# Optimizing Gaussian Processes
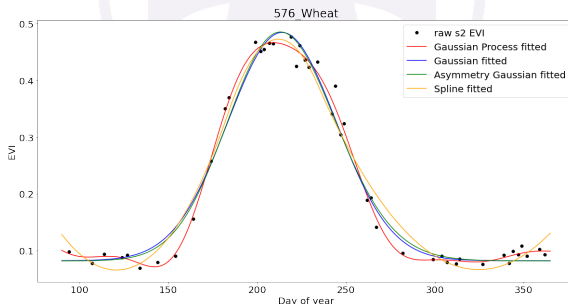
**Honours Research Project**

Michael Ciccotosto-Camp - 44302913

THE UNIVERSITY
OF QUEENSLAND
AUSTRALIA

## Problem Setting and Motivation

- This project focuses on the problem of time series prediction.
- Given a data set of $n$ observations $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{n}$, where each input $x_i \in \mathbb{R}_{>0}$ is a time value and $y_i \in \mathbb{R}$ is a output or experimental observation that acts a function of time, the goal of time series prediction is to try and best predict a value $y_\star$ at time $x_\star$.
- The idea of studying time series prediction came from a research group from the Gatton campus, lead by Andries Potgieter, analysing crop growth from previous seasons to forecast when certain phenological stages will take place in the current harvest.



576_Wheat

**Introduction to Gaussian Processes**

- Originally, Potgieter's team surveyed a number of different parameteric models to carry out forecasting. However, the parameteric models we serverely limited in their ability to inform when key phenological stages would take place. After seeing the success of applying GPs to other remote sensing tasks, the team investigated the use of GPs in their own research to find that they could produce much higher resolution predictions from which they could infer a far richer phenological timeline.

- A Gaussian Process (GP) is a collection of random variables with index set $I$, such that every finite subset of random variables has a joint Gaussian distribution.

- A GP is completely characterized by a mean function $m : X \rightarrow \mathbb{R}$ and a kernel $k : X \times X \rightarrow \mathbb{R}$.

$$m(\boldsymbol{x}) = \mathbb{E}\left[f(\boldsymbol{x})\right]$$
$$k(\boldsymbol{x}, \boldsymbol{x'}) = \mathbb{E}\left[(f(\boldsymbol{x}) - m(\boldsymbol{x}))(f(\boldsymbol{x'}) - m(\boldsymbol{x'}))\right].$$

- In this context, think of the kernel as a function that provides some notion of similarity between points.

## Predictions

- The aim of GPs is to find a suitable mean function, $m$, for which we can then predict inputs from outside the observed values, $X_\star$. This requires an understanding of the function $f$.

- Adopting the notation $(K_{WW'})_{i,j} \triangleq k(w_i, w'_j)$, when attempting to model our value function we usually do not have access to the value function itself but a noisy version thereof, $y = f(x) + \varepsilon$ where $\varepsilon \sim \mathbb{N}(0, \sigma_n^2)$ meaning the prior on the noisy values becomes $\text{cov}(y) = K_{XX} + \sigma_n^2 I$.

- Using the assumption that our data can be modelled as a Gaussian process, we can write out the new distribution of the observed noisy values along the points at which we wish to test the underlying function as

$$\begin{bmatrix} y \\ f_\star \end{bmatrix} \sim \mathbb{N}\left(0, \quad \begin{bmatrix} K_{XX} + \sigma_n^2 \mathbb{I}_{n \times n} & K_{X_\star X}^\intercal \\ K_{X_\star X} & K_{X_\star X_\star} \end{bmatrix}\right).$$

- The mean and covariance can then be computed as

$$\overline{f_\star} \triangleq K_{X_\star X} \left[ K_{XX} + \sigma_n^2 \mathbb{I}_{n \times n} \right]^{-1} y$$

$$\text{cov}(f_\star) = K_{X_\star X_\star} - K_{X_\star X} \left[ K_{XX} + \sigma_n^2 \mathbb{I}_{n \times n} \right]^{-1} K_{X_\star X}^\intercal.$$

## Unoptimized GPR

---

**Algorithm 1:** Unoptimized GPR

**input** : Observations $\boldsymbol{X}, \boldsymbol{y}$ and a test input $\boldsymbol{x}_\star$.

**output:** A prediction $\overline{f_\star}$ with its corresponding variance $\mathbb{V}[f_\star]$.

1 $\boldsymbol{L} = \text{cholesky}\left(\boldsymbol{K_{XX}} + \sigma_n^2 \mathbb{I}_{n \times n}\right)$

2 $\alpha = \text{lin-solve}\left(\boldsymbol{L}^\mathsf{T}, \text{lin-solve}\left(\boldsymbol{L}, \boldsymbol{y}\right)\right)$

3 $\overline{f_\star} = \boldsymbol{K}_{\boldsymbol{x}_\star \boldsymbol{X}} \alpha$

4 $\boldsymbol{v} = \text{lin-solve}\left(\boldsymbol{L}, \boldsymbol{K}_{\boldsymbol{x}_\star \boldsymbol{X}}\right)$

5 $\mathbb{V}[f_\star] = \boldsymbol{K}_{\boldsymbol{x}_\star \boldsymbol{x}_\star} - \boldsymbol{v}^\mathsf{T} \boldsymbol{v}$

6 **return** $\overline{f_\star}, \mathbb{V}[f_\star]$

---

## Problems with Unoptimized GPR

---

**Algorithm 2:** Unoptimized GPR

---

**input** : Observations $\boldsymbol{X}, \boldsymbol{y}$ and a test input $\boldsymbol{x}_\star$.

**output:** A prediction $\overline{f_\star}$ with its corresponding variance $\mathbb{V}[f_\star]$.

1   $\boldsymbol{L} = \text{cholesky}\left(\boldsymbol{K_{XX}} + \sigma_n^2 \mathbb{I}_{n \times n}\right)$

2   $\alpha = \text{lin-solve}\left(\boldsymbol{L}^\mathsf{T}, \text{lin-solve}\left(\boldsymbol{L}, \boldsymbol{y}\right)\right)$

3   $\overline{f_\star} = \boldsymbol{K_{x_\star X}}\alpha$

4   $\boldsymbol{v} = \text{lin-solve}\left(\boldsymbol{L}, \boldsymbol{K_{x_\star X}}\right)$

5   $\mathbb{V}[f_\star] = \boldsymbol{K_{x_\star x_\star}} - \boldsymbol{v}^\mathsf{T}\boldsymbol{v}$

6   **return** $\overline{f_\star}, \mathbb{V}[f_\star]$

---

- Line 1 can be incredibly slow as computing $\boldsymbol{K_{XX}}$ and performing a Cholesky decomposition scale poorly as the number of inputs, $n$, grows.

## Nystrom Approximation

- One technique to speed up the computation of $K_{XX}$ is to use a Nystrom approximation.
- The Nystrom method we seek a matrix $Q \in \mathbb{R}^{n \times k}$ that satisfies $\|A - QQ^*A\|_F \leq \varepsilon$, where $A \in \mathbb{R}^{n \times n}$ is positive semi definite matrix, to form the rank$-k$ approximation

$$
\begin{aligned}
A &\simeq QQ^*A \\
&\simeq Q \left( Q^*AQ \right) Q^* \\
&= Q \left( Q^*AQ \right) \left( Q^*AQ \right)^\dagger \left( Q^*AQ \right) Q^* \\
&\simeq (AQ) \left( Q^*AQ \right)^\dagger \left( Q^*A \right).
\end{aligned}
$$

- A matrix $Q$ that satisfies the above conditions can be built using through a very popular column sampling technique.
- As the name suggests, the matrix $Q$ essentially samples and rescales columns from $A$ using a probability distribution $\{p_i\}_{i=1}^n$.
- When $Q$ is constructed in this manner, it is usually referred to as a sketching matrix and denoted $S$.

**Random Fourier Feature Approximation**

- The other technique investigated to speed up the computation of $K_{XX}$ is the Random Fourier Feature (RFF) approximation.
- The main idea is instead of using a kernel function to implicitly lift data into a higher dimensional feature space, an explicit feature map $\varphi : \mathbb{R}^d \to \mathbb{R}^D$ could be used to approximate $k$ as $k(x, y) = \langle \Phi(x), \Phi(y) \rangle_{\mathbb{R}^N} \simeq \langle \varphi(x), \varphi(y) \rangle_{\mathbb{R}^D}$ where $D$ is chosen so that $n \gg D$. Once $\varphi(x_i)$ has been computed for each $x_i$, every entry of the Gram matrix can be swiftly approximated as $K_{ij} = K_{ji} \simeq \langle \varphi(x_i), \varphi(y_j) \rangle_{\mathbb{R}^D}$.
- The RFF technique hinges on Bochners theorem which characterises positive definite functions (namely kernels) and states that any positive definite functions can be represented as

$$k(x - y) = k(x - y) = \int_{\mathbb{C}^d} \exp(i \langle \omega, x - y \rangle) \mu_k(d\omega)$$

where $\mu_k$ is a positive finite measure on the frequencies of $\omega$.
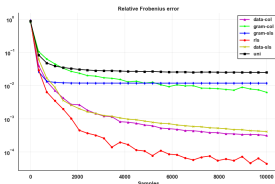
- The integral $k(\mathbf{x} - \mathbf{y}) = \int_{\mathbb{C}^d} \exp(i\langle \boldsymbol{\omega}, \mathbf{x} - \mathbf{y}\rangle)\, \mu_k(d\boldsymbol{\omega})$ can then be approximated via the following Monte Carlo estimate

$$
\begin{aligned}
k(\mathbf{x} - \mathbf{y}) &= \int_{\mathbb{C}^d} \exp(i\langle \boldsymbol{\omega}, \mathbf{x} - \mathbf{y}\rangle)\, p(\boldsymbol{\omega})\, d\boldsymbol{\omega} \\
&= \mathbb{E}_{\boldsymbol{\omega} \sim p(\cdot)}\left(\exp(i\langle \boldsymbol{\omega}, \mathbf{x} - \mathbf{y}\rangle)\right) \\
&\simeq \frac{1}{D} \sum_{j=1}^{D} \exp(i\langle \boldsymbol{\omega}_j, \mathbf{x} - \mathbf{y}\rangle) \\
&= \sum_{j=1}^{D} \left(\frac{1}{\sqrt{D}} \exp(i\langle \boldsymbol{\omega}_j, \mathbf{x}\rangle)\right) \overline{\left(\frac{1}{\sqrt{D}} \exp(i\langle \boldsymbol{\omega}_j, \mathbf{y}\rangle)\right)} \\
&= \langle \varphi(\mathbf{x}), \varphi(\mathbf{y})\rangle_{\mathbb{C}^D}
\end{aligned}
$$
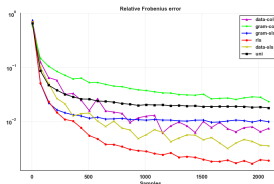
where $\boldsymbol{\omega}_i \overset{\text{iid}}{\sim} p(\cdot)$ using the feature map

$$
\varphi(\mathbf{x}) = \frac{1}{\sqrt{D}}\left[z(\boldsymbol{\omega}_1, \mathbf{x}), z(\boldsymbol{\omega}_2, \mathbf{x}), \ldots, z(\boldsymbol{\omega}_D, \mathbf{x})\right]^{\mathsf{T}}
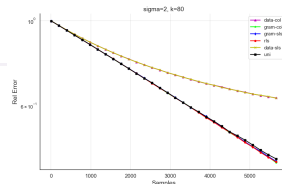$$

where for convenience of notation $z(\boldsymbol{\omega}, \mathbf{x}) = \exp(i\langle \boldsymbol{\omega}, \mathbf{x}\rangle)$.
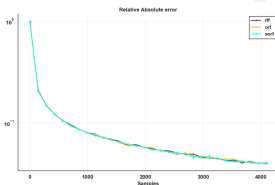
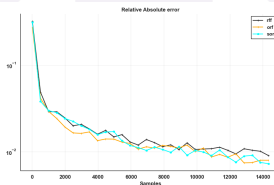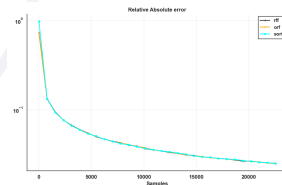(a) 3D-Spatial network

(b) Abalone

(c) Temperature

Figure: Comparison of Nystrom methods for various datasets.



(a) 3D-Spatial network

(b) Abalone

(c) Wine

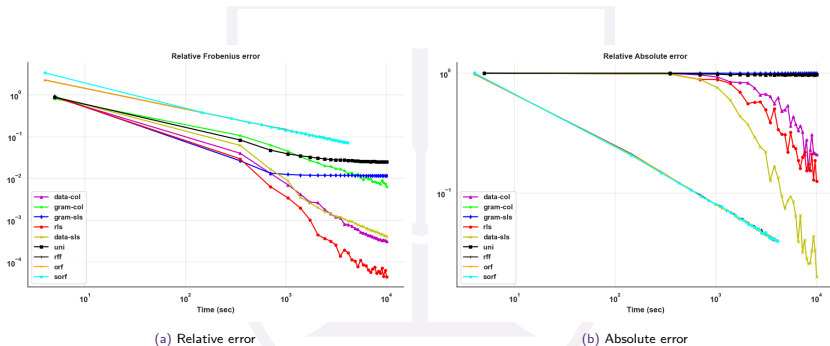Figure: Comparison of RFF methods for various datasets.

(a) Relative error

(b) Absolute error

Figure: Comparison between Nystrom and RFF approximations for the 3D-Spatial network data.

## Moving Forward

- We saw that the Nystrom technique is better at producing approximations of the Gram matrix, $K_{XX}$, with smaller *relative errors* while the RFF technique is better at producing approximations with smaller *absolute errors*.

- The question reamins which is better for GP prediction, lower relative error or absolute error (or perhaps some combination of the two)? This will ultimately determined which of the two techniques are more useful in Machine Learning applications.

- Recall, the other bottle neck in the GPR algorithm was the Cholesky decomposition.

- This can be avoided by employing faster linear system solvers. The two solvers we are interested in are the *Conjugate Gradient* (CG) and *Minimum Residual* (MINRES) methods used for solving positive semi-definite and symmeteric indefinite systems respectively.

- While MINRES can be applied to a wider class of linear systems, in certain scenarios it may perform better than the more widely used CG method. It will be interesting to see which of the two produced better predictions in the context of GPs.