

CS 559: Machine Learning Fundamentals & Applications

Lecture 7: KNN, Kernel Method, and Gaussian Process





KNN: k-Nearest Neighbors

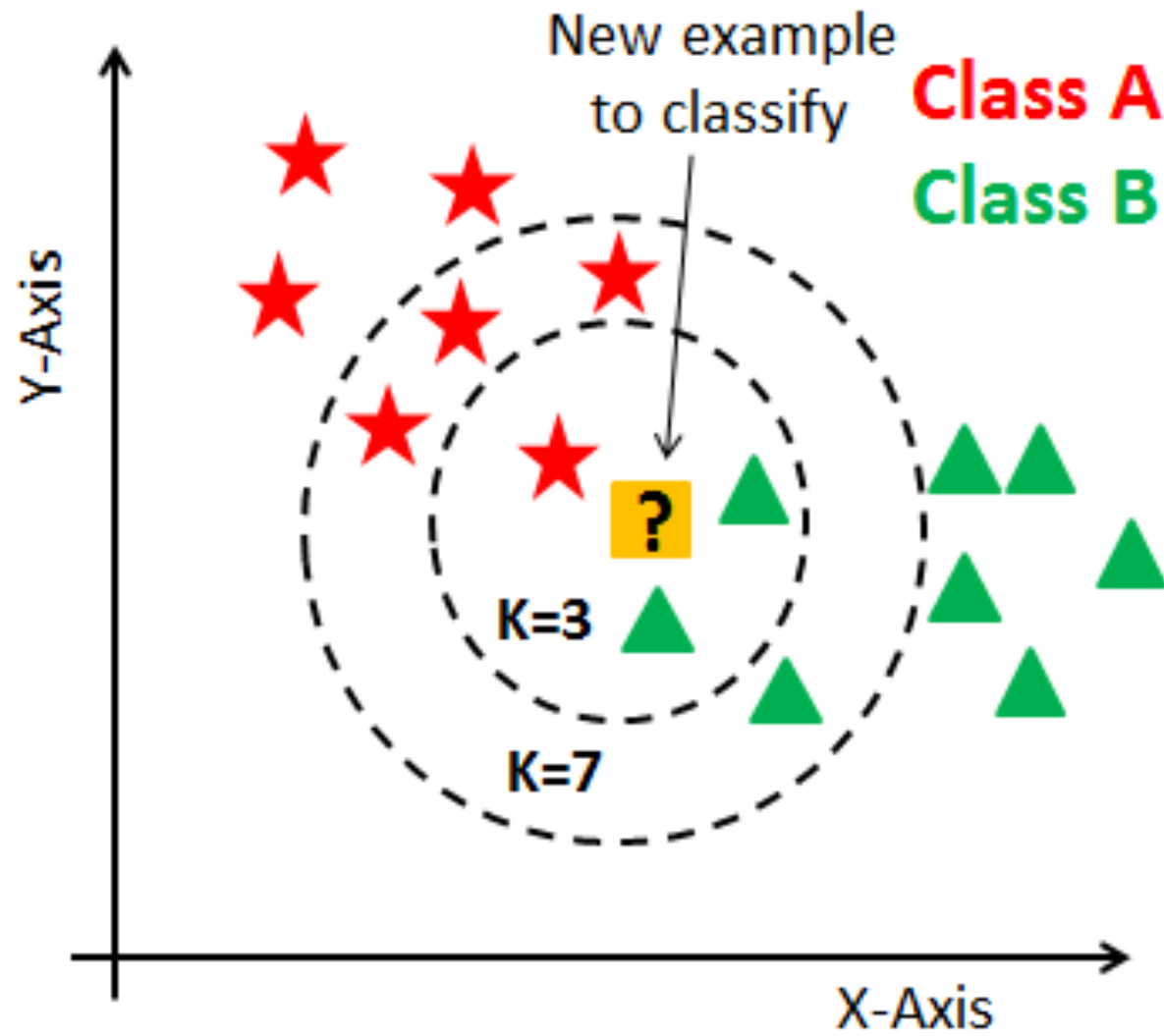
KNN: Non-parametric Models



Non-parametric model: Using instance-based learning, is characterizes by memorizing the training dataset. When the cost is 0 during the learning process, we call it lazy learning.

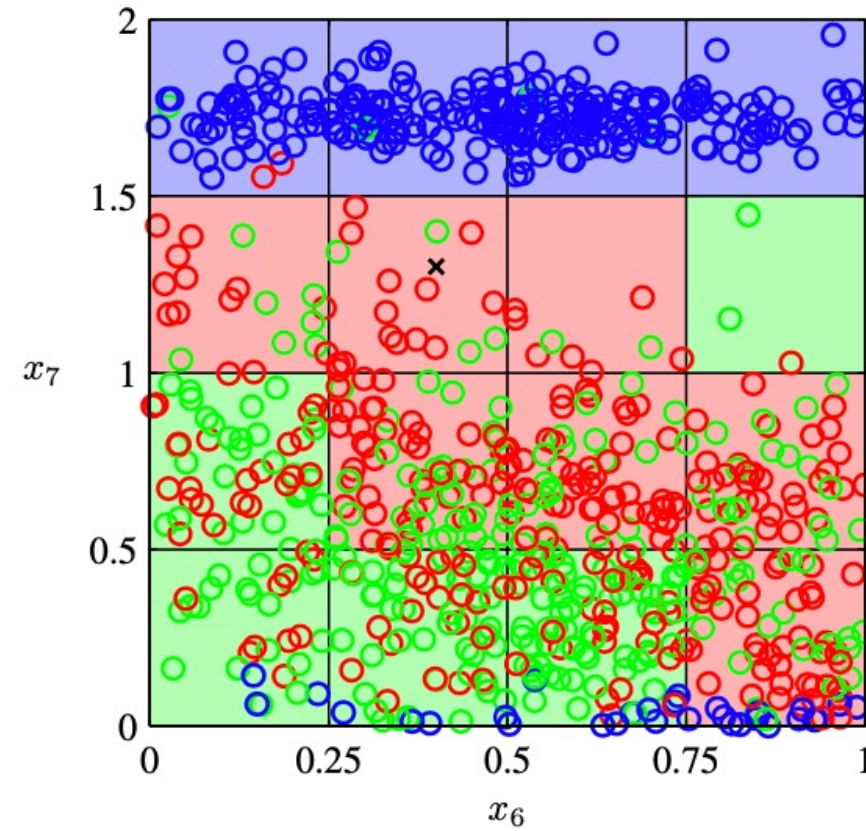
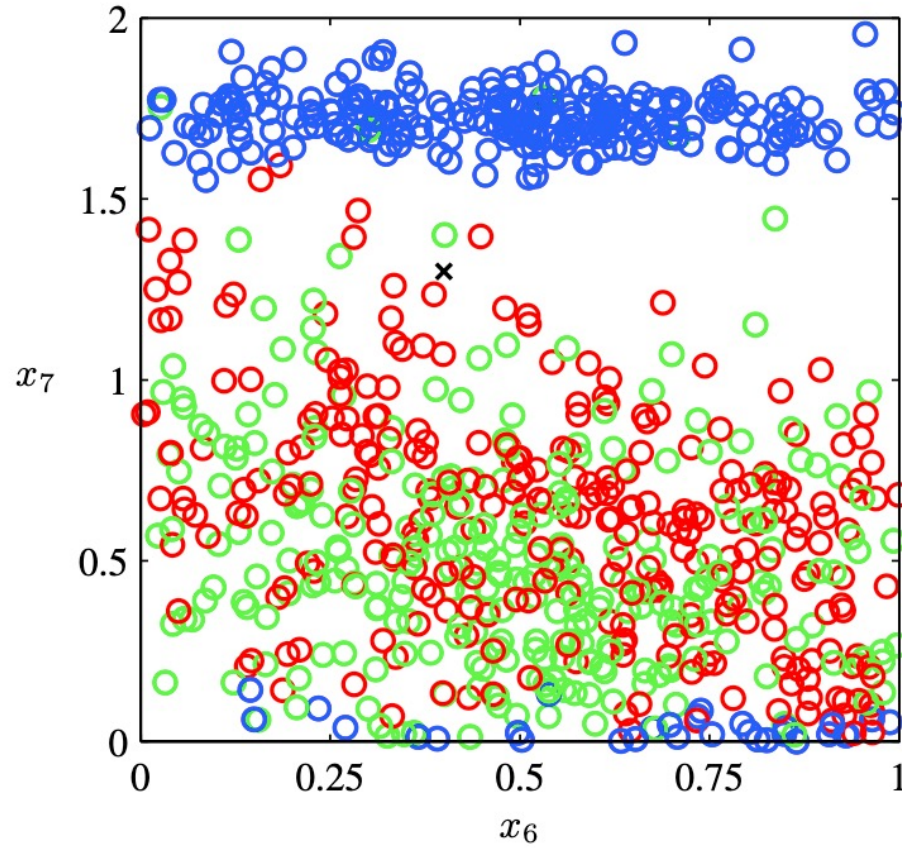
Typical non-parametric models are KNN, decision trees, random forest, and kernel SVM.

KNN



Curse of Dimensionality

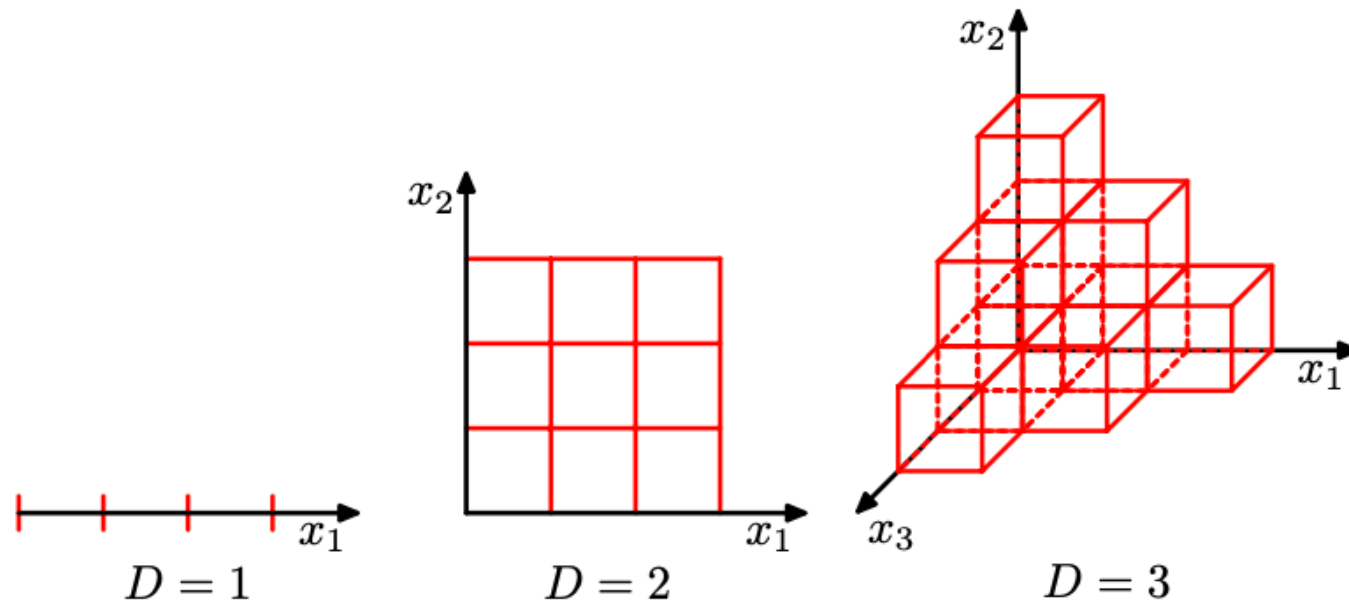
The KNN is very susceptible to overfitting due to the **curse of dimensionality**.



Curse of Dimensionality



- Consider if we are going to classify in a simple approach, there are several problems. The most severe is when we have a larger number of input variables – input spaces in higher dimensionality. The addition of one variable increases the dimension exponentially.
- The problem arises because we need an exponentially large quantity of training data in order to ensure that all cells are not empty.



Curse of Dimensionality



- For the simple consideration, let the model be $y(\mathbf{x}, \mathbf{w}) = \sum \mathbf{w} \mathbf{x}$,

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^D w_i x_i + \sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j + \dots$$

- As D increases, the number of coefficients grows proportional to D^M .
- The KNN requires all points to be close to every dimension axis.
- An additional dimension needs to make the point to be closer to the new axis.



KNN(X,k)

1. Load the training and test data
2. Choose the value of k
3. For each point in test data:
 4. Find the Euclidean distance to all training data points.
 5. Store the distances in a list and sort it.
 6. Choose the first k points.
7. Assign a class to the test point based on the majority of classes present in the chosen points.
8. End



Gaussian Distribution

Conditional Gaussian Distribution

Suppose \mathbf{x} is a D -dimensional vector with Gaussian distribution $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and that partition into two disjoint subsets \mathbf{x}_a (M components) and \mathbf{x}_b ($D-M$ components),

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}$$

and

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$$

where $\boldsymbol{\Sigma}_{ab} = \boldsymbol{\Sigma}_{ba}^T$.

Let the *precision matrix* be

$$\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$$

so that

$$\begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}^{-1} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}.$$

Using the identity of a **portioned matrix**

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} & -(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B}\mathbf{D}^{-1} \end{pmatrix},$$

the precision matrix is equivalent to

$$\boldsymbol{\Lambda} = \begin{pmatrix} (\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba})^{-1} & -(\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba})^{-1}\boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1} \\ -\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}(\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba})^{-1} & \boldsymbol{\Sigma}_{bb}^{-1} + \boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}(\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba})^{-1}\boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1} \end{pmatrix}$$

(Eq. 1)



Conditional Gaussian Distribution

By the definition, the conditional probability $p(\mathbf{x}_a|\mathbf{x}_b)$ can be expressed as

$$p(\mathbf{x}_a|\mathbf{x}_b) = \frac{p(\mathbf{x}_a \cap \mathbf{x}_b)}{p(\mathbf{x}_a)}$$

evaluated from the joint distribution $p(\mathbf{x}) = p(\mathbf{x}_a, \mathbf{x}_b)$ by fixing \mathbf{x}_b to the observed value and normalizing the resulting expression to obtain a valid probability distribution over \mathbf{x}_a .

Instead, we can obtain the solution using the **quadratic form** in the exponent of the Gaussian distribution and **reinstating the normalization** at the end.

Starting with the Gaussian,

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{(a|b)}, \boldsymbol{\Sigma}_{(a|b)}) = \frac{1}{(2\pi)^{\frac{D}{2}} |\boldsymbol{\Sigma}_{(a|b)}|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_{(a|b)} - \boldsymbol{\mu}_{(a|b)})^T \boldsymbol{\Sigma}_{(a|b)}^{-1} (\mathbf{x}_{(a|b)} - \boldsymbol{\mu}_{(a|b)}) \right\}.$$

The quadratic term, $-\frac{1}{2} (\mathbf{x}_{(a|b)} - \boldsymbol{\mu}_{(a|b)})^T \boldsymbol{\Sigma}_{(a|b)}^{-1} (\mathbf{x}_{(a|b)} - \boldsymbol{\mu}_{(a|b)})$, in the **conditional probability** can be expressed as

$$-\frac{1}{2} (\mathbf{x}_{(a|b)} - \boldsymbol{\mu}_{(a|b)})^T \boldsymbol{\Sigma}_{(a|b)}^{-1} (\mathbf{x}_{(a|b)} - \boldsymbol{\mu}_{(a|b)}) = -\frac{1}{2} \mathbf{x}_{(a|b)}^T \boldsymbol{\Sigma}_{(a|b)}^{-1} \mathbf{x}_{(a|b)} + \mathbf{x}_{(a|b)}^T \boldsymbol{\Sigma}_{(a|b)}^{-1} \boldsymbol{\mu}_{(a|b)} + C$$

(Eq. 2)

where the constant C contains the terms that independent of \mathbf{x} .

Then $\boldsymbol{\mu}$ can be obtained from $\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$ term.

Conditional Gaussian Distribution



In the same manner, consider the **conditional distribution**. Eq. 2 can be expressed by following Eq. 1. as

$$\begin{aligned} &= -\frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &\quad - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b). \end{aligned}$$

which can be expand as

$$\begin{aligned} &= -\frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{aa} \mathbf{x}_a + \mathbf{x}_a^T \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{ab} \mathbf{x}_b + \frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{ab} \boldsymbol{\mu}_b + \frac{1}{2}\mathbf{x}_b^T \boldsymbol{\Lambda}_{ab} \boldsymbol{\mu}_a - \frac{1}{2}\mathbf{x}_b^T \boldsymbol{\Lambda}_{ba} \mathbf{x}_a + \frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{ba} \boldsymbol{\mu}_b + \dots \end{aligned} \tag{Eq. 3}$$

Conditional Gaussian Distribution



In the same manner, consider the **conditional distribution**. Eq. 2 can be expressed by following Eq. 1. as

$$\begin{aligned} &= -\frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &\quad - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b). \end{aligned}$$

which can be expand as

$$= -\frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{aa} \mathbf{x}_a + \mathbf{x}_a^T \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{ab} \mathbf{x}_b + \frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{ab} \boldsymbol{\mu}_b + \frac{1}{2}\mathbf{x}_b^T \boldsymbol{\Lambda}_{ab} \boldsymbol{\mu}_a - \frac{1}{2}\mathbf{x}_b^T \boldsymbol{\Lambda}_{ba} \mathbf{x}_a + \frac{1}{2}\mathbf{x}_b^T \boldsymbol{\Lambda}_{ba} \boldsymbol{\mu}_b + \dots \quad (\text{Eq. 3})$$

The \mathbf{x}_a second-order term in Eq. 3 shows that

$$\boldsymbol{\Sigma}_{(a|b)} = \boldsymbol{\Lambda}_{aa}^{-1}. \quad (\text{Eq. 4})$$

Conditional Gaussian Distribution



In the same manner, consider the **conditional distribution**. Eq. 2 can be expressed by following Eq. 1. as

$$\begin{aligned} &= -\frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &\quad - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b). \end{aligned}$$

which can be expand as

$$\begin{aligned} &= -\frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{aa} \mathbf{x}_a + \mathbf{x}_a^T \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{ab} \mathbf{x}_b + \frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{ab} \boldsymbol{\mu}_b + \frac{1}{2}\mathbf{x}_b^T \boldsymbol{\Lambda}_{ab} \boldsymbol{\mu}_a - \frac{1}{2}\mathbf{x}_b^T \boldsymbol{\Lambda}_{ba} \mathbf{x}_a + \frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{ba} \boldsymbol{\mu}_b + \dots \end{aligned} \quad (\text{Eq. 3})$$

The \mathbf{x}_a second-order term in Eq. 3 shows that

$$\boldsymbol{\Sigma}_{(a|b)} = \boldsymbol{\Lambda}_{aa}^{-1}. \quad (\text{Eq. 4})$$

The terms are linear in \mathbf{x}_a are

$$\begin{aligned} &\mathbf{x}_a^T \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{ab} \mathbf{x}_b + \frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{ab} \boldsymbol{\mu}_b - \frac{1}{2}\mathbf{x}_b^T \boldsymbol{\Lambda}_{ba} \mathbf{x}_a + \frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{ba} \boldsymbol{\mu}_b \\ &= \mathbf{x}_a^T \{\boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ba}(\mathbf{x}_b - \boldsymbol{\mu}_b)\}. \end{aligned}$$

Conditional Gaussian Distribution



The comparison between (Eq. 2) and (Eq. 3) shows that

$$\begin{aligned} \mathbf{x}_{(a|b)}^T \boldsymbol{\Sigma}_{(a|b)}^{-1} \boldsymbol{\mu}_{(a|b)} &= \mathbf{x}_a^T \{ \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ba} (\mathbf{x}_b - \boldsymbol{\mu}_b) \} \\ \boldsymbol{\mu}_{(a|b)} &= \boldsymbol{\Sigma}_{(a|b)} \{ \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ba} (\mathbf{x}_b - \boldsymbol{\mu}_b) \} \\ \boldsymbol{\mu}_{(a|b)} &= \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ba} (\mathbf{x}_b - \boldsymbol{\mu}_b) \end{aligned}$$

(Eq. 5)

Using (Eq. 1), (Eq. 4) and (Eq. 5) can be expressed in terms of the mean and covariance of the conditional distribution

$$\begin{aligned} \boldsymbol{\Sigma}_{(a|b)} &= \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba} \\ \boldsymbol{\mu}_{(a|b)} &= \boldsymbol{\mu}_a - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} (\mathbf{x}_b - \boldsymbol{\mu}_b). \end{aligned}$$

Sequential Learning



Considering the result of MLE for the mean,

$$\begin{aligned}\mu_{ML}^{(N)} &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} \mathbf{x}_N + \frac{1}{N} \sum_{n=1}^{N-1} \mathbf{x}_n \\ &= \frac{1}{N} \mathbf{x}_N + \frac{(N-1)}{N} \mu_{ML}^{(N-1)} \\ &= \mu_{ML}^{(N-1)} + \frac{1}{N} \left(\mathbf{x}_N - \mu_{ML}^{(N-1)} \right),\end{aligned}$$

the contribution from data points gets smaller as N increases. However, this solution is the equivalent solution as a general MLE and we need a general derivable formulation.

Sequential Learning

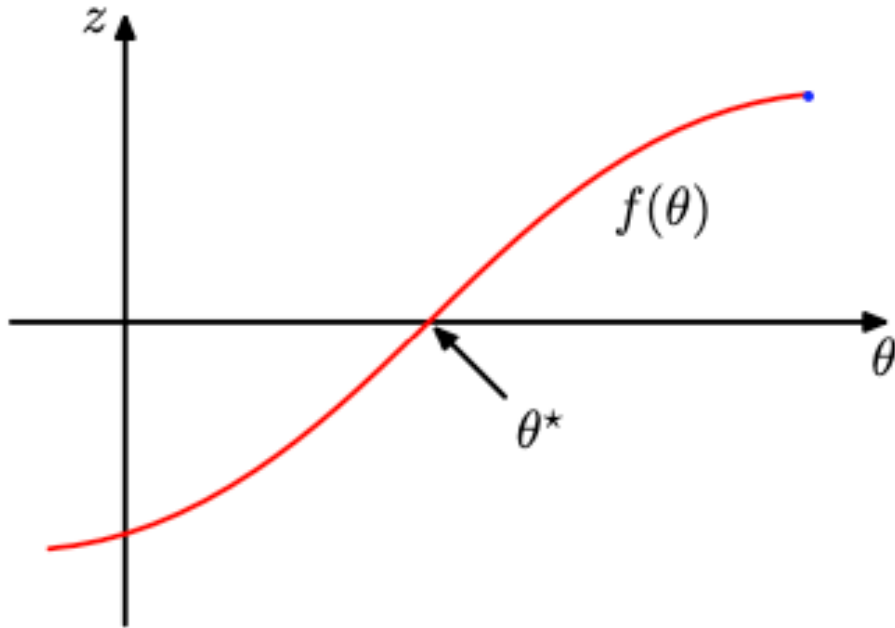


The *Robbins-Monro* algorithm allows for the construction of a more general formulation of sequential learning.

Consider a pair of random variables θ and z in a joint distribution $p(z, \theta)$.

The conditional expectation, $\mathbb{E}[z|\theta]$, defines a deterministic function $f(\theta)$ that is given by

$$f(\theta) = \mathbb{E}[z|\theta] = \int z p(z|\theta) dz .$$



Sequential Learning



The goal is to find the root θ^* : $f(\theta^*) = 0$.

Suppose we have a large data set of observation.

We can find the regression model directly but the sequential estimation scheme for θ^* is not easy.

The following conditions are given by Robbins and Monro.

1. The conditional variance of z is finite: $\mathbb{E}[(z - f)^2 | \theta] < \infty$.
2. A sequence of successive estimates of θ^* is defined as $\theta^{(N)} = \theta^{(N-1)} - a_{N-1}z(\theta^{(N-1)})$.
3. The coefficients $\{a_N\}$ satisfy the conditions.

$$\begin{aligned}\lim_{N \rightarrow \infty} a_N &= 0 \\ \sum_{N=1}^{\infty} a_N &= \infty \\ \sum_{N=1}^{\infty} a_N^2 &< \infty\end{aligned}$$

The 2nd condition shows that the sequence will converge to the root with probability 1.

The 3-1 condition shows that the successive correlations decrease in magnitude so that the process can converge to a limiting value. The 3-2 condition ensures the algorithm does not converge short of the root. The 3-3 condition ensures that the accumulated noise has finite variance and does not spoil convergence.

Sequential Learning



By the definition, the MLE solution θ_{ML} is a stationary point of the negative log-likelihood function and satisfies

$$\left. \frac{\partial}{\partial \theta} \left\{ -\frac{1}{N} \sum_{n=1}^N \ln p(x_n | \theta) \right\} \right|_{\theta_{ML}} = 0.$$

By taking the limit $N \rightarrow \infty$, the equation above becomes

$$-\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N \frac{\partial}{\partial \theta} \ln p(x_n | \theta) = \mathbb{E}_x \left[-\frac{\partial}{\partial \theta} \ln p(x_n | \theta) \right].$$

We can apply the Robbins-Monro procedure,

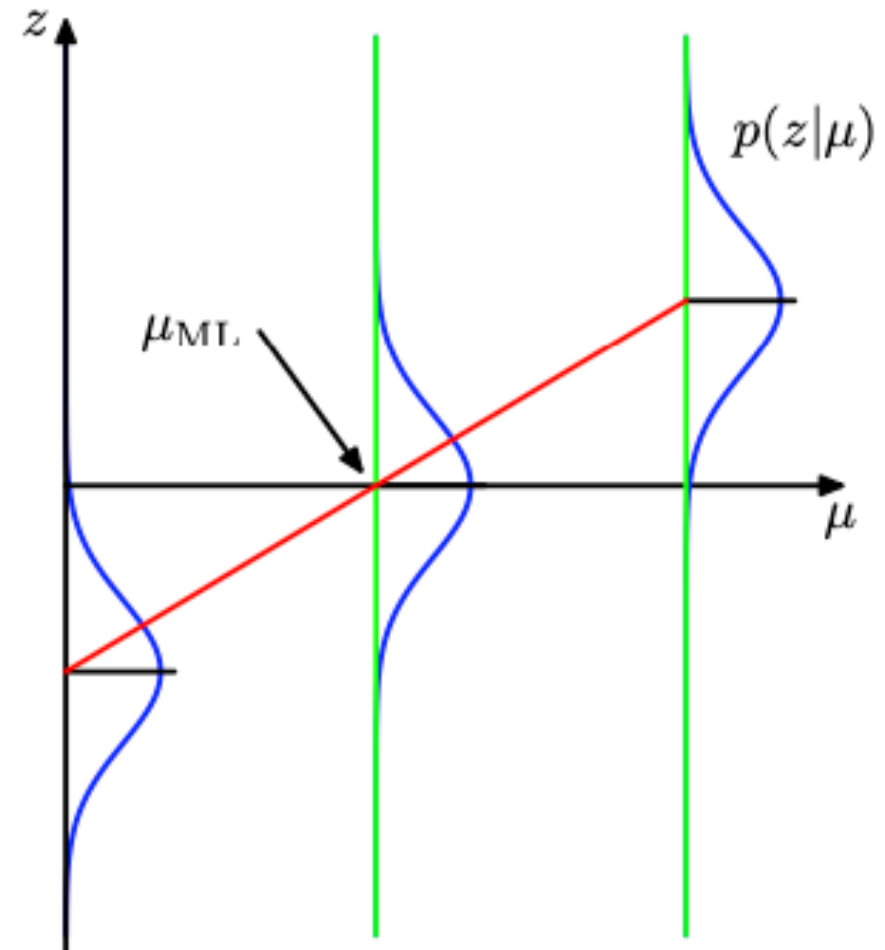
$$\theta^{(N)} = \theta^{(N-1)} - a_{N-1} \frac{\partial}{\partial \theta^{N-1}} [-\ln p(x_N | \theta^{N-1})].$$

Since $\theta^{(N)} = \mu_{ML}^{(N)}$,

$$z = \frac{\partial}{\partial \mu_{ML}} [-\ln p(x | \mu_{ML}, \sigma^2)] = -\frac{1}{\sigma^2} (x - \mu_{ML}).$$

The Gaussian z distribution has the mean of $-(\mu - \mu_{ML})/\sigma^2$.

This allows us to choose the coefficients a_N in the form σ^2/N .





Kernel Method



Dual Presentation

Many linear parametric models can be re-cast into an equivalent *dual representation* that based on linear combinations of a *kernel function* evaluated at the training data points.

For models are based on a fixed nonlinear *feature space* mapping $\phi(\mathbf{x})$, the kernel function is given by

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

Dual Presentation

Consider a linear regression model

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

where $\lambda \geq 0$. Setting $\nabla J(\mathbf{w}) = 0$, the solution of \mathbf{w} is

$$\mathbf{w} = -\frac{1}{\lambda} \sum_{n=1}^N \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\} \phi(\mathbf{x}_n) = \sum_{n=1}^N a_n \phi(\mathbf{x}_n) = \Phi^T \mathbf{a}$$

where $\mathbf{a} = (a_1, \dots, a_N)^T$ and a_n is defined as

$$a_n = -\frac{1}{\lambda} \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}.$$

The least-squares algorithm can be reformulated using a *dual representation* with the parameter vector \mathbf{w} by substituting $\mathbf{w} = \Phi^T \mathbf{a}$ into $J(\mathbf{w})$:

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a}$$

(Eq. 6)

where $\mathbf{t} = (t_1, \dots, t_N)^T$.



Dual Presentation

Eq. 6 can be simplified using the *Gram* matrix $\mathbf{K} = \Phi\Phi^T$ which is an $N \times N$ symmetric matrix with elements

$$K_{nm} = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$$

where $k(x, x')$ is the *kernel function*:

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}.$$

(Eq. 7)

Setting the gradient of $J(\mathbf{a})$ w.r.t. to \mathbf{a} to 0,

$$\begin{aligned} \nabla_{\mathbf{a}} J(\mathbf{a}) &= 0 = \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{K} \mathbf{t} + \lambda \mathbf{K} \mathbf{a} \\ \mathbf{a} &= (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}. \end{aligned}$$

(Eq. 8)



Dual Presentation

Substituting Eq. 8 back to the linear model, the prediction for a new \mathbf{x} can be found as

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \mathbf{a}^T \Phi \phi(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

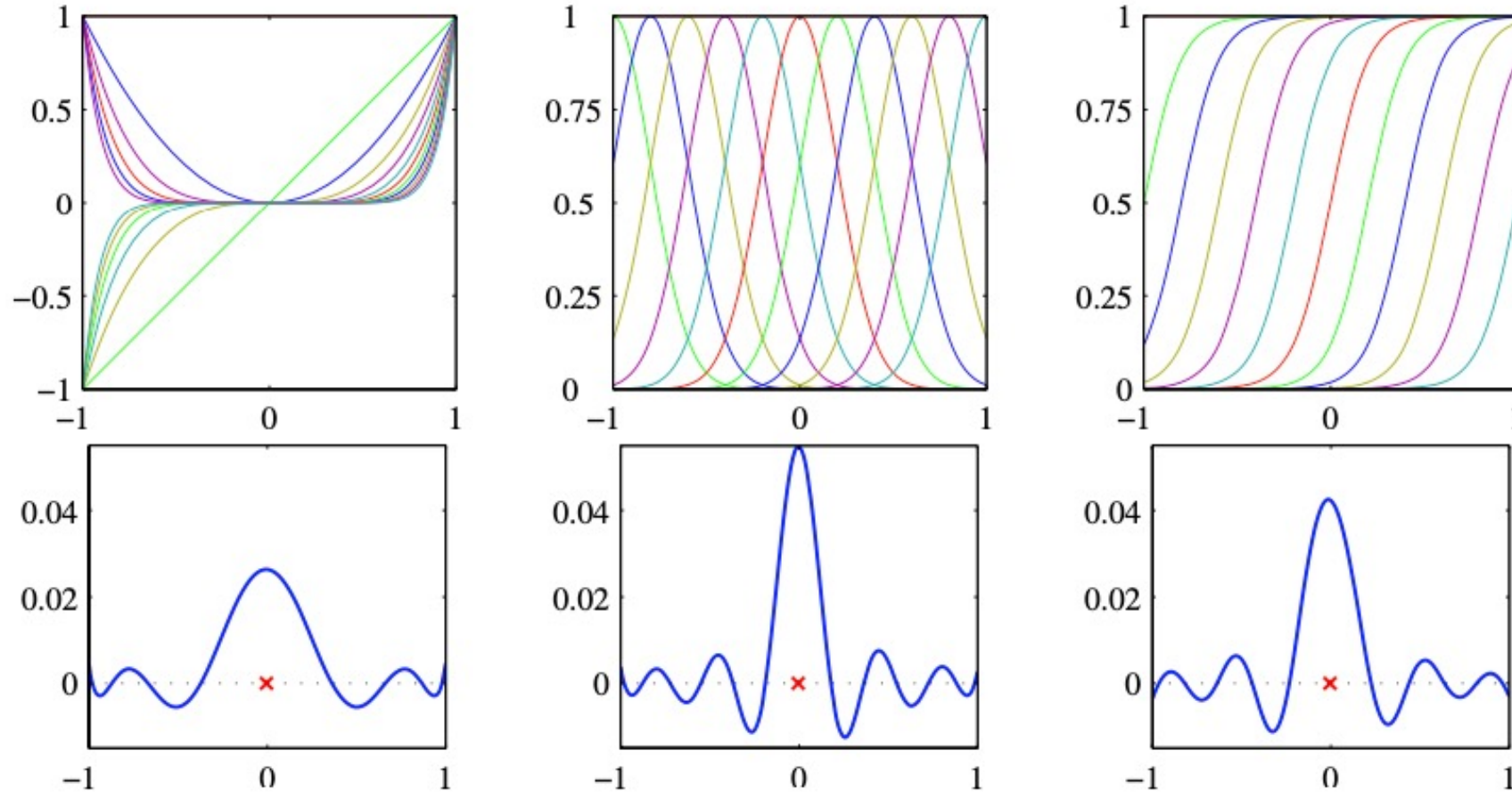
(Eq. 9)

where the vector $\mathbf{k}(\mathbf{x})$ has elements $k_n(\mathbf{x}) = k(\mathbf{x}_n, \mathbf{x})$.

Eq. 9 expresses the solution completely in $k(\mathbf{x}, \mathbf{x}')$ where \mathbf{a} can be expressed as a linear combination of $\phi(\mathbf{x})$. This will allow us to recover the solution in terms of \mathbf{w} .

Kernel Construction

To construct the valid kernel functions, we can map a feature space $\phi(\mathbf{x})$ and find the corresponding kernel validity.





Kernel Construction

In 1-D input space, the kernel function is defined as

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x').$$

Or we can construct kernel functions directly using the valid kernel that is a scalar product in some feature space.

For example, consider the following

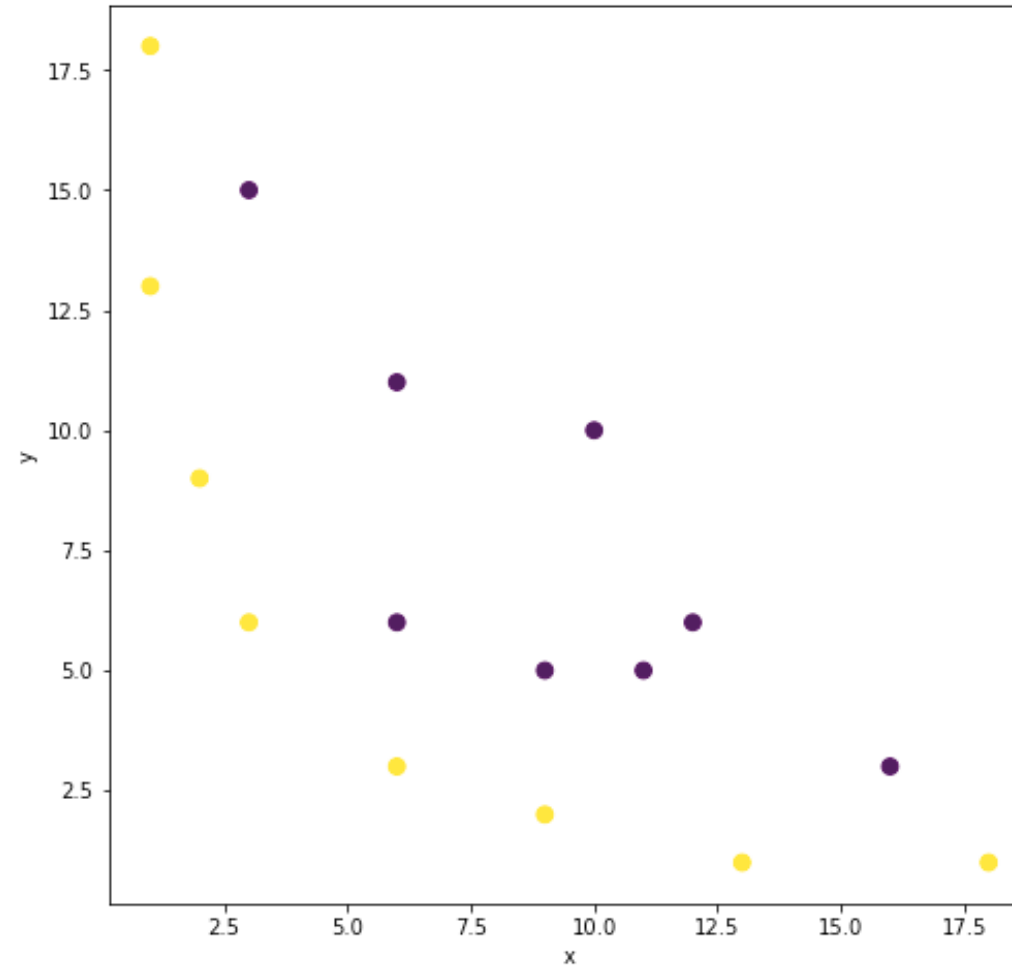
$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2.$$

Let the input space in 2-D be $\mathbf{x} = (x_1, x_2)$.

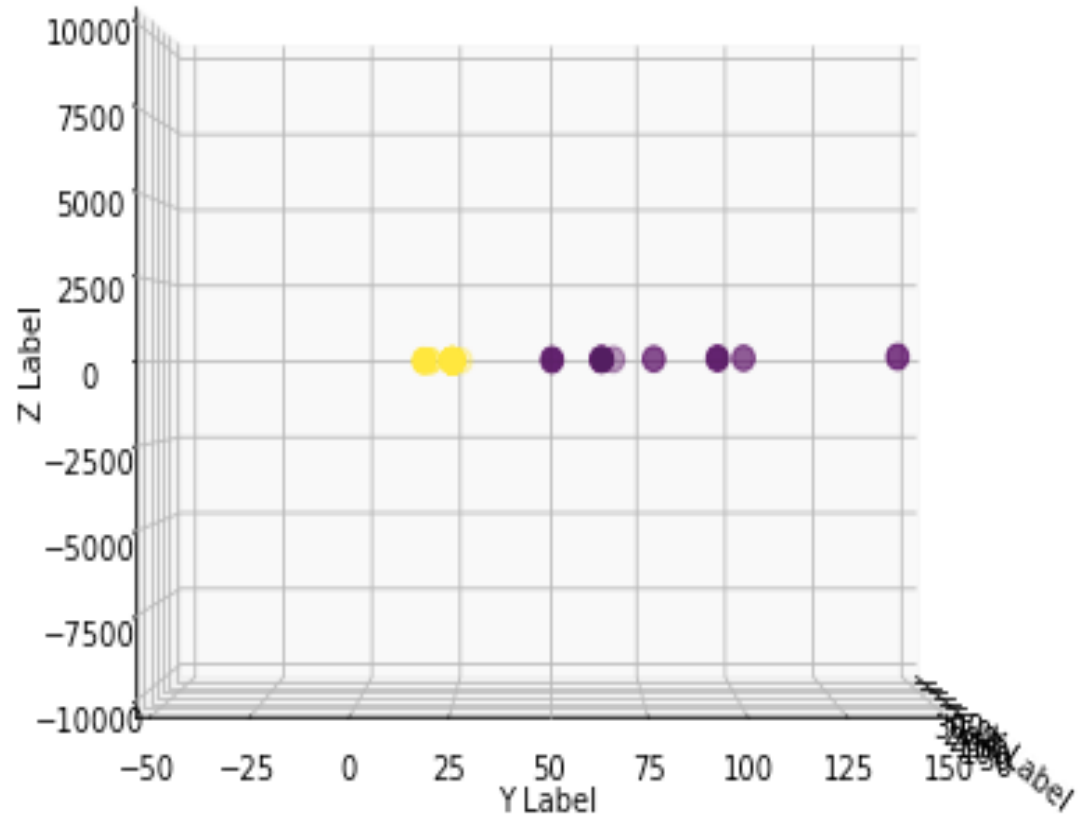
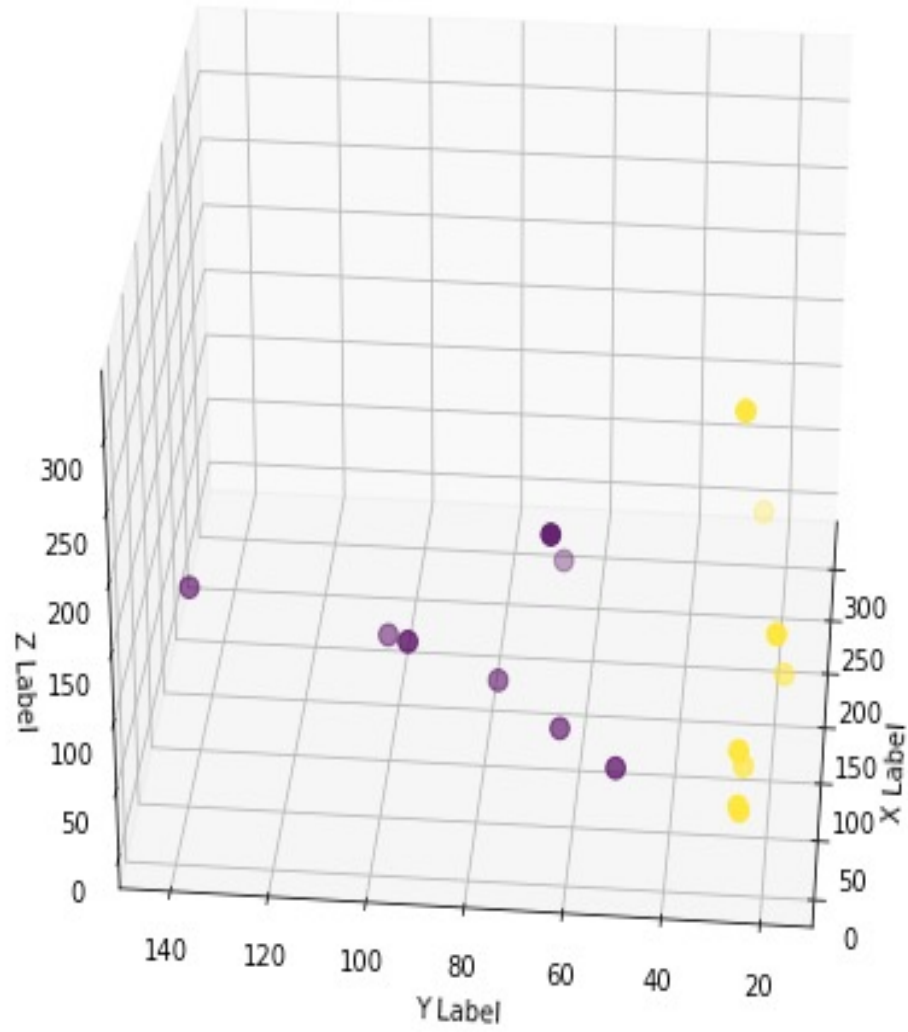
If we expand out the terms and identify the corresponding nonlinear feature mapping

$$\begin{aligned} k(\mathbf{x}, \mathbf{z}) &= (\mathbf{x}^T \mathbf{z})^2 = (x_1 z_1 + x_2 z_2)^2 \\ &= x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2 \\ &= (x_1^2, \sqrt{2}x_1 x_2, x_2^2)(z_1^2, \sqrt{2}z_1 z_2, z_2^2)^T \\ &= \phi(\mathbf{x})^T \phi(\mathbf{z}). \end{aligned}$$

Kernel Construction



Kernel Construction





Kernel Construction

As long as the Gram matrix \mathbf{K} has elements $k(\mathbf{x}_n, \mathbf{x}_m)$ are positive semidefinite for all possible choices of the set $\{\mathbf{x}_n\}$, the validity of kernel is confirmed.

Positive semidefinite matrices are the scalars $\mathbf{x}^T \mathbf{M} \mathbf{x}$ and $\mathbf{x}^* \mathbf{M} \mathbf{x}$ are positive or zero where \mathbf{x}^* is a conjugate transpose of \mathbf{x} .

Check if matrix \mathbf{A} is a positive definite.

$$\mathbf{A} = \begin{bmatrix} 9 & -15 \\ -15 & 25 \end{bmatrix}$$

Let \mathbf{x} be a 2×1 vector, $\mathbf{x} = [x_1, x_2]$.

$$\begin{aligned} \mathbf{x}^T \mathbf{A} \mathbf{x} &= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \begin{bmatrix} 9 & -15 \\ -15 & 25 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &= 9x_1^2 - 15x_1x_2 - 15x_1x_2 + 25x_2^2 \\ &= (3x_1 - 5x_2)^2 \end{aligned}$$

If $\mathbf{x} = [5, 3]$, $\mathbf{x}^T \mathbf{A} \mathbf{x} = 0$.

If $\mathbf{x} = [5, 2]$, $\mathbf{x}^T \mathbf{A} \mathbf{x} = 25$.



Kernel Construction

Another commonly used kernel function is a “Gaussian” kernel also known as RBF and is defined as

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{|\mathbf{x} - \mathbf{x}'|^2}{2\sigma^2}\right)$$

where $|\mathbf{x} - \mathbf{x}'|^2 = \mathbf{x}^T \mathbf{x} + (\mathbf{x}')^T \mathbf{x}' - 2\mathbf{x}^T \mathbf{x}'$ and this gives

$$\begin{aligned} k(\mathbf{x}, \mathbf{x}') &= \exp\left(-\frac{\mathbf{x}^T \mathbf{x}}{2\sigma^2}\right) \exp\left(\frac{\mathbf{x}^T \mathbf{x}'}{\sigma^2}\right) \exp\left(-\frac{\mathbf{x}'^T \mathbf{x}'}{2\sigma^2}\right) \\ &= \exp\left\{-\frac{1}{2\sigma^2} (k(\mathbf{x}, \mathbf{x}) + k(\mathbf{x}', \mathbf{x}') - 2k(\mathbf{x}, \mathbf{x}'))\right\}. \end{aligned}$$



Gaussian Processes

Gaussian Processes



By extending the roles of kernels to probabilistic discriminative models, we can see how kernels arise in a Bayesian setting and lead to the framework of Gaussian processes.

Recall linear regression model $y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$.

- A prior distribution over \mathbf{w} induced a corresponding prior distribution over functions $y(\mathbf{x}, \mathbf{w})$.
- Then the posterior distribution over \mathbf{w} is evaluated and obtained the corresponding posterior distribution over regression functions that predict new input data \mathbf{x} .

In the Gaussian process, we define the prior probability distribution over functions directly.

- for a finite training set, we only need to consider the values of the function at the discrete set of input values.



Linear Regression Revisit

Consider a model $y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$ and a prior distribution over \mathbf{w} given by

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}).$$

Since we are interested in the function evaluation at specific value of \mathbf{x} , we are also interested in the joint distribution of the function values $y(\mathbf{x}_1), \dots, y(\mathbf{x}_N)$,

$$\mathbf{y} = \Phi \mathbf{w}.$$

The mean and covariance of \mathbf{y} are

$$\mathbb{E}[\mathbf{y}] = \Phi \mathbb{E}[\mathbf{w}] = \mathbf{0}$$

$$\text{cov}[\mathbf{y}] = \mathbb{E}[\mathbf{y}\mathbf{y}^T] = \Phi \mathbb{E}[\mathbf{w}\mathbf{w}^T] \Phi^T = \frac{1}{\alpha} \Phi \Phi^T = \mathbf{K}$$

where \mathbf{K} is the Gram matrix with elements

$$K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m) = \frac{1}{\alpha} \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m).$$

Gaussian Processes



In general, a GP is defined as a probability distribution over functions $y(\mathbf{x})$ s.t. the set of $y(\mathbf{x})$ values are evaluated at an arbitrary set of points \mathbf{x} jointly have a Gaussian distribution.

The stochastic process is the joint distribution over N variables y_1, \dots, y_N that are specified by the mean and the covariance.

- We do not know the mean of $y(\mathbf{x})$.
- But we can take it to be zero symmetrically.
- This allows the mean of the prior over weight values $p(\mathbf{w}|\alpha)$ be zero in the basis function viewpoint.
- Then, we give the covariance of $y(\mathbf{x})$ at any two values of \mathbf{x} that is given by the kernel function $\mathbb{E}[y(\mathbf{x}_n), y(\mathbf{x}_m)] = k(\mathbf{x}_n, \mathbf{x}_m)$.

The **Ornstein-Uhlenbeck process** allows defining the kernel function directly.



Gaussian Process for regression

Let the observed target values be

$$t_n = y_n + \epsilon_n$$

where $y_n = y(\mathbf{x}_n)$ and ϵ_n is the random noise variable chosen i.i.d,

$$p(t_n|y_n) = \mathcal{N}(t_n|y_n, \beta^{-1}).$$

The joint distribution of \mathbf{t} conditioned on \mathbf{y} is an isotropic Gaussian form

$$p(\mathbf{t}|\mathbf{y}) = \mathcal{N}(\mathbf{t}|\mathbf{y}, \beta^{-1}\mathbf{I}_N).$$

The marginal distribution $p(\mathbf{y})$ is given by a Gaussian

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}).$$

The marginal distribution $p(\mathbf{t})$ conditioned on inputs is

$$p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{y})p(\mathbf{y})d\mathbf{y} = \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C})$$

where the covariance matrix \mathbf{C} has elements

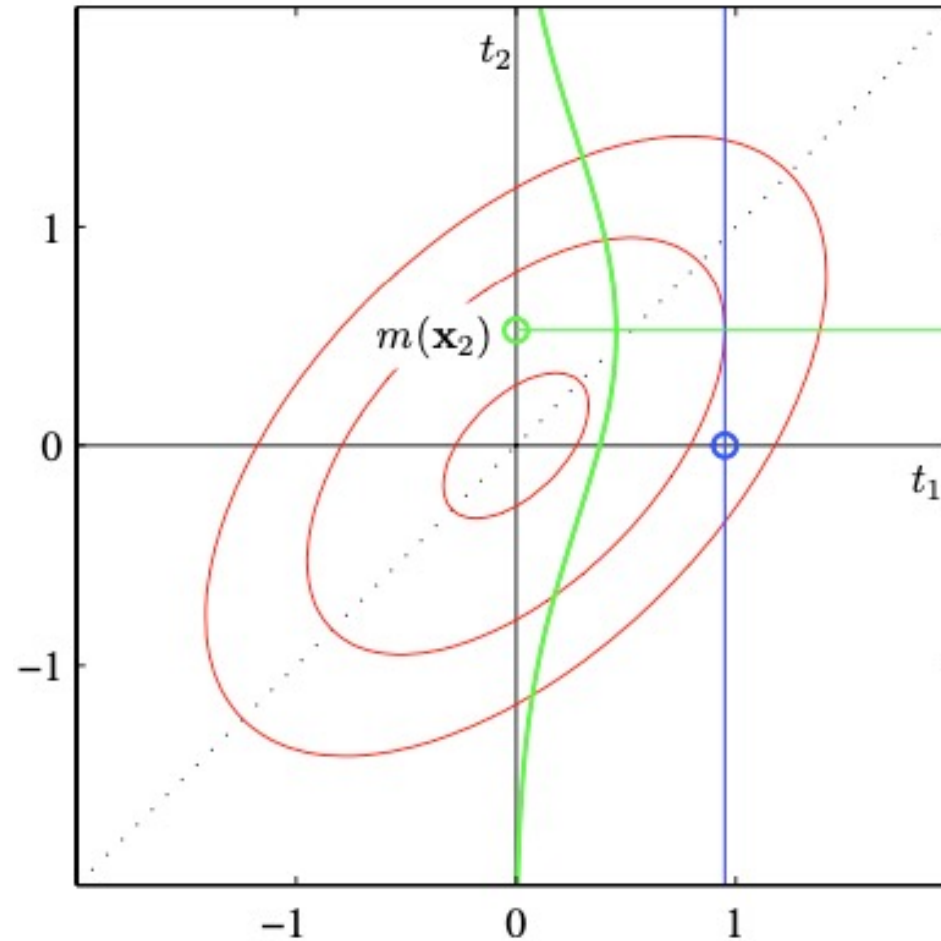
$$\mathbf{C} = \begin{pmatrix} k(\mathbf{x}_n, \mathbf{x}_n) + \beta^{-1} & k(\mathbf{x}_n, \mathbf{x}_m) \\ k(\mathbf{x}_m, \mathbf{x}_n) & k(\mathbf{x}_m, \mathbf{x}_m) + \beta^{-1} \end{pmatrix}$$
$$C(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \beta^{-1}\delta_{nm}.$$

(Eq. 10)

Gaussian Processes

Suppose the training data set $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$ and $\mathbf{t}_N = (t_1, \dots, t_N)^T$.

If the new data has the input vector \mathbf{x}_{N+1} with the target variable t_{N+1} , the predictive distribution $p(t_{N+1}|\mathbf{t}_N)$.



Gaussian Processes



The joint distribution $p(\mathbf{t}_{N+1})$ where the target vector is $\mathbf{t}_{N+1} = (t_1, \dots, t_N, t_{N+1})^T$ can be expressed as

$$p(\mathbf{t}_{N+1}) = \mathcal{N}(\mathbf{t}_{N+1} | \mathbf{0}, \mathbf{C}_{N+1})$$

where \mathbf{C}_{N+1} is the covariance matrix with elements given by (Eq. 10).

Note that the joint distribution is Gaussian, we can apply the Gaussian distribution discussed above.

The covariance matrix partition is then

$$\mathbf{C}_{N+1} = \begin{pmatrix} \mathbf{C}_N & k(\mathbf{x}_n, \mathbf{x}_{N+1}) \\ k^T(\mathbf{x}_n, \mathbf{x}_{N+1}) & k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1} \end{pmatrix}$$

for $n = 1, \dots, N$.

Using the result from (Eq. 5), the mean and covariance of the conditional distribution of $p(t_{N+1} | \mathbf{t})$ are

$$\begin{aligned} m(\mathbf{x}_{N+1}) &= \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t} \\ \sigma^2(\mathbf{x}_{N+1}) &= k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1} - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}. \end{aligned}$$



Learning the hyperparameters

The hyperparameters' learning is based on the evaluation of likelihood function $p(\mathbf{t}|\boldsymbol{\theta})$ and $\boldsymbol{\theta}$ can be estimated by maximizing the log likelihood function.

The maximization of the log likelihood can be done by using the gradient-based optimization algorithms.

Start with the standard form for a multivariate Gaussian distribution

$$\ln p(\mathbf{t}|\boldsymbol{\theta}) = -\frac{1}{2}\ln|\mathbf{C}_N| - \frac{1}{2}\mathbf{t}^T \mathbf{C}_N^{-1} \mathbf{t} - \frac{N}{2}\ln(2\pi)$$

and set the gradient w.r.t. θ_i equal to 0 for the θ_i estimation.

Learning the hyperparameters

The derivative of the inverse of a matrix \mathbf{A}^{-1} is

$$\frac{\partial}{\partial x} \mathbf{A}^{-1} = -\mathbf{A}^{-1} \left(\frac{\partial \mathbf{A}}{\partial x} \right) \mathbf{A}^{-1}.$$

(Eq. 11)

Using (Eq. 11), the gradient of \mathbf{C}_N^{-1} is

$$\frac{\partial}{\partial \theta_i} \mathbf{C}_N^{-1} = -\mathbf{C}_N^{-1} \left(\frac{\partial \mathbf{C}_N}{\partial \theta_i} \right) \mathbf{C}_N^{-1}.$$

(Eq. 12)

Using the derivative of $\log|\mathbf{A}|$,

$$\frac{\partial}{\partial x} \ln |\mathbf{A}| = \text{Tr} \left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \right),$$

the derivative of $\ln|\mathbf{C}_N|$ is

$$\frac{\partial}{\partial \theta_i} \ln |\mathbf{C}_N| = \text{Tr} \left(\mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_i} \right).$$

(Eq. 13)

Combining (Eq. 12) and (Eq. 13), the derivation of the Gaussian distribution is

$$\frac{\partial}{\partial \theta_i} \ln p(\mathbf{t}|\boldsymbol{\theta}) = -\frac{1}{2} \text{Tr} \left(\mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_i} \right) + \frac{1}{2} \mathbf{t}^T \mathbf{C}_N^{-1} \left(\frac{\partial \mathbf{C}_N}{\partial \theta_i} \right) \mathbf{C}_N^{-1} \mathbf{t}.$$

Gaussian Processes



Because $\ln p(\mathbf{t}|\boldsymbol{\theta})$ is generally a nonconvex function, there can be multiple maxima.

The estimation is straightforward:

1. a prior over $\boldsymbol{\theta}$
2. log posterior via gradient-based methods.

While it is a Bayesian treatment, finding the exact marginalization is not possible. Instead, we need to approximate the marginalization.

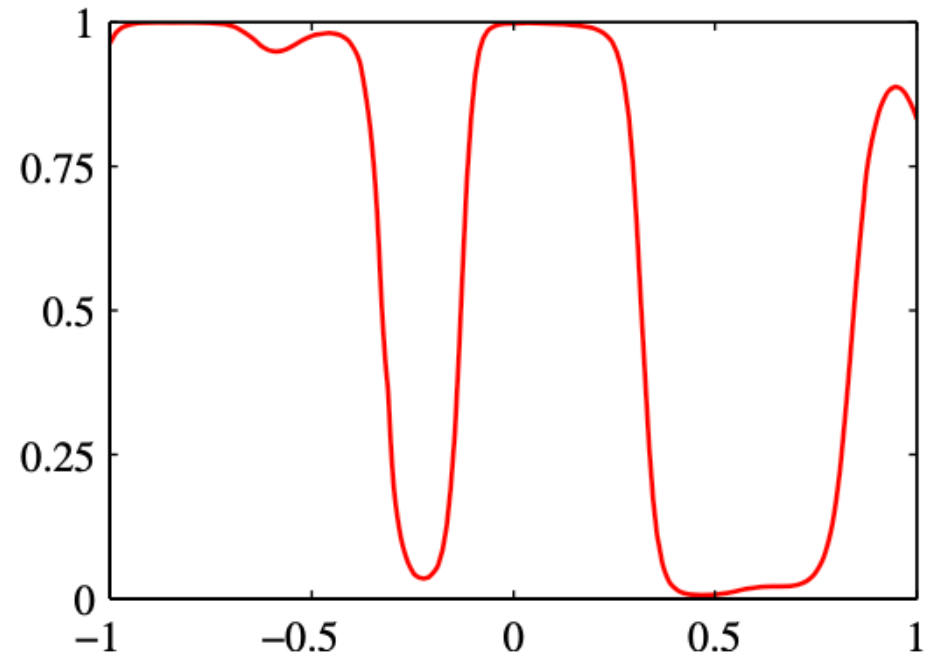
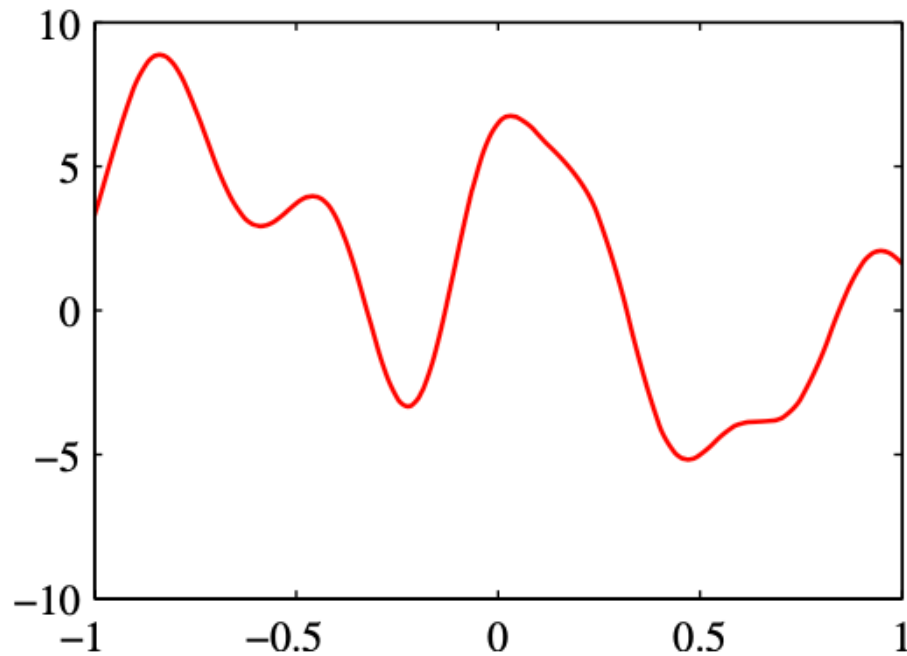
Gaussian Process for Classification

Consider the binary class problem with a target variable $t \in \{0,1\}$.

Define a Gaussian process over a function $a(\mathbf{x})$ and transform the function using a sigmoid function, $y = \sigma(a)$. Obtain a non-Gaussian stochastic process over functions $y(\mathbf{x})$ where $y \in (0,1)$.

The Bernoulli distribution over the target variable t is given as

$$p(t|a) = \sigma(a)^t (1 - \sigma(a))^{1-t}.$$





Gaussian Process for Classification

Similar to the regression problem, we introduce a Gaussian process prior over the vector \mathbf{a}_{N+1} and a non-Gaussian process over \mathbf{t}_{N+1} .

The Gaussian process prior for \mathbf{a}_{N+1} :

$$p(\mathbf{a}_{N+1}) = \mathcal{N}(\mathbf{a}_{N+1} | \mathbf{0}, \mathbf{C}_{N+1})$$

The covariance matrix is

$$C(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \nu \delta_{nm}$$

where ν is a fixed noise term.

The predictive distribution is

$$p(t_{N+1} = 1 | \mathbf{t}_N) = \int p(t_{N+1} = 1 | a_{N+1}) p(a_{N+1} | \mathbf{t}_N) da_{N+1}$$

where $p(t_{N+1} = 1 | a_{N+1}) = \sigma(a_{N+1})$.

Gaussian Process for Classification

The posterior distribution over a_{N+1} using Bayes' theorem given by

$$p(a_{N+1}|\mathbf{t}_N) = \int p(a_{N+1}, \mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N = \frac{1}{p(\mathbf{t}_N)} \int p(a_{N+1}, \mathbf{a}_N) p(\mathbf{t}_N | a_{N+1}, \mathbf{a}_N) d\mathbf{a}_N$$
$$\frac{1}{p(\mathbf{t}_N)} \int p(a_{N+1} | \mathbf{a}_N) p(\mathbf{a}_N) p(\mathbf{t}_N | \mathbf{a}_N) d\mathbf{a}_N = \int p(a_{N+1} | \mathbf{a}_N) p(\mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N$$

where $p(\mathbf{t}_N | a_{N+1}, \mathbf{a}_N) = p(\mathbf{t}_N | \mathbf{a}_N)$.

The conditional distribution $p(a_{N+1} | \mathbf{a}_N)$ is obtained by invoking the result seen from the regression,

$$p(a_{N+1} | \mathbf{a}_N) = \mathcal{N}(a_{N+1} | \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{a}_N, c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}).$$

Gaussian Process for Classification

The prior $p(\mathbf{a}_N)$ is given by a zero-mean Gaussian process with \mathbf{C}_N and the data term is given by:

$$p(\mathbf{t}_N | \mathbf{a}_N) = \prod_{n=1}^N \sigma(a_n)^{t_n} (1 - \sigma(a_n))^{(1-t_n)} = \prod_{n=1}^N e^{a_n t_n} \sigma(-a_n).$$

The approximation can be obtained by Taylor expanding the log of $p(\mathbf{t}_N | \mathbf{a}_N)$, $\Psi(\mathbf{a}_N)$, is

$$\begin{aligned} \Psi(\mathbf{a}_N) &= \ln p(\mathbf{a}_N) + \ln p(\mathbf{t}_N | \mathbf{a}_N) = \\ &= -\frac{1}{2} \mathbf{a}_N^T \mathbf{C}_N^{-1} \mathbf{a}_N - \frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{C}_N| + \mathbf{t}_N^T \mathbf{a}_N - \sum_{n=1}^N \ln(1 + e^{a_n}) + \text{const.} \end{aligned}$$

To find the mode of the posterior distribution, we need to evaluate the gradient of $\Psi(\mathbf{a}_N)$ that is given by

$$\nabla \Psi(\mathbf{a}_N) = \mathbf{t}_N - \boldsymbol{\sigma}_N - \mathbf{C}_N^{-1} \mathbf{a}_N$$

where $\boldsymbol{\sigma}_N$ is a vector with elements $\sigma(a_n)$.

Gaussian Process for Classification

The second derivative of $\Psi(\mathbf{a}_N)$ is

$$\nabla\nabla\Psi(\mathbf{a}_N) = -\mathbf{W}_N - \mathbf{C}_N^{-1}$$

where \mathbf{W}_N is a diagonal matrix with elements $\sigma(a_n)(1 - \sigma(a_n))$ in the range of $(1, 1/4)$ and a positive definite matrix. Since \mathbf{C}_N^{-1} is also positive definite matrix, the Hessian matrix $-\nabla\nabla\Psi(\mathbf{a}_N)$ is also positive definite.

Using the Newton-Raphson formula, the iterative update equation is

$$\mathbf{a}_N^{new} = \mathbf{C}_N(\mathbf{I} + \mathbf{W}_N\mathbf{C}_N)^{-1}\{\mathbf{t}_N - \boldsymbol{\sigma}_N + \mathbf{W}_N\mathbf{a}_N\}.$$

At the mode, the gradient $\nabla\Psi(\mathbf{a}_N)$ will vanish and \mathbf{a}_N^* will satisfy

$$\mathbf{a}_N^* = \mathbf{C}_N(\mathbf{t}_N - \boldsymbol{\sigma}_N).$$

Once the mode of the posterior is found, the Hessian matrix can be evaluated:

$$\mathbf{H} = -\nabla\nabla\Psi(\mathbf{a}_n) = \mathbf{W}_N + \mathbf{C}_N^{-1}.$$

This defines the Gaussian approximation of the posterior distribution $p(\mathbf{a}_N|\mathbf{t}_N)$ given by

$$q(\mathbf{a}_N) = \mathcal{N}(\mathbf{a}_N|\mathbf{a}_N^*, \mathbf{H}^{-1}).$$



Gaussian Process for Classification

The combination of $p(a_{N+1}|\mathbf{a}_N)$ and $q(\mathbf{a}_N)$ can obtain the result

$$\begin{aligned}\mathbb{E}[a_{N+1}|\mathbf{t}_N] &= \mathbf{k}^T(\mathbf{t}_N - \boldsymbol{\sigma}_N) \\ \text{var}[a_{N+1}|\mathbf{t}_N] &= c - \mathbf{k}^T(\mathbf{W}_N^{-1} + \mathbf{C}_N)^{-1}\mathbf{k}.\end{aligned}$$

Then the determination of $\boldsymbol{\theta}$ of the covariance function can be found by maximizing the likelihood $p(\mathbf{t}_N|\boldsymbol{\theta})$:

$$p(\mathbf{t}_N|\boldsymbol{\theta}) = \int p(\mathbf{t}_N|\mathbf{a}_N)p(\mathbf{a}_N|\boldsymbol{\theta})d\mathbf{a}_N.$$

The approximation of log likelihood function is

$$\ln p(\mathbf{t}_N|\boldsymbol{\theta}) = \Psi(\mathbf{a}_N^*) - \frac{1}{2}\ln|\mathbf{W}_N + \mathbf{C}_N^{-1}| + \frac{N}{2}\ln(2\pi)$$

where $\Psi(\mathbf{a}_N^*) = \ln p(\mathbf{a}_N^*|\boldsymbol{\theta}) + \ln p(\mathbf{t}_N|\mathbf{a}_N^*)$.

Gaussian Process for Classification

The derivative w.r.t. θ_j is then

$$\frac{\partial \ln p(\mathbf{t}_N | \boldsymbol{\theta})}{\partial \theta_j} = \frac{1}{2} \mathbf{a}_N^{*-1} \mathbf{C}_N^{-1} \left(\frac{\partial \mathbf{C}_N}{\partial \theta_j} \right) \mathbf{C}_N^{-1} \mathbf{a}_N^* - \frac{1}{2} \text{Tr} \left\{ (\mathbf{I} + \mathbf{C}_N \mathbf{W}_N)^{-1} \mathbf{W}_N \left(\frac{\partial \mathbf{C}_N}{\partial \theta_j} \right) \right\}.$$

This leaves the following contribution to the derivative w.r.t. a component θ_j of $\boldsymbol{\theta}$

$$-\frac{1}{2} \sum_{n=1}^N \frac{\partial \ln |W_N + C_N^{-1}|}{\partial a_n^*} \left(\frac{\partial a_n^*}{\partial \theta_j} \right) = -\frac{1}{2} \sum_{n=1}^N [(I + C_N W_N)^{-1} C_N]_{nn} \sigma_n^* (1 - \sigma_n^*) (1 - 2\sigma_n^*) \left(\frac{\partial a_n^*}{\partial \theta_j} \right)$$

where $\sigma_n^* = \sigma(a_n^*)$.

The differentiation of $\mathbf{a}_N^* = \mathbf{C}_N(\mathbf{t}_N - \boldsymbol{\sigma}_N)$ w.r.t. θ_j gives

$$\frac{\partial a_n^*}{\partial \theta_j} = \frac{\partial \mathbf{C}_N}{\partial \theta_j} (\mathbf{t}_N - \boldsymbol{\sigma}_N) - \mathbf{C}_N \mathbf{W}_N \left(\frac{\partial a_n^*}{\partial \theta_j} \right).$$

Rearranging then gives

$$\frac{\partial a_n^*}{\partial \theta_j} = (\mathbf{I} + \mathbf{W}_N \mathbf{C}_N)^{-1} \left(\frac{\partial \mathbf{C}_N}{\partial \theta_j} \right) (\mathbf{t}_N - \boldsymbol{\sigma}_N).$$

The evaluation of the gradient of the log likelihood function can be used to determine the value for $\boldsymbol{\theta}$.

Gaussian Process for Classification

