

STA3010 Regression Analysis

Feng Yin

The Chinese University of Hong Kong (Shenzhen)

yinfeng@cuhk.edu.cn

May 7, 2020

Overview

- 1 Background
- 2 Kernel Methods
 - Dual Representation of RLS
 - Constructing Valid Kernels
- 3 Gaussian Processes (GP) for Machine Learning
 - GP for Regression
 - GP for Classification
 - Hyper-parameter Optimization
 - Two GP Examples
- 4 Relationships with Other Models
- 5 Summary and Selected Recent Advances
- 6 Appendix

Outline

- 1 Background
- 2 Kernel Methods
 - Dual Representation of RLS
 - Constructing Valid Kernels
- 3 Gaussian Processes (GP) for Machine Learning
 - GP for Regression
 - GP for Classification
 - Hyper-parameter Optimization
 - Two GP Examples
- 4 Relationships with Other Models
- 5 Summary and Selected Recent Advances
- 6 Appendix

Background

We start with a general regression model

$$y_i = f(\mathbf{x}_i) + e_i, \quad i = 1, 2, \dots, n, \quad (1)$$

where

- n is the number of the data points;
- \mathbf{x}_i is the i -th input of dimension p ;
- y_i is the i -th output (assumed to be scalar);
- $f(\mathbf{x}) : \mathbb{R}^p \rightarrow \mathbb{R}$ is the underlying regression function to be approximated;
- e_i is an error term.

Background

In general, a machine learning algorithm involves the following two phases:

Learning phase: Assume a (parametric or nonparametric) model of $f(\mathbf{x})$ and train the model given a set of training data, denoted by $\mathcal{D} \triangleq \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\} = \{\mathbf{X}, \mathbf{y}\}$.

Prediction phase: Given a novel data point with input \mathbf{x}_* , we aim to predict the output y_* .

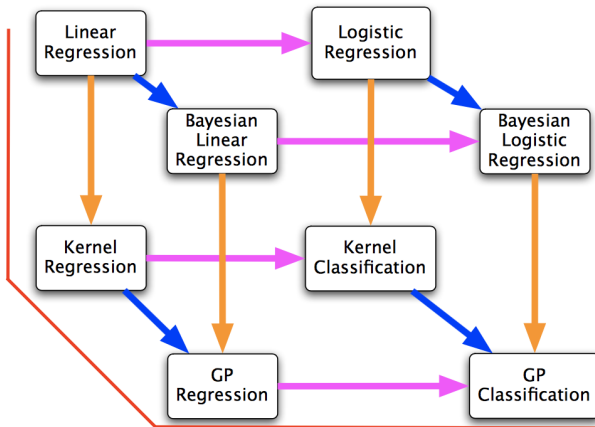
Two options prevail:

- 1 Discard the training data and make predictions purely based on the model that has been trained \rightarrow **parametric model**.
- 2 Use the training data during the prediction phase \rightarrow **nonparametric model**.

In this presentation, we mainly focus on the **second option**.

Background

“Roadmap” of this presentation:



Outline

- 1 Background
- 2 **Kernel Methods**
 - Dual Representation of RLS
 - Constructing Valid Kernels
- 3 Gaussian Processes (GP) for Machine Learning
 - GP for Regression
 - GP for Classification
 - Hyper-parameter Optimization
 - Two GP Examples
- 4 Relationships with Other Models
- 5 Summary and Selected Recent Advances
- 6 Appendix

Regularized Least-Squares in Canonical Form

Let us approximate the underlying function by a linear regression model

$$f(\mathbf{x}_i) = \mathbf{w}^T \phi(\mathbf{x}_i), \quad (2)$$

where \mathbf{w} is an $m \times 1$ vector of weighting factors of a fixed nonlinear feature space mapping $\phi(\mathbf{x}_i) : \mathbb{R}^p \rightarrow \mathbb{R}^m$, which returns a vector of m basis function values evaluated at point \mathbf{x}_i .

The classic **regularized least-squares (RLS)** method for solving the weighting factors \mathbf{w} can be expressed as:

$$\mathbf{w}_{RLS} = \arg \min_{\mathbf{w}} J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^n \left(\mathbf{w}^T \phi(\mathbf{x}_i) - y_i \right)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}, \quad (3)$$

where λ is the regularization parameter.

Regularized Least-Squares in Dual Form

The point estimate of the weighting factors can be formulated as

$$\mathbf{w} = \sum_{i=1}^n a_i \phi(\mathbf{x}_i) = \Phi^T \mathbf{a}, \quad (4)$$

where Φ is the design matrix of size $n \times m$, whose i th row is given by $\phi(\mathbf{x}_i)^T$ and the vector $\mathbf{a} = [a_1, a_2, \dots, a_n]^T$ with the i th entry defined by

$$a_i \triangleq -\frac{1}{\lambda} \left(\mathbf{w}^T \phi(\mathbf{x}_i) - y_i \right). \quad (5)$$

Regularized Least-Squares in Dual Form

The dual representation of the above RLS can be obtained by substituting $\mathbf{w} = \Phi^T \mathbf{a}$ into the original cost function $J(\mathbf{w})$ and expanding the square term, namely,

$$J(\mathbf{a}) \triangleq \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{y} + \frac{1}{2} \mathbf{y}^T \mathbf{y} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}, \quad (6)$$

where $\mathbf{y} = [y_1, y_2, \dots, y_n]^T$ and $\mathbf{K} = \Phi \Phi^T$ is called the **Gram matrix or kernel matrix**, which is an $n \times n$ symmetric matrix with the (i, j) -th entry (the i th row and j th column of \mathbf{K})

$$[\mathbf{K}]_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) \triangleq \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j), \quad (7)$$

where the term **kernel function** $k(\mathbf{x}, \mathbf{x}')$ is introduced eventually.

Regularized Least-Squares in Dual Form

Taking the derivative of $J(\mathbf{a})$ w.r.t. \mathbf{a} and set it equal to zero, yields

$$\mathbf{a}_{RLS} = (\mathbf{K} + \lambda \mathbf{I}_n)^{-1} \mathbf{y}. \quad (8)$$

Inserting this result back into the linear parametric regression model, we obtain the following prediction for a new input \mathbf{x}_* as follows:

$$y_* = \phi(\mathbf{x}_*)^T \mathbf{w}_{RLS} = \phi(\mathbf{x}_*)^T \Phi^T \mathbf{a}_{RLS} = \mathbf{k}(\mathbf{x}_*)^T (\mathbf{K} + \lambda \mathbf{I}_n)^{-1} \mathbf{y} \quad (9)$$

where we have defined the vector $\mathbf{k}(\mathbf{x}_*)$ with entries $k_i(\mathbf{x}_*) \triangleq k(\mathbf{x}_i, \mathbf{x}_*)$.

Remark: The dual form rewrites the RLS solution entirely in terms of the kernel function $k(\mathbf{x}, \mathbf{x}')$.

Regularized Least-Squares in Dual Form

The cons and pros of the dual representation are the following:

- **Pros:** The prediction **works directly with the kernel function and avoids the selection of fixed basis functions $\phi(\mathbf{x})$** , which allows us to use feature spaces of higher, even infinite, dimension.
- **Cons:** In the dual representation, we need to **invert an $n \times n$ matrix** for computing **\mathbf{a}** rather than **invert a smaller $m \times m$ matrix** for **\mathbf{w}** in the original problem.

In summary, larger model flexibility is gained with higher computational complexity.

Constructing Valid Kernels

According to [J. Shawe-Taylor and N. Cristianini, 04], a **necessary and sufficient condition** for a function $k(\mathbf{x}, \mathbf{x}')$ to be a valid kernel is that the corresponding kernel matrix \mathbf{K} should be **positive semidefinite (PSD)** for all possible choices of $\mathbf{x} \in \mathcal{X}$.

For clarity, we give the kernel matrix \mathbf{K} explicitly as follows:

$$\mathbf{K}(\mathbf{X}, \mathbf{X}; \boldsymbol{\theta}_h) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1; \boldsymbol{\theta}_h) & k(\mathbf{x}_1, \mathbf{x}_2; \boldsymbol{\theta}_h) & \dots & k(\mathbf{x}_1, \mathbf{x}_n; \boldsymbol{\theta}_h) \\ k(\mathbf{x}_2, \mathbf{x}_1; \boldsymbol{\theta}_h) & k(\mathbf{x}_2, \mathbf{x}_2; \boldsymbol{\theta}_h) & \dots & k(\mathbf{x}_2, \mathbf{x}_n; \boldsymbol{\theta}_h) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1; \boldsymbol{\theta}_h) & k(\mathbf{x}_n, \mathbf{x}_2; \boldsymbol{\theta}_h) & \dots & k(\mathbf{x}_n, \mathbf{x}_n; \boldsymbol{\theta}_h) \end{bmatrix}, \quad (10)$$

where $\boldsymbol{\theta}_h$ is a set of parameters that control the characteristic of the kernel function. $\boldsymbol{\theta}_h$ is often called **hyper-parameters** and later \mathbf{K} or $\mathbf{K}(\mathbf{X}, \mathbf{X})$ is short for $\mathbf{K}(\mathbf{X}, \mathbf{X}; \boldsymbol{\theta}_h)$.

Constructing Valid Kernels

We have two options for constructing valid kernels:

- 1 **Naive solution**: Choose a feature space mapping $\phi(\mathbf{x})$ and use the definition

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') = \sum_{i=1}^m \phi_i(\mathbf{x}) \phi_i(\mathbf{x}'), \quad (11)$$

where $\phi_i(\mathbf{x})$, $i = 1, 2, \dots, m$ are the basis functions.

- 2 **More elegant solution**: Construct valid kernel functions without having to construct the feature space mapping $\phi(\mathbf{x})$ first. **Build complex but valid kernels out of elementary kernels under valid operations.**

Constructing Valid Kernels

According to [C. Bishop, 06], some fundamental properties can be exploited for constructing a new valid kernel $k(\mathbf{x}, \mathbf{x}')$, they are:

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})f(\mathbf{x}') \quad (12a)$$

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad (12b)$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \quad (12c)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(k_1(\mathbf{x}, \mathbf{x}')\right) \quad (12d)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \quad (12e)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') \cdot k_2(\mathbf{x}, \mathbf{x}') \quad (12f)$$

where $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$ are both known valid kernels; $f(\mathbf{x}) : \mathbb{R}^p \rightarrow \mathbb{R}$ is any function; $c \geq 0$ is a constant.

Elementary Kernels: SE Kernel

Squared exponential (SE) kernel (a.k.a. radial basis function kernel or Gaussian kernel), has the form:

$$k_{SE}(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2l^2} \right). \quad (13)$$

The SE kernel is often deemed as the default kernel for Gaussian Processes and SVM.

The set of hyper-parameters to be tuned is $\theta_h = [\sigma^2, l]^T$, where

- The length scale l determines the length of the 'wiggles' in the function.
- The variance σ^2 determines the average distance of the function away from its mean.

Elementary Kernels: SE Kernel

- **Simple task I:** Show that the SE kernel is a valid kernel.

Answer:

- 1 Expand $\|\mathbf{x} - \mathbf{x}'\|^2$ to be $\mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T \mathbf{x}' + (\mathbf{x}')^T \mathbf{x}'$.
 - 2 Express $k_{SE}(\mathbf{x}, \mathbf{x}') = \sigma^2 \cdot f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$, where $\sigma^2 > 0$, $l^2 > 0$,
 $f(\mathbf{x}) = \exp\left[\frac{-\mathbf{x}^T \mathbf{x}}{2l^2}\right]$ and $k_1(\mathbf{x}, \mathbf{x}') = \exp\left[\frac{\mathbf{x}^T \mathbf{x}'}{l^2}\right]$.
 - 3 $k_2(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$, i.e., the linear kernel is a valid kernel, and so is $k_2(\mathbf{x}, \mathbf{x}')/l^2$.
 - 4 Apply the properties shown beforehand.
- **Simple task II:** Show that the SE kernel corresponds to infinite number of basis function, i.e., the dimension of $\phi(\mathbf{x})$ is infinite.

Answer: Use Taylor expansion, details can be found in appendix.

Elementary Kernels: RQ Kernel

Rational quadratic (RQ) kernel is of the form:

$$k_{RQ}(\mathbf{x}, \mathbf{x}') = \sigma^2 \left(1 + \frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\alpha l^2} \right)^{-\alpha}. \quad (14)$$

This kernel can be seen as a scale mixture (an infinite sum) of SE kernel functions with different characteristic length-scales. When $\alpha \rightarrow \infty$, RQ kernel \rightarrow SE kernel, see Appendix.

The set of hyper-parameters to be tuned is $\boldsymbol{\theta}_h = [\sigma^2, l, \alpha]^T$.

Elementary Kernels: Periodic Kernel

Periodic kernel is of the form:

$$k_{PER}(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left(-\frac{\sin^2 \left(\pi \|\mathbf{x} - \mathbf{x}'\| / p \right)}{l^2} \right), \quad (15)$$

which allows one to **model functions which repeat themselves exactly**.

The set of hyper-parameters to be tuned is $\theta_h = [\sigma^2, p, l]^T$, where

- The period p simply determines the distance between repetitions of the function.

Elementary Kernels: Local Periodic Kernel

Locally periodic kernel is of the form:

$$k_{LP}(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left(-\frac{\sin^2 \left(\pi \|\mathbf{x} - \mathbf{x}'\| / p \right)}{l_1^2} \right) \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2l_2^2} \right). \quad (16)$$

An SE kernel times a periodic kernel results in functions which are periodic, but can slowly vary over time.

The set of hyper-parameters to be tuned is $\theta_h = [\sigma^2, p, l_1, l_2]^T$, where

- The length scales l_1 and l_2 together determine the length of the 'wiggles' in the function.

Stationary and Isotropic Kernels

Before we move on, let us quickly recap some background knowledge:

- **Stationary kernel** function to be a function of $\boldsymbol{\tau} = \mathbf{x} - \mathbf{x}'$.
- **Isotropic kernel** function to be a function of $r = \|\mathbf{x} - \mathbf{x}'\|$.

It is obvious that the above four examples are not only stationary but also isotropic.

Stationary kernel and its spectral density are Fourier duals:

$$k(\boldsymbol{\tau}) = \int S(\mathbf{s}) e^{j2\pi \mathbf{s}^T \boldsymbol{\tau}} d\boldsymbol{\tau} \quad (17)$$

$$S(\mathbf{s}) = \int k(\boldsymbol{\tau}) e^{-j2\pi \mathbf{s}^T \boldsymbol{\tau}} d\boldsymbol{\tau} \quad (18)$$

Example: $\tau = t$ representing time and $s = f$ representing frequency.

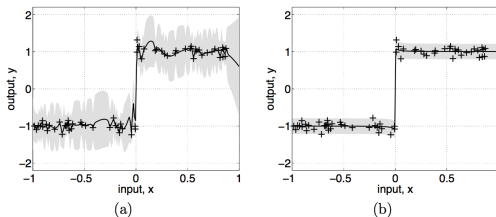
Advanced Kernels: Neural Network Kernel

Lastly, we introduce a famous **non-stationary kernel**, namely the neural network kernel:

$$k_{NN}(\mathbf{x}, \mathbf{x}') = \frac{2}{\pi} \sin^{-1} \left(\frac{2\tilde{\mathbf{x}}\Sigma\tilde{\mathbf{x}}'}{\sqrt{(1 + 2\tilde{\mathbf{x}}\Sigma\tilde{\mathbf{x}})(1 + 2\tilde{\mathbf{x}}'\Sigma\tilde{\mathbf{x}}')}} \right), \quad (19)$$

where $\tilde{\mathbf{x}} \triangleq [1, \mathbf{x}^T]^T$ is an augmented input vector. Often $\Sigma = \text{diag}(\sigma_1^2, \sigma_2^2)$ is a diagonal matrix and thus $\boldsymbol{\theta}_h = [\sigma_1^2, \sigma_2^2]^T$.

Non-stationary kernel is capable of providing better modeling in terms of abrupt change and saturation effect in the data, see the following example.



- ① C. Bishop, Pattern Recognition and Machine Learning, Springer, 2006.
- ② C. Rasmussen, Gaussian Process for Machine Learning, MIT press, 2006.
- ③ J. Shawe-Taylor and N. Cristianini, Kernel Methods for Pattern Analysis, Cambridge University Press, 2004.
- ④ A. G. Wilson, Covariance Kernels for Fast Automatic Pattern Discovery and Extrapolation with Gaussian Processes, Ph.D. Thesis, University of Cambridge, 2014.
- ⑤ D. Duvenaud, The Kernel Cookbook.
<http://www.cs.toronto.edu/~duvenaud/cookbook/>

Outline

- 1 Background
- 2 Kernel Methods
 - Dual Representation of RLS
 - Constructing Valid Kernels
- 3 **Gaussian Processes (GP) for Machine Learning**
 - GP for Regression
 - GP for Classification
 - Hyper-parameter Optimization
 - Two GP Examples
- 4 Relationships with Other Models
- 5 Summary and Selected Recent Advances
- 6 Appendix

Parametric vs Nonparametric Models [Z. Ghahramani, 15]

- **Parametric models** assume some finite set of parameters \mathbf{w} . Given the parameters, future prediction, y_* , is independent of the observed data, therefore \mathbf{w} captures everything there is to know about the data.
- **Nonparametric models** assume that the data distribution cannot be defined in terms of such a finite set of parameters. But they can be defined by assuming an **infinite dimensional \mathbf{w}** . Usually we think of \mathbf{w} as a function. The amount of information that \mathbf{w} can capture about the observed data can grow as the amount of data grows.
- **Nonparametric models**: A really large parametric model in practice [Teh, 11].

Bayesian Nonparametric Models

A Bayesian nonparametric model is a [Bayesian model](#) introduced on an [infinite-dimensional parameter space](#).

The parameter space is typically chosen as the set of all possible solutions for a given learning problem, e.g.,

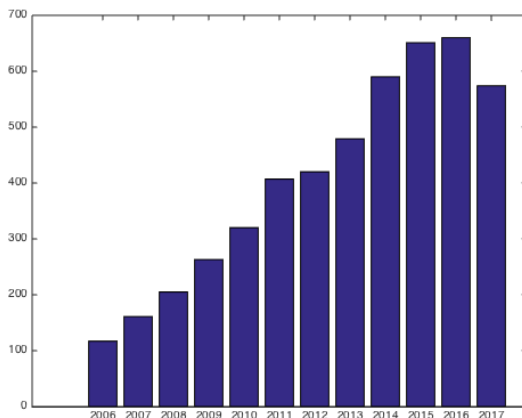
- for [regression problem](#), the parameter space can be the set of continuous functions;
- for [density estimation problem](#), the parameter space can consist of all densities.

[Examples:](#)

Parametric	Non-parametric	Application
polynomial regression	Gaussian processes	function approx.
logistic regression	Gaussian process classifiers	classification
mixture models, k-means	Dirichlet process mixtures	clustering

Some examples of Bayesian non-parametric models [Z. Ghahramani, 15].

Gaussian Process (GP)

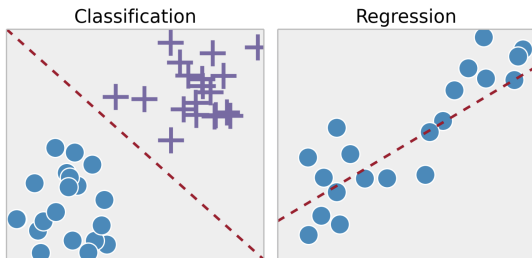


Number of publications on GP in engineering fields,
<https://www.engineeringvillage.com>

Gaussian Process

In the following, Gaussian process will be viewed from a different angle and used as advanced Bayesian tools for

- nonlinear regression
- classification
- clustering (less seen)



Definition [Rasmussen and Williams, 06]

A Gaussian process is a collection of random variables, any finite number of which have Gaussian distributions.

We focus on real-valued Gaussian process, which is completely specified by a mean function and a covariance function (a.k.a. kernel function), i.e.,

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}_h)), \quad (20)$$

where

- $m(\mathbf{x})$ is the mean function, often set to zero in practice, especially when there is no prior knowledge available.
- $k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}_h)$ is the covariance/kernel function controlled by the hyper-parameters $\boldsymbol{\theta}_h$.

Gaussian Process for Regression

We consider the general regression model

$$y = f(\mathbf{x}) + e, \quad (21)$$

where

- the output y is continuous-valued;
- the underlying regression function $f(\mathbf{x})$ is assumed to be a real-valued zero-mean Gaussian process;
- and additive i.i.d. Gaussian noise with variance σ^2 is assumed;
- the noise terms are mutually independent.

The set of all unknown parameters is denoted by $\boldsymbol{\theta} = [\boldsymbol{\theta}_h, \sigma^2]^T$.

Gaussian Process for Regression

Given a training dataset $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$ where $\mathbf{y} = [y_1, y_2, \dots, y_n]^T$ contains the training outputs and $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$ contains the training inputs, we aim to predict $\mathbf{y}_* = [y_{*,1}, y_{*,2}, \dots, y_{*,n_*}]^T$ given test inputs $\mathbf{X}_* = [\mathbf{x}_{*,1}, \mathbf{x}_{*,2}, \dots, \mathbf{x}_{*,n_*}]$ with a posterior distribution $p(\mathbf{y}_* | \mathcal{D}, \mathbf{X}_*; \theta)$.

Special note on the notation: sometimes we write \mathbf{X} and/or \mathbf{x}_* in the conditional probabilities just to clarify which input correspond to which output, but they are essentially deterministic quantities and can be omitted in the conditional probability expressions.

Why GP is a Bayesian Nonparametric Model?

- A Gaussian prior is selected for deriving the posterior distribution of a desired metric, where Bayes theorem naturally comes into play.
- $f(\mathbf{x})$ can be seen as an infinite dimensional parameters \mathbf{w} .
- As the number of the observed data (\mathbf{x}_i, y_i) grows, we get to know the model $f(\mathbf{x})$ better.
- Consequently, the posterior distribution out of the GP model is refined as the number of observed data grows!

Gaussian Process for Regression

In light of the definition of Gaussian process given beforehand, the joint prior distribution of the **training output** \mathbf{y} and **test output** \mathbf{y}_* can be written explicitly as:

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{y}_* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_n & \mathbf{K}(\mathbf{X}, \mathbf{X}_*) \\ \mathbf{K}(\mathbf{X}_*, \mathbf{X}) & \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) + \sigma^2 \mathbf{I}_{n_*} \end{bmatrix} \right), \quad (22)$$

where

- $\mathbf{K}(\mathbf{X}, \mathbf{X}_*)$ is an $n \times n_*$ matrix of the covariance evaluated at all pairs of training and test points.
- $\mathbf{K}(\mathbf{X}_*, \mathbf{X})$, $\mathbf{K}(\mathbf{X}, \mathbf{X})$, $\mathbf{K}(\mathbf{X}_*, \mathbf{X}_*)$ that are of size $n_* \times n$, $n \times n$, and $n_* \times n_*$, respectively.
- $\mathbf{K}(\mathbf{X}_*, \mathbf{X}) = \mathbf{K}(\mathbf{X}, \mathbf{X}_*)^T$.

Gaussian Process for Regression

To derive the posterior distribution $p(\mathbf{y}_*|\mathcal{D}, \mathbf{X}_*, \boldsymbol{\theta})$, we require:

Conditional Gaussian Distribution

If the two random variables \mathbf{x} and \mathbf{y} are jointly multivariate Gaussian as

$$p(\mathbf{x}, \mathbf{y}) \sim \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}\right), \quad (23)$$

then it is easy to derive the following conditional probabilities

$$p(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}(\boldsymbol{\mu}_{x|y}, \Sigma_{x|y}), \quad p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}(\boldsymbol{\mu}_{y|x}, \Sigma_{y|x}), \quad (24)$$

where

$$\boldsymbol{\mu}_{x|y} = \boldsymbol{\mu}_x + \Sigma_{xy}\Sigma_{yy}^{-1}(\mathbf{y} - \boldsymbol{\mu}_y), \quad \boldsymbol{\mu}_{y|x} = \boldsymbol{\mu}_y + \Sigma_{yx}\Sigma_{xx}^{-1}(\mathbf{x} - \boldsymbol{\mu}_x), \quad (25)$$

$$\Sigma_{x|y} = \Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx}, \quad \Sigma_{y|x} = \Sigma_{yy} - \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy}. \quad (26)$$

Gaussian Process for Regression

Applying the above result to our GP regression model, we thus obtain

$$p(\mathbf{y}_*|\mathcal{D}, \mathbf{X}_*; \boldsymbol{\theta}_h) \sim \mathcal{N}(\bar{\mathbf{m}}, \bar{\mathbf{V}}), \quad (27)$$

where

$$\bar{\mathbf{m}} \triangleq \mathbf{K}(\mathbf{X}_*, \mathbf{X}) [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_n]^{-1} \mathbf{y}, \quad (28)$$

$$\bar{\mathbf{V}} \triangleq \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) + \sigma^2 \mathbf{I}_{n_*} - \mathbf{K}(\mathbf{X}_*, \mathbf{X}) [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_n]^{-1} \mathbf{K}(\mathbf{X}, \mathbf{X}_*). \quad (29)$$

Remark: for ease of our narration on the inference, the parameters $\boldsymbol{\theta}$ is assumed to be known! Later, we will mention how to tune these parameters.

Gaussian Process for Regression

For the simplest case where there is only one test input \mathbf{x}_* , the posterior distribution is given by

$$p(y_*|\mathcal{D}, \mathbf{x}_*; \boldsymbol{\theta}) \sim \mathcal{N}(\bar{m}, \bar{v}). \quad (30)$$

where the posterior mean and variance are respectively

$$\bar{m} \triangleq \mathbf{k}_* [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_n]^{-1} \mathbf{y}, \quad (31)$$

$$\bar{v} \triangleq k(\mathbf{x}_*, \mathbf{x}_*) + \sigma^2 - \mathbf{k}_*^T [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_n]^{-1} \mathbf{k}_*. \quad (32)$$

Here, $\mathbf{k}_* \triangleq [k(\mathbf{x}_*, \mathbf{x}_1), k(\mathbf{x}_*, \mathbf{x}_2), \dots, k(\mathbf{x}_*, \mathbf{x}_n)]^T$ is the vector of covariances between the test point and the n training points in \mathbf{X} .

Gaussian Process for Regression

Some interpretations of the posterior mean and posterior variance (for one test point case) are as follows:

- The posterior mean can be written as $\bar{m} = \beta^T \mathbf{y}$, which is a linear combination of the observations, often referred as **linear predictor**.
- The posterior variance

$$\bar{v} \triangleq k(\mathbf{x}_*, \mathbf{x}_*) + \sigma^2 - \mathbf{k}_*^T [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_n]^{-1} \mathbf{k}_* \quad (33)$$

is the **difference** between the **variance of the prior** and the **variance explained from the training data \mathcal{D}** .

Gaussian Process for Classification

- In classification problems, we wish to assign a new input \mathbf{x}_* to one of C classes, $\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_C$.
- We focus on the **probabilistic classification based on Gaussian processes**, where test predictions take the form of class probabilities.
- In the sequel, we narrow down our focus to **discriminative approach** for **binary classification**, i.e., the case $C = 2$.

Gaussian Process for Classification

We start with a representative **parametric** binary classification model, namely the **linear logistic regression model**:

$$p(y = +1|\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{x}^T \mathbf{w}), \quad (34)$$

$$p(y = -1|\mathbf{x}, \mathbf{w}) = 1 - p(y = +1|\mathbf{x}, \mathbf{w}), \quad (35)$$

where $\sigma(z) = \frac{1}{1+\exp(-z)}$ is the logistic function. However, $\sigma(z)$ may also take other Sigmoidal type functions.

For **symmetric functions** satisfying $\sigma(z) = 1 - \sigma(-z)$, the likelihood can be written compactly as

$$p(y_i = +1|\mathbf{x}_i, \mathbf{w}) = \sigma(y_i f_i), \quad (36)$$

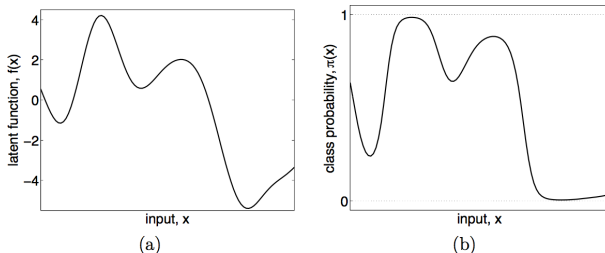
where $f_i \triangleq f(\mathbf{x}_i) = \mathbf{x}_i^T \mathbf{w}$.

Gaussian Process for Classification

Instead of working with the weighting factors \mathbf{w} , Gaussian process classification introduces

$$\pi(\mathbf{x}) \triangleq p(y = +1|f(\mathbf{x})) = \sigma(f(\mathbf{x})), \quad (37)$$

where $f(\mathbf{x})$ is assumed to be a Gaussian process.



Panel (a) shows a sample latent function $f(x)$ drawn from a Gaussian process; Panel (b) shows the transformation of this sample function through $\sigma(f(x))$.

Gaussian Process for Classification

In GP classification, we are particularly interested in the class probability $\bar{\pi}_*$ for novel test point \mathbf{x}_* given the dataset \mathcal{D} , concretely,

$$\bar{\pi}_* \triangleq p(y_* = +1 | \mathbf{x}_*, \mathcal{D}) = \int \sigma(f_*) p(f_* | \mathbf{x}_*, \mathcal{D}) df_*, \quad (38)$$

where

$$p(f_* | \mathbf{x}_*, \mathcal{D}) = p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(f_* | \mathbf{f}, \mathbf{x}_*, \mathbf{X}) p(\mathbf{f} | \mathbf{X}, \mathbf{y}) d\mathbf{f}, \quad (39)$$

and $f_* \triangleq f(\mathbf{x}_*)$, $\mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)]^T$.

- 1 The probability $p(f_* | \mathbf{f}, \mathbf{x}_*, \mathbf{X})$ can be computed using similar technique applied for GP regression.
- 2 The probability $p(\mathbf{f} | \mathbf{X}, \mathbf{y}) = p(\mathbf{y} | \mathbf{f}, \mathbf{X}) p(\mathbf{f} | \mathbf{X}) / p(\mathbf{y} | \mathbf{X})$ is the posterior over the latent variables \mathbf{f} .

Gaussian Process for Classification

Difficulties with GP classification:

- The non-Gaussian distribution $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ makes the integral in Eq.(39) analytically intractable in general.
- The integral in Eq.(38) can be also analytically intractable for some Sigmoid function.

Solution:

- ① Approximate $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ with a Gaussian distribution using e.g.,
 - Laplace approximation [Williams and Baber, 1998], see Appendix
 - Expectation propagation [Minka, 2001]
 - Variational approximation [Opper, 2009]
- ② or Markov chain Monte Carlo (MCMC) [Neal, 1999]

Gaussian Process for Classification

Laplace approximation of $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ is done by Taylor expansion of $\log p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ around the maximum of the posterior, yielding

$$q(\mathbf{f}|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{f}|\hat{\mathbf{f}}, \mathbf{A}^{-1}) \quad (40)$$

where

$$\hat{\mathbf{f}} = \arg \max_{\mathbf{f}} p(\mathbf{f}|\mathbf{X}, \mathbf{y}), \quad (41)$$

$$\mathbf{A} = -\nabla_{\mathbf{f}} \nabla_{\mathbf{f}} \log p(\mathbf{f}|\mathbf{X}, \mathbf{y})|_{\mathbf{f}=\hat{\mathbf{f}}}, \quad (42)$$

are respectively the **maximum of the posterior** and the **Hessian of the negative log posterior at the maximum**. Details see Appendix.

The Hessian \mathbf{A} can be further expressed as $\mathbf{A} = \mathbf{W} + \mathbf{K}^{-1}$, where \mathbf{K} is short for $\mathbf{K}(\mathbf{X}, \mathbf{X})$ and \mathbf{W} is a diagonal matrix with

$$\mathbf{W} \triangleq -\nabla_{\mathbf{f}} \nabla_{\mathbf{f}} \log p(\mathbf{y}|\mathbf{f})|_{\mathbf{f}=\hat{\mathbf{f}}}. \quad (43)$$

Gaussian Process for Classification

Now, we come to the inference part.

The exact posterior mean for a novel test point \mathbf{x}_* is

$$\mathbb{E}_p[f_*|\mathbf{x}_*, \mathcal{D}] = \int f_* p(f_*|\mathbf{x}_*, \mathcal{D}) df_* \quad (44)$$

$$= \int \mathbb{E}[f_*|\mathbf{f}, \mathbf{X}, \mathbf{x}_*] p(\mathbf{f}|\mathbf{X}, \mathbf{y}) d\mathbf{f} \quad (45)$$

$$= \int \mathbf{k}(\mathbf{x}_*)^T \mathbf{K}^{-1} \mathbf{f} p(\mathbf{f}|\mathbf{X}, \mathbf{y}) d\mathbf{f} \quad (46)$$

$$= \mathbf{k}(\mathbf{x}_*)^T \mathbf{K}^{-1} \mathbb{E}[\mathbf{f}|\mathbf{X}, \mathbf{y}], \quad (47)$$

where line 1 to 2 is due to $p(f_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X}) p(\mathbf{f}|\mathbf{X}, \mathbf{y}) d\mathbf{f}$.

We obtain the Laplace approximated posterior mean as

$$\mathbb{E}_q[f_*|\mathbf{x}_*, \mathcal{D}] = \mathbf{k}(\mathbf{x}_*)^T \mathbf{K}^{-1} \hat{\mathbf{f}}. \quad (48)$$

Gaussian Process for Classification

The **exact posterior variance** for a novel test point \mathbf{x}_* is

$$\text{Var}_p[f_*|\mathbf{x}_*, \mathcal{D}] = \int (f_* - \mathbb{E}_p[f_*|\mathbf{x}_*, \mathcal{D}])^2 p(f_*|\mathbf{x}_*, \mathcal{D}) df_* \quad (49)$$

$$\begin{aligned} &= \mathbb{E}_{p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X})} \left[(f_* - \mathbb{E}[f_*|\mathbf{f}, \mathbf{X}, \mathbf{x}_*])^2 \right] \\ &+ \mathbb{E}_{p(\mathbf{f}|\mathbf{X}, \mathbf{y})} \left[(\mathbb{E}[f_*|\mathbf{f}, \mathbf{X}, \mathbf{x}_*] - \mathbb{E}[f_*|\mathbf{y}, \mathbf{X}, \mathbf{x}_*])^2 \right] \end{aligned} \quad (50)$$

We obtain the **Laplace approximated posterior variance** as

$$\begin{aligned} \text{Var}_q[f_*|\mathbf{x}_*, \mathcal{D}] &= \mathbb{E}_{p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X})} \left[(f_* - \mathbb{E}[f_*|\mathbf{f}, \mathbf{X}, \mathbf{x}_*])^2 \right] \\ &+ \mathbb{E}_{q(\mathbf{f}|\mathbf{X}, \mathbf{y})} \left[(\mathbb{E}[f_*|\mathbf{f}, \mathbf{X}, \mathbf{x}_*] - \mathbb{E}[f_*|\mathbf{y}, \mathbf{X}, \mathbf{x}_*])^2 \right] \end{aligned} \quad (51)$$

$$\begin{aligned} &= k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}(\mathbf{x}_*)^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}_*) \\ &+ \mathbf{k}(\mathbf{x}_*)^T \mathbf{K}^{-1} (\mathbf{K}^{-1} + \mathbf{W})^{-1} \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}_*) \end{aligned} \quad (52)$$

$$= k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}(\mathbf{x}_*)^T (\mathbf{K} + \mathbf{W}^{-1})^{-1} \mathbf{k}(\mathbf{x}_*). \quad (53)$$

Detailed derivations of the above steps can be found in the Appendix.



Gaussian Process for Classification

Given an approximation of the posterior distribution $p(f_*|\mathbf{y}, \mathbf{X}, \mathbf{x}_*)$, i.e., $q(f_*|\mathbf{y}, \mathbf{X}, \mathbf{x}_*)$ which is Gaussian distributed with its mean and variance derived in the previous two slides, we have the probability of having class 1 approximately equal to:

$$\begin{aligned}\bar{\pi}_* &\triangleq p(y_* = +1|\mathbf{x}_*, \mathcal{D}) = \int \sigma(f_*) p(f_*|\mathbf{x}_*, \mathcal{D}) df_* \\ &\approx \int \sigma(f_*) q(f_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) df_*. \end{aligned} \tag{54}$$

When $\sigma(z)$ takes the logistic function, the above one-dimensional integral has to be evaluated numerically.

- ① T. P. Minka, A Family of Algorithms for Approximate Bayesian Inference. PhD thesis, Massachusetts Institute of Technology, 2001.
- ② C. K. I. Williams and D. Barber, Bayesian Classification with Gaussian Processes. IEEE Transactions on Pattern Analysis and Machine Intelligence, 20(12): 1342–1351, 1998.
- ③ R. M. Neal, Regression and Classification using Gaussian Process Priors. In Bayesian Statistics 6, pages 475–501. Oxford University Press. 1999.
- ④ M. Opper and C. Archambeau, The variational Gaussian approximation revisited. Neural computation, 21(3): 786–792, 2009.

Hyper-parameter Optimization: ML Method

The predictive performance of GP regression and classification depends on the goodness of a set of hyper-parameters.

The following methods can be used for tuning the hyper-parameters:

- ① deterministic methods including:
 - maximum likelihood (ML)
 - cross-validation
 - alignment score (mainly for GP classification)
- ② stochastic methods including:
 - hybrid Monte-Carlo (HMC) [Rasmussen 96, Neal 97].
 - slice sampling

Hyper-parameter Optimization: ML Method

The dominant method for tuning the hyper-parameters is via **maximizing the marginal likelihood function** $p(\mathbf{y}; \boldsymbol{\theta})$, where

$$p(\mathbf{y}; \boldsymbol{\theta}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f}. \quad (55)$$

Recall that

- ① for the **GP regression**, we have $y_i = f(\mathbf{x}_i) + n_i$, where n_i is assumed to be Gaussian noise;
 - ② and for the **GP classification**, we have $y_i = \sigma(f(\mathbf{x}_i))$, and $\sigma(z)$ is a non-linear mapping,
- and $f(\mathbf{x})$ is assumed to be a Gaussian process in both cases.

Hyper-parameter Optimization: ML Method

For [GP regression](#), due to the Gaussian assumption on the noise, the log-marginal likelihood function can be obtained in closed form as

$$\log p(\mathbf{y}; \boldsymbol{\theta}) = -\frac{1}{2} \left\{ \log \det(\mathbf{C}(\boldsymbol{\theta})) + \mathbf{y}^T \mathbf{C}^{-1}(\boldsymbol{\theta}) \mathbf{y} + n \log(2\pi) \right\} \quad (56)$$

where $\mathbf{C}(\boldsymbol{\theta}) \triangleq \mathbf{K}(\mathbf{X}, \mathbf{X}; \boldsymbol{\theta}_h) + \sigma^2 \mathbf{I}_n$, and $\boldsymbol{\theta} = [\boldsymbol{\theta}_h^T, \sigma^2]^T$.

The hyper-parameter is tuned equivalently by [minimizing the negative log-likelihood function](#). In the GP society, this is often solved using the [gradient based algorithms](#), such as LFGS-Newton, which requires the gradient:

$$\frac{\partial -\log p(\mathbf{y}; \boldsymbol{\theta})}{\partial \theta_i} = \frac{1}{2} \text{tr} \left(\mathbf{C}^{-1}(\boldsymbol{\theta}) \frac{\partial \mathbf{C}(\boldsymbol{\theta})}{\partial \theta_i} \right) - \frac{1}{2} \mathbf{y}^T \mathbf{C}^{-1}(\boldsymbol{\theta}) \frac{\partial \mathbf{C}(\boldsymbol{\theta})}{\partial \theta_i} \mathbf{C}^{-1}(\boldsymbol{\theta}) \mathbf{y}. \quad (57)$$

Hyper-parameter Optimization: ML Method

For **GP classification**, due to the nonlinear mapping $\sigma(y_i f_i)$, the marginal likelihood can not be obtained in closed form in general.

Since

$$p(\mathbf{y}; \boldsymbol{\theta}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f} = \int \exp(\Psi(\mathbf{f})) d\mathbf{f}, \quad (58)$$

where $\Psi(\mathbf{f}) \triangleq \log p(\mathbf{y}|\mathbf{f}) + \log p(\mathbf{f})$. Perform **Taylor expansion of $\Psi(\mathbf{f})$ at the maximum of the posterior $\hat{\mathbf{f}}$** , an **approximated** log-likelihood function is then obtained as

$$\log q(\mathbf{y}; \boldsymbol{\theta}) = -\frac{1}{2}\hat{\mathbf{f}}^T \mathbf{K}(\mathbf{X}, \mathbf{X}; \boldsymbol{\theta}_h)^{-1} \hat{\mathbf{f}} + \log p(\mathbf{y}|\hat{\mathbf{f}}) - \frac{1}{2} \log \det(\mathbf{B}), \quad (59)$$

where $\det(\mathbf{B}) = \det(\mathbf{I}_n + \mathbf{W}^{\frac{1}{2}} \mathbf{K}(\mathbf{X}, \mathbf{X}; \boldsymbol{\theta}_h) \mathbf{W}^{\frac{1}{2}})$.

Gradient can be obtained similarly as the regression case.

Hyper-parameter Optimization: ML Method

We know that the gradient based algorithms are iterative and at each iteration, in general,

- The computational complexity for evaluating the gradient of the negative marginal likelihood function requires the matrix inverse $\mathbf{C}^{-1}(\boldsymbol{\theta})$ or $\mathbf{K}^{-1}(\boldsymbol{\theta})$, scales as $\mathcal{O}(n^3)$.
- Once $\mathbf{C}^{-1}(\boldsymbol{\theta})$ or $\mathbf{K}^{-1}(\boldsymbol{\theta})$ is obtained, the computation of the derivatives scales as $\mathcal{O}(n^2)$ per hyper-parameter.

When the dimension of \mathbf{y} is high and the number of hyper-parameters to be optimized is also high, better strategy for parameter optimization is sought after.

Hyper-parameter Optimization: Monte Carlo Method

Another way of dealing with the hyper-parameters is to use Monte Carlo method. Principally, the posterior distribution is calculated by

$$p(\mathbf{y}_*|\mathcal{D}) = \int_{\boldsymbol{\theta} \in \Theta} p(\mathbf{y}_*|\mathcal{D}, \boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta}. \quad (60)$$

The Monte Carlo method uses [sampling methods](#) to get an approximated posterior distribution via

$$p(\mathbf{y}_*|\mathcal{D}) \approx \frac{1}{N_{MC}} \sum_{i=1}^{N_{MC}} p(\mathbf{y}_*|\mathcal{D}, \boldsymbol{\theta}_i), \quad (61)$$

where $\boldsymbol{\theta}_i$, $i = 1, 2, \dots, N_{MC}$ are samples drawn from the posterior distribution $p(\boldsymbol{\theta}|\mathcal{D})$.

Hyper-parameter Optimization: Monte Carlo Method

Some remarks are:

- The Monte Carlo approximation of the posterior distribution is a Gaussian mixture;
- The approximation accuracy will in general increase as we take more samples from $p(\theta|\mathcal{D})$;
- The method that we use to sample from $p(\theta|\mathcal{D})$ impacts the approximation accuracy.
- Low complex and efficient sampling method for high-dimensional space is still under development. Some classic work can be found in [Duane 87, Rasmussen 96, Neal 97].

Summary

- Bayesian nonparametric model enables good model flexibility;
- Gaussian noise assumption and factorizing property of the likelihood function largely simplify the derivations;
- Computational complexity for training the hyper-parameters is quite high.

- ① C. E. Rasmussen, Evaluation of Gaussian Processes and Other Methods for Nonlinear Regression. PhD thesis, Dept. of Computer Science, University of Toronto. 1996.
- ② R. M. Neal, Monte Carlo Implementation of Gaussian Process Models for Bayesian Regression and Classification. Technical Report 9702, Department of Statistics, University of Toronto. 1997.

GPR Example: Mauna Loa Atmospheric Carbon Dioxide

Review: Given the training dataset $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$ and a novel test data point \mathbf{x}_* , the GPR method does the following to predict y_* :

- 1 **Select a suitable kernel**, stationary or non-stationary, according to the prior knowledge about the data.
- 2 **Maximize the log-marginal likelihood function** in Eq.(56) to get
 - an optimal set of parameters θ ,
 - the inverse of the covariance matrix $\mathbf{C}(\theta) \triangleq \mathbf{K}(\mathbf{X}, \mathbf{X}; \theta_h) + \sigma^2 \mathbf{I}_n$.
- 3 **Compute the posterior/predictive mean and variance** according to Eq.(31) and Eq.(32).

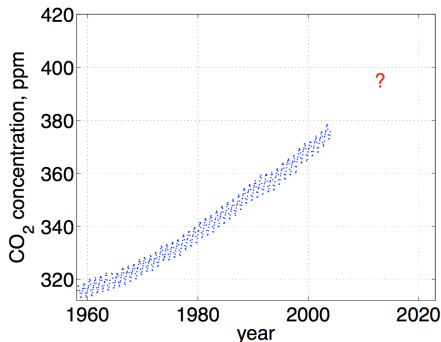
GPR Example: Data Description

The data [Keeling and Whorf, 2004] consists of monthly average atmospheric CO₂ concentrations (unit: per million by volume (ppmv)) derived from in-situ air samples collected at the Mauna Loa Observatory, Hawaii, between 1958 and 2003.

Patterns seen from the data:

- long term rising trend
- seasonal variation
- smaller irregularities
- noise

Task: predict beyond 2003.



GPR Example: Kernel Selection

Rasmussen and Williams suggest a composite kernel:

$$k(x, x'; \theta_h) = k_1(x, x') + k_2(x, x') + k_3(x, x') + k_4(x, x') \quad (62)$$

where

- long-term (smooth) rising trend addressed by the SE kernel

$$k_1(x, x') = \theta_1^2 \exp \left[\frac{-(x-x')^2}{\theta_2^2} \right]$$

- seasonal trend addressed by the LP kernel

$$k_2(x, x') = \theta_3^2 \exp \left[\frac{-2 \sin^2(\pi(x-x'))}{\theta_5^2} \right] \times \exp \left[\frac{-(x-x')^2}{2\theta_4^2} \right]$$

- smaller irregularities/anomaly addressed by the RQ kernel

$$k_3(x, x') = \theta_6^2 \left(1 + \frac{(x-x')^2}{2\theta_8\theta_7^2} \right)^{-\theta_8}$$

- and noise addressed by another SE kernel

$$k_4(x, x') = \theta_9^2 \exp \left[\frac{-(x-x')^2}{\theta_{10}^2} \right]$$

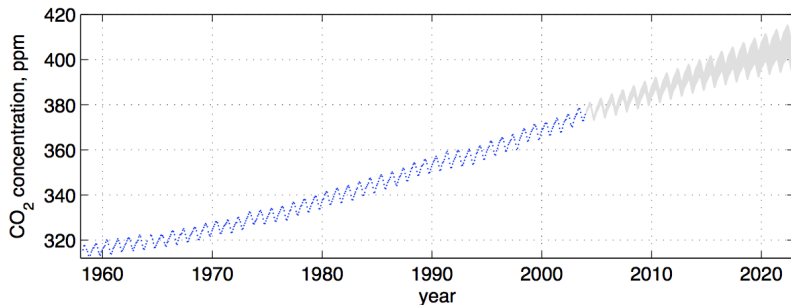
GPR Example: Hyper-parameter Optimization

Hyperparameter optimization:

- Empirical mean of the data (341 ppm) is first subtracted from the data to have a **zero-mean Gaussian process**.
- The parameters $\theta = [\theta_1, \theta_2, \dots, \theta_{10}, \sigma^2]^T$ are tuned by maximizing the log-marginal likelihood using **conjugate gradient method**.
- **Random restarts** are tried and pick the run with the best log-marginal likelihood, which is $\log p(\mathbf{y}|\theta) = -108.5$.

GPR Example: Inference

Prediction performance using the above composite kernel:



The [545 observations](#) of monthly averages of the atmospheric concentration of CO₂ made between 1958 and the end of 2003, together with [95% predictive confidence region](#) for a GP regression model, 20 years into the future.

GPC Example: Binary Handwritten Digit Classification

Review: Given the training dataset $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$ and a novel test data point \mathbf{x}_* , the **binary GPC** method does the following to obtain $p(y_* = 1|\mathcal{D})$ and the final label:

- ① Select a suitable kernel, stationary or non-stationary, according to the prior knowledge about the data.
- ② Select an initial guess of the hyper-parameters $\boldsymbol{\theta}_h^{(0)}$. Do the following steps until some convergence condition is met:
 - Update the maximum of the posterior $\hat{\mathbf{f}}$ at the current set of hyper-parameters $\boldsymbol{\theta}_h^{(i)}$ in the i -th iteration according to Eq.(41).
 - Maximize the approximated log-marginal likelihood function to update $\boldsymbol{\theta}_h^{(i+1)}$ according to (59).
- ③ Compute the Laplace approximated posterior mean and variance according to Eq.(48) and Eq.(53).
- ④ If $p(y_* = 1|\mathcal{D}) > 0.5$, we go for label 1; otherwise label 0 for \mathbf{x}_* . Or alternatively, check the sign of the posterior mean.

GPC Example: Data Description

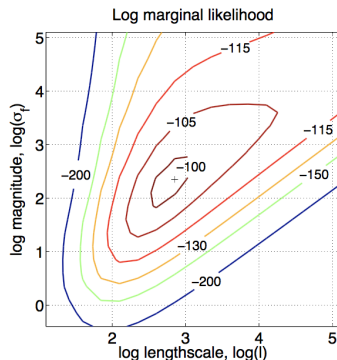
We consider the discrimination of images of the digit 3 from images of the digit 5 as an example of binary classification.

Data description:

- US Postal Service (USPS) database of handwritten digits as segmented 16×16 greyscale images.
- The images are normalized so that the intensity of the pixels lies in $[-1, 1]$.
- The 3s vs. 5s data has 767 training cases, split into 406/361, while the test set has 773 cases split into 418/355.

GPC Example: Hyper-parameter Optimization

- The SE kernel is selected, $k(x, x'; \theta_h) = \theta_f^2 \exp(-(x - x')^2 / l^2)$.
- An optimal set of hyper-parameters $\theta_h = [\theta_f, l]^T$ is obtained by maximizing the Laplace approximated log-marginal likelihood function using the conjugate gradient method.

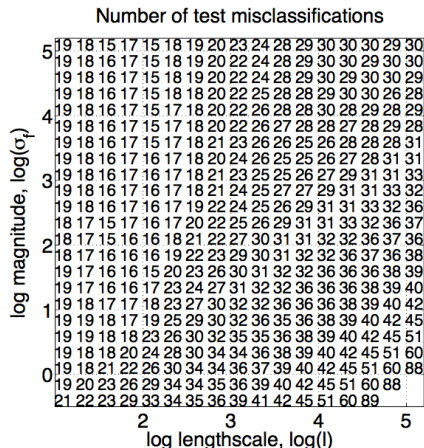


At the optimum:

- $\log(l) = 2.85$ and $\log(\sigma_f) = 2.35$
- $\log q(\mathbf{y}; \theta) = -99$

GPC Example: Inference

The predictive performance in terms of the test errors:



the number of test errors (out of 773) when predicting using the sign of the posterior mean.

① MATLAB:

- built-in *Gaussian Process Regression Models* in Statistic and Machine Learning Toolbox, 2017b.
<https://cn.mathworks.com/help/stats/gaussian-process-regression-models.html>
- GPML Matlab Code version 4.0, written by C. Rasmussen,
<http://www.gaussianprocess.org/gpml/code/matlab/doc/>

② PYTHON:

- Gaussian Processes package,
http://scikit-learn.org/stable/modules/gaussian_process.html

③ JULIA:

- Gaussian Processes package,
<https://github.com/STOR-i/GaussianProcesses.jl>

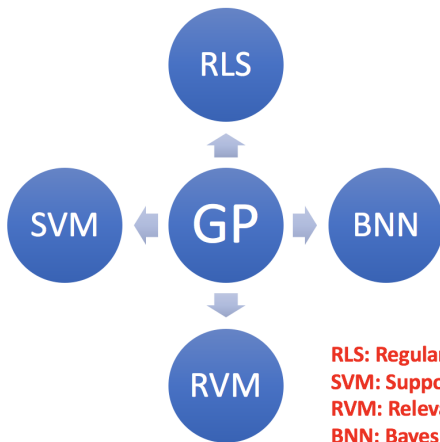
- 1 C. D. Keeling and T. P. Whorf, Atmospheric CO₂ Records from Sites in the SIO Air Sampling Network. In Trends: A Compendium of Data on Global Change. Carbon Dioxide Information Analysis Center, Oak Ridge National Laboratory, Oak Ridge, Tenn., U.S.A., 2004.

Outline

- 1 Background
- 2 Kernel Methods
 - Dual Representation of RLS
 - Constructing Valid Kernels
- 3 Gaussian Processes (GP) for Machine Learning
 - GP for Regression
 - GP for Classification
 - Hyper-parameter Optimization
 - Two GP Examples
- 4 Relationships with Other Models
- 5 Summary and Selected Recent Advances
- 6 Appendix

Overview

The relationships between the Gaussian Process (GP) models and other models are summarized in the figure below:



RLS: Regularized Least-Squares
SVM: Support Vector Machine
RVM: Relevance Vector Machine
BNN: Bayesian Neural Network

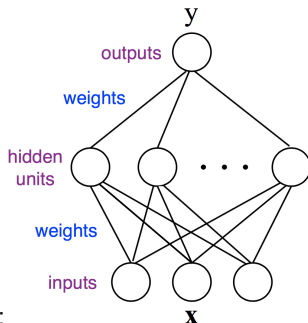
Relationship with Bayesian Neural Networks

Consider a neural network with **one hidden layer of H units** that takes an input $\mathbf{x} \in \mathbb{R}^p$. The mapping function can be written as:

$$f(\mathbf{x}) = b + \sum_{j=1}^H v_j h(\mathbf{x}; \mathbf{u}_j), \quad (63)$$

where

- v_j 's are the hidden-to-output weights;
- $h(\mathbf{x}; \mathbf{u})$ is the hidden unit transfer/activation function (**assume to be bounded**) which depends on the input-to-hidden weights \mathbf{u} ;
- b is a bias term.



Relationship with Bayesian Neural Networks

Assuming that

- b and v 's have independent zero-mean distributions of variance σ_b^2 and σ_v^2 , respectively;
- the weights \mathbf{u}_j for each hidden unit be i.i.d..

Stacking all weights to obtain \mathbf{w} , we have

$$\mathbb{E}_{\mathbf{w}} [f(\mathbf{x})] = 0 \quad (64)$$

$$\mathbb{E}_{\mathbf{w}} [f(\mathbf{x})f(\mathbf{x}')] = \sigma_b^2 + \sigma_v^2 \sum_{j=1}^H \mathbb{E}_{\mathbf{u}} [h(\mathbf{x}; \mathbf{u}_j)h(\mathbf{x}'; \mathbf{u}_j)]. \quad (65)$$

Remark: The sum is over H i.i.d. RVs. As the transfer function is bounded, all moments of the distribution will be bounded and applying the CLT shows that $f(\mathbf{x})$ converges to a Gaussian process as $H \rightarrow \infty$.

Relationship with Bayesian Neural Networks

For the following Sigmoid transfer functions, the corresponding kernel function is well known:

- $h(\mathbf{x}; \mathbf{u}) = \Phi\left(u_0 + \sum_{i=1}^p u_i x_i\right)$, where $\Phi(z) = \frac{2}{\sqrt{\pi}} \int_0^z \exp(-t^2) dt$ is the error function.
- Assume that \mathbf{u} is drawn from a zero-mean Gaussian distribution with covariance matrix Σ . Furthermore, let $\tilde{\mathbf{x}} = [1, x_1, x_2, \dots, x_p]^T$.
- The neural network kernel is derived then as

$$k_{NN}(\mathbf{x}, \mathbf{x}') \triangleq \mathbb{E}_{\mathbf{u}} \left[h(\mathbf{x}, \mathbf{u}) h(\mathbf{x}', \mathbf{u}) \right] \quad (66)$$

$$= \int \Phi(\mathbf{u}^T \tilde{\mathbf{x}}) \Phi(\mathbf{u}^T \tilde{\mathbf{x}}') \mathcal{N}(\mathbf{u}; \mathbf{0}, \Sigma) d\mathbf{u} \quad (67)$$

$$= \frac{2}{\pi} \sin^{-1} \left(\frac{2 \tilde{\mathbf{x}}^T \Sigma \tilde{\mathbf{x}}'}{\sqrt{(1 + 2(\tilde{\mathbf{x}})^T \Sigma \tilde{\mathbf{x}})(1 + 2(\tilde{\mathbf{x}}')^T \Sigma \tilde{\mathbf{x}}')}} \right) \quad (68)$$

Relationship with Support Vector Machine

This section aims to compare the Gaussian process classification with the support vector machine (SVM).

For the non-separable case, the soft margin SVM computes

$$\mathbf{w}_{SVM} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(1 - y_i f_i, 0) \quad (69)$$

where

- $f_i \triangleq \mathbf{w}^T \phi(\mathbf{x}_i)$;
- $\max(z, 0)$ return z , if $z > 0$ or zero otherwise;
- $C > 0$ is a regularization term.

Relationship with Support Vector Machine

This convex optimization problem can be solved using QP methods, e.g., the sequential minimal optimization (SMO) algorithm [Platt, 99] and yields a solution of the form $\sum_i^n \alpha_i \phi(\mathbf{x}_i)$.

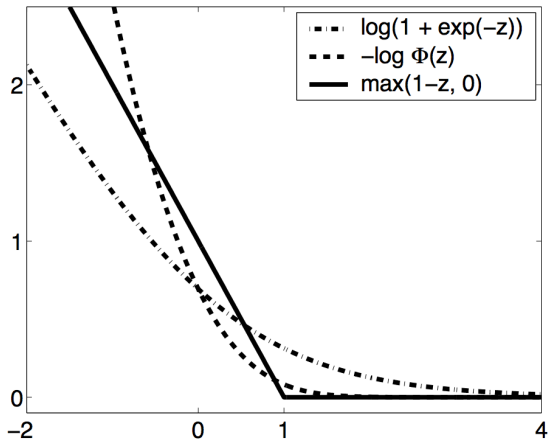
Then, we have $\|\mathbf{w}\|^2 = \boldsymbol{\alpha}^T K \boldsymbol{\alpha} = \mathbf{f}^T K^{-1} \mathbf{f}$, where we introduced $K \boldsymbol{\alpha} = \mathbf{f}$. Thus the soft margin objective function can be re-written as

$$\arg \min_{\mathbf{f}} \frac{1}{2} \mathbf{f}^T K^{-1} \mathbf{f} + C \sum_{i=1}^n (1 - y_i f_i)_+ \quad (70)$$

For the binary GP classifier, to obtain the MAP value $\hat{\mathbf{f}}$ of $p(\mathbf{f}|\mathbf{y})$ we solve the following minimization problem:

$$\arg \min_{\mathbf{f}} \frac{1}{2} \mathbf{f}^T K^{-1} \mathbf{f} - \sum_{i=1}^n \log p(y_i | f_i). \quad (71)$$

Relationship with Support Vector Machine



A comparison between different error functions.

- ① J. C. Platt, Fast Training of Support Vector Machines Using Sequential Minimal Optimization. In Schölkopf, Burges, and Smola, editors, *Advances in Kernel Methods*, pages 185–208. MIT Press, 1999.

Outline

- 1 Background
- 2 Kernel Methods
 - Dual Representation of RLS
 - Constructing Valid Kernels
- 3 Gaussian Processes (GP) for Machine Learning
 - GP for Regression
 - GP for Classification
 - Hyper-parameter Optimization
 - Two GP Examples
- 4 Relationships with Other Models
- 5 **Summary and Selected Recent Advances**
- 6 Appendix

Summary

- We introduced Bayesian nonparametric models based on Gaussian processes for nonlinear regression and classification.
- The performance depends on the selection of a kernel function.
- Excellent performance is gained at the cost of high computational complexity, $\mathcal{O}(n^3)$, which prohibits the classic algorithms to be used for big data applications.
- Assumptions on the Gaussian likelihood and factoring property of the likelihood function are unrealistic for many applications.

Selected Recent Advances

- Low-complex Gaussian processes (for big data)
 - ① M. Todescato *et.al.*, "Efficient Spatio-Temporal Gaussian Regression via Kalman Filtering", <https://arxiv.org/abs/1705.01485>, 2017.
 - ② J. Hensman and N. Lawrence, "Gaussian Processes for Big Data", <http://www.auai.org/uai2013>, 2013.
 - ③ M. Deisenroth, "Distributed Gaussian Processes", <https://arxiv.org/pdf/1502.02843>, 2015.
 - ④ S. Ambikasaran *et.al.*, Fast Direct Methods for Gaussian Processes, <https://arxiv.org/pdf/1403.6015.pdf>, 2015.
- More sophisticated GP model
 - ① A. Damianou and N. Lawrence, "Deep Gaussian Process", Proceedings of Machine Learning Research, 2013.
 - ② M van der Wilk, "Convolutional Gaussian Processes", <https://arxiv.org/abs/1709.01894>, 2017.

Some of the above listed papers will be covered in SRIBD reading group.

Outline

- 1 Background
- 2 Kernel Methods
 - Dual Representation of RLS
 - Constructing Valid Kernels
- 3 Gaussian Processes (GP) for Machine Learning
 - GP for Regression
 - GP for Classification
 - Hyper-parameter Optimization
 - Two GP Examples
- 4 Relationships with Other Models
- 5 Summary and Selected Recent Advances
- 6 **Appendix**

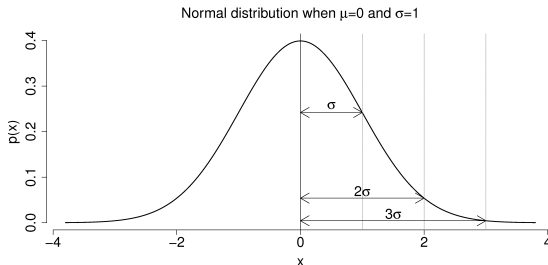
Gaussian Distribution: Univariate Case

The probability density function (pdf) is

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[\frac{-(x - \mu)^2}{2\sigma^2} \right]. \quad (72)$$

The mean and variance are

$$\mathbb{E}(x) = \mu, \quad \text{var}(x) = \mathbb{E}[(x - \mathbb{E}(x))^2] = \sigma^2. \quad (73)$$



Gaussian Distribution: Multivariate Case

The probability density function (pdf) is

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{p/2} \det(\Sigma)^{1/2}} \exp \left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right] \quad (74)$$

The mean and covariance matrix are

$$\mathbb{E}(\mathbf{x}) = \boldsymbol{\mu}, \quad \text{Cov}(\mathbf{x}) = \mathbb{E} \left[(\mathbf{x} - \mathbb{E}(\mathbf{x}))(\mathbf{x} - \mathbb{E}(\mathbf{x}))^T \right] = \Sigma. \quad (75)$$

Conditional Gaussian Distribution

If the two random variables $\mathbf{x} \in \mathbb{R}^{d_x}$ and $\mathbf{y} \in \mathbb{R}^{d_y}$ are jointly Gaussian with the following joint distribution:

$$p(\mathbf{x}, \mathbf{y}) \sim \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}\right), \quad (76)$$

then it is easy to derive the following **conditional probabilities**:

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_{x|y}, \Sigma_{x|y}), \quad p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_{y|x}, \Sigma_{y|x}), \quad (77)$$

where

$$\boldsymbol{\mu}_{x|y} = \boldsymbol{\mu}_x + \Sigma_{xy}\Sigma_{yy}^{-1}(\mathbf{y} - \boldsymbol{\mu}_y), \quad \boldsymbol{\mu}_{y|x} = \boldsymbol{\mu}_y + \Sigma_{yx}\Sigma_{xx}^{-1}(\mathbf{x} - \boldsymbol{\mu}_x), \quad (78)$$

and

$$\Sigma_{x|y} = \Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx}, \quad \Sigma_{y|x} = \Sigma_{yy} - \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy}. \quad (79)$$

Marginal Gaussian Distribution

Following the previous slide where the joint Gaussian distribution $p(\mathbf{x}, \mathbf{y})$ was defined.

The **marginal distributions** out of it are still Gaussian, i.e.,

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_x, \Sigma_{xx}), \quad (80)$$

and

$$p(\mathbf{y}) = \int p(\mathbf{x}, \mathbf{y}) d\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_y, \Sigma_{yy}). \quad (81)$$

An Important Gaussian Identities

The product of two Gaussians gives another (un-normalized) Gaussian

$$\mathcal{N}(\mathbf{x}|\mathbf{a}, \mathbf{A}) \cdot \mathcal{N}(\mathbf{x}|\mathbf{b}, \mathbf{B}) = \mathbf{Z}^{-1} \mathcal{N}(\mathbf{x}|\mathbf{c}, \mathbf{C}), \quad (82)$$

where $\mathbf{c} = \mathbf{C}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b})$ and $\mathbf{C} = (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1}$ and the normalizing factor is

$$\mathbf{Z}^{-1} = (2\pi)^{-p/2} \det(\mathbf{A} + \mathbf{B})^{-1/2} \exp \left[-\frac{1}{2}(\mathbf{a} - \mathbf{b})^T (\mathbf{A} + \mathbf{B})^{-1} (\mathbf{a} - \mathbf{b}) \right]. \quad (83)$$

Eigendecomposition of Symmetric, PD Matrices

For **real, symmetric, positive definite (PD)** covariance matrix Σ , which is of size $n \times n$, its eigenvector equation can be expressed as

$$\Sigma \mathbf{u}_i = \lambda_i \mathbf{u}_i, \quad i = 1, 2, \dots, n$$

where

- the eigenvalues are real;
- and the eigenvectors can be chosen to form an orthonormal set, i.e., $\mathbf{u}_i^T \mathbf{u}_j = 0$, if $i \neq j$; otherwise $\mathbf{u}_i^T \mathbf{u}_j = 1$, if $i = j$.

As a consequence, the covariance matrix and its inverse matrix can be re-expressed as

$$\Sigma = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^T, \quad \Sigma^{-1} = \sum_{i=1}^n \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T. \quad (84)$$

Cholesky Decomposition

A **symmetric, positive definite matrix** \mathbf{A} can be decomposed into a product of a lower triangular matrix \mathbf{L} and its transpose (upper triangular)

$$\mathbf{L}\mathbf{L}^T = \mathbf{A}, \quad (85)$$

where \mathbf{L} is called the Cholesky factor.

To solve the linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$, we follow the steps below:

- first solve $\mathbf{L}\mathbf{y} = \mathbf{b}$ by forward substitution;
- then $\mathbf{L}^T\mathbf{x} = \mathbf{y}$ by back substitution;
- write the solution as $\mathbf{x} = \mathbf{L}^T \setminus (\mathbf{L} \setminus \mathbf{b})$, where the backslash notation $\mathbf{A} \setminus \mathbf{b}$ is the vector \mathbf{x} which solves $\mathbf{A}\mathbf{x} = \mathbf{b}$.

An useful result: the determinant of a symmetric, positive definite (PD) matrix can be calculated by

$$\det(\mathbf{A}) = \prod_{i=1}^n L_{ii}^2, \quad \text{or} \quad \log \det(\mathbf{A}) = 2 \sum_{i=1}^n \log L_{ii}. \quad (86)$$

The limit of a rational quadratic is squared exponential due to

$$\lim_{\alpha \rightarrow \infty} \left(1 + \frac{x^2}{2\alpha}\right)^{-\alpha} = \exp\left(-\frac{x^2}{2}\right), \quad (87)$$

which explains that the limit of the RQ kernel in Eq.(14) for $\alpha \rightarrow \infty$ becomes the SE kernel in Eq.(13).

Let Ω be the set of all possible outcomes of an experiment and \mathcal{F} be a σ -field of subsets of Ω which contains all the events of interest, then μ is a countably additive measure if it is real and non-negative and for all mutually disjoint sets $A_1, A_2, \dots \in \mathcal{F}$, we have

$$\mu \left(\bigcup_{i=1}^{\infty} A_i \right) = \sum_{i=1}^{\infty} \mu(A_i). \quad (88)$$

- ① If $\mu(\Omega) < \infty$, then μ is called a finite measure;
- ② if $\mu(\Omega) = 1$, then μ is called a probability measure;
- ③ The Lebesgue measure defines a uniform measure over subsets of Euclidean space.

Integration

We give the integration of a function $f(\mathbf{x}) : \mathbb{R}^p \rightarrow \mathbb{R}$ with respect to a measure μ by

$$\int f(\mathbf{x}) d\mu(\mathbf{x}), \quad (89)$$

where it is assumed that $f(\cdot)$ is measurable, i.e., for any Borel-measurable set $A \in \mathbb{R}$, $f^{-1}(A) \in \mathcal{B}^d$.

We classify the following two cases:

- 1 When μ is a Lebesgue measure, $\int f(\mathbf{x}) d\mu(\mathbf{x}) = \int f(\mathbf{x}) d\mathbf{x}$;
- 2 When μ is a probability measure, $\int f(\mathbf{x}) d\mu(\mathbf{x}) = \int f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$, where the non-negative function $p(\mathbf{x})$ is called the density of the measure if for all $A \in \mathcal{B}^d$ we have $\mu(A) = \int_A p(\mathbf{x}) d\mathbf{x}$. Note that not all probability measures have densities.

Kernel and PSD Kernel

A general name for a function $k(\cdot, \cdot)$ of two arguments mapping a pair of inputs $\mathbf{x} \in \mathcal{X}, \mathbf{x}' \in \mathcal{X}$ into \mathbb{R} is a *kernel*. This term arises in the theory of integral operators, where the operator T_k is defined as

$$(T_k f)(\mathbf{x}) = \int_{\mathcal{X}} k(\mathbf{x}, \mathbf{x}') f(\mathbf{x}') d\mu(\mathbf{x}'), \quad (90)$$

where μ denotes a measure.

A *kernel* is said to be positive semidefinite if

$$\int k(\mathbf{x}, \mathbf{x}') f(\mathbf{x}) f(\mathbf{x}') d\mu(\mathbf{x}) d\mu(\mathbf{x}') \geq 0, \quad (91)$$

for all $f \in L_2(\mathcal{X}, \mu)$. Equivalently a kernel function which gives rise to PSD Gram matrices for any choice of $n \in \mathbb{N}$ and \mathcal{D} is positive semi-definite.

Expansion of SE Kernel

For $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^p$, the polynomial kernel $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^d$, where d is the polynomial order, a positive integer. The kernel function can be expressed as:

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') \quad (92)$$

where

$$\phi(\mathbf{x}) = \left\{ \sqrt{\frac{d!}{n_1! n_2! \cdots n_p!}} x_1^{n_1} x_2^{n_2} \cdots x_p^{n_p} \right\}, \quad n_i \geq 0, n_i \in \mathbb{N}, \sum_{i=1}^p n_i = d \quad (93)$$

For example, for $p = d = 2$, $(\mathbf{x}^T \mathbf{x}')^2 = \phi(\mathbf{x})^T \phi(\mathbf{x}')$, where according to the above formula,

$$\phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)^T. \quad (94)$$

Expansion of SE Kernel (Con'd)

For the SE kernel,

$$\begin{aligned}k_{SE}(\mathbf{x}, \mathbf{x}') &= \sigma^2 \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2l^2}\right) \\&= \sigma^2 \exp\left(\frac{-\|\mathbf{x}\|^2}{2}\right) \exp\left(\frac{-\|\mathbf{x}'\|^2}{2}\right) \exp\left(\mathbf{x}^T \mathbf{x}'\right) \\&= \sigma^2 \exp\left(\frac{-\|\mathbf{x}\|^2}{2}\right) \exp\left(\frac{-\|\mathbf{x}'\|^2}{2}\right) \sum_{j=0}^{\infty} \frac{(\mathbf{x}^T \mathbf{x}')^j}{j!}\end{aligned}\quad (95)$$

According to the results shown above for the polynomial kernel, we have

$$\phi(\mathbf{x}) = \sigma \exp\left(\frac{-\|\mathbf{x}\|^2}{2}\right) \left\{ \sqrt{\frac{j!}{n_1! n_2! \dots n_p!}} x_1^{n_1} x_2^{n_2} \dots x_p^{n_p} \right\} \quad (96)$$

for all $j = 0, 1, 2, \dots, \infty$, $\sum_{i=1}^p n_i = j$.

Expansion of SE Kernel (Con'd)

Example: for $k=2$, when

- $j = 0$, we have term 1
- $j = 1$, we have two terms x_1, x_2
- $j = 2$, we have three terms $x_1^2, x_2^2, \sqrt{2}x_1x_2$
- $j = 3$, we have four terms $x_1^3, x_2^3, \sqrt{3}x_1^2x_2, \sqrt{3}x_2^2x_1$
- $j = 4, \dots$
- \vdots

Thus, we get

$$\phi(\mathbf{x}) = \left\{ 1, x_1, x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2, x_1^3, x_2^3, \sqrt{3}x_1^2x_2, \sqrt{3}x_2^2x_1, \dots \right\} \\ \times \sigma \exp\left(\frac{-\|\mathbf{x}\|^2}{2}\right) \quad (97)$$

Laplace Approximation

Laplace approximation of $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ given by

$$q(\mathbf{f}|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{f}|\hat{\mathbf{f}}, \mathbf{A}^{-1}) \quad (98)$$

where

$$\hat{\mathbf{f}} = \arg \max_{\mathbf{f}} p(\mathbf{f}|\mathbf{X}, \mathbf{y}), \quad (99)$$

$$\mathbf{A} = -\nabla_{\mathbf{f}} \nabla_{\mathbf{f}} \log p(\mathbf{f}|\mathbf{X}, \mathbf{y})|_{\mathbf{f}=\hat{\mathbf{f}}}. \quad (100)$$

are respectively the maximum of the posterior and the Hessian of the negative log posterior at the maximum.

Laplace Approximation (Con'd)

Due to the Bayes theorem, the maximum of the posterior $\hat{\mathbf{f}}$ can be obtained by

$$\begin{aligned}\hat{\mathbf{f}} &= \arg \max_{\mathbf{f}} p(\mathbf{f}|\mathbf{X}, \mathbf{y}) \\ &\equiv \arg \max_{\mathbf{f}} \log p(\mathbf{f}|\mathbf{X}, \mathbf{y}) \\ &\equiv \arg \max_{\mathbf{f}} \log p(\mathbf{f}) + \log p(\mathbf{y}|\mathbf{f}) \\ &\equiv \arg \max_{\mathbf{f}} -\frac{1}{2}\mathbf{f}^T K(\mathbf{X}, \mathbf{X})^{-1}\mathbf{f} + \sum_{i=1}^n \log \sigma(y_i f_i)\end{aligned}\quad (101)$$

Since the cost function is concave in terms of the latent variables \mathbf{f} , the MAP estimate $\hat{\mathbf{f}}$ is the unique global maximum.

Laplace Approximation (Con'd)

Due to the Bayes theorem, we know that

$$\log p(\mathbf{f}|\mathbf{X}, \mathbf{y}) = \log p(\mathbf{f}) + \log p(\mathbf{y}|\mathbf{f}) + \text{Const} \quad (102)$$

where the constant terms is irrespective of \mathbf{f} .

Therefore, the Hessian \mathbf{A} can be further expressed as $\mathbf{A} = \mathbf{W} + \mathbf{K}^{-1}$, where \mathbf{K} is short for $\mathbf{K}(\mathbf{X}, \mathbf{X})$ and \mathbf{W} is a diagonal matrix defined below:

$$\begin{aligned} \mathbf{W} &\triangleq -\nabla_{\mathbf{f}} \nabla_{\mathbf{f}} \log p(\mathbf{y}|\mathbf{f})|_{\mathbf{f}=\hat{\mathbf{f}}} \\ &= -\left(\nabla_{\mathbf{f}} \nabla_{\mathbf{f}} \sum_{i=1}^n \log \sigma(y_i f_i) \right) |_{\mathbf{f}=\hat{\mathbf{f}}}. \end{aligned} \quad (103)$$

Matrix Identities

The **matrix inversion lemma** states that

① for the matrix inverse,

$$\left(\mathbf{Z} + \mathbf{U}\mathbf{W}\mathbf{V}^T\right)^{-1} = \mathbf{Z}^{-1} - \mathbf{Z}^{-1}\mathbf{U}\left(\mathbf{W}^{-1} + \mathbf{V}^T\mathbf{Z}^{-1}\mathbf{U}\right)^{-1}\mathbf{V}^T\mathbf{Z}^{-1}, \quad (104)$$

② and for the matrix determinant,

$$\det\left(\mathbf{Z} + \mathbf{U}\mathbf{W}\mathbf{V}^T\right) = \det(\mathbf{Z})\det(\mathbf{W})\det(\mathbf{W}^{-1} + \mathbf{V}^T\mathbf{Z}^{-1}\mathbf{U}). \quad (105)$$

where the matrices \mathbf{Z} is of size $n \times n$, \mathbf{W} is of size $m \times m$, \mathbf{U} and \mathbf{V} are both of size $n \times m$. It is assumed that the relevant matrix inverses all exist.

- ① C. Rasmussen, Gaussian Process for Machine Learning, MIT press, 2006.
- ② A. Shashua, "Introduction to Machine Learning", 2008, available on arXiv