

M5MS10

Machine Learning

Spring 2018

Lecture 6

Dr Ben Calderhead
b.calderhead@imperial.ac.uk



M5MS10

Machine Learning

Spring 2018

Lecture 6

Dr Ben Calderhead
b.calderhead@imperial.ac.uk

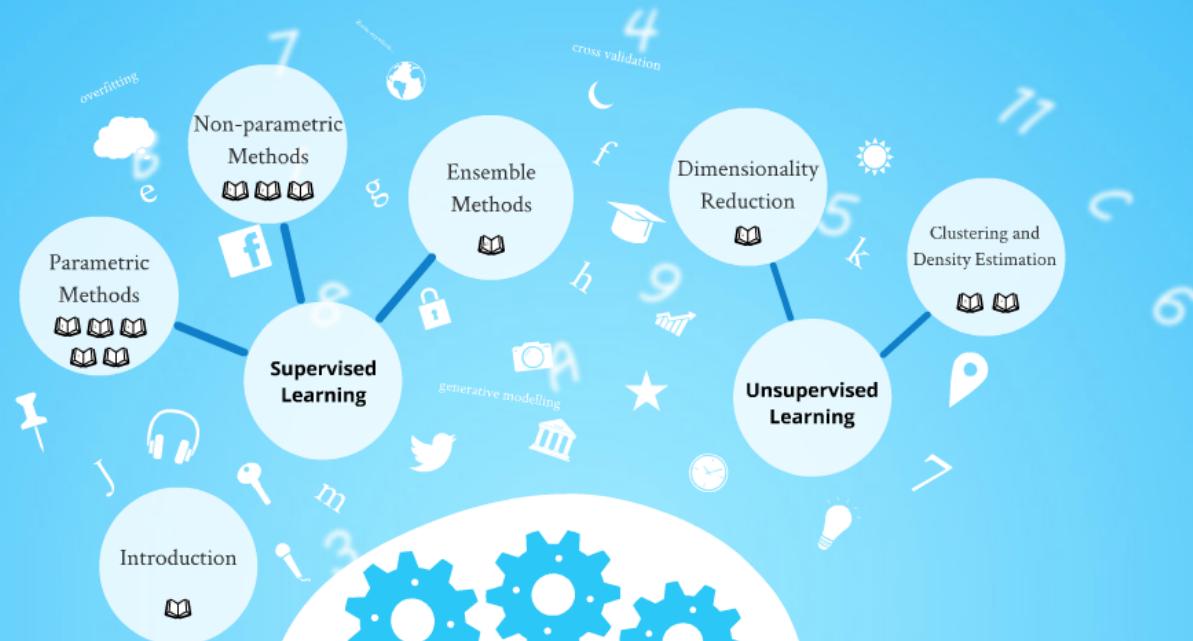


IS10 Learning

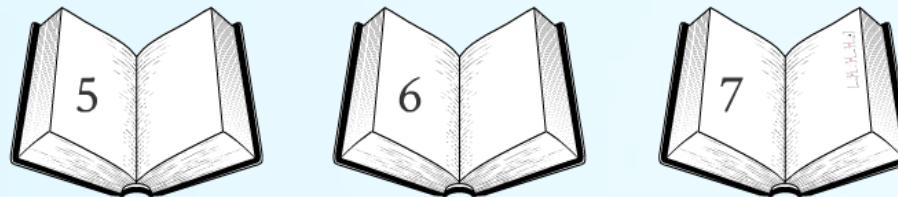
2018

re 6

Gelderhead
imperial.ac.uk

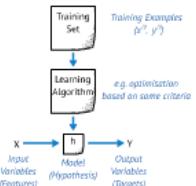


Non-parametric Methods



Introduction to Gaussian Processes

Nonparametric Regression



We'll consider continuous output variables, i.e. *regression*.

Gaussian Processes

A **Gaussian process** is an infinite collection of random variables, any *finite* number of which have a (consistent) multivariate Gaussian distribution.

We can think of this as a generalisation of a multivariate Gaussian to *infinite* dimensions.

A Gaussian process defines a distribution over functions.

- We can consider a function to approximately be an infinitely long vector.

Gaussian Process versus Gaussian Distribution

Gaussian distribution: $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$



Gaussian process: $f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$



Gaussian Process versus Gaussian Distribution

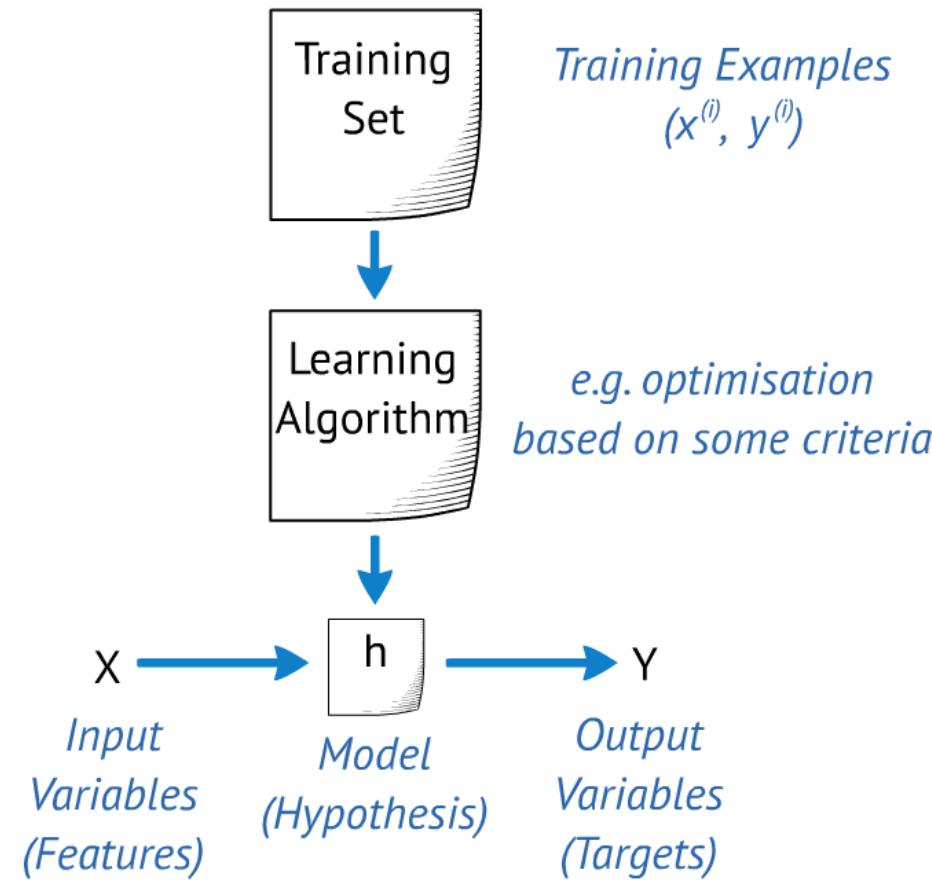
A Gaussian distribution is a distribution over a variable of a given dimension.

- It is fully specified by a mean vector and a covariance matrix:
 $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$.
- For each index i denoting each dimension of \mathbf{x} , the mean and variance is defined directly.

A GP is a distribution over functions and so the index set is of infinite dimension.

- It is fully specified by a mean function and a covariance function:
 $f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$.
- For each dimension of f , the mean and (co-)variance is specified via a function indexed by \mathbf{x} , where e.g. $\mathbf{x} \in \mathbb{R}$.

Nonparametric Regression



We'll consider continuous output variables, i.e. *regression*.

Gaussian Processes

A **Gaussian process** is an infinite collection of random variables, any *finite* number of which have a (consistent) multivariate Gaussian distribution.

We can think of this as a generalisation of a multivariate Gaussian to *infinite* dimensions.

A Gaussian process defines a distribution over functions.

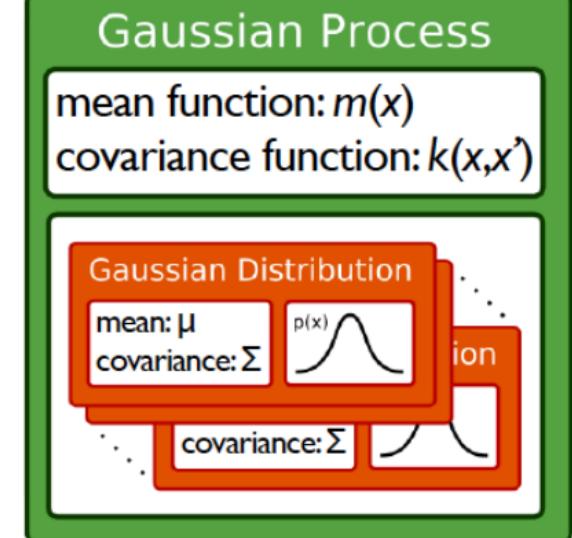
- ▶ We can consider a function to approximately be an infinitely long vector.

Gaussian Process versus Gaussian Distribution

Gaussian distribution: $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$



Gaussian process: $f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$



Gaussian Process versus Gaussian Distribution

A Gaussian distribution is a distribution over a variable of a given dimension.

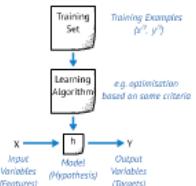
- ▶ It is fully specified by a mean vector and a covariance matrix:
 $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.
- ▶ For each index i denoting each dimension of \mathbf{x} , the mean and variance is defined directly.

A GP is a distribution over functions and so the index set is of infinite dimension.

- ▶ It is fully specified by a mean function and a covariance function: $f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$.
- ▶ For each dimension of f , the mean and (co-)variance is specified via a function indexed by \mathbf{x} , where e.g. $\mathbf{x} \in \mathbb{R}$.

Introduction to Gaussian Processes

Nonparametric Regression



We'll consider continuous output variables, i.e. *regression*.

Gaussian Processes

A **Gaussian process** is an infinite collection of random variables, any *finite* number of which have a (consistent) multivariate Gaussian distribution.

We can think of this as a generalisation of a multivariate Gaussian to *infinite* dimensions.

A Gaussian process defines a distribution over functions.

- We can consider a function to approximately be an infinitely long vector.

Gaussian Process versus Gaussian Distribution

Gaussian distribution: $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$



Gaussian process: $f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$



Gaussian Process versus Gaussian Distribution

A Gaussian distribution is a distribution over a variable of a given dimension.

- It is fully specified by a mean vector and a covariance matrix: $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$.
- For each index i denoting each dimension of \mathbf{x} , the mean and variance is defined directly.

A GP is a distribution over functions and so the index set is of infinite dimension.

- It is fully specified by a mean function and a covariance function: $f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$.
- For each dimension of f , the mean and (co-)variance is specified via a function indexed by \mathbf{x} , where e.g. $\mathbf{x} \in \mathbb{R}$.

Multivariate Gaussian Distributions

Multivariate Gaussian

Although the Gaussian process is an *infinite dimensional* object, it turns out that we only ever need to deal with finitely many points at any one time.

Multivariate Gaussian distributions are all we require!

A vector-valued random variable $\mathbf{x} \in \mathbb{R}^n$ is said to have a multivariate normal (or Gaussian) distribution with mean $\mu \in \mathbb{R}^n$ and covariance matrix $\Sigma \in \mathbb{S}_{++}^n$ if

$$p(\mathbf{x}|\mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma (\mathbf{x} - \mu)\right)$$

where \mathbb{S}_{++}^n denotes the space of $n \times n$ positive-definite symmetric matrices.

Multivariate Gaussian

Gaussian random variables are extremely useful in machine learning and statistics for the following main reasons.

- ▶ **Measurement error** can often be considered to be the accumulation of a large number of small independent random perturbations - by the *Central Limit Theorem*, summations of independent random variables will tend towards being Gaussian distributed.
- ▶ They are **convenient** because many of the integrals involving Gaussian distributions that arise in practice may be written as simple closed form solutions.

Multivariate Gaussian

The following properties of multivariate Gaussians are vital for manipulating Gaussian processes analytically.

Consider a random vector $\mathbf{x} \in \mathbb{R}^n$ with $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$. We can then consider the variable \mathbf{x} that has been partitioned into two sets $\mathbf{x}_A = [x_1, \dots, x_r]^T \in \mathbb{R}^r$ and $\mathbf{x}_B = [x_{r+1}, \dots, x_n]^T \in \mathbb{R}^{(n-r)}$. We therefore have

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_A \\ \mathbf{x}_B \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_A \\ \mu_B \end{pmatrix} \quad \Sigma = \begin{pmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{pmatrix}$$

We note that $\Sigma_{AB} = \Sigma_{BA}^T$, since $\Sigma = E[(\mathbf{x} - \mu)(\mathbf{x} - \mu)^T] = \Sigma^T$.

Multivariate Gaussian

The Gaussian density function is correctly **normalised**.

$$\int_{\mathbf{x}} p(\mathbf{x}|\mu, \Sigma) d\mathbf{x} = 1$$

Multivariate Gaussian

8.1.8 Product of gaussian densities

Let $\mathcal{N}_k(\mathbf{m}, \Sigma)$ denote a density of \mathbf{x} , then

$$\mathcal{N}_k(\mathbf{m}_1, \Sigma_1) \cdot \mathcal{N}_k(\mathbf{m}_2, \Sigma_2) = c_v \mathcal{N}_k(\mathbf{m}_v, \Sigma_v) \quad (371)$$

$$\begin{aligned} c_v &= \mathcal{N}_m(\mathbf{m}_2, (\Sigma_1 + \Sigma_2)) \\ &= \frac{1}{\sqrt{\det(2\pi(\Sigma_1 + \Sigma_2))}} \exp\left[-\frac{1}{2}(\mathbf{m}_1 - \mathbf{m}_2)^T (\Sigma_1 + \Sigma_2)^{-1}(\mathbf{m}_1 - \mathbf{m}_2)\right] \\ m_v &= (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1}(\Sigma_1^{-1}\mathbf{m}_1 + \Sigma_2^{-1}\mathbf{m}_2) \\ \Sigma_v &= (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1} \end{aligned}$$

but note that the product is not normalized as a density of \mathbf{x} .

[Matrix Cookbook...](#)

Multivariate Gaussian

The marginal densities

$$\begin{aligned} p(\mathbf{x}_A) &= \int_{\mathbf{x}_B} p(\mathbf{x}_A, \mathbf{x}_B | \mu, \Sigma) d\mathbf{x}_B \\ p(\mathbf{x}_B) &= \int_{\mathbf{x}_A} p(\mathbf{x}_A, \mathbf{x}_B | \mu, \Sigma) d\mathbf{x}_A \end{aligned}$$

are also Gaussian with the following form,

$$\begin{aligned} \mathbf{x}_A &\sim \mathcal{N}(\mu_A, \Sigma_{AA}) \\ \mathbf{x}_B &\sim \mathcal{N}(\mu_B, \Sigma_{BB}) \end{aligned}$$

Multivariate Gaussian

The **conditional densities**

$$\begin{aligned} p(\mathbf{x}_A | \mathbf{x}_B) &= \frac{p(\mathbf{x}_A, \mathbf{x}_B | \mu, \Sigma)}{\int_{\mathbf{x}_A} p(\mathbf{x}_A, \mathbf{x}_B | \mu, \Sigma) d\mathbf{x}_A} \\ p(\mathbf{x}_B | \mathbf{x}_A) &= \frac{p(\mathbf{x}_A, \mathbf{x}_B | \mu, \Sigma)}{\int_{\mathbf{x}_B} p(\mathbf{x}_A, \mathbf{x}_B | \mu, \Sigma) d\mathbf{x}_B} \end{aligned}$$

are also Gaussian with the following form,

$$\begin{aligned} \mathbf{x}_A | \mathbf{x}_B &\sim \mathcal{N}(\mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(\mathbf{x}_B - \mu_B), \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}) \\ \mathbf{x}_B | \mathbf{x}_A &\sim \mathcal{N}(\mu_B + \Sigma_{BA}\Sigma_{AA}^{-1}(\mathbf{x}_A - \mu_A), \Sigma_{BB} - \Sigma_{BA}\Sigma_{AA}^{-1}\Sigma_{AB}) \end{aligned}$$

Multivariate Gaussian

Although the Gaussian process is an *infinite dimensional* object, it turns out that we only ever need to deal with finitely many points at any one time.

Multivariate Gaussian distributions are all we require!

A vector-valued random variable $\mathbf{x} \in \mathbb{R}^n$ is said to have a multivariate normal (or Gaussian) distribution with mean $\boldsymbol{\mu} \in \mathbb{R}^n$ and covariance matrix $\boldsymbol{\Sigma} \in \mathbf{S}_{++}^n$ if

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{2/n} |\boldsymbol{\Sigma}|^{1/n}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma} (\mathbf{x} - \boldsymbol{\mu})\right)$$

where \mathbf{S}_{++}^n denotes the space of $n \times n$ positive-definite symmetric matrices.

Multivariate Gaussian

Gaussian random variables are extremely useful in machine learning and statistics for the following main reasons.

- ▶ **Measurement error** can often be considered to be the accumulation of a large number of small independent random perturbations - by the *Central Limit Theorem*, summations of independent random variables will tend towards being Gaussian distributed.
- ▶ They are **convenient** because many of the integrals involving Gaussian distributions that arise in practice may be written as simple closed form solutions.

Multivariate Gaussian

The following properties of multivariate Gaussians are vital for manipulating Gaussian processes analytically.

Consider a random vector $\mathbf{x} \in \mathbb{R}^n$ with $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. We can then consider the variable \mathbf{x} that has been partitioned into two sets $\mathbf{x}_A = [x_1, \dots, x_r]^T \in \mathbb{R}^r$ and $\mathbf{x}_B = [x_{r+1}, \dots, x_n]^T \in \mathbb{R}^{(n-r)}$. We therefore have

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_A \\ \mathbf{x}_B \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_A \\ \boldsymbol{\mu}_B \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{AA} & \boldsymbol{\Sigma}_{AB} \\ \boldsymbol{\Sigma}_{BA} & \boldsymbol{\Sigma}_{BB} \end{pmatrix}$$

We note that $\boldsymbol{\Sigma}_{AB} = \boldsymbol{\Sigma}_{BA}^T$, since $\boldsymbol{\Sigma} = E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] = \boldsymbol{\Sigma}^T$.

Multivariate Gaussian

The Gaussian density function is correctly **normalised**.

$$\int_{\mathbf{x}} p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x} = 1$$

Multivariate Gaussian

8.1.8 Product of gaussian densities

Let $\mathcal{N}_x(\mathbf{m}, \Sigma)$ denote a density of x , then

$$\mathcal{N}_x(\mathbf{m}_1, \Sigma_1) \cdot \mathcal{N}_x(\mathbf{m}_2, \Sigma_2) = c_c \mathcal{N}_x(\mathbf{m}_c, \Sigma_c) \quad (371)$$

$$c_c = \mathcal{N}_{\mathbf{m}_1}(\mathbf{m}_2, (\Sigma_1 + \Sigma_2))$$

$$= \frac{1}{\sqrt{\det(2\pi(\Sigma_1 + \Sigma_2))}} \exp \left[-\frac{1}{2}(\mathbf{m}_1 - \mathbf{m}_2)^T (\Sigma_1 + \Sigma_2)^{-1} (\mathbf{m}_1 - \mathbf{m}_2) \right]$$

$$\mathbf{m}_c = (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1} (\Sigma_1^{-1} \mathbf{m}_1 + \Sigma_2^{-1} \mathbf{m}_2)$$

$$\Sigma_c = (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1}$$

but note that the product is not normalized as a density of x .

Matrix Cookbook...

Multivariate Gaussian

The marginal densities

$$p(\mathbf{x}_A) = \int_{\mathbf{x}_B} p(\mathbf{x}_A, \mathbf{x}_B | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}_B$$

$$p(\mathbf{x}_B) = \int_{\mathbf{x}_A} p(\mathbf{x}_A, \mathbf{x}_B | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}_A$$

are also Gaussian with the following form,

$$\mathbf{x}_A \sim \mathcal{N}(\boldsymbol{\mu}_A, \boldsymbol{\Sigma}_{AA})$$

$$\mathbf{x}_B \sim \mathcal{N}(\boldsymbol{\mu}_B, \boldsymbol{\Sigma}_{BB})$$

Multivariate Gaussian

The conditional densities

$$p(\mathbf{x}_A | \mathbf{x}_B) = \frac{p(\mathbf{x}_A, \mathbf{x}_B | \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\int_{\mathbf{x}_A} p(\mathbf{x}_A, \mathbf{x}_B | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}_A}$$

$$p(\mathbf{x}_B | \mathbf{x}_A) = \frac{p(\mathbf{x}_A, \mathbf{x}_B | \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\int_{\mathbf{x}_B} p(\mathbf{x}_A, \mathbf{x}_B | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}_B}$$

are also Gaussian with the following form,

$$\begin{aligned}\mathbf{x}_A | \mathbf{x}_B &\sim \mathcal{N}(\boldsymbol{\mu}_A + \boldsymbol{\Sigma}_{AB} \boldsymbol{\Sigma}_{BB}^{-1} (\mathbf{x}_B - \boldsymbol{\mu}_B), \boldsymbol{\Sigma}_{AA} - \boldsymbol{\Sigma}_{AB} \boldsymbol{\Sigma}_{BB}^{-1} \boldsymbol{\Sigma}_{BA}) \\ \mathbf{x}_B | \mathbf{x}_A &\sim \mathcal{N}(\boldsymbol{\mu}_B + \boldsymbol{\Sigma}_{BA} \boldsymbol{\Sigma}_{AA}^{-1} (\mathbf{x}_A - \boldsymbol{\mu}_A), \boldsymbol{\Sigma}_{BB} - \boldsymbol{\Sigma}_{BA} \boldsymbol{\Sigma}_{AA}^{-1} \boldsymbol{\Sigma}_{AB})\end{aligned}$$

Multivariate Gaussian Distributions

Multivariate Gaussian

Although the Gaussian process is an *infinite dimensional* object, it turns out that we only ever need to deal with finitely many points at any one time.

Multivariate Gaussian distributions are all we require!

A vector-valued random variable $\mathbf{x} \in \mathbb{R}^n$ is said to have a multivariate normal (or Gaussian) distribution with mean $\mu \in \mathbb{R}^n$ and covariance matrix $\Sigma \in \mathbb{S}_{++}^n$ if

$$p(\mathbf{x}|\mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma (\mathbf{x} - \mu)\right)$$

where \mathbb{S}_{++}^n denotes the space of $n \times n$ positive-definite symmetric matrices.

Multivariate Gaussian

Gaussian random variables are extremely useful in machine learning and statistics for the following main reasons.

- ▶ **Measurement error** can often be considered to be the accumulation of a large number of small independent random perturbations - by the *Central Limit Theorem*, summations of independent random variables will tend towards being Gaussian distributed.
- ▶ They are **convenient** because many of the integrals involving Gaussian distributions that arise in practice may be written as simple closed form solutions.

Multivariate Gaussian

The following properties of multivariate Gaussians are vital for manipulating Gaussian processes analytically.

Consider a random vector $\mathbf{x} \in \mathbb{R}^n$ with $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$. We can then consider the variable \mathbf{x} that has been partitioned into two sets $\mathbf{x}_A = [x_1, \dots, x_r]^T \in \mathbb{R}^r$ and $\mathbf{x}_B = [x_{r+1}, \dots, x_n]^T \in \mathbb{R}^{(n-r)}$. We therefore have

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_A \\ \mathbf{x}_B \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_A \\ \mu_B \end{pmatrix} \quad \Sigma = \begin{pmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{pmatrix}$$

We note that $\Sigma_{AB} = \Sigma_{BA}^T$, since $\Sigma = E[(\mathbf{x} - \mu)(\mathbf{x} - \mu)^T] = \Sigma^T$.

Multivariate Gaussian

The Gaussian density function is correctly **normalised**.

$$\int_{\mathbf{x}} p(\mathbf{x}|\mu, \Sigma) d\mathbf{x} = 1$$

Multivariate Gaussian

8.1.8 Product of gaussian densities

Let $\mathcal{N}_k(\mathbf{m}, \Sigma)$ denote a density of \mathbf{x} , then

$$\mathcal{N}_k(\mathbf{m}_1, \Sigma_1) \cdot \mathcal{N}_k(\mathbf{m}_2, \Sigma_2) = c_v \mathcal{N}_k(\mathbf{m}_v, \Sigma_v) \quad (371)$$

$$\begin{aligned} c_v &= \mathcal{N}_m(\mathbf{m}_2, (\Sigma_1 + \Sigma_2)) \\ &= \frac{1}{\sqrt{\det(2\pi(\Sigma_1 + \Sigma_2))}} \exp\left[-\frac{1}{2}(\mathbf{m}_1 - \mathbf{m}_2)^T (\Sigma_1 + \Sigma_2)^{-1}(\mathbf{m}_1 - \mathbf{m}_2)\right] \\ m_v &= (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1}(\Sigma_1^{-1}\mathbf{m}_1 + \Sigma_2^{-1}\mathbf{m}_2) \\ \Sigma_v &= (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1} \end{aligned}$$

but note that the product is not normalized as a density of \mathbf{x} .

[Matrix Cookbook...](#)

Multivariate Gaussian

The **marginal densities**

$$\begin{aligned} p(\mathbf{x}_A) &= \int_{\mathbf{x}_B} p(\mathbf{x}_A, \mathbf{x}_B | \mu, \Sigma) d\mathbf{x}_B \\ p(\mathbf{x}_B) &= \int_{\mathbf{x}_A} p(\mathbf{x}_A, \mathbf{x}_B | \mu, \Sigma) d\mathbf{x}_A \end{aligned}$$

are also Gaussian with the following form,

$$\begin{aligned} \mathbf{x}_A &\sim \mathcal{N}(\mu_A, \Sigma_{AA}) \\ \mathbf{x}_B &\sim \mathcal{N}(\mu_B, \Sigma_{BB}) \end{aligned}$$

Multivariate Gaussian

The **conditional densities**

$$\begin{aligned} p(\mathbf{x}_A | \mathbf{x}_B) &= \frac{p(\mathbf{x}_A, \mathbf{x}_B | \mu, \Sigma)}{\int_{\mathbf{x}_A} p(\mathbf{x}_A, \mathbf{x}_B | \mu, \Sigma) d\mathbf{x}_A} \\ p(\mathbf{x}_B | \mathbf{x}_A) &= \frac{p(\mathbf{x}_A, \mathbf{x}_B | \mu, \Sigma)}{\int_{\mathbf{x}_B} p(\mathbf{x}_A, \mathbf{x}_B | \mu, \Sigma) d\mathbf{x}_B} \end{aligned}$$

are also Gaussian with the following form,

$$\begin{aligned} \mathbf{x}_A | \mathbf{x}_B &\sim \mathcal{N}(\mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(\mathbf{x}_B - \mu_B), \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}) \\ \mathbf{x}_B | \mathbf{x}_A &\sim \mathcal{N}(\mu_B + \Sigma_{BA}\Sigma_{AA}^{-1}(\mathbf{x}_A - \mu_A), \Sigma_{BB} - \Sigma_{BA}\Sigma_{AA}^{-1}\Sigma_{AB}) \end{aligned}$$

Gaussian Process Priors

Gaussian Process Prior

We have quite a lot of freedom when choosing the [mean](#) and [covariance](#) functions that will define our prior over functions.

Usually the mean is set to be the zero function.

The covariance function is much more interesting however, since it defines how the points in our function covary with one another.

The only requirement of a covariance function is that it is [positive semi-definite](#).

Covariance Functions

If k , k_1 and k_2 are positive semi-definite, then we can create new [valid covariance functions](#) according the to following functions, which are also positive semi-definite.

1. $ak(\mathbf{x}, \mathbf{x}')/R$, for $a \geq 0$
2. $k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$
3. $k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}')$
4. $P(k(\mathbf{x}, \mathbf{x}'))$, where $P(\cdot)$ is some polynomial function with positive coefficients
5. $\exp(k(\mathbf{x}, \mathbf{x}'))$
6. $f(\mathbf{x})k(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$
7. $k(\phi(\mathbf{x}), \phi(\mathbf{x}'))$

Squared Exponential Covariance Function

The squared exponential covariance function is a common choice,

$$k(x, x') = \sigma_0 \exp \left\{ -\frac{1}{2} \left(\frac{x - x'}{\lambda} \right)^2 \right\}$$

Intuitively, we see that function variables close to each other in the input space are highly correlated, whereas those far from each other are uncorrelated.

There are hyperparameters that change the output of the covariance function: σ controls the amplitude, λ controls the lengthscale.

The squared exponential produces very smooth functions, since it is [infinitely differentiable](#).

Matern Covariance Function

The Matern exponential covariance function is another common choice,

$$k(x, x') = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}|x - x'|}{\lambda} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu}|x - x'|}{\lambda} \right)$$

where K_ν is a modified Bessel function.

There are hyperparameters that change the output of the covariance function: ν controls the roughness of the sample functions, λ controls the lengthscale.

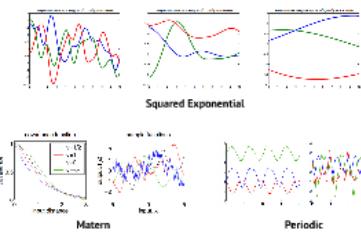
With $\nu = \infty$ we recover the squared exponential covariance, with $\nu = \frac{1}{2}$ we obtain an Ornstein-Uhlenbeck process.

Periodic Covariance Function

We can even model periodic functions,

$$k(x, x') = \exp \left(\frac{-2 \sin^2(x - x')}{\lambda^2} \right)$$

Covariance Function Comparison



Gaussian Process Prior

We have quite a lot of freedom when choosing the **mean** and **covariance** functions that will define our prior over functions.

Usually the mean is set to be the zero function.

The covariance function is much more interesting however, since it defines how the points in our function covary with one another.

The only requirement of a covariance function is that it is **positive semi-definite**.

Covariance Functions

If k , k_1 and k_2 are positive semi-definite, then we can create new **valid covariance functions** according the to following functions, which are also positive semi-definite.

1. $ak(\mathbf{x}, \mathbf{x}')R$, for $a \geq 0$
2. $k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$
3. $k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}')$
4. $P(k(\mathbf{x}, \mathbf{x}'))$, where $P()$ is some polynomial function with positive coefficients
5. $\exp(k(\mathbf{x}, \mathbf{x}'))$
6. $f(\mathbf{x})k(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$
7. $k(\phi(\mathbf{x}), \phi(\mathbf{x}'))$

Squared Exponential Covariance Function

The squared exponential covariance function is a common choice,

$$k(x, x') = \sigma_0 \exp \left\{ -\frac{1}{2} \left(\frac{x - x'}{\lambda} \right)^2 \right\}$$

Intuitively, we see that function variables close to each other in the input space are highly correlated, whereas those far from each other are uncorrelated.

There are hyperparameters that change the output of the covariance function: σ controls the amplitude, λ controls the lengthscale.

The squared exponential produces very smooth functions, since it is *infinitely differentiable*.

Matern Covariance Function

The Matern exponential covariance function is another common choice,

$$k(x, x') = \frac{2^{1-v}}{\Gamma(v)} \left(\frac{\sqrt{2v}|x - x'|}{\lambda} \right)^v K_v \left(\frac{\sqrt{2v}|x - x'|}{\lambda} \right)$$

where K_v is a modified Bessel function.

There are hyperparameters that change the output of the covariance function: v controls the roughness of the sample functions, λ controls the lengthscales.

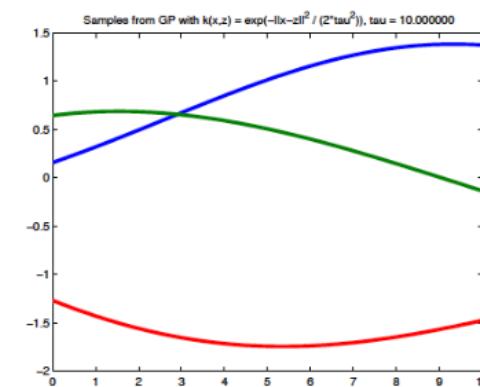
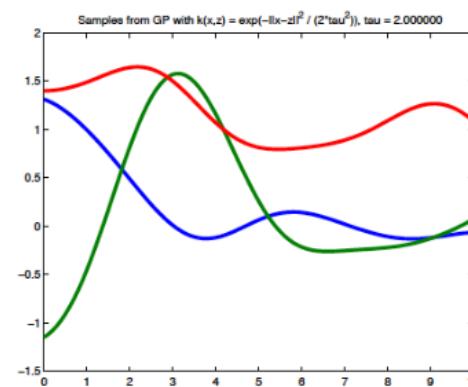
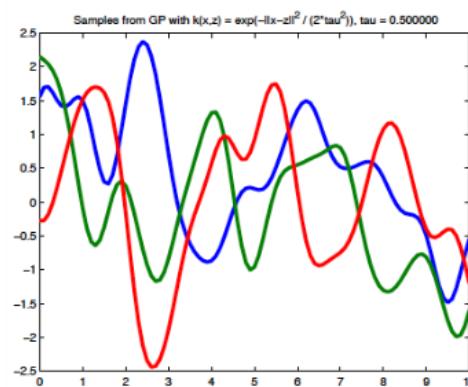
With $v = \infty$ we recover the squared exponential covariance, with $v = \frac{1}{2}$ we obtain an Ornstein-Uhlenbeck process.

Periodic Covariance Function

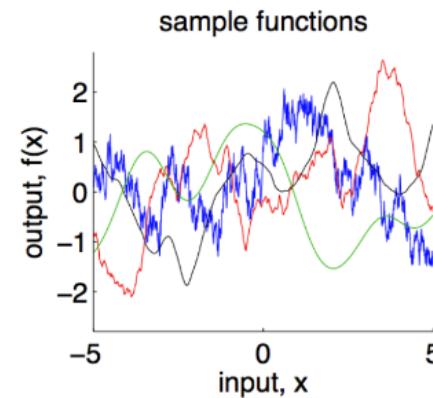
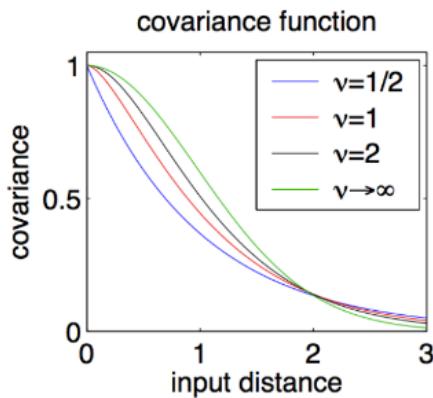
We can even model periodic functions,

$$k(x, x') = \exp\left(\frac{-2 \sin^2(x - x')}{\lambda^2}\right)$$

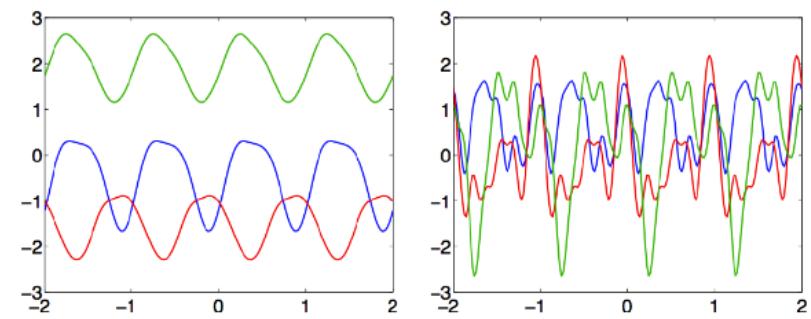
Covariance Function Comparison



Squared Exponential



Matern



Periodic

Gaussian Process Priors

Gaussian Process Prior

We have quite a lot of freedom when choosing the [mean](#) and [covariance](#) functions that will define our prior over functions.

Usually the mean is set to be the zero function.

The covariance function is much more interesting however, since it defines how the points in our function covary with one another.

The only requirement of a covariance function is that it is [positive semi-definite](#).

Covariance Functions

If k , k_1 and k_2 are positive semi-definite, then we can create new [valid covariance functions](#) according the to following functions, which are also positive semi-definite.

1. $ak(\mathbf{x}, \mathbf{x}')/R$, for $a \geq 0$
2. $k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$
3. $k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}')$
4. $P(k(\mathbf{x}, \mathbf{x}'))$, where $P(\cdot)$ is some polynomial function with positive coefficients
5. $\exp(k(\mathbf{x}, \mathbf{x}'))$
6. $f(\mathbf{x})k(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$
7. $k(\phi(\mathbf{x}), \phi(\mathbf{x}'))$

Squared Exponential Covariance Function

The squared exponential covariance function is a common choice,

$$k(x, x') = \sigma_0 \exp \left\{ -\frac{1}{2} \left(\frac{x - x'}{\lambda} \right)^2 \right\}$$

Intuitively, we see that function variables close to each other in the input space are highly correlated, whereas those far from each other are uncorrelated.

There are hyperparameters that change the output of the covariance function: σ controls the amplitude, λ controls the lengthscale.

The squared exponential produces very smooth functions, since it is [infinitely differentiable](#).

Matern Covariance Function

The Matern exponential covariance function is another common choice,

$$k(x, x') = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}|x - x'|}{\lambda} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu}|x - x'|}{\lambda} \right)$$

where K_ν is a modified Bessel function.

There are hyperparameters that change the output of the covariance function: ν controls the roughness of the sample functions, λ controls the lengthscale.

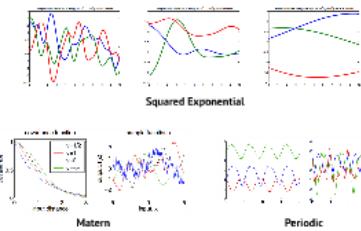
With $\nu = \infty$ we recover the squared exponential covariance, with $\nu = \frac{1}{2}$ we obtain an Ornstein-Uhlenbeck process.

Periodic Covariance Function

We can even model periodic functions,

$$k(x, x') = \exp \left(\frac{-2 \sin^2(x - x')}{\lambda^2} \right)$$

Covariance Function Comparison



Making Predictions with Gaussian Processes

Gaussian Process Regression

We can therefore model each datapoint as some function of the inputs plus noise, such that

$$y_i = f(x_i) + \epsilon_i$$

where we can define our prior over functions as a Gaussian process,

$$f(\mathbf{x}) \sim \mathcal{GP}(m = \mathbf{0}, k)$$

and we can assume Gaussian noise for the likelihood,

$$\rho(\mathbf{y}|\mathbf{f}) \sim \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I})$$

Conveniently, because everything is Gaussian, we can integrate over all the functions to obtain the marginal likelihood,

$$\rho(\mathbf{y}) = \int \rho(\mathbf{y}|\mathbf{f})\rho(\mathbf{f})d\mathbf{f} = \mathcal{N}(\mathbf{0}, \mathbf{K} + \sigma^2 \mathbf{I})$$

where \mathbf{K} is the matrix generated by the covariance function at the given inputs.

Gaussian Process Regression

Let's assume we now have a training set with N measurements \mathbf{y}_N observed with input variables \mathbf{x}_N . We then wish to make T predictions of \mathbf{y}_T for some other input variable \mathbf{x}_T .

We employ the notation \mathbf{K}_N to denote the matrix obtained by evaluating the covariance function at the N inputs \mathbf{x}_N , and \mathbf{K}_{NT} denotes the cross covariance between the inputs \mathbf{x}_N and \mathbf{x}_T .

We can consider the marginal likelihood at the joint training and test points,

$$\rho(\mathbf{y}, \mathbf{y}_T) = \mathcal{N}(\mathbf{0}, \mathbf{K}_{N+T} + \sigma^2 \mathbf{I})$$

where

$$\mathbf{K}_{N+T} = \begin{pmatrix} \mathbf{K}_N & \mathbf{K}_{NT} \\ \mathbf{K}_{TN} & \mathbf{K}_T \end{pmatrix}$$

Gaussian Process Regression

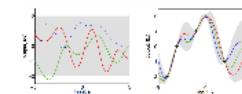
Finally, given this partitioned Gaussian distribution, we can condition the test variables on the training variables, such that

$$\rho(\mathbf{y}_T | \mathbf{y}) = \mathcal{N}(\mu_T, \Sigma_T)$$

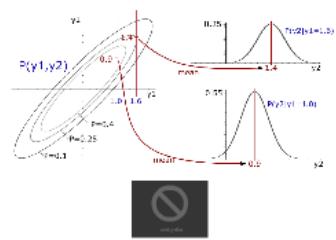
where

$$\mu_T = \mathbf{K}_{TN} (\mathbf{K}_N + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\Sigma_T = \mathbf{K}_T - \mathbf{K}_{TN} (\mathbf{K}_N + \sigma^2 \mathbf{I})^{-1} \mathbf{K}_{NT}$$



Gaussian Process Regression



Choosing Hyperparameters

Other non-Bayesian methods (such as Support Vector Machines) require the use of cross validation to choose the appropriate parameters for obtaining good predictive performance.

The advantage of Gaussian process approach is that we can choose the hyperparameters and covariances directly from the data!

In particular we can minimise the negative log marginal likelihood with respect to all hyperparameters θ ,

$$\mathcal{L} = -\log \rho(\mathbf{y}|\theta) = -\frac{1}{2} \log \det \mathbf{C}(\theta) + \frac{1}{2} \mathbf{y}^T \mathbf{C}(\theta) \mathbf{y} + \frac{N}{2} \log(2\pi)$$

where $\mathbf{C}(\theta) = \mathbf{K} + \sigma^2 \mathbf{I}$.

Relation to Linear Regression

Let's consider modelling functions defined by a weighted sum of a finite set of basis function, such that

$$f(x) = \sum_{m=1}^M \theta_m \phi_m(x) = \theta^T \phi(x)$$

We can place a Gaussian prior on the weights, so that $\rho(\theta) = \mathcal{N}(\mathbf{0}, \Sigma_\theta)$.

This defines a Gaussian process with a covariance function given by,

$$K(x, x') = E[f(x), f(x')] = \phi(x)^T \Sigma_\theta \phi(x')$$

Relation to Linear Regression

There is in fact a direct correspondence between covariance functions and basis functions given by [Mercer's theorem](#), which states that we can always decompose a covariance function into eigenfunctions and eigenvectors,

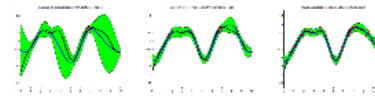
$$K(x, x') = \sum_{i=1}^{\infty} \lambda_i \varphi_i(x) \varphi_i(x')$$

If the sum is finite then it can be defined analogously using linear regression.

Often the sum is infinite, and there are no closed form solutions for the eigenvectors and eigenvalues.

We see how this is another form of the *kernel trick* - the covariance function effectively transforms the input space into a high or even infinite dimensional feature space, yet we can still compute in the lower dimensional space.

Example: Increasing Data



Summary

Gaussian processes are a very powerful approach for performing nonparametric supervised learning.

We can use the kernel trick to allow us to do linear regression effectively using an infinite number of basis functions, while retaining analytic tractability.

The main challenge is the computational scaling; it is cubic in the number of observations.

Gaussian Process Regression

We can therefore model each datapoint as some function of the inputs plus noise, such that

$$y_i = f(x_i) + \epsilon_i$$

where we can define our prior over functions as a Gaussian process,

$$\mathbf{f}(\mathbf{x}) \sim \mathcal{GP}(m = \mathbf{0}, k)$$

and we can assume Gaussian noise for the likelihood,

$$p(\mathbf{y}|\mathbf{f}) \sim \mathcal{N}(\mathbf{f}, \sigma^2)$$

Conveniently, because everything is Gaussian, we can integrate over all the functions to obtain the marginal likelihood,

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f} = \mathcal{N}(\mathbf{0}, \mathbf{K} + \sigma^2 \mathbf{I})$$

where \mathbf{K} is the matrix generated by the covariance function at the given inputs.

Gaussian Process Regression

Let's assume we now have a training set with N measurements \mathbf{y}_N observed with input variables \mathbf{x}_N . We then wish to make T predictions of \mathbf{y}_T for some other input variable \mathbf{x}_T .

We employ the notation \mathbf{K}_N to denote the matrix obtained by evaluating the covariance function at the N inputs \mathbf{x}_N , and \mathbf{K}_{NT} denotes the cross covariance between the inputs \mathbf{x}_N and \mathbf{x}_T .

We can consider the marginal likelihood at the joint training and test points,

$$p(\mathbf{y}, \mathbf{y}_T) = \mathcal{N}(\mathbf{0}, \mathbf{K}_{N+T} + \sigma^2 \mathbf{I})$$

where

$$\mathbf{K}_{N+T} = \begin{pmatrix} \mathbf{K}_N & \mathbf{K}_{NT} \\ \mathbf{K}_{TN} & \mathbf{K}_T \end{pmatrix}$$

Gaussian Process Regression

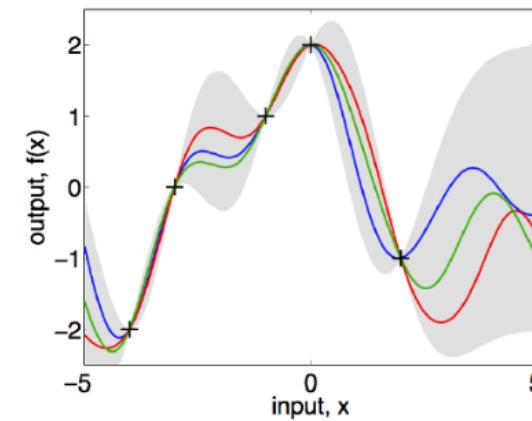
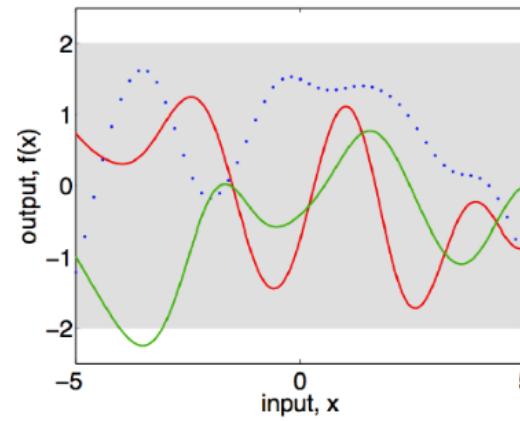
Finally, given this partitioned Gaussian distribution, we can condition the test variables on the training variables, such that

$$p(\mathbf{y}_T | \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_T, \boldsymbol{\Sigma}_T)$$

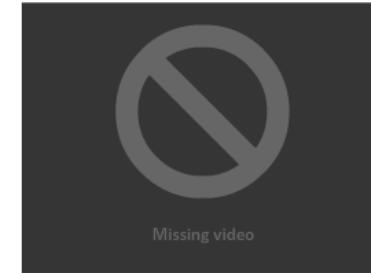
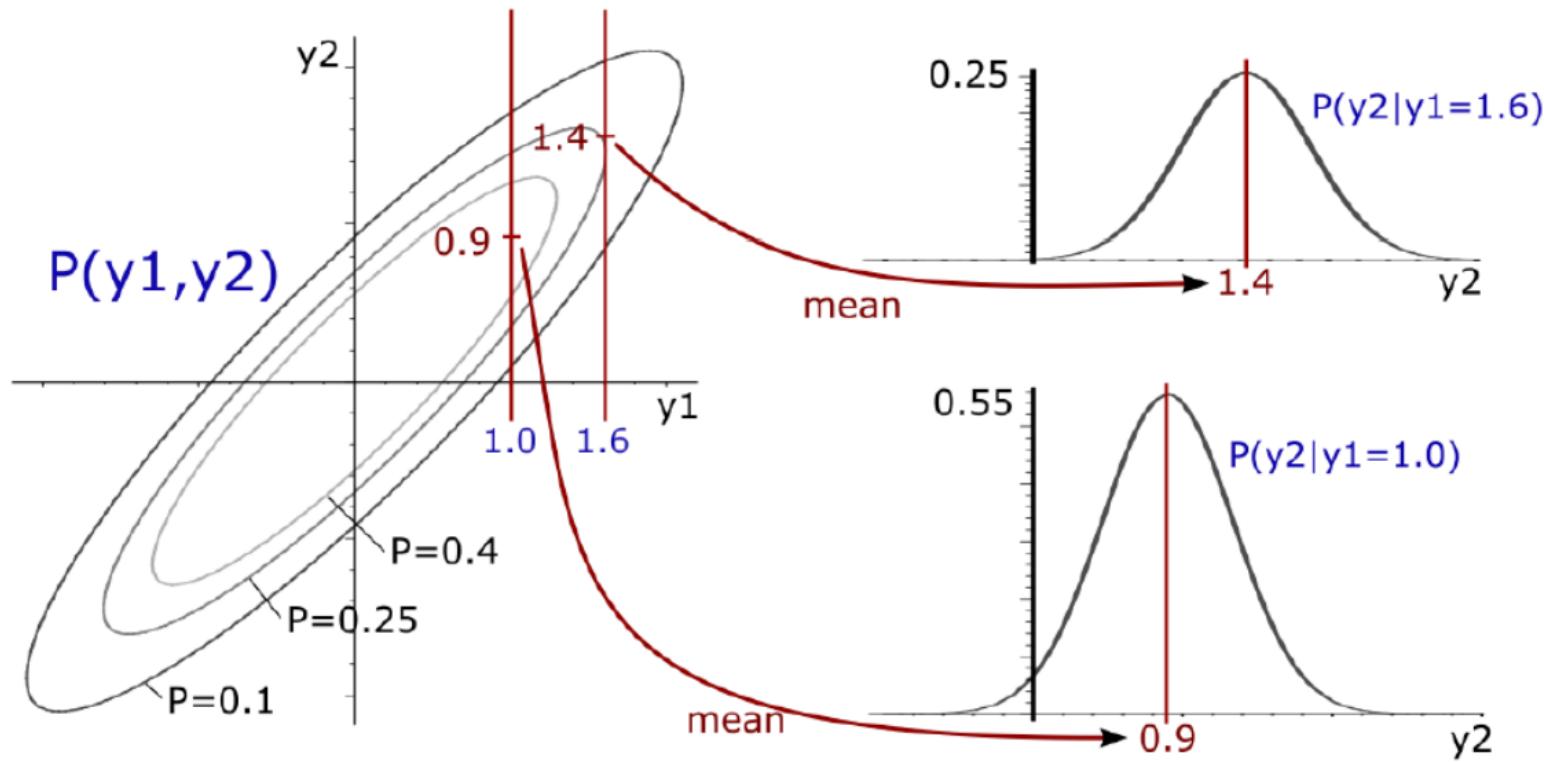
where

$$\boldsymbol{\mu}_T = \mathbf{K}_{TN} (\mathbf{K}_N + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\boldsymbol{\Sigma}_T = \mathbf{K}_T - \mathbf{K}_{TN} (\mathbf{K}_N + \sigma^2 \mathbf{I})^{-1} \mathbf{K}_{NT}$$



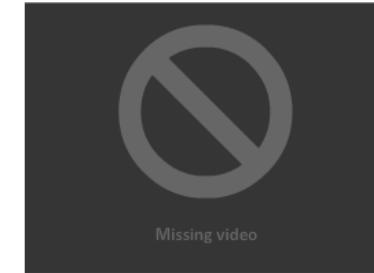
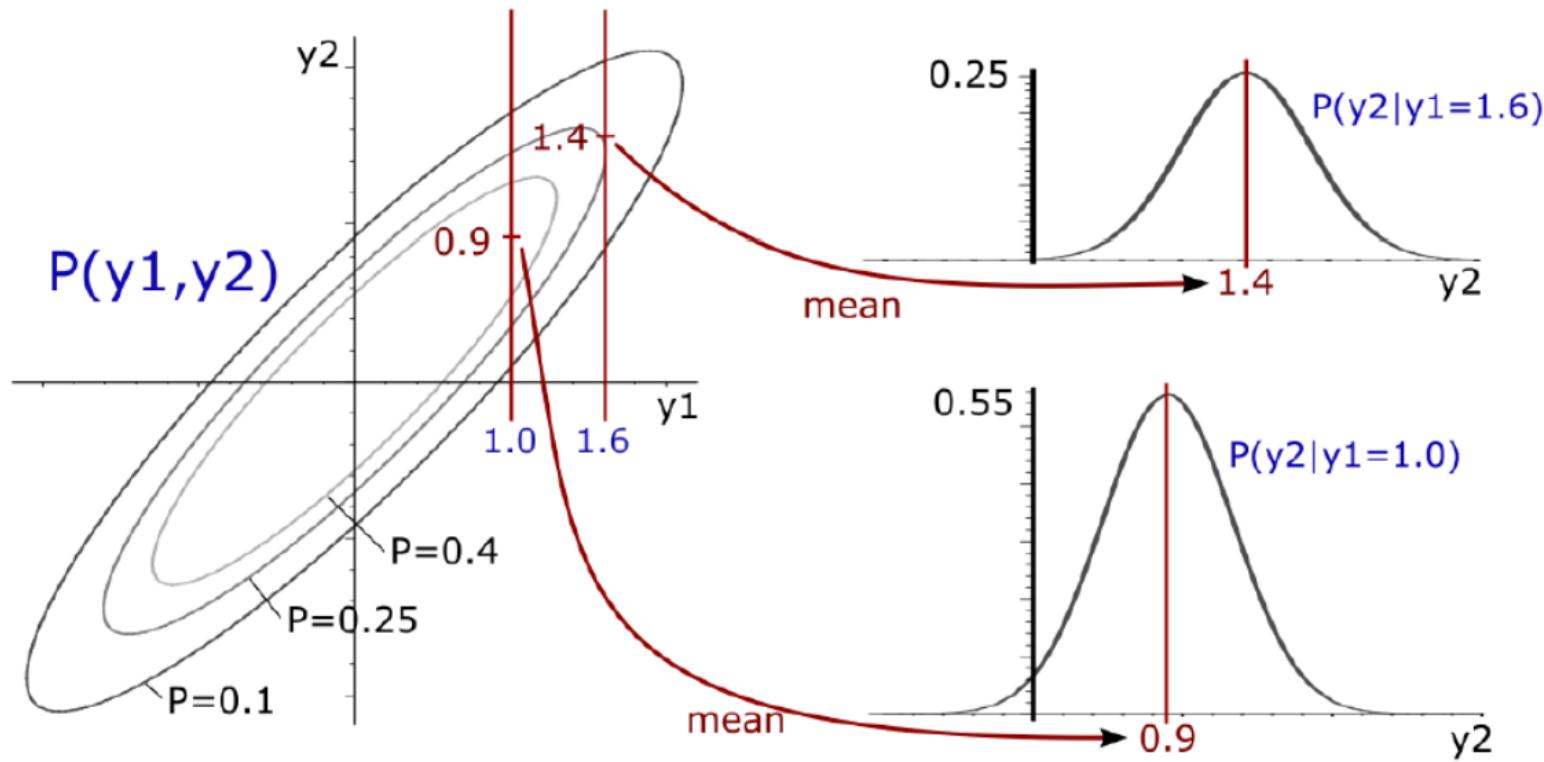
Gaussian Process Regression





Missing video

Gaussian Process Regression



Choosing Hyperparameters

Other non-Bayesian methods (such as Support Vector Machines) require the use of cross validation to choose the appropriate parameters for obtaining good predictive performance.

The advantage of Gaussian process approach is that we can choose the hyperparameters and covariances directly from the data!

In particular we can minimise the negative log marginal likelihood with respect to all hyperparameters θ ,

$$\mathcal{L} = -\log p(\mathbf{y}|\theta) = \frac{1}{2} \log \det \mathbf{C}(\theta) + \frac{1}{2} \mathbf{y}^T \mathbf{C}(\theta) \mathbf{y} + \frac{N}{2} \log(2\pi)$$

where $\mathbf{C}(\theta) = \mathbf{K} + \sigma^2 \mathbf{I}$.

Relation to Linear Regression

Let's consider modelling functions defined by a weighted sum of a finite set of basis function, such that

$$f(x) = \sum_{m=1}^M \theta_m \phi_m(x) = \boldsymbol{\theta}^T \boldsymbol{\phi}(x)$$

We can place a Gaussian prior on the weights, so that
 $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}})$.

This defines a *Gaussian process* with a covariance function given by,

$$K(x, x') = E[f(x), f(x')] = \boldsymbol{\phi}(x)^T \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \boldsymbol{\phi}(x')$$

Relation to Linear Regression

There is in fact a direct correspondence between covariance functions and basis functions given by **Mercer's theorem**, which states that we can always decompose a covariance function into eigenfunctions and eigenvectors,

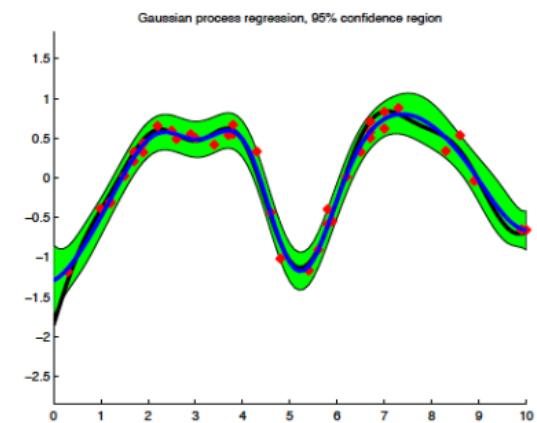
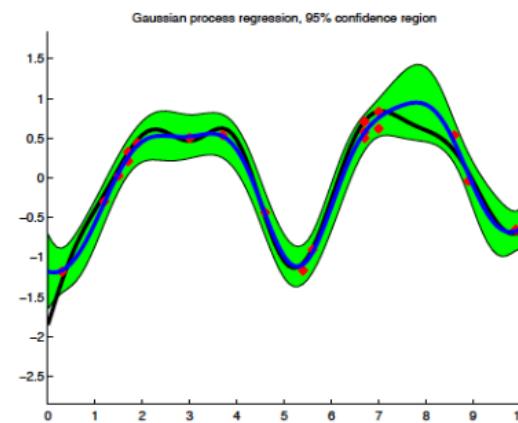
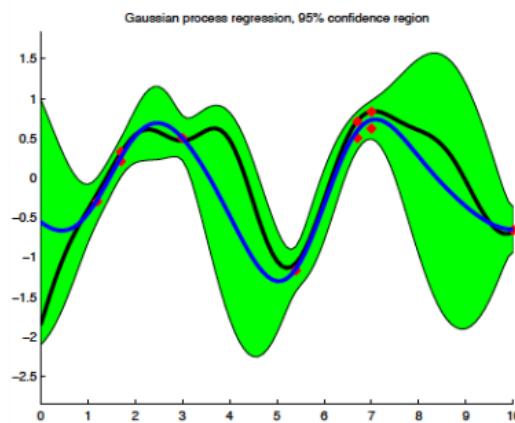
$$K(x, x') = \sum_{i=1}^{\infty} \lambda_i \varphi_i(x) \varphi_i(x')$$

If the sum is finite then it can be defined analogously using linear regression.

Often the sum is infinite, and there are no closed form solutions for the eigenvectors and eigenvalues.

We see how this is another form of the *kernel trick* - the covariance function effectively transforms the input space into a high or even infinite dimensional feature space, yet we can still compute in the lower dimensional space.

Example: Increasing Data



Summary

Gaussian processes are a very powerful approach for performing nonparametric supervised learning.

We can use the kernel trick to allow us to do linear regression effectively using an infinite number of basis functions, while retaining analytic tractability.

The main challenge is the computational scaling; it is cubic in the number of observations.

Making Predictions with Gaussian Processes

Gaussian Process Regression

We can therefore model each datapoint as some function of the inputs plus noise, such that

$$y_i = f(x_i) + \epsilon_i$$

where we can define our prior over functions as a Gaussian process,

$$f(\mathbf{x}) \sim \mathcal{GP}(m = \mathbf{0}, k)$$

and we can assume Gaussian noise for the likelihood,

$$\rho(\mathbf{y}|\mathbf{f}) \sim \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I})$$

Conveniently, because everything is Gaussian, we can integrate over all the functions to obtain the marginal likelihood,

$$\rho(\mathbf{y}) = \int \rho(\mathbf{y}|\mathbf{f})\rho(\mathbf{f})d\mathbf{f} = \mathcal{N}(\mathbf{0}, \mathbf{K} + \sigma^2 \mathbf{I})$$

where \mathbf{K} is the matrix generated by the covariance function at the given inputs.

Gaussian Process Regression

Let's assume we now have a training set with N measurements \mathbf{y}_N observed with input variables \mathbf{x}_N . We then wish to make T predictions of \mathbf{y}_T for some other input variable \mathbf{x}_T .

We employ the notation \mathbf{K}_N to denote the matrix obtained by evaluating the covariance function at the N inputs \mathbf{x}_N , and \mathbf{K}_{NT} denotes the cross covariance between the inputs \mathbf{x}_N and \mathbf{x}_T .

We can consider the marginal likelihood at the joint training and test points,

$$\rho(\mathbf{y}, \mathbf{y}_T) = \mathcal{N}(\mathbf{0}, \mathbf{K}_{N+T} + \sigma^2 \mathbf{I})$$

where

$$\mathbf{K}_{N+T} = \begin{pmatrix} \mathbf{K}_N & \mathbf{K}_{NT} \\ \mathbf{K}_{TN} & \mathbf{K}_T \end{pmatrix}$$

Gaussian Process Regression

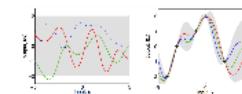
Finally, given this partitioned Gaussian distribution, we can condition the test variables on the training variables, such that

$$\rho(\mathbf{y}_T | \mathbf{y}) = \mathcal{N}(\mu_T, \Sigma_T)$$

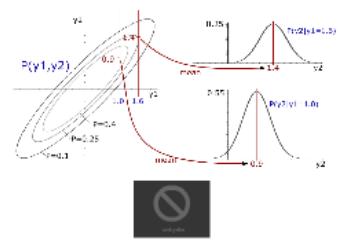
where

$$\mu_T = \mathbf{K}_{TN} (\mathbf{K}_N + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\Sigma_T = \mathbf{K}_T - \mathbf{K}_{TN} (\mathbf{K}_N + \sigma^2 \mathbf{I})^{-1} \mathbf{K}_{NT}$$



Gaussian Process Regression



Choosing Hyperparameters

Other non-Bayesian methods (such as Support Vector Machines) require the use of cross validation to choose the appropriate parameters for obtaining good predictive performance.

The advantage of Gaussian process approach is that we can choose the hyperparameters and covariances directly from the data!

In particular we can minimise the negative log marginal likelihood with respect to all hyperparameters θ ,

$$\mathcal{L} = -\log \rho(\mathbf{y}|\theta) = -\frac{1}{2} \log \det \mathbf{C}(\theta) + \frac{1}{2} \mathbf{y}^T \mathbf{C}(\theta) \mathbf{y} + \frac{N}{2} \log(2\pi)$$

where $\mathbf{C}(\theta) = \mathbf{K} + \sigma^2 \mathbf{I}$.

Relation to Linear Regression

Let's consider modelling functions defined by a weighted sum of a finite set of basis function, such that

$$f(x) = \sum_{m=1}^M \theta_m \phi_m(x) = \boldsymbol{\theta}^T \phi(x)$$

We can place a Gaussian prior on the weights, so that $\rho(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \Sigma_\theta)$.

This defines a Gaussian process with a covariance function given by,

$$K(x, x') = E[f(x), f(x')] = \phi(x)^T \Sigma_\theta \phi(x')$$

Relation to Linear Regression

There is in fact a direct correspondence between covariance functions and basis functions given by [Mercer's theorem](#), which states that we can always decompose a covariance function into eigenfunctions and eigenvectors,

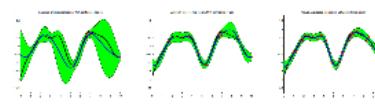
$$K(x, x') = \sum_{i=1}^{\infty} \lambda_i \varphi_i(x) \varphi_i(x')$$

If the sum is finite then it can be defined analogously using linear regression.

Often the sum is infinite, and there are no closed form solutions for the eigenvectors and eigenvalues.

We see how this is another form of the *kernel trick* - the covariance function effectively transforms the input space into a high or even infinite dimensional feature space, yet we can still compute in the lower dimensional space.

Example: Increasing Data



Summary

Gaussian processes are a very powerful approach for performing nonparametric supervised learning.

We can use the kernel trick to allow us to do linear regression effectively using an infinite number of basis functions, while retaining analytic tractability.

The main challenge is the computational scaling; it is cubic in the number of observations.

M5MS10

Machine Learning

Spring 2018

Lecture 6

Dr Ben Calderhead
b.calderhead@imperial.ac.uk

