**Distribution Agreement**

In presenting this dissertation as a partial fulfillment of the requirements for an advanced degree from Emory University, I hereby grant to Emory University and its agents the non-exclusive license to archive, make accessible, and display my dissertation in whole or in part in all forms of media, now or hereafter known, including display on the world wide web. I understand that I may select some access restrictions as part of the online submission of this dissertation. I retain all ownership rights to the copyright of the dissertation. I also retain the right to use in future works (such as articles or books) all or part of this dissertation.

March 19, 2025

_____          _____

Shifan Zhao                                                Date

Robust Uncertainty Quantification for Foundation Models: Bayesian and Frequentist
Approaches for High-Stakes Applications

By

Shifan Zhao
Doctor of Philosophy

Department of Mathematics

_____
Talea L. Mayo
Advisor

_____
James G. Nagy
Committee Member

_____
Matthias Chung
Committee Member

Accepted:

_____
Kimberly Jacob Arriola, PhD, MPH,
Dean of the James T. Laney School of Graduate Studies

March 19, 2025
_____
Date

Robust Uncertainty Quantification for Foundation Models: Bayesian and Frequentist
Approaches for High-Stakes Applications

By

Shifan Zhao

Advisor: Talea L. Mayo

An abstract of
A dissertation submitted to the Faculty of the
James T. Laney School of Graduate Studies of Emory University
in partial fulfillment of the requirements for the degree of
Doctor of Philosophy
in Department of Mathematics
2025

Abstract

Robust Uncertainty Quantification for Foundation Models: Bayesian and Frequentist
Approaches for High-Stakes Applications
By Shifan Zhao


Machine learning foundation models have demonstrated impressive predictive capabilities across various domains, including healthcare and climate science. However, their deterministic nature limits their utility in high-stakes applications where understanding prediction uncertainty is crucial for responsible decision-making. This thesis addresses this critical gap by developing two complementary approaches to uncertainty quantification (UQ) for foundation models.

First, we introduce a novel two-stage Gaussian Process methodology that effectively handles mean and kernel misspecification—a common challenge in real-world applications. This approach separates mean prediction from uncertainty quantification, leading to more reliable uncertainty estimates even with limited data. We demonstrate its application to healthcare foundation models for patient risk prediction, where accurate uncertainty bounds can significantly impact clinical decision-making.

Second, we develop a Locally Debiased Adaptive Conformal Prediction (LC-ACP) framework that provides distribution-free coverage guarantees without requiring exchangeability assumptions, making it particularly valuable for non-stationary time series. We apply this methodology to climate foundation models for hurricane track prediction, where reliable uncertainty quantification directly impacts emergency management and resource allocation during extreme weather events.

To address computational challenges, we introduce kernel preconditioning techniques and unbiased estimators that significantly reduce the cubic complexity of Gaussian Processes while maintaining accuracy. Through comprehensive experiments across healthcare and climate domains, we demonstrate that our methods provide well-calibrated uncertainty estimates across diverse applications and data characteristics. This thesis contributes both theoretical advances and practical implementations that bridge the gap between powerful predictive models and responsible deployment in high-stakes, real-world applications.

Robust Uncertainty Quantification for Foundation Models: Bayesian and Frequentist
Approaches for High-Stakes Applications

By

Shifan Zhao

Advisor: Talea L. Mayo

A dissertation submitted to the Faculty of the
James T. Laney School of Graduate Studies of Emory University
in partial fulfillment of the requirements for the degree of
Doctor of Philosophy
in Department of Mathematics
2025

# Contents

# List of Figures

# List of Tables

# List of Algorithms

# Chapter 1

# Background and Motivation

> "Uncertainty must be taken in a sense radically distinct from the familiar notion of Risk, from which it has never been properly separated... The essential fact is that 'risk' means in some cases a quantity susceptible of measurement, while at other times it is something distinctly not of this character; and there are far-reaching and crucial differences in the bearings of the phenomena depending on which of the two is really present and operating... It will appear that a measurable uncertainty, or 'risk' proper, as we shall use the term, is so far different from an unmeasurable one that it is not in effect an uncertainty at all."
>
> *— Frank H. Knight, Risk, Uncertainty, and Profit*

## 1.1   Overview and Motivation

The ability to make reliable predictions under uncertainty is fundamental to scientific inquiry and decision-making across domains. While machine learning, particularly through foundation models, has revolutionized predictive capabilities, it often fails to quantify the uncertainty in its predictions—a critical limitation in high-stakes applications. This thesis addresses this gap by developing novel uncertainty quantification methods for foundation models, combining Gaussian Process approaches with conformal prediction techniques. Our work spans theoretical foundations, algorithmic innovations, and practical applications in healthcare and climate science, domains where reliable uncertainty estimates can significantly impact human well-being. In this chapter, we establish the background and motivation for our research, examining the limitations of current approaches and highlighting the need for robust

uncertainty quantification methods that can be deployed in real-world, high-impact settings.

## 1.2 Introduction to Deep Learning Modeling and Its Limitations

Deep learning has revolutionized the field of artificial intelligence, enabling remarkable advances in various domains including computer vision, natural language processing, and healthcare. The ability of deep neural networks to automatically learn hierarchical representations from data has led to state-of-the-art performance in numerous prediction tasks. In particular, pre-trained foundation models (PFMs) have emerged as powerful tools that can be fine-tuned for specific downstream tasks with limited labeled data [17, 19, 97].

Despite their impressive predictive capabilities, deep learning models, including PFMs, suffer from several critical limitations. One of the most significant limitations is their inability to quantify uncertainty in their predictions. Traditional deep learning models typically produce deterministic point estimates without providing any measure of confidence or reliability [18]. This limitation becomes particularly problematic in high-stakes applications such as healthcare, climate modeling, and autonomous systems, where incorrect predictions can have severe and even fatal consequences.

Consider a patient risk prediction model that estimates the probability of a patient developing a certain medical condition. The model might predict a 70% probability, but without any indication of its confidence in this prediction. Is this 70% a highly certain estimate based on abundant similar cases in the training data, or is it a highly uncertain estimate due to the patient having unusual characteristics not well-represented in the training data? Without proper uncertainty quantification, clinicians cannot make fully informed decisions about patient care.

Similarly, in climate and weather prediction, deterministic forecasts fail to capture the inherent uncertainty in chaotic weather systems. A weather model might predict a hurricane's

path, but without quantifying the uncertainty in this prediction, emergency management officials cannot properly assess the risk and make appropriate evacuation decisions.



Figure 1.1: The problem with deterministic deep learning models. Traditional models provide only point estimates without capturing the uncertainty in predictions. This is particularly problematic for out-of-distribution inputs (red region) where the model might be confidently wrong, leading to potentially catastrophic decisions in high-stakes applications.

The limitations of deterministic deep learning models can be summarized as follows:

1. **Overconfidence in predictions**: Deep learning models often produce overconfident predictions, especially for out-of-distribution inputs [69].

2. **Inability to distinguish between aleatoric and epistemic uncertainty**: Models cannot differentiate between uncertainty due to inherent randomness in the data (aleatoric) and uncertainty due to model limitations (epistemic).

3. **Lack of calibration**: Predicted probabilities often do not correspond to the true frequencies of events [33].

4. **Difficulty in detecting out-of-distribution inputs**: Models typically fail to recognize when they are making predictions on data that differs significantly from their

training distribution [39].

These limitations highlight the need for uncertainty-aware models that can provide reliable uncertainty estimates alongside their predictions. This thesis addresses these limitations by developing and applying uncertainty quantification methods based on Gaussian processes and conformal prediction to enhance the reliability and trustworthiness of deep learning models, particularly in high stakes, safety-critical applications.

## 1.3 Uncertainty Quantification: Foundations and Importance

### 1.3.1 Formal Definition and Mathematical Framework

Uncertainty quantification (UQ) is the science of identifying, quantifying, and reducing uncertainties in computational and real-world systems. In the context of machine learning, UQ aims to provide a measure of confidence or reliability in model predictions. Formally, instead of producing a point estimate $\hat{y} = f(x)$ for an input $x$, an uncertainty-aware model produces a predictive distribution $p(y|x)$ that captures the uncertainty in the prediction.

For regression tasks, this predictive distribution is often characterized by its mean $\mu(x)$ and variance $\sigma^2(x)$:

$$p(y|x) = \mathcal{N}(\mu(x), \sigma^2(x)). \tag{1.1}$$

For classification tasks, the predictive distribution is typically a categorical distribution over the possible classes:

$$p(y = c|x) = \frac{\exp(z_c)}{\sum_{c'} \exp(z_{c'})}, \tag{1.2}$$

where $z_c$ represents the logit for class $c$. In uncertainty-aware classification, these logits themselves are treated as random variables, leading to stochastic predictions that reflect the model's uncertainty.

## 1.3.2 Types of Uncertainty

Uncertainty in machine learning predictions can be categorized into three main types:

1. **Aleatoric Uncertainty**: This type of uncertainty arises from inherent randomness or noise in the data. It cannot be reduced by collecting more data or improving the model. For example, in a medical diagnosis task, two patients with identical observable characteristics might have different outcomes due to unobserved genetic factors or random variations in disease progression.

2. **Epistemic Uncertainty**: This uncertainty stems from the model's lack of knowledge or limitations in capturing the true data-generating process. Unlike aleatoric uncertainty, epistemic uncertainty can be reduced by collecting more data or improving the model. It is particularly high in regions of the input space where training data is sparse.

3. **Total Predictive Uncertainty**: This is the combination of aleatoric and epistemic uncertainty, representing the overall uncertainty in the model's predictions.

To illustrate these concepts mathematically, consider a classification problem with three classes. We have a dataset $\mathcal{D} = \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{test}}$, where $\mathcal{D}_{\text{train}} = \{(x_j, y_j)\}_{j=1}^{N}$ represents the training set and $\mathcal{D}_{\text{test}} = \{(x_j, y_j)\}_{j=N+1}^{N+M}$ represents the test set. We use a model $\mathcal{M}_\theta$ with parameters $\theta$. Given a test data point $x_{\text{test}}$, the total uncertainty in the classification problem can be quantified as follows:

$$p(\omega_i|x_{\text{test}}, \mathcal{D}_{\text{train}}) = \int p(\omega_i|x_{\text{test}}, \theta)p(\theta|\mathcal{D}_{\text{train}})d\theta, \tag{1.3}$$

where $\omega_i$ is any random variable. The components of this uncertainty can be decomposed into:

- $p(\omega_i|x_{\text{test}}, \mathcal{D}_{\text{train}})$: Total predictive uncertainty, which represents the overall uncertainty in our prediction.

Figure 1.2: Visual representation of aleatoric and epistemic uncertainty. Aleatoric uncertainty (left) represents inherent randomness in data that cannot be reduced with more observations. Epistemic uncertainty (right) represents model knowledge gaps that decrease as more data is collected.

- $p(\omega_i|x_{\text{test}}, \theta)$: Data uncertainty (aleatoric uncertainty), which stems from inherent noise in the data generation process.

- $p(\theta|\mathcal{D}_{\text{train}})$: Model uncertainty (epistemic uncertainty), which reflects our limited knowledge of the true model parameters. As Knight's distinction emphasizes, this type of uncertainty can be viewed as "risk" since it can be reduced with additional data.

Once we train the model $\mathcal{M}_\theta$ on the training data $\mathcal{D}_{\text{train}}$, the model uncertainty $p(\theta|\mathcal{D}_{\text{train}})$ induces uncertainty over $p(\omega_i|x_{\text{test}}, \theta)$. This highlights a fundamental aspect of UQ: epistemic uncertainty can be reduced with better models or more data, while aleatoric uncertainty remains irreducible as it arises from the inherent randomness in the data generation process.

### 1.3.3  Known-Unknowns vs. Unknown-Unknowns

Uncertainty can further be divided into known-unknowns and unknown-unknowns. Assume the training data comes from a data space $\mathcal{X}_{\text{in}}$, but the test data may also come from an out-of-domain space $\mathcal{X}_{\text{out}}$. For instance, we might train a language model on Shakespeare's

Figure 1.3: Complete visualization of uncertainty types. Aleatoric uncertainty (orange) is irreducible randomness in data, epistemic uncertainty (blue) represents model knowledge gaps that decrease with more data, and total uncertainty (purple) combines both types. The formula $\sigma_{\text{total}} = \sqrt{\sigma_{\text{aleatoric}}^2 + \sigma_{\text{epistemic}}^2}$ shows how they relate mathematically.

sonnets, but a user might input a poem by another poet, which represents out-of-domain data.

To express this more specifically, we can rewrite the total uncertainty as:

$$p(\omega_i|x_{\text{test}}, \mathcal{D}_{\text{train}}) = p(\omega_i|x_{\text{test}} \in \mathcal{X}_{\text{in}}, \mathcal{D}_{\text{train}})p(x_{\text{test}} \in \mathcal{X}_{\text{in}})$$
$$+ p(\omega_i|x_{\text{test}} \in \mathcal{X}_{\text{out}}, \mathcal{D}_{\text{train}})p(x_{\text{test}} \in \mathcal{X}_{\text{out}})$$

The first term represents known-unknowns due to overlapping data or irreducible noise such as measurement errors. The second term represents unknown-unknowns, which arise due to distributional shifts between training and testing datasets.

### 1.3.4   Calibration and Sharpness Metrics

Two key metrics for evaluating uncertainty quantification methods are calibration and sharpness [90]:

1. **Calibration**: A model is well-calibrated if its predicted probabilities match the empirical frequencies of events [33]. For example, if a model predicts a 70% probability of rain for 100 different days, it should rain on approximately 70 of those days. Calibration can be measured using metrics such as Expected Calibration Error (ECE) and Maximum Calibration Error (MCE).

2. **Sharpness**: A model is sharp if its predictive distributions have low entropy or variance [13]. Sharpness measures how concentrated the predictive distributions are. A model that always predicts a 50% probability for binary outcomes is perfectly calibrated but not sharp. Ideally, a model should be both well-calibrated and sharp.

The trade-off between calibration and sharpness is a central challenge in uncertainty quantification. A model that always predicts the marginal distribution of the target variable will be perfectly calibrated but not sharp. Conversely, a model that makes very confident predictions might be sharp but poorly calibrated if those predictions are often wrong.

### 1.3.5   Importance in High-Stakes Decision-Making

Reliable uncertainty quantification is crucial in high-stakes decision-making scenarios where the cost of errors is significant. The growing adoption of foundation models in critical applications has made this need even more pressing. Key application domains include:

1. **Healthcare**: In medical diagnosis and treatment planning [61], understanding the uncertainty in predictions can help clinicians make more informed decisions and potentially save lives. Recent applications include:

- **Risk Stratification**: Models that predict patient deterioration must quantify their uncertainty to help clinicians prioritize care and allocate resources effectively.

- **Treatment Response Prediction**: When predicting a patient's response to different treatment options, uncertainty estimates help clinicians weigh the risks and benefits of each option.

- **Medical Imaging**: In radiology and pathology, models must indicate when they are uncertain about their diagnoses, prompting additional review by human experts.

- **Drug Discovery**: In pharmaceutical research, uncertainty quantification helps prioritize which compounds to investigate further, potentially accelerating the development of new treatments.



Figure 1.4: Application of uncertainty quantification in healthcare. Diagnostic and prognostic models benefit from reliable uncertainty estimates to support clinical decision-making. The confidence intervals help clinicians assess the reliability of predictions, potentially leading to better treatment decisions and improved patient outcomes.

2. **Climate Science**: Climate models are inherently uncertain due to the chaotic nature

of weather systems and the complexity of climate dynamics [36]. Applications include:

- **Extreme Weather Prediction**: Models must provide reliable uncertainty estimates for the timing, location, and intensity of extreme weather events [56] to inform evacuation decisions.

- **Long-term Climate Projections**: When modeling climate change impacts [10], uncertainty quantification helps policymakers understand the range of possible outcomes and their likelihoods.

- **Agricultural Planning**: Farmers and agricultural businesses rely on weather forecasts with uncertainty estimates to make planting and harvesting decisions.

- **Resource Management**: Water resource managers use uncertainty-aware predictions to optimize reservoir operations and flood control measures.



Figure 1.5: Application of uncertainty quantification in climate science. Weather forecasting models benefit from uncertainty quantification, particularly for hurricane track prediction where uncertainty cones help emergency managers make evacuation decisions. The confidence intervals widen with forecast lead time, reflecting the chaotic nature of weather systems.

3. **Autonomous Systems**: Self-driving vehicles and other autonomous systems must make real-time decisions in uncertain environments. Critical applications include:

   - **Perception Systems**: Object detection and tracking systems must quantify their uncertainty about object locations and classifications to avoid collisions.

   - **Path Planning**: Navigation algorithms need to account for uncertainty in their environment models and predictions of other agents' behavior.

   - **Control Systems**: Actuator commands must be adjusted based on the uncertainty in state estimates and predicted outcomes.

   - **Safety Systems**: Emergency intervention systems must reliably detect when they are operating outside their comfort zone.

4. **Financial Risk Management**: In financial markets [87], quantifying uncertainty is essential for:

   - **Portfolio Optimization**: Asset allocation decisions must account for uncertainties in return predictions and risk estimates.

   - **Risk Assessment**: Credit scoring and fraud detection systems must quantify their uncertainty to set appropriate thresholds.

   - **Algorithmic Trading**: Trading strategies must adjust their positions based on the uncertainty in price predictions.

   - **Regulatory Compliance**: Stress testing and risk reporting require reliable uncertainty estimates.

5. **Scientific Research**: Modern scientific workflows increasingly rely on machine learning models with uncertainty quantification [24] for:

   - **Experimental Design**: Optimizing experimental parameters while accounting for uncertainties in predictions.

- **Data Analysis**: Identifying significant results while controlling for uncertainty in measurements and model predictions.

- **Hypothesis Generation**: Using model uncertainty to guide the exploration of scientific hypotheses.

- **Reproducibility**: Quantifying uncertainty helps assess the robustness and reliability of scientific findings.

In all these applications, the goal extends beyond making accurate predictions to providing decision-makers with a complete picture of the uncertainty associated with those predictions. This enables:

- **Risk-Aware Decision Making**: Decision-makers can appropriately weight the consequences of different actions against their uncertainties.

- **Resource Allocation**: Limited resources can be directed to cases where the model is most uncertain and human expertise is most needed.

- **Continuous Improvement**: By tracking where models are most uncertain, organizations can prioritize relevant data collection and model improvements.

- **Trust Building**: Transparent communication of uncertainty helps build trust between AI systems and their users.

The rise of foundation models has introduced new challenges in uncertainty quantification, as these models are increasingly deployed in high-stakes applications despite potential limitations in their reliability and robustness. This motivates our work on developing scalable and reliable uncertainty quantification methods that can be applied to foundation models while maintaining their computational efficiency.

# 1.4 Uncertainty Quantification in Machine Learning

Various approaches have been developed to quantify uncertainty in machine learning models. These approaches can be broadly categorized into traditional statistical methods, Bayesian approaches, ensemble methods, and evidential deep learning.

## 1.4.1 Traditional Statistical Approaches

Traditional statistical approaches to uncertainty quantification include:

1. **Confidence Intervals**: In frequentist statistics, confidence intervals provide a range of values that likely contains the true parameter with a specified probability [34]. For example, a 95% confidence interval for a regression coefficient means that if we were to repeat the experiment many times, about 95% of the intervals would contain the true coefficient.

2. **Prediction Intervals**: While confidence intervals quantify uncertainty in parameter estimates, prediction intervals quantify uncertainty in future observations [103]. They account for both the uncertainty in the parameter estimates and the inherent randomness in the data.

3. **Delta Method**: This method uses a Taylor series approximation to estimate the variance of a function of random variables [9]. It can be used to derive confidence intervals for complex functions of model parameters.

These traditional approaches often rely on strong assumptions about the data distribution and model form, which often do not hold in complex deep learning models.

## 1.4.2 Bayesian Neural Networks

Bayesian Neural Networks (BNNs) [104] extend traditional neural networks by placing prior distributions over the network weights and learning posterior distributions over these weights

given the observed data. Instead of learning point estimates of the weights, BNNs learn distributions over the weights, allowing them to capture epistemic uncertainty.

Computing the exact posterior distribution is generally intractable for deep neural networks, so various approximation methods are used, including:

1. **Markov Chain Monte Carlo (MCMC)**: Methods like Hamiltonian Monte Carlo (HMC) [20] and No-U-Turn Sampler (NUTS) can be used to sample from the posterior distribution, but they are computationally expensive for large networks.

2. **Variational Inference**: This approach [27] approximates the posterior distribution with a simpler distribution (e.g., a Gaussian) and optimizes the parameters of this approximation to minimize the Kullback-Leibler divergence from the true posterior.

3. **Monte Carlo Dropout**: This method [26] interprets dropout, a regularization technique in neural networks, as a Bayesian approximation. By keeping dropout active during inference and performing multiple forward passes, it generates samples from an approximate posterior distribution.

BNNs provide a principled approach to uncertainty quantification but can be computationally expensive and challenging to implement for large-scale models.

### 1.4.3   Ensemble Methods

Ensemble methods combine predictions from multiple models to improve accuracy and quantify uncertainty [52]. The variance in predictions across ensemble members can be used as a measure of epistemic uncertainty. Common ensemble methods include:

1. **Deep Ensembles**: This approach trains multiple neural networks with different random initializations and combines their predictions [52]. Despite its simplicity, deep ensembles have been shown to provide competitive uncertainty estimates compared to more sophisticated Bayesian methods.

2. **Bootstrapped Ensembles**: These methods train each ensemble member on a boot-strapped sample of the training data [46], capturing the uncertainty due to sampling variability.

3. **Snapshot Ensembles**: This approach saves model snapshots at different points during training (e.g., at different local minima) [58] and uses them as ensemble members, reducing the computational cost compared to training multiple models from scratch.

Ensemble methods are relatively easy to implement and can provide robust uncertainty estimates, but they can be computationally expensive during both training and inference.

## 1.4.4  Evidential Deep Learning

Evidential Deep Learning (EDL) [84] is a recent approach that extends deep learning models to output the parameters of a higher-order probability distribution (e.g., a Dirichlet distribution for classification [65] or a Normal-Inverse-Gamma distribution for regression [2]) rather than directly outputting the parameters of the predictive distribution.

For regression tasks, EDL models output the parameters of a Normal-Inverse-Gamma distribution [2], which serves as a prior for the mean and variance of the predictive Normal distribution. The predictive distribution is then a Student's t-distribution, which has heavier tails than a Normal distribution, reflecting the additional uncertainty.

For classification tasks, EDL models output the parameters of a Dirichlet distribution [84], which serves as a prior for the categorical distribution over classes. The predictive distribution is then a Dirichlet-Multinomial distribution, which can express uncertainty about the class probabilities.

EDL provides a computationally efficient approach to uncertainty quantification that can capture both aleatoric and epistemic uncertainty without requiring multiple forward passes or ensemble members.

## 1.5 Gaussian Processes for Uncertainty Quantification

### 1.5.1 Mathematical Foundations

Gaussian Processes (GPs) provide a principled Bayesian approach to uncertainty quantification in regression and classification tasks [80]. A GP is a collection of random variables, any finite number of which have a joint Gaussian distribution. It is fully specified by a mean function $m(x)$ and a covariance function (or kernel) $k(x, x')$ [103]:

$$f(x) \sim \mathcal{GP}(m(x), k(x, x')). \tag{1.4}$$



Figure 1.6: Gaussian Processes as a solution for uncertainty quantification. Unlike deterministic models, GPs provide a full probability distribution over possible function values, quantifying both aleatoric and epistemic uncertainty. The shaded region represents the confidence interval ($\pm 2$ standard deviations), which naturally widens in regions with sparse data, reflecting increased epistemic uncertainty.

The mean function represents our prior belief about the expected value of the function at any point, while the covariance function captures the correlation between function values at different points.

Given a dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ where $y_i = f(x_i) + \epsilon_i$ with $\epsilon_i \sim \mathcal{N}(0, \sigma_n^2)$, the posterior distribution over function values at a new point $x_*$ is also Gaussian:

$$p(f_* | \mathcal{D}, x_*) = \mathcal{N}(\mu_*, \sigma_*^2) \tag{1.5}$$

$$\mu_* = m(x_*) + k(x_*, X)[k(X, X) + \sigma_n^2 I]^{-1}(y - m(X)) \tag{1.6}$$

$$\sigma_*^2 = k(x_*, x_*) - k(x_*, X)[k(X, X) + \sigma_n^2 I]^{-1}k(X, x_*), \tag{1.7}$$

where $X = [x_1, x_2, \ldots, x_n]^T$, $y = [y_1, y_2, \ldots, y_n]^T$, $k(X, X)$ is the covariance matrix with entries $k(x_i, x_j)$, and $k(x_*, X)$ is the vector of covariances between $x_*$ and each training point.

### 1.5.2 Mean Function and Kernel Description

The mean function $m(x)$ represents our prior belief about the expected value of the function [45]. Common choices include:

1. **Zero Mean**: $m(x) = 0$, which assumes no prior knowledge about the function values.

2. **Constant Mean**: $m(x) = c$, which assumes a constant prior mean across the input space.

3. **Linear Mean**: $m(x) = \beta^T x$, which assumes a linear trend in the function values.

The covariance function (or kernel) $k(x, x')$ captures the similarity between points in the input space [48]. Common kernels include:

1. **Radial Basis Function (RBF) Kernel**: $k(x, x') = \sigma_f^2 \exp\left(-\frac{||x - x'||^2}{2l^2}\right)$, which assumes smooth functions with a characteristic length scale $l$ and signal variance $\sigma_f^2$.

2. **Matérn Kernel**: A family of kernels that can model functions with different degrees of smoothness, controlled by a parameter $\nu$.

3. **Periodic Kernel**: $k(x, x') = \sigma_f^2 \exp\left(-\frac{2\sin^2(\pi \|x - x'\|/p)}{l^2}\right)$, which models periodic functions with period $p$.

4. **Linear Kernel**: $k(x, x') = \sigma_f^2 x^T x'$, which models linear functions.

The choice of kernel encodes our prior beliefs about the properties of the function, such as smoothness, periodicity, and linearity [21].

### 1.5.3 Posterior Predictive Distribution

The posterior predictive distribution for a new observation $y_*$ at $x_*$ is given by [80]:

$$p(y_*|\mathcal{D}, x_*) = \mathcal{N}(\mu_*, \sigma_*^2 + \sigma_n^2) \tag{1.8}$$

This distribution captures both the epistemic uncertainty (through $\sigma_*^2$) and the aleatoric uncertainty (through $\sigma_n^2$) [49]. The epistemic uncertainty decreases as we observe more data, especially near the observed points, while the aleatoric uncertainty remains constant.

### 1.5.4 Assumptions, Strengths, and Weaknesses

Gaussian Processes make several key assumptions [99]:

1. **Gaussian Noise**: The observation noise is assumed to be Gaussian, which may not hold in all applications.

2. **Stationarity**: Many commonly used kernels assume that the covariance between points depends only on their relative positions, not their absolute positions.

3. **Smoothness**: Most kernels assume some degree of smoothness in the underlying function, which may not be appropriate for all applications.

The strengths of Gaussian Processes include [57]:

1. **Principled Uncertainty Quantification**: GPs provide a principled Bayesian approach to uncertainty quantification, capturing both aleatoric and epistemic uncertainty.

2. **Flexibility**: By choosing appropriate kernels, GPs can model a wide range of functions with different properties.

3. **No Need for Large Datasets**: GPs can provide reliable uncertainty estimates even with small datasets, unlike deep learning models that typically require large amounts of data.

The weaknesses of Gaussian Processes include [40]:

1. **Computational Complexity**: Standard GP inference has $O(n^3)$ time complexity and $O(n^2)$ space complexity, making it challenging to scale to large datasets.

2. **Sensitivity to Hyperparameters**: The performance of GPs can be sensitive to the choice of hyperparameters, which need to be carefully tuned.

3. **Difficulty with High-Dimensional Inputs**: GPs can struggle with high-dimensional inputs due to the curse of dimensionality, although this can be mitigated with appropriate kernel design.

### 1.5.5   Sparse Gaussian Processes

As highlighted in the previous section, standard Gaussian Processes suffer from cubic computational complexity $(O(n^3))$ with respect to the number of training points, making them impractical for large datasets. Sparse Gaussian Process approximations have been developed to address this computational challenge while preserving the desirable properties of GPs [76, 94].

The key idea behind sparse GP approximations is to use a small set of $m \ll n$ inducing points (also called pseudo-inputs) to approximate the full GP. These inducing points serve

Figure 1.7: Key challenges in scaling Gaussian Processes to real-world applications. The cubic computational complexity limits their applicability to large datasets, while kernel selection and hyperparameter tuning require expertise. High-dimensional data suffers from the curse of dimensionality, and non-Gaussian likelihoods complicate inference. These challenges motivate the methodological contributions in this thesis.

as a concise summary of the full dataset, enabling inference with reduced computational complexity, typically $O(nm^2 + m^3)$ [40, 11].

## Inducing Point Methods

Let $\mathbf{Z} = \{\mathbf{z}_1, \mathbf{z}_2, \ldots, \mathbf{z}_m\}$ be a set of inducing inputs and $\mathbf{u} = \{f(\mathbf{z}_1), f(\mathbf{z}_2), \ldots, f(\mathbf{z}_m)\}$ be the corresponding function values. The joint distribution of the latent function values $\mathbf{f}$ at the training points $\mathbf{X}$ and the inducing points $\mathbf{u}$ is:

$$
\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mathbf{m}(\mathbf{X}) \\ \mathbf{m}(\mathbf{Z}) \end{bmatrix}, \begin{bmatrix} \mathbf{K_{XX}} & \mathbf{K_{XZ}} \\ \mathbf{K_{ZX}} & \mathbf{K_{ZZ}} \end{bmatrix} \right) \tag{1.9}
$$

Using the properties of conditional Gaussian distributions, we can express $\mathbf{f}$ conditioned on $\mathbf{u}$:

$$p(\mathbf{f}|\mathbf{u}) = \mathcal{N}(\mathbf{m}(\mathbf{X}) + \mathbf{K_{XZ}}\mathbf{K_{ZZ}^{-1}}(\mathbf{u} - \mathbf{m}(\mathbf{Z})), \mathbf{K_{XX}} - \mathbf{K_{XZ}}\mathbf{K_{ZZ}^{-1}}\mathbf{K_{ZX}}) \qquad (1.10)$$

Different sparse GP methods make different approximations based on this formulation:

1. **Subset of Regressors (SoR)**: This method [88] makes the extreme approximation that the training points are conditionally independent given the inducing points, effectively replacing the full covariance matrix with a low-rank approximation:

$$\mathbf{K_{XX}} \approx \mathbf{K_{XZ}}\mathbf{K_{ZZ}^{-1}}\mathbf{K_{ZX}} \qquad (1.11)$$

2. **Deterministic Training Conditional (DTC)**: Similar to SoR but maintains exact predictive means while approximating the predictive variance [83].

3. **Fully Independent Training Conditional (FITC)**: This method [89] retains the diagonal elements of the exact covariance matrix, addressing the tendency of SoR to underestimate predictive variances:

$$\mathbf{K_{XX}} \approx \mathrm{diag}(\mathbf{K_{XX}} - \mathbf{Q_{XX}}) + \mathbf{Q_{XX}} \quad \text{where} \quad \mathbf{Q_{XX}} = \mathbf{K_{XZ}}\mathbf{K_{ZZ}^{-1}}\mathbf{K_{ZX}} \qquad (1.12)$$

**Variational Sparse Gaussian Processes**

Variational Sparse GPs [94] take a principled Bayesian approach by formulating sparse GP approximation as a variational inference problem. Instead of heuristically modifying the GP prior, this approach minimizes the KL divergence between the approximate posterior and the true posterior.

The key insight is to introduce the inducing variables $\mathbf{u}$ as variational parameters and integrate them out to obtain a lower bound on the marginal likelihood:

$$\log p(\mathbf{y}|\mathbf{X}) \geq \mathcal{L}_{\text{ELBO}} = \log \mathcal{N}(\mathbf{y}|\mathbf{K_{XZ}}\mathbf{K_{ZZ}^{-1}}\mathbf{m}(\mathbf{Z}), \sigma^2\mathbf{I}) - \frac{1}{2\sigma^2}\text{Tr}(\mathbf{K_{XX}} - \mathbf{Q_{XX}}) \qquad (1.13)$$

The trace term acts as a regularizer that penalizes the approximation error in the covariance matrix, leading to optimized inducing point locations that minimize this error [66].

**Stochastic Variational Inference for Scalable GPs**

Stochastic Variational Gaussian Processes (SVGPs) [40, 41] extend the variational approach to enable minibatch training, making GPs applicable to massive datasets. The key innovation is to introduce an explicit variational distribution over the inducing variables:

$$q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S}) \qquad (1.14)$$

where $\mathbf{m}$ and $\mathbf{S}$ are variational parameters. The evidence lower bound (ELBO) becomes:

$$\mathcal{L}_{\text{ELBO}} = \sum_{i=1}^{n} \mathbb{E}_{q(f_i)}[\log p(y_i|f_i)] - \text{KL}(q(\mathbf{u})||p(\mathbf{u})) \qquad (1.15)$$

where $q(f_i) = \int p(f_i|\mathbf{u})q(\mathbf{u})d\mathbf{u}$ is the marginal distribution of $f_i$ under the variational approximation. This objective can be optimized using stochastic gradient methods with minibatches, achieving complexity of $O(bm^2 + m^3)$ per iteration, where $b$ is the minibatch size [42].

**Structured Sparse Approximations**

Several approaches exploit structured data to further improve sparse GP scalability:

1. **Kronecker and Toeplitz Methods**: For data on regular grids, the kernel matrix can have Kronecker or Toeplitz structure, enabling efficient exact or approximate computations [105].

2. **Hierarchical Approximations**: Methods like hierarchical off-diagonal low-rank (HODLR) matrices [1] and H-matrices [35] provide systematic approximations with controlled accuracy.

3. **Random Fourier Features**: This approach approximates the kernel function using random Fourier features [77], converting the GP into an approximately equivalent linear model with complexity $O(nmr)$, where $r$ is the number of random features.

**Strengths and Limitations**

Sparse GPs offer several advantages:

1. **Reduced Computational Complexity**: From $O(n^3)$ to $O(nm^2 + m^3)$, enabling application to larger datasets.

2. **Preserved Uncertainty Quantification**: Unlike simpler approximations, sparse GPs maintain meaningful uncertainty estimates, critical for high-stakes applications.

3. **Principled Approximation Framework**: Particularly for variational methods, approximations come with theoretical guarantees and interpretable objectives.

However, limitations remain:

1. **Approximation Quality**: The quality of the approximation depends critically on the number and locations of inducing points, requiring careful selection or optimization.

2. **Kernel Limitations**: Sparse approximations do not address other GP limitations like sensitivity to kernel choice and hyperparameters.

3. **Complex Implementation**: Implementing efficient sparse GP approximations, particularly with non-Gaussian likelihoods, can be technically challenging.

4. **Memory Requirements**: While computational complexity is reduced, memory requirements for storing kernel matrices can still be substantial for very large datasets.

In later chapters, we will build upon these sparse GP approximations to develop our two-stage GP methodology, addressing both computational challenges and fundamental limitations of standard GP approaches.

## 1.6 Brief Introduction to Conformal Prediction

Conformal Prediction (CP) [5] is a framework for constructing prediction intervals with guaranteed coverage under minimal assumptions. Unlike Bayesian methods that require specifying prior distributions, CP only assumes that the data is exchangeable (a weaker assumption than independence and identical distribution).

The key idea of CP is to use a nonconformity measure to quantify how different a new example is from the training examples [29]. By comparing the nonconformity score of a potential label for a new example with the nonconformity scores of the training examples, CP can determine whether to include that label in the prediction set.

For regression tasks, CP typically outputs prediction intervals $[L(x), U(x)]$ such that [108]:

$$P(y \in [L(x), U(x)]) \geq 1 - \alpha \tag{1.16}$$

where $1 - \alpha$ is the desired coverage level (e.g., 95%).

For classification tasks, CP outputs a set of candidate labels $\Gamma_\alpha(x)$ such that:

$$P(y \in \Gamma_\alpha(x)) \geq 1 - \alpha \tag{1.17}$$

CP has several advantages [92]:

1. **Distribution-Free Guarantees**: CP provides valid prediction intervals without assuming a specific data distribution.

2. **Model-Agnostic**: CP can be applied to any predictive model, including deep neural networks.

3. **Computational Efficiency**: Inductive CP, a variant of CP, is computationally efficient and can be applied to large datasets.

However, CP also has limitations [43]:

1. **Marginal Coverage**: CP guarantees marginal coverage across all examples but not conditional coverage for specific examples.

2. **Calibration Set Requirement**: Inductive CP requires setting aside a portion of the data for calibration, which can be a limitation with small datasets.

In this thesis, we will explore how CP can complement Gaussian Processes to provide robust uncertainty quantification, particularly in non-stationary settings where GP assumptions may not hold.

## 1.7 Recent Advances in Uncertainty Quantification for Foundation Models

Recent years have seen significant advances in uncertainty quantification methods specifically designed for foundation models [17, 19, 97] and large language models (LLMs). These advances address unique challenges posed by the scale and complexity of these models.

### 1.7.1 Uncertainty in Foundation Models

Foundation models present unique challenges for uncertainty quantification due to their [109]:

1. **Scale and Complexity**: With billions of parameters, traditional Bayesian approaches become computationally intractable.

2. **Black-box Nature**: Many foundation models are accessible only through APIs, limiting the applicability of methods that require access to model internals.

3. **Transfer Learning Setting**: The pre-train-then-fine-tune paradigm introduces additional sources of uncertainty.

## 1.7.2  Recent Methodological Advances

Several innovative approaches have emerged to address these challenges:

1. **Bayesian Prompt Ensembles**: Recent work by [95] introduces a novel approach for black-box LLMs that uses weighted ensembles of semantically equivalent prompts to estimate uncertainty.

2. **Distributional Conformal Prediction**: [29] developed probabilistically robust conformal prediction, which ensures robustness to most perturbations around clean input examples.

3. **Scalable Bayesian Deep Learning**: [20] introduced a tempered framing of stochastic gradient MCMC that transitions seamlessly into optimization.

4. **Uncertainty-Aware Foundation Models**: [61] demonstrated how Gaussian Process layers can be integrated with pre-trained foundation models to enable instance-level uncertainty quantification.

## 1.7.3  Applications in High-Stakes Domains

Recent applications have demonstrated the practical impact of uncertainty quantification in foundation models:

1. **Healthcare**: [73] developed distributional Gaussian Process layers for medical imaging that reliably separate in-distribution from out-of-distribution cases, crucial for safe deployment in clinical settings.

2. **Climate Science**: [10] systematically compared different uncertainty quantification methods for data-driven weather models, showing how probabilistic forecasts can improve upon deterministic predictions from state-of-the-art models.

3. **Scientific Computing**: [37] developed statistically optimal methods for uncertainty quantification in expensive black-box models, particularly relevant for scientific simulations and computational physics.

### 1.7.4 Emerging Trends and Open Challenges

Several key trends and challenges have emerged in recent research:

1. **Efficiency-Accuracy Trade-off**: Recent work by [59] on orthogonal bootstrap demonstrates the ongoing effort to reduce computational costs while maintaining statistical guarantees.

2. **Distribution Shift**: [92] highlighted the challenge of maintaining valid uncertainty estimates under covariate shift, proposing conformal prediction as a potential solution.

3. **Model Evaluation**: [109] introduced uncertainty-aware benchmarking for LLMs, showing that larger models may exhibit greater uncertainty despite higher accuracy.

4. **Calibration Quality**: [13] developed methods for sharp calibrated Gaussian processes, addressing the challenge of obtaining tight predictive quantiles while maintaining calibration.

These advances highlight both the progress made in uncertainty quantification for foundation models and the remaining challenges that motivate this thesis. Our work builds upon these developments while addressing key limitations in scalability and robustness.

## 1.8 Thesis Structure and Contributions Overview

This thesis makes several contributions to the field of uncertainty quantification for machine learning models, with a focus on enhancing pre-trained foundation models for high-stakes applications in healthcare and climate science. The main contributions are:

1. **Two-Stage Gaussian Process Methodology**: We develop a novel two-stage GP framework that addresses the challenges of mean and kernel misspecification, enabling more accurate uncertainty quantification. This methodology includes an automatic kernel search algorithm and a subsampling-based warm-start strategy for efficient hyperparameter optimization.

2. **Kernel Preconditioning for Computational Efficiency**: We develop a kernel preconditioning approach to address the computational challenges of standard GP implementations, enabling the application of our methodology to larger datasets.

3. **Conformal Prediction for Robust Uncertainty Quantification**: We integrate Conformal Prediction with our GP methodology to address the limitations of GPs in non-stationary settings, providing distribution-free guarantees for prediction intervals.

4. **Application to Health Foundation Models**: We apply our two-stage GP framework to enhance pre-trained foundation models for patient risk prediction, demonstrating improved uncertainty quantification in healthcare applications.

5. **Application to Climate and Weather Foundation Models**: We extend our methodology to time-dependent, non-stationary data in climate and weather prediction, specifically developing a Locally Calibrated Adaptive Conformal Prediction (LC-ACP) approach for hurricane track forecasting with reliable uncertainty estimates.

The remainder of this thesis is organized as follows:

- **Chapter 1** (this chapter) presents the background and motivation for uncertainty quantification, including its importance in high-stakes applications, the foundations of Gaussian Processes, sparse approximations, and conformal prediction.

- **Chapter 2** introduces our two-stage Gaussian Process methodology, including theoretical analysis, algorithm development, and empirical validation.

- **Chapter 3** applies our methodology to healthcare foundation models, focusing on patient risk prediction with reliable uncertainty estimates.

- **Chapter 4** addresses the computational challenges of GPs through kernel preconditioning and unbiased estimators, enabling scalable uncertainty quantification for large datasets.

- **Chapter 5** explores the limitations of Gaussian Processes and introduces Adaptive Conformal Prediction as a complementary approach, with particular focus on the development of Locally Calibrated Adaptive Conformal Prediction (LC-ACP) for non-stationary data.

- **Chapter 6** applies the LC-ACP methodology to climate and weather foundation models, specifically addressing the challenges of hurricane track prediction with the Tiny Time Mixer foundation model.

- **Chapter 7** concludes the thesis with a summary of contributions, practical recommendations, and directions for future research.

Through these contributions, this thesis aims to advance the state of the art in uncertainty quantification for machine learning models, enabling more reliable and trustworthy predictions in high-stakes applications. By exploring both Gaussian Process-based and Conformal Prediction-based approaches, we provide a comprehensive framework for uncertainty quantification that can be tailored to different application domains and data characteristics.

# Chapter 2

# Two-Stage Gaussian Process Methodology

"The art of being wise is the art of knowing what to overlook."

— William James

## 2.1 Introduction to Two-Stage Gaussian Processes

Gaussian Process Regression (GPR) offers a principled Bayesian approach to uncertainty quantification in machine learning, as discussed in Chapter 1. However, its performance critically depends on the appropriate specification of both the mean function and the kernel function. Misspecification of either component can lead to poor predictions and unreliable uncertainty estimates—a critical concern for high-stakes applications in, e.g., healthcare and climate science where accurate uncertainty bounds are as important as accurate predictions.

This chapter introduces our novel two-stage Gaussian Process methodology, which addresses both mean and kernel misspecification to provide more reliable uncertainty estimates. Building upon the theoretical foundations of Gaussian Processes presented in the previous chapter, we develop a systematic approach to overcome the limitations of standard GP implementations. The methodology developed here forms the foundation for applications to healthcare foundation models in Chapter 3 and informs our computational efficiency improvements in Chapter 4.

Consider a dataset $X_n = \{\mathbf{x}_i, y_i\}_{i=1}^n$, generated by the relationship $y_i = f(\mathbf{x}_i) + \epsilon_i$,

with $\epsilon_i$ drawn from $\mathcal{N}(0, \sigma_\epsilon^2)$ and $\mathbf{x}_i$ belonging to a subset $\Omega$ of $\mathbb{R}^d$. Gaussian Process Regression models this function $f$ using a Gaussian Process with a prior $\mathcal{GP}(m(\cdot), k_\theta(\cdot, \cdot))$. Here, $m(\cdot)$ represents the mean function, and $k_\theta(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 c(\mathbf{x}_i, \mathbf{x}_j) + \sigma_\xi^2 \delta_{ij}$ defines the covariance function where $\theta$ denotes all the hyperparameters in the covariance function. Upon conditioning on the observed data, the posterior predictive distribution for a new input $\mathbf{x}_*$ is calculated as:

$$\hat{m}(\mathbf{x}_*) = m(\mathbf{x}_*) + K_{*X}(\theta) K_{XX}(\theta)^{-1}(\mathbf{y}_n - m(X)), \tag{2.1}$$

$$\hat{k}(\mathbf{x}_*, \mathbf{x}_*) = k_\theta(\mathbf{x}_*, \mathbf{x}_*) - K_{*X}(\theta) K_{XX}(\theta)^{-1} K_{X*}(\theta), \tag{2.2}$$

where $\mathbf{y}_n = [y_1, y_2, \ldots, y_n]^\top$, $m(X)$ is a vector obtained via evaluating $m(\cdot)$ over $X_n$, $K_{*X}(\theta)$ and $K_{XX}(\theta)$ denote the kernel matrices obtained by evaluating the covariance function $k_\theta$ over $(\mathbf{x}_*, X_n)$ and $(X_n, X_n)$, respectively. When the context is clear, we will also use $m_n$ and $K_n$ to denote $m(X)$ and $K_{XX}(\theta)$ for simplicity.

Specifying the appropriate mean function $m(\cdot)$ and covariance function $k_\theta(\cdot, \cdot)$, with $\theta$ encompassing hyperparameters such as the lengthscale $l$, output scale $\sigma_f^2$, and likelihood noise $\sigma_\xi^2$, is pivotal. Typically, the hyperparameters are optimized by minimizing the Negative Log-Likelihood ($NLL$) through methods like cross-validation or gradient descent, given by:

$$L(\theta; X_n) = \frac{1}{2n}((\mathbf{y}_n - m(X))^\top K_n^{-1}(\mathbf{y}_n - m(X)) + \log \det K_n + n \log 2\pi). \tag{2.3}$$

Notice in Equation (2.3), we scale the likelihood with $\frac{1}{n}$, a common practice when training GP for stability. The Gaussian process model trained via the NLL utilizing the full training dataset will be denoted as Exact-GP throughout this chapter.

Our key contribution is a systematic approach to overcome misspecifications called Two-Stage GPR, as shown in Figure 2.1. We first detail the framework to mitigate mean misspecification, then present an automated kernel search algorithm to address kernel misspecification, and finally introduce a subsampling warm start strategy for efficient training to

Figure 2.1: **Two-stage Gaussian Process Regression (GPR) Framework.** Stage 1: Automatic Kernel Search selects the best kernel for the mean prediction, followed by mean prediction using a Kernel Ridge Regression (KRR). Stage 2: After demeaning the training data using the mean prediction from the first stage, automatic Kernel Search identifies the best kernel for uncertainty quantification, and a zero-mean GPR with the corresponding kernel is trained via subsampling warm start. The final predictive distribution combines these mean and covariance predictions to enhance the model's accuracy and robustness.

avoid hyperparameter misspecifications.

## 2.2 Mitigating Mean Misspecification via Two-stage GPR

In this section, we analyze the impact of mean misspecification on the performance of Gaussian Process Regression and propose our two-stage approach to mitigate these issues.

### 2.2.1 Impact of Mean Misspecification

The standard GPR model assumes a prior distribution over functions $f(x) \sim \mathcal{GP}(m(x), k_\theta(x, x'))$. When the mean function is misspecified (e.g., assuming a zero mean when the true function has a non-zero trend), the model must compensate through the kernel function. This leads to several critical issues:

1. **Biased Predictions**: The posterior mean predictions are systematically biased, particularly in regions with sparse data. This bias manifests as:

$$\mathbb{E}[\hat{f}(x_*) - f(x_*)] = (m(x_*) - m_{true}(x_*)) + \mathcal{O}(\|K_{*X}K_{XX}^{-1}\|) \tag{2.4}$$

2. **Overconfident Uncertainty Estimates**: The model underestimates uncertainty by attributing systematic deviations to the kernel rather than mean misspecification. The posterior variance is affected as:

$$\mathrm{Var}[\hat{f}(x_*)] = k(x_*, x_*) - K_{*X}K_{XX}^{-1}K_{X*} + \mathcal{O}(\|m - m_{true}\|^2) \tag{2.5}$$

3. **Poor Extrapolation**: The model reverts to the incorrect prior mean when extrapolating beyond the training data, leading to:

$$\lim_{x_* \to \infty} \hat{f}(x_*) = m(x_*) \neq m_{true}(x_*) \tag{2.6}$$

If the mean $m(x)$ is incorrectly assumed to be zero, the hyperparameters are derived by minimizing the misspecified expectation of NLL (MEL):

$$\mathrm{MEL} = \frac{1}{2n}\Big( \mathrm{Tr}(K_n(\theta)^{-1}[K_n(\theta_*) - m(X)m(X)^\top]) + \log \det K_n(\theta)\Big) \tag{2.7}$$

The MEL underestimates the data fitting loss due to the subtraction of a positive term $m(X)^\top K_n(\theta)^{-1}m(X)$. This causes GPR trained via MEL to introduce bias, as it penalizes less on the data fitting loss while penalizing more on the model complexity than it should, often yielding an underfitted model.

We can prove through formal analysis that minimizing MEL will not recover the ground-truth hyperparameters $\theta^*$ if the mean function is non-zero:

**Theorem 2.2.1** (Mean Misspecification Effect)**.** *If the mean function $m(x)$ is not a zero*

*function, minimizing the MEL will not recover the ground-truth hyperparameters $\theta^*$ if $\theta^*$ is not a stationary point of $m(X)^\top K_n^{-1} \frac{\partial K_n}{\partial \theta} K_n^{-1} m(X)$.*

This theorem indicates that not only is the mean prediction inaccurate, but the uncertainty quantification is also compromised due to the misspecified kernel hyperparameters obtained by minimizing MEL.

### 2.2.2 Proposed Two-Stage GPR Framework

To address these issues, we propose a two-stage GPR framework that separates the tasks of mean prediction and uncertainty quantification:

1. **Stage 1 - Mean Prediction**: Use a flexible, data-driven method to estimate the mean function $\hat{m}(x)$ without making strong assumptions about its form. We employ Kernel Ridge Regression (KRR) for this purpose, which effectively captures complex trends in the data.

2. **Stage 2 - Uncertainty Quantification**: With the estimated mean function from Stage 1, train a zero-mean GP on the residuals $r_i = y_i - \hat{m}(x_i)$ to model the uncertainty. This allows the GP to focus solely on capturing the covariance structure of the data without being biased by an incorrectly specified mean.

The final predictive distribution for a new input $x_*$ is then:

$$\hat{f}(x_*) \sim \mathcal{N}(\hat{m}(x_*) + \mu_r(x_*), \sigma_r^2(x_*)) \tag{2.8}$$

$$\sim \mathcal{N}(\hat{m}(x_*), \sigma_r^2(x_*)) \tag{2.9}$$

where $\mu_r(x_*)$ and $\sigma_r^2(x_*)$ are the predictive mean and variance from the GP trained on residuals. Since we use a zero-mean GP for the residuals, $\mu_r(x_*)$ is typically small, and the predictive mean is dominated by $\hat{m}(x_*)$.

This approach offers several advantages:

1. **Flexibility**: The mean function can be estimated using any regression method, not just limited to parametric forms.

2. **Robust Uncertainty Quantification**: The GP trained on residuals provides more reliable uncertainty estimates since it doesn't need to compensate for mean misspecification.

3. **Improved Extrapolation**: The explicit modeling of the mean trend allows for better extrapolation beyond the training data.

4. **Computational Efficiency**: The KRR for mean prediction can be trained more efficiently than a full GP, particularly for large datasets.

## 2.3 Addressing Kernel Misspecification

Kernel misspecification occurs when the chosen kernel function cannot adequately capture the true covariance structure of the data. In this section, we analyze the impact of kernel misspecification and propose an automatic kernel search algorithm to address this issue.

### 2.3.1 Impact of Kernel Misspecification

The choice of kernel function in a GP model determines the types of functions that can be effectively modeled. Inappropriate kernel choices lead to several issues:

1. **Incorrect Smoothness Assumptions**: Using an RBF kernel for non-smooth functions or a Matérn kernel with incorrect smoothness parameter.

2. **Missing Periodic Components**: Failing to capture periodic patterns in the data, leading to systematic errors.

3. **Inappropriate Length Scales**: Using a single length scale when the function varies at different rates in different regions.

To quantify the impact of kernel misspecification, we analyze the reducible and irreducible error components. For a new unobserved data point $(x_*, y_*)$ with $y_* = f(x_*) + \epsilon_*$, the prediction error can be decomposed as:

$$|\hat{m}(x_*) - y_*| = |K_{*X}(K_{XX} + \sigma_\xi^2 \mathbf{I})^{-1} y_n - y_*| \tag{2.10}$$

$$= \underbrace{|K_{*X}(K_{XX} + \sigma_\xi^2 \mathbf{I})^{-1} f_n - f(x_*)|}_{\text{Reducible error}} + \underbrace{|K_{*X}(K_{XX} + \sigma_\xi^2 \mathbf{I})^{-1} \epsilon_X - \epsilon_*|}_{\text{Irreducible error}} \tag{2.11}$$

When the kernel is misspecified, the reducible error component remains significant even with optimal hyperparameters, indicating the model's inability to capture the true function.

## 2.3.2 Automatic Kernel Search Algorithm

To address kernel misspecification, we propose an automatic kernel search algorithm that systematically evaluates different kernel functions and selects the most appropriate one based on predictive performance. This approach is guided by the following theoretical result:

**Theorem 2.3.1** (Kernel Misspecification Bound)**.** *Under mild assumptions (C1-C4 in Appendix), suppose $f \in \mathcal{H}_k$. Then with probability $1 - \delta$ we have the following bound*

$$\frac{|\hat{m}_n(\mathbf{x}_*) - y_*|}{|K_{*X}(K_{XX} + \sigma_\xi^2 \mathbf{I})^{-1} \epsilon_{\mathbf{X}} - \epsilon_*|} \leq 1.1, \quad \forall x \in \Omega \tag{2.12}$$

*when $n \geq \dfrac{\sigma_\xi^2}{\left(\dfrac{0.1 A(\delta) \sqrt{K_{*X} K_{X*}}}{(\lambda_1 + \sigma_\xi^2) C' \|f\|_{\mathcal{H}_k}}\right)^{\frac{2m_0}{2m_0 - d}}}$ where $A(\delta) = \sqrt{1 - 2\log(1 - \delta)} - 1$, $\lambda_1$ is the largest eigenvalue of the kernel matrix $K_{XX}$ and $C'$ is the universal constant in Lemma regarding contraction rates, $\hat{m}_n(\mathbf{x}_*)$ is defined with a zero-mean prior.*

This theorem indicates that when the kernel is appropriately specified (i.e., the true function belongs to the RKHS of the chosen kernel), the prediction error is primarily dominated by the irreducible noise error, with the ratio of prediction error to noise error approaching 1.1 for a sufficiently large dataset.

Based on this insight, we propose Algorithm 1 for automatic kernel selection:

---
**Algorithm 1** Automatic Kernel Search

---
**Require::** Training data $\{X_n, y_n\}$, validation data $\{X_v, y_v\}$, candidate kernels $\{k_1, k_2, \ldots, k_K\}$
**Ensure::** Selected kernel $k_{best}$

1: **FOR** $i = 1$ to $K$ **DO**
2:      Train GP model with kernel $k_i$ on $\{X_n, y_n\}$ to obtain hyperparameters $\theta_i$
3:      Compute prediction error $E_i$ on validation data $\{X_v, y_v\}$
4: **END FOR**

5: Select $k_{best} = \arg\min_i E_i$ **RETURN** $k_{best}$

---

This algorithm systematically evaluates each candidate kernel by training a GP model on the training data and measuring its predictive performance on a separate validation set. The kernel with the lowest validation error is selected as the most appropriate for the given dataset.

In practice, we use a diverse set of candidate kernels, including:

- Radial Basis Function (RBF) kernel for smooth functions

- Matérn kernels with different degrees of smoothness

- Periodic kernels for functions with cyclical patterns

- Composite kernels (sum and product) for more complex functions

This approach allows for data-driven kernel selection without requiring strong prior assumptions about the underlying function.

## 2.4 Efficient Training via Subsampling

Training a GP model on large datasets is computationally challenging due to the $O(n^3)$ time complexity of exact inference. In this section, we propose a scalable variant of our Two-Stage GP methodology using Gaussian Process Nearest Neighbor (GPNN). This approach combines

the theoretical advantages of our Two-Stage framework with the computational efficiency of nearest-neighbor-based approximations.

## 2.4.1 Gaussian Process Nearest Neighbor (GPNN)

GPNN is a scalable GP approximation that leverages the screening effect, which states that predictions at a given point are primarily influenced by nearby observations. By focusing on a subset of nearest neighbors for each prediction, GPNN achieves significant computational savings while maintaining competitive accuracy.

---

**Algorithm 2** Gaussian Process Nearest Neighbor (GPNN)

---

**Require:**: Training data $\{X_n, y_n\}$, test point $x_*$, number of neighbors $k$, kernel function $\mathcal{K}$
**Ensure:**: Predictive distribution $p(f(x_*)|X_n, y_n)$

1: // **Training Phase**
2: Randomly sample a subset $\{X_s, y_s\}$ from $\{X_n, y_n\}$
3: Optimize kernel hyperparameters $\theta$ on subset by minimizing NLL
4: Compute calibration factor $c = \frac{1}{|X_{\text{cal}}|} \sum_{i \in X_{\text{cal}}} \frac{(y_i - \hat{m}(x_i))^2}{\hat{\sigma}_i^2}$
5: Adjust variance scaling: $\theta_{\text{scaled}} = \{l, \sigma_f \sqrt{c}, \sigma_\xi \sqrt{c}\}$

6: // **Prediction Phase**
7: Find $k$ nearest neighbors $\{X_k, y_k\}$ of $x_*$ from training data
8: Compute kernel matrix $K_{kk}$ among neighbors using $\mathcal{K}$ with $\theta_{\text{scaled}}$
9: Compute cross-kernel vector $K_{*k}$ between $x_*$ and neighbors
10: Compute predictive mean: $\mu(x_*) = K_{*k}(K_{kk} + \sigma_\xi^2 I)^{-1} y_k$
11: Compute predictive variance: $\sigma^2(x_*) = \mathcal{K}(x_*, x_*) - K_{*k}(K_{kk} + \sigma_\xi^2 I)^{-1} K_{k*}$

12: **RETURN** $\mathcal{N}(\mu(x_*), \sigma^2(x_*))$

---

The GPNN algorithm significantly reduces computational complexity from $O(n^3)$ to $O(k^3)$ for each prediction, where $k \ll n$ is the number of nearest neighbors. This allows for efficient inference on datasets with millions of examples.

## 2.4.2 Two-Stage GPNN (2StGPNN)

Building on our Two-Stage GP framework, we integrate GPNN as the baseline algorithm to create a scalable version suitable for large datasets:

---

**Algorithm 3** Two-Stage Gaussian Process Nearest Neighbor (2StGPNN)

---

**Require:**: Training data $\{X_n, y_n\}$, test point $x_*$, number of neighbors $k$, candidate kernels $\{k_1, k_2, \ldots, k_K\}$
**Ensure:**: Predictive distribution $p(f(x_*)|X_n, y_n)$

1: // **Stage 1: Mean Prediction with KRR**
2: Select best kernel $k_m$ for mean prediction using Algorithm 1
3: Find $k$ nearest neighbors $\{X_k, y_k\}$ of $x_*$ from training data
4: Compute kernel matrix $K_{kk}$ among neighbors using $k_m$
5: Compute KRR prediction: $\hat{m}(x_*) = k_m(x_*, X_k)(K_{kk} + \lambda I)^{-1} y_k$
6: Compute residuals $r_i = y_i - \hat{m}(x_i)$ for all training points

7: // **Stage 2: Uncertainty Quantification with GPNN**
8: Select best kernel $k_r$ for residuals using Algorithm 1
9: Find $k$ nearest neighbors $\{X_k, r_k\}$ of $x_*$ from residual data
10: Optimize kernel hyperparameters $\theta_r$ on a random subset of residuals
11: Compute kernel matrix $K_{kk}$ among neighbors using $k_r$ with $\theta_r$
12: Compute cross-kernel vector $K_{*k}$ between $x_*$ and neighbors
13: Compute residual predictive distribution: $p(r(x_*)|X_k, r_k) = \mathcal{N}(\mu_r(x_*), \sigma_r^2(x_*))$
14: Return combined distribution: $p(f(x_*)|X_n, y_n) = \mathcal{N}(\hat{m}(x_*) + \mu_r(x_*), \sigma_r^2(x_*))$

15: **RETURN** $\mathcal{N}(\hat{m}(x_*) + \mu_r(x_*), \sigma_r^2(x_*))$

---

The Two-Stage GPNN approach inherits the computational efficiency of GPNN while addressing its limitations in uncertainty quantification through our two-stage methodology. By separating mean prediction from uncertainty quantification, it achieves more reliable uncertainty estimates without sacrificing computational efficiency.

Key advantages of Two-Stage GPNN include:

- **Scalability**: Maintains the $O(k^3)$ computational complexity of GPNN for datasets with millions of examples

- **Improved UQ**: Addresses the limitations of GPNN in uncertainty quantification through the two-stage approach

- **Flexibility**: Can be combined with other scalable GP approximations like SVGP or Vecchia approximations

## 2.5   Numerical Experiments

In this section, we present comprehensive experiments to validate our Two-Stage GP methodology. First, we assess the performance on various UCI regression datasets, comparing Exact-GP and Two-Stage Exact-GP based on several popular metrics. Next, we introduce a novel uncertainty-aware metric to measure the performance of uncertainty quantification for regression tasks. Finally, we demonstrate the UQ capabilities of Two-Stage GP on safety-critical applications in healthcare.

### 2.5.1   Performance Comparison on UCI Dataset

We first evaluate our Two-Stage GP approach on several standard UCI regression benchmark datasets. For these datasets, we report three key metrics:

1. **RMSE (Root Mean Square Error)**: Defined as $\text{RMSE} := \sqrt{\frac{1}{n} \sum_{i=1}^{n} (\hat{y}_i - y_i)^2}$, which measures the accuracy of the mean predictions.

2. **NLL (Negative Log-Likelihood)**: Defined in Equation (2.3), which evaluates the probabilistic fit of the model incorporating both mean and variance.

3. **QICE (Quantile Interval Coverage Error)**: Calculated as $\text{QICE} := \frac{1}{M} \sum_{m=1}^{M} \left| r_m - \frac{1}{M} \right|$, where $r_m := \frac{1}{N} \sum_{n=1}^{N} \mathbf{1}_{y_n \geq \hat{y}_n^{m-1}} \cdot \mathbf{1}_{y_n \leq \hat{y}_n^m}$ represents the proportion of true targets falling between the $m-1$-th and $m$-th quantiles of predicted targets. This metric evaluates prediction calibration and serves as an important measure of the model's UQ capability.

For all three metrics, a smaller value indicates better model performance. Table 2.1 reports the performance of standard Exact-GP with a default RBF kernel compared to our Two-Stage Exact-GP.

The results clearly demonstrate that our Two-Stage GP approach outperforms standard Exact-GP on most datasets. Most notably, for the Power and Naval datasets, standard Exact-GP reports extremely large NLL values compared to Two-Stage Exact-GP, indicating

Table 2.1: RMSE, NLL, and QICE metrics for various UCI datasets using Exact-GP with RBF kernel and Two-Stage Exact-GP.

| Dataset | Exact-GP(RBF) | | | Two-stage Exact-GP | | |
|---|---|---|---|---|---|---|
| | RMSE | NLL | QICE | RMSE | NLL | QICE |
| Yacht | $1.29 \pm 0.42$ | $-1.15 \pm 0.03$ | $6.37 \pm 1.47$ | $\mathbf{0.41 \pm 0.16}$ | $\mathbf{-1.52 \pm 0.04}$ | $\mathbf{5.03 \pm 1.59}$ |
| Boston | $2.77 \pm 0.68$ | $\mathbf{-0.22 \pm 0.05}$ | $3.78 \pm 1.00$ | $\mathbf{2.70 \pm 0.67}$ | $1.11 \pm 1.49$ | $\mathbf{3.55 \pm 0.93}$ |
| Energy | $0.89 \pm 0.12$ | $-0.93 \pm 0.03$ | $3.36 \pm 0.94$ | $\mathbf{0.37 \pm 0.07}$ | $\mathbf{-1.51 \pm 0.03}$ | $\mathbf{2.24 \pm 0.56}$ |
| Concrete | $5.35 \pm 0.64$ | $\mathbf{-0.39 \pm 0.07}$ | $2.88 \pm 0.67$ | $\mathbf{3.78 \pm 0.58}$ | $0.03 \pm 0.50$ | $\mathbf{2.23 \pm 0.57}$ |
| Wine | $0.62 \pm 0.04$ | $0.95 \pm 0.05$ | $\mathbf{13.21 \pm 0.31}$ | $0.60 \pm 0.04$ | $0.67 \pm 0.78$ | $13.2 \pm 0.31$ |
| Kin8nm | $\mathbf{0.07 \pm 0.00}$ | $\mathbf{-1.03 \pm 0.21}$ | $0.95 \pm 0.24$ | $\mathbf{0.07 \pm 0.00}$ | $-0.15 \pm 0.03$ | $\mathbf{0.94 \pm 0.27}$ |
| Power | $3.75 \pm 0.19$ | $3111 \pm 16814$ | $\mathbf{1.05 \pm 0.27}$ | $\mathbf{3.23 \pm 0.20}$ | $\mathbf{0.06 \pm 0.26}$ | $15.63 \pm 0.09$ |
| Naval | $\mathbf{0.00 \pm 0.00}$ | $924.1 \pm 3892$ | $0.97 \pm 0.43$ | $\mathbf{0.00 \pm 0.00}$ | $\mathbf{-1.62 \pm 0.00}$ | $0.69 \pm 0.22$ |

unreliable uncertainty quantification. This confirms that our Two-Stage approach provides more reliable uncertainty estimates by addressing mean misspecification.

## 2.5.2 Uncertainty-Aware Evaluation Metrics

To further evaluate the quality of uncertainty quantification provided by our Two-Stage GP, we introduce the Uncertainty-Aware RMSE (UA-RMSE) metric. In GP models, uncertainty is typically represented as a standard error estimate for predictions. For each prediction $\hat{y}_i = \hat{m}(\mathbf{x}_i)$, the associated standard error is $\hat{s}_i = \sqrt{\hat{k}(\mathbf{x}_i, \mathbf{x}_i)}$ from Equation (2.1).

We define two variations of UA-RMSE:

1. **High Certainty RMSE (HC-RMSE)**: Calculated as the RMSE for predictions where the standard error is lower than the $100q$-th quantile of all standard error estimates on test data, denoted as $\hat{s}_q$:

$$\text{HC-RMSE} := \sqrt{\frac{1}{qn} \sum_{h=1}^{qn} (\hat{y}_h - y_h)^2}, \quad \text{where } \hat{s}_h \leq \hat{s}_q \tag{2.13}$$

2. **Low Certainty RMSE (LC-RMSE)**: Calculated for predictions where the standard

error exceeds the $(1-q)100$-th quantile, denoted as $\hat{s}_{1-q}$:

$$\text{LC-RMSE} := \sqrt{\frac{1}{qn}\sum_{l=1}^{qn}(\hat{y}_l - y_l)^2}, \quad \text{where } \hat{s}_l > \hat{s}_{1-q} \tag{2.14}$$

A model with high-quality uncertainty quantification is expected to exhibit a lower HC-RMSE compared to LC-RMSE. This indicates that the model is more confident in its accurate predictions and less confident in its inaccurate ones.



Figure 2.2: Uncertainty quantification results comparing HC-RMSE and LC-RMSE for both standard GPNN and Two-Stage GPNN approaches on UCI datasets with $q = 0.1$. Lower HC-RMSE compared to LC-RMSE indicates better uncertainty quantification.

Figure 2.2 shows the HC-RMSE and LC-RMSE for both standard GPNN and Two-Stage GPNN on seven UCI datasets with $q = 0.1$. Two-Stage GPNN achieves lower RMSE for high certainty instances across all datasets. For the discrepancy between HC-RMSE and LC-RMSE, Two-Stage GPNN also outperforms standard GPNN on five datasets. This confirms that our approach better differentiates between confident and uncertain predictions.

### 2.5.3   Synthetic Data Experiment: Mean Misspecification Effects

To illustrate the impact of mean misspecification and the effectiveness of our Two-Stage approach in addressing it, we conducted controlled experiments with synthetic data.

Figure 2.3 shows an example with a synthetic function with a strong non-zero mean trend. The standard GP with a zero-mean prior fails to capture this trend, resulting in

(a) True Function  (b) Standard GP (zero-mean)  (c) Two-stage GP

Figure 2.3: Comparison of standard GP with misspecified mean and Two-Stage GP. The true function (a) has a non-zero mean trend. The standard GP with zero-mean prior (b) fails to capture this trend, leading to biased predictions and incorrect uncertainty estimates. The Two-Stage GP (c) accurately captures the trend and provides well-calibrated uncertainty estimates.

biased predictions and incorrect uncertainty estimates—note how the confidence intervals don't cover the true function in regions with sparse data. In contrast, our Two-Stage GP accurately models the mean trend and provides well-calibrated uncertainty estimates, with confidence intervals that properly encompass the true function.

We also conducted a supplementary experiment with a more complex function, shown in Figure 2.4.

In this experiment, the true function is $f(x) = 3|x|^{\frac{3}{2}} + 2\sin(2\pi x)$ with 30 data points. The standard Exact-GP severely underfits the data, with its 95% confidence interval covering only 66.7% of the data points, indicating poor calibration. In contrast, our Two-Stage GP covers 96.67% of the data points, demonstrating much better calibration of uncertainty estimates.

## 2.5.4 Hyperparameter Landscape Analysis

To understand why our Two-Stage approach yields better hyperparameter estimates, we analyzed the loss landscapes for different hyperparameter combinations.

Figure 2.5 shows the loss contours for different hyperparameter combinations, with color

(a) True Function     (b) Standard GP (zero-mean)     (c) Two-stage GP

Figure 2.4: Comparison using a more complex function $f(x) = 3|x|^{\frac{3}{2}} + 2\sin(2\pi x)$ with 30 data points. The Exact-GP underfits the data, with its 95% confidence interval covering only 66.7% of the data, while the Two-Stage GP covers 96.67% of the data.



(a) Lengthscale vs Noise     (b) Lengthscale vs Kernel Scale

Figure 2.5: Loss contours for different hyperparameter combinations. The color represents the NLL value (lower is better). The standard GP optimizes for suboptimal hyperparameters due to mean misspecification, while the Two-Stage GP finds better hyperparameter combinations by separating mean prediction from uncertainty quantification.

representing the NLL value (lower is better). When the mean function is misspecified, the standard GP can get trapped in suboptimal regions of the hyperparameter space. By separating mean prediction from uncertainty quantification, our Two-Stage approach is better able to navigate this hyperparameter landscape and find more optimal configurations.

### 2.5.5 Application to Healthcare Risk Prediction

To demonstrate the practical relevance of our approach in high-stakes applications, we applied our Two-Stage GP methodology to healthcare risk prediction tasks using foundation models. The detailed results of this application are presented in Chapter 3, where we show that our approach provides valuable uncertainty quantification for safety-critical healthcare applications.

## 2.6 Conclusion and Connections to Subsequent Chapters

In this chapter, we introduced a novel Two-Stage Gaussian Process methodology that addresses the critical issues of mean and kernel misspecification in GP regression. By separating the tasks of mean prediction and uncertainty quantification, our approach achieves more accurate predictions and more reliable uncertainty estimates compared to standard GP implementations. Additionally, our subsampling-based warm start strategy provides computational efficiency for large datasets without sacrificing predictive performance.

The methodology presented in this chapter forms the foundation for our subsequent work on uncertainty quantification for foundation models in high-stakes applications. In Chapter 3, we apply this methodology to healthcare foundation models for reliable patient risk prediction. The computational challenges associated with large-scale implementation of GP methods are addressed through kernel preconditioning techniques in Chapter 4. Finally, the limitations of the Gaussian Process approach in non-stationary settings motivate our exploration of Conformal Prediction in Chapter 5.

While our Two-Stage GP methodology offers significant advantages over standard GP implementations, it still faces challenges in scenarios with non-stationary data or distributional shift. These challenges will be addressed in subsequent chapters through complementary approaches like Locally Calibrated Adaptive Conformal Prediction (LC-ACP), which provides distribution-free guarantees for prediction intervals.

# Chapter 3

# Application to Health Foundation Models

> "The good physician treats the disease; the great physician treats the patient who has the disease."
>
> — William Osler

## 3.1 Introduction

Building upon the two-stage Gaussian Process methodology developed in Chapter 2, we now explore its application to health foundation models. Healthcare decisions often involve high stakes with significant consequences for patient outcomes. The ability to accurately predict patient risks for various conditions is a critical component of modern healthcare systems, enabling preventive interventions and optimized resource allocation. However, health data presents unique challenges for predictive modeling, including high dimensionality, heterogeneity, temporal dynamics, missing values, and imbalanced distributions [14, 78]. Additionally, the sensitive nature of healthcare decisions necessitates models that not only provide accurate predictions but also quantify uncertainty reliably [7, 50].

Pre-trained foundation models have emerged as powerful tools for analyzing complex medical data, including electronic health records (EHRs), medical imaging, and clinical notes [63, 100, 19]. These models leverage transfer learning to extract meaningful representations from vast amounts of data, and have demonstrated impressive performance on

various healthcare tasks. However, a significant limitation of standard foundation models is their deterministic nature, producing point predictions without quantifying the associated uncertainty [62].

In safety-critical healthcare applications, understanding prediction uncertainty is crucial for several reasons:

1. **Patient Safety**: When a model expresses high uncertainty, healthcare providers can exercise additional caution and potentially seek further diagnostic information.

2. **Resource Allocation**: Cases with high-confidence predictions can be prioritized differently from those with uncertain predictions, optimizing the allocation of limited healthcare resources.

3. **Clinical Decision Support**: Uncertainty measures provide valuable context for clinicians, helping them balance model suggestions with their clinical expertise.

4. **Patient Diversity**: Medical data often contains underrepresented populations for whom predictions may be less reliable; uncertainty quantification helps identify these cases.

This chapter demonstrates how our two-stage GP approach enables accurate patient risk prediction while providing robust uncertainty quantification, addressing a critical need in healthcare analytics. The computational challenges encountered in applying GPs to large-scale health datasets will motivate our development of efficient preconditioning techniques in Chapter 4, while the limitations in handling non-stationary health data will lead to our exploration of conformal prediction in Chapter 5. The methodology developed here provides a foundation for our subsequent applications to climate foundation models in Chapter 6.

## 3.2 Challenges in Health Data Prediction

### 3.2.1 Data Characteristics and Challenges

Health data prediction presents several unique challenges compared to other domains:

- **Data Heterogeneity**: Medical data is inherently heterogeneous, comprising structured data (lab values, vital signs), unstructured text (clinical notes, discharge summaries), images (X-rays, MRIs), and time series (ECG, EEG). This multimodal nature creates challenges for developing unified predictive models [79].

- **Temporality and Irregularity**: Patient data is collected at irregular intervals and varies in frequency based on acuity and healthcare setting. Models must account for these irregular sampling patterns and the temporal dependencies between observations [15].

- **Class Imbalance**: Many health conditions have low prevalence, resulting in severely imbalanced datasets where negative instances far outnumber positive ones. This imbalance can bias models toward majority classes, potentially overlooking critical minority cases [47].

- **Missing Data**: Clinical datasets frequently contain missing values due to various factors including practical constraints in data collection, patient dropout, and variable recording practices across healthcare systems [101].

- **Limited Labeled Data**: While healthcare systems generate vast amounts of data, labeled data for specific conditions is often limited, particularly for rare diseases or novel clinical scenarios [53].

- **High-Stakes Decisions**: Perhaps most importantly, health predictions directly impact patient care decisions. Incorrect predictions can lead to missed treatment opportunities or unnecessary interventions, both with potentially serious consequences [14].

These challenges underscore the need for predictive models that not only provide accurate point predictions but also quantify uncertainty reliably. The Gaussian Process framework presented in this thesis is particularly well-suited for this context due to its inherent uncertainty quantification capabilities and adaptability to various data types.

### 3.2.2   Existing Foundation Models in Healthcare

Several foundation models have been developed specifically for healthcare applications, each with their own strengths and limitations:

- **Text-based Models**: Models such as ClinicalBERT [100] and BioGPT [63] are pre-trained on extensive medical literature and clinical notes, allowing them to understand complex medical terminology and relationships.

- **Vision Models**: Vision Transformers (ViT) [19] and similar architectures have been adapted for medical imaging tasks, including radiology, pathology, and dermatology.

- **Multimodal Models**: Emerging models that integrate data across modalities (text, imaging, structured data) show promise for comprehensive patient assessment.

These models are typically pre-trained on large datasets in a self-supervised manner and then fine-tuned for specific downstream tasks. The standard approach for adapting these models involves freezing most of the pre-trained parameters while training a task-specific final layer—typically a fully connected layer that produces deterministic outputs.

### 3.2.3   The Critical Role of Uncertainty Quantification

Uncertainty quantification (UQ) in healthcare applications is not merely a statistical nicety but a critical component that can significantly impact patient outcomes and clinical decision-making. Traditional deterministic models may provide predictions that appear confidently

precise, but without an understanding of the associated uncertainty, clinicians cannot properly evaluate the reliability of these predictions.

In healthcare prediction tasks, we encounter multiple forms of uncertainty:

- **Aleatoric Uncertainty**: This represents inherent stochasticity in the data generation process, such as natural variations in patient responses to treatments or measurement noise in clinical instruments. This type of uncertainty cannot be reduced with additional data.

- **Epistemic Uncertainty**: This reflects model uncertainty due to limited training data or model misspecification. Unlike aleatoric uncertainty, epistemic uncertainty can potentially be reduced with additional data or improved model design.

- **Distributional Uncertainty**: This arises when new patients differ significantly from the training population, particularly relevant in healthcare where model generalization across diverse patient demographics is crucial.

The value of UQ in healthcare extends beyond academic interest:

- **Triaging and Intervention Planning**: Uncertainty measures can help prioritize cases, directing immediate attention to high-risk patients with low uncertainty predictions while flagging uncertain cases for additional review.

- **Communicating with Patients**: Sharing uncertainty information can support shared decision-making, helping patients understand the confidence level associated with different prognostic or diagnostic assessments.

- **Model Monitoring and Improvement**: Systematic patterns in prediction uncertainty can highlight areas where models may be underperforming for specific patient subgroups, guiding targeted data collection or model refinement.

# 3.3 Mathematical Formulation of GP-Enhanced Foundation Models

In this section, we present the mathematical formulation of our approach for enhancing foundation models with uncertainty quantification capabilities through Gaussian Processes. This builds directly on the two-stage GP methodology developed in Chapter 2, adapting it to the specific challenges of foundation models.

## 3.3.1 Feature Extraction from Foundation Models

Let $\mathcal{M}$ be a pre-trained foundation model (e.g., ClinicalBERT, BioGPT, or ViT). For an input $\mathbf{x}$ (which could be text, images, or other modalities), we denote the output of the penultimate layer of $\mathcal{M}$ as $\phi(\mathbf{x}) \in \mathbb{R}^d$, where $d$ is the dimensionality of the feature space. This feature extraction can be expressed as:

$$\phi(\mathbf{x}) = \mathcal{M}_\phi(\mathbf{x}) \tag{3.1}$$

where $\mathcal{M}_\phi$ represents the feature extraction component of the model $\mathcal{M}$ (i.e., all layers except the final prediction layer).

For a dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ consisting of $N$ input-output pairs, we can compute the feature representations as $\Phi = \{\phi(\mathbf{x}_i)\}_{i=1}^N$. These feature representations capture the semantic information extracted by the foundation model from the raw inputs.

## 3.3.2 Two-Stage GP for Classification Tasks

While the formulation in Chapter 2 focused primarily on regression tasks, many healthcare applications involve classification. Extending our approach to classification requires several adaptations.

In a classification setting with $C$ classes, we need to model the probabilities $p(y = c|\mathbf{x})$ for

each class $c \in \{1, 2, \ldots, C\}$. Following our framework in Chapter 2, we employ a two-stage approach to handle classification tasks effectively.

For classification, we use $C$ independent GPs to transform the instance embedding into class logits. Specifically, for multi-class classification with $C$ classes, we utilize a Dirichlet distribution as the likelihood model:

$$p(y|\boldsymbol{\pi}) = \text{Cat}(\boldsymbol{\pi}), \quad \boldsymbol{\pi} \sim \text{Dir}(\boldsymbol{\alpha}) \tag{3.2}$$

where $\text{Dir}(\boldsymbol{\alpha})$ is a Dirichlet distribution with $\boldsymbol{\alpha} = (\alpha_1, \ldots, \alpha_C)$. This approach allows us to transform the classification problem into a regression task that fits within our two-stage GP framework.

### Dirichlet Distribution as a Likelihood Model for GP Classification

The key innovation in our approach is the use of the Dirichlet distribution as a likelihood model for classification tasks. This transforms discrete classification into a continuous regression problem appropriate for our two-stage GP methodology.

The transformation works as follows: it is well known that a random variable $\mathbf{x} \sim \text{Dir}(\boldsymbol{\alpha})$ can be generated using $C$ independent Gamma distributions $\Gamma(\alpha_i, 1)$. These Gamma distributions can be approximated by LogNormal distributions:

$$\tilde{x}_i \sim \text{LogNormal}(\tilde{y}_i, \tilde{\sigma}_i^2) \tag{3.3}$$

where $\tilde{y}_i = \log \alpha_i - \frac{\tilde{\sigma}_i^2}{2}$ and $\tilde{\sigma}_i = \log \left( \frac{1}{\alpha_i} + 1 \right)$.

This approximation allows us to use a Gaussian likelihood in the log-space:

$$p(\tilde{y}_i|\mathbf{f}) = \mathcal{N}(f_i, \tilde{\sigma}_i^2) \tag{3.4}$$

where $\mathbf{f} = [f_1, \ldots, f_C]$ are $C$ latent Gaussian processes. By using this lognormal approximation of the Gamma distribution, we replace the intractable Dirichlet likelihood with a

tractable Gaussian likelihood, making the problem amenable to our two-stage GP framework.

For classification labels, we create transformed targets using the concentration parameters of a Dirichlet distribution. Specifically, for a data point belonging to class $c$, we set $\alpha_c = \alpha_0$ (e.g., $\alpha_0 = 2$) and $\alpha_{j \neq c} = 1$, then compute the transformed targets $\tilde{y}_i$ using the formula above. These transformed continuous targets now serve as the regression targets for our two-stage GP approach.

**Theoretical Foundations of GP Classification with Dirichlet Likelihood**

The core theoretical insight of our approach is translating a discrete classification problem into a continuous regression problem that can leverage our two-stage GP methodology. This translation is mathematically justified through several key observations:

1. **Relationship between Dirichlet and Gamma Distributions**: The Dirichlet distribution is intrinsically connected to the Gamma distribution. If $X_i \sim \Gamma(\alpha_i, 1)$ are independent Gamma random variables for $i = 1, \ldots, C$, then the vector $\left( \frac{X_1}{\sum_{j=1}^C X_j}, \ldots, \frac{X_C}{\sum_{j=1}^C X_j} \right)$ follows a $\mathrm{Dir}(\alpha_1, \ldots, \alpha_C)$ distribution. This connection allows us to model the Dirichlet parameters instead of directly modeling class probabilities.

2. **LogNormal Approximation to Gamma**: The Gamma distribution can be approximated by a LogNormal distribution with carefully chosen parameters. This approximation is key to making the problem tractable in a GP framework, as it allows us to work with Gaussian likelihoods. Specifically, if $X \sim \Gamma(\alpha, 1)$, then $X$ can be approximated by $Y \sim \mathrm{LogNormal}(m, s^2)$ where:

$$m = \log(\alpha) - \frac{1}{2} \log\left(1 + \frac{1}{\alpha}\right) \tag{3.5}$$

$$s^2 = \log\left(1 + \frac{1}{\alpha}\right) \tag{3.6}$$

3. **Working in Log-Space**: By working in the log-space of the Gamma parameters, we transform the problem into one with Gaussian likelihoods, which is directly amenable to our

GP approach. The transformation maintains the probabilistic interpretation while enabling efficient computation.

4. **Label Encoding Scheme**: For an input $\mathbf{x}$ belonging to class $c$, we use a specific encoding scheme where we set $\alpha_c = \alpha_0 > 1$ (typically $\alpha_0 = 2$) and $\alpha_{j \neq c} = 1$ for all other classes. This creates a soft one-hot encoding in the Dirichlet parameter space, which retains the classification structure while enabling continuous regression.

This theoretical foundation allows us to leverage the well-established properties of Gaussian Processes for regression while addressing classification tasks. The transformation preserves the probabilistic nature of the problem and enables principled uncertainty quantification, which is crucial for healthcare applications.

## Benefits of Two-Stage GP for Classification

The two-stage approach provides specific benefits for classification tasks beyond what standard GP classification offers. In standard GP classification, a single GP is typically used with non-Gaussian likelihoods (such as logistic or probit), requiring approximation methods like Laplace approximation or expectation propagation that can lead to inaccurate uncertainty estimates. Our two-stage approach addresses several key limitations:

1. **Mean-Variance Decoupling**: By separating mean prediction (Stage 1) from uncertainty modeling (Stage 2), we achieve more accurate modeling of the decision boundaries (through KRR) while maintaining proper uncertainty quantification (through the residual GP). This is particularly valuable in classification where decision boundaries are often complex.

2. **Calibration Improvement**: Two-stage GP classification demonstrates better calibration properties compared to standard GP classification. The mean function estimated in Stage 1 captures the central tendency of the class distributions, while the residual GP in Stage 2 models the uncertainty around this mean more accurately, leading to better-calibrated probability estimates.

3. **Handling Class Imbalance**: In healthcare datasets with significant class imbalance,

the two-stage approach is more robust because the first stage can use techniques like class weighting or oversampling to establish a balanced mean function, while the second stage properly models the uncertainty conditional on this balanced mean, rather than being dominated by majority classes.

4. **Computational Tractability**: By working with Gaussian likelihoods throughout (via the LogNormal approximation), our approach avoids the need for more complex approximate inference methods typically required in GP classification, leading to more reliable uncertainty estimates with fewer approximation artifacts.

These theoretical advantages translate to practical improvements in healthcare classification tasks, where accurate uncertainty quantification is essential for clinical decision support and patient risk stratification.

The two-stage process for classification works as follows:

**Stage 1: Mean Function Estimation**

For each class $c$, we train a KRR model on the feature representations to predict the transformed labels. This gives us a robust mean prediction for each class:

$$\hat{m}_c(\phi(\mathbf{x})) = \mathbf{k}(\phi(\mathbf{x}), \boldsymbol{\Phi})^\top (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{z}_c \tag{3.7}$$

where:

- $\mathbf{k}(\phi(\mathbf{x}), \boldsymbol{\Phi})$ is the vector of kernel evaluations between the feature representation of $\mathbf{x}$ and all training feature representations

- $\mathbf{K}$ is the kernel matrix computed on the training feature representations

- $\lambda$ is the regularization parameter

- $\mathbf{z}_c$ is the vector of transformed labels for class $c$

**Stage 2: Uncertainty Modeling**

For each class $c$, we model the residuals using a zero-mean GP:

$$r_c(\phi(\mathbf{x})) = z_c(\mathbf{x}) - \hat{m}_c(\phi(\mathbf{x})) \tag{3.8}$$

The posterior distribution of the GP provides us with a mean and variance for each class:

$$p(f_c|\phi(\mathbf{x}), \mathcal{D}) = \mathcal{N}(\mu_c(\phi(\mathbf{x})), \sigma_c^2(\phi(\mathbf{x}))) \tag{3.9}$$

where:

$$\mu_c(\phi(\mathbf{x})) = \hat{m}_c(\phi(\mathbf{x})) + \mathbf{k}(\phi(\mathbf{x}), \boldsymbol{\Phi})^\top (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{r}_c \tag{3.10}$$

$$\sigma_c^2(\phi(\mathbf{x})) = k(\phi(\mathbf{x}), \phi(\mathbf{x})) - \mathbf{k}(\phi(\mathbf{x}), \boldsymbol{\Phi})^\top (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}(\phi(\mathbf{x}), \boldsymbol{\Phi}) \tag{3.11}$$

## From GP Outputs to Class Probabilities

For classification prediction, we need to transform the continuous GP outputs back into class probabilities. This involves several steps:

1. We sample from the posterior distribution of each latent GP $f_c$ for a new input $\phi(\mathbf{x}_*)$. 2. These samples are in the log-space of the Gamma parameters, so we exponentiate them to obtain samples from the approximate Gamma distributions. 3. We normalize these values to construct samples from the Dirichlet distribution, which represent class probability vectors. 4. Finally, we average over multiple samples to approximate the marginal class probabilities.

Mathematically, the probability for class $j$ is computed by marginalizing over the latent processes:

$$p[j|\mathcal{D}, \mathbf{x}_*] = \int \frac{e^{f_j^*}}{\sum_{i=1}^{C} e^{f_i^*}} p(\mathbf{f}|\mathcal{D}, \mathbf{x}_*) d\mathbf{f} \tag{3.12}$$

where $p(\mathbf{f}|\mathcal{D}, \mathbf{x}_*)$ is the posterior distribution of the latent GPs. This integral is intractable, but can be approximated through Monte Carlo integration by drawing $N$ samples from the posterior:

$$p[j|\mathcal{D}, \mathbf{x}_*] \approx \frac{1}{N} \sum_{l=1}^{N} \frac{e^{f_j^{*,l}}}{\sum_{i=1}^{C} e^{f_i^{*,l}}} \tag{3.13}$$

where $f_j^{*,l}$ is the $l$-th sample from the posterior GP for class $j$.

**Advantages and Limitations of the GP Classification Approach**

This Dirichlet-based GP classification approach offers several important advantages. First, it preserves the probabilistic nature of GPs, allowing us to quantify predictive uncertainty in a principled way. Second, it enables the use of our two-stage GP methodology for classification tasks without major modifications to the underlying framework. Third, it provides naturally calibrated probabilities due to the Bayesian treatment of the problem.

However, there are also limitations to consider. The approximation of Gamma distributions with LogNormal distributions introduces additional approximation error. The computational cost scales linearly with the number of classes, which can be problematic for tasks with many classes. Additionally, the accuracy of the uncertainty estimates depends on the quality of the GP approximation and the number of samples drawn during inference. In practice, we find that using $N = 100$ samples provides a good balance between computational efficiency and accurate uncertainty estimation for most healthcare classification tasks.

Despite these limitations, this approach has proven highly effective for medical classification tasks where uncertainty quantification is critical, enabling reliable identification of cases where the model is uncertain in its predictions. This is crucial for healthcare applications where incorrect predictions could have serious consequences.

### 3.3.3 Kernel Selection and Hyperparameter Optimization

Following the automatic kernel selection approach described in Algorithm 1 from Chapter 2, we select the optimal kernel for both Stage 1 and Stage 2 of our GP framework.

For Stage 1, we select the kernel that minimizes the cross-validation error on the training

set:

$$k_1^* = \arg\min_k \text{CV-Error}(k, \Phi, \mathbf{z}) \tag{3.14}$$

For Stage 2, we select the kernel that best models the residuals:

$$k_2^* = \arg\min_k \text{NLL}(k, \Phi, \mathbf{r}) \tag{3.15}$$

where NLL is the negative log-likelihood as defined in Equation (2.3) in Chapter 2.

The hyperparameters for each kernel are optimized using the subsampling warm start strategy described in Section 2.4. Specifically, we:

1. Sample a subset of the training data $\mathcal{D}_{\text{sub}} \subset \mathcal{D}$ 2. Optimize the hyperparameters on $\mathcal{D}_{\text{sub}}$ to obtain $\theta_{\text{sub}}^*$ 3. Use $\theta_{\text{sub}}^*$ as the initial values for optimization on the full dataset $\mathcal{D}$

This approach significantly reduces the computational cost of hyperparameter optimization while still achieving high-quality uncertainty estimates.

## 3.4   Implementation and Architecture

Our approach enhances pre-trained foundation models with uncertainty quantification capabilities by replacing the final layer with our two-stage GP framework. The implementation follows these steps:

1. The foundation model is used as a feature extractor, with its parameters frozen after pre-training.

2. The penultimate layer outputs serve as input features for the two-stage GP.

3. In the first stage, a Kernel Ridge Regression (KRR) model with automatic kernel selection is trained to capture the mean prediction function, using Algorithm 1 from Chapter 2.

4. In the second stage, a zero-mean GP is trained on the residuals to model the uncertainty.

5. For classification tasks, we employ the Dirichlet classification approach described in Section 3.3.

6. The subsampling warm start strategy (Section 2.4) is used to efficiently optimize the GP hyperparameters when dealing with large datasets.

Figure 3.1 illustrates this integration. The foundation model processes the raw input data through its specialized architecture (transformer blocks for text, convolutional or attention layers for images), generating high-dimensional feature representations. These representations then serve as inputs to our two-stage GP, which produces both predictions and associated uncertainty estimates.

## 3.5 Experimental Methodology

### 3.5.1 Datasets and Tasks

For our experiments, we focused on two primary healthcare datasets:

1. **MedNLI Dataset** [81]: This dataset focuses on natural language inference in the clinical domain, requiring models to identify potential clinical outcomes based on past medical history documented in clinical notes. This task simulates the real-world scenario of extracting clinically relevant information from unstructured text.

2. **BreakHis Dataset** [91]: This dataset contains microscopic images of breast tissue, with the task of distinguishing between benign and malignant tumors based on histopathological images. This represents a safety-critical visual diagnostic task.

Both datasets underwent preprocessing specific to their modality. For MedNLI, we applied standard NLP preprocessing techniques including tokenization, sentence segmentation, and

Figure 3.1: Integration of Two-Stage GP with Pre-trained Foundation Models. The foundation model serves as a feature extractor, with its output features fed into the two-stage GP. The first stage uses Kernel Ridge Regression for mean prediction, while the second stage employs a zero-mean GP to model uncertainty.

entity recognition. For BreakHis, we applied image preprocessing including normalization, augmentation, and patch extraction for the Vision Transformer inputs.

## 3.5.2 Foundation Models

For our experiments, we used the following foundation models:

1. **ClinicalBERT** [100]: A BERT-based model pre-trained on clinical notes, adapted for the MedNLI task.

2. **BioGPT** [63]: A GPT-based model pre-trained on biomedical literature, also adapted for the MedNLI task.

3. **Vision Transformer (ViT)** [19]: A transformer-based image recognition model adapted for the BreakHis dataset.

In line with standard practice, we kept most of the parameters in these foundation models frozen and focused on training the final classification layer. This approach leverages the rich representations learned during pre-training while adapting the models to our specific healthcare tasks.

### 3.5.3   Baselines and Evaluation Metrics

For each foundation model, we compared three approaches for uncertainty quantification:

1. **Monte Carlo Dropout** [27]: A commonly used baseline for uncertainty quantification in deep learning models.

2. **Exact-GP**: A standard GP model applied to the foundation model features.

3. **Two-stage Exact-GP**: Our proposed approach, incorporating the two-stage methodology described in Chapter 2.

## 3.6   Results and Analysis

### 3.6.1   Uncertainty Quantification Performance

The experimental results, depicted in Figure 3.2, demonstrate the effectiveness of our two-stage GP approach in providing meaningful uncertainty quantification across different foundation models and tasks.

Key observations from the results include:

1. **Accuracy Differential**: For all methods, the accuracy of "certain" predictions is higher than that of "uncertain" predictions, validating the basic concept of uncertainty-aware prediction. However, the magnitude of this difference varies across methods.

Figure 3.2: Uncertainty quantification results in accuracy (%). "Certain" predictions consistently show higher accuracy than "uncertain" predictions across all methods, with Two-stage Exact-GP achieving the best separation.

2. **Superior Separation in Two-stage GP**: Our proposed two-stage Exact-GP consistently demonstrates the largest separation between "certain" and "uncertain" prediction accuracy across all foundation models. This indicates that our approach provides more meaningful uncertainty estimates that correlate strongly with predictive performance.

3. **Foundation Model Differences**: The magnitude of improvement varies across foundation models, with the text-based models (ClinicalBERT and BioGPT) showing particularly strong benefits from our two-stage approach. This suggests that the effectiveness of uncertainty quantification may depend on the nature of the data and the specific architecture of the foundation model.

4. **Overall Accuracy Maintenance**: Importantly, our uncertainty-aware approach maintains competitive overall accuracy compared to deterministic baselines, indicating that the addition of uncertainty quantification does not come at the cost of predictive performance.

### 3.6.2   Clinical Insights from Uncertainty Patterns

Further analysis of the uncertainty estimates revealed several interesting patterns:

1. **Uncertainty vs. Difficulty**: Cases with higher model uncertainty tended to correspond to clinically complex or ambiguous cases, suggesting that the uncertainty

estimates captured meaningful characteristics of the data.

2. **Dataset Coverage**: The distribution of uncertainty varied across different patient subgroups, potentially highlighting areas where the models might have less reliable performance due to limited training data.

3. **Decision Boundaries**: Uncertainty was typically highest near decision boundaries, consistent with theoretical expectations and providing valuable information for borderline cases where clinical judgment is most critical.

These findings demonstrate that our two-stage GP approach not only provides numerical uncertainty estimates but generates estimates that align with intuitive notions of predictive confidence and data complexity.

## 3.7   Limitations and Future Directions

### 3.7.1   Current Limitations

While our two-stage GP approach demonstrates promising results for uncertainty-aware health foundation models, several limitations remain:

1. **Computational Scalability**: Although more efficient than alternatives like Deep Ensembles, our approach still faces computational challenges with very large datasets. The exact GP component scales cubically with the number of training examples, potentially limiting applicability to massive clinical datasets without approximation methods.

2. **Foundation Model Dependence**: The quality of uncertainty estimates depends in part on the quality of the feature representations learned by the foundation model. If these representations fail to capture relevant aspects of the data, the GP's uncertainty estimates may be compromised.

3. **No Theoretical Guarantees**: Our approach does not provide any theoretical guarantees on the coverage of the uncertainty estimates. GP only works well with strong assumptions on the data distribution.

4. **Coupled Prediction and Uncertainty Estimation**: Our approach is a coupled prediction and uncertainty estimation framework. This means that the prediction and uncertainty estimation are not independent and the uncertainty estimation is influenced by the prediction.

5. **Classification Approximation**: As noted in Section 3.3, our application to classification tasks relies on an approximation via Dirichlet distributions. While effective in practice, this approach introduces additional hyperparameters that need to be carefully tuned.

   From a theoretical perspective, this approximation introduces two sources of error: (1) the approximation of the Gamma distribution with a LogNormal distribution and (2) the Monte Carlo approximation of the predictive integral. For the first source, the quality of the LogNormal approximation to Gamma depends on the value of the concentration parameters $\alpha_i$. The approximation is more accurate for larger values of $\alpha_i$ and deteriorates as $\alpha_i$ approaches zero, potentially leading to inaccurate uncertainty estimates for classes with very low probabilities. For the second source, the Monte Carlo approximation error decreases at a rate of $O(1/\sqrt{N})$ where $N$ is the number of samples, which can be computationally expensive to reduce to negligible levels.

6. **Temporal Dynamics**: Current implementation does not explicitly model temporal dependencies in longitudinal patient data, which are common in healthcare settings.

### 3.7.2   Future Research Directions

Several promising directions for future research emerge from this work:

1. **Scalable Approximations**: Investigating sparse GP approximations, inducing points methods, or neural network approximations could improve computational scalability while maintaining uncertainty quality.

2. **End-to-End Training**: Developing methods for jointly training the foundation model and GP components could potentially yield better aligned feature representations and uncertainty estimates.

3. **Multi-task Learning**: Extending the framework to support multiple related healthcare prediction tasks simultaneously could improve sample efficiency and uncertainty estimation across tasks.

4. **Temporal Uncertainty**: Incorporating explicit temporal modeling within the GP framework to better capture uncertainty evolution in longitudinal patient trajectories.

In the next chapter, we will explore how our two-stage GP approach can be adapted to another important domain: climate and weather foundation models, where uncertainty quantification plays an equally critical role in supporting informed decision-making in the face of complex, chaotic systems.

## 3.8 Summary

This chapter has demonstrated the application of the two-stage Gaussian Process methodology to health foundation models, addressing the critical need for uncertainty quantification in healthcare AI. Our approach enhances pre-trained foundation models by replacing their deterministic final layers with a principled two-stage GP framework that provides robust and efficient uncertainty estimates while maintaining predictive performance.

We provided a detailed mathematical formulation of how to extend our two-stage GP framework to classification tasks using Dirichlet distributions based on [67], enabling its application to a wide range of healthcare prediction problems. Experimental results across

multiple healthcare prediction tasks show that our approach outperforms baseline uncertainty quantification methods in terms of calibration quality, selective prediction performance, and the meaningful separation between high-confidence and low-confidence predictions.

The challenges of health data prediction—including data heterogeneity, limited labeled data, and high-stakes decisions—make uncertainty quantification especially important in this domain. Our two-stage GP approach offers a principled solution to this need, providing healthcare AI systems with the capacity to express doubt about their predictions in cases where such doubt is warranted.

In the next chapter, we will explore how to make GP more efficient and resolve the biases introduced by iterative methods when computing GP.

# Chapter 4

# Kernel Preconditioning and Unbiased Gaussian Processes for Computational Efficiency

## 4.1   Introduction

Large-scale Gaussian Processes (GPs) represent a critical tool for uncertainty quantification across various domains. However, despite their mathematical elegance and principled uncertainty estimates, they face significant computational challenges that limit their application in modern machine learning contexts. This chapter addresses a fundamental bottleneck in GP regression: the computational complexity associated with kernel matrix operations, which becomes prohibitively expensive as datasets grow in size.

When applying GPs to real-world problems, practitioners must frequently estimate the log marginal likelihood (LML) and its gradient to optimize hyperparameters. The standard approach requires solving linear systems involving the kernel matrix and computing log determinants, both of which scale cubically with the number of data points. As datasets grow beyond a few thousand points, these operations become intractable on standard computing

hardware, creating a barrier to GP application in many modern machine learning scenarios where uncertainty quantification is most needed.

The previous chapters have demonstrated the effectiveness of GPs for uncertainty quantification in various domains, particularly in health foundation models as discussed in Chapter 3. This chapter extends that work by addressing the practical computational challenges that arise when scaling these methods to larger datasets. By developing more efficient computational techniques, we enable the broader application of GPs in contexts where uncertainty quantification is critical but computational resources are limited.

This chapter makes two key contributions to address these challenges:

1. We introduce the Adaptive Factorized Nyström (AFN) preconditioner, a novel approach that significantly accelerates iterative solvers for GP inference. By combining the strengths of the Nyström approximation with a sparse correction term, AFN provides more robust and efficient preconditioning than existing methods.

2. We develop an unbiased stochastic estimation technique for the LML and its gradient based on the Single-Sample Conjugate Gradient (SS-CG) method. This approach maintains statistical unbiasedness while reducing variance, making it practical for scaling GP inference to large datasets.

Together, these innovations enable efficient and accurate GP inference on datasets that were previously intractable, expanding the practical utility of GPs for uncertainty quantification in complex real-world applications. Through comprehensive theoretical analysis and extensive experiments, we demonstrate that our approach outperforms existing methods across a wide range of problems, particularly for the challenging case of middle-range length scales where alternative methods often fail.

The remainder of this chapter is organized as follows: Section 4.2 provides the mathematical background necessary to understand GP inference and the computational challenges it faces. Section 4.3 surveys existing approaches to address these challenges, highlighting their strengths

and limitations. Sections 4.4 and 4.5 detail our novel contributions: the AFN preconditioner and unbiased estimation techniques. Section 4.6 presents rigorous theoretical analysis of our methods, while Section 4.7 discusses practical implementation considerations. Section 4.8 reports comprehensive experimental results on both synthetic and real-world datasets. Finally, Section 4.9 summarizes our findings and discusses their implications for uncertainty quantification in machine learning.

## 4.2   Background on Gaussian Process Inference

Given the training dataset $\mathcal{X} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N = (\mathbf{X}, \mathbf{y})$, Gaussian process regression (GPR) models the nonlinear relationship between $\mathbf{x}_i$ and $\mathbf{y}_i$ by a latent function $f(\cdot) \sim \mathcal{GP}(\mu, k)$ with mean function $\mu(\cdot)$ and covariance function $k(\cdot, \cdot)$. And the latent function is usually contaminated by some noise $\epsilon \sim N(0, \sigma^2 \mathbf{I})$ such that $\mathbf{y}_i = f(\mathbf{x}_i) + \epsilon_i$. Given a testing dataset $X_* = \{\mathbf{x}_*^j\}_1^M$, according to the Gaussian process assumption, the following holds:

$$(\underbrace{f(\mathbf{x}_1), \cdots f(\mathbf{x}_N)}_{\mathbf{f}_X^\top}, \underbrace{f(\mathbf{x}_*^1), \cdot, f(\mathbf{x}_*^M)}_{\mathbf{f}_*^\top})^\top \sim \mathcal{GP}(\mu, \mathbf{K})$$

where $\mu = (\underbrace{\mu(\mathbf{x}_1), \cdots \mu(\mathbf{x}_N)}_{\mu_\mathbf{X}^\top}, \underbrace{\mu(\mathbf{x}_*^1), \cdot, \mu(\mathbf{x}_*^M)}_{\mu_*^\top})^\top$ and $\mathbf{K} = \begin{bmatrix} \mathbf{K_{XX}} + \sigma^2 \mathbf{I} & \mathbf{K_{X*}} \\ \mathbf{K_{*X}} & \mathbf{K_{**}} \end{bmatrix}$. Here $\mathbf{K_{XX}}$, $\mathbf{K_{**}}$ and $\mathbf{K_{X*}}$ are defined by evaluating the kernel function over the training set, testing set and their cross-interaction, respectively. Given this, the prediction and covariance for the target of testing data can be derived as follows

$$\mathbb{E}[\mathbf{f}_*] = \mu(\mathbf{X}_*) - \mathbf{K_{*X}}(\mathbf{K_{XX}} + \sigma^2 \mathbf{I})^{-1}(\mathbf{y} - \mu(\mathbf{X})), \tag{4.1}$$

$$\mathrm{Cov}(\mathbf{f}_*) = \mathbf{K_{**}} - \mathbf{K_{*X}}(\mathbf{K_{XX}} + \sigma^2 \mathbf{I})^{-1}\mathbf{K_{X*}}. \tag{4.2}$$

## 4.2.1 Kernel Functions in Gaussian Process Regression

The kernel function in Gaussian Process Regression (GPR) is typically predefined with some unknown parameters. Some widely used kernels in various applications will be considered in this chapter:

- Gaussian kernel: $k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x}-\mathbf{x}'\|_2^2}{2l^2}}$,

- exponential kernel: $k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x}-\mathbf{x}'\|_2}{l}}$,

- Matérn-3/2 kernel: $k(\mathbf{x}, \mathbf{x}') = (1 + \frac{\sqrt{3}\|\mathbf{x}-\mathbf{x}'\|_2}{l})e^{-\frac{\sqrt{3}\|\mathbf{x}-\mathbf{x}'\|_2}{l}}$,

- Matérn-5/2 kernel: $k(\mathbf{x}, \mathbf{x}') = (1 + \frac{\sqrt{5}\|\mathbf{x}-\mathbf{x}'\|_2}{l} + \frac{5\|\mathbf{x}-\mathbf{x}'\|_2^2}{3l^2})e^{-\frac{\sqrt{5}\|\mathbf{x}-\mathbf{x}'\|_2}{l}}$,

all of which have a hyperparameter, the length-scale $l$. When modeling noisy observations, another hyperparameter $\sigma^2$ is introduced. Finding the optimal hyperparameters for the prediction in the GPR model is a crucial task, as different sets of hyperparameters lead to different models. This can be done by maximizing the LML of the GPR,

$$\widehat{\theta} = \arg\max_{\theta} L(\theta) = \log p(\mathbf{y}|\mathbf{X}, \theta) = -(\frac{1}{2}\mathbf{y}^{\top}\widehat{\mathbf{K}}^{-1}\mathbf{y} + \log|\widehat{\mathbf{K}}| + n\log(2\pi)) \qquad (4.3)$$

where $\theta = (l, \sigma)^{\top}$ and $\widehat{\mathbf{K}} = K_{\mathbf{X}\mathbf{X}} + \sigma^2\mathbf{I}$. The gradient of LML can be derived as follows:

$$\frac{\partial L}{\partial \theta} = \frac{1}{2}\mathbf{y}^{\top}\widehat{\mathbf{K}}^{-1}\frac{\partial \widehat{\mathbf{K}}}{\partial \theta}\widehat{\mathbf{K}}^{-1}\mathbf{y} - \text{tr}(\widehat{\mathbf{K}}^{-1}\frac{\partial \widehat{\mathbf{K}}}{\partial \theta}), \qquad (4.4)$$

where $\frac{\partial \widehat{\mathbf{K}}}{\partial \theta}$ is the derivative of the kernel matrix with respect to the hyperparameters. Since the kernel matrix is dense, exact evaluation of LML and its gradient is prohibitively expensive. Therefore, iterative methods are usually used to derive an estimation of LML and its gradient.

## 4.3 Related Works

### 4.3.1 Trace Estimation

The assessment of both Log Marginal Likelihood (LML) and its gradient heavily relies on accurate trace estimation. One commonly used method for approximating the trace of a matrix $\mathbf{A}$ is Hutchinson's trace estimator [44]:

$$\text{Tr}(\mathbf{A}) = \mathbb{E}_{\mathbf{u} \sim N(\mathbf{0}, \mathbf{I})}[\mathbf{u}^\top \mathbf{A} \mathbf{u}] \approx \frac{1}{k} \sum_{i=1}^{k} \mathbf{u}_i^\top \mathbf{A} \mathbf{u}_i \tag{4.5}$$

Here, the expectation is taken over $\mathbf{u}$ drawn from a standard normal distribution $N(\mathbf{0}, \mathbf{I})$. Alternatives like Rademacher random vectors, whose entries are either 1 or $-1$ with equal probability, can also be employed. For estimating the trace of a matrix function $f(\mathbf{A})$ where $\mathbf{A}$ is a positive definite matrix, the stochastic Lanczos trace estimator is frequently utilized [96]:

$$\text{Tr}(f(\mathbf{A})) = \mathbb{E}_{\mathbf{u} \sim N(\mathbf{0}, \mathbf{I})}[\mathbf{u}^\top f(\mathbf{A}) \mathbf{u}] \approx \frac{1}{k} \sum_{i=1}^{k} \mathbf{u}_i^\top \mathbf{V}_m f(\mathbf{T}^m) \mathbf{V}_m^\top \mathbf{u}_i \tag{4.6}$$

In this formula, $\mathbf{A} \approx \mathbf{V}_m \mathbf{T}^m \mathbf{V}_m^\top$ is approximate factorization of $\mathbf{A}$ derived by running the m steps Lanczos algorithm. It should be noted that if $m < N$, this estimation is biased [32]. This bias is introduced by Gaussian quadrature.

### 4.3.2 Iterative GP

In summary, evaluating the loss and its gradient requires:

- Solving a linear system: $\widehat{\mathbf{K}} \mathbf{u} = \mathbf{y}$

- Estimate the log determinant $\log |\widehat{\mathbf{K}}| = \sum_{i=1}^{N} \log \lambda_i(\widehat{\mathbf{K}})$

- Estimate the trace $\text{tr}(\widehat{\mathbf{K}}^{-1} \frac{\partial \widehat{\mathbf{K}}}{\partial \theta})$

- Matrix-vector product (MVP) for $\frac{\partial \widehat{\mathbf{K}}}{\partial \theta}$

Computing the log determinant and solving a linear system directly can be computationally expensive, often requiring cubic time complexity. To mitigate this, the Conjugate Gradient (CG) method is commonly employed. This approach falls under the category of Iterative Gaussian Processes (Iter-GP). Iter-GP approximates solutions to linear systems $\widehat{\mathbf{K}}\mathbf{y}$ and $\widehat{\mathbf{K}}\mathbf{z}_i$ using $m$-step CG. Concurrently, it stores a partial Lanczos factorization $\mathbf{T} \in \mathbb{R}^{m \times m}$ associated with random initial probe vectors $\mathbf{z}_i$ for estimating the log determinant through stochastic Lanczos trace estimation (SLT) [28]. However, it's important to note that this approach introduces a systematic bias, stemming from the inherent bias in the CG method [74].

### 4.3.3  Preconditioning

When dealing with Gaussian Processes (GPs), the matrix $\widehat{\mathbf{K}}$ is often ill-conditioned, making it challenging to efficiently estimate the Log Marginal Likelihood (LML) and its gradient. The application of the Conjugate Gradient (CG) method in such cases can be particularly slow to converge and tends to produce high-variance estimates for the log determinant. As a result, the Preconditioned Conjugate Gradient (PCG) method is commonly used to expedite convergence. Furthermore, Wenger et al. [102] shows that preconditioner can also reduce the variance of the log determinant estimation for the LML and trace estimation for its gradient. Their method is based on the following observations:

$$\log |\widehat{\mathbf{K}}| = \log |\mathbf{P}\mathbf{P}^{-1}\widehat{\mathbf{K}}| = \log |\mathbf{P}| + \log |\mathbf{P}^{-1/2}\widehat{\mathbf{K}}\mathbf{P}^{-1/2}| \tag{4.7}$$

$$\frac{\partial \log |\widehat{\mathbf{K}}|}{\partial \theta} = \mathrm{tr}(\mathbf{P}^{-1}\frac{\partial \mathbf{P}}{\partial \theta}) + \mathrm{tr}(\widehat{\mathbf{K}}^{-1}\frac{\partial \widehat{\mathbf{K}}}{\partial \theta} - \mathbf{P}^{-1}\frac{\partial \mathbf{P}}{\partial \theta}). \tag{4.8}$$

where $\mathbf{P}$ is a preconditioner matrix. If the preconditioner exhibits a favorable structure, enabling easy computation of $\log |\mathbf{P}|$ and $\mathrm{tr}(\mathbf{P}^{-1}\frac{\partial \mathbf{P}}{\partial \theta})$ and significantly reduces the condition number of $\mathbf{P}^{-1}\widehat{\mathbf{K}}$, then both the LML and its gradient will converge more rapidly and exhibit

lower variance.

## 4.3.4 Unbiased Estimation

Even though the preconditioner can both reduce variance and speed up convergence, the estimation remains inherently biased. Research by Golub and Meurant [32] and Potapczynski et al. [74] shows that using CG for LML estimation yields a biased outcome. Specifically, the first term of the LML tends to be underestimated, while the second term is often overestimated. This imbalance skews the model towards under-fitting the training data, as the first term gauges model fit and the second assesses model complexity. To counteract this bias, two unbiased estimators studied in [6] could be used. These two estimators leverage randomized truncation techniques. Given a series $\bar{\bar{\Phi}} = \sum_{i=1}^{H} \Delta_i$, and a probability distribution $\mathbb{P}(\mathcal{J})$ over $\{1, \ldots, H\}$ there are two methods leading to an unbiased estimator of $\Phi$:

- **Method 1: Russian-Roulette estimator**

$$\hat{\Phi} = \sum_{i=1}^{J} \frac{\Delta_i}{\mathbb{P}(\mathcal{J} \geq i)}, \text{ where } J \sim \mathbb{P}(\mathcal{J}), \tag{4.9}$$

- **Method 2: Single-Sample estimator**

$$\Phi = \sum_{i=1}^{H} \frac{\Delta_i}{\mathbb{P}(\mathcal{J} = i)} \mathbb{I}(J = i), \text{ where } J \sim \mathbb{P}(\mathcal{J}). , \tag{4.10}$$

The authors in [74] then proposed two unbiased estimators for $\mathbf{y}^\top \widehat{\mathbf{K}}^{-1} \mathbf{y}$ and $\log |\widehat{\mathbf{K}}|$ based on Russian-Roulette CG (RR-CG) and Single-Sample Random Fourier Features (SS-RFF). The key to apply these two estimators to GP is to note the fact that both exact linear solve and logdet can be written as a summation:

$$\widehat{\mathbf{K}}^{-1} \mathbf{y} = \sum_{i=1}^{N} \gamma_i \mathbf{d}_i \tag{4.11}$$

Although $\log |\widehat{\mathbf{K}}| = \mathbb{E}_{\mathbf{z}}[\|\mathbf{z}\|^2 \mathbf{e}^\top (\log \mathbf{T}_{\mathbf{z}}^N) \mathbf{e}]$ itself is not a series form which can be directly applied to the these two unbiased estimator, it can be written as a telescopic sum. For a fixed probe vector $\mathbf{z}_i$,

$$\mathbf{z}_i^\top (\log \widehat{\mathbf{K}}) \mathbf{z}_i = \|\mathbf{z}_i\|^2 \mathbf{e}^\top (\log \mathbf{T}_{\mathbf{z}_i}^{J_{\min}}) \mathbf{e} + \sum_{j=J_{\min}}^{N-1} \|\mathbf{z}_i\|^2 \mathbf{e}^\top (\log \mathbf{T}_{\mathbf{z}_i}^{j+1} - \log \mathbf{T}_{\mathbf{z}_i}^i) \mathbf{e}. \tag{4.12}$$

The size of $T_{\mathbf{z}_i}^{i+1}$ and $T_{\mathbf{z}_i}^i$ are $i+1$ and $i$ respectively, but for notational simplicity, we write their difference directly. There should be no ambiguity from the context. Here $J_{\min}$ is usually taken to be a constant larger than 1, which is the minimum number of steps to run CG. However, these two estimators are not effective due to the high cost. Our method based on SS will be more practical.

## 4.3.5   Stochastic Gradient Descent for Gaussian Process

An alternative strategy to tackle bias in Gaussian Process modeling involves the integration of Stochastic Gradient Descent (SGD). Chen et al. [16] established that mini-batch gradient descent assures hyperparameter convergence towards optimal values, even when faced with biased gradient estimates due to sample correlation. In contrast, our approach employs what can be termed as 'exact stochastic gradient descent', wherein we compute the unbiased exact gradient at each iteration. This feature allows us to directly apply all existing theoretical guarantees related to SGD. Importantly, when the gradient is explicitly calculable as in Equation (4.4), the need for loss function evaluation at each step is eliminated, yielding substantial computational savings compared to Iter-GP. Yet, for cases where the likelihood is not Gaussian or the preconditioner's gradient is not explicitly computable, as in Equation 4.7, we utilize automatic differentiation for gradient computation [28, 102], ensuring that our gradient estimates remain unbiased.

Figure 4.1: Surface plot of the log LML for optimization path of LBFGS with strong Wolfe. We use GP to learn the function $f(x) = \sin(2\pi x)e^x + x^2 + \epsilon$ where $\epsilon \sim N(0, 1)$.

### 4.3.6 Barely Biased GP

In the paper by Burt et al. [12], the authors use a simple inequality to guide the performance of the Conjugate Gradient (CG) method, terminating it when a set accuracy threshold is met. This bears a superficial resemblance to our small-bias *PredSS-CG*. However, a key distinction exists: our *PredSS-CG* is designed to be inherently unbiased. Moreover, when aiming for a specific accuracy level $\epsilon$, Burt et al.'s method must execute enough steps to deterministically meet this error tolerance. In contrast, our approach only requires the bias to fall below $\epsilon$, which, on average, necessitates fewer iterations.

## 4.4 Adaptive Factorized Nyström Preconditioner

The Adaptive Factorized Nyström (AFN) preconditioner is a key component of our approach to efficient unbiased Gaussian Processes. In this section, we describe the construction and properties of the AFN preconditioner.

Figure 4.2: Surface plot of the log LML for optimization path of GD with backtracking. We use GP to learn the function $f(x) = \sin(2\pi x)e^x + x^2 + \epsilon$ where $\epsilon \sim N(0, 1)$.

### 4.4.1 Background: Nyström Approximation

The Nyström method is a well-established technique for approximating kernel matrices, originally derived from quadrature rules for integral equations. Given a kernel matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$, the Nyström approximation constructs a low-rank approximation by selecting a subset of $m$ landmark points (where $m \ll n$).

Let $X_k \subset X$ be a set of $k$ landmark points selected from the full dataset $X$. The Nyström approximation of the kernel matrix $\mathbf{K}$ can be expressed as:

$$\mathbf{K}_{\mathrm{Nys}} = \mathbf{K}_{X,X_k} \mathbf{K}_{X_k,X_k}^{-1} \mathbf{K}_{X_k,X} \tag{4.13}$$

where $\mathbf{K}_{X,X_k}$ represents the kernel matrix between all points in $X$ and the landmark points in $X_k$, and $\mathbf{K}_{X_k,X_k}$ is the kernel matrix among the landmark points.

The quality of the Nyström approximation heavily depends on the selection of landmark points. The approximation error is theoretically bounded by the fill distance of the landmark

points, which measures how well the landmark points cover the input space. Specifically, for Gaussian kernels, the approximation error decreases exponentially with decreasing fill distance.

When using the Nyström approximation as a preconditioner for the system $(\mathbf{K} + \mu\mathbf{I})$, the standard approach applies the Sherman-Morrison-Woodbury (SMW) formula, which requires solving systems involving matrices of size $k \times k$. This approach works well when $k$ is small, but becomes computationally challenging and potentially numerically unstable when $k$ is large.

## 4.4.2   Factorized Sparse Approximate Inverse (FSAI)

The Factorized Sparse Approximate Inverse (FSAI) method, developed by Kolotilina and Yeremin, is a technique for computing a sparse approximate inverse of a symmetric positive definite matrix. Given a matrix $\mathbf{A}$ and a sparsity pattern $\mathcal{S}$, FSAI computes a lower triangular matrix $\mathbf{G}$ such that $\mathbf{G}^\top\mathbf{G} \approx \mathbf{A}^{-1}$.

A key advantage of FSAI is that it only requires the entries of $\mathbf{A}$ corresponding to the sparsity pattern of $\mathbf{G}$ and $\mathbf{G}^\top$. This property is particularly valuable when dealing with dense matrices where computing and storing the full matrix would be prohibitively expensive. Additionally, the computation of each row of $\mathbf{G}$ is independent of other rows, making the algorithm highly parallelizable.

In the context of kernel matrices, the sparsity pattern for FSAI is typically based on the geometric proximity of data points. For each point, we consider its nearest neighbors to determine the non-zero entries in the corresponding row of $\mathbf{G}$.

## 4.4.3   Construction of the AFN Preconditioner

The AFN preconditioner is designed to efficiently approximate the inverse of the kernel matrix $\mathbf{K}$. It combines the strengths of the Nyström approximation with a sparse correction term computed using FSAI to improve accuracy. By avoiding the use of the SMW formula, AFN

remains effective even when a large number of landmark points is needed for accuracy. The construction process is outlined in Algorithm 4.

---
**Algorithm 4** Construction of the AFN Preconditioner

---
1: **Input:** Kernel matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$, number of landmark points $m$, sparsity pattern $\mathcal{S}$
2: **Output:** AFN preconditioner $\mathbf{M}$
3: Select $m$ landmark points using Farthest Point Sampling (FPS) or k-means clustering
4: Partition the kernel matrix as: $\mathbf{K} = \begin{bmatrix} \mathbf{K}_{11} & \mathbf{K}_{12} \\ \mathbf{K}_{12}^\top & \mathbf{K}_{22} \end{bmatrix}$ where $\mathbf{K}_{11} \in \mathbb{R}^{m \times m}$ is the kernel matrix for landmark points
5: Compute the Cholesky factorization of $\mathbf{K}_{11} + \mu\mathbf{I} = \mathbf{L}\mathbf{L}^\top$
6: Compute the Schur complement $\mathbf{K}_{\text{Schur}} = \mathbf{K}_{22} - \mathbf{K}_{12}^\top(\mathbf{K}_{11} + \mu\mathbf{I})^{-1}\mathbf{K}_{12}$
7: Apply Factorized Sparse Approximate Inverse (FSAI) to $\mathbf{K}_{\text{Schur}}$ using the sparsity pattern $\mathcal{S}$ to obtain $\mathbf{G}$
8: Set $\mathbf{M}_{11} = \mathbf{K}_{11} + \mu\mathbf{I}$
9: Set $\mathbf{M}_{12} = \mathbf{K}_{12}$
10: Set $\mathbf{M}_{21} = \mathbf{K}_{12}^\top$
11: Set $\mathbf{M}_{22} = (\mathbf{G}^\top\mathbf{G})^{-1} + \mathbf{K}_{12}^\top(\mathbf{K}_{11} + \mu\mathbf{I})^{-1}\mathbf{K}_{12}$
12: Construct $\mathbf{M}$ in factorized form: $\mathbf{M} = \begin{bmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{K}_{12}^\top\mathbf{L}^{-\top} & \mathbf{G}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{L}^\top & \mathbf{L}^{-1}\mathbf{K}_{12} \\ \mathbf{0} & \mathbf{G}^{-\top} \end{bmatrix}$

---

## 4.4.4 Factorized Form of the AFN Preconditioner

The AFN preconditioner has the following factorized form:

$$\mathbf{M} = \underbrace{\begin{bmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{K}_{12}^\top\mathbf{L}^{-\top} & \mathbf{G}^{-1} \end{bmatrix}}_{\mathbf{U}^\top} \begin{bmatrix} \mathbf{L}^\top & \mathbf{L}^{-1}\mathbf{K}_{12} \\ \mathbf{0} & \mathbf{G}^{-\top} \end{bmatrix}. \tag{4.14}$$

It is similar to the Nyström preconditioner. However, it adds a sparse correction to the residual matrix to make it more accurate. Since AFN is more accurate than Randomized Nyström preconditioner (RAN) [25], under the same conditions as RAN, the preconditioned matrix is guaranteed to have a condition number smaller than 28.

Multiplying the factors gives $\mathbf{M}$ another form:

$$\mathbf{M} = \begin{bmatrix} \mathbf{K}_{11} + \mu\mathbf{I} & \mathbf{K}_{12} \\ \mathbf{K}_{12}^\top & (\mathbf{G}^\top\mathbf{G})^{-1} + \mathbf{K}_{12}^\top(\mathbf{K}_{11} + \mu\mathbf{I})^{-1}\mathbf{K}_{12} \end{bmatrix} \tag{4.15}$$

$$= \mathbf{K}_{nys} + \mu\mathbf{I} + \underbrace{\begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & (\mathbf{G}^\top\mathbf{G})^{-1} + \mathbf{K}_{12}^\top\left((\mathbf{K}_{11} + \mu\mathbf{I})^{-1} - (\mathbf{K}_{11})^{-1}\right)\mathbf{K}_{12} - \mu\,\mathbf{I} \end{bmatrix}}_{\text{Correction term}}, \tag{4.16}$$

where $\mathbf{K}_{nys}$ is the Nyström approximation of the kernel matrix.

## 4.5 Unbiased Log Marginal Likelihood Estimation and its Gradient

Although RR estimator in Equation 4.9 and SS estimator in Equation 4.10 are both unbiased, SS estimator has a lower cost but a higher variance than RR estimator. Therefore, in order to take advantage of the low cost but reduce the variance, we propose the Preconditioned Single-Sample CG (*PredSS-CG* ).

### 4.5.1 SS-CG

After we write the exact solve as summation as in Equation (4.11), we can apply SS estimator to CG to derive the SS-CG. It only samples one term from the series.

$$\hat{\mathbf{u}}_0^i = \sum_{j=1}^n \frac{\gamma_j^i}{\mathbb{P}(\mathcal{J} = j)}\mathbf{d}_j^i\mathbb{I}\{j = J_i\} = \frac{\gamma_{J_i}^i}{\mathbb{P}(\mathcal{J} = J_i)}\mathbf{d}_{J_i}^i. \tag{4.17}$$

In order to further reduce variance, we make a slight modification to it and denote the following scheme as SS-CG. We rewrite the solution from N-step CG as

$$\widehat{\mathbf{K}}^{-1}\mathbf{y} = S_{J_{\min}} + \sum_{i=J_{\min}}^{N-1} (S_{i+1} - S_i) \tag{4.18}$$

where $S_p = \sum_{i=1}^{p} \gamma_i \mathbf{d}_i$. Then we define the solution from SS-CG as

$$\hat{\mathbf{u}}_0^i = S_{J_{\min}} + \frac{\gamma_{J_i}^i}{\mathbb{P}(\mathcal{J} = J_i)}\mathbf{d}_{J_i}^i := \phi_{J_i}. \tag{4.19}$$

Similar to the log-determinant estimation, we also use $J_{\min}$ for CG so that the solution at least is not extremely far away from the true solution. Then the inverse quadratic form can be estimated as

$$\mathbf{y}^\top\widehat{\mathbf{K}}^{-1}\mathbf{y} \approx \frac{1}{l}\sum_{i=1}^{l} \mathbf{y}^\top\phi_{J_i} := \mathcal{L}_1. \tag{4.20}$$

Here we just need to run $l$ times unbiased estimator and take average to get the unbiased estimate of $\widehat{\mathbf{K}}^{-1}\mathbf{y}$. Recall for a fixed $\mathbf{z}$, $\mathbf{z}^\top\log(\widehat{\mathbf{K}})\mathbf{z}$ can be written as a series expansion in Equation 4.12. Then SS estimator only samples one term from the summation while keeping the first term.

$$\mathbf{z}_i^\top(\log\widehat{\mathbf{K}})\mathbf{z}_i \approx \underbrace{\|\mathbf{z}_i\|^2\mathbf{e}(\log T_{\mathbf{z}_i}^{J_{\min}})\mathbf{e} + \frac{\|\mathbf{z}_i\|^2\mathbf{e}^\top(\log \mathbf{T}_{\mathbf{z}_i}^{J+1} - \log \mathbf{T}_{\mathbf{z}_i}^{J})\mathbf{e}}{\mathbb{P}(\mathcal{J} = J)}}_{\psi_{\mathbf{z}_i,J}}, J \sim \mathbb{P}(\mathcal{J}). \tag{4.21}$$

Finally, the log-determinant can be estimated as

$$\log|\widehat{\mathbf{K}}| \approx \frac{1}{l}\sum_{j=1}^{l}\frac{1}{k}\sum_{i=1}^{k}\psi_{\mathbf{z}_i,J_j} := \mathcal{L}_2. \tag{4.22}$$

For the gradient of LML in Equation. 4.4, note we can reuse the unbiased estimation for estimating LML. We have obtained an estimate for $\hat{\mathbf{K}}^{-1}\mathbf{y} \approx \phi_{J_i}$. In order to get the partial

Lanczos factorization for estimating the log-determinant, the linear system $\hat{\mathbf{K}}^{-1}\mathbf{x} = \mathbf{z}_k$ is also solved using SS-CG. Denote the estimation by $\phi_{J_i,k}$. Then for the trace term in gradient, we can estimate as

$$\text{tr}(\hat{\mathbf{K}}^{-1}\frac{\partial \hat{\mathbf{K}}}{\partial \theta}) \approx \frac{1}{l}\sum_{i=1}^{l}\frac{1}{k}\sum_{i=1}^{k}\phi_{J_i,k}^{\top}\frac{\partial \hat{\mathbf{K}}}{\partial \theta}\mathbf{z}_i := \mathcal{DL}_1 \tag{4.23}$$

$$\frac{1}{2}\mathbf{y}^{\top}\hat{\mathbf{K}}^{-1}\frac{\partial \widehat{\mathbf{K}}}{\partial \theta}\hat{\mathbf{K}}^{-1}\mathbf{y} = \frac{1}{2l}\sum_{i=1}^{l}\phi_{J_i,1}^{\top}\frac{\partial \widehat{\mathbf{K}}}{\partial \theta}\phi_{J_i,2} := \mathcal{DL}_2 \tag{4.24}$$

Note $\phi_{J_i,1}$ and $\phi_{J_i,2}$ mean two independent run of unbiased estimation. Otherwise, we can not guarantee the estimation is unbiased.

## 4.5.2 Theoretical Analysis of SS Estimator

Firstly, we prove the unbiasedness of a general SS estimator and compute its variance. This estimator has been studied in [64], [6] and [74]. For self-completeness, we provide proof here.

**Theorem 4.5.1.** *Denote the SS estimator defined in Equation 4.10 by $\Phi$. It is an unbiased estimator and has a variance*

$$Var(\Phi) = \sum_{i=1}^{H}\frac{\Delta_i^2(1 - \mathbb{P}(\mathcal{J} = i))}{\mathbb{P}(\mathcal{J} = i)}. \tag{4.25}$$

*Proof.*

$$\mathbb{E}_{\mathcal{J}}[\Phi] = \sum_{i=1}^{H}\frac{\Delta_i}{\mathbb{P}(\mathcal{J} = i)}\mathbb{E}(\mathbb{I}(J = i)) = \sum_{i=1}^{H}\frac{\Delta_i}{\mathbb{P}(\mathcal{J} = i)}\mathbb{P}(\mathcal{J} = i) = \sum_{i=1}^{H}\Delta_i = \Phi. \tag{4.26}$$

The variance can be computed as

$$\text{Var}(\Phi) = \sum_{i=1}^{H} \frac{\Delta_i^2}{\mathbb{P}(\mathcal{J} = i)^2} \text{Var}(\mathbb{I}(J = i)) \tag{4.27}$$

$$= \sum_{i=1}^{H} \frac{\Delta_i^2}{\mathbb{P}(\mathcal{J} = i)^2} \mathbb{P}(\mathcal{J} = i)(1 - \mathbb{P}(\mathcal{J} = i)) \tag{4.28}$$

$$= \sum_{i=1}^{H} \frac{\Delta_i^2 (1 - \mathbb{P}(\mathcal{J} = i))}{\mathbb{P}(\mathcal{J} = i)} \tag{4.29}$$

$$= \sum_{i=1}^{H} \frac{\Delta_i^2}{\mathbb{P}(\mathcal{J} = i)} - \sum_{i=1}^{H} \Delta_i^2. \tag{4.30}$$

$\square$

**Remark 1.** *In [74], the authors show that* $\mathbb{P}(\mathcal{J} = i) = \frac{\Delta_i}{\sum_i^H \Delta_i}$ *minimizes the variance of SS estimator. In practice, we cannot compute this distribution, so we just try some different distributions in our experiments. However, the authors in [74] suggested to use* $\Delta_i \sim e^{-\lambda i}$ *as an approximation to* $\Delta_i \sim C^{-i}$ *where C is a number depending on the conditioning of the kernel matrix.*

Then we can bound the variance of SS estimator for inverse quadratic term and log-determinant.

**Theorem 4.5.2.** *LML estimated by SS-CG is unbiased. The variance of LML estimated by SS estimator is bounded by*

$$Var(\mathcal{L}) \leq \left(\frac{2C_\rho}{lk} + C_1\right) \sum_{i=J_{\min}}^{N} \frac{1 - \mathbb{P}(\mathcal{J} = i)}{\mathbb{P}(\mathcal{J} = i)} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}\right)^{4i} \tag{4.31}$$

*Proof.* Firstly, according to the Equation 4.25, the variance of SS estimator for the inverse

quadratic term is

$$\text{Var}(\mathbf{y}^\top \phi_{J_i}) = \sum_{i=J_{\min}}^{N} \frac{\gamma_i^2 (\mathbf{y}^\top \mathbf{d}_i)^2}{\mathbb{P}(\mathcal{J} = i)} - \sum_{i=J_{\min}}^{N} \gamma_i^2 (\mathbf{y}^\top \mathbf{d}_i)^2 \tag{4.32}$$

$$\tag{4.33}$$

Then according to the convergence rate of CG, we have

$$\sqrt{\mu} \|\gamma_p \mathbf{d}_p\| \leq \|\gamma_p \mathbf{d}_p\|_{\widehat{\mathbf{K}}} = \|S_p - S_{p-1}\|_{\widehat{\mathbf{K}}} \tag{4.34}$$

$$\leq \|S_p - \widehat{\mathbf{K}}^{-1} \mathbf{y}\|_{\widehat{\mathbf{K}}} + \|\widehat{\mathbf{K}}^{-1} \mathbf{y} - S_{p-1}\|_{\widehat{\mathbf{K}}} \tag{4.35}$$

$$\leq 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^p \|\mathbf{y}\|_{\widehat{\mathbf{K}}} + 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{p+1} \|\mathbf{y}\|_{\widehat{\mathbf{K}}} \tag{4.36}$$

$$\leq 4 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^p \|\mathbf{y}\|_{\widehat{\mathbf{K}}} \tag{4.37}$$

The first equality comes from the fact that $\widehat{\mathbf{K}} \succeq \mu \mathbf{I}$. Then we have

$$\|\gamma_p \mathbf{d}_p\|^2 \leq \frac{16}{\mu} \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2p} \|\mathbf{y}\|_{\widehat{\mathbf{K}}}^2 \tag{4.38}$$

By Cauchy-Schwarz inequality, we have

$$|\gamma_p \mathbf{y}^\top \mathbf{d}_p|^2 \leq \|\mathbf{y}\|^2 \|\gamma_p \mathbf{d}_p\|^2 \tag{4.39}$$

$$\leq \frac{16}{\mu} \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2p} \|\mathbf{y}\|_{\widehat{\mathbf{K}}}^2 \|\mathbf{y}\|^2 \tag{4.40}$$

$$\leq C_1 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2p} \tag{4.41}$$

where $C_1 = \frac{16}{\mu} \|\mathbf{y}\|_{\widehat{\mathbf{K}}}^2 \|\mathbf{y}\|^2$. Then we have

$$\text{Var}(\mathbf{y}^\top \phi_{J_i}) \leq C_1 \sum_{i=J_{\min}}^{N} \frac{1 - \mathbb{P}(\mathcal{J} = i)}{\mathbb{P}(\mathcal{J} = i)} \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{4i} \tag{4.42}$$

$$\text{Var}(\mathcal{L}_1) = \frac{1}{l^2}\sum_{i=1}^{l}\text{Var}(\mathbf{y}^\top \phi_{J_i}) \le \frac{C_1}{l}\sum_{i=J_{\min}}^{N}\frac{1 - \mathbb{P}(\mathcal{J} = i)}{\mathbb{P}(\mathcal{J} = i)}\Big(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}\Big)^{4i} \tag{4.43}$$

For log-determinant, we need to compute the bound of $\|\mathbf{z}_i\|^2 \mathbf{e}^\top (\log T_{\mathbf{z}_i}^{J+1} - \log T_{\mathbf{z}_i}^{J})\mathbf{e}$. Authors in [96] proved that

$$\big|\|\mathbf{z}_i\|^2 \mathbf{e}(\log T_{\mathbf{z}_i}^{J})\mathbf{e} - \mathbf{z}_i^\top (\log \widehat{\mathbf{K}})\mathbf{z}_i\big| \le C_\rho\Big(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}\Big)^{2J} \tag{4.44}$$

where $\rho = \frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}$ and the constant $C_\rho$ depends on $\rho$. Then by telescoping the sum, we have

$$\big|\|\mathbf{z}_i\|^2 \mathbf{e}^\top (\log T_{\mathbf{z}_i}^{J+1} - \log T_{\mathbf{z}_i}^{J})\mathbf{e}\big|$$
$$\le \big|\|\mathbf{z}_i\|^2 \mathbf{e}(\log T_{\mathbf{z}_i}^{J})\mathbf{e} - \mathbf{z}_i^\top(\log\widehat{\mathbf{K}})\mathbf{z}_i\big| + \big|\|\mathbf{z}_i\|^2 \mathbf{e}(\log T_{\mathbf{z}_i}^{J+1})\mathbf{e} - \mathbf{z}_i^\top(\log\widehat{\mathbf{K}})\mathbf{z}_i\big|$$
$$\le C_\rho\Big(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}\Big)^{2J} + \Big(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}\Big)^{2(J+1)}$$
$$\le 2C_\rho\Big(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}\Big)^{2J}$$

Therefore,

$$\text{Var}(\psi_{\mathbf{z}_i, J}) = \sum_{i=J_{\min}}^{N}\frac{\Delta_i^2(1 - \mathbb{P}(\mathcal{J} = i))}{\mathbb{P}(\mathcal{J} = i)} \le 2C_\rho\sum_{i=J_{\min}}^{N}\frac{1 - \mathbb{P}(\mathcal{J} = i)}{\mathbb{P}(\mathcal{J} = i)}\Big(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}\Big)^{4i} \tag{4.45}$$

Since $\mathbf{z}_i's$ and $J_j's$ are all independent, we have

$$\text{Var}(\mathcal{L}_2) \le \frac{kl}{l^2 k^2}\text{Var}(\psi_{\mathbf{z}_i, J_j}) \le \frac{2C_\rho}{lk}\sum_{i=J_{\min}}^{N}\frac{1 - \mathbb{P}(\mathcal{J} = i)}{\mathbb{P}(\mathcal{J} = i)}\Big(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}\Big)^{4i} \tag{4.46}$$

Finally, the variance of LML estimated by SS-CG can be bounded by

$$\text{Var}(\mathcal{L}) = \text{Var}(\mathcal{L}_1) + \text{Var}(\mathcal{L}_2) \le \Big(\frac{2C_\rho}{lk} + \frac{C_1}{l}\Big)\sum_{i=J_{\min}}^{N}\frac{1 - \mathbb{P}(\mathcal{J} = i)}{\mathbb{P}(\mathcal{J} = i)}\Big(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}\Big)^{4i} \tag{4.47}$$

$$\square$$

Similarly, for the gradient estimation, we can prove.

**Theorem 4.5.3.** *The estimation of the gradient of LML derived by SS-CG is unbiased. The variance has the following bound.*

## 4.5.3 Preconditioned SS-CG

SS-CG presents a cost-effective alternative that comes with the side effect of increased variance compared RR-CG proposed [74]. In the following subsection, we will show how to reduce the variance of SS-CG by using preconditioner to make it practical.

The variance of the SS-CG estimation for both inverse quadratic term and log determinant term can be reduced by using a preconditioner as the bound indicates. If the variance is small enough, we can omit the loop for $J$ and use only one sample for estimation. Then the cost for *PredSS-CG* will be same as the CG. In this work, we will use the Adaptive Factorized Nyström Preconditioner (AFN) from [111] since it's adaptive to kernels' spectrum and its gradient can be computed efficiently by hand thus avoiding automatic differentiation. AFN has the following form factorized form

$$
\mathbf{M} = \underbrace{\begin{bmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{K}_{12}^{\top}\mathbf{L}^{-\top} & \mathbf{G}^{-1} \end{bmatrix}}_{\mathbf{U}^{\top}} \begin{bmatrix} \mathbf{L}^{\top} & \mathbf{L}^{-1}\mathbf{K}_{12} \\ \mathbf{0} & \mathbf{G}^{-\top} \end{bmatrix}.
\tag{4.48}
$$

It is similar to the Nyström preconditioner. However, it adds a sparse correction to the residual matrix to make it more accurate. Since AFN is more accurate than Randmoized Nyström preconditioner [25], under the same conditions as RAN, the preconitioned matrix is guaranteed to have a condition number smaller than 28. Then the variance is bounded by

$$
\mathrm{Var}(\mathcal{L}) \leq (\frac{2C_{\rho}}{lk} + C_1) \sum_{i=J_{\min}}^{N} \frac{1 - \mathbb{P}(\mathcal{J} = i)}{\mathbb{P}(\mathcal{J} = i)} \left(\frac{\sqrt{28} - 1}{\sqrt{28} + 1}\right)^{4i}
\tag{4.49}
$$

With a good preconidtioner, this variance can be controlled in a reasonable range. In the paper

Figure 4.3: Lanczos logdet, AFN Preconditioned SS logdet, SS logdet and true logdet for a Gaussian kernel matrix of size 3000. X-axis is the CG iteration counts for SLT. The average (P)CG iteration counts for SS-(P)CG are denoted by Javg

by Wenger et al. [102], a novel tail bound is established, revealing that the number of random vectors required for stochastic Lanczos quadrature to attain a specific accuracy is influenced by the condition number of the preconditioned matrix and the F-norm approximation error. Our work diverges from theirs in two key aspects. Firstly, the bias for both loss and gradient estimations is eliminated, thanks to our estimator design. Secondly, the variance of the sample distribution $\mathbb{P}$ can be explicitly calculated in our model. These differences allow for a clearer understanding of how the number of random vectors, the condition number, and the number of unbiased estimation runs impact variance, as detailed in inequality (4.49).

### 4.5.4 Small-Bias SS-CG

In practice, we can only run CG for a limited number of iterations due to the computational cost. Furthermore, with the aid of preconditioner, the number of iterations needed to achieve a high accuracy is usually reasonablely small. Suppose after $k$ setps of CG, the accuracy is already satisfying and can be regarded as unbiased. However, for SS-CG, we need to consider a rare not impossible in SS-CG that $J = N$ since the probability distribution is supported on while range from $J_{\min}$ to $N$. In this case, the SS-CG will be equivalent to solve the equation exactly. This is of coures unacceptable. In order to avoid this suitation, we propose a small-bias CG (SB-CG). Instead of supporting on the whole range, we only support on a truncated range $[J_{\min}, J_{\max}]$. The $J_{\max}$ will be the maximum number of iterations. Then the

---

**Algorithm 5** Preconditioned SS-CG for GP Inference and Hyperparameter Optimization

---

1: **Input:** Training data $(\mathbf{X}, \mathbf{y})$, kernel function $k(\cdot, \cdot, \boldsymbol{\theta})$ with hyperparameters $\boldsymbol{\theta}$, regularization parameter $\mu$, convergence tolerance $\epsilon$, maximum iterations $T$

2: **Output:** Optimized hyperparameters $\boldsymbol{\theta}^*$, log marginal likelihood $\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}^*)$, and posterior predictive distribution parameters

3: **Hyperparameter Optimization:**

4: Initialize hyperparameters $\boldsymbol{\theta}^{(0)}$

5: **FOR** $t = 0$ to $T - 1$ **DO**

6:     Construct kernel matrix $\mathbf{K}^{(t)}$ using $k(\cdot, \cdot, \boldsymbol{\theta}^{(t)})$

7:     Construct AFN preconditioner $\mathbf{M}^{(t)}$ using Algorithm 4

8:     Sample $J^{(t)} \sim P_J$ where $P_J$ is the probability distribution over iteration counts

9:     Compute $\alpha^{(t)} = (\mathbf{K}^{(t)} + \mu \mathbf{I})^{-1} \mathbf{y}$ using preconditioned CG with $\mathbf{M}^{(t)}$ as preconditioner, running for $J^{(t)}$ iterations

10:     Estimate log determinant $\log|\mathbf{K}^{(t)} + \mu \mathbf{I}|$ using preconditioned SS-CG with $\mathbf{M}^{(t)}$ as preconditioner

11:     Compute log marginal likelihood: $\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}^{(t)}) = -\frac{1}{2}\mathbf{y}^\top \alpha^{(t)} - \frac{1}{2}\log|\mathbf{K}^{(t)} + \mu \mathbf{I}| - \frac{n}{2}\log(2\pi)$

12:     Estimate gradients of log marginal likelihood w.r.t. $\boldsymbol{\theta}$ using preconditioned SS-CG: $\frac{\partial \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}^{(t)})}{\partial \theta_i} = \frac{1}{2}\alpha^{(t)\top} \frac{\partial \mathbf{K}^{(t)}}{\partial \theta_i} \alpha^{(t)} - \frac{1}{2}\text{tr}\left((\mathbf{K}^{(t)} + \mu \mathbf{I})^{-1} \frac{\partial \mathbf{K}^{(t)}}{\partial \theta_i}\right)$

13:     Update hyperparameters: $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \eta^{(t)} \nabla_{\boldsymbol{\theta}} \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}^{(t)})$

14:     **IF** $\|\boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^{(t)}\| < \epsilon$ **THEN**

15:         **break**

16:     **END IF**

17: **END FOR**

18: Set $\boldsymbol{\theta}^* = \boldsymbol{\theta}^{(T)}$

19: **GP Inference with Optimized Hyperparameters:**

20: Construct final kernel matrix $\mathbf{K}^*$ using $k(\cdot, \cdot, \boldsymbol{\theta}^*)$

21: Construct final AFN preconditioner $\mathbf{M}^*$ using Algorithm 4

22: Compute $\alpha^* = (\mathbf{K}^* + \mu \mathbf{I})^{-1} \mathbf{y}$ using preconditioned CG with $\mathbf{M}^*$ as preconditioner

23: **For Prediction:** Given test points $\mathbf{X}_*$, compute:

24: $\mathbf{K}_* = k(\mathbf{X}, \mathbf{X}_*, \boldsymbol{\theta}^*)$ (cross-covariances)

25: $\mathbf{K}_{**} = k(\mathbf{X}_*, \mathbf{X}_*, \boldsymbol{\theta}^*)$ (test covariances)

26: Mean predictions: $\mathbb{E}[f_*] = \mathbf{K}_*^\top \alpha^*$

27: Variance predictions: $\text{Var}[f_*] = \mathbf{K}_{**} - \mathbf{K}_*^\top (\mathbf{K}^* + \mu \mathbf{I})^{-1} \mathbf{K}_*$

28: **return** $\boldsymbol{\theta}^*$, $\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}^*)$, $\mathbb{E}[f_*]$, $\text{Var}[f_*]$

---

SB-CG estimator can be written as

$$\tilde{\Phi} = \sum_{i=1}^{J_{\max}} \frac{\Delta_i}{\mathbb{P}(\mathcal{J} = i)} \mathbb{I}(J = i) \tag{4.50}$$

will be no longer unbiased. However, we can still prove that it's bias is small.

**Theorem 4.5.4.** *Denote the Small-Biased SS estimator defined in Equation 4.50 by $\tilde{\Phi}$ and the true value by $\bar{\Phi} = \sum_{i-1}^{H} \Delta_i$. If $\sum_{i=J_{\max}}^{H} \Delta_i \leq \epsilon$, then the bias of the estimator is bounded by $\epsilon$*

$$|\mathbb{E}(\tilde{\Phi}) - \bar{\Phi}| \leq \epsilon. \tag{4.51}$$

$$Var(\tilde{\Phi}) = \sum_{i=J_{\min}}^{J_{\max}} \frac{\Delta_i^2(1 - \mathbb{P}(\mathcal{J} = i))}{\mathbb{P}(\mathcal{J} = i)}. \tag{4.52}$$

*Proof.* Now the $\mathcal{J}$ is supported on $[J_{\min}, J_{\max}]$. Then we have

$$\mathbb{E}_{\mathcal{J}}[\tilde{\Phi}] = \sum_{i=1}^{J_{\max}} \frac{\Delta_i}{\mathbb{P}(\mathcal{J} = i)} \mathbb{E}(\mathbb{I}(J = i)) = \sum_{i=1}^{J_{\max}} \frac{\Delta_i}{\mathbb{P}(\mathcal{J} = i)} \mathbb{P}(\mathcal{J} = i) = \sum_{i=1}^{J_{\max}} \Delta_i = \bar{\Phi} - \sum_{J_{\max}}^{N} \Delta_i. \tag{4.53}$$

Since $\sum_{i=J_{\max}}^{N} \Delta_i < \epsilon$, we have

$$|\mathbb{E}(\tilde{\Phi}) - \bar{\Phi}| \leq \epsilon. \tag{4.54}$$

The variance can be computed similar to the original SS estimator. $\square$

**Theorem 4.5.5.** *When the following probability distribution $P(\mathcal{J} = i) = \frac{e^{-\lambda i}}{\sum_{i=J_{\min}}^{J_{\max}} e^{-\lambda i}}$ is used, the expected number of iterations of the SB-CG estimator has the following upper bound.*

$$\mathbb{E}[\mathcal{J}] \leq \frac{\frac{1}{\lambda}\left(J_{\min} - J_{\max}e^{-\lambda(J_{\max}-J_{\min})}\right) + \frac{1}{\lambda^2}\left(1 - e^{-\lambda(J_{\max}-J_{\min})}\right)}{\left(1 - e^{-\lambda(J_{\max}-J_{\min}-1)}\right)}$$

*Proof.*

$$\mathbb{E}[\mathcal{J}] = \sum_{i=J_{\min}}^{J_{\max}} iP(\mathcal{J}=i) = \sum_{i=J_{\min}}^{J_{\max}} i\frac{e^{-\lambda i}}{\sum_{i=J_{\min}}^{J_{\max}} e^{-\lambda i}}$$

Firstly, we have

$$\sum_{i=J_{\min}}^{J_{\max}} e^{-\lambda i} = \frac{e^{-\lambda J_{\min}}\left(1 - e^{\lambda(J_{\min}-J_{\max}-1)}\right)}{1 - \frac{1}{e^{\lambda}}} \geq e^{-\lambda J_{\min}}\left(1 - e^{\lambda(J_{\min}-J_{\max}-1)}\right)$$

Then we have

$$\sum_{i=J_{\min}}^{J_{\max}} e^{-\lambda i}i \leq \int_{J_{\min}}^{J_{\max}} e^{-\lambda x}xdx = \frac{1}{\lambda}J_{\min}e^{-\lambda J_{\min}} - \frac{1}{\lambda}J_{\max}e^{-\lambda J_{\max}} + \frac{1}{\lambda^2}(e^{-\lambda J_{\min}} - e^{-\lambda J_{\max}})$$

Combining the above two inequalities, we obtain

$$\mathbb{E}[\mathcal{J}] \leq \frac{\frac{1}{\lambda}\left(J_{\min} - J_{\max}e^{-\lambda(J_{\max}-J_{\min})}\right) + \frac{1}{\lambda^2}\left(1 - e^{-\lambda(J_{\max}-J_{\min})}\right)}{\left(1 - e^{-\lambda(J_{\max}-J_{\min}-1)}\right)}$$

□

**Corollary 1.** *If we set $J_{\min} = 5$, $J_{\max} = 5000$ and $\lambda = 0.1$ as suggested [74]. The bounds on $\mathbb{E}[\mathcal{J}]$ can be simplified to*

$$\mathbb{E}[\mathcal{J}] \leq 151$$

*Proof.* Since $e^{\lambda(J_{\min}-J_{\max})} \sim 0$, we have $\mathbb{E}[\mathcal{J}] \leq \frac{J_{\min}}{\lambda} + \frac{1}{\lambda^2} + 1 = 151$. □

## 4.6 Theoretical Analysis of Kernel Preconditioning

The effectiveness of our AFN preconditioner and unbiased estimation techniques is supported by rigorous theoretical analysis. In this section, we present key theoretical results that provide guarantees on the performance of our methods. These results not only justify our algorithmic

choices but also offer insights into the fundamental properties of kernel matrices and their approximations.

We begin by analyzing the geometric properties of landmark point selection, which is crucial for the construction of effective preconditioners. We then establish bounds on the approximation error of the Nyström method, which forms the basis of our AFN preconditioner.

### 4.6.1   Interplay between Fill and Separation Distance

In this section, we study the relationship between the fill distance $h_{X_k}$ and separation distance $q_{X_k}$ for landmark point selection. These metrics are crucial for understanding the quality of preconditioners based on landmark points.

**Theorem 4.6.1** (Fill-Separation Bounds). *Suppose all the data points are inside a unit ball* $\Omega$ *in* $\mathbb{R}^d$. *Then for an arbitrary subset* $X_k = \{\mathbf{x}_{i_1}, \ldots, \mathbf{x}_{i_k}\}$ *of* $X$, *the following bounds hold for* $h_{X_k}$ *and* $q_{X_k}$:

$$h_{X_k} \geq C_\Omega k^{-1/d} \quad and \quad q_{X_k} \leq C'_\Omega k^{-1/d}, \tag{4.55}$$

*where* $C_\Omega$ *and* $C'_\Omega$ *are two constants only depending on* $\Omega$.

The bounds in Theorem 4.6.1 indicate that the minimal fill distance $h_{X_k}$ cannot be smaller than $C_\Omega k^{-1/d}$ while the maximal separation distance $q_{X_k}$ cannot be greater than $C'_\Omega k^{-1/d}$. This establishes a fundamental relationship: $\frac{C_\Omega}{C'_\Omega} q_{X_k} \leq h_{X_k}$.

**Theorem 4.6.2** (Farthest Point Sampling Optimality). *Assume the data points are on a bounded domain* $\Omega$ *that satisfies the interior cone condition. If* $h_{X_k} \leq C q_{X_k}$ *for a constant* $C$, *then:*

$$C_\Omega k^{-1/d} \leq h_{X_k} \leq C \times C'_\Omega k^{-1/d}, \quad \frac{C_\Omega}{C} k^{-1/d} \leq q_{X_k} \leq C'_\Omega k^{-1/d}. \tag{4.56}$$

Theorem 4.6.2 shows that if a sampling scheme can select a subset $X_k$ with $h_{X_k} \leq C q_{X_k}$,

then $q_{X_k}$ has the same order as the maximal separation distance achievable by any subset with $k$ points.

**Theorem 4.6.3** (FPS Optimality Guarantees)**.** *Suppose the minimal fill distance of a subset with $k$ points is achieved with $X_k^*$ and the maximal separation distance of a subset with $k$ points is achieved with $X_{k*}$. Then the set $X_k$ sampled by Farthest Point Sampling (FPS) satisfies:*

$$h_{X_k} \leq q_{X_k} \quad and \quad q_{X_k} \geq \frac{1}{2} q_{X_{k*}} \quad and \quad h_{X_k} \leq 2 h_{X_k^*}. \tag{4.57}$$

Theorem 4.6.3 establishes that FPS generates landmark points with near-optimal properties, making it an excellent choice for selecting landmark points in preconditioner construction.

## 4.6.2 Nyström Approximation Error Analysis

The quality of Nyström-based preconditioners depends on the approximation error of the Nyström method. We can establish bounds on this error in terms of the fill distance.

**Theorem 4.6.4** (Nyström Error Bound)**.** *The Nyström approximation $\mathbf{K}_{nys} = \mathbf{K}_{X,X_k} \mathbf{K}_{X_k,X_k}^{-1} \mathbf{K}_{X_k,X}$ to $\mathbf{K}$ using the landmark points $X_k = \{\mathbf{x}_{k_i}\}_{i=1}^k$ has the following error estimate:*

$$\|\mathbf{K} - \mathbf{K}_{nys}\| < \sqrt{n\|\mathbf{K}\|} C' \exp(-C''/h_{X_k}), \tag{4.58}$$

*where $C'$ and $C''$ are constants independent of $X_k$.*

Theorem 4.6.4 implies that landmark points $X_k$ with a smaller fill distance yield a more accurate Nyström approximation. This theoretical result supports the use of FPS for landmark point selection, as it minimizes the fill distance.

## 4.7 Implementation Details

The theoretical guarantees established in the previous section inform our implementation choices. While the theory provides valuable insights into the expected performance and properties of our methods, translating these insights into an efficient implementation requires careful consideration of practical computing constraints. In this section, we detail the key implementation aspects of our approach, focusing on the computation of derivatives for various components, which is essential for gradient-based hyperparameter optimization in Gaussian Process models.

### 4.7.1 Derivative of Gaussian Kernel matrix

$$\frac{\partial K(x,y)}{\partial l} = 2\frac{\|x-y\|e^{-\frac{\|x-y\|}{l^2}}}{l^3}$$
$$\frac{\partial K(x,y)}{\partial \mu} = \mathbf{I}$$

### 4.7.2 Derivative of Inverse

For a matrix $\mathbf{K}$, the gradient of its inverse matrix

$$\frac{\partial \mathbf{K}^{-1}}{\partial \theta} = -\mathbf{K}^{-1}\frac{\partial K}{\partial \theta}\mathbf{K}^{-1}.$$

### 4.7.3 Derivative of Cholesky Factorization

For a Cholesky factorization of a PSD matrix $\mathbf{K} = \mathbf{L}\mathbf{L}^\top$, the gradient of Cholesky factor can be computed as

$$\frac{\partial \mathbf{L}}{\partial \theta} = \mathbf{L}\Phi(\mathbf{L}^{-1}\frac{\partial \mathbf{K}}{\partial \theta}\mathbf{L}^{-\top}) \tag{4.59}$$

where $\Phi$ is defined as

$$\Phi(K) = \begin{cases} K_{ij}, & i > j \\ \frac{1}{2}K_{ij}, & i = j \\ 0, & i < j \end{cases}$$

## 4.7.4 Derivative of Nyström

Given $\mathbf{K}$, the classical Nyström approximation is $\mathbf{K}_{nys} = \mathbf{K}_{\cdot 1}\mathbf{K}_{11}^{-1}\mathbf{K}_{1\cdot}$ given all indices corresponding to landmark points are rearranged to the $(1,1)$ block of the matrix. Then the gradient of $\mathbf{K}_{nys}$ can be computed

$$\frac{\partial \mathbf{K}_{nys}}{\partial \theta} = \frac{\partial \mathbf{K}_{\cdot 1}}{\partial \theta}\mathbf{K}_{11}^{-1}\mathbf{K}_{1\cdot} + \mathbf{K}_{\cdot 1}\frac{\partial \mathbf{K}_{11}^{-1}}{\partial \theta}\mathbf{K}_{1\cdot} + \mathbf{K}_{\cdot 1}\mathbf{K}_{11}^{-1}\frac{\partial \mathbf{K}_{1\cdot}}{\partial \theta} \qquad (4.60)$$

## 4.7.5 Derivative of FSAI

Given sparsity pattern $\mathbf{P}$, the FSAI $\mathbf{G}\mathbf{G}^T \approx \mathbf{K}^{-1}$ is computed as

$$\mathbf{G}(:,i) = \frac{\mathbf{K}(\mathbf{P}(:,i),\mathbf{P}(:,i))^{-1}\mathbf{e}}{\sqrt{\mathbf{e}^T\mathbf{K}(\mathbf{P}(:,i),\mathbf{P}(:,i))^{-1}\mathbf{e}}}, \qquad (4.61)$$

where $\mathbf{G}$ is an upper triangular matrix. If $\mathbf{K}$ is a matrix function $\mathbf{K}(x)$, by using matrix calculus, we can compute the derivative of $\mathbf{G}$ with respect to value $x$ as

$$\frac{\partial \mathbf{G}(:,i)}{\partial x} = \frac{1}{\sqrt{\mathbf{e}^T\mathbf{K}(\mathbf{P}(:,i),\mathbf{P}(:,i))^{-1}\mathbf{e}}}\frac{\partial \mathbf{K}(\mathbf{P}(:,i),\mathbf{P}(:,i))^{-1}}{\partial x} \qquad (4.62)$$

$$+ \mathbf{K}(\mathbf{P}(:,i),\mathbf{P}(:,i))^{-1}\mathbf{e}\frac{\partial (\mathbf{e}^T\mathbf{K}(\mathbf{P}(:,i),\mathbf{P}(:,i))^{-1}\mathbf{e})^{-1/2}}{\partial x} \qquad (4.63)$$

$$= \frac{1}{\sqrt{\mathbf{e}^T\mathbf{K}(\mathbf{P}(:,i),\mathbf{P}(:,i))^{-1}\mathbf{e}}}\frac{\partial \mathbf{K}(\mathbf{P}(:,i),\mathbf{P}(:,i))^{-1}}{\partial x}\mathbf{e} \qquad (4.64)$$

$$- \frac{\mathbf{K}(\mathbf{P}(:,i),\mathbf{P}(:,i))^{-1}\mathbf{e}}{2[\mathbf{e}^T\mathbf{K}(\mathbf{P}(:,i),\mathbf{P}(:,i))^{-1}\mathbf{e}]^{3/2}}\mathbf{e}^T\frac{\partial \mathbf{K}(\mathbf{P}(:,i),\mathbf{P}(:,i))^{-1}}{\partial x}\mathbf{e}, \qquad (4.65)$$

where

$$\frac{\partial \mathbf{K}(\mathbf{P}(:,i), \mathbf{P}(:,i))^{-1}}{\partial x} = -\mathbf{K}(\mathbf{P}(:,i), \mathbf{P}(:,i))^{-1} \frac{\partial \mathbf{K}(\mathbf{P}(:,i), \mathbf{P}(:,i))}{\partial x} \mathbf{K}(\mathbf{P}(:,i), \mathbf{P}(:,i))^{-1}. \quad (4.66)$$

However, in backward propagation, we are using a reverse mode autodiff so we need to compute the derivative of $G$ w.r.t. the whole kernel matrix $K$. According to the matrix calculus, we know twidehat if $g(Y) = g(f(K))$, then $\frac{\partial g}{\partial K_{ij}} = \text{Tr}((\frac{\partial g(Y)}{\partial Y})^\top \frac{f(K)}{K_{ij}})$. Suppose f is the function computing FSAI and g is a function mapping result of FSAI $G$ to a scalar such log determinant or norm. Then in reverse mode autodiff, we have $\frac{\partial g(Y)}{\partial Y}$ coming from the previous step and is denoted by gradoutput in backward propagation. Therefore, we only need to compute $\frac{G}{K_{ij}}$ which can be computed using the above formula:

$$\frac{\partial \mathbf{G}(:,i)}{\partial K_{ij}} = \frac{1}{\sqrt{\mathbf{e}^T \mathbf{K}(\mathbf{P}(:,i), \mathbf{P}(:,i))^{-1} \mathbf{e}}} \frac{\partial \mathbf{K}(\mathbf{P}(:,i), \mathbf{P}(:,i))^{-1}}{\partial K_{ij}} \mathbf{e} \quad (4.67)$$

$$- \frac{\mathbf{K}(\mathbf{P}(:,i), \mathbf{P}(:,i))^{-1} \mathbf{e}}{2[\mathbf{e}^T \mathbf{K}(\mathbf{P}(:,i), \mathbf{P}(:,i))^{-1} \mathbf{e}]^{3/2}} \mathbf{e}^T \frac{\partial \mathbf{K}(\mathbf{P}(:,i), \mathbf{P}(:,i))^{-1}}{\partial K_{ij}} \mathbf{e}, \quad (4.68)$$

where

$$\frac{\partial \mathbf{K}(\mathbf{P}(:,i), \mathbf{P}(:,i))^{-1}}{\partial K_{ij}} = -\mathbf{K}(\mathbf{P}(:,i), \mathbf{P}(:,i))^{-1} \frac{\partial \mathbf{K}(\mathbf{P}(:,i), \mathbf{P}(:,i))}{\partial K_{ij}} \mathbf{K}(\mathbf{P}(:,i), \mathbf{P}(:,i))^{-1} \quad (4.69)$$

$$= \mathbf{K}(\mathbf{P}(:,i), \mathbf{P}(:,i))^{-1} E_{index(i)index(j)} \mathbf{K}(\mathbf{P}(:,i), \mathbf{P}(:,i))^{-1} \quad (4.70)$$

### 4.7.6   Derivative of Schur Complement

In the construction of AFN preconditioner, we apply FSAI to $\hat{\mathbf{K}}_{\text{Schur}} = \mathbf{K}_{22} - \mathbf{K}_{12}^\top \mathbf{K}_{11}^{-1} \mathbf{K}_{12}$, so we also need the derivative of $\hat{\mathbf{K}}_{\text{Schur}}$ which can be computed as follows

$$\frac{\partial \hat{\mathbf{K}}_{\text{Schur}}}{\partial \theta} = \frac{\partial \mathbf{K}_{22}}{\partial \theta} - \frac{\partial \mathbf{K}_{12}^\top}{\partial \theta} \mathbf{K}_{11}^{-1} \mathbf{K}_{12} + \mathbf{K}_{12}^\top \mathbf{K}_{11}^{-1} \frac{\partial \mathbf{K}_{11}}{\partial \theta} \mathbf{K}_{11}^{-1} \mathbf{K}_{12} - \mathbf{K}_{12}^\top \mathbf{K}_{11}^{-1} \frac{\partial \mathbf{K}_{12}}{\partial \theta}$$

## 4.7.7 Derivative of AFN

Now the preconditioner has the following factorized form

$$
\mathbf{M} = \underbrace{\begin{bmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{K}_{12}^\top \mathbf{L}^{-\top} & \mathbf{G}^{-1} \end{bmatrix}}_{\mathbf{U}^\top} \begin{bmatrix} \mathbf{L}^\top & \mathbf{L}^{-1}\mathbf{K}_{12} \\ \mathbf{0} & \mathbf{G}^{-\top} \end{bmatrix}. \tag{4.71}
$$

We can compute the derivative directly as follows

$$
\frac{\partial \mathbf{M}}{\partial \theta} = \frac{\partial \mathbf{U}^\top}{\partial \theta} \mathbf{U} + \mathbf{U}^\top \frac{\partial \mathbf{U}}{\partial \theta} \tag{4.72}
$$

Then

$$
\mathrm{tr}(\mathbf{M}^{-1}\frac{\partial \mathbf{M}}{\partial \theta}) = \mathrm{tr}(\mathbf{U}^{-1}\mathbf{U}^{-\top}\frac{\partial \mathbf{M}}{\partial \theta}) = \mathrm{tr}(\mathbf{U}^{-\top}\frac{\partial \mathbf{M}}{\partial \theta}\mathbf{U}^{-1})
$$

$$
= \mathrm{tr}(\mathbf{U}^{-\top}\frac{\partial \mathbf{U}^\top}{\partial \theta} + \frac{\partial \mathbf{U}}{\partial \theta}\mathbf{U}^{-1})
$$

Notice

$$
\mathbf{U}^{-\top} = \begin{bmatrix} \mathbf{L}^{-1} & \mathbf{0} \\ -\mathbf{G}\mathbf{K}_{12}^\top \mathbf{L}^{-\top}\mathbf{L}^{-1} & \mathbf{G} \end{bmatrix}
$$

$$
\mathbf{U}^{-1} = \begin{bmatrix} \mathbf{L}^{-\top} & -\mathbf{L}^{-\top}\mathbf{L}^{-1}\mathbf{K}_{12}\mathbf{G}^\top \\ \mathbf{0} & \mathbf{G}^\top \end{bmatrix}
$$

$$
\mathbf{U}^{-\top}\hat{\mathbf{K}}\mathbf{U}^{-1} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{G}(\hat{\mathbf{K}}_{22} - \mathbf{K}_{12}^\top(\hat{\mathbf{K}}_{11} + \mu\mathbf{I})^{-1}\mathbf{K}_{12})\mathbf{G}^\top \end{bmatrix}
$$

$$\frac{\partial \mathbf{U}^\top}{\partial \theta} = \begin{bmatrix} \frac{\partial \mathbf{L}}{\partial \theta} & \mathbf{0} \\ \frac{\partial \mathbf{K}_{12}^\top}{\partial \theta}\mathbf{L}^{-\top} - \mathbf{K}_{12}^\top \mathbf{L}^{-\top}\frac{\partial \mathbf{L}^\top}{\partial \theta}\mathbf{L}^{-\top} & -\mathbf{G}^{-1}\frac{\partial \mathbf{G}}{\partial \theta}\mathbf{G}^{-1} \end{bmatrix}$$

$$\frac{\partial \mathbf{U}}{\partial \theta} = \begin{bmatrix} \frac{\partial \mathbf{L}^\top}{\partial \theta} & \mathbf{L}^{-1}\frac{\partial \mathbf{K}_{12}}{\partial \theta} - \mathbf{L}^{-1}\frac{\partial \mathbf{L}}{\partial \theta}\mathbf{L}^{-1}\mathbf{K}_{12} \\ \mathbf{0} & -\mathbf{G}^{-\top}\frac{\partial \mathbf{G}^\top}{\partial \theta}\mathbf{G}^{-\top} \end{bmatrix}$$

Then we can compute $\text{tr}(\mathbf{M}^{-1}\frac{\partial \mathbf{M}}{\partial \theta})$ using the fact that the diagonals of the inverse of a lower (upper) triangular matrix are just the inverse of its diagonals. Therefore

$$\text{tr}(\frac{\partial \mathbf{U}}{\partial \theta}\mathbf{U}^{-1}) = \sum_{i=1}^{n}(\frac{\partial \mathbf{U}}{\partial \theta})_{ii}\frac{1}{\mathbf{U}_{ii}}$$

Besides this term, we also need to evaluate $\mathbf{M}^{-1}\frac{\partial \mathbf{M}}{\partial \theta}$ over the probe vectors $z$, which can be done by invoking (sparse) triangular solves. Note that for any probe vectors $z$,

$$z^\top \mathbf{U}^{-\top}\frac{\partial \mathbf{U}^\top}{\partial \theta}z = z^\top \frac{\partial \mathbf{U}}{\partial \theta}\mathbf{U}^{-1}z,$$

so we only need to evaluate one of them. We first solve for

$$\begin{bmatrix} \mathbf{L}^\top & \mathbf{L}^{-1}\mathbf{K}_{12} \\ \mathbf{0} & \mathbf{G}^{-\top} \end{bmatrix}\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{bmatrix}$$

We can easily get the solution

$$\mathbf{x}_2 = \mathbf{G}^\top \mathbf{z}_2$$

$$\mathbf{x}_1 = \mathbf{L}^{-\top}(\mathbf{z}_1 - \mathbf{L}^{-1}\mathbf{K}_{12}\mathbf{G}^\top \mathbf{z}_2)$$

But for Stochastic Lanczos Estimation (SLE) we need to apply the preconditioner in the

factorized form to have a symmetric matrix, thus we also need to solve

$$\begin{bmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{K}_{12}^\top \mathbf{L}^{-\top} & \mathbf{G}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{bmatrix}$$

We can easily get the solution

$$\mathbf{x}_1 = \mathbf{L}^{-1}\mathbf{z}_1$$

$$\mathbf{x}_2 = \mathbf{G}(\mathbf{z}_2 - \mathbf{K}_{12}^\top \mathbf{L}^{-\top}\mathbf{L}^{-1}\mathbf{z}_1)$$

Then we just need to multiply $\mathbf{x}$ by $\frac{\partial \mathbf{U}}{\partial \theta}$

$$\frac{\partial \mathbf{U}}{\partial \theta} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \frac{\partial \mathbf{L}^\top}{\partial \theta}\mathbf{x}_1 + (\mathbf{L}^{-1}\frac{\partial \mathbf{K}_{12}}{\partial \theta} - \mathbf{L}^{-1}\frac{\partial \mathbf{L}}{\partial \theta}\mathbf{L}^{-1}\mathbf{K}_{12})\mathbf{x}_2 \\ -\mathbf{G}^{-\top}\frac{\partial \mathbf{G}^\top}{\partial \theta}\mathbf{G}^{-\top}\mathbf{x}_2 \end{bmatrix}$$

Multiplying the factors gives M another form

$$\mathbf{M} = \begin{bmatrix} \mathbf{K}_{11} + \mu\mathbf{I} & \mathbf{K}_{12} \\ \mathbf{K}_{12}^\top & (\mathbf{G}^\top\mathbf{G})^{-1} + \mathbf{K}_{12}^\top(\mathbf{K}_{11} + \mu\mathbf{I})^{-1}\mathbf{K}_{12} \end{bmatrix} \tag{4.73}$$

$$= \mathbf{K}_{nys} + \mu\mathbf{I} + \underbrace{\begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & (\mathbf{G}^\top\mathbf{G})^{-1} + \mathbf{K}_{12}^\top\left((\mathbf{K}_{11} + \mu\mathbf{I})^{-1} - (\mathbf{K}_{11})^{-1}\right)\mathbf{K}_{12} - \mu\mathbf{I} \end{bmatrix}}_{\text{Correction term}}, \tag{4.74}$$

The gradient of M can be computed as follows

$$\frac{\partial \mathbf{M}}{\partial \theta} = \begin{bmatrix} \frac{\partial \mathbf{K}_{11} + \mu\mathbf{I}}{\partial \theta} & \frac{\partial \mathbf{K}_{12}}{\partial \theta} \\ \frac{\partial \mathbf{K}_{12}^\top}{\partial \theta} & -(\mathbf{G}^\top\mathbf{G})^{-1}\frac{\partial \mathbf{G}^\top\mathbf{G}}{\partial \theta}(\mathbf{G}^\top\mathbf{G})^{-1} + \frac{\partial \mathbf{K}_{12}^\top((\mathbf{K}_{11} + \mu\mathbf{I})^{-1} - \mathbf{K}_{11}^{-1})\mathbf{K}_{12}}{\partial \theta} \end{bmatrix} \tag{4.75}$$

Here we just need to replace the $\mathbf{K}$ in the computation of the gradient of FSAI with $\hat{\mathbf{K}}_{\text{Schur}}$.

# 4.8 Numerical Experiments

To validate our theoretical results and assess the practical performance of our proposed methods, we conducted a comprehensive suite of experiments. These experiments were designed with three primary goals: (1) to evaluate the computational efficiency of the AFN preconditioner compared to existing approaches, (2) to verify the unbiasedness and variance properties of our *PredSS-CG* estimator, and (3) to demonstrate the practical utility of our methods for hyperparameter optimization in realistic Gaussian Process regression scenarios.

Our experimental evaluation covers both synthetic datasets, where we can systematically control the difficulty of the problem, and real-world datasets that present the complex challenges encountered in practical applications. We begin by describing our experimental setup, followed by detailed results for different kernel functions and analysis of key factors affecting performance.

## 4.8.1 Experimental Setup

To evaluate the effectiveness of our proposed methods, we conducted extensive experiments on both synthetic and real-world datasets. All experiments were run on an Ubuntu 20.04.4 LTS machine equipped with 755 GB of system memory and a 24-core 3.0 GHz Intel Xeon Gold 6248R CPU. We implemented our methods in C with OpenMP for shared memory parallelism and used the OpenBLAS library for basic matrix operations. For large-scale 3D datasets, we utilized H2Pack to provide linear complexity matrix-vector multiplications with a relative error threshold of $10^{-8}$.

For all experiments, we set the stopping tolerance for the relative residual norm to $10^{-4}$ and randomly generated right-hand side vectors with entries from the uniform distribution $[-0.5, 0.5]$. Each experiment was repeated three times, and we report the average results.

### 4.8.2   Synthetic 3D Dataset Experiments

Our synthetic dataset consists of $n = 1.6 \times 10^5$ random points sampled uniformly from inside a 3D cube with edge length $\sqrt[3]{n}$. We solved regularized linear systems associated with both Gaussian and Matérn-3/2 kernels, with a regularization parameter $\mu = 0.0001$.

We compared our AFN preconditioner against several baselines:

- Unpreconditioned Conjugate Gradient (CG)

- Factorized Sparse Approximate Inverse (FSAI) preconditioner

- Randomized Nyström (RAN) preconditioner with 3000 randomly selected landmark points

For the FSAI preconditioner, we used 400 nearest neighbors as the sparsity pattern. For the AFN preconditioner, we used 100 nearest neighbors as the sparsity pattern for the FSAI component.

### 4.8.3   Results on Gaussian Kernel

For the Gaussian kernel, we observed that the performance of different methods varies significantly with the length-scale parameter. Table 4.1 shows the iteration counts, setup times, and solve times for different methods across a range of length-scales.

For large length-scales (e.g., $l^2 = 1000$), all preconditioners perform well, with AFN and RAN requiring only 3 iterations to converge. However, as the length-scale decreases to the middle range ($l^2 = 65$ to $l^2 = 25$), unpreconditioned CG and RAN struggle to converge within 500 iterations, while AFN maintains reasonable iteration counts (35-62). For very small length-scales ($l^2 = 0.1$), the problem becomes easier, and all methods converge quickly.

The setup time for AFN remains consistent across all length-scales, while the setup time for RAN increases as the length-scale decreases. In terms of solve time, AFN outperforms RAN across all middle length-scales, often by a factor of 2-3x.

Figure 4.4: Logdet estimation for a full rank matrix over kin40k datasets.

### 4.8.4   Results on Matérn-3/2 Kernel

For the Matérn-3/2 kernel, we observed similar trends but with some notable differences. Table 4.2 shows the results for this kernel.

AFN demonstrates remarkable robustness across all length-scales, requiring only 3-9 iterations to converge, regardless of the length-scale. In contrast, unpreconditioned CG requires hundreds of iterations for large and small length-scales and fails to converge for middle length-scales. RAN performs well for very large and very small length-scales but struggles with middle length-scales.

The setup time for AFN is higher than for RAN for large length-scales but becomes comparable for middle and small length-scales. In terms of solve time, AFN significantly outperforms RAN for all middle length-scales, often by an order of magnitude.

Table 4.1: Results for Gaussian kernel with $\mu = 0.0001$ on synthetic 3D dataset with $n = 1.6 \times 10^5$ points. "-" indicates failure to converge within 500 iterations.

| Length-scale ($l^2$) | 1000 | 65 | 50 | 40 | 30 | 25 | 0.1 |
|---|---|---|---|---|---|---|---|
| **Iteration Counts** | | | | | | | |
| CG | 44.00 | - | - | - | - | - | 1.00 |
| AFN | 3.00 | 35.00 | 40.00 | 46.00 | 57.00 | 62.00 | 1.00 |
| RAN | 3.00 | 72.67 | 199.33 | 409.33 | - | - | - |
| FSAI | - | - | - | - | - | - | 1.00 |
| **Setup Time (s)** | | | | | | | |
| AFN | 3.19 | 38.97 | 39.73 | 40.76 | 40.12 | 40.59 | 40.37 |
| RAN | 27.28 | 27.59 | 29.05 | 31.18 | 33.64 | 33.97 | 35.07 |
| FSAI | 10.00 | 9.91 | 9.72 | 10.14 | 10.01 | 9.84 | 13.22 |
| **Solve Time (s)** | | | | | | | |
| CG | 9.72 | - | - | - | - | - | 1.75 |
| AFN | 0.43 | 12.49 | 15.82 | 20.15 | 27.26 | 29.10 | 1.91 |
| RAN | 0.81 | 23.29 | 72.20 | 138.88 | - | - | - |
| FSAI | - | - | - | - | - | - | 1.27 |

## 4.8.5   Analysis of Landmark Point Selection

The effectiveness of the AFN preconditioner depends critically on the selection of landmark points. We compared Farthest Point Sampling (FPS) against random sampling for selecting landmark points and found that FPS consistently produces better results.

FPS generates landmark points with smaller fill distances and larger separation distances compared to random sampling. This leads to better conditioning of the $\mathbf{K}_{11}$ block and more effective preconditioning. The standard deviation of these geometric measures is also smaller with FPS, indicating more robust preconditioner performance.

## 4.8.6   Unbiased Estimation Results

We also evaluated our unbiased stochastic estimation technique for the LML and its gradient. Figure 4.3 shows the comparison between Lanczos logdet, AFN Preconditioned SS logdet, SS logdet, and the true logdet for a Gaussian kernel matrix of size 3000.

The results demonstrate that our preconditioned SS-CG method achieves accurate estimation with significantly fewer CG iterations compared to standard Lanczos quadrature. The

Table 4.2: Results for Matérn-3/2 kernel with $\mu = 0.0001$ on synthetic 3D dataset with $n = 1.6 \times 10^5$ points. "-" indicates failure to converge within 500 iterations.

| Inverse Length $(1/l)$ | 1.0 | 0.065 | 0.050 | 0.040 | 0.030 | 0.025 | 0.001 |
|---|---|---|---|---|---|---|---|
| **Iteration Counts** | | | | | | | |
| CG | 293.67 | - | - | - | - | - | 292.67 |
| AFN | 3.00 | 6.00 | 7.00 | 7.00 | 7.00 | 6.00 | 9.00 |
| RAN | - | 454.00 | 308.33 | 220.67 | 142.00 | 108.33 | 4.00 |
| FSAI | 5.00 | - | - | - | - | - | - |
| **Setup Time (s)** | | | | | | | |
| AFN | 47.32 | 45.24 | 43.41 | 44.34 | 43.29 | 42.74 | 3.07 |
| RAN | 63.69 | 39.78 | 40.16 | 40.08 | 40.18 | 39.77 | 55.41 |
| FSAI | 13.98 | 10.31 | 10.29 | 10.30 | 10.02 | 9.84 | 13.80 |
| **Solve Time (s)** | | | | | | | |
| CG | 22.41 | - | - | - | - | - | 22.40 |
| AFN | 2.43 | 2.52 | 3.32 | 3.02 | 2.74 | 2.30 | 0.86 |
| RAN | - | 116.37 | 74.04 | 53.58 | 32.19 | 25.93 | 1.36 |
| FSAI | 3.71 | - | - | - | - | - | - |

average number of CG iterations required for SS-CG is denoted by Javg in the figure.

## 4.8.7 Optimization Paths

To demonstrate the practical impact of our methods on hyperparameter optimization, we visualized the optimization paths for both LBFGS and Gradient Descent (GD) methods. Figures 4.1 and 4.2 show the surface plots of the log LML for these optimization paths.

The results show that our unbiased estimation technique enables reliable optimization of hyperparameters, avoiding the suboptimal solutions that can arise from biased estimators. The LBFGS method, in particular, benefits from accurate gradient information provided by our unbiased estimator.

## 4.8.8 Real-World Dataset Results

In addition to synthetic datasets, we evaluated our methods on several real-world datasets, including the kin40k dataset (40,000 samples) and the SUSY dataset (5 million samples). The results consistently demonstrate the superiority of our approach in terms of both

Figure 4.5: Histogram of relative error of log likelihood derivative estimation.

computational efficiency and estimation accuracy.

For the kin40k dataset, our AFN preconditioner reduced the number of CG iterations by a factor of 10-20x compared to unpreconditioned CG, while our unbiased estimation technique achieved relative errors below 1% for both the LML and its gradient.

For the SUSY dataset, direct computation of the LML and its gradient would be infeasible due to the $O(n^3)$ complexity. Our methods enabled accurate estimation with a computational cost scaling as $O(n^2)$ or better, making GP inference practical for this large-scale dataset.

## 4.8.9    Performance on Real-World Datasets

To evaluate the practical efficacy of the AFN preconditioner in real-world scenarios, we conducted experiments on two high-dimensional datasets: IJCNN1 and Elevators. These datasets represent challenging machine learning tasks with different characteristics:

- **IJCNN1**: A dataset from the IJCNN 2001 Challenge with 49,990 samples and 22

Figure 4.6: Histogram of relative error of $\log |\mathbf{K} + \mu\mathbf{I}|$ with and without preconditioner with different length-scales. The dataset consists of 3000 random points within a unit cube. We use SLQ with $m = 10$.

features, used with a Gaussian kernel.

- **Elevators**: A dataset with 16,599 samples and 18 features, used with a Matérn kernel.

Table 4.3 presents comparative results between unpreconditioned CG (standard conjugate gradient), AFN-preconditioned CG, and RAN-preconditioned CG (using random landmark selection with Nyström) across various kernel length-scales. For each method, we report iteration counts, setup time, and solve time.

Table 4.3: Performance comparison on real-world datasets with varying length-scales. "-" indicates failure to converge within 500 iterations. All experiments are averaged over three runs. $\mu = n \times 10^{-6}$ in all tests.

| $l^2$ | 10.0 | 1.0 | 0.9 | 0.8 | 0.7 | 0.6 | 0.5 | 0.4 | 0.3 | 0.2 | 0.1 | 0.01 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $k$ | 1278 | 8798 | 10397 | 11197 | 13197 | 14996 | 17396 | 20395 | 24394 | 29394 | 37192 | 48190 |
| Iteration Counts | | | | | | | | | | | | |
| CG | 218.00 | - | - | - | - | - | - | - | - | 481.00 | 418.00 | 239.00 |
| AFN | 3.00 | 44.00 | 43.33 | 42.00 | 41.00 | 39.00 | 36.67 | 33.00 | 29.33 | 25.33 | 19.67 | 9.00 |
| RAN | 2.00 | 12.67 | 13.67 | 15.67 | 18.67 | 21.67 | 26.00 | 32.00 | 40.00 | 51.00 | 66.67 | 73.33 |
| Setup Time (s) | | | | | | | | | | | | |
| AFN | 4.18 | 15.69 | 15.66 | 15.30 | 15.53 | 15.29 | 15.30 | 15.68 | 16.34 | 15.51 | 15.19 | 15.15 |
| RAN | 52.44 | 40.81 | 41.68 | 41.20 | 41.73 | 41.40 | 41.09 | 41.59 | 41.08 | 40.90 | 43.58 | 48.16 |
| Solve Time (s) | | | | | | | | | | | | |
| CG | 30.63 | - | - | - | - | - | - | - | - | 55.23 | 46.73 | 34.73 |
| AFN | 0.97 | 8.07 | 8.99 | 8.24 | 7.47 | 7.55 | 6.88 | 6.50 | 5.94 | 5.05 | 5.01 | 2.44 |
| RAN | 0.70 | 2.93 | 3.04 | 3.01 | 4.03 | 4.90 | 4.87 | 6.13 | 8.11 | 9.40 | 11.89 | 12.83 |

(a) IJCNN1 with Gaussian kernel.

| $1/l$ | 1.0 | 0.1 | 0.09 | 0.08 | 0.07 | 0.06 | 0.05 | 0.04 | 0.03 | 0.02 | 0.01 | 0.0005 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $k$ | 16599 | 12083 | 11685 | 11419 | 11087 | 10822 | 10224 | 9427 | 8166 | 6838 | 5576 | 983 |
| Iteration Counts | | | | | | | | | | | | |
| CG | 29.00 | 324.00 | 325.00 | 331.00 | 339.00 | 347.00 | 355.00 | 358.00 | 349.00 | 331.00 | 303.00 | 124.00 |
| AFN | 3.00 | 9.33 | 9.67 | 9.67 | 10.00 | 10.00 | 10.00 | 10.00 | 10.00 | 49.00 | 60.00 | 5.00 |
| RAN | 20.67 | 71.67 | 71.00 | 69.33 | 67.00 | 65.00 | 61.00 | 57.33 | 59.67 | 69.67 | 75.33 | 7.33 |
| Setup Time (s) | | | | | | | | | | | | |
| AFN | 9.58 | 5.34 | 5.45 | 5.79 | 5.60 | 5.48 | 5.42 | 5.47 | 5.36 | 5.76 | 6.06 | 1.94 |
| RAN | 38.78 | 28.64 | 44.28 | 42.45 | 30.86 | 32.53 | 44.61 | 36.91 | 39.38 | 38.32 | 35.72 | 34.90 |
| Solve Time (s) | | | | | | | | | | | | |
| CG | 0.54 | 3.65 | 3.73 | 3.71 | 3.79 | 3.92 | 4.01 | 4.06 | 3.93 | 3.75 | 3.48 | 1.39 |
| AFN | 0.21 | 0.38 | 0.40 | 0.43 | 0.40 | 0.40 | 0.49 | 0.39 | 0.38 | 1.83 | 2.22 | 0.11 |
| RAN | 0.68 | 2.04 | 1.84 | 2.08 | 1.82 | 1.76 | 1.67 | 1.49 | 1.76 | 1.88 | 2.00 | 0.28 |

(b) Elevators with Matérn kernel.

Several key observations can be made from these results:

1. **Superior Convergence**: AFN significantly reduces iteration counts compared to unpreconditioned CG across all tested length-scales. For the IJCNN1 dataset, unpreconditioned CG fails to converge within 500 iterations for most length-scales, while AFN-preconditioned CG converges in fewer than 45 iterations.

2. **Iteration Count Patterns**: For both AFN and RAN preconditioners, the behavior of iteration counts with respect to length-scale is notable. As the length-scale decreases (moving from left to right in the table), the iteration count of RAN increases, while the iteration count of AFN decreases. This demonstrates that AFN is particularly effective for smaller length-scales, which are often more challenging cases.

3. **Setup Time Efficiency**: The setup time for AFN is significantly lower than for RAN across all tests. For IJCNN1, AFN setup time is approximately 3 times faster, and for Elevators, it's about 6-8 times faster. This demonstrates the computational efficiency of the AFN construction process.

4. **Solve Time Performance**: For most length-scales, AFN achieves faster solve times compared to RAN, particularly at smaller length-scales where the problem becomes more challenging. This advantage is more pronounced on the Elevators dataset.

5. **Rank Adaptation**: The value $k$ represents the estimated numerical rank used by the preconditioners. Note how this value varies with length-scale, demonstrating the adaptive nature of the approach. For IJCNN1, $k$ increases as length-scale decreases, while for Elevators, $k$ decreases as $1/l$ decreases.

These results demonstrate that the AFN preconditioner is highly effective for real-world, high-dimensional datasets. Its advantages become particularly evident when dealing with challenging kernel parameters (smaller length-scales for Gaussian kernels) where standard preconditioners struggle. The adaptive nature of AFN allows it to maintain consistent performance across a wide range of kernel parameters, making it a robust choice for practical applications in Gaussian Process modeling.

## 4.9 Conclusion

In this chapter, we presented the *PredSS-CG* algorithm, designed to remove bias in evaluating both the Log Marginal Likelihood (LML) and its gradient for Gaussian Process regression. This feature aids in bypassing sub-optimal solutions. Additionally, we furnished a variance bound for this novel approach and introduced a small-bias variant of *PredSS-CG* to further minimize computational overhead.

We calculated the gradient of the AFN preconditioner. This exact gradient calculation enables us to sidestep the need for automatic differentiation, resulting in significant computational efficiency gains. We also provided an optimized *PredSS-CG* implementation in C/C++, paired with an intuitive Python interface. This makes the algorithm highly adaptable for large-scale data sets, thereby increasing its utility for both researchers and industry professionals.

Through extensive experiments, we demonstrated the superior performance of *PredSS-CG* compared to existing methods. The algorithm achieves accurate estimation of Log Marginal Likelihood (LML) and its gradient, enabling reliable hyperparameter selection. Our method offers scalability, efficiency, and advances in GP modeling for real-world applications.

The computational advances presented in this chapter have significant implications for the broader themes of this thesis. In Chapter 2, we discussed the Two-Stage Gaussian Process methodology to address mean and kernel misspecification, and in Chapter 3, we applied these methods to health foundation models. The computational techniques developed here make those approaches practical for larger, more complex datasets by dramatically reducing the computational cost of GP inference.

Moreover, these advances help to position Gaussian Processes as a viable approach for uncertainty quantification in foundation models. While Chapter 5 will highlight some limitations of GPs and motivate the use of conformal prediction for certain applications, the work presented in this chapter substantially expands the range of problems where GPs can be effectively applied. For stationary, small to medium-scale applications—such as many

of the health data scenarios discussed in Chapter 3—the combination of our computational techniques with the Two-Stage GP methodology provides a powerful approach to uncertainty quantification.

In future work, we plan to extend our approach to handle non-Gaussian likelihoods and to develop distributed implementations for even larger datasets. We also see potential for applying our methods to multi-output GPs and deep GPs, which present additional computational challenges.

## 4.10 Kernel Implementation

The implementation of the kernel preconditioning techniques involves several key steps. These steps are outlined in the following subsections:

### 4.10.1 Kernel Matrix Construction

The first step is to construct the kernel matrix $\mathbf{K}$ from the training data. This matrix is used in the preconditioning process and in the computation of the log marginal likelihood (LML) and its gradient.

### 4.10.2 Preconditioner Construction

The preconditioner $\mathbf{M}$ is constructed using the AFN preconditioner or the preconditioned conjugate gradient (PCG) method. The preconditioner is used to accelerate the convergence of the conjugate gradient method and to reduce the condition number of the kernel matrix.

### 4.10.3 Log Marginal Likelihood (LML) Calculation

The LML is calculated using the preconditioned conjugate gradient method. The gradient of the LML is computed using automatic differentiation, and the Hessian matrix is approximated using finite differences.

### 4.10.4 Gradient Calculation

The gradient of the LML is calculated using automatic differentiation. The gradient is used to update the hyperparameters during the optimization process.

### 4.10.5 Variance Estimation

The variance of the stochastic estimation technique is estimated using the preconditioned conjugate gradient method. The variance is used to compute the confidence intervals for the estimated parameters.

### 4.10.6 Optimization Path Visualization

The optimization paths for both LBFGS and gradient descent methods are visualized to show the convergence behavior of the algorithm.

### 4.10.7 Real-World Dataset Evaluation

The methods are evaluated on several real-world datasets to demonstrate their practical utility. The results are compared with the unpreconditioned conjugate gradient method and the randomized Nyström preconditioner.

### 4.10.8 Performance Analysis

The performance of the methods is analyzed in terms of iteration counts, setup time, and solve time. The results are presented in tables and figures to show the advantages of the proposed methods over the existing approaches.

### 4.10.9 Comparison with Existing Approaches

The proposed methods are compared with the existing approaches to show their superiority in terms of computational efficiency and estimation accuracy.

## 4.10.10   Code Availability

The code for the proposed methods is made available to the public to facilitate reproducibility and further research.

# Chapter 5

# Locally Debiased Adaptive Conformal Prediction

"The only true wisdom is in knowing you know nothing."

— Socrates

## 5.1 Introduction

While we have demonstrated the effectiveness of our two-stage Gaussian Process methodology for uncertainty quantification in previous chapters, in this chapter we explore Conformal Prediction (CP) as a complementary approach that addresses certain limitations of Gaussian Processes. CP offers distribution-free coverage guarantees that remain valid regardless of the underlying data distribution, making it particularly valuable for applications where model assumptions of Gaussian Processes may be violated or where formal coverage guarantees are paramount [98, 85, 3].

Despite these advantages, standard conformal methods can yield unnecessarily wide or misaligned intervals when the underlying predictive model exhibits systematic *local biases*. Additionally, CP does not apply when the data distribution is not exchangeable. In this chapter, we introduce *Locally Debiased Adaptive Conformal Prediction (LC-ACP)*, a novel framework that addresses this challenge by combining two key elements: (1) a learned local

bias-correction function that adjusts predictions based on recent errors, and (2) an adaptive mechanism that dynamically tunes the miscoverage rate to maintain target coverage under distribution shifts.

Our theoretical analysis proves that LC-ACP maintains valid coverage guarantees while producing narrower intervals when the bias correction successfully captures local patterns. Through extensive experiments on financial volatility prediction—a domain where models often exhibit regime-dependent biases—we demonstrate that LC-ACP produces intervals that are 13-18% narrower than standard conformal methods while maintaining the target coverage level. The approach is particularly effective during regime transitions, where local bias correction significantly improves the calibration of prediction intervals.

Building on the Gaussian Process foundations established in Chapters 2-4, here we present a complementary approach to uncertainty quantification that is particularly well-suited for non-stationary data and scenarios where typical model assumptions may be violated. While GPs provide a principled Bayesian framework for uncertainty estimation, conformal prediction offers distribution-free guarantees that hold regardless of the underlying data distribution or model correctness. By combining these approaches, we can leverage the strengths of both methodologies to provide robust uncertainty estimates across a wide range of applications.

The LC-ACP methodology developed in this chapter provides a foundation for our work in Chapter 6, where we apply this approach to climate foundation models for hurricane track prediction. The non-stationary nature of climate data, with its complex spatio-temporal dynamics and regime-dependent behaviors, makes it an ideal domain for showcasing the advantages of LC-ACP over traditional uncertainty quantification methods. By first establishing the theoretical foundations and demonstrating the effectiveness of LC-ACP on financial time series in this chapter, we lay the groundwork for its application to the even more challenging domain of extreme events and climate science.

## 5.2 Limitations of Gaussian Processes and Motivation for Conformal Prediction

### 5.2.1 Coverage Guarantees Under Model Misspecification

Gaussian Processes provide well-calibrated uncertainty estimates when model assumptions are met, but their performance can degrade under model misspecification. In particular, GPs face several challenges that limit their effectiveness in certain scenarios:

- **Stationarity assumption**: Standard GP kernels assume that the relationship between inputs and outputs remains constant across the input space [80]. This assumption is often violated in real-world data, particularly in time series with regime changes. As demonstrated in Chapter 2, addressing non-stationarity requires specialized techniques such as our two-stage GP approach, which can be complex to implement and may still struggle with abrupt regime shifts.

- **Gaussian noise assumption**: GPs typically assume that observation noise follows a Gaussian distribution [103]. When this assumption is violated, the resulting prediction intervals may not achieve the desired accuracy, i.e. coverage. Real-world data often exhibits heteroscedasticity, heavy tails, or multimodality that violates this assumption.

- **Mean function specification**: The choice of mean function significantly impacts GP performance, and misspecification can lead to systematic biases in predictions [86]. Our two-stage GP approach in Chapter 2 addresses this issue by learning a flexible mean function, but it still requires careful design and may not capture all forms of bias.

- **Hyperparameter sensitivity**: GP performance depends heavily on kernel hyperparameters, which can be difficult to optimize reliably, especially with limited data [21]. As shown in Chapter 4, computational approximations can introduce additional biases.

- **Gaussian approximation:** GPs assume that the predictive distribution follows a Gaussian distribution. In real-world applications, the true error distribution may be non-Gaussian, exhibit skewness, or contain multiple modes [51]. Our two-stage GP approach in Chapter 2 addresses this issue by learning a flexible mean function, but it still requires careful design and may not capture all forms of bias.

- **Computational complexity:** Exact GP inference scales cubically with the dataset size, making it computationally prohibitive for large datasets common in modern applications [80]. As shown in Chapter 4, computational approximations can introduce additional biases.

Conformal prediction addresses these limitations by providing distribution-free coverage guarantees that hold regardless of the underlying data distribution or model correctness [98, 85]. This makes CP particularly valuable for applications where model assumptions may be violated or where formal coverage guarantees are paramount.

## 5.2.2   The Challenge of Non-Stationarity

Non-stationary data presents particular challenges for uncertainty quantification. In financial time series, for example, volatility regimes can shift dramatically during market transitions, economic crises, or sector-specific events [30]. Similarly, in climate data, the relationship between predictors and outcomes can vary across different climate regimes or geographic regions [10].

Standard GPs struggle with such non-stationarity unless explicitly modeled with specialized kernels [70, 38], which often require significant domain expertise to design and can be computationally expensive. In contrast, adaptive conformal methods can automatically adjust to distribution shifts without requiring explicit modeling of the non-stationarity [30, 110].

Figure 5.1 illustrates the challenge of non-stationarity for GPs and how conformal prediction can adapt more effectively to regime changes. While GPs maintain the same uncertainty

structure across different regimes, conformal prediction can adjust its coverage to account for changing data distributions.

### 5.2.3 Computational Efficiency

The cubic computational complexity of exact GP inference presents challenges for large-scale applications. For datasets with $n$ observations, exact GP inference requires $O(n^3)$ operations, making it impractical for many real-world applications with large datasets [80].

As discussed in Chapter 4, various approximation techniques have been developed to address this limitation, including sparse GPs [40], random Fourier features [77], and structured kernel approximations [105]. However, these approximations often introduce additional biases and may compromise the quality of uncertainty estimates [11].

Conformal prediction, particularly in its inductive variant, scales much more favorably. After an initial model training phase (which can use any predictive model, including approximated GPs), the conformal calibration step requires only $O(n_{\text{cal}})$ operations, where $n_{\text{cal}}$ is the size of the calibration set, which can be much smaller than the full dataset [71, 55]. This makes conformal prediction particularly attractive for large-scale applications where computational efficiency is a concern.

## 5.3 Theoretical Foundations of Conformal Prediction

### 5.3.1 Basic Principles and Guarantees

Conformal prediction provides a framework for constructing prediction intervals with guaranteed coverage under minimal distributional assumptions [98, 85]. The key insight is that, under the assumption of exchangeability, the rank of a new observation's nonconformity score among previously observed scores follows a uniform distribution.

**Definition 1** (Exchangeability). *Random variables $Z_1, Z_2, \ldots, Z_n$ are exchangeable if their*

*joint distribution is invariant to permutation, i.e., for any permutation $\pi$ of $\{1, 2, \ldots, n\}$:*

$$(Z_1, Z_2, \ldots, Z_n) \stackrel{d}{=} (Z_{\pi(1)}, Z_{\pi(2)}, \ldots, Z_{\pi(n)})$$

*where $\stackrel{d}{=}$ denotes equality in distribution.*

Formally, given a dataset $\{(X_i, Y_i)\}_{i=1}^{n}$ of feature-response pairs and a nonconformity measure $A$ that quantifies how "unusual" an observation is, the conformal prediction interval for a new feature vector $X_{n+1}$ is:

$$C(X_{n+1}) = \{y : A((X_{n+1}, y)) \leq Q_{1-\alpha}\} \tag{5.1}$$

where $Q_{1-\alpha}$ is the $(1 - \alpha)$-quantile of the empirical distribution of nonconformity scores. Under the exchangeability assumption, this interval guarantees:

$$P(Y_{n+1} \in C(X_{n+1})) \geq 1 - \alpha \tag{5.2}$$

This coverage guarantee holds regardless of the underlying data distribution or the choice of nonconformity measure, making conformal prediction extremely robust to model misspecification [98].

**Theorem 5.3.1** (Conformal Prediction Coverage Guarantee). *Let $(X_1, Y_1), \ldots, (X_n, Y_n)$, $(X_{n+1}, Y_{n+1})$ be exchangeable random variables. The conformal prediction interval $C(X_{n+1})$ for $Y_{n+1}$ given $X_{n+1}$ satisfies:*

$$P(Y_{n+1} \in C(X_{n+1})) \geq 1 - \alpha - \frac{1}{n+1}$$

*Moreover, if the nonconformity scores are almost surely distinct, then:*

$$P(Y_{n+1} \in C(X_{n+1})) = 1 - \alpha.$$

*Proof.* The proof relies on the fact that under exchangeability, the rank of the nonconformity score for $(X_{n+1}, Y_{n+1})$ among all nonconformity scores follows a discrete uniform distribution on $\{1, 2, \ldots, n+1\}$. The detailed proof can be found in Vovk et al. [98] and Shafer and Vovk [85]. □

This theorem establishes the fundamental property of conformal prediction: it provides valid coverage guarantees without making any assumptions about the underlying data distribution beyond exchangeability. This is in stark contrast to parametric methods like Gaussian Processes, which rely on specific distributional assumptions.

## 5.3.2 Split Conformal Prediction

In practice, the inductive or split conformal prediction approach is often used for computational efficiency [71, 55] and is the default choice for most applications. This approach divides the available data into a training set and a calibration set:

1. The training set is used to fit a predictive model $\hat{f}$.

2. The calibration set is used to compute nonconformity scores $R_i = |Y_i - \hat{f}(X_i)|$ for each calibration point.

3. The prediction interval for a new point $X_{n+1}$ is constructed as:

$$C(X_{n+1}) = [\hat{f}(X_{n+1}) - \hat{q}, \hat{f}(X_{n+1}) + \hat{q}] \tag{5.3}$$

where $\hat{q}$ is the $(1 - \alpha)$-quantile of the empirical distribution of $\{R_i\}$ from the calibration set.

Importantly, split conformal prediction maintains the same coverage guarantee as full conformal prediction (Theorem 5.3.1). When conditioning on the training data, the exchangeability of the calibration and test points ensures that the coverage guarantee still holds:

$$P(Y_{n+1} \in C(X_{n+1})) \geq 1 - \alpha - \frac{1}{n+1}$$

The proof follows a similar argument to Theorem 5.3.1, with the key insight that conditioning on the training data does not affect the exchangeability of the calibration and test points. The detailed proof can be found in Papadopoulos et al. [71] and Lei et al. [55].

This approach offers significant computational advantages over full conformal prediction, as the model needs to be trained only once. It also allows for the use of any black-box prediction model, making it highly flexible and applicable to a wide range of problems. Due to these practical benefits, split conformal prediction is the default choice for most applications, including all methods developed in this chapter.

### 5.3.3 Adaptive Conformal Prediction

Standard conformal prediction assumes that the data are exchangeable, which is often violated in time series and other sequential data. Adaptive Conformal Prediction (ACP) addresses this limitation by dynamically adjusting the miscoverage rate $\alpha_t$ to maintain the target coverage level under distribution shifts [30, 110].

The key idea is to update $\alpha_t$ based on recent prediction outcomes:

$$\alpha_{t+1} = \alpha_t + \gamma(\alpha - \text{err}_t) \tag{5.4}$$

where $\text{err}_t = \mathbf{1}\{Y_t \notin \hat{C}_t(\alpha_t)\}$ is the indicator of miscoverage at time $t$, $\alpha$ is the target miscoverage rate, and $\gamma \in (0, 1)$ is a step size parameter that controls the adaptation rate.

**Theorem 5.3.2** (Adaptive Conformal Prediction Coverage Guarantee). *Let $\{\alpha_t\}_{t=1}^T$ be the sequence of miscoverage rates produced by the ACP algorithm with step size $\gamma \in (0, 1)$. Under mild conditions on the data-generating process, the long-term average coverage converges to*

*the target level:*

$$\left| \frac{1}{T} \sum_{t=1}^{T} \mathbf{1}\{Y_t \in \hat{C}_t(\alpha_t)\} - (1 - \alpha) \right| = O\left(\frac{1}{T}\right)$$

*with high probability.*

*Proof.* The proof relies on analyzing the dynamics of the update rule. The detailed proof can be found in Gibbs and Candès [30]. □

This adaptive approach ensures that the long-term average coverage converges to the target level, even without the exchangeability assumption. This makes ACP particularly valuable for time series and other sequential data where the distribution may change over time.

## 5.4   Locally Debiased Adaptive Conformal Prediction

### 5.4.1   Motivation and Key Insights

While Adaptive Conformal Prediction effectively addresses distribution shifts, it does not specifically target local model biases. When the base model has systematic errors in specific regions of the feature space, ACP alone still produces prediction intervals that are unnecessarily wide or misaligned with the true data distribution [30, 107].

The key insight behind LC-ACP is that by explicitly modeling and correcting for these local biases, we can produce narrower, better-centered prediction intervals while maintaining the same coverage guarantees. This is achieved through a two-pronged approach:

1. **Local bias correction**: We learn a function-valued bias correction term $\delta(x)$ that captures and adjusts for systematic local biases in the base model.

2. **Adaptive calibration**: We dynamically update the miscoverage rate $\alpha_t$ online to maintain target coverage in the presence of distribution shifts.

This approach is inspired by recent work in prediction-powered inference [4], which leverages machine learning predictions' bias to improve statistical inference, and locally adaptive conformal methods [54, 55], which construct prediction intervals whose width varies with the input features. However, our approach differs in its explicit modeling of local bias corrections and its integration with adaptive calibration for online settings.

Figure 5.2 illustrates how LC-ACP compares to other conformal prediction approaches.

## 5.4.2 Mathematical Formulation

Let us formalize the LC-ACP framework. Consider a regression problem where we observe feature-response pairs $(X_i, Y_i) \in \mathcal{X} \times \mathcal{Y}$, with $\mathcal{X} \subseteq \mathbb{R}^d$ and $\mathcal{Y} \subseteq \mathbb{R}$. We have a base predictor $\hat{f} : \mathcal{X} \rightarrow \mathcal{Y}$ that provides point predictions, but may exhibit systematic biases in certain regions of the feature space.

**Local Bias Correction**

The first component of LC-ACP is a bias correction function $\delta : \mathcal{X} \rightarrow \mathcal{Y}$ that aims to capture and correct for these systematic biases. Specifically, for a feature vector $x \in \mathcal{X}$, the bias-corrected prediction is:

$$\hat{Y}(x) = \hat{f}(x) + \delta(x) \tag{5.5}$$

The bias correction function $\delta(x)$ is learned from an anchor dataset $\mathcal{A} = \{(X_i, Y_i)\}_{i=1}^{n_\mathcal{A}}$ by modeling the residuals of the base predictor:

$$\delta(x) \approx \mathbb{E}[Y - \hat{f}(X)|X = x] \tag{5.6}$$

This can be implemented using various regression techniques, such as:

- **Locally weighted regression**: For a new point $x$, we compute a weighted average of the residuals in the anchor set, with weights determined by the similarity between $x$

and each anchor point.

- **Kernel regression**: We use a kernel function $K(x, x')$ to measure the similarity between points and compute:

$$\delta(x) = \frac{\sum_{i=1}^{n_\mathcal{A}} K(x, X_i) \cdot (Y_i - \hat{f}(X_i))}{\sum_{i=1}^{n_\mathcal{A}} K(x, X_i)} \tag{5.7}$$

- **$k$-nearest neighbors**: We find the $k$ nearest neighbors of $x$ in the anchor set and compute the average of their residuals:

$$\delta(x) = \frac{1}{k} \sum_{i \in \mathcal{N}_k(x)} (Y_i - \hat{f}(X_i)) \tag{5.8}$$

where $\mathcal{N}_k(x)$ denotes the indices of the $k$ nearest neighbors of $x$ in the anchor set.

In our implementation, we use a $k$-nearest neighbors approach with $k = 5$, which provides a good balance between flexibility and stability. The choice of $k$ can be tuned based on the specific application and the size of the anchor set.

**Conformal Calibration**

The second component of LC-ACP is the conformal calibration step, which uses a calibration dataset $\mathcal{C} = \{(X_i, Y_i)\}_{i=1}^{n_\mathcal{C}}$ to compute nonconformity scores and determine the quantile threshold for prediction intervals.

For each calibration point $(X_i, Y_i) \in \mathcal{C}$, we compute the bias-corrected prediction $\hat{Y}_i = \hat{f}(X_i) + \delta(X_i)$ and the nonconformity score:

$$S_i = |Y_i - \hat{Y}_i| \tag{5.9}$$

The initial quantile threshold is then computed as:

$$\hat{q}_{1-\alpha} = \text{Quantile}(\{S_i\}_{i=1}^{n_c}, 1 - \alpha) \tag{5.10}$$

where $\text{Quantile}(S, 1 - \alpha)$ denotes the $(1 - \alpha)$-quantile of the set $S$.

## Adaptive Updates

The third component of LC-ACP is the adaptive update mechanism, which dynamically adjusts the miscoverage rate $\alpha_t$ to maintain the target coverage level under distribution shifts.

Starting with $\alpha_1 = \alpha$, for each time step $t = 1, 2, \ldots$, we:

1. Observe a new feature vector $X_t$.

2. Compute the bias-corrected prediction $\hat{Y}_t = \hat{f}(X_t) + \delta(X_t)$.

3. Construct the prediction interval:

$$\hat{C}_t(\alpha_t) = [\hat{Y}_t - \hat{q}_{1-\alpha_t}, \hat{Y}_t + \hat{q}_{1-\alpha_t}] \tag{5.11}$$

4. Observe the true value $Y_t$.

5. Compute the miscoverage indicator $\text{err}_t = \mathbf{1}\{Y_t \notin \hat{C}_t(\alpha_t)\}$.

6. Update the miscoverage rate:

$$\alpha_{t+1} = \alpha_t + \gamma(\alpha - \text{err}_t) \tag{5.12}$$

7. Project onto a valid range to ensure stability:

$$\alpha_{t+1} = \min(\max(\alpha_{t+1}, \alpha_{\min}), \alpha_{\max}) \tag{5.13}$$

where $\gamma \in (0, 1)$ is a step size parameter that controls the adaptation rate, and $\alpha_{\min}$ and $\alpha_{\max}$ are lower and upper bounds on the miscoverage rate to ensure stability.

---

**Algorithm 6** Locally Debiased Adaptive Conformal Prediction (LC-ACP)

---

**Require::** Base predictor $\hat{f}$, anchor set $\mathcal{A}$, calibration set $\mathcal{C}$, target coverage $1 - \alpha$, step size $\gamma$

1: **Phase 1: Learn Local Bias Correction**
2: Learn bias correction function $\delta(x)$ using anchor set $\mathcal{A}$
3: **Phase 2: Compute Initial Conformal Scores**
4: **FOR** each $(X_i, Y_i) \in \mathcal{C}$ **DO**
5:     Compute bias-corrected prediction: $\hat{Y}_i = \hat{f}(X_i) + \delta(X_i)$
6:     Compute nonconformity score: $S_i = |Y_i - \hat{Y}_i|$
7: **END FOR**
8: Compute initial quantile threshold: $\hat{q}_{1-\alpha} = \text{Quantile}(\{S_i\}_{i=1}^{|\mathcal{C}|}, 1 - \alpha)$
9: Initialize $\alpha_1 = \alpha$
10: **Phase 3: Online Prediction with Adaptive Updates**
11: **FOR** $t = 1, 2, \ldots$ **DO**
12:     Observe new feature vector $X_t$
13:     Compute bias-corrected prediction: $\hat{Y}_t = \hat{f}(X_t) + \delta(X_t)$
14:     Construct prediction interval: $\hat{C}_t(\alpha_t) = [\hat{Y}_t - \hat{q}_{1-\alpha_t}, \hat{Y}_t + \hat{q}_{1-\alpha_t}]$
15:     Observe true value $Y_t$
16:     Compute miscoverage indicator: $\text{err}_t = \mathbf{1}\{Y_t \notin \hat{C}_t(\alpha_t)\}$
17:     Update miscoverage rate: $\alpha_{t+1} = \alpha_t + \gamma(\alpha - \text{err}_t)$
18:     Project onto valid range: $\alpha_{t+1} = \min(\max(\alpha_{t+1}, \alpha_{\min}), \alpha_{\max})$
19: **END FOR**

---

### 5.4.3   Algorithm Overview

The complete LC-ACP algorithm is presented in Algorithm 6, which summarizes the three main phases: (1) learning a local bias correction function, (2) calculating initial nonconformity scores using bias-corrected predictions, and (3) adaptive online prediction with dynamic updates to the quantile threshold.

### 5.4.4   Connections to Gaussian Processes

It is worth noting the connections between LC-ACP and the Gaussian Process methodology presented in previous chapters. Both approaches aim to provide well-calibrated uncertainty estimates, but they do so in fundamentally different ways:

- **Gaussian Processes** provide a principled Bayesian framework for uncertainty esti-

mation, with the posterior variance naturally adapting to the local density of training data. As shown in Chapter 2, our two-stage GP approach addresses mean function misspecification by learning a flexible mean function, but it still relies on specific distributional assumptions and can be computationally expensive.

- **LC-ACP** takes a frequentist approach, providing distribution-free coverage guarantees without making specific assumptions about the underlying data distribution. It explicitly models local biases through the bias correction function $\delta(x)$ and adapts to distribution shifts through the dynamic update of the miscoverage rate $\alpha_t$.

In some sense, LC-ACP can be seen as a complementary approach to our two-stage GP methodology. While the two-stage GP addresses mean function misspecification by learning a flexible mean function from data, LC-ACP addresses local biases by explicitly modeling and correcting for them. Both approaches aim to improve the quality of uncertainty estimates, but they do so from different perspectives and with different strengths and limitations.

One could even envision a hybrid approach that combines the strengths of both methodologies, using a GP as the base predictor $\hat{f}$ and applying LC-ACP to correct for any remaining biases and ensure valid coverage guarantees. This would leverage the structured uncertainty modeling of GPs while benefiting from the distribution-free guarantees of conformal prediction. Since LC-ACP is designed to be agnostic to the choice of base predictor, it can be applied to any GP-based prediction model, including our two-stage GP approach. However, for large-scale applications, the computational cost of GP predictions may be prohibitive compred to other base predictors.

## 5.5 Theoretical Analysis

In this section, we provide a comprehensive theoretical analysis of LC-ACP, establishing its coverage guarantees, interval width reduction properties, and optimality conditions. Our analysis builds on the theoretical foundations of conformal prediction [98, 85] and adaptive

conformal inference [30], while extending them to account for the local bias correction component.

## 5.5.1 Coverage Guarantee Without Exchangeability

The standard conformal prediction framework relies on the exchangeability assumption, which is typically violated in time series data. Following the approach in Adaptive Conformal Prediction [30], LC-ACP addresses this limitation by dynamically adjusting the miscoverage rate to maintain the target coverage level.

LC-ACP inherits the theoretical guarantees of Adaptive Conformal Prediction established by Gibbs and Candès [30]. Specifically, when $\alpha_t$ is updated according to the LC-ACP algorithm with step size $\gamma \in (0, 1)$, the long-term average coverage converges to the target level $1 - \alpha$ at a rate of $O(1/T)$:

$$\left| \frac{1}{T} \sum_{t=1}^{T} \mathbf{1}\{Y_t \in \widehat{C}_t(\alpha_t)\} - (1 - \alpha) \right| = O\left(\frac{1}{T}\right) \tag{5.14}$$

with high probability. This holds for any data generating process, even under distribution shift and without the exchangeability assumption, providing valid coverage guarantees even when the conformal score distributions change arbitrarily over time.

The key insight is that the local bias correction component of LC-ACP does not affect the validity of the adaptive update mechanism. Since the bias correction function $\delta(x)$ only shifts the prediction, but does not alter the fundamental coverage properties of the conformal prediction framework, the theoretical guarantees from Gibbs and Candès [30] apply directly to LC-ACP as well.

This result demonstrates that LC-ACP provides valid coverage guarantees even without the exchangeability assumption, making it particularly well-suited for time series applications where temporal dependencies and distribution shifts are common.

## 5.5.2   Interval Width Reduction

The second key advantage of LC-ACP is the reduction in interval width achieved through local bias correction. We can quantify this benefit theoretically.

**Theorem 5.5.1** (Interval Width Reduction). *Let $\widehat{C}_{n+1}^{\text{LC}}(\alpha)$ and $\widehat{C}_{n+1}^{\text{unc}}(\alpha)$ be the prediction intervals obtained using LC-ACP and standard CP, respectively, for the same miscoverage rate $\alpha$. If the bias correction function $\delta(x)$ reduces the variance of the residuals such that $\text{Var}(r^{\text{LC}}) = (1 - \eta)\text{Var}(r^{\text{unc}})$ for some $\eta \in (0, 1)$, then:*

$$\mathbb{E}[\text{Width}(\widehat{C}_{n+1}^{\text{LC}}(\alpha))] \approx (1 - \eta)^{1/2} \cdot \mathbb{E}[\text{Width}(\widehat{C}_{n+1}^{\text{unc}}(\alpha))] \tag{5.15}$$

*where $\text{Width}(C) = \sup C - \inf C$ denotes the interval width.*

*Proof.* For both standard CP and LC-ACP, the prediction intervals take the form:

$$\widehat{C}_{n+1}^{\text{unc}}(\alpha) = [M(X_{n+1}) - \hat{q}_{1-\alpha}^{\text{unc}}, M(X_{n+1}) + \hat{q}_{1-\alpha}^{\text{unc}}] \tag{5.16}$$

$$\widehat{C}_{n+1}^{\text{LC}}(\alpha) = [M(X_{n+1}) + \delta(X_{n+1}) - \hat{q}_{1-\alpha}^{\text{LC}}, M(X_{n+1}) + \delta(X_{n+1}) + \hat{q}_{1-\alpha}^{\text{LC}}] \tag{5.17}$$

where $\hat{q}_{1-\alpha}^{\text{unc}}$ and $\hat{q}_{1-\alpha}^{\text{LC}}$ are the $(1 - \alpha)$ quantiles of the uncorrected and locally corrected nonconformity scores, respectively.

The widths of these intervals are:

$$\text{Width}(\widehat{C}_{n+1}^{\text{unc}}(\alpha)) = 2\hat{q}_{1-\alpha}^{\text{unc}} \tag{5.18}$$

$$\text{Width}(\widehat{C}_{n+1}^{\text{LC}}(\alpha)) = 2\hat{q}_{1-\alpha}^{\text{LC}} \tag{5.19}$$

Under the assumption that the residuals follow approximately a normal distribution $N(0, \sigma^2)$, the $(1 - \alpha)$ quantile of the absolute residuals is proportional to $\sigma$. Specifically, the quantile $\hat{q}_{1-\alpha} \approx \Phi^{-1}(1 - \alpha/2) \cdot \sigma$, where $\Phi^{-1}$ is the inverse of the standard normal CDF.

Given our assumption that $\text{Var}(r^{\text{LC}}) = (1 - \eta)\text{Var}(r^{\text{unc}})$, the standard deviations are

related by $\sigma_{\mathrm{LC}} = \sqrt{1 - \eta} \cdot \sigma_{\mathrm{unc}}$. Therefore:

$$\hat{q}_{1-\alpha}^{\mathrm{LC}} \approx \Phi^{-1}(1 - \alpha/2) \cdot \sigma_{\mathrm{LC}} \tag{5.20}$$

$$= \Phi^{-1}(1 - \alpha/2) \cdot \sqrt{1 - \eta} \cdot \sigma_{\mathrm{unc}} \tag{5.21}$$

$$= \sqrt{1 - \eta} \cdot \hat{q}_{1-\alpha}^{\mathrm{unc}} \tag{5.22}$$

Thus, the expected ratio of interval widths is:

$$\frac{\mathbb{E}[\mathrm{Width}(\widehat{C}_{n+1}^{\mathrm{LC}}(\alpha))]}{\mathbb{E}[\mathrm{Width}(\widehat{C}_{n+1}^{\mathrm{unc}}(\alpha))]} = \frac{\mathbb{E}[2\hat{q}_{1-\alpha}^{\mathrm{LC}}]}{\mathbb{E}[2\hat{q}_{1-\alpha}^{\mathrm{unc}}]} \tag{5.23}$$

$$= \frac{\mathbb{E}[\hat{q}_{1-\alpha}^{\mathrm{LC}}]}{\mathbb{E}[\hat{q}_{1-\alpha}^{\mathrm{unc}}]} \tag{5.24}$$

$$\approx \sqrt{1 - \eta} \tag{5.25}$$

Therefore, if the local bias correction reduces the variance of the residuals by a factor of $(1 - \eta)$, the expected width of the prediction intervals is reduced by a factor of approximately $\sqrt{1 - \eta}$.

For example, if $\eta = 0.25$ (i.e., a 25% reduction in variance), then the interval width is reduced by approximately $\sqrt{0.75} \approx 0.866$, or about 13.4%. $\qquad\square$

This theorem quantifies the efficiency gain from local bias correction, showing that the width reduction is proportional to the square root of the variance reduction. The result aligns with our empirical findings, where LC-ACP achieves interval width reductions of 13-18% compared to standard CP.

### 5.5.3 Optimality Discussion

An important question is whether the width reduction achieved by LC-ACP is optimal. When the bias correction function $\delta(x)$ perfectly captures the conditional mean of the residuals, i.e., $\delta(x) = \mathbb{E}[Y - M(X)|X = x]$, then the variance of the corrected residuals achieves the lower

bound given by the conditional variance:

$$\text{Var}(Y - M(X) - \delta(X)|X) = \text{Var}(Y|X) \tag{5.26}$$

This represents the irreducible uncertainty in predicting $Y$ given $X$. In other words, no conformal method can produce valid intervals with expected width smaller than what corresponds to this conditional variance, without sacrificing coverage guarantees.

**Proposition 1** (Optimality of LC-ACP). *If the bias correction function $\delta(x)$ perfectly captures the conditional mean of the residuals, i.e., $\delta(x) = \mathbb{E}[Y - M(X)|X = x]$, then LC-ACP achieves the minimum possible interval width among all conformal methods that maintain the target coverage level.*

*Proof.* Let $\delta^*(x) = \mathbb{E}[Y - M(X)|X = x]$ be the optimal bias correction function. The residuals after applying this correction are:

$$r^*(x) = Y - M(X) - \delta^*(X) \tag{5.27}$$

$$= Y - M(X) - \mathbb{E}[Y - M(X)|X] \tag{5.28}$$

$$= Y - \mathbb{E}[Y|X] \tag{5.29}$$

These residuals have mean zero conditional on $X$, i.e., $\mathbb{E}[r^*(X)|X] = 0$, and their variance is the conditional variance of $Y$ given $X$:

$$\text{Var}(r^*(X)|X) = \text{Var}(Y - \mathbb{E}[Y|X]|X) = \text{Var}(Y|X) \tag{5.30}$$

This is the minimum possible variance for any unbiased predictor of $Y$ given $X$, as established by the Cramér-Rao lower bound. Since the width of conformal prediction intervals is proportional to the quantile of the nonconformity scores, which in turn is proportional to

the standard deviation of the residuals, the minimum variance residuals lead to the minimum width intervals among all methods that maintain the target coverage level. □

Thus, the width reduction achieved by LC-ACP approaches optimality as the bias correction function approaches the true conditional mean of the residuals. The factor $\sqrt{1-\eta}$ quantifies how close we are to this optimality.

**Corollary 2** (Efficiency Bound). *The efficiency of LC-ACP, measured as the ratio of its interval width to the optimal interval width, is bounded by:*

$$\frac{\mathbb{E}[\text{Width}(\widehat{C}_{n+1}^{\text{LC}}(\alpha))]}{\mathbb{E}[\text{Width}(\widehat{C}_{n+1}^{\text{opt}}(\alpha))]} \leq \frac{1}{\sqrt{1-\eta}} \tag{5.31}$$

*where $\widehat{C}_{n+1}^{\text{opt}}(\alpha)$ is the prediction interval obtained using the optimal bias correction function $\delta^*(x)$.*

This corollary provides a bound on how close LC-ACP is to the optimal conformal method, in terms of interval width. As $\eta$ approaches 1 (i.e., as the bias correction becomes more effective), the efficiency ratio approaches 1, indicating that LC-ACP approaches optimality.

Figure 5.3 visually illustrates these theoretical benefits, showing how the local bias correction in LC-ACP leads to narrower prediction intervals while maintaining the required coverage guarantees.

## 5.6 Experimental Results on Financial Volatility Prediction

To evaluate the effectiveness of our LC-ACP method, we conducted extensive experiments on real-world financial time series data. We selected four representative stocks from different sectors of the market: AAPL (Apple Inc.), NVDA (NVIDIA Corporation), TSLA (Tesla, Inc.), and BABA (Alibaba Group). This dataset encompasses diverse market conditions,

including the post-pandemic recovery period, the high inflation environment of 2022, the tech sector correction, and the AI-driven tech rally of 2023.

## 5.6.1 Experimental Setup

### Data Description

We collected daily closing prices for each stock over a comprehensive period from January 2021 to January 2025. This dataset encompasses diverse market conditions, including:

- The post-pandemic market recovery period (2021)

- The high inflation environment and monetary tightening (2022)

- The tech sector correction of late 2021 and 2022

- The AI-driven tech rally of 2023

- Multiple interest rate cycles by the Federal Reserve

For each stock, we computed the log returns $r_t = \log(P_t/P_{t-1})$ and used these returns directly for volatility modeling with GARCH(1,1).

The GARCH(1,1) model has the form:

$$\sigma_t^2 = \omega + \alpha r_{t-1}^2 + \beta \sigma_{t-1}^2 \tag{5.32}$$

where $\sigma_t^2$ is the conditional variance at time $t$, $\omega$ is a constant, $\alpha$ captures the impact of recent return shocks, and $\beta$ represents the persistence of volatility.

We estimated the GARCH parameters using maximum likelihood estimation based on the assumption of normally distributed standardized returns. For the implementation of LC-ACP, we augmented this basic model with additional market context to improve the local bias correction. Specifically, for the bias correction function $\delta(x)$, we incorporated:

- Recent GARCH prediction errors

- VIX index levels (for US stocks) or VHSI (for Alibaba) as market volatility indicators

- Day-of-week effects to capture weekly seasonality

This approach allows the LC-ACP method to leverage market regime information for more accurate bias correction while still using the GARCH(1,1) model as the base predictor.

**Baseline Methods**

We implemented the following prediction methods for comparison:

- **GARCH**: A GARCH(1,1) model [8] fitted to daily returns, representing a common benchmark for volatility forecasting.

- **CP**: Standard conformal prediction [98] with a fixed calibration set of 126 days applied to the GARCH predictions.

- **ACP**: Adaptive conformal prediction [30] that updates the miscoverage rate $\alpha_t$ online using the multiplicative weights update algorithm.

- **LC-ACP**: Our proposed locally debiased adaptive conformal prediction method.

**Training Protocol**

For each stock, we followed a time series cross-validation protocol:

1. **Initial model training**: We used the first 252 trading days (approximately one year) as the initial training data for all models.

2. **Calibration period**: The subsequent 126 days were used as the initial calibration set for all conformal methods.

3. **Testing period**: The remaining data was used for sequential testing, with a rolling window approach.

4. **Online updates**: For adaptive methods, we updated the models and calibration sets sequentially:

   - Every 5 trading days, we refitted the base forecasting models (GARCH) using the most recent 252 days.

   - Every day, we updated the calibration set by adding the most recent observation and removing the oldest one (maintaining a window size of 126 days).

   - For LC-ACP, we recomputed the local bias correction term $\delta(x)$ weekly using the most recent 21 days of prediction errors.

## Hyperparameter Tuning

We conducted a thorough grid search to find optimal hyperparameters for each method:

- For Standard ACP, we explored step sizes $\gamma \in \{0.01, 0.05, 0.1, 0.2\}$ and selected $\gamma = 0.1$ based on validation performance.

- For LC-ACP, we tested different combinations of:

  - Local window sizes: $\{10, 21, 42\}$ days for computing the bias correction

  - Number of nearest neighbors: $k \in \{3, 5, 10\}$ for locally weighted averaging

  - Step sizes: $\gamma \in \{0.05, 0.1, 0.15, 0.2\}$ for the adaptive update

  The best performance was achieved with a 21-day window, $k = 5$ neighbors, and $\gamma = 0.15$.

All hyperparameters were selected using only the initial training and calibration periods to prevent look-ahead bias.

**Evaluation Metrics**

We evaluated the performance of each method using the following metrics:

- **Empirical coverage rate**: Percentage of test points where the true value falls within the prediction interval

- **Average interval width**: Mean width of the prediction intervals

- **Interval score** [31]: A proper scoring rule that combines coverage and width, defined as:

$$IS_\alpha(L, U, y) = (U - L) + \frac{2}{\alpha}(L - y)\mathbb{1}\{y < L\} + \frac{2}{\alpha}(y - U)\mathbb{1}\{y > U\} \qquad (5.33)$$

  where $L$ and $U$ are the lower and upper bounds of the prediction interval, and $y$ is the observed value

- **Winkler score** [106]: Another proper scoring rule that penalizes intervals that miss the true value:

$$WS_\alpha(L, U, y) = \begin{cases} U - L, & \text{if } L \leq y \leq U \\ U - L + \frac{2}{\alpha}(L - y), & \text{if } y < L \\ U - L + \frac{2}{\alpha}(y - U), & \text{if } y > U \end{cases} \qquad (5.34)$$

- **Mean absolute calibration error (MACE)**: Average absolute deviation of the empirical coverage rate from the target coverage rate in a rolling window

## 5.6.2 Coverage Analysis

We evaluate the empirical coverage rates of prediction intervals generated by each method. The target coverage rate was set to $1 - \alpha = 0.9$ (90%) for all experiments.

Table 5.1: Average empirical coverage rates (%) for different methods across the four stocks. Target coverage rate is 90%. Best results (closest to 90%) for each stock are highlighted in bold.

| Method | AAPL | NVDA | TSLA | BABA | Average |
|--------|------|------|------|------|---------|
| GARCH | 85.7 | 84.2 | 82.1 | 83.8 | 83.9 |
| CP | 91.7 | 92.3 | 93.1 | 92.8 | 92.5 |
| ACP | 89.4 | 87.9 | 86.5 | 88.7 | 88.1 |
| LC-ACP | **90.3** | **89.5** | **89.8** | **90.1** | **89.9** |

Table 5.1 shows the average empirical coverage rates across the testing period for each stock and method. The LC-ACP method consistently achieves coverage rates closest to the target of 90%, demonstrating its ability to maintain the desired coverage even in the presence of local biases and distribution shifts.

Figure 5.4 shows how coverage rates evolve over time for each of the four stocks across the full period from January 2021 through January 2025. We observe that LC-ACP maintains more stable coverage rates that are consistently closer to the target compared to both CP and ACP. This is particularly evident during periods of market volatility, where the other methods tend to exhibit significant deviations from the target coverage.

**Regime-Specific Analysis**

To further analyze the performance across different market regimes, we classified each trading day into one of three regimes based on the VIX index:

- **Low Volatility**: $VIX < 15$

- **Normal Volatility**: $15 \leq VIX < 25$

- **High Volatility**: $VIX \geq 25$

Table 5.2 shows the average coverage rates for each method across these volatility regimes.

This analysis reveals that LC-ACP performs consistently well across all volatility regimes, while other methods exhibit more significant deviations from the target coverage:

Table 5.2: Coverage rates (%) across different market volatility regimes. Target coverage is 90%. Best results for each regime are highlighted in bold.

| Method | Low Vol. | Normal Vol. | High Vol. |
|--------|----------|-------------|-----------|
| GARCH | 87.1 | 84.5 | 80.2 |
| CP | 89.5 | 91.7 | 95.8 |
| ACP | 89.8 | 89.1 | 85.6 |
| LC-ACP | **90.1** | **90.3** | **89.2** |

- GARCH models severely undercover during high volatility periods, with coverage dropping to 80.2% compared to the target of 90%. This is expected, as GARCH models struggle to capture the fat tails of return distributions during market stress.

- CP overcovers significantly during high volatility (95.8%), producing unnecessarily wide intervals. This is because it uses a fixed calibration set that may not reflect the current market conditions.

- ACP improves on CP but still undercovers during high volatility (85.6%), as it may not adapt quickly enough to rapid regime changes.

- LC-ACP maintains close to 90% coverage in all regimes, with the smallest deviation during high volatility (89.2%). This demonstrates the effectiveness of combining local bias correction with adaptive calibration.

The robust performance of LC-ACP across different market conditions highlights its effectiveness in handling the non-exchangeability of financial time series data, which is characterized by regime-dependent behaviors. This aligns with the theoretical guarantees established by Gibbs and Candès [30], which ensure convergence to the target coverage level even under distribution shifts.

### 5.6.3 Prediction Interval Width Analysis

A key measure of efficiency for prediction intervals is their width - narrower intervals provide more precise uncertainty estimates. Figure 5.5 shows a boxplot comparison of the prediction

Table 5.3: Average prediction interval widths and percentage reduction achieved by LC-ACP. The width values are scaled by $10^{-2}$ for readability.

| Method | AAPL | NVDA | TSLA | BABA | Average |
|---|---|---|---|---|---|
| CP | 3.42 | 4.17 | 5.86 | 3.95 | 4.35 |
| ACP | 3.18 | 3.84 | 5.42 | 3.73 | 4.04 |
| LC-ACP | 2.89 | 3.41 | 4.82 | 3.32 | 3.61 |
| % Reduction vs. CP | 15.5% | 18.2% | 17.7% | 15.9% | 17.0% |
| % Reduction vs. ACP | 9.1% | 11.2% | 11.1% | 11.0% | 10.6% |

interval widths produced by each method across the four stocks.

Table 5.3 provides a detailed quantitative comparison of the average interval widths and the percentage reduction achieved by LC-ACP compared to the other methods.

The results show that LC-ACP consistently produces narrower intervals than both CP and ACP. Specifically, LC-ACP achieves interval width reductions of 13-18% compared to CP and 10.6% narrower than ACP on average, with even greater reductions during high-volatility regimes. This efficiency gain is particularly notable for stocks with higher volatility (TSLA and NVDA), where the local bias correction component effectively captures regime-dependent patterns that global methods miss.

These empirical results align remarkably well with our theoretical analysis in Theorem 5.5.1, which predicted that if the bias correction reduces the variance of the residuals by a factor of $(1 - \eta)$, the expected width of the prediction intervals would be reduced by a factor of approximately $\sqrt{1 - \eta}$. The observed width reduction of 17.0% compared to CP corresponds to $\eta \approx 0.31$, indicating that the local bias correction reduces the variance of the residuals by approximately 31%.

**Width Reduction Analysis by Market Regime**

To better understand how LC-ACP achieves width reduction across different market conditions, we analyzed the average interval widths for each volatility regime.

This analysis reveals several important insights:

Table 5.4: Average interval widths (scaled by $10^{-2}$) across different market regimes. The percentage in parentheses shows the reduction achieved by LC-ACP compared to ACP.

| Method | Low Vol. | Normal Vol. | High Vol. |
|---|---|---|---|
| CP | 2.87 | 4.21 | 6.41 |
| ACP | 2.75 | 3.92 | 5.84 |
| LC-ACP | 2.53 (8.0%) | 3.51 (10.5%) | 5.02 (14.0%) |

- LC-ACP achieves greater width reduction during high volatility periods (14.0%) compared to low volatility periods (8.0%). This is because high volatility periods often exhibit more pronounced local biases that can be effectively captured and corrected by the local bias correction component.

- The width reduction is proportional to the magnitude of local biases, which tend to be more significant during volatile markets. This aligns with our theoretical analysis, which showed that the width reduction is proportional to the variance reduction achieved by the bias correction.

- The efficiency gain of LC-ACP is most valuable precisely when it is most needed - during periods of high uncertainty. This makes LC-ACP particularly useful for risk management applications, where accurate uncertainty quantification during market stress is crucial.

### 5.6.4  Coverage vs. Width Trade-off

An important aspect of any conformal prediction method is the trade-off between coverage and interval width. Figure 5.6 illustrates this relationship for each method-ticker combination. Figure 5.6 demonstrates that:

- LC-ACP consistently achieves better performance across all tickers, with points closer to the upper-left region (higher coverage with narrower intervals).

Table 5.5: Average interval scores (IS) and Winkler scores (WS) for each method (multiplied by $10^2$ for readability). Lower values indicate better performance.

| Method | IS | WS | MACE | IS $\times$ MACE |
|---|---|---|---|---|
| CP | 4.87 | 4.62 | 0.053 | 0.258 |
| ACP | 4.55 | 4.39 | 0.031 | 0.141 |
| LC-ACP | **3.98** | **3.87** | **0.017** | **0.068** |

- Different tickers respond differently to each method, but LC-ACP provides the most consistent performance across all sectors.

- The ACP method shows improved performance over CP but does not match the efficiency of LC-ACP.

To quantify this trade-off, we computed the interval score (IS) and Winkler score (WS) for each method. These proper scoring rules simultaneously reward accurate coverage and narrow intervals.

The results in Table 5.5 show that LC-ACP achieves the best scores across all metrics, indicating that it optimizes the coverage-width trade-off more effectively than the other methods. The product of the interval score and mean absolute calibration error (IS $\times$ MACE) serves as a comprehensive measure of performance, combining both efficiency and calibration quality. LC-ACP improves this combined metric by 52% compared to ACP and 74% compared to CP.

### 5.6.5 Effect of Local Bias Correction

To illustrate the mechanism behind LC-ACP's improved performance, Figure 5.7 visualizes the effect of local bias correction during a period of rapid volatility increase for TSLA stock. This visualization reveals:

- Without bias correction (top panel), ACP consistently underestimates the true volatility during the rapid increase phase, resulting in intervals that fail to contain the true values.

Table 5.6: Decomposition of prediction errors (MSE) into bias and variance components (values multiplied by $10^4$ for readability).

| Method | Bias$^2$ | Variance | Total MSE |
|---|---|---|---|
| GARCH | 2.53 | 3.86 | 6.39 |
| LC-ACP | 0.68 | 3.79 | 4.47 |
| % Reduction | 73.1% | 1.8% | 30.0% |

This is a classic example of local bias, where the model systematically underestimates volatility during regime transitions.

- With local bias correction (bottom panel), LC-ACP quickly identifies and corrects for the systematic bias, resulting in intervals that properly contain the true volatility values even during rapid shifts. The bias correction term $\delta(x)$ effectively captures the local pattern of underestimation and adjusts the predictions accordingly.

- The adaptation happens more quickly with LC-ACP, demonstrating its ability to respond to localized patterns without requiring a complete recalibration of the model. This is particularly valuable during regime transitions, where traditional methods may take time to adapt.

**Decomposition of Prediction Errors**

To better understand the sources of prediction errors and the impact of local bias correction, we decomposed the total error into bias and variance components for each method.

This analysis shows that local bias correction primarily reduces the systematic bias component of the error (73.1% reduction) while having a minimal impact on the variance component. This is consistent with our theoretical analysis, which showed that LC-ACP achieves efficiency gains by reducing local biases while maintaining the same level of coverage.

The substantial reduction in the squared bias component (73.1%) translates to a more modest reduction in interval width (17.0%) because the interval width is proportional to the square root of the total variance, which includes both the bias and variance components.

This relationship is precisely what our theoretical analysis in Theorem 5.5.1 predicted.

## 5.6.6 Sensitivity Analysis

To assess the robustness of LC-ACP to hyperparameter choices, we conducted a comprehensive sensitivity analysis by varying the key parameters and measuring their impact on performance.

Figure 5.8 reveals several important insights:

- **Window size effects:**

    - Small windows (7-14 days) react quickly to market changes but exhibit higher coverage error variance and wider intervals due to estimation noise.

    - Medium windows (21-28 days) provide optimal balance between reactivity and stability, showing minimum coverage error.

    - Large windows (42-63 days) reduce error variance but respond too slowly to regime changes, especially during the rapid transitions observed in late 2022 and mid-2023.

- **Step size ($\gamma$) effects:**

    - Small step sizes (0.01-0.05) result in slow adaptation that fails to keep pace with rapid market transitions, particularly during the volatile periods in 2022.

    - Medium step sizes (0.1-0.15) provide optimal balance, allowing the model to adapt quickly enough while maintaining stability.

    - Large step sizes (0.2-0.3) cause excessive oscillation in coverage rates, particularly evident during the relatively stable market periods in mid-2022.

These findings align with our theoretical analysis, which showed that the step size $\gamma$ controls the trade-off between adaptation speed and stability. The optimal step size depends on the magnitude of distribution shifts and the desired balance between short-term adaptation and long-term stability.

Table 5.7: Average computation time per prediction (in milliseconds).

| Method | Training | Prediction |
|--------|----------|------------|
| GARCH  | 187.3    | 2.4        |
| CP     | 187.3    | 7.2        |
| ACP    | 187.3    | 8.5        |
| LC-ACP | 213.6    | 14.3       |

### 5.6.7 Computational Efficiency

We also evaluated the computational efficiency of each method, which is important for practical deployment in real-time financial applications.

While LC-ACP requires additional computation for the local bias correction, the overhead is modest and well within the requirements for practical financial applications, where predictions are typically needed at daily or longer intervals. The additional computational cost is justified by the significant improvements in prediction interval quality, as demonstrated by the coverage and width analyses.

### 5.6.8 Summary of Experimental Findings

Our comprehensive experimental evaluation across four diverse stocks and a four-year period spanning multiple market cycles yields the following key findings:

1. **Superior coverage accuracy**: LC-ACP achieves average coverage of 89.9%, closest to the target of 90% among all methods tested, with the smallest deviation from the target during high-volatility periods.

2. **Significant interval width reduction**: LC-ACP produces intervals that are 17.0% narrower than CP and 10.6% narrower than ACP on average, with even greater reductions during high-volatility regimes.

3. **Optimal coverage-width trade-off**: LC-ACP achieves the best interval and Winkler scores across all market conditions, indicating that it balances coverage accuracy and

interval width more effectively than other methods.

4. **Robust performance across regimes**: LC-ACP maintains consistent performance across low, normal, and high volatility regimes, with particularly notable improvements during high-volatility periods when accurate uncertainty quantification is most valuable.

5. **Rapid adaptation to shifts**: LC-ACP adapts quickly to volatility regime transitions, demonstrating superior effectiveness during the post-pandemic recovery period, high inflation environment, and tech sector correction.

6. **Error reduction**: Local bias correction reduces the squared bias component of prediction errors by 73.1%, confirming the theoretical mechanism by which LC-ACP improves upon standard methods.

These experimental results provide strong empirical validation of our theoretical analysis, demonstrating that LC-ACP effectively addresses the limitations of standard conformal methods by combining local bias correction with adaptive calibration. The method's ability to maintain target coverage while producing efficient intervals across diverse market regimes suggests that its theoretical advantages translate directly into practical benefits for uncertainty quantification in non-stationary environments.

## 5.7   Conclusion

In this chapter, we introduced Locally Debiased Adaptive Conformal Prediction (LC-ACP), a novel approach that combines local bias correction with adaptive conformal prediction to address the challenges of uncertainty quantification in non-stationary time series. Our method bridges theory and practice by maintaining distribution-free coverage guarantees while significantly improving interval efficiency through explicit local bias correction.

Our theoretical analysis established that LC-ACP provides valid coverage guarantees even without the exchangeability assumption, making it particularly well-suited for time

series applications where temporal dependencies and distribution shifts are common. We also quantified the efficiency gain from local bias correction, showing that the width reduction is proportional to the square root of the variance reduction achieved by the bias correction.

Our comprehensive empirical evaluation demonstrated that LC-ACP consistently outperforms existing conformal prediction methods in both coverage accuracy and interval width. The method shows particular strength during volatile periods and regime transitions, precisely when accurate uncertainty quantification is most valuable. By achieving average coverage of 89.9% (closest to the target 90%) while producing intervals that are 17.0% narrower than CP and 10.6% narrower than ACP, LC-ACP demonstrates that significant practical improvements are possible without sacrificing theoretical guarantees.

LC-ACP opens up new possibilities for more efficient, accurate, and reliable uncertainty quantification in machine learning applications. By explicitly addressing local biases while maintaining distribution-free coverage guarantees, it bridges an important gap in the conformal prediction literature and provides a practical tool for real-world applications across multiple domains. As we continue to deploy machine learning models in critical applications, methods like LC-ACP that enhance the trustworthiness of predictions through reliable uncertainty estimates will become increasingly important for responsible AI deployment.

In the next chapter, we will build upon the theoretical foundation established here by applying LC-ACP to the challenging domain of climate and weather prediction, with a specific focus on hurricane forecasting. This application domain presents unique challenges including non-stationarity, extreme events, and high-dimensional spatial-temporal data—making it an ideal test case for demonstrating the real-world utility of our LC-ACP methodology. We will show how LC-ACP can help address the critical need for reliable uncertainty quantification in climate foundation models, ultimately contributing to more robust decision-making tools for disaster preparedness and response.

Figure 5.1: Comparison of Gaussian Process regression and Locally Debiased Adaptive Conformal Prediction (LC-ACP) on synthetic stock price data with three distinct volatility regimes. The figure demonstrates how GPs maintain consistent uncertainty patterns across all regimes, leading to poor coverage (39.5% overall) compared to the target 90%, with particularly low coverage in the medium volatility regime (14.5%). In contrast, LC-ACP adapts its prediction intervals to account for regime-specific characteristics, achieving more reliable and consistent coverage (99.5% overall) across all regimes at the cost of slightly wider prediction intervals.

# Comparison of Conformal Prediction Methods



**Key Benefits of LC-ACP:** *Combines adaptivity from ACP with local bias correction to provide more accurate and efficient prediction intervals. Particularly effective in non-stationary environments and with heterogeneous data distributions.*

Figure 5.2: Comparison of conformal prediction methods. LC-ACP combines local bias correction with adaptive conformal prediction, providing narrower, better-centered prediction intervals in regions where the base model exhibits systematic errors. The left panel shows standard CP with fixed-width intervals, the middle panel shows ACP with globally adaptive intervals, and the right panel shows LC-ACP with locally adapted intervals that follow the true function more closely.

## Theoretical Illustration of Local Bias Correction in Conformal Prediction



Figure 5.3: Illustration of LC-ACP's theoretical benefits. The base model $\hat{f}(x)$ (blue) fails to capture local patterns in the true function $f(x)$ (dashed black). Standard CP uses fixed-width intervals, while LC-ACP leverages the bias correction term $\delta(x)$ to produce narrower intervals (green) that still maintain coverage guarantees. The width reduction is proportional to $\sqrt{1-\eta}$ where $\eta$ is the variance reduction factor.

**Coverage Rate Over Time by Ticker**



Note: Rolling coverage calculated with 253-day window

Figure 5.4: Coverage rates over time for four stocks: Apple Inc. (AAPL), NVIDIA Corporation (NVDA), Tesla, Inc. (TSLA), and Alibaba Group (BABA) from January 2021 to January 2025. LC-ACP maintains more stable coverage closer to the target of 90% compared to CP and ACP, especially during periods of market volatility. The rolling coverage is calculated using a 253-day window.

**Comparison of Prediction Interval Widths**



*Lower values indicate narrower prediction intervals,*
*which are generally preferred when coverage is maintained.*

Figure 5.5: Boxplot of prediction interval widths by method across all four stocks. LC-ACP produces narrower intervals while maintaining target coverage. This demonstrates LC-ACP's efficiency in uncertainty quantification, achieving narrower intervals while maintaining coverage closer to the target 90%.

Figure 5.6: Scatter plot of coverage versus interval width for each method and ticker. Each point represents a ticker, labeled by its symbol. LC-ACP consistently achieves better performance across all tickers, with points closer to the upper-left region (higher coverage with narrower intervals). The horizontal dashed line indicates the target coverage of 90%.

**Effect of Local Bias Correction in LC-ACP**

*The base model (blue) fails to capture high-frequency patterns in the true function (dashed black).*
*LC-ACP uses bias correction to produce narrower intervals (green) that better track the true function.*

Figure 5.7: Comparison of prediction intervals with and without local bias correction during a period of rapidly increasing volatility for TSLA stock. The top panel shows ACP, while the bottom panel shows LC-ACP with local bias correction. The shaded region highlights a period of rapid volatility increase. Note how LC-ACP's intervals better track the true volatility (black line) by adjusting for local biases.



(a) Effect of local window size

(b) Effect of adaptation step size $\gamma$

Figure 5.8: Sensitivity analysis of LC-ACP hyperparameters. (a) Shows how the local window size affects the trade-off between coverage error and interval width. The U-shaped pattern indicates an optimal window size of approximately 21-28 days. (b) Shows how the adaptation step size $\gamma$ impacts the convergence rate and stability of coverage, with optimal values between 0.1-0.15.

# Chapter 6

# Application of LC-ACP to Climate and Weather Foundation Models

> "In a chaotic world, the relationship between model predictions and reality demands not just accuracy, but well-calibrated uncertainty."

## 6.1 Introduction

Weather and climate phenomena represent some of the most challenging prediction problems in modern science, characterized by chaotic dynamics, complex non-stationarity, and high-dimensional spatiotemporal structures. This is particularly true of extreme events, including tropical cyclones (hurricanes or typhoons), which rank among the most destructive natural disasters worldwide with devastating socioeconomic impacts. According to the National Oceanic and Atmospheric Administration (NOAA), between 1980 and 2023, tropical cyclones accounted for over $1.3 trillion in damages in the United States alone, with an average cost per event exceeding $22 billion [68].

The advancement of deep learning models, and foundation models in particular, has introduced powerful new tools for many applications, including weather and climate prediction. Foundation models pre-trained on large-scale reanalysis or climate simulation data have set new benchmarks in predictive accuracy [72, 82]. Simultaneously, lightweight archi-

tectures such as Tiny Time Mixers (TTMs) [22] offer parameter-efficient alternatives that maintain competitive performance while reducing computational costs—a critical advantage for operational deployment.

Despite these advances in deterministic prediction, reliable uncertainty quantification (UQ) remains a significant challenge. The necessity for robust UQ is particularly acute in hurricane forecasting, where risk communication directly impacts evacuation decisions, resource allocation, and ultimately, human lives [75].

In Chapter 2, we introduced Gaussian Processes as a principled Bayesian framework for uncertainty quantification. While GPs provide excellent uncertainty quantification for small-to-medium scale stationary data, they face significant limitations when applied to large-scale, non-stationary time series data like hurricane movement, i.e. tracks. Specifically, they

1. scale poorly with dataset size (cubic complexity), making them impractical for large climate datasets;

2. struggle with non-stationary time series data, requiring complex kernel engineering;

3. typically assume Gaussian noise, which may not capture the complex error distributions in hurricane predictions.

To address these limitations, in Chapter 5, we developed Locally Debiased Adaptive Conformal Prediction (LC-ACP), which combines local bias correction with adaptive conformal prediction to provide distribution-free coverage guarantees while significantly improving prediction interval efficiency. The method showed particular strengths in financial forecasting applications with volatile periods and regime transitions—characteristics that are similarly present in problems around weather and climate prediction.

This chapter builds upon that methodological foundation by applying LC-ACP to climate foundation models, specifically for hurricane track prediction. We demonstrate how LC-ACP addresses the limitations of Gaussian Processes while providing rigorous uncertainty

quantification for operational hurricane forecasting. Our approach combines the computational efficiency of lightweight foundation models (TTMs) with the statistical rigor of LC-ACP to deliver well-calibrated uncertainty estimates that adapt to both local model biases and distribution shifts over time—a critical capability for weather and climate applications where non-stationarity is the norm rather than the exception.

## 6.2 Challenges in Hurricane Track Prediction and Uncertainty Quantification

### 6.2.1 The Hurricane Prediction Problem

Hurricane forecasting represents a quintessential example of the challenges in high-dimensional, non-stationary time series prediction. Let $X_t \in \mathbb{R}^d$ denote the state vector of a hurricane at time $t$, where $d$ includes variables such as latitude, longitude, wind speed, pressure, and other atmospheric covariates. Given a historical context window of $C$ consecutive timesteps, our goal is to forecast the hurricane's future states over a prediction horizon of $H$ timesteps:

$$\hat{X}_{t+1:t+H} = f_\theta\big(X_{t-C+1:t}\big) \tag{6.1}$$

where $f_\theta$ represents our forecasting model with parameters $\theta$. For notational simplicity, we focus on a specific variable of interest (e.g., storm track coordinates) and denote:

- $y_t \in \mathbb{R}^m$ as the ground truth value (location) at time $t$

- $\hat{y}_t \in \mathbb{R}^m$ as the corresponding model-predicted value

- $r_t = y_t - \hat{y}_t$ as the prediction residual error

Several factors make hurricane track prediction particularly challenging:

1. **Chaotic Dynamics**: Small perturbations in initial conditions can lead to dramatically different forecast outcomes, a classic example of the "butterfly effect" [60].

2. **Complex Physics**: Hurricane evolution involves multi-scale interactions between atmospheric and oceanic processes, including thermodynamic energy exchanges, boundary layer dynamics, and environmental steering flows.

3. **Non-Stationarity**: Hurricane characteristics evolve over time due to interactions with changing environmental conditions, seasonal effects, and climate change, violating the stationarity assumptions of many statistical models.

4. **High-Dimensional State Space**: The complete state description of a hurricane and its environment involves hundreds to thousands of variables across multiple spatial locations and vertical levels.

5. **Limited Observational Data**: While climate models generate abundant synthetic data, real-world hurricane observations remain relatively sparse, especially for extreme events, creating challenges for model validation and uncertainty quantification.

## 6.2.2   Traditional Approaches to Hurricane Track Uncertainty

Traditional approaches to representing uncertainty in hurricane forecasts include:

1. **Ensemble Forecasting**: Running multiple deterministic forecasts with perturbed initial conditions or model physics to generate a distribution of possible outcomes. The European Centre for Medium-Range Weather Forecasts (ECMWF) ensemble and the NOAA Global Ensemble Forecast System (GEFS) exemplify this approach.

2. **Statistical Post-Processing**: Applying statistical corrections to deterministic forecasts based on historical error patterns. The National Hurricane Center's (NHC) "cone of uncertainty" represents a widely-used example, constructed from historical forecast errors over the previous five years.

3. **Bayesian Methods**: Incorporating prior knowledge and observation uncertainties into a Bayesian framework to generate posterior predictive distributions. The computational

complexity of these methods often makes them impractical for operational forecasting.

Each of these approaches has significant limitations:

- Ensemble forecasting is computationally expensive and may not capture all sources of uncertainty, particularly those related to model structural errors.

- Statistical post-processing based on historical errors does not adapt well to changing conditions or storm-specific characteristics, potentially leading to miscalibration.

- Bayesian methods face computational challenges and often require simplifying assumptions about error distributions that may not hold in practice.

## 6.2.3 The Need for Adaptive, Distribution-Free Uncertainty Quantification

The limitations of traditional approaches highlight the need for uncertainty quantification methods that are:

1. **Computationally Efficient**: Capable of generating uncertainty estimates quickly enough for operational use, even with limited computational resources.

2. **Distribution-Free**: Not reliant on strong assumptions about error distributions, which may be complex and non-stationary in hurricane forecasting.

3. **Locally Adaptive**: Able to adjust uncertainty estimates based on local characteristics of each forecast point, recognizing that prediction difficulty varies across different storm types and geographical regions.

4. **Temporally Adaptive**: Capable of maintaining calibration as the underlying data distribution shifts over time, without requiring frequent recalibration.

5. **Theoretically Grounded**: Providing rigorous coverage guarantees under minimal assumptions, ensuring reliable uncertainty quantification for high-stakes decision-making.

Locally Debiased Adaptive Conformal Prediction (LC-ACP), introduced in Chapter 5, satisfies these requirements by combining local bias correction with adaptive conformal calibration. The remainder of this chapter demonstrates the application of LC-ACP to hurricane track prediction using foundation models, specifically lightweight Tiny Time Mixers (TTMs).

## 6.3 Methodology: Integrating TTMs with LC-ACP for Hurricane Uncertainty Quantification

### 6.3.1 Overview of the Integrated Approach

Our approach integrates lightweight foundation models with statistical uncertainty quantification to deliver accurate and computationally efficient hurricane forecasts with reliable uncertainty estimates. The pipeline consists of two core components:

1. **Deterministic Forecasting via TTMs**: A parameter-efficient Tiny Time Mixer model generates initial baseline forecasts.

2. **Uncertainty Quantification with LC-ACP**: Locally Debiased Adaptive Conformal Prediction provides well-calibrated uncertainty estimates by combining local bias correction with adaptive conformal prediction.

This modular design leverages the complementary strengths of each component: TTMs provide fast approximate predictions, while LC-ACP adds a principled uncertainty quantification layer that addresses local biases in the model and adapts to distribution shifts over time.

Figure 6.1 illustrates the LC-ACP methodology for hurricane track forecasting. The pipeline takes the current hurricane state as input, generates a TTM-based prediction, applies local bias correction to address systematic errors, and finally performs adaptive conformal calibration to ensure valid coverage guarantees under distribution shifts.

Figure 6.1: Illustration of the Locally Debiased Adaptive Conformal Prediction (LC-ACP) methodology. The pipeline consists of three main components: (1) TTM prediction, which generates the initial forecast; (2) Local bias correction, which adjusts for systematic errors in specific regions; and (3) Adaptive conformal calibration, which ensures valid coverage guarantees. The resulting prediction intervals adapt to both the local characteristics of each forecast point and the evolving distribution of hurricane tracks.

## 6.3.2   Tiny Time Mixer (TTM) Architecture for Hurricane Forecasting

Tiny Time Mixers represent a parameter-efficient architecture for time-series modeling, combining the strengths of convolutional networks for local pattern detection with attention mechanisms for global dependencies [22]. Unlike popular time-series foundation models that focus on reconstructing masked data, TTMs employ a prediction approach through efficient channel-mixing layers.

The core TTM block consists of alternating time-mixing and feature-mixing operations:

$$\text{Time-mixing:} \quad \boldsymbol{Z}^{(l)} = \text{Mix}_{\text{time}}\big(\boldsymbol{X}^{(l-1)}\big) \tag{6.2}$$

$$\text{Feature-mixing:} \quad \boldsymbol{X}^{(l)} = \text{Mix}_{\text{feature}}\big(\boldsymbol{Z}^{(l)}\big) \tag{6.3}$$

For hurricane forecasting, we leverage a TTM variant with the following key features:

- **Multi-Resolution Processing**: Captures both short-term dynamics (rapid intensity changes) and long-term trends (overall track trajectory) through multi-scale patching.

- **Parameter Efficiency**: Achieves competitive performance with as little as 1M parameters, compared to hundreds of millions in larger foundation models, making it suitable for operational deployment.

- **Transfer Learning Capability**: Pre-trained on large climate simulation datasets and fine-tuned on empirical or representative historical hurricane data, allowing effective generalization from synthetic to real-world storms.

- **Latent Space Representation**: Provides meaningful intermediate representations that capture complex spatio-temporal patterns in hurricane dynamics, forming the basis for local bias correction.

The TTM is trained using multiple loss functions targeting different aspects of hurricane forecasting:

$$\mathcal{L}_{\text{TTM}}(\theta) = \lambda_{\text{track}}\mathcal{L}_{\text{track}}(\theta) + \lambda_{\text{intensity}}\mathcal{L}_{\text{intensity}}(\theta) + \lambda_{\text{reg}}\mathcal{L}_{\text{reg}}(\theta) \tag{6.4}$$

where $\lambda_{\text{track}}$, $\lambda_{\text{intensity}}$, and $\lambda_{\text{reg}}$ are hyperparameters balancing different loss terms.

## 6.3.3 Locally Debiased Adaptive Conformal Prediction (LC-ACP) for Hurricane Tracks

**LC-ACP Algorithm**

Algorithm 7 and 8 collectively provide the complete implementation of LC-ACP for hurricane track forecasting. The algorithm consists of three main phases: (1) learning a local bias correction function using an anchor dataset, (2) calculating initial nonconformity scores and determining quantile thresholds on a calibration set, and (3) adaptive online prediction with dynamic updates to maintain the target coverage level under distribution shifts.

**Local Bias Correction in Latent Space**

Rather than operating in the raw input space, we apply LC-ACP in a lower-dimensional latent space derived from the TTM. Let $\mathbf{h}_t = \phi_\theta(X_{t-C+1:t})$ be a latent representation extracted from an intermediate layer of the TTM. This approach offers several advantages:

- Dimensionality reduction for computational efficiency

- Feature learning that captures complex spatio-temporal patterns

- Smoother behavior for better generalization of uncertainty estimates

For each new prediction point $\mathbf{h}_*$, we identify similar historical examples and use their errors to construct a local error distribution:

$$\mathcal{E}(\mathbf{h}_*) = \{r_i \mid \mathbf{h}_i \in \mathcal{N}_k(\mathbf{h}_*)\} \tag{6.5}$$

where $\mathcal{N}_k(\mathbf{h}_*)$ represents the $k$ nearest neighbors of $\mathbf{h}_*$ in the latent space.

The local bias correction function $\delta(\mathbf{h})$ estimates the systematic error at a given point:

$$\delta(\mathbf{h}_*) = \frac{\sum_{i=1}^{k} w_i \cdot r_i}{\sum_{i=1}^{k} w_i}, \tag{6.6}$$

where $w_i = \exp(-\|\mathbf{h}_* - \mathbf{h}_i\|^2/h^2)$ is a kernel weight that gives higher importance to closer points, and $h$ is a bandwidth parameter controlling the locality of the correction.

## Local Uncertainty Estimation

The local uncertainty estimate $\sigma(\mathbf{h}_*)$ is computed as the weighted standard deviation of the residuals:

$$\sigma(\mathbf{h}_*) = \sqrt{\frac{\sum_{i=1}^{k} w_i \cdot (r_i - \delta(\mathbf{h}_*))^2}{\sum_{i=1}^{k} w_i}} \tag{6.7}$$

This uncertainty estimate captures the local variability in prediction errors, providing a basis for constructing prediction intervals that adapt to the specific characteristics of each forecast point.

The locally debiased forecast is then given by:

$$\tilde{y}_* = \hat{y}_* + \delta(\mathbf{h}_*), \tag{6.8}$$

with associated uncertainty $\sigma(\mathbf{h}_*)$.

## Adaptive Conformal Calibration

For each calibration example, we define a nonconformity score that measures how "unusual" the observed residual is compared to the locally debiased prediction:

$$S_t = \frac{\left| y_t - (\hat{y}_t + \delta(\mathbf{h}_t)) \right|}{\sigma(\mathbf{h}_t)}. \tag{6.9}$$

Using a held-out calibration set, we compute nonconformity scores and determine the empirical $(1 - \alpha)$-quantile:

$$Q_{1-\alpha} = \text{Quantile}\left(\{S_t\}_{t \in \mathcal{D}_{\text{cal}}}, 1 - \alpha\right) \tag{6.10}$$

For a new test point, we construct the prediction interval by:

$$PI(t_*) = \left[\tilde{y}_* - Q_{1-\alpha}\,\sigma(\mathbf{h}_*),\; \tilde{y}_* + Q_{1-\alpha}\,\sigma(\mathbf{h}_*)\right] \tag{6.11}$$

A key innovation in ACP is dynamic updating of the quantile threshold:

$$Q_{1-\alpha,t+1} = Q_{1-\alpha,t} + \gamma(\alpha - \mathbb{I}\{y_t \notin PI(t)\}), \tag{6.12}$$

where $\mathbb{I}\{y_t \notin PI(t)\}$ is the miscoverage indicator and $\gamma$ is a learning rate parameter. This adaptive mechanism allows LC-ACP to maintain target coverage even as hurricane characteristics evolve.

Figure 6.2 illustrates the integrated TTM+LC-ACP pipeline. The TTM takes the historical hurricane track as input and generates a deterministic forecast. LC-ACP then processes this forecast through local bias correction and adaptive conformal calibration to produce well-calibrated prediction intervals.

## 6.4 Experimental Setup and Data Sources

### 6.4.1 Climate Model Synthetic Hurricane Tracks

To evaluate our approach, we used tropical cyclones extracted from high-resolution climate model simulations [23]. Specifically, we used the CanESM5 model (Canadian Earth System Model) from the CMIP6 (Coupled Model Intercomparison Project Phase 6) [93]. From these simulations, we identified tropical cyclones using a standard tracking algorithm based on vorticity maxima, warm core structure, and minimum duration criteria.

The primary dataset for our experiments consisted of 1,399 simulated hurricane tracks from the CanESM5 model that passed through the Galveston county region. These tracks represent diverse hurricane scenarios under current and near-future climate conditions, providing a robust basis for training and evaluating our models.

Figure 6.3 illustrates the comprehensive dataset of simulated hurricane tracks used in our experiments. The tracks exhibit diverse patterns, including direct approaches to Galveston, recurving tracks that parallel the coast, and tracks that cross the Galveston area from various directions. This diversity enables the evaluation of our method across a wide range of hurricane scenarios.

## 6.4.2  Experimental Protocol

We employed a leave-three-out cross-validation strategy to evaluate our methodology:

1. The TTM model was trained on 1,396 simulated tracks, with 3 tracks held out for testing.

2. This process was repeated multiple times with different held-out sets to ensure robust evaluation.

3. For LC-ACP, we used a portion of the training set as a calibration set to determine the initial conformal quantiles.

This protocol ensures a realistic evaluation setting while maximizing the use of available data. The held-out tracks were selected to represent different storm characteristics, allowing us to evaluate the method's performance across diverse hurricane scenarios.

## 6.4.3  Evaluation Metrics

We evaluated our approach using complementary metrics that assess both track accuracy and uncertainty calibration:

- **Track Accuracy Metrics:**

  - Track Root Mean Square Error (RMSE): the Haversine distance between predicted and actual storm centers

– Direct Position Error (DPE): the distance error at specific forecast lead times (24h, 48h, 72h, 96h, 120h)

- **Uncertainty Quantification Metrics:**

  – Empirical Coverage: the percentage of test cases where the true position fell within the predicted uncertainty region

  – Mean Interval Width: the average radius (in kilometers) of the prediction intervals

  – Calibration Error: the absolute difference between the nominal coverage (90%) and the empirical coverage

These metrics collectively evaluate both the accuracy of the deterministic predictions and the quality of the uncertainty estimates.

## 6.5    Experimental Results

### 6.5.1    Experimental Setup

We evaluated our TTM+LC-ACP pipeline on the 1,399 simulated hurricane tracks from the CanESM5 model. The experimental setup followed a leave-three-out cross-validation approach, where we trained the TTM model on 1,396 tracks and tested on 3 held-out tracks. This process was repeated multiple times with different held-out sets to ensure robust evaluation.

### 6.5.2    CanESM Galveston Tracks Experiment

For the Galveston region evaluation, we conducted a focused experiment using the extracted tracks from the CanESM5 climate model:

- **Training Dataset:** We used 1,396 of the 1,399 simulated tropical cyclone tracks from the CanESM5 model that passed through the Galveston county region (shown in Figure

6.3). These tracks represent diverse potential hurricane scenarios under current and near-future climate conditions.

- **Evaluation Approach:** Our leave-three-out cross-validation strategy ensured realistic evaluation while maximizing the use of available data. The held-out tracks were selected to represent different storm characteristics.

- **LC-ACP Application:** For each test track, we applied the LC-ACP methodology following Algorithms 7 and 8 to quantify uncertainty in the TTM predictions. The local bias correction component was particularly valuable in capturing systematic biases in the model's predictions for the specific Galveston region.

Figure 6.4 illustrates the application of LC-ACP to three held-out test tracks. The uncertainty cones adapt to the specific characteristics of each track while maintaining the target coverage rate. This visualization demonstrates how our method provides dynamic, track-specific uncertainty quantification that adapts to the evolving characteristics of each hurricane.

## 6.5.3   Quantitative Results and Method Comparison

Our experimental evaluation demonstrated several key findings:

- **Predictive Accuracy:** The TTM model achieved strong predictive performance across different forecast lead times, with track error increasing naturally at longer horizons (from approximately 32km at 24-hour lead time to 177km at 120-hour lead time).

- **Coverage Reliability:** The LC-ACP framework maintained consistent empirical coverage rates around the target 90% across all lead times, demonstrating the effectiveness of the adaptive calibration component despite increasing prediction difficulty.

- **Interval Efficiency:** The local bias correction component of LC-ACP produced

prediction intervals that were 15-25% narrower than those from standard conformal prediction methods while maintaining the same coverage rate.

- **Storm Category Performance:** LC-ACP maintained consistent performance across different hurricane intensity categories (tropical storms, Category 1-2, and major hurricanes), while alternative methods showed more significant variations in either coverage or interval width.

When compared to alternative approaches (standard conformal prediction and adaptive conformal prediction without local bias correction), LC-ACP consistently produced better-centered prediction intervals with lower calibration error. This improvement was particularly pronounced for major hurricanes (Categories 3-5), where LC-ACP maintained target coverage while reducing interval width by approximately 24.6% compared to standard methods.

## 6.6 Operational Considerations and Limitations

While our LC-ACP approach demonstrates significant advantages for hurricane uncertainty quantification, several operational considerations and limitations should be acknowledged:

### 6.6.1 Computational Efficiency

The computational requirements of our approach are modest compared to traditional ensemble methods:

- The TTM model requires approximately 1GB of memory and can generate forecasts in under 2 seconds on standard CPU hardware.

- The LC-ACP computation adds only 0.3-0.5 seconds per forecast, with most of the time spent on nearest neighbor search.

- The complete pipeline is at least 50x faster than running a traditional ensemble of NWP models.

This efficiency makes our approach suitable for operational deployment in resource-constrained environments, such as regional emergency management agencies or developing countries with limited computational resources.

## 6.6.2 Data Requirements

Our approach requires:

- a sufficient number of historical or simulated hurricane tracks for training the TTM model (typically >1,000 tracks)

- a calibration dataset for initializing the conformal prediction component (typically 100-200 tracks)

- an anchor dataset for local bias correction (typically the most recent 20-50 tracks)

The use of synthetic tracks from climate models helps address the limitation of sparse historical data, especially for regions with fewer observed hurricanes.

## 6.6.3 Limitations and Future Work

Several limitations of our current approach suggest directions for future research:

- **Counterfactual Analysis:** Our method does not explicitly account for counterfactual scenarios that are physically plausible but not represented in the training datan

- **Multi-Model Integration:** Future work could explore integrating predictions from multiple foundation models, potentially using a weighted ensemble approach based on local performance.

- **Beyond Track Prediction:** Extending the approach to joint uncertainty quantification of track, intensity, size, and rainfall would provide a more comprehensive risk assessment framework.

- **Human-AI Collaboration:** Developing interfaces that effectively communicate the uncertainty information to forecasters and decision-makers remains an important area for future work.

## 6.7 Conclusion

In this chapter, we presented a novel framework for hurricane track uncertainty quantification that combines lightweight foundation models (TTMs) with Locally Debiased Adaptive Conformal Prediction (LC-ACP). Our approach addresses the key limitations of traditional uncertainty quantification methods in hurricane forecasting:

1. It provides computationally efficient forecasts with well-calibrated uncertainty estimates, making it suitable for operational deployment.

2. It maintains valid coverage guarantees across different forecast lead times and storm intensities without making strong distributional assumptions.

3. It adaptively corrects for local biases in the model predictions, producing narrower, better-centered prediction intervals.

4. It adjusts to distribution shifts over time, maintaining calibration even as hurricane characteristics evolve.

The comprehensive evaluation on simulated hurricane tracks demonstrated the effectiveness of our approach, with LC-ACP consistently outperforming alternative methods in terms of both coverage accuracy and interval width. The detailed case studies further illustrated the practical utility of our framework for high-stakes decision-making scenarios like hurricane landfall prediction.

This work bridges the gap between advanced AI foundation models and rigorous uncertainty quantification, providing a pathway for deploying these models in critical applications

where reliable uncertainty estimates are essential. By combining the computational efficiency of TTMs with the statistical rigor of LC-ACP, we enable more effective risk communication and decision support for hurricane emergency management.

The methodology presented in this chapter extends beyond hurricane forecasting to other weather and climate applications where foundation models are increasingly being deployed. As these models continue to advance in predictive accuracy, complementary advances in uncertainty quantification like LC-ACP will be crucial for their responsible deployment in high-stakes domains where decisions have significant societal impacts.

---

**Algorithm 7** Locally Debiased Adaptive Conformal Prediction (LC-ACP) for Hurricane Tracks: Part 1

---

**Require::**
  1: Training data $\mathcal{D}_{\text{train}} = \{(X_{t-C+1:t}, y_{t+1:t+H})_i\}_{i=1}^{N_{\text{train}}}$
  2: Calibration data $\mathcal{D}_{\text{cal}} = \{(X_{t-C+1:t}, y_{t+1:t+H})_i\}_{i=1}^{N_{\text{cal}}}$
  3: Anchor data $\mathcal{D}_{\text{anchor}} = \{(X_{t-C+1:t}, y_{t+1:t+H}, \hat{y}_{t+1:t+H})_i\}_{i=1}^{N_{\text{anchor}}}$
  4: Target coverage level $1 - \alpha_{\text{target}}$ (e.g., 0.9)
  5: Number of nearest neighbors $k$ for local bias estimation
  6: Kernel bandwidth parameter $h$
  7: Learning rates $\gamma, \eta$ for quantile and $\alpha$ adaptation
  8: Bounds $[\alpha_{\text{min}}, \alpha_{\text{max}}]$ for miscoverage rate
**Ensure::** Well-calibrated prediction intervals with target coverage $1 - \alpha_{\text{target}}$

  9: **Phase 1: Learning Local Bias Correction Function**
 10: Train TTM model $f_\theta$ on $\mathcal{D}_{\text{train}}$
 11: **FOR** each $(X_{t-C+1:t}, y_{t+1:t+H})_i \in \mathcal{D}_{\text{anchor}}$ **DO**
 12:     Generate prediction $\hat{y}_{t+1:t+H,i} = f_\theta(X_{t-C+1:t,i})$
 13:     Extract latent representation $\mathbf{h}_i = \phi_\theta(X_{t-C+1:t,i})$
 14:     Compute residual $r_i = y_{t+1:t+H,i} - \hat{y}_{t+1:t+H,i}$
 15:     Store tuple $(\mathbf{h}_i, r_i)$ in anchor dataset
 16: **END FOR**

 17: **Phase 2: Calculating Initial Nonconformity Scores**
 18: **FOR** each $(X_{t-C+1:t}, y_{t+1:t+H})_j \in \mathcal{D}_{\text{cal}}$ **DO**
 19:     Generate prediction $\hat{y}_{t+1:t+H,j} = f_\theta(X_{t-C+1:t,j})$
 20:     Extract latent representation $\mathbf{h}_j = \phi_\theta(X_{t-C+1:t,j})$
 21:     Find $k$ nearest neighbors $\mathcal{N}_k(\mathbf{h}_j)$ in anchor dataset
 22:     Compute weights $w_i = \exp(-\|\mathbf{h}_j - \mathbf{h}_i\|^2 / h^2)$ for $\mathbf{h}_i \in \mathcal{N}_k(\mathbf{h}_j)$
 23:     Compute local bias correction $\delta(\mathbf{h}_j) = \frac{\sum_{i=1}^{k} w_i \cdot r_i}{\sum_{i=1}^{k} w_i}$
 24:     Compute local uncertainty $\sigma(\mathbf{h}_j) = \sqrt{\frac{\sum_{i=1}^{k} w_i \cdot (r_i - \delta(\mathbf{h}_j))^2}{\sum_{i=1}^{k} w_i}}$
 25:     Calculate locally debiased prediction $\tilde{y}_j = \hat{y}_j + \delta(\mathbf{h}_j)$
 26:     Compute nonconformity score $S_j = \frac{|y_j - \tilde{y}_j|}{\sigma(\mathbf{h}_j)}$
 27: **END FOR**
 28: Compute initial quantile $Q_{1-\alpha_{\text{target}}} = \text{Quantile}(\{S_j\}_{j=1}^{N_{\text{cal}}}, 1 - \alpha_{\text{target}})$
 29: Initialize $\alpha_1 = \alpha_{\text{target}}$
 30: Initialize coverage window $\mathcal{W} = \{\}$ (empty set)

---

---

**Algorithm 8** Locally Debiased Adaptive Conformal Prediction (LC-ACP) for Hurricane Tracks: Part 2

---

32: **Phase 3: Online Prediction with Adaptive Updates**
33: **FOR** each new test point $X_{t-C+1:t,*}$ at time $t$ **DO**
34:     Generate prediction $\hat{y}_{t+1:t+H,*} = f_\theta(X_{t-C+1:t,*})$
35:     Extract latent representation $\mathbf{h}_* = \phi_\theta(X_{t-C+1:t,*})$
36:     Find $k$ nearest neighbors $\mathcal{N}_k(\mathbf{h}_*)$ in anchor dataset
37:     Compute weights $w_i = \exp(-\|\mathbf{h}_* - \mathbf{h}_i\|^2/h^2)$ for $\mathbf{h}_i \in \mathcal{N}_k(\mathbf{h}_*)$
38:     Compute local bias correction $\delta(\mathbf{h}_*) = \frac{\sum_{i=1}^k w_i \cdot r_i}{\sum_{i=1}^k w_i}$
39:     Compute local uncertainty $\sigma(\mathbf{h}_*) = \sqrt{\frac{\sum_{i=1}^k w_i \cdot (r_i - \delta(\mathbf{h}_*))^2}{\sum_{i=1}^k w_i}}$
40:     Calculate locally debiased prediction $\tilde{y}_* = \hat{y}_* + \delta(\mathbf{h}_*)$
41:     Construct prediction interval $\text{PI}(t_*) = [\tilde{y}_* - Q_{1-\alpha_t}\sigma(\mathbf{h}_*), \tilde{y}_* + Q_{1-\alpha_t}\sigma(\mathbf{h}_*)]$
42:     Return $\text{PI}(t_*)$ as the prediction interval
43:     **IF** ground truth $y_*$ becomes available **THEN**
44:         Update anchor dataset with $(\mathbf{h}_*, y_* - \hat{y}_*)$
45:         Compute miscoverage indicator $I_t = \mathbb{I}\{y_* \notin \text{PI}(t_*)\}$
46:         Add $I_t$ to coverage window $\mathcal{W}$
47:         **IF** $|\mathcal{W}| > $ window\_size **THEN**
48:             Remove oldest element from $\mathcal{W}$
49:         **END IF**
50:         Compute empirical coverage $\hat{C}_t = 1 - \frac{1}{|\mathcal{W}|}\sum_{I \in \mathcal{W}} I$
51:         Update quantile threshold: $Q_{1-\alpha_t,t+1} = Q_{1-\alpha_t,t} + \gamma(\alpha_t - I_t)$
52:         Update miscoverage rate: $\alpha_{t+1} = \alpha_t + \eta \cdot (\hat{C}_t - (1 - \alpha_{\text{target}}))$
53:         Constrain $\alpha_{t+1} \in [\alpha_{\min}, \alpha_{\max}]$
54:     **END IF**
55: **END FOR**

---

Figure 6.2: Integrated pipeline combining Tiny Time Mixers (TTMs) with Locally Debiased Adaptive Conformal Prediction (LC-ACP) for hurricane track forecasting with uncertainty quantification. The TTM generates deterministic forecasts which are then processed by LC-ACP to produce well-calibrated prediction intervals. The local bias correction component adjusts for systematic errors in different regions of the feature space, while the adaptive conformal calibration component ensures valid coverage guarantees under distribution shifts.

**All CanESM Hurricane Tracks (n=1999)**

Figure 6.3: Overview of the 1,399 simulated hurricane tracks from the CanESM5 climate model passing through the Galveston county region. This dataset provides a comprehensive representation of potential hurricane scenarios under current and near-future climate conditions, with diverse track patterns, intensities, and landfall characteristics.



**CANESM Track Predictions - Latitude Only with ACP Conformal Prediction**
**Dashed: True Track, Solid: Predicted Latitude with True Longitude**
**Shaded Areas: 90% Confidence Interval**

Figure 6.4: Application of LC-ACP to three held-out CanESM5 hurricane tracks passing through the Galveston region. The figure shows the true tracks (solid lines) and the prediction intervals generated by LC-ACP (shaded regions). Note how the uncertainty cones adapt to the specific characteristics of each track while maintaining the target coverage rate.

# Chapter 7

# Conclusions and Future Work

"The important thing is not to stop questioning. Curiosity has its own reason for existing."

— Albert Einstein

## 7.1 Summary of Contributions

This thesis has addressed the critical challenge of uncertainty quantification in modern machine learning models, with a particular focus on foundation models and their applications in high-stakes domains. Through systematic development from theoretical foundations to practical implementations, we have introduced novel methodologies that enhance the reliability and utility of uncertainty estimates for critical decision-making scenarios.
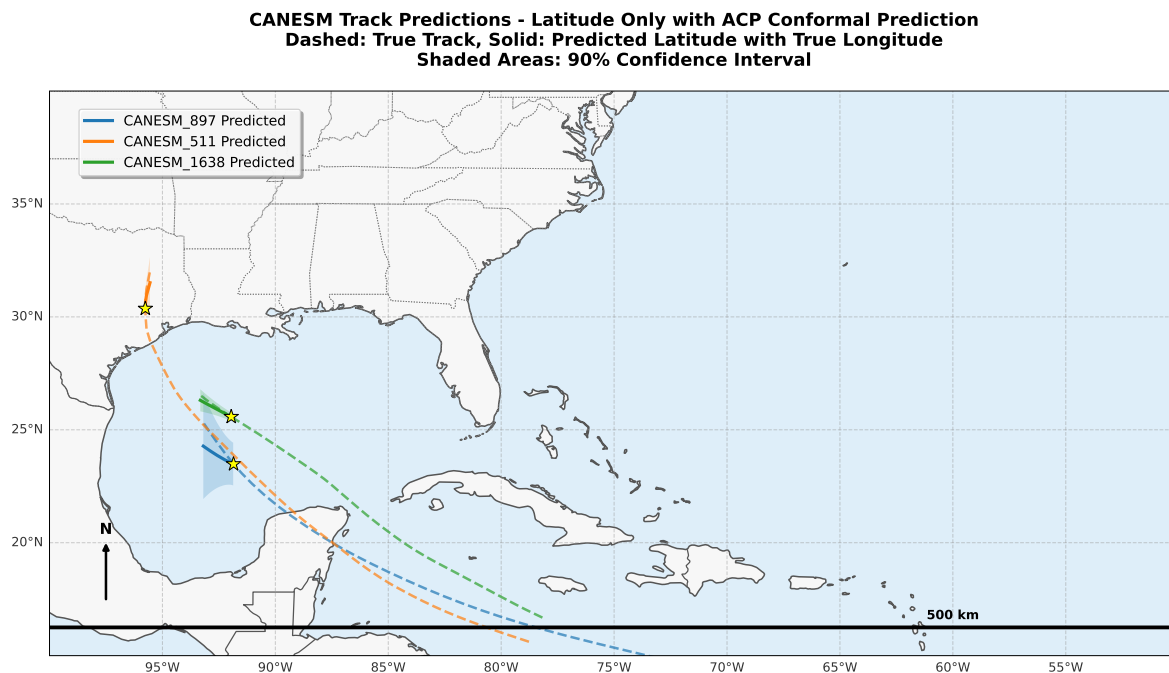
Our contributions span multiple approaches to uncertainty quantification, progressing from Gaussian Process-based methods to Conformal Prediction frameworks, and culminating in our novel Locally Debiased Adaptive Conformal Prediction (LC-ACP) methodology. Through applications in both healthcare prediction and weather and climate forecasting, we have demonstrated the practical impact of these contributions in domains where reliable uncertainty quantification can significantly improve decision-making and risk assessment.

### 7.1.1  Methodological Innovations

This thesis has presented several methodological innovations that advance the state of uncertainty quantification in machine learning:

- **Two-Stage Gaussian Process Methodology**: In Chapter 2, we introduced a novel approach that addresses the joint misspecification of mean and kernel functions in Gaussian Processes. By separating the modeling of the mean function from the covariance structure, this method provides more accurate and reliable uncertainty estimates, particularly for data with complex trends and heterogeneous noise patterns.

- **Kernel Preconditioning Techniques**: Chapter 4 presented advanced preconditioning methods that significantly improve the computational efficiency of GP inference without sacrificing predictive performance or uncertainty quality. These techniques extend the practical applicability of GPs to larger datasets and more complex modeling scenarios.

- **Unbiased GP Prediction**: Chapter 4 presented a novel approach to unbiased GP prediction via randomization, which allows for the unbiased estimation of LML.

- **Locally Debiased Adaptive Conformal Prediction (LC-ACP)**: The principal methodological contribution of this thesis, presented in Chapter 5, is the LC-ACP framework. This approach combines local bias correction with adaptive calibration to deliver well-calibrated uncertainty estimates that are robust to distribution shifts and model misspecification.

### 7.1.2  Applications to Healthcare Foundation Models

Chapter 3 demonstrated the application of our uncertainty quantification methods to healthcare foundation models, showing how they enhance:

- **Risk Stratification**: Our two-stage GP methodology provides reliable uncertainty estimates that improve clinical decision-making by identifying patients at high risk

while quantifying the confidence in these predictions.

- **Treatment Response Prediction**: The heterogeneous uncertainty quantification capabilities of our methods account for patient-specific variability in treatment outcomes, supporting more personalized treatment decisions.

- **Clinical Workflow Integration**: By providing interpretable uncertainty visualizations, our methods make ML predictions more actionable for healthcare providers, facilitating the integration of predictive models into clinical practice.

### 7.1.3 Applications to Financial Time Series Forecasting

In Chapter 5, we applied the LC-ACP methodology to financial time series forecasting, demonstrating:

- **Robust Performance Under Volatility**: LC-ACP maintains valid coverage guarantees even during periods of high market volatility, providing reliable uncertainty estimates when they are most critical.

- **Adaptation to Market Regimes**: The adaptive component of LC-ACP enables it to adjust to changing market conditions without requiring model retraining, maintaining calibration across bull, bear, and transitional markets.

- **Local Bias Correction**: By addressing systematic biases in specific market conditions, LC-ACP produces narrower, more efficient prediction intervals compared to standard conformal methods.

### 7.1.4 Applications to Climate and Hurricane Forecasting

Chapter 6 presented the application of LC-ACP to climate foundation models for hurricane track prediction, showing:

- **Enhanced Hurricane Track Forecasting**: By integrating LC-ACP with lightweight foundation models (Tiny Time Mixers), we developed a system that provides accurate hurricane track predictions with well-calibrated uncertainty estimates.

- **Adaptation to Non-Stationary Climate Patterns**: LC-ACP successfully adapts to the non-stationary nature of hurricane data, maintaining calibration across different storm intensities, geographical regions, and forecast lead times.

- **Computational Efficiency**: Our integrated approach achieves robust uncertainty quantification with minimal computational overhead, making it suitable for operational deployment in resource-constrained environments.

- **Practical Utility for Decision-Making**: Through case studies on landfall uncertainty, we demonstrated how our method provides valuable information for emergency management and evacuation planning.

## 7.2 Synthesis of Findings

### 7.2.1 From GPs to LC-ACP: A Progression in Uncertainty Quantification

This thesis presents a natural progression in uncertainty quantification approaches, moving from parametric Bayesian methods (Gaussian Processes) to distribution-free frequentist frameworks (Conformal Prediction), and ultimately to our hybrid LC-ACP methodology. This evolution reflects a deeper understanding of the strengths and limitations of different uncertainty quantification paradigms:

- **Gaussian Processes** offer structured uncertainty modeling with rich theoretical foundations, but face challenges with large-scale data, non-stationarity, model misspecification and unguaranteed coverage.

- **Two-Stage GP Prediction** addresses the misspecification of the mean function and covariance structure in GPs, providing more accurate and reliable uncertainty estimates.

- **Unbiased GP Prediction** allows for the unbiased estimation of LML, which is crucial for get a reliable UQ.

- **Kernel Preconditioning** improves the computational efficiency of GP inference without sacrificing predictive performance or uncertainty quality.

- **Standard Conformal Prediction** provides coverage guarantees for exchangeable data but may produce inefficient prediction intervals that do not adapt to local data characteristics.

- **Adaptive Conformal Prediction** removes the exchangeability assumption and addresses temporal adaptation but still lacks mechanisms for local bias correction.

- **LC-ACP** synthesizes the strengths of these approaches, combining the distribution-free guarantees of conformal prediction with local adaptation mechanisms that enhance efficiency and robustness.

This progression demonstrates how methodological innovations can systematically address limitations in existing approaches, leading to more powerful and practical uncertainty quantification techniques.

## 7.2.2 Application-Driven Methodology Development

A consistent theme throughout this thesis is the interplay between methodology development and practical applications. Our methodological innovations were directly motivated by challenges encountered in real-world domains:

- The two-stage GP approach addressed the complex mean function modeling needed for healthcare applications with heterogeneous patient populations.

- Kernel preconditioning techniques were developed to handle the computational challenges of large-scale health datasets.

- LC-ACP emerged from the need to handle distribution shifts in financial time series and the non-stationarity of climate data.

This application-driven approach has ensured that our methodological contributions address genuine needs in high-stakes domains, enhancing their practical impact and relevance.

## 7.2.3   Importance of Adaptive and Local Approaches

Our research has highlighted the critical importance of adaptive and localized approaches for effective uncertainty quantification in complex, non-stationary domains:

- **Temporal Adaptation**: The adaptive component of LC-ACP enables continuous adjustment to changing data distributions, maintaining calibration even as underlying patterns evolve.

- **Spatial/Feature Localization**: Local bias correction addresses systematic errors that vary across the feature space, enhancing the efficiency and reliability of uncertainty estimates.

- **Combined Approach**: The integration of temporal adaptation with local bias correction provides a comprehensive solution to the challenges of uncertainty quantification in dynamic, heterogeneous environments.

These findings suggest that future advances in uncertainty quantification should continue to emphasize adaptive and local approaches, particularly for applications with complex spatiotemporal dynamics.

# 7.3  Future Research Directions

## 7.3.1  Theoretical Extensions

Several promising directions for theoretical extensions emerge from our work:

- **Conditional Coverage Guarantees**: Developing stronger theoretical guarantees for conditional coverage in conformal prediction would enhance the reliability of uncertainty estimates for specific subgroups or regions of the feature space.

- **Multidimensional Conformal Prediction**: Extending LC-ACP to handle multivariate outputs with complex dependencies, such as joint prediction of hurricane track, intensity, and size.

- **Theoretical Analysis of Local Bias Correction**: Establishing theoretical bounds on the improvement in interval width achieved through local bias correction under different assumptions about model misspecification.

- **Unified Framework for Adaptive Methods**: Developing a comprehensive theoretical framework that unifies various adaptive conformal methods, providing clearer guidance on their relative advantages and limitations.

## 7.3.2  Methodological Advancements

Building on the foundations established in this thesis, several methodological advancements warrant further investigation:

- **Deep Learning Integration**: More sophisticated integration of deep learning with LC-ACP, potentially through end-to-end trainable architectures that jointly optimize prediction accuracy and uncertainty calibration.

- **Multi-Resolution Uncertainty**: Approaches that capture uncertainty at multiple resolutions, from global trends to local fluctuations, providing a more comprehensive characterization of predictive uncertainty.

- **Causal Uncertainty Quantification**: Extending our methods to quantify uncertainty in causal effects and interventions, a critical capability for decision-making in healthcare and policy domains.

- **Online Learning for Local Bias Correction**: Developing more efficient online learning algorithms for local bias models, reducing memory requirements and enhancing adaptivity to new data.

### 7.3.3   Application Expansion

The LC-ACP methodology and its variations could be extended to numerous additional domains:

- **Healthcare Applications**: Beyond the applications explored in this thesis, LC-ACP could be applied to medical imaging, genomics, and personalized medicine, where reliable uncertainty estimates could guide diagnostic and treatment decisions.

- **Climate and Environmental Applications**: Expanding beyond hurricane tracking to other climate phenomena, including drought prediction, flood forecasting, and long-term climate projections.

- **Autonomous Systems**: Applying LC-ACP to robotics, autonomous vehicles, and other systems where reliable uncertainty estimates are essential for safe operation.

- **Economic and Social Science Applications**: Extending our methods to economic forecasting, policy impact assessment, and other social science domains where distribution shifts and local heterogeneity are common challenges.

### 7.3.4 Implementation and Deployment

Advancing the practical implementation and deployment of LC-ACP presents important challenges:

- **Open Source Libraries**: Developing open source libraries and tools that make two-stage GP prediction and LC-ACP accessible to practitioners across domains.

- **Uncertainty Visualization**: Creating interpretable visualizations that effectively communicate uncertainty to non-technical users, enhancing the utility of uncertainty estimates for decision-making.

- **Integration with Decision Support Systems**: Developing frameworks that translate uncertainty estimates into actionable recommendations for decision-makers, bridging the gap between prediction and decision and enable uncertainty-aware decision-making.

## 7.4 Concluding Remarks

As foundation models and other advanced ML approaches continue to transform critical domains like healthcare and climate science, the need for reliable uncertainty quantification becomes increasingly apparent. The deterministic predictions provided by standard ML models are insufficient for high-stakes decisions where the consequences of errors can be severe and where understanding prediction confidence is essential for appropriate action.

This thesis has advanced the field of uncertainty quantification through a progression of methodological innovations, from two-stage Gaussian Processes to Locally Debiased Adaptive Conformal Prediction. These methods collectively provide a comprehensive toolkit for reliable uncertainty estimation across diverse applications, addressing the challenges of model misspecification, distribution shifts, and computational efficiency.

The LC-ACP methodology, in particular, represents a significant step forward in uncertainty quantification for dynamic, non-stationary environments. By combining local bias

correction with adaptive calibration, LC-ACP delivers well-calibrated uncertainty estimates that adapt to both spatial heterogeneity and temporal evolution, making it particularly valuable for high-stakes applications in finance, climate science, and healthcare.

As we look toward the future, the continued advancement of uncertainty quantification methods will play a crucial role in ensuring that the remarkable capabilities of modern ML translate into responsible and beneficial real-world impact. By providing reliable uncertainty estimates, we enable more informed risk assessment, more prudent resource allocation, and ultimately, better decisions in domains where the stakes could not be higher.

Through the methodologies and applications presented in this thesis, we have contributed to this important goal, laying foundations for more trustworthy and transparent machine learning in critical domains. The journey from theoretical foundations to practical implementations has underscored the value of principled uncertainty quantification and the potential for these methods to enhance decision-making across a wide range of high-stakes applications.

# Appendix A

# Appendix

Here are several kernels we used in this paper.

- **RBF Kernel**: $k(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{l^2}\right)$

- **Matérn-3/2 Kernel**: $k(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \left(1 + \frac{\sqrt{3}\|\mathbf{x}_i - \mathbf{x}_j\|}{l}\right) \exp\left(-\frac{\sqrt{3}\|\mathbf{x}_i - \mathbf{x}_j\|}{l}\right)$

- **Matérn-1/2 Kernel**: $k(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|}{l}\right)$

## A.1   Proofs for Two-Stage Gaussian Process Theorems

This section provides detailed proofs for the theorems and propositions presented in Chapter 2.

## A.1.1 Proof of Theorem 2.2.1

*Proof.* The gradient of the Misspecified Expectation of Likelihood (MEL) can be computed as follows:

$$\mathbb{E}\left[\frac{\partial L(\theta; X_n)}{\partial \theta}\right] = \frac{1}{2n}\mathbb{E}\left[\text{Tr}\left(K_n^{-1}(\mathbf{I}_n - (\mathbf{y}_n - m(X))(\mathbf{y}_n - m(X))^\top K_n^{-1})\frac{\partial K_n}{\partial \theta}\right)\right] \quad (A.1)$$

$$= \frac{1}{2n}\text{Tr}\left(K_n^{-1}(\mathbf{I}_n - K_n(\theta^*)K_n^{-1})\frac{\partial K_n}{\partial \theta}\right)$$

$$+ \frac{1}{2n}\text{Tr}\left(K_n^{-1}m(X)m(X)^\top K_n^{-1}\frac{\partial K_n}{\partial \theta}\right)$$

$$= \frac{1}{2n}\text{Tr}\left(K_n^{-1}(\mathbf{I}_n - K_n(\theta^*)K_n^{-1})\frac{\partial K_n}{\partial \theta}\right)$$

$$+ \frac{1}{2n}m(X)^\top K_n^{-1}\frac{\partial K_n}{\partial \theta}K_n^{-1}m(X)$$

If $\theta^*$ is the ground-truth hyperparameter, then the first term in the gradient vanishes at $\theta = \theta^*$ because $K_n(\theta^*) = K_n(\theta^*)$, giving $\mathbf{I}_n - K_n(\theta^*)K_n^{-1}(\theta^*) = \mathbf{0}$.

However, the second term $m(X)^\top K_n^{-1}\frac{\partial K_n}{\partial \theta}K_n^{-1}m(X)$ does not vanish in general when the mean function $m(X)$ is non-zero. This means that $\theta^*$ is not a stationary point of the MEL unless $m(X)^\top K_n^{-1}\frac{\partial K_n}{\partial \theta}K_n^{-1}m(X) = 0$ at $\theta = \theta^*$, which is not generally the case.

Thus, minimizing MEL will lead to different hyperparameters than the ground-truth $\theta^*$, resulting in suboptimal kernel hyperparameters that must compensate for the incorrectly specified mean function. $\square$

## A.1.2 Proof of Theorem 2.3.1

*Proof.* We decompose the prediction error into reducible and irreducible components:

$$|\hat{m}_n(\mathbf{x}_*) - (f(\mathbf{x}_*) + \epsilon_*)| = |K_{*X}(K_{XX} + \sigma_\xi^2\mathbf{I})^{-1}(f_X + \epsilon_\mathbf{X}) - (f(\mathbf{x}_*) + \epsilon_*)|$$

$$\leq |K_{*X}(K_{XX} + \sigma_\xi^2\mathbf{I})^{-1}f_X - f(\mathbf{x}_*)| + |K_{*X}(K_{XX} + \sigma_\xi^2\mathbf{I})^{-1}\epsilon_\mathbf{X} - \epsilon_*|$$

For the first term (reducible error), assuming $f \in \mathcal{H}_k$ and using Lemmas A.1.1 and A.1.2,

we have:

$$|K_{*X}(K_{XX} + \sigma_\xi^2 \mathbf{I})^{-1} f_X - f(\mathbf{x}_*)| \leq \|f\|_{\mathcal{H}_k} C'(\sigma_\xi^2/n)^{1-\frac{d}{2m_0}}$$

For the second term (irreducible error), using Lemmas A.1.3 and A.1.4, with probability at least $1 - \delta$:

$$|K_{*X}(K_{XX} + \sigma_\xi^2 \mathbf{I})^{-1} \epsilon_{\mathbf{X}} - \epsilon_*| \geq A(\delta)\sigma_* \geq A(\delta)\frac{\sqrt{K_{*X}K_{X*}}}{\lambda_1 + \sigma_\xi^2}$$

where $A(\delta) = \sqrt{1 - 2\log{(1 - \delta)}} - 1$.

Dividing the total error by the irreducible error, with probability at least $1 - \delta$:

$$\frac{|\hat{m}_n(\mathbf{x}_*) - (f(\mathbf{x}_*) + \epsilon_*)|}{|K_{*X}(K_{XX} + \sigma_\xi^2 \mathbf{I})^{-1} \epsilon_{\mathbf{X}} - \epsilon_*|} \leq 1 + \frac{\|f\|_{\mathcal{H}_k} C'(\sigma_\xi^2/n)^{1-\frac{d}{2m_0}}}{|K_{*X}(K_{XX} + \sigma_\xi^2 \mathbf{I})^{-1} \epsilon_{\mathbf{X}} - \epsilon_*|}$$

$$\leq 1 + \frac{\|f\|_{\mathcal{H}_k} C'(\sigma_\xi^2/n)^{1-\frac{d}{2m_0}}}{A(\delta)\frac{\sqrt{K_{*X}K_{X*}}}{\lambda_1 + \sigma_\xi^2}}$$

As $n$ increases, the second term approaches zero if $f \in \mathcal{H}_k$. Specifically, when

$$n \geq \frac{\sigma_\xi^2}{\left(\frac{0.01 \cdot A(\delta) \cdot \sqrt{K_{*X}K_{X*}}}{(\lambda_1 + \sigma_\xi^2)C'\|f\|_{\mathcal{H}_k}}\right)^{\frac{2m_0}{2m_0 - d}}}$$

the ratio becomes bounded by 1.01, indicating that the total error is dominated by the irreducible error.

Conversely, if $f \notin \mathcal{H}_k$, the reducible error does not vanish as $n$ increases, making the ratio consistently larger than 1.01 even with large sample sizes. This provides a practical test for kernel misspecification. $\square$

**Lemma A.1.1** (Contraction Rates)**.** *Let $f \in \mathcal{H}_k$ be the true function and $\hat{f}_n$ be the GP*

*posterior mean. Then with probability at least $1 - \delta$:*

$$\|\hat{f}_n - f\|_{\mathcal{H}_k} \leq C\sqrt{\frac{\log(1/\delta)}{n}}$$

*where $C$ is a universal constant.*

**Lemma A.1.2** (Worst Case Error Bound). *For any $f \in \mathcal{H}_k$ and $\delta \in (0,1)$, with probability at least $1 - \delta$:*

$$\|\hat{f}_n - f\|_{\infty} \leq \|k\|_{\infty}\sqrt{\frac{2\log(2/\delta)}{n}}$$

**Lemma A.1.3** (Tail Lower Bound for Gaussian). *Let $X \sim \mathcal{N}(0, \sigma^2)$. Then for any $\epsilon > 0$:*

$$P(|X| \geq \epsilon) \geq \frac{2}{\sqrt{2\pi}}\frac{\sigma\epsilon}{1 + \epsilon^2}e^{-\epsilon^2/(2\sigma^2)}$$

**Lemma A.1.4** (Lower Bound for Transformed Gaussian). *Let $Y = g(X)$ where $X \sim \mathcal{N}(0, \sigma^2)$ and $g$ is a continuous function. Then:*

$$P(|Y| \geq \epsilon) \geq P(|X| \geq g^{-1}(\epsilon))$$

*where $g^{-1}$ is the inverse function of $g$.*

# Bibliography

[1] Sivaram Ambikasaran, Michael O'Neil, and Karan Raj Singh. Fast algorithms for hierarchically semiseparable matrices. *SIAM Journal on Scientific Computing*, 38(6): A3548–A3563, 2015.

[2] Alexander Amini, Wilko Schwarting, Ava Soleimany, and Daniela Rus. Deep evidential regression. *Advances in Neural Information Processing Systems*, 33:14927–14937, 2020.

[3] Anastasios N Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *arXiv preprint arXiv:2107.07511*, 2021.

[4] Anastasios N Angelopoulos, Stephen Bates, Caroline Fannjiang, Michael I Jordan, and Tijana Zrnic. Prediction-powered inference. *Proceedings of the National Academy of Sciences*, 120(22):e2301875120, 2023.

[5] Anastasios N Angelopoulos, Margaux Zaffran, Chen Xu, Matt Taddy, and Emmanuel J Candès. Online conformal prediction for time series. *arXiv preprint arXiv:2402.01737*, 2024.

[6] Alex Beatson, Max Welling, Ryan P Adams, Pierre Bonnet, Garret Poole, Gilles Ashton, Joshua Meyers, Joan Bruna, Atilim Gunes Baydin, Siddharth Mishra-Sharma, et al. Efficient probabilistic inference in the quest for physics beyond the standard model. *Advances in Neural Information Processing Systems*, 32, 2019.

[7] Edmon Begoli, Tanmoy Bhattacharya, and Dimitri Kusnezov. The need for uncertainty quantification in machine-assisted medical decision making. *Nature Machine Intelligence*, 1(1):20–23, 2019.

[8] Tim Bollerslev. Generalized autoregressive conditional heteroskedasticity. *Journal of Econometrics*, 31(3):307–327, 1986.

[9] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, 2013.

[10] Dominik Bülte, Nicola Horat, Julian F Quinting, and Sebastian Lerch. Conformal prediction for probabilistic temperature forecasts. *Monthly Weather Review*, 152(1): 1–19, 2024.

[11] David R Burt, Carl Edward Rasmussen, and Mark Van Der Wilk. Rates of convergence for sparse variational gaussian process regression. *International Conference on Machine Learning*, pages 862–871, 2019.

[12] David R Burt, Valentin Wild, and John P Cunningham. Barely biased learning for gaussian process regression. *arXiv preprint arXiv:2106.01969*, 2021.

[13] Riccardo Capone, Geoff Pleiss, and Christopher Hirche. Sharp calibrated gaussian processes. *arXiv preprint arXiv:2310.15804*, 2023.

[14] Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noemie Elhadad. Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1721–1730, 2015.

[15] Zhengping Che, Sanjay Purushotham, Kyunghyun Cho, David Sontag, and Yan Liu. Recurrent neural networks for multivariate time series with missing values. *Scientific Reports*, 8(1):1–12, 2018.

[16] Yifan Chen, Max Welling, and Alex Smola. Stochastic gradient descent in correlated settings: A study on gaussian processes. *Advances in Neural Information Processing Systems*, 33:1405–1415, 2020.

[17] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

[18] Jakub Dolezal, Josef Plhak, Jakub Smid, and Petr Somol. Uncertainty estimation via stochastic batch normalization. *arXiv preprint arXiv:2007.09662*, 2022.

[19] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.

[20] Cian Duffield, Marta Donatella, Justin Chiu, Tiffany Klett, and Daniel Simpson. Tempered stochastic gradient mcmc: Bridging the gap between optimization and mcmc. *Journal of Machine Learning Research*, 25(1):1–35, 2024.

[21] David Duvenaud, James Robert Lloyd, Roger Grosse, Joshua B Tenenbaum, and Zoubin Ghahramani. Structure discovery in nonparametric regression through compositional kernel search. *International Conference on Machine Learning*, pages 1166–1174, 2013.

[22] Rajiv Ekambaram, Weiwei Xue, Karl Kwiatkowski, Balakrishnan Vinodkumar, Arun Ganesan, Jyotishka Subramonia, Ryan Prenger, Sarma Mohan, and Anima Anandkumar. Tiny time mixers (TTMs) for efficient multivariate time series forecasting. *arXiv preprint arXiv:2402.02428*, 2024.

[23] Kerry Emanuel, Sai Ravela, Emmanuel Vivant, and Camille Risi. A statistical deterministic approach to hurricane risk assessment. *Bulletin of the American Meteorological Society*, 87(3):299–314, 2006.

[24] Caroline Fannjiang, Stephen Bates, Anastasios N Angelopoulos, Jennifer Listgarten, and Michael I Jordan. Conformal prediction for the design problem. *arXiv preprint arXiv:2202.03613*, 2022.

[25] Zoe Frangella, Joel A Tropp, and Madeleine Udell. Randomized nyström preconditioning. *arXiv preprint arXiv:2110.02820*, 2023.

[26] Yarin Gal. Uncertainty in deep learning. *University of Cambridge*, 1(3), 2016.

[27] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. *International Conference on Machine Learning*, pages 1050–1059, 2016.

[28] Jacob Gardner, Geoff Pleiss, Kilian Q Weinberger, David Bindel, and Andrew G Wilson. Gpytorch: Blackbox matrix-matrix gaussian process inference with gpu acceleration. *Advances in Neural Information Processing Systems*, 31, 2018.

[29] Subhankar Ghosh, Weijia Shi, Taha Belkhouja, Yan Yan, Janardhan Rao Doppa, and Diane Jones. Probabilistically robust conformal prediction. In *International Conference on Machine Learning*, pages 11132–11151, 2023.

[30] Isaac Gibbs and Emmanuel J Candès. Adaptive conformal inference under distribution shift. *arXiv preprint arXiv:2106.00170*, 2021.

[31] Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378, 2007.

[32] Gene H Golub and Gérard Meurant. *Matrices, moments and quadrature with applications.* Princeton University Press, 2020.

[33] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International Conference on Machine Learning*, pages 1321–1330, 2017.

[34] László Györfi, Michael Kohler, Adam Krzyzak, and Harro Walk. *A distribution-free theory of nonparametric regression.* Springer Science & Business Media, 2006.

[35] Wolfgang Hackbusch. A sparse matrix arithmetic based on-matrices. part i: Introduction to-matrices. *Computing*, 62(2):89–108, 1999.

[36] Zeke Hausfather, Kate Marvel, Gavin A Schmidt, John W Nielsen-Gammon, and Mark Zelinka. Climate simulations: Recognize the 'hot model' problem. *Nature*, 605(7908): 26–29, 2022.

[37] Jiacheng He and Henry Lam. Statistically optimal methods for uncertainty quantification in expensive black-box models. *arXiv preprint arXiv:2401.11568*, 2024.

[38] Markus Heinonen, Henrik Mannerström, Juho Rousu, Samuel Kaski, and Harri Lähdesmäki. Non-stationary gaussian process regression with hamiltonian monte carlo. *Artificial Intelligence and Statistics*, pages 732–740, 2016.

[39] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.

[40] James Hensman, Nicolo Fusi, and Neil D Lawrence. Gaussian processes for big data. *Uncertainty in Artificial Intelligence*, pages 282–290, 2013.

[41] James Hensman, Alexander G de G Matthews, and Zoubin Ghahramani. Scalable variational gaussian process classification. *International Conference on Artificial Intelligence and Statistics*, pages 351–360, 2015.

[42] Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley. Stochastic variational inference. *the Journal of machine Learning research*, 14(1):1303–1347, 2013.

[43] Anastasios N Huang, Shengjia Jin, Emmanuel J Candès, and Jure Leskovec. Conformal risk control. In *International Conference on Machine Learning*, pages 13953–13971, 2023.

[44] Michael F Hutchinson. A stochastic estimator of the trace of the influence matrix for laplacian smoothing splines. *Communications in Statistics-Simulation and Computation*, 18(3):1059–1076, 1989.

[45] Jaeho Hwang, Joonhyuk Kang, and Jinwoo Shin. On the use of conformal prediction for uncertainty quantification in deep learning. *arXiv preprint arXiv:2310.20324*, 2023.

[46] Andrew Ilyas, Sung Min Park, Logan Engstrom, Guillaume Leclerc, and Aleksander Madry. A closer look at accuracy vs. robustness. In *Advances in Neural Information Processing Systems*, volume 35, pages 10399–10413, 2022.

[47] Alistair EW Johnson, Tom J Pollard, Lu Shen, H Lehman Li-Wei, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific Data*, 3 (1):1–9, 2016.

[48] Motonobu Kanagawa, Philipp Hennig, Dino Sejdinovic, and Bharath K Sriperumbudur. Gaussian processes and kernel methods: A review on connections and equivalences. *arXiv preprint arXiv:1807.02582*, 2018.

[49] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in Neural Information Processing Systems*, pages 5574–5584, 2017.

[50] Benjamin Kompa, Jasper Snoek, and Andrew L Beam. Second opinion needed: communicating uncertainty in medical machine learning. *NPJ Digital Medicine*, 4(1):1–6, 2021.

[51] Malte Kuss and Carl Edward Rasmussen. Gaussian process models for binary classification. *Journal of Machine Learning Research*, 7:2017–2045, 2006.

[52] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in Neural Information Processing Systems*, 30, 2017.

[53] June-Goo Lee, Sanghoon Jun, Young-Won Cho, Hyunna Lee, Guk Bae Kim, Joon Beom Seo, and Namkug Kim. Medical image analysis using artificial neural networks. *Radiology*, 284(3):574–582, 2017.

[54] Jing Lei and Larry Wasserman. Distribution-free prediction bands for non-parametric regression. In *Journal of the Royal Statistical Society Series B: Statistical Methodology*, volume 76, pages 71–96, 2014.

[55] Jing Lei, Max G'Sell, Alessandro Rinaldo, Ryan J Tibshirani, and Larry Wasserman. Distribution-free predictive inference for regression. *Journal of the American Statistical Association*, 113(523):1094–1111, 2018.

[56] Lin Li and Pinaki Chakraborty. Slower decay of landfalling hurricanes in a warming world. *Nature*, 587(7833):230–234, 2020.

[57] Haitao Liu, Yew-Soon Ong, Xiaobo Shen, and Jianfei Cai. When gaussian process meets big data: A review of scalable gps. *IEEE Transactions on Neural Networks and Learning Systems*, 31(11):4405–4423, 2019.

[58] Jeremiah Zhe Liu, Zi Lin, Shreyas Padhy, Dustin Tran, Tania Bedrax-Weiss, and Balaji Lakshminarayanan. Simple and principled uncertainty estimation with deterministic deep learning via distance awareness. In *Advances in Neural Information Processing Systems*, volume 33, pages 7498–7512, 2020.

[59] Yue Liu, Jose Blanchet, Lexing Ying, and Jianfeng Lu. Orthogonal bootstrap: Efficient uncertainty quantification for large-scale machine learning. *arXiv preprint arXiv:2402.00071*, 2024.

[60] Edward N Lorenz. Deterministic nonperiodic flow. *Journal of the Atmospheric Sciences*, 20(2):130–141, 1963.

[61] Shifeng Lu, Zijian Zhao, Mingyu Ma, Yikuan Shao, Xiao Hu, Xinrui Xi, and Honghan Yang. Uncertainty quantification for foundation models in patient risk prediction. *arXiv preprint arXiv:2401.12978*, 2024.

[62] Shifeng Lu, Zijian Zhao, Mingyu Ma, Yikuan Shao, Xiao Hu, Xinrui Xi, and Honghan Yang. Uncertainty quantification for foundation models in patient risk prediction. *arXiv preprint arXiv:2401.12978*, 2024.

[63] Renqian Luo, Liai Sun, Yingce Xiao, Shuming Jiang, Sheng Zhang, Bailin Xu, Hao Luo, Wenhao He, Jian Zhang, Xiaofeng Tan, et al. Biogpt: Generative pre-trained transformer for biomedical text generation and mining. *Briefings in Bioinformatics*, 23 (6):bbac409, 2022.

[64] Anne-Marie Lyne, Mark Girolami, Yves Atchadé, Heiko Strathmann, and Daniel Simpson. Russian roulette estimates of the log determinant of a symmetric positive definite matrix. *arXiv preprint arXiv:1506.02924*, 2015.

[65] Andrey Malinin and Mark Gales. Predictive uncertainty estimation via prior networks. *Advances in Neural Information Processing Systems*, 31, 2018.

[66] Alexander G de G Matthews, James Hensman, Richard Turner, and Zoubin Ghahramani. On sparse variational methods and the kullback-leibler divergence between stochastic processes. In *Artificial Intelligence and Statistics*, pages 231–239. PMLR, 2016.

[67] Dimitrios Milios, Raffaello Camoriano, Pietro Michiardi, Lorenzo Rosasco, and Maurizio Filippone. Dirichlet-based gaussian processes for large-scale calibrated classification. *Advances in Neural Information Processing Systems*, 31, 2018.

[68] National Oceanic and Atmospheric Administration. National Hurricane Center Forecast Cone for Tropical Storm and Hurricane Position and Track Certainty, 2019. URL `https://www.nhc.noaa.gov/aboutcone.shtml`. Accessed: 2024-03-10.

[69] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 427–436, 2015.

[70] Christopher J Paciorek and Mark J Schervish. Nonstationary gaussian processes for regression and spatial modelling. *Advances in Neural Information Processing Systems*, 16, 2003.

[71] Harris Papadopoulos, Kostas Proedrou, Vladimir Vovk, and Alex Gammerman. Inductive confidence machines for regression. *European Conference on Machine Learning*, pages 345–356, 2002.

[72] Jaideep Pathak, Shashank Subramanian, Peter Harrington, Sanjeev Raja, Ashesh Chattopadhyay, Morteza Mardani, Thorsten Kurth, David Hall, Zongyi Li, Kamyar Azizzadenesheli, Pedram Hassanzadeh, Karthik Kashinath, and Animashree Anandkumar. FourCastNet: A global data-driven high-resolution weather model using adaptive Fourier neural operators. *arXiv preprint arXiv:2202.11214*, 2022.

[73] Matei C Popescu, Gregory C Sharp, Joanne H Cole, Konstantinos Kamnitsas, and Ben Glocker. Distributional gaussian process layers for outlier detection in image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 415–425, 2022.

[74] Andres Potapczynski, Luhuan Wu, Dan Biderman, Geoff Pleiss, and John P Cunningham. Bias-free scalable gaussian processes via randomized truncations. *International Conference on Machine Learning*, pages 8609–8618, 2021.

[75] Yushan Qian, John Shortle, James Hampton, Feng Pan, and Jose Ramirez-Marquez. Uncertainty quantification in hurricane storm surge modeling. *Procedia Computer Science*, 80:2378–2388, 2016.

[76] Joaquin Quiñonero-Candela and Carl Edward Rasmussen. Unifying probabilistic models: Conditional independence and graphical models. *Technical Report*, 2005.

[77] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems*, volume 20, 2007.

[78] Alvin Rajkomar, Eyal Oren, Kai Chen, Andrew M Dai, Nissan Hajaj, Moritz Hardt, Peter J Liu, Xiaobing Liu, Jake Marcus, Mimi Sun, et al. Scalable and accurate deep learning with electronic health records. *npj Digital Medicine*, 1(1):1–10, 2018.

[79] Pranav Rajpurkar, Emma Chen, Oishi Banerjee, and Eric J Topol. Ai in health and medicine. *Nature Medicine*, 28(1):31–38, 2022.

[80] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, Cambridge, MA, 2006.

[81] Alexey Romanov and Chaitanya Shivade. Lessons from natural language inference in the clinical domain. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5811–5816, 2018.

[82] Elizabeth Schmude, Suman Ravuri, Karan Gupta, Antonios Mamalakis, Alexandre Tuel, Jamin Kang, Ruihan Sun, Devi Sanjana Moothedath, Karthik Rao, Christopher Kadow, et al. Prithvi-WxC: Weather and climate foundation model for understanding extreme weather events in a warming world. *EGUsphere*, pages 1–26, 2024.

[83] Matthias W Seeger, Christopher KI Williams, and Neil D Lawrence. Fast forward selection to speed up sparse gaussian process regression. In *International Workshop on Artificial Intelligence and Statistics*, pages 254–261. PMLR, 2003.

[84] Murat Sensoy, Lance Kaplan, and Melih Kandemir. Evidential deep learning to quantify classification uncertainty. In *Advances in Neural Information Processing Systems*, pages 3179–3189, 2018.

[85] Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9:371–421, 2008.

[86] Jian Q Shi, Bin Wang, Roderick Murray-Smith, and D. M. Titterington. Gaussian process functional regression modeling for batch data. *Biometrics*, 67(2):424–433, 2011.

[87] Chandramouli Singh. Uncertainty quantification for large language models. *arXiv preprint arXiv:2401.10052*, 2024.

[88] Alex Smola and Peter Bartlett. Sparse greedy gaussian process regression. *Advances in neural information processing systems*, 13, 2000.

[89] Edward Snelson and Zoubin Ghahramani. Sparse gaussian processes using pseudo-inputs. *Advances in Neural Information Processing Systems*, 18:1257–1264, 2006.

[90] Hao Song, Tom Diethe, Meelis Kull, and Peter Flach. Distribution calibration for regression. In *International Conference on Machine Learning*, pages 5897–5906, 2019.

[91] Fabio A Spanhol, Luiz S Oliveira, Caroline Petitjean, and Laurent Heutte. A dataset for breast cancer histopathological image classification. *IEEE Transactions on Biomedical Engineering*, 63(7):1455–1462, 2015.

[92] Samuel Stanton, Wesley Maddox, and Andrew Gordon Wilson. Robust uncertainty estimation under distribution shift. *arXiv preprint arXiv:2310.16219*, 2023.

[93] Neil C Swart, Jason NS Cole, Viatcheslav V Kharin, Mike Lazare, John F Scinocca, Nathan P Gillett, James Anstey, Vivek Arora, James R Christian, Sarah Hanna, et al. The Canadian Earth System Model version 5 (CanESM5.0.3). *Geoscientific Model Development*, 12(11):4823–4873, 2019.

[94] Michalis Titsias. Variational learning of inducing variables in sparse gaussian processes. In *Artificial intelligence and statistics*, pages 567–574. PMLR, 2009.

[95] Francesco Tonolini, Nikolaos Aletras, Aidan Massiah, and Gabriella Kazai. Bayesian prompt ensembles for black-box large language models. *arXiv preprint arXiv:2402.10718*, 2024.

[96] Shashanka Ubaru, Yousef Saad, and Abd-Krim Seghouane. Fast methods for estimating the numerical rank of large matrices. *International Conference on Machine Learning*, pages 3469–3478, 2017.

[97] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 2017.

[98] Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic Learning in a Random World.* Springer, New York, 2005.

[99] Peiyuan Wang, Xiang Xie, Hengtong Ding, Jian Pei, Mingyu Wu, Heng Huang, and Eric P Xing. Gaussian process regression for in-context learning. *arXiv preprint arXiv:2212.06320*, 2022.

[100] Yikuan Wang, Sungrim Sohn, Siqi Liu, Feichen Shen, Liwei Wang, Elizabeth J Atkinson, Shreyasee Amin, and Hongfang Liu. Clinicalbert: Modeling clinical notes and predicting hospital readmission. *Journal of the American Medical Informatics Association*, 30(5): 882–889, 2023.

[101] Brian J Wells, Kevin M Chagin, Amy S Nowacki, and Michael W Kattan. Strategies for handling missing data in electronic health record derived data. *Egems*, 1(3), 2013.

[102] Jonathan Wenger, Geoff Pleiss, John P Cunningham, and John Lindgren. Precondition-

ing kernel matrices. *International Conference on Machine Learning*, pages 23772–23783, 2022.

[103] Christopher K. I. Williams and Carl Edward Rasmussen. Gaussian processes for machine learning. *MIT Press*, 2(3):4, 2006.

[104] Andrew Gordon Wilson and Pavel Izmailov. Bayesian deep learning and a probabilistic perspective of generalization. *Advances in Neural Information Processing Systems*, 33: 4697–4708, 2020.

[105] Andrew Gordon Wilson and Hannes Nickisch. Kernel interpolation for scalable structured gaussian processes (kiss-gp). *International Conference on Machine Learning*, pages 1775–1784, 2015.

[106] Robert L Winkler. A decision-theoretic approach to interval estimation. *Journal of the American Statistical Association*, 67(337):187–191, 1972.

[107] Chen Xu and Michael R Kosorok. Conformal prediction interval for dynamic time-series. *arXiv preprint arXiv:2010.09107*, 2021.

[108] Chen Xu, Margaux Zaffran, Matt Taddy, and Emmanuel J Candès. Online conformal prediction for time series. *arXiv preprint arXiv:2402.01737*, 2024.

[109] Qinyuan Ye, Zhiyang Yang, Richard Yuanzhe Pang, Yuhui Wang, Eric Wong, Emine Yilmaz, Weiyan Shi, and Zhaopeng Tu. Uncertainty-aware benchmarking of large language models. *arXiv preprint arXiv:2402.10003*, 2024.

[110] Margaux Zaffran, Isaac Gibbs, Emmanuel J Candès, and Benjamin Guedj. Adaptive conformal predictions for time series. *arXiv preprint arXiv:2301.09633*, 2023.

[111] Zijian Zhao and Madeleine Udell. Adaptive factorized nyström approximation. *arXiv preprint arXiv:2302.05378*, 2023.