# COMP4131: Data Modelling and Analysis
## Lecture 7: Gaussian Process Regression and Classification

Daokun Zhang

University of Nottingham Ningbo China

*daokun.zhang@nottingham.edu.cn*

31 March, 2025

# Overview

1 Gaussian Process

2 Gaussian Process Regression

3 Gaussian Process Classification

# Gaussian Process

## Linear Regression

Recall the probabilistic analysis of the standard linear regression model with Gaussian noise

$$f(\mathbf{x}) = \mathbf{w}^{\mathrm{T}} \phi(\mathbf{x}), \quad y = f(\mathbf{x}) + \varepsilon,$$

where $\varepsilon$ is the additive noise that follows an i.i.d Gaussian distribution with zero mean and variance $\sigma_n^2$:

$$\varepsilon \sim \mathcal{N}\left(0, \sigma_n^2\right).$$

Given the training data $\mathbf{X} = \{\mathbf{x}_1, \ldots, \mathbf{x}_N\}$ and $\mathbf{y} = \{y_1, \cdots, y_N\}$, we can use maximum likelihood estimation (MLE) to find the solution of $\mathbf{w}$ as

$$\mathbf{w}^* = \left(\mathbf{\Phi}^{\mathrm{T}}\mathbf{\Phi}\right)^{-1}\mathbf{\Phi}^{\mathrm{T}}\mathbf{y},$$

where $\mathbf{\Phi} = [\phi(\mathbf{x}_1), \cdots, \phi(\mathbf{x}_N)]^{\mathrm{T}}$ is an $N \times M$ matrix.

# Bayesian Linear Regression

In the Bayesian formalism, we can specify a prior over the parameters, expressing our beliefs about the parameters before we look at the observations. We put a Gaussian prior with mean vector $\boldsymbol{\mu_w}$ and covariance matrix $\boldsymbol{\Sigma_w}$ on the weights

$$\mathbf{w} \sim \mathcal{N}\left(\boldsymbol{\mu_w}, \boldsymbol{\Sigma_w}\right).$$

Due to the randomness in $\mathbf{w}$, $f(x) = \mathbf{w}^{\mathrm{T}}\phi(\mathbf{x})$ is no longer a deterministic but a random function. Let $f_p$ denote the linear function value $f(\mathbf{x}_n)$ at sample $\mathbf{x}_n$, then $f_1, \cdots f_N$ follow a multivariate Gaussian distribution with

$$\mathbb{E}\left[f_p\right] = \boldsymbol{\mu_w}^{\mathrm{T}}\phi(\mathbf{x}_p),$$
$$\mathrm{Cov}\left(f_n, f_m\right) = \phi\left(\mathbf{x}_n\right)^{\mathrm{T}}\boldsymbol{\Sigma_w}\phi\left(\mathbf{x}_m\right).$$

Now we get a Gaussian process $f(\mathbf{x})$ with mean function $m(\mathbf{x}) = \boldsymbol{\mu_w}^{\mathrm{T}}\phi(\mathbf{x})$ and covariance function $k(\mathbf{x}, \mathbf{x}') = \phi\left(\mathbf{x}\right)^{\mathrm{T}}\boldsymbol{\Sigma_w}\phi\left(\mathbf{x}'\right)$.

# Gaussian Process

**Definition**: *A Gaussian process is a collection of random variables, any finite number of which have a joint Gaussian distribution.*

**Implications**

- "a collection of random variables" means a set that has
  - a finite number of random variables
  - an infinite number of random variables
- "any finite number of which have a joint Gaussian distribution" means:
  - multiple random variables follow a joint multivariate Gaussian distribution
  - after marginalization, a single random variable follows a univariate Gaussian distribution

# Gaussian Process

A Gaussian process is completely specified by its mean function and covariance function. We define mean function $m(\mathbf{x})$ and the covariance function $k(\mathbf{x}, \mathbf{x}')$ as

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

and write the Gaussian process as

$$f(\mathbf{x}) \sim \mathcal{GP}\left(m(\mathbf{x}), k\left(\mathbf{x}, \mathbf{x}'\right)\right).$$

Usually, for notational simplicity we will take the mean function to be zero, i.e., set $m(\mathbf{x}) = 0$.

**Example**: For Bayesian linear regression model $f(\mathbf{x}) = \mathbf{w}^{\mathrm{T}}\phi(\mathbf{x})$ with prior $\mathbf{w} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{\Sigma_w}\right)$, we have the mean and covariance function

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})] = \phi(\mathbf{x})^{\mathrm{T}}\mathbb{E}[\mathbf{w}] = 0,$$
$$k\left(\mathbf{x}, \mathbf{x}'\right) = \mathbb{E}\left[f(\mathbf{x})f\left(\mathbf{x}'\right)\right] = \phi(\mathbf{x})^{\mathrm{T}}\mathbb{E}\left[\mathbf{w}\mathbf{w}^{\mathrm{T}}\right]\phi\left(\mathbf{x}'\right) = \phi(\mathbf{x})^{\mathrm{T}}\mathbf{\Sigma_w}\phi\left(\mathbf{x}'\right).$$

# Kernel Trick

The covariance function $k(\mathbf{x}, \mathbf{x}')$ can be extended to any kernel functions.

For all $\mathbf{x}$ and $\mathbf{x}'$ in the input space $\mathcal{X}$, the function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is called a kernel function, if there exists a feature map $\varphi : \mathcal{X} \to \mathcal{V}$ that satisfies

$$k\left(\mathbf{x}, \mathbf{x}'\right) = \left\langle \varphi(\mathbf{x}), \varphi\left(\mathbf{x}'\right) \right\rangle_{\mathcal{V}},$$

where $\mathcal{V}$ is a inner product space, and $\langle \cdot, \cdot \rangle_{\mathcal{V}}$ is the inner product operation defined by $\mathcal{V}$.

An alternative definition can be formulated by the *positive semidefinite (PSD)* property: For any points $(\mathbf{x}_1, \ldots, \mathbf{x}_n)$ in $\mathcal{X}$, and all choices of $n$ real-valued coefficients $(c_1, \ldots, c_n)$, the function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is a kernel function, if

$$\sum_{i=1}^{n} \sum_{j=1}^{n} k\left(\mathbf{x}_i, \mathbf{x}_j\right) c_i c_j \geq 0.$$

The spanned matrix, $\boldsymbol{K} \in \mathbb{R}^n \times \mathbb{R}^n$ with its $ij$-th entry $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$, is called Gram matrix.

# Kernel Trick

Many kernels can be chosen for various application scenarios

- Fisher kernel
- Polynomial kernel
- Radial basis function kernel (RBF)
- String kernels
- Graph kernels

The Radial basis function kernel (RBF) is also called squared-exp or Gaussian kernel, which is formulated as

$$k_{\mathrm{SE}}\left(\mathbf{x}_i, \mathbf{x}_j\right) = \sigma_f^2 \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\ell^2}\right).$$

# Gaussian Process Regression

# Prediction with Noise-free Observations

We first consider the ideal case where the observation $y$ can be described by a Gaussian process $f(\mathbf{x}) \sim \mathcal{GP}\left(0, k\left(\mathbf{x}, \mathbf{x}'\right)\right)$ in a noise-free manner, i.e., $y = f(x)$. For a set of test samples with size $n_*$, and feature matrix $\boldsymbol{X}_*$, its function value vector $\mathbf{f}_*$ follows a multivariate Gaussian distribution

$$\mathbf{f}_* \sim \mathcal{N}\left(\mathbf{0}, \boldsymbol{K}\left(\boldsymbol{X}_*, \boldsymbol{X}_*\right)\right),$$

where $\boldsymbol{K}(\cdot, \cdot)$ is the covariance matrix between two collections of input feature vectors.
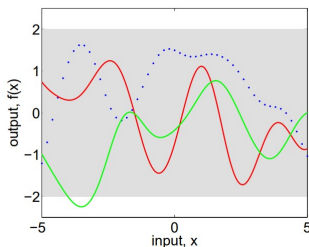
If we are given a set of training samples with size $n$, and feature matrix $\boldsymbol{X}$, its function value vector $\mathbf{f}$ should follow a multivariate Gaussian distribution jointly with $\mathbf{f}_*$

$$\left[\begin{array}{c} \mathbf{f} \\ \mathbf{f}_* \end{array}\right] \sim \mathcal{N}\left(\mathbf{0}, \left[\begin{array}{cc} \boldsymbol{K}(\boldsymbol{X}, \boldsymbol{X}) & \boldsymbol{K}(\boldsymbol{X}, \boldsymbol{X}_*) \\ \boldsymbol{K}(\boldsymbol{X}_*, \boldsymbol{X}) & \boldsymbol{K}(\boldsymbol{X}_*, \boldsymbol{X}_*) \end{array}\right]\right).$$

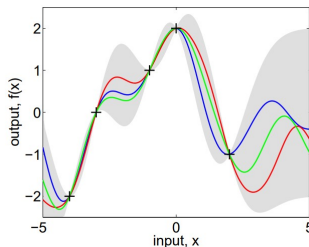## Prediction with Noise-free Observations

Given the observations of **f**, the posterior distribution of $\mathbf{f}_*$ is also a Gaussian distribution

$$\mathbf{f}_* | \boldsymbol{X}_*, \boldsymbol{X}, \mathbf{f} \sim \mathcal{N} \left( \boldsymbol{K}(\boldsymbol{X}_*, \boldsymbol{X}) \, \boldsymbol{K}(\boldsymbol{X}, \boldsymbol{X})^{-1} \mathbf{f} \right.$$
$$\left. \boldsymbol{K}(\boldsymbol{X}_*, \boldsymbol{X}_*) - \boldsymbol{K}(\boldsymbol{X}_*, \boldsymbol{X}) \, \boldsymbol{K}(\boldsymbol{X}, \boldsymbol{X})^{-1} \boldsymbol{K}(\boldsymbol{X}, \boldsymbol{X}_*) \right).$$



(a), prior

(b), posterior

The prior and posterior distributions of $\mathbf{f}_*$ with the RBF kernel function $k(\mathbf{x}_p, \mathbf{x}_q) = \exp(-\frac{1}{2} \|\mathbf{x}_p - \mathbf{x}_q\|_2^2)$.

# Prediction using Noisy Observations

It is typical for more realistic modelling situations that we do not have access to function values themselves, but only noisy versions thereof $y = f(\mathbf{x}) + \varepsilon$.

Assuming additive independent identically distributed Gaussian noise $\varepsilon$ with zero mean and variance $\sigma_n^2$, the prior on the noisy observations becomes another Gaussian distribution with mean $\mathbb{E}[y] = \mathbb{E}[f(\mathbf{x}) + \varepsilon] = 0$, and covariance

$$\mathrm{cov}\left(y_p, y_q\right) = k\left(\mathbf{x}_p, \mathbf{x}_q\right) + \sigma_n^2 \delta_{pq} \text{ or } \mathrm{cov}(\mathbf{y}) = K(\boldsymbol{X}, \boldsymbol{X}) + \sigma_n^2 \boldsymbol{I},$$

where $\delta_{pq}$ is a Kronecker delta symbol which is one if and only if $p = q$ and zero otherwise.

# Gaussian Process Regression

We can write the joint distribution of the observed target values $\mathbf{y}$ and the function values $\mathbf{f}_*$ of test samples under the prior as

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I} & \mathbf{K}(\mathbf{X}, \mathbf{X}_*) \\ \mathbf{K}(\mathbf{X}_*, \mathbf{X}) & \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right).$$
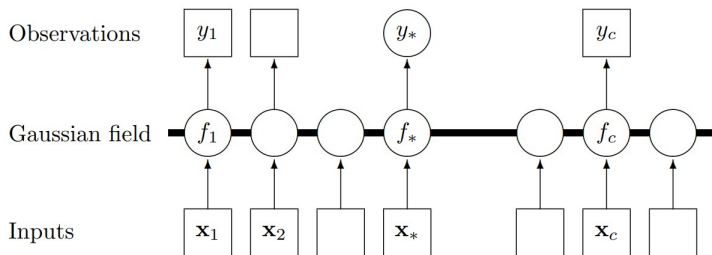
Deriving the conditional distribution, we arrive at the key predictive equations for Gaussian process regression

$$\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_* \sim \mathcal{N} \left( \bar{\mathbf{f}}_*, \operatorname{cov}(\mathbf{f}_*) \right), \text{ where}$$

$$\bar{\mathbf{f}}_* \triangleq \mathbb{E}\left[\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*\right] = \mathbf{K}(\mathbf{X}_*, \mathbf{X}) \left[ \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I} \right]^{-1} \mathbf{y},$$

$$\operatorname{cov}(\mathbf{f}_*) = \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) - \mathbf{K}(\mathbf{X}_*, \mathbf{X}) \left[ \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I} \right]^{-1} \mathbf{K}(\mathbf{X}, \mathbf{X}_*).$$

# Gaussian Process Regression



Graphical model (chain graph) for a Gaussian process for regression.
Squares represent observed variables and circles represent unknowns.

# Gaussian Process Regression

By using a compact form of the notation setting $\boldsymbol{K} = \boldsymbol{K}(\boldsymbol{X}, \boldsymbol{X})$ and $\mathbf{k}_* = \boldsymbol{K}(\boldsymbol{X}, \mathbf{x}_*)$ for a test sample $\mathbf{x}_*$, the prediction $f_*$ for the single test sample $\mathbf{x}_*$ has mean $\bar{f}_*$ and variance $\mathbb{V}[f_*]$ as

$$\bar{f}_* = \mathbf{k}_*^\top \left( \boldsymbol{K} + \sigma_n^2 \boldsymbol{I} \right)^{-1} \mathbf{y},$$

$$\mathbb{V}[f_*] = k\left(\mathbf{x}_*, \mathbf{x}_*\right) - \mathbf{k}_*^\top \left( \boldsymbol{K} + \sigma_n^2 \boldsymbol{I} \right)^{-1} \mathbf{k}_*.$$

Another way to look at the expression of the mean $\bar{f}_*$ is to see it as a linear combination of $N$ kernel functions, each one centered on a training point, by writing

$$\bar{f}_* = \sum_{i=1}^{N} \alpha_i k\left(\mathbf{x}_i, \mathbf{x}_*\right),$$

where $\boldsymbol{\alpha} = \left( \boldsymbol{K} + \sigma_n^2 \boldsymbol{I} \right)^{-1} \mathbf{y}$.

# Gaussian Process Regression

**input**: $X$ (inputs), $\mathbf{y}$ (targets), $k$ (covariance function), $\sigma_n^2$ (noise level), $\mathbf{x}_*$ (test input)

2:  $L := \text{cholesky}(K + \sigma_n^2 I)$

     $\boldsymbol{\alpha} := L^\top \backslash (L \backslash \mathbf{y})$

4:  $\bar{f}_* := \mathbf{k}_*^\top \boldsymbol{\alpha}$ $\qquad\qquad\qquad\qquad \Big\}$ predictive mean

     $\mathbf{v} := L \backslash \mathbf{k}_*$

6:  $\mathbb{V}[f_*] := k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^\top \mathbf{v}$ $\qquad \Big\}$ predictive variance

     $\log p(\mathbf{y}|X) := -\frac{1}{2}\mathbf{y}^\top \boldsymbol{\alpha} - \sum_i \log L_{ii} - \frac{n}{2}\log 2\pi$

8:  **return**: $\bar{f}_*$ (mean), $\mathbb{V}[f_*]$ (variance), $\log p(\mathbf{y}|X)$ (log marginal likelihood)

Algorithm for Gaussian process regression, where the predictions (the target value's mean and variance) at the test sample $\mathbf{x}_*$ and the log marginal likelihood on training set are returned.

# Varying the Hyperparameters

Setting the RBF kernel $k(\mathbf{x}_p, \mathbf{x}_q) = \sigma_f^2 \exp(-\frac{1}{2\ell^2}\|\mathbf{x}_p - \mathbf{x}_q\|_2^2)$ as the covariance function, the covariance between the observations $y_p$ and $y_q$ becomes
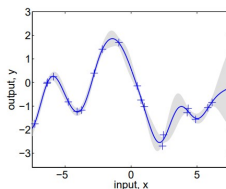
$$\text{cov}\,(y_p, y_q) = \sigma_f^2 \exp(-\frac{1}{2\ell^2}\|\mathbf{x}_p - \mathbf{x}_q\|_2^2) + \sigma_n^2 \delta pq,$$

which implies that the model relies on three hyperparameters:
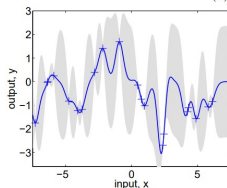
- signal variance $\sigma_f^2$
- noise variance $\sigma_n^2$
- length-scale $\ell$

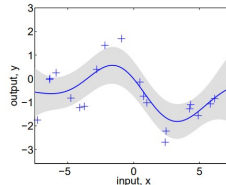The log marginal likelihood can be used to tune the hyperparameters.

# Varying the Hyperparameters



(a), $\ell = 1$

(b), $\ell = 0.3$

(c), $\ell = 3$

Small length-scale $\ell = 0.3 \rightarrow$ a sharply varying function $f$ with small noise.
Large length-scale $\ell = 3 \rightarrow$ a slowly varying function $f$ with large noise.
Moderate length-scale $\ell = 1 \rightarrow$ exact fitting of the underlying function $f$.

# Gaussian Process Classification

# Binary Gaussian Process Classifier

For binary classification with label candidate set $\{+1, -1\}$, linear regression uses the sigmoid function to squash the linear regression output into the label probability $p(y = +1|\mathbf{x}) = \sigma(\mathbf{w}^{\mathrm{T}}\phi(\mathbf{x}))$.

In Gaussian process, the latent function value $f(\mathbf{x})$ is used to construct the linear regression output.

Similarly, for Gaussian process classification, we can also construct the label probability via the sigmoid function

$$p(y = +1|f(\mathbf{x})) = \sigma(f(\mathbf{x})).$$

Different from $\mathbf{w}^{\mathrm{T}}\phi(\mathbf{x})$, $f(\mathbf{x})$ is Gaussian random variable, we need to marginalize it to get the posterior probability

$$p(y = +1|\mathbf{x}) = \int p(y = +1|f)p(f|\mathbf{x})\mathrm{d}f.$$

# Binary Gaussian Process Classifier

Given a training set $\{\boldsymbol{X}, \mathbf{y}\}$, for a new test sample $\mathbf{x}_*$, Gaussian process classifier predicts its label by estimating the posterior probability

$$\bar{\pi}_* \triangleq p(y_* = +1 | \boldsymbol{X}, \mathbf{y}, \mathbf{x}_*) = \int p(y_* = +1 | f_*) \, p(f_* | \boldsymbol{X}, \mathbf{y}, \mathbf{x}_*) \, \mathrm{d}f_*$$

$$= \int \sigma(f_*) p(f_* | \boldsymbol{X}, \mathbf{y}, \mathbf{x}_*) \, \mathrm{d}f_*.$$

To achieve the estimation, we need first calculate the posterior distribution of the latent variable $f_*$ corresponding to the test sample $\mathbf{x}_*$

$$p(f_* | \boldsymbol{X}, \mathbf{y}, \mathbf{x}_*) = \int p(f_* | \boldsymbol{X}, \mathbf{x}_*, \mathbf{f}) \, p(\mathbf{f} | \boldsymbol{X}, \mathbf{y}) \mathrm{d}\mathbf{f},$$

where $p(\mathbf{f} | \boldsymbol{X}, \mathbf{y}) = p(\mathbf{y} | \mathbf{f}) p(\mathbf{f} | \boldsymbol{X}) / p(\mathbf{y} | \boldsymbol{X})$ is the posterior over the latent variables and $p(\mathbf{y} | \mathbf{f})$ is computed by applying the elementwise sigmoid function to $\mathbf{f}$.

# Laplace Approximation for Binary GP Classifier

In the regression case, $\mathbf{y}$ and $f_*$ jointly follow a multivariate Gaussian distribution, the conditional probability distribution $p(f_*|\mathbf{X}, \mathbf{y}, \mathbf{x}_*)$ can be computed analytically.

However, in the classification case, due to the computational intractability of the distribution $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$, the posterior distribution $p(f_*|\mathbf{X}, \mathbf{y}, \mathbf{x}_*)$ cannot be computed analytically with the intractable integral.

To solve this problem, we resort to the Laplace approximation that approximates the distribution $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ with a Gaussian distribution $q(\mathbf{f}|\mathbf{X}, \mathbf{y})$, obtained by doing a second order Taylor expansion of $\log p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ around the maximum of the posterior

$$q(\mathbf{f}|\mathbf{X}, \mathbf{y}) = \mathcal{N}\left(\mathbf{f}|\hat{\mathbf{f}}, \mathbf{A}^{-1}\right) \propto \exp\left(-\frac{1}{2}(\mathbf{f} - \hat{\mathbf{f}})^\top \mathbf{A}(\mathbf{f} - \hat{\mathbf{f}})\right),$$

where $\hat{\mathbf{f}} = \arg\max_{\mathbf{f}} p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ and $\mathbf{A} = -\nabla\nabla \log p(\mathbf{f}|\mathbf{X}, \mathbf{y})|_{\mathbf{f}=\hat{\mathbf{f}}}$ is the Hessian matrix of the negative log posterior at that point.

# Laplace Approximation for Binary GP Classifier

By Bayes' rule, the posterior over the latent variables $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ is given by

$$p(\mathbf{f}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}, \mathbf{f}|\mathbf{X})}{p(\mathbf{y}|\mathbf{X})} = \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{X})}{p(\mathbf{y}|\mathbf{X})},$$

but $p(\mathbf{y}|\mathbf{X})$ is independent of $\mathbf{f}$, we only need to consider the un-normalized posterior $p(\mathbf{y}, \mathbf{f}|\mathbf{X}) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{X})$ when maximizing w.r.t. $\mathbf{f}$.

Taking the logarithm on the un-normalized posterior, we get

$$\begin{aligned}
\mathbf{\Psi}(\mathbf{f}) &\triangleq \log p(\mathbf{y}|\mathbf{f}) + \log p(\mathbf{f}|\mathbf{X}) \\
&= \log p(\mathbf{y}|\mathbf{f}) - \frac{1}{2}\mathbf{f}^{\mathrm{T}}\mathbf{K}^{-1}\mathbf{f} - \frac{1}{2}\log|\mathbf{K}| - \frac{n}{2}\log 2\pi,
\end{aligned}$$

where $p(\mathbf{f}|\mathbf{X})$ is a Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{K})$, $\mathbf{K}$ is short for the kernel matrix $\mathbf{K}(\mathbf{X}, \mathbf{X})$, and $|\mathbf{K}|$ is the determinant of the kernel matrix $\mathbf{K}$.

## Laplace Approximation for Binary GP Classifier

Differentiating $\Psi(\mathbf{f})$ w.r.t. $\mathbf{f}$, we obtain

$$\nabla \Psi(\mathbf{f}) = \nabla \log p(\mathbf{y}|\mathbf{f}) - \mathbf{K}^{-1}\mathbf{f},$$
$$\nabla\nabla \Psi(\mathbf{f}) = \nabla\nabla \log p(\mathbf{y}|\mathbf{f}) - \mathbf{K}^{-1},$$

where $\nabla\nabla \log p(\mathbf{y}|\mathbf{f})$ is diagonal, since the distribution for $y_i = +1$ or $-1$ depends only on $f_i$, not on $f_{j \neq i}$.

We can apply the Newton-Raphson method to find the latent function value vector $\hat{\mathbf{f}}$ that maximizes $\Psi(\mathbf{f})$, and obtain the approximated probability distribution as

$$q(\mathbf{f}|\mathbf{X}, \mathbf{y}) = \mathcal{N}\left(\hat{\mathbf{f}}, \left(\mathbf{K}^{-1} + \mathbf{W}\right)^{-1}\right),$$

where $\mathbf{W} \triangleq -\nabla\nabla \log p(\mathbf{y}|\mathbf{f})|_{\mathbf{f}=\hat{\mathbf{f}}}$.

# Laplace Approximation for Binary GP Classifier

With $q(\mathbf{f}|\mathbf{X}, \mathbf{y})$ approximating $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$, we can get the approximation for the posterior $p(f_*|\mathbf{X}, \mathbf{y}, \mathbf{x}_*)$ as

$$q(f_*|\mathbf{X}, \mathbf{y}, \mathbf{x}_*) = \int p(f_*|\mathbf{X}, \mathbf{x}_*, \mathbf{f}) \, q(\mathbf{f}|\mathbf{X}, \mathbf{y}) \mathrm{d}\mathbf{f},$$

where both $p(f_*|\mathbf{X}, \mathbf{x}_*, \mathbf{f})$ and $q(\mathbf{f}|\mathbf{X}, \mathbf{y})$ are Gaussian distributions. It can be proved that $q(f_*|\mathbf{X}, \mathbf{y}, \mathbf{x}_*)$ is also a Gaussian distribution with mean and variance as

$$\mathbb{E}_q[f_*|\mathbf{X}, \mathbf{y}, \mathbf{x}_*] = \mathbf{k}_*^{\mathrm{T}} \boldsymbol{K}^{-1} \hat{\mathbf{f}},$$

$$\mathbb{V}_q[f_*|\mathbf{X}, \mathbf{y}, \mathbf{x}_*] = k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^{\mathrm{T}} \left(\boldsymbol{K} + \boldsymbol{W}^{-1}\right)^{-1} \mathbf{k}_*.$$

Finally, the label probability for the test sample $\mathbf{x}_*$ can computed as

$$\bar{\pi}_* \triangleq p(y_* = +1|\mathbf{X}, \mathbf{y}, \mathbf{x}_*) := \int \sigma(f_*) q(f_*|\mathbf{X}, \mathbf{y}, \mathbf{x}_*) \, \mathrm{d}f_*.$$

Numeric integration techniques are generally used to calculate the one-dimensional integral.

# Binary Gaussian Process Classifier

**input**: $K$ (covariance matrix), $\mathbf{y}$ ($\pm 1$ targets), $p(\mathbf{y}|\mathbf{f})$ (likelihood function)
2: $\mathbf{f} := \mathbf{0}$               initialization
   **repeat**               Newton iteration
4:     $W := -\nabla\nabla \log p(\mathbf{y}|\mathbf{f})$
      $L := \text{cholesky}(I + W^{\frac{1}{2}} K W^{\frac{1}{2}})$
6:     $\mathbf{b} := W\mathbf{f} + \nabla \log p(\mathbf{y}|\mathbf{f})$
      $\mathbf{a} := \mathbf{b} - W^{\frac{1}{2}} L^{\top} \backslash (L \backslash (W^{\frac{1}{2}} K \mathbf{b}))$
8:     $\mathbf{f} := K\mathbf{a}$
   **until** convergence
10: $\log q(\mathbf{y}|X, \theta) := -\frac{1}{2}\mathbf{a}^{\top}\mathbf{f} + \log p(\mathbf{y}|\mathbf{f}) - \sum_i \log L_{ii}$
   **return**: $\hat{\mathbf{f}} := \mathbf{f}$ (post. mode), $\log q(\mathbf{y}|X, \theta)$ (approx. log marg. likelihood)

Training algorithm for the binary Gaussian process classifier with Laplace approximation. The mode $\hat{\mathbf{f}}$ for the posterior probability distribution $p(\mathbf{f}|\boldsymbol{X}, \mathbf{y})$ is returned, as well as the approximated log marginal likelihood, which can be use to tune hyperparameters.

# Binary Gaussian Process Classifier

> **input**: $\hat{\mathbf{f}}$ (mode), $X$ (inputs), $\mathbf{y}$ ($\pm 1$ targets), $k$ (covariance function),
> $\qquad\qquad\qquad\qquad p(\mathbf{y}|\mathbf{f})$ (likelihood function), $\mathbf{x}_*$ test input
>
> 2: $W := -\nabla\nabla \log p(\mathbf{y}|\hat{\mathbf{f}})$
> $\quad L := \text{cholesky}(I + W^{\frac{1}{2}} K W^{\frac{1}{2}})$
> 4: $\bar{f}_* := \mathbf{k}(\mathbf{x}_*)^\top \nabla \log p(\mathbf{y}|\hat{\mathbf{f}})$
> $\quad \mathbf{v} := L \backslash \left(W^{\frac{1}{2}} \mathbf{k}(\mathbf{x}_*)\right)$
> 6: $\mathbb{V}[f_*] := k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^\top \mathbf{v}$
> $\quad \bar{\pi}_* := \int \sigma(z) \mathcal{N}(z|\bar{f}_*, \mathbb{V}[f_*]) dz$
> 8: **return**: $\bar{\pi}_*$ (predictive class probability (for class 1))

Prediction algorithm for the binary Gaussian classifier with Laplace approximation.

## Multi-Class Gaussian Process Classifier

For the multi-class classification problem, where the class label $y$ belongs to the label candidate set $\{1, \cdots, C\}$ with size $C$, we need to construct a Gaussian process for each class $c \in \{1, \cdots, C\}$, with the latent function value $f^c(\mathbf{x})$ evaluated at $\mathbf{x}$.

The conditional label probability $p\left(y = c | f^1(\mathbf{x}), \cdots, f^C(\mathbf{x})\right)$ can be constructed by the softmax function

$$p\left(y = c | f^1(\mathbf{x}), \cdots, f^C(\mathbf{x})\right) = \frac{\exp\left(f^c(\mathbf{x})\right)}{\sum_{c'=1}^{C} \exp\left(f^{c'}(\mathbf{x})\right)}.$$

Following the same procedure as the binary classification case, given a test sample $\mathbf{x}_*$, we can estimate the posterior probabilities $p(y_* = c | \boldsymbol{X}, \mathbf{y}, \mathbf{x}_*)$ to predict $\mathbf{x}_*$'s class label.

# References

Williams CK, Rasmussen CE. **Gaussian Processes for Machine Learning**. Cambridge, MA: MIT press; 2006.

- Section 1 - Section 3

# The End