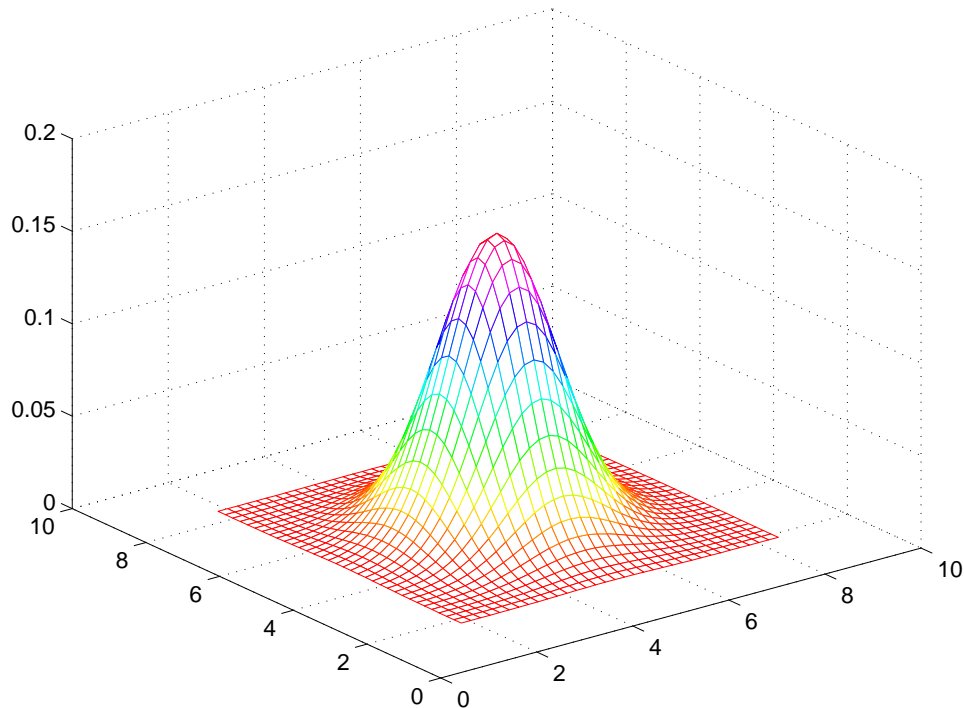


**University of Cambridge
Engineering Part IIB**

**Paper 4F10: Statistical Pattern
Processing**

Handout 8: Gaussian Processes



Bill Byrne
bill.byrne@eng.cam.ac.uk
Michaelmas 2012

Outline

The next two lectures cover a range of topics motivating and describing **Gaussian Processes**. The following areas will be discussed:

- Regression problems
- Linear regression
 - Bayesian linear regression
- Regression/interpolation with basis functions
- Gaussian process regression
 - form of covariance function
- Classification with regression models
 - Bayesian classification
- Classification with Gaussian Processes

Regression

The majority of the techniques so far described have concentrated on classification rather than regression. Prediction/interpolation are standard, commonly occurring, regression problems.

Here rather than mapping from an observation \mathbf{x} to a class $y \in \{-1, 1\}$, we are interested in predicting a continuous value given a set of observations. Examples of this are:

- stock market prediction
- speech enhancement

The first part of these two lectures on Gaussian processes will look at regression. The application of Gaussian processes to classification will then be considered. For both cases supervised training will be assumed.

The regression (or prediction) is assumed to have two parts:

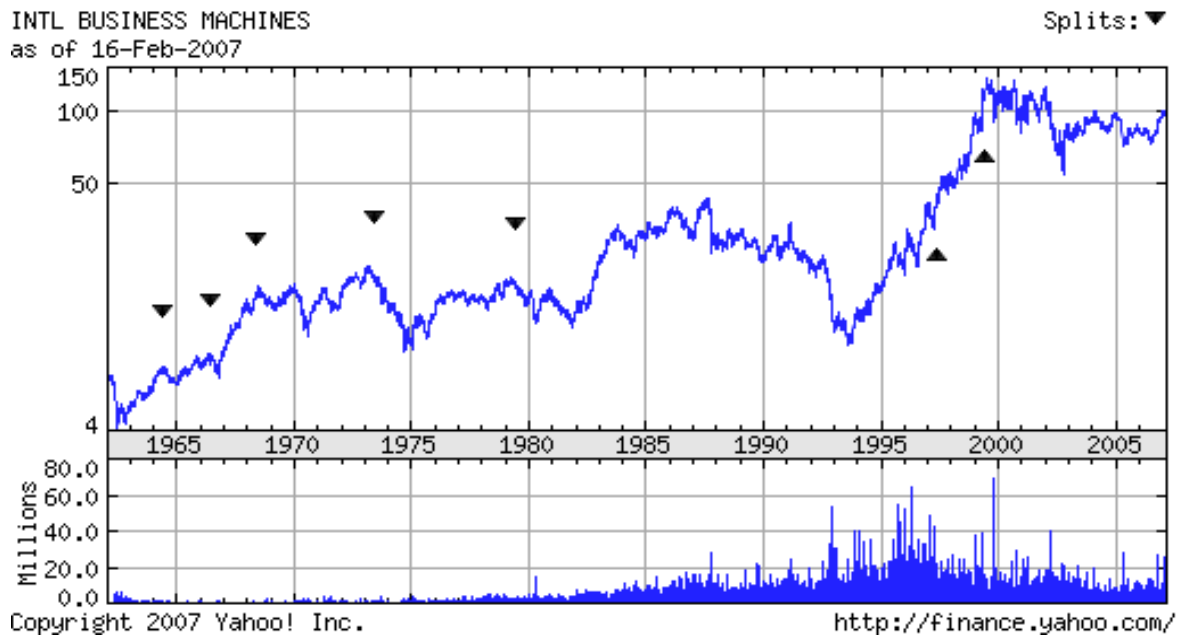
- **systematic variation**: can be accurately predicted if the underlying process is known ($f(\mathbf{x})$)
- **random variation**: the inherent “unpredictability” of the system.

These two parts will be modelled separately

$$y = f(\mathbf{x}) + \epsilon$$

where ϵ will be assumed to be Gaussian distributed, with a mean of zero and variance σ_n^2 .

Stock Market Prediction



- Would like to make predictions and actions
- What value will IBM shares have today given the past performance?
 - IBM is now at \$120.82 - \$92.68 last year though not a linear rise!
- As well as the predicted value I am probably interested in the **confidence** of the prediction ...

Gaussian Process application

- Gaussian process approximations of **stochastic differential equations**.
- Prediction of **deformed and annealed microstructures** using Bayesian neural networks and Gaussian processes.
- **Biomarker discovery in microarray gene expression data** with Gaussian processes.
- Gaussian process model based **predictive control**.
- Gaussian processes for **multiuser detection in cdma receivers**.
- Learning a Gaussian process prior for **automatically generating music playlists**.
- GPPS: A Gaussian process **positioning system for cellular networks**.

For a wide-range of information and papers on Gaussian processes see

<http://www.gaussianprocess.org/>

Basic Probability Revisited (again)

These lectures will make extensive use of the following standard probability concepts:

- **continuous random variables**: the RV can take any value within a, possibly, infinite, range. Associated **probability density function** satisfies

$$\int p(x)dx = 1; \quad p(x) \geq 0$$

- **joint** probability $p(x, y)$
- **joint** independence $p(x, y) = p(x)p(y)$
- **conditional** probability $p(x|y)$
- **Bayes' rule**

$$p(x|y) = \frac{p(x)p(y|x)}{p(y)}$$

- **marginal** probability

$$p(x) = \int_y p(x, y)dy$$

where the integral is over all possible values of y .

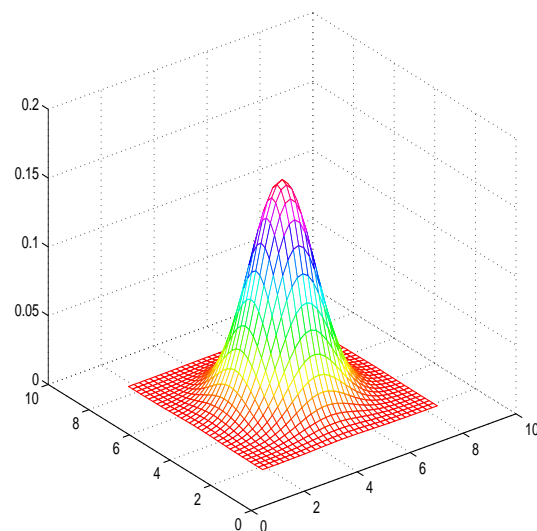
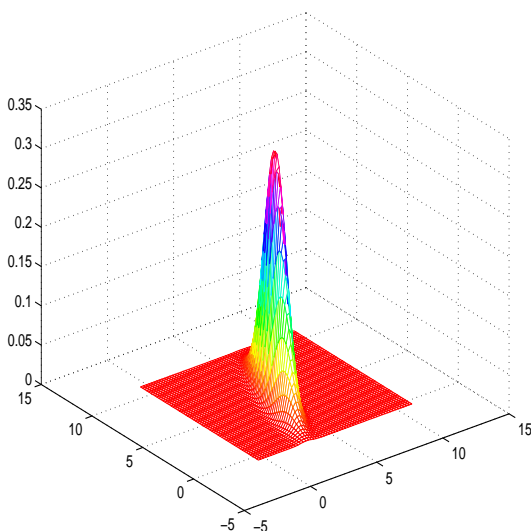
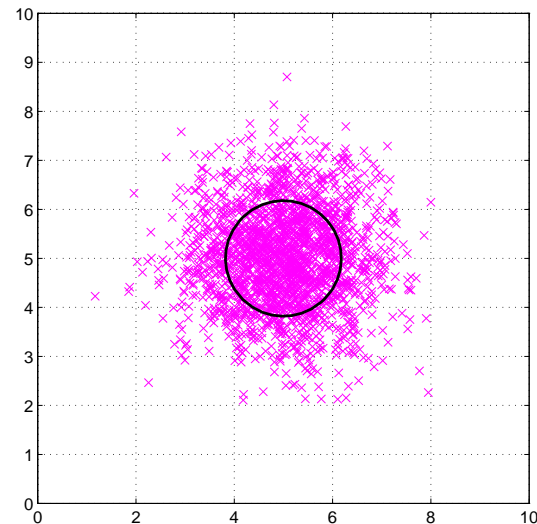
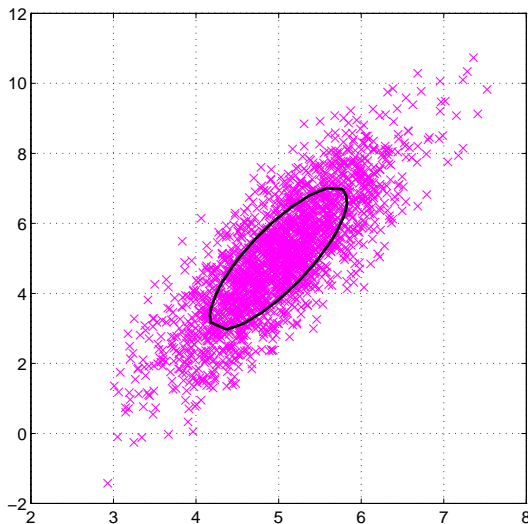
Multivariate Gaussian Distribution

For a d -dimensional data feature vector, \mathbf{x} , the multivariate Gaussian distribution is

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

$$\Sigma = \begin{pmatrix} 3 & 1 \\ 1 & 0.5 \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$



Gaussian Distributions

Conditional Gaussian distribution:

Consider the joint distribution between vectors \mathbf{x}_1 and \mathbf{x}_2

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix} \right)$$

We are interested in the conditional distribution, which itself is Gaussian

$$\mathbf{x}_1 | \mathbf{x}_2 \sim \mathcal{N}(\boldsymbol{\mu}_{1|2}, \boldsymbol{\Sigma}_{1|2})$$

where

$$\begin{aligned} \boldsymbol{\mu}_{1|2} &= \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ \boldsymbol{\Sigma}_{1|2} &= \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \end{aligned}$$

Product of Gaussian distributions:

Consider the two distributions

$$p_1(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \quad p_2(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$$

The product is an **un-normalised Gaussian**

$$p_1(\mathbf{x})p_2(\mathbf{x}) \propto \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

where

$$\begin{aligned} \boldsymbol{\mu} &= \boldsymbol{\Sigma} (\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2) \\ \boldsymbol{\Sigma} &= (\boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1})^{-1} \end{aligned}$$

Linear Interpolation

Consider the supervised training data of n samples, each with an observation \mathbf{x}_i and output y_i . The regression function $f(\mathbf{x})$ is **linear** if defined as

$$f(\mathbf{x}) = \mathbf{w}'\mathbf{x}$$

and the target value has Gaussian noise so that

$$y = f(\mathbf{x}) + \epsilon$$

where $\epsilon \sim \mathcal{N}(0, \sigma_n^2)$.

For a given value of \mathbf{w} the likelihood of the outputs can be expressed as

$$p(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{w}) = \prod_{i=1}^n p(y_i | \mathbf{x}_i, \mathbf{w}) = \mathcal{N}(\mathbf{y}; \mathbf{X}'\mathbf{w}, \sigma_n^2 \mathbf{I})$$

where \mathbf{I} is the $n \times n$ identity matrix and

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_n \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

The value of \mathbf{w} can then be found by maximising this likelihood. This is equivalent to minimising the least squares cost function

$$\min_{\mathbf{w}} \left\{ \sum_{i=1}^n (y_i - \mathbf{w}'\mathbf{x}_i)^2 \right\} = \min_{\mathbf{w}} \left\{ \|\mathbf{y} - \mathbf{X}'\mathbf{w}\|^2 \right\}$$

which can be solved using the pseudo-inverse

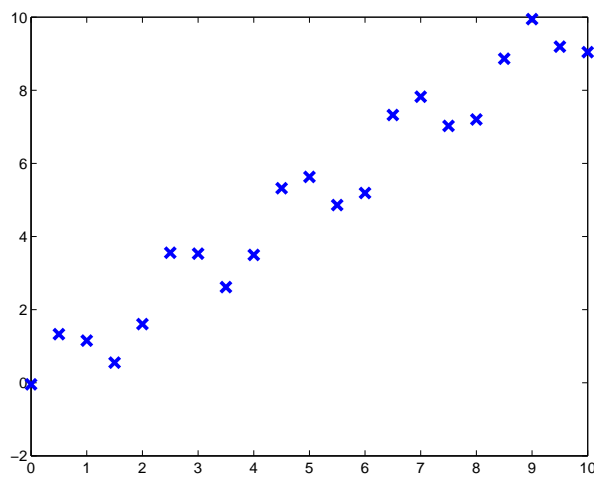
$$\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}')^{-1} \mathbf{X}\mathbf{y}$$

Example of Linear Regression

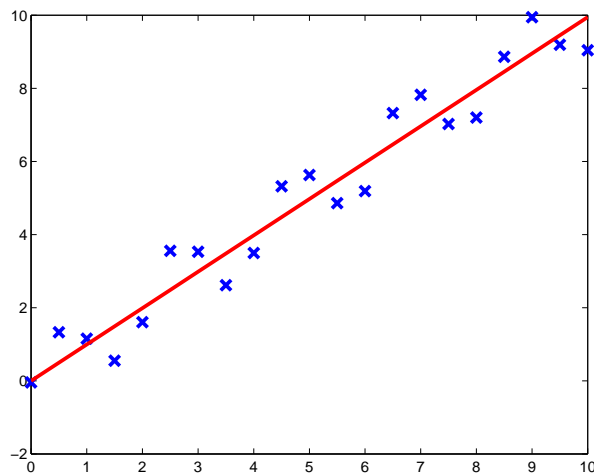
Generate data using

$$y = x + \sin(3x) + \epsilon$$

where $\epsilon \sim \mathcal{N}(0, 0.01)$.



A set of points generated using this function



Linear regression (interpolation/prediction) using this data

$$f(x) = 0.9948x$$

Bayesian Parameter Estimation

The parameter estimation schemes have assumed that there is a **single estimate** of the model parameters that can be obtained from the given training data. As this is only one possible set of training data (from all possible sets), then **error bars** are sometimes used to reflect this uncertainty.

In contrast a Bayesian scheme for estimating the model parameters, estimates a **distribution over the model parameters**. This distribution represents the uncertainty.

These two approaches are often referred to as

- **Frequentist**: linear regression prediction from \mathbf{x} is

$$\hat{y} = \mathbf{x}'\hat{\mathbf{w}}$$

where $\hat{\mathbf{w}}$ is the ML estimate (possibly include error bars)

- **Bayesian**: the distribution over the prediction is

$$p(y|\mathbf{x}, \mathbf{y}, \mathbf{X}) = \int p(y|\mathbf{x}, \mathbf{w})p(\mathbf{w}|\mathbf{y}, \mathbf{X})d\mathbf{w}$$

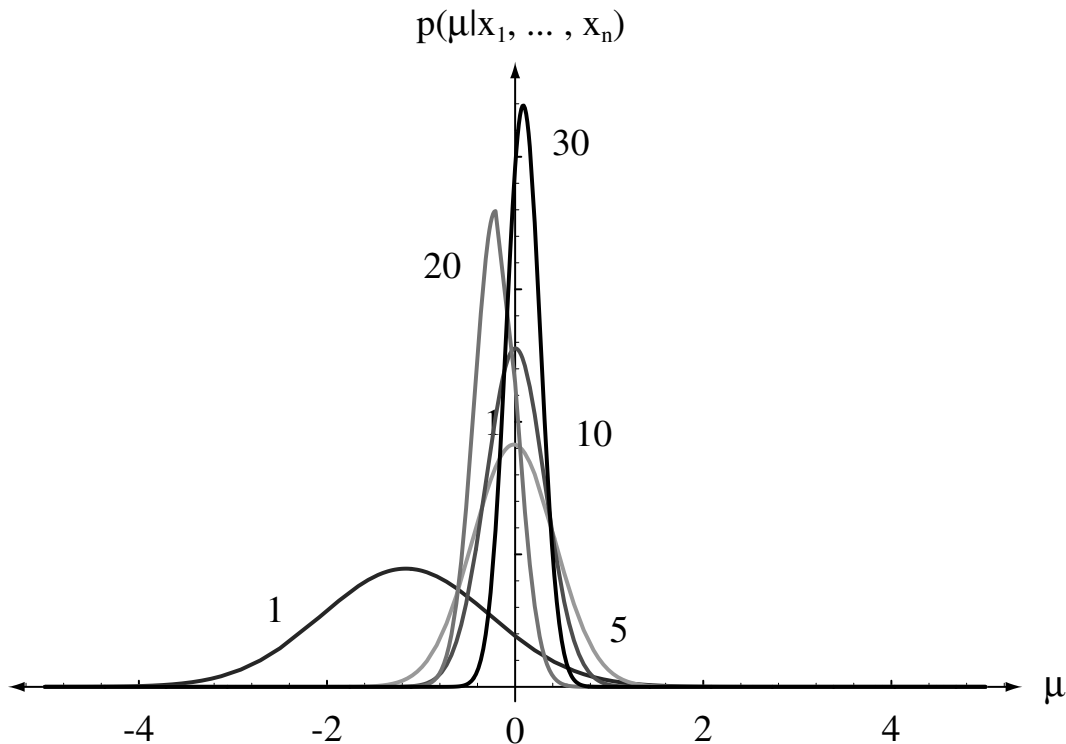
where $p(\mathbf{w}|\mathbf{y}, \mathbf{X})$ is the distribution over the model parameters given the training data. The point estimate can be selected using for example the mean of the output distribution.

Bayesian approaches are becoming increasingly popular as the additional computational load of the **marginalisation** (integration) is becoming less of an issue through the use of approximations (and more computers!).

Gaussian Mean “Estimation”

As an aside Bayesian estimation can be applied to model parameters. Consider the Bayesian estimation of the mean.

The figure below (from DHS) shows how the posterior distribution for estimating the mean of a Gaussian distribution as the number of data samples increases.



- The posterior becomes more sharply peaked (reduced variance) and the MAP estimate (the mode of the distribution) gradually changes.
- As the number of samples increases the MAP estimate will tend towards the ML estimate.

Bayesian Linear Regression

A **Bayesian** version of linear regression may be used. Here a distribution over the weight vector is estimated and then marginalised over to yield the prediction. “Standard” linear regression can be written as

$$p(y|\mathbf{x}, \mathbf{X}, \mathbf{y}) = \int p(y|\mathbf{x}, \mathbf{w})p(\mathbf{w}|\mathbf{X}, \mathbf{y})d\mathbf{w}$$

where

$$\begin{aligned} p(y|\mathbf{x}, \mathbf{w}) &= \mathcal{N}(y; \mathbf{w}'\mathbf{x}, \sigma_n^2) \\ p(\mathbf{w}|\mathbf{X}, \mathbf{y}) &= \delta(\mathbf{w} - \hat{\mathbf{w}}) \end{aligned}$$

This can then be written as

$$p(y|\mathbf{x}, \mathbf{X}, \mathbf{y}) = \mathcal{N}(y; \hat{\mathbf{w}}'\mathbf{x}, \sigma_n^2)$$

If only the deterministic part of the regression progress is considered

$$f(\mathbf{x}) = \hat{\mathbf{w}}'\mathbf{x}$$

More generally the **full** distribution can be used. From the training data, \mathbf{y} and \mathbf{X} , the distribution of the weights is

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

the denominator is independent of the weight vector so

$$p(\mathbf{w}|\mathbf{y}, \mathbf{X}) \propto p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})$$

Weight Vector Posterior Distribution

If the weights vector has a Gaussian prior then

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \mathbf{w}) &= \mathcal{N}(\mathbf{y}; \mathbf{X}'\mathbf{w}, \sigma_n^2\mathbf{I}) \\ \mathbf{w} &\sim \mathcal{N}(\mathbf{0}, \sigma_w^2\mathbf{I}) \end{aligned}$$

The distribution of the weights can then be written as

$$\log(p(\mathbf{w}|\mathbf{y}, \mathbf{X})) \propto -\frac{1}{2\sigma_n^2}(\mathbf{y} - \mathbf{X}'\mathbf{w})'(\mathbf{y} - \mathbf{X}'\mathbf{w}) - \frac{1}{2\sigma_w^2}\mathbf{w}'\mathbf{w}$$

As both are Gaussian distributed

$$p(\mathbf{w}|\mathbf{y}, \mathbf{X}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w)$$

where

$$\begin{aligned} \boldsymbol{\mu}_w &= \frac{1}{\sigma_n^2} \left(\frac{1}{\sigma_n^2} \mathbf{X} \mathbf{X}' + \frac{1}{\sigma_w^2} \mathbf{I} \right)^{-1} \mathbf{X} \mathbf{y} \\ \boldsymbol{\Sigma}_w &= \left(\frac{1}{\sigma_n^2} \mathbf{X} \mathbf{X}' + \frac{1}{\sigma_w^2} \mathbf{I} \right)^{-1} \end{aligned}$$

Note: as $\sigma_w^2 \rightarrow \infty$, then $\boldsymbol{\mu}_w \rightarrow \hat{\mathbf{w}}$ the ML estimate.

This has so far only considered the [training](#).

[What about regression/prediction?](#)

Bayesian Linear Regression (cont)

The joint distribution of all the weights is Gaussian. Regression can be based on this weight posterior.

Given an observation x , and the training data \mathbf{X} and \mathbf{y} , what is the distribution of the output value y .

- First examine the **deterministic part**: $f(\mathbf{x})$

$$\begin{aligned} p(f(\mathbf{x})|\mathbf{x}, \mathbf{X}, \mathbf{y}) &= \int p(f(\mathbf{x})|\mathbf{x}, \mathbf{w})p(\mathbf{w}|\mathbf{X}, \mathbf{y})d\mathbf{w} \\ &= \mathcal{N}(f(\mathbf{x}); \mathbf{x}'\boldsymbol{\mu}_w, \mathbf{x}'\boldsymbol{\Sigma}_w\mathbf{x}) \end{aligned}$$

since $f(\mathbf{x})$ is a linear transform of the weight vector \mathbf{w} . This distribution over the function $f(\mathbf{x})$ includes all the information that can be extracted from the training data and the observation, \mathbf{x} .

- The **distribution of the target value** y is easy to find - this is the deterministic part, $f(\mathbf{x})$, plus the noise ϵ , as these are independent and Gaussian

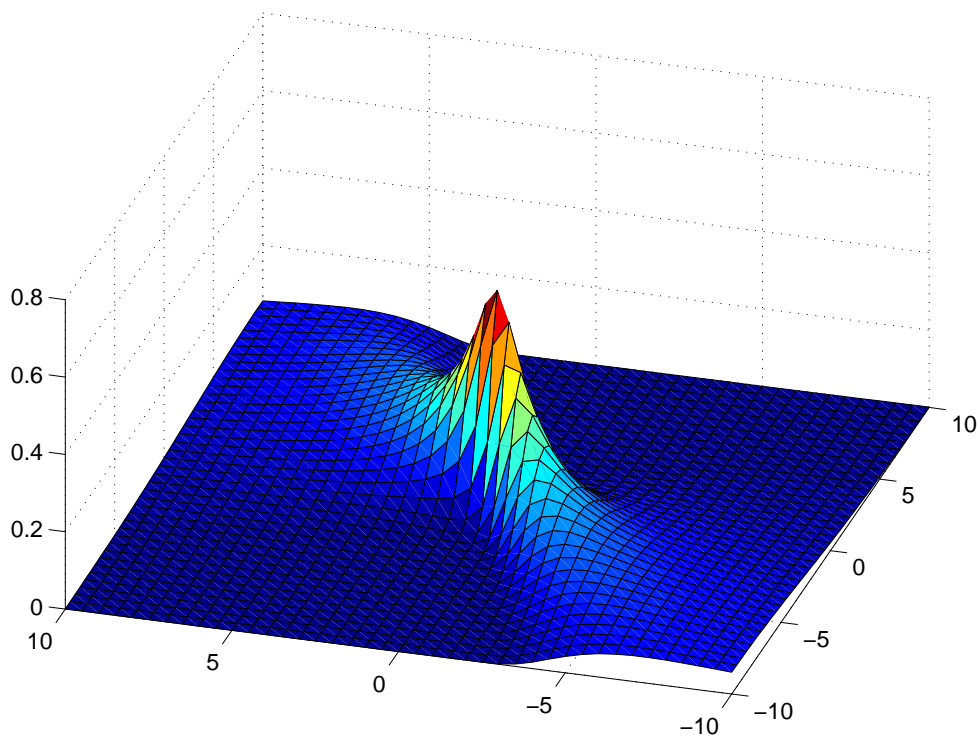
$$\begin{aligned} p(y|\mathbf{x}, \mathbf{X}, \mathbf{y}) &= \int p(y|f(\mathbf{x}))p(f(\mathbf{x})|\mathbf{X}, \mathbf{y}, \mathbf{x})df(\mathbf{x}) \\ &= \mathcal{N}(y; \mathbf{x}'\boldsymbol{\mu}_w, \mathbf{x}'\boldsymbol{\Sigma}_w\mathbf{x} + \sigma_n^2) \end{aligned}$$

Thus the posterior distribution of the weights are Gaussian as is the prediction of the observations. This is not surprising as the relationship the output and the function is linear. If this relationship is non-linear then the posterior and prediction will no longer be Gaussian.

Bayesian Linear Regression Example

Consider a simple example with

- 1-dimensional observations, $x, \sigma_n^2 = 0.5$
- From training data $\mu_w = 1, \sigma_w^2 = 0.5$



$$p(y|x, \mathbf{X}, \mathbf{y}) = \mathcal{N}(y; x\mu_w, x^2\sigma_w^2 + \sigma_n^2)$$

Above plot shows distribution over the input/output-space

- Regression forced to go through zero
 - variance at this point σ_n^2
- Variance increases as we move away from the origin
- The increase in variance depends on σ_w^2
 - decreases as $\sum_{i=1}^n x_i^2$ increases

Marginal Likelihood

In addition to using the training data for prediction it will also be useful to look at the **marginal likelihood** (compare this to the evidence)

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w}$$

This is useful as the values of \mathbf{w} are not observed, whereas \mathbf{X} and \mathbf{y} are. Effectively \mathbf{w} is treated as a **latent variable**.

Using the distributions defined on the previous slide, it is clear that the marginal likelihood is also Gaussian distributed

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}) &= \int \mathcal{N}(\mathbf{y}; \mathbf{X}'\mathbf{w}, \sigma_n^2\mathbf{I})\mathcal{N}(\mathbf{w}; \mathbf{0}, \sigma_w^2\mathbf{I})d\mathbf{w} \\ &= \mathcal{N}(\mathbf{y}; \mathbf{0}, \sigma_w^2\mathbf{X}'\mathbf{X} + \sigma_n^2\mathbf{I}) \end{aligned}$$

Note this has made use of the following matrix equality

$$\frac{1}{\sigma_n^2}\mathbf{I} - \frac{1}{\sigma_n^2}\mathbf{X}'\left(\frac{1}{\sigma_n^2}\mathbf{X}\mathbf{X}' + \frac{1}{\sigma_w^2}\mathbf{I}\right)^{-1}\mathbf{X}\frac{1}{\sigma_n^2} = \left(\sigma_w^2\mathbf{X}'\mathbf{X} + \sigma_n^2\mathbf{I}\right)^{-1}$$

The marginal likelihood can be used as a mechanism for finding the parameters σ_w^2 and σ_n^2 .

Interpolation using Basis Functions

Linear interpolation is limited as the relationship between the observations and the output (or target value) is constrained to be linear. An alternative is to perform interpolation using a set of basis functions.

Consider a **basis function** of the form $\phi(\|\mathbf{x} - \mathbf{x}_i\|)$, where $\phi()$ is some non-linear function and $\|\mathbf{x} - \mathbf{x}_i\|$ is a **distance** of the vector \mathbf{x} from the prototype vector \mathbf{x}_i .

For the case of n training examples each being used as a prototype, the mapping can be defined as

$$f(\mathbf{x}) = \sum_{i=1}^n w_i \phi(\|\mathbf{x} - \mathbf{x}_i\|) = \phi(\mathbf{x})' \mathbf{w}$$

where

$$\phi(\mathbf{x}) = \left[\phi(\|\mathbf{x} - \mathbf{x}_1\|) \quad \dots \quad \phi(\|\mathbf{x} - \mathbf{x}_n\|) \right]'$$

The output value is again considered to be

$$y = f(\mathbf{x}) + \epsilon$$

The values for \mathbf{w} needs to estimated. Following the standard linear interpolation example, for the training data

$$\Phi = \begin{bmatrix} \phi(\|\mathbf{x}_1 - \mathbf{x}_1\|) & \dots & \phi(\|\mathbf{x}_1 - \mathbf{x}_n\|) \\ \vdots & \ddots & \vdots \\ \phi(\|\mathbf{x}_n - \mathbf{x}_1\|) & \dots & \phi(\|\mathbf{x}_n - \mathbf{x}_n\|) \end{bmatrix} = \begin{bmatrix} \phi(\mathbf{x}_1)' \\ \vdots \\ \phi(\mathbf{x}_n)' \end{bmatrix}$$

So for the training data

$$\mathbf{y} = \Phi \mathbf{w} + \epsilon$$

Interpolation using Basis Functions (2)

If the inverse Φ^{-1} exists then the ML estimate is

$$\hat{\mathbf{w}} = \Phi^{-1} \mathbf{y}$$

It has been shown that for a large class of functions $\phi(\cdot)$ if the set of points $\mathbf{x}_1, \dots, \mathbf{x}_n$ is distinct then Φ^{-1} exists.

With the weights defined above the function $f(\mathbf{x})$ represents a continuous function that passes through **all** the target points, y_i - this is **exact** interpolation.

There are number of standard forms of $\phi(\cdot)$ used:

Non-linearity	Form ($z = (\ \mathbf{x} - \mathbf{x}_p\)/\sigma$)
Gaussian	$\exp(-z^2/2)$
Exponential	$\exp(-z)$
Quadratic	$z^2 + \alpha z + \beta$
Inverse Quadratic	$\frac{1}{1+z^2}$
Thin plate spine	$z^\alpha \log(z)$
Trigonometric	$\sin(z)$

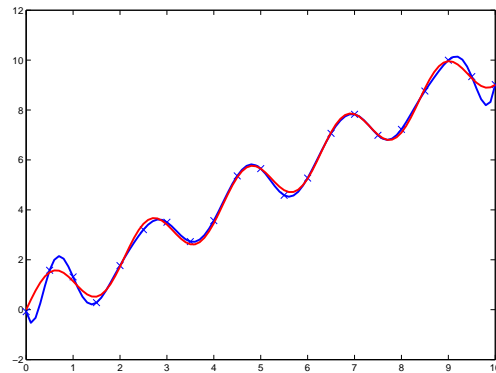
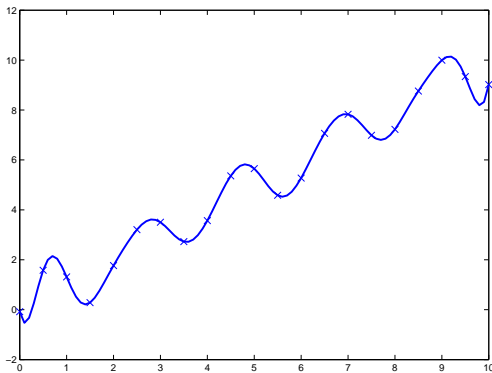
σ is a variable that controls the smoothness of the interpolation function.

Some of the functions are **localised basis functions**. These have the property that $\phi(z) \rightarrow 0$ as $z \rightarrow \infty$.

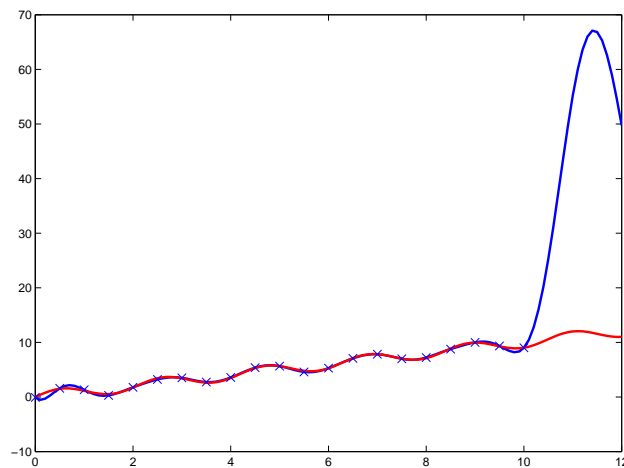
Basis Function Example

Interpolation was performed using a Gaussian basis function

$$\phi(\|x - x_i\|) = \exp\left(-\frac{1}{2}(x - x_i)^2\right)$$



Using the Gaussian basis finds the non-linearities and models the function well (exactly models the training points!).



However prediction is not too good!

Useful to use a Bayesian approach to directly get distribution

Bayesian Regression with Basis Functions

The previous analysis of Bayesian Linear regression can be applied to regression with basis functions.

Again, let the prior distribution for the weight vector be Gaussian distributed

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{0}, \sigma_w^2 \mathbf{I})$$

We need the posterior weight distribution. Now

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{y}; \Phi \mathbf{w}, \sigma_n^2 \mathbf{I})$$

From the linear case (equating Φ and \mathbf{X}')

$$p(\mathbf{w} | \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w)$$

where

$$\boldsymbol{\mu}_w = \frac{1}{\sigma_n^2} \left(\frac{1}{\sigma_n^2} \Phi' \Phi + \frac{1}{\sigma_w^2} \mathbf{I} \right)^{-1} \Phi' \mathbf{y}$$

$$\boldsymbol{\Sigma}_w = \left(\frac{1}{\sigma_n^2} \Phi' \Phi + \frac{1}{\sigma_w^2} \mathbf{I} \right)^{-1}$$

Thus a Bayesian form of **non-linear regression** may be addressed in the same fashion as linear regression, by using basis functions.

Training Output Covariance Matrix

It is interesting (and related to Gaussian Processes) to look at the covariance matrix of the training outputs, \mathbf{y} . The covariance matrix of \mathbf{y} , $\Sigma_{\mathbf{y}}$, can then be written as

$$\Sigma_{\mathbf{y}} = \mathcal{E} \{ \mathbf{y} \mathbf{y}' \} - \mathcal{E} \{ \mathbf{y} \} \mathcal{E} \{ \mathbf{y} \}'$$

By inspection the mean is zero (\mathbf{w} has zero mean prior, as is ϵ)

$$\begin{aligned} \Sigma_{\mathbf{y}} &= \mathcal{E} \{ \mathbf{y} \mathbf{y}' \} \\ &= \mathcal{E} \{ (\Phi \mathbf{w} + \epsilon)(\Phi \mathbf{w} + \epsilon)' \} \\ &= \sigma_{\mathbf{w}}^2 \Phi \Phi' + \sigma_{\mathbf{n}}^2 \mathbf{I} \end{aligned}$$

(same form as obtained for the marginal likelihood)

The target values are all jointly Gaussian under these assumptions. The same process can be applied to the weights posterior distribution. These are both examples of **Gaussian Processes**.

Interestingly element i, j element of the covariance matrix, $\Sigma_{\mathbf{y}}$, has the form

$$\sigma_{\mathbf{w}}^2 \phi(\mathbf{x}_i)' \phi(\mathbf{x}_j) + \sigma_{\mathbf{n}}^2 \delta_{ij}$$

where

$$\delta_{ij} = \begin{cases} 1, & i = j \\ 0, & \text{otherwise} \end{cases}$$

Gaussian Processes

The definition is:

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

A Gaussian Process is completely specified by:

- **mean function:**

$$m(\mathbf{x}) = \mathcal{E}\{f(\mathbf{x})\}$$

- **covariance function:**

$$k(\mathbf{x}, \mathbf{x}_j) = \mathcal{E}\{(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}_j) - m(\mathbf{x}_j))\}$$

A Gaussian Process for two elements can be written as

$$\begin{bmatrix} f(\mathbf{x}_i) \\ f(\mathbf{x}_j) \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} m(\mathbf{x}_i) \\ m(\mathbf{x}_j) \end{bmatrix}, \begin{bmatrix} k(\mathbf{x}_i, \mathbf{x}_i) & k(\mathbf{x}_i, \mathbf{x}_j) \\ k(\mathbf{x}_i, \mathbf{x}_j) & k(\mathbf{x}_j, \mathbf{x}_j) \end{bmatrix} \right)$$

For simplicity the mean function will be taken as zero. Note for linear regression $f(\mathbf{x}) = \mathbf{w}'\mathbf{x}$.

Initially GPs will be considered for **regression**. This will have a form similar to Bayesian regression. The distribution of the output (including the random variation)

$$p(y|\mathbf{X}, \mathbf{x}, \mathbf{y}) = \int p(y|f(\mathbf{x}))p(f(\mathbf{x})|\mathbf{X}, \mathbf{x}, \mathbf{y})df(\mathbf{x})$$

where $f(\mathbf{x})$ is the systematic prediction.

Regression using Gaussian Processes

We would like to do prediction given that the outputs are all jointly Gaussian.

Initially let's examine obtaining the **best** prediction process $f(\mathbf{x})$ given the training data.

- Assume that the prediction process is jointly Gaussian distributed with the training outputs.
- As in basis function regression let the covariance function be a function of the training observations \mathbf{X} (and regression observation \mathbf{x}).

The **joint distribution** can then be written as

$$\begin{bmatrix} f(\mathbf{x}) \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} k(\mathbf{x}, \mathbf{x}) & \mathbf{k}(\mathbf{x}, \mathbf{X})' \\ \mathbf{k}(\mathbf{x}, \mathbf{X}) & \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I} \end{bmatrix}\right)$$

The **conditional distribution** can then be expressed as

$$f(\mathbf{x}) | \mathbf{X}, \mathbf{y}, \mathbf{x} \sim \mathcal{N}(\mu_f, \sigma_f^2)$$

where (using equalities from slide 7)

$$\begin{aligned} \mu_f &= \mathbf{k}(\mathbf{x}, \mathbf{X})' [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{y} \\ \sigma_f^2 &= k(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x}, \mathbf{X})' [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{k}(\mathbf{x}, \mathbf{X}) \end{aligned}$$

and

$$\mathbf{k}(\mathbf{x}, \mathbf{X}) = \begin{bmatrix} k(\mathbf{x}, \mathbf{x}_1) \\ \vdots \\ k(\mathbf{x}, \mathbf{x}_n) \end{bmatrix}, \quad \mathbf{K}(\mathbf{X}, \mathbf{X}) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_n, \mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_1, \mathbf{x}_n) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$$

Basis Function Regression (Revisited)

Consider the linear model for regression

$$f(\mathbf{x}) = \sum_{i=1}^n w_i \phi(\|\mathbf{x} - \mathbf{x}_i\|) = \mathbf{w}'\boldsymbol{\phi}(\mathbf{x})$$

where

$$\boldsymbol{\phi}(\mathbf{x}) = \left[\phi(\|\mathbf{x} - \mathbf{x}_1\|) \quad \dots \quad \phi(\|\mathbf{x} - \mathbf{x}_n\|) \right]'$$

with a Gaussian prior distribution $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma_w^2 \mathbf{I})$ then

$$\mathcal{E}\{f(\mathbf{x})\} = \mathcal{E}\{\mathbf{w}'\}\boldsymbol{\phi}(\mathbf{x}) = 0$$

$$\begin{aligned} \mathcal{E}\{f(\mathbf{x})f(\mathbf{x})\} &= \boldsymbol{\phi}(\mathbf{x})'\mathcal{E}\{\mathbf{w}\mathbf{w}'\}\boldsymbol{\phi}(\mathbf{x}) \\ &= \sigma_w^2 \boldsymbol{\phi}(\mathbf{x})'\boldsymbol{\phi}(\mathbf{x}) \end{aligned}$$

This has the same form as the **kernels** previously seen where $\boldsymbol{\phi}(\mathbf{x})$ defines the **feature-space**.

$$\mathcal{E}\{f(\mathbf{x}_i)f(\mathbf{x}_j)\} = k(\mathbf{x}_i, \mathbf{x}_j) = \sigma_w^2 \langle \boldsymbol{\phi}(\mathbf{x}_i), \boldsymbol{\phi}(\mathbf{x}_j) \rangle$$

The distribution of the training outputs ($p(\mathbf{y}|\mathbf{X})$) can then be written in terms of kernel operations

$$\boldsymbol{\Sigma}_y = \sigma_w^2 \boldsymbol{\Phi}\boldsymbol{\Phi}' + \sigma_n^2 \mathbf{I} = \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}$$

again

$$\mathbf{K}(\mathbf{X}, \mathbf{X}) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_n, \mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_1, \mathbf{x}_n) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$$

Covariance Functions and Kernels

The form of $\mathbf{K}(\mathbf{X}, \mathbf{X})$ is determined by the covariance function. However it has the same form as the **Gram matrix** used when describing kernels with SVM.

The Gram matrix associated with a covariance functions **must** be positive semi-definite (PSD) i.e. for a real matrix \mathbf{K}

$$\mathbf{v}'\mathbf{K}\mathbf{v} \geq 0$$

and these matrices will also be symmetric, $k(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_j, \mathbf{x}_i)$. It is possible to have kernels that are not PSD, but Covariance Functions **must** be PSD.

In the Gaussian Process literature covariance functions are split into two broad groups:

- **stationary**: are a function of $\mathbf{x}_i - \mathbf{x}_j$. These are invariant to translations in the input space. The various forms of basis interpolation will fall under this category.
- **non-stationary**: these are functions which do vary with translation.

The most important sub-set of these covariance functions are **dot-product** covariance functions. These are only functions of $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$.

- **linear**: $k(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$
- **polynomial**, order d : $k(\mathbf{x}_i, \mathbf{x}_j) = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle + 1)^d$

Squared Exponential Covariance Function

One standard stationary covariance function is given by

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma_1^2}\right)$$

This has **free parameters** (or hyper-parameters)

- **length scale** - σ_1
- **signal variance** - α

This is often used as the covariance function over the systematic prediction function, $f(\mathbf{x})$. The covariance function for the observation would be

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \exp\left(\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma_1^2}\right)$$

There will normally be some prediction noise. In this case the covariance function will become

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma_1^2}\right) + \sigma_n^2 \delta_{ij}$$

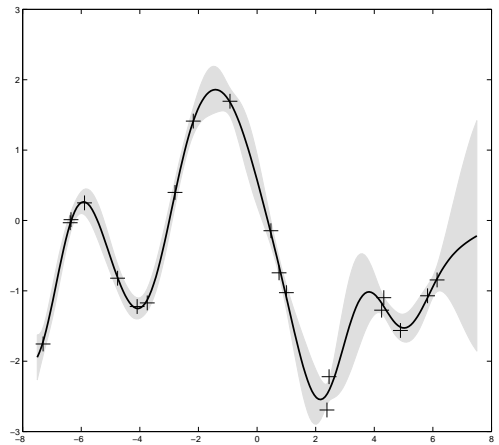
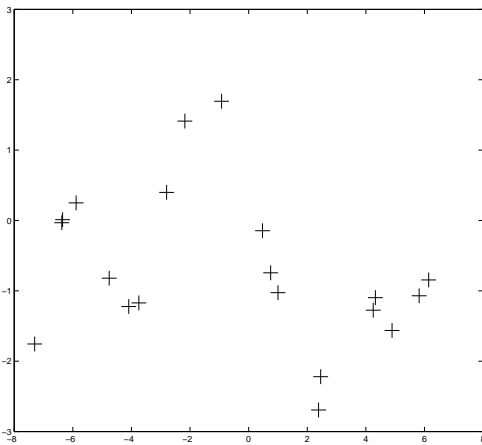
where

$$\delta_{ij} = \begin{cases} 1, & i = j \\ 0, & \text{otherwise} \end{cases}$$

Example Gaussian Process

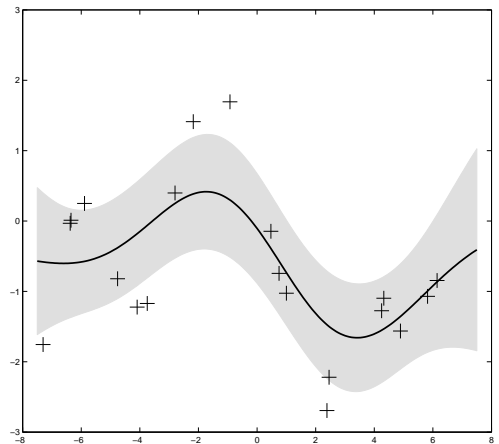
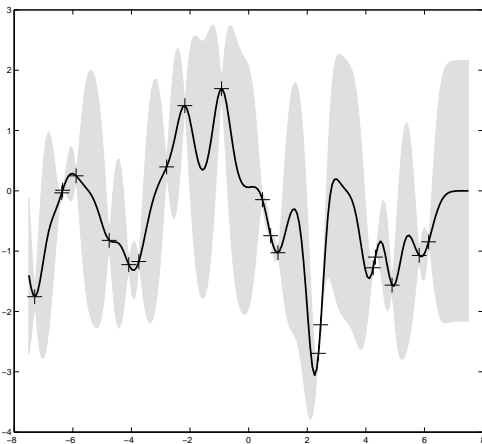
Consider a covariance function of the form

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma_1^2}\right) + \sigma_n^2 \delta_{ij}$$



Points (drawn from right)

$\sigma_1 = 1.0, \alpha = 1.0, \sigma_n = 0.1$



$\sigma_1 = 0.3, \alpha = 1.08, \sigma_n = 5.0 \times 10^{-5}$ $\sigma_1 = 3.0, \alpha = 1.16, \sigma_n = 0.89$

The grey region shows the ± 2 standard deviations (95% confidence) - $f(\mathbf{x})$ only (noise std, σ_n subtracted)

Gaussian Process Practicalities

It is interesting to note that it is possible to write the mean of the prediction

$$\mu_{\mathbf{f}} = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i)$$

where

$$\boldsymbol{\alpha} = [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{y}$$

In general **all** values of α_i will be non-zero.

Thus regression/prediction is expressed in a **non-sparse** combination of the kernel between the “test” observation and all the training points. This can be contrasted with the SVM where the decision boundary is only a function of the support vectors (thus it is highly sparse).

In the same fashion as SVM decision boundary Gaussian process regression is limited to the space spanned within the feature-space ($\phi(\mathbf{x})$) of all the training points.

One of the major problems with implementing Gaussian Processes is the **computational cost**. This is dominated by the need to invert $\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}$. This has dimensionality of the number of training examples. For large training sets a number of approximate approaches can be used.

Learning Hyper-parameters

So far it has been assumed that the free parameters, the **hyper-parameters** of the covariance function, θ , are known. It would be preferable to be able to train these.

One approach to do this is to maximise the **marginal likelihood**. This is (remembering the Bayesian interpretation of linear regression and assuming zero mean)

$$p(\mathbf{y}|\mathbf{X}, \theta) = \mathcal{N}(\mathbf{y}; \mathbf{0}, \mathbf{K}_\theta(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I})$$

where θ are the hyper-parameters.

It is possible to write

$$\log(p(\mathbf{y}|\mathbf{X}, \theta)) = -\frac{1}{2} \mathbf{y}' \Sigma_y^{-1} \mathbf{y} - \frac{1}{2} \log(|\Sigma_y|) - \frac{n}{2} \log(2\pi)$$

where

$$\Sigma_y = \mathbf{K}_\theta(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}$$

This expression may be considered to have two parts

- **data fit**: $-\frac{1}{2} \mathbf{y}' \Sigma_y^{-1} \mathbf{y}$
- **complexity penalty**: $-\frac{1}{2} \log(|\Sigma_y|)$

(the last term is a constant). Using this expression it is also possible to estimate the noise variance, σ_n^2 as well as θ .

Hyper-parameter Optimisation

Matrix Calculus (Reference)

We will need the gradient which will make use of

$$\begin{aligned}\frac{\partial}{\partial \theta} \mathbf{K}^{-1} &= -\mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial \theta} \mathbf{K}^{-1} \\ \frac{\partial}{\partial \theta} \log(|\mathbf{K}|) &= \text{tr} \left(\mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial \theta} \right)\end{aligned}$$

where

$$\text{tr}(\mathbf{K}) = \sum_{i=1}^n k_{ii}$$

assuming that \mathbf{K} is an $n \times n$ matrix.

Hyper-parameter Derivatives

We now need the partial derivatives of the marginal likelihood with respect to the hyper-parameters. Using the matrix calculus equalities

$$\begin{aligned}\frac{\partial}{\partial \theta_j} \log(p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})) &= \frac{1}{2} \mathbf{y}' \mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial \theta_j} \mathbf{K}^{-1} \mathbf{y} - \frac{1}{2} \text{tr} \left(\mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial \theta_j} \right) \\ &= \frac{1}{2} \text{tr} \left((\mathbf{K}^{-1} \mathbf{y} \mathbf{y}' \mathbf{K}^{-1} - \mathbf{K}^{-1}) \frac{\partial \mathbf{K}}{\partial \theta_j} \right)\end{aligned}$$

(using \mathbf{K} as a shorthand for $\mathbf{K}_\theta(\mathbf{X}, \mathbf{X})$ and the equality $\mathbf{y}' \mathbf{A} \mathbf{y} = \text{tr}(\mathbf{y} \mathbf{y}' \mathbf{A})$). The computational cost of this is dominated by the inversion of \mathbf{K} which has cost $\mathcal{O}(n^3)$

It is not possible to guarantee that the marginal likelihood doesn't suffer from **multiple local maxima**. Each of these will correspond to different interpretations of the data.

Classification with Regression Models

The examples of Gaussian processes given have been examined for regression. They may also be used for classification.

To use the output from Gaussian process regression, requires mapping from points in the range $-\infty$ to ∞ into “probabilities” in the range 0 to 1. There are a number of ways that this can be done. Initially considering just the **binary** case:

- **logistic regression**: here (considering linear regression)

$$P(y = +1|\mathbf{x}, \hat{\mathbf{w}}) = \frac{1}{1 + \exp(-\hat{\mathbf{w}}'\mathbf{x})}$$

- **probit regression**: again considering linear regression

$$P(y = +1|\mathbf{x}, \hat{\mathbf{w}}) = \int_{-\infty}^{\hat{\mathbf{w}}'\mathbf{x}} \mathcal{N}(z; 0, 1) dz$$

Using the sum to one constraint

$$P(y = -1|\mathbf{x}, \hat{\mathbf{w}}) = 1 - P(y = +1|\mathbf{x}, \hat{\mathbf{w}})$$

These appear to be noise free - in the sense that they appear in the same form as “noise-less” regression $y = f(\mathbf{x})$. It can be shown that probit regression is related to a 1/0 probability

$$P(y = +1|\mathbf{x}, \hat{\mathbf{w}}) = \begin{cases} 1, & \hat{\mathbf{w}}'\mathbf{x} + \epsilon \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

with Gaussian noise (see examples paper).

How to apply the same concepts to Gaussian Processes?

Bayesian Classification

In the same fashion as Bayesian regression Bayesian classification involves integrating (marginalising) over the model parameter distribution.

Before approaching more general Gaussian processes, consider the simple case of $f(\mathbf{x}) = \mathbf{w}'\mathbf{x}$, then classification is based on

$$P(y = +1|\mathbf{x}, \mathbf{X}, \mathbf{y}) = \int P(y = +1|\mathbf{x}, \mathbf{w})p(\mathbf{w}|\mathbf{X}, \mathbf{y})d\mathbf{w}$$

In the same fashion as prediction we need $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$.

Again using the Gaussian prior for the weights

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) \propto p(\mathbf{w})P(\mathbf{y}|\mathbf{X}, \mathbf{w})$$

The targets are conditionally independent, so

$$\log(p(\mathbf{w}|\mathbf{X}, \mathbf{y})) \propto -\frac{1}{\sigma_w^2}\mathbf{w}'\mathbf{w} + \sum_{i=1}^n \log(P(y_i|\mathbf{x}_i, \mathbf{w}))$$

Unlike the regression case this is not Gaussian distributed.
The posterior of the weights given the training data does not have a simple analytical solution.

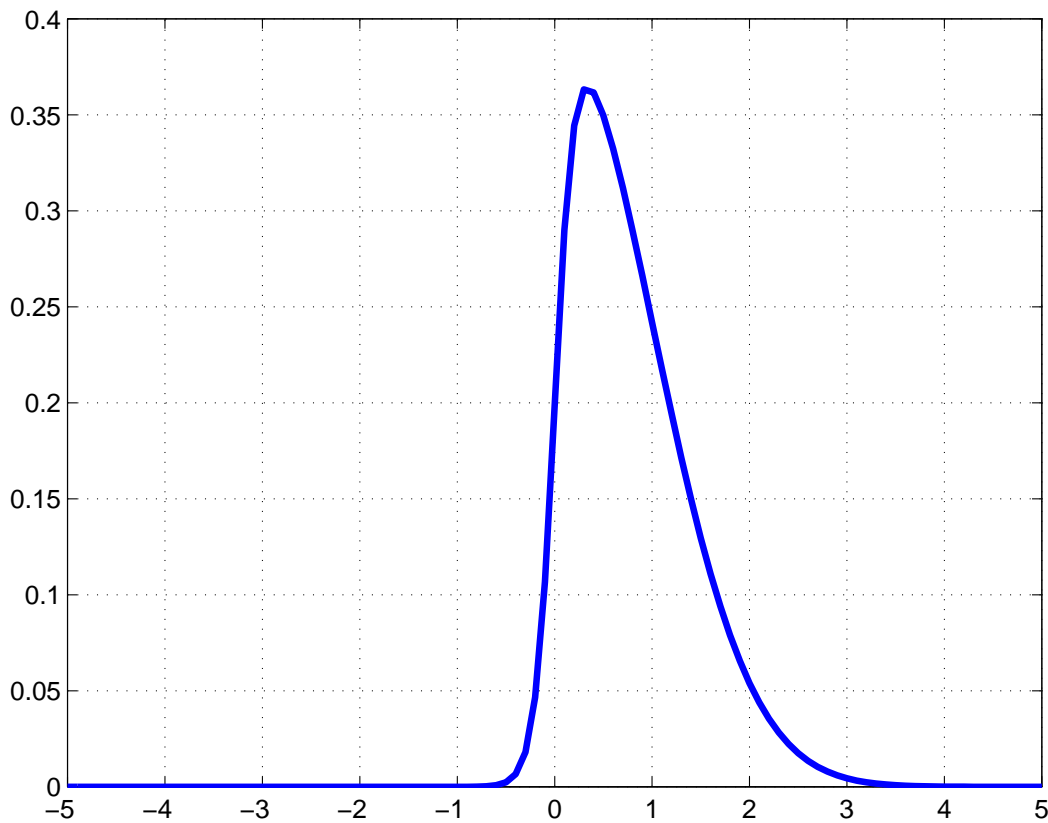
Posterior Distribution Example

Taking a simple example of what the posterior looks like for $f(x) = wx$. Here

- weight prior is Gaussian $p(w) = \mathcal{N}(w; 0, 1)$
- single training observation $y_1 = 1, x_1 = 10$

Posterior (using logistic regression) is given by

$$p(w|y_1, x_1) \propto \mathcal{N}(w; 0, 1) \left(\frac{1}{1 + \exp(-10w)} \right)$$



Weight posterior is non-Gaussian

Gaussian Process Classification

For the more general classification, the function $f(\mathbf{x})$ will be modelled using a Gaussian process. For the training data the distribution of the vector of training functions $\mathbf{f}(\cdot)$

$$\mathbf{f}(\cdot)|\mathbf{X} = \begin{bmatrix} f(\mathbf{x}_1) \\ \vdots \\ f(\mathbf{x}_n) \end{bmatrix} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}(\mathbf{X}, \mathbf{X}))$$

In contrast to the regression case even if $\mathbf{f}(\cdot)|\mathbf{X}$ is Gaussian distributed, the output targets, \mathbf{y} will **not** be Gaussian distributed.

In a similar fashion to classification with linear prediction, the outputs are conditionally independent given the systematic prediction

$$P(\mathbf{y}|\mathbf{f}(\cdot)) = \prod_{i=1}^n P(y_i|f(\mathbf{x}_i))$$

If logistic regression is used then (note targets 1 and 0)

$$\begin{aligned} \log(P(\mathbf{y}|\mathbf{f}(\cdot))) &= \sum_{i=1}^n \left[y_i \log \left(\frac{1}{1 + \exp(-f(\mathbf{x}_i))} \right) \right. \\ &\quad \left. + (1 - y_i) \log \left(1 - \frac{1}{1 + \exp(-f(\mathbf{x}_i))} \right) \right] \end{aligned}$$

The non-linear relationship between the targets and the systematic prediction is again clear. The classification posterior is then given by

$$P(y|\mathbf{X}, \mathbf{x}, \mathbf{y}) = \int P(y|f(\mathbf{x}))p(f(\mathbf{x})|\mathbf{X}, \mathbf{y}, \mathbf{x})df(\mathbf{x})$$

Gaussian Process Classification (cont)

Now considering the general process procedure. The distributions over the functions is again assumed to be Gaussian. Binary classification proceeds in two steps

- **function posterior distribution**: first the distribution over the functions given the training data and the observation is required.

$$p(f(\mathbf{x})|\mathbf{X}, \mathbf{y}, \mathbf{x}) = \int p(f(\mathbf{x})|\mathbf{X}, \mathbf{x}, \mathbf{f}(\cdot))p(\mathbf{f}(\cdot)|\mathbf{X}, \mathbf{y})d\mathbf{f}(\cdot)$$

As the relationship between the function and the target is no longer linear

$$p(\mathbf{f}(\cdot)|\mathbf{X}, \mathbf{y}) \propto P(\mathbf{y}|\mathbf{f}(\cdot))p(\mathbf{f}(\cdot)|\mathbf{X})$$

will not be Gaussian distributed even if the function (as here) is modelled using a Gaussian process.

- **marginalise over distribution**: consider logistic regression

$$P(y = +1|\mathbf{x}, \mathbf{X}, \mathbf{y}) = \int \left(\frac{1}{1 + \exp(-f(\mathbf{x}))} \right) p(f(\mathbf{x})|\mathbf{X}, \mathbf{y}, \mathbf{x})df(\mathbf{x})$$

Both these integrals are, in general, not analytical solvable:

- the marginalisation over the distribution is possible using standard numerical techniques as it is **1-dimensional** integral.
- **finding the function posterior is more complicated.**

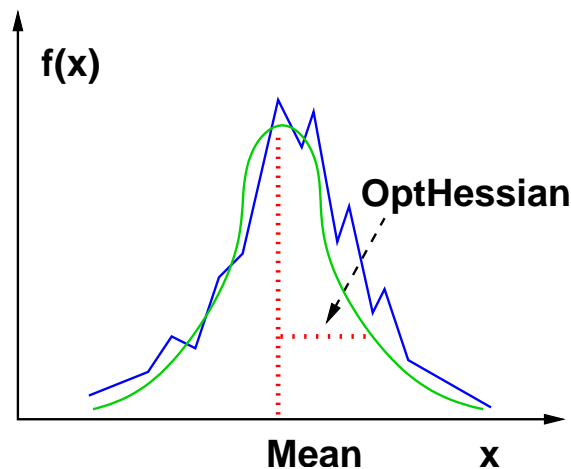
Laplace's Approximation

It is necessary to find the function posterior distribution

$$p(f(\mathbf{x})|\mathbf{X}, \mathbf{y}, \mathbf{x}) = \int p(f(\mathbf{x})|\mathbf{X}, \mathbf{x}, \mathbf{f}(\cdot))p(\mathbf{f}(\cdot)|\mathbf{X}, \mathbf{y})d\mathbf{f}(\cdot)$$

As the function, $f(\mathbf{x})$ is modelled using a Gaussian process then $p(f(\mathbf{x})|\mathbf{X}, \mathbf{x}, \mathbf{f}(\cdot))$ is Gaussian distributed.

A general approximation technique to modelling posterior distributions is to use **Laplace's approximation**.



Here a Gaussian distribution is used:

$$p(\mathbf{f}(\cdot)|\mathbf{X}, \mathbf{y}) \approx q(\mathbf{f}(\cdot)|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{f}(\cdot); \boldsymbol{\mu}_f, \boldsymbol{\Sigma}_f)$$

where

$$\begin{aligned} \boldsymbol{\mu}_f &= \arg \max_{\mathbf{f}} \{p(\mathbf{f}(\cdot)|\mathbf{X}, \mathbf{y})\} \\ \boldsymbol{\Sigma}_f^{-1} &= -\nabla^2 \log(p(\mathbf{f}(\cdot)|\mathbf{X}, \mathbf{y})) \Big|_{\mathbf{f}(\cdot)=\boldsymbol{\mu}_f} \end{aligned}$$

($\boldsymbol{\Sigma}_f^{-1}$ is the **Hessian** evaluated at $\boldsymbol{\mu}_f$)

Laplace's Approximation (cont)

Applying Laplace's approximation to Gaussian processes

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

As the denominator is not a function of $\mathbf{f}(\cdot)$ then need to find

$$\max_{\mathbf{f}(\cdot)} \{\log(P(\mathbf{y}|\mathbf{f}(\cdot))) + \log(p(\mathbf{f}(\cdot)|\mathbf{X}))\} = \max_{\mathbf{f}(\cdot)} \{\mathcal{F}(\mathbf{f}(\cdot))\}$$

This may be expressed as

$$\max_{\mathbf{f}(\cdot)} \left\{ \log(P(\mathbf{y}|\mathbf{f}(\cdot))) - \frac{1}{2}\mathbf{f}(\cdot)\mathbf{K}^{-1}\mathbf{f}(\cdot) - \frac{1}{2}\log(|\mathbf{K}|) \right\}$$

($\mathbf{K} = \mathbf{K}(\mathbf{X}, \mathbf{X})$). Differentiating this yields

$$\begin{aligned} \nabla \mathcal{F}(\mathbf{f}(\cdot)) &= \nabla \log(P(\mathbf{y}|\mathbf{f}(\cdot))) - \mathbf{K}^{-1}\mathbf{f}(\cdot) \\ \nabla^2 \mathcal{F}(\mathbf{f}(\cdot)) &= \nabla^2 \log(P(\mathbf{y}|\mathbf{f}(\cdot))) - \mathbf{K}^{-1} \end{aligned}$$

Note each of the observations are assumed independent. Taking the example of logistic regression it is simple to show (and used for logistic regression/classification training previously)

$$\frac{\partial}{\partial f(\mathbf{x}_i)} \log(P(y_i|f(\mathbf{x}_i))) = y_i - \frac{1}{1 + \exp(-f(\mathbf{x}_i))}$$

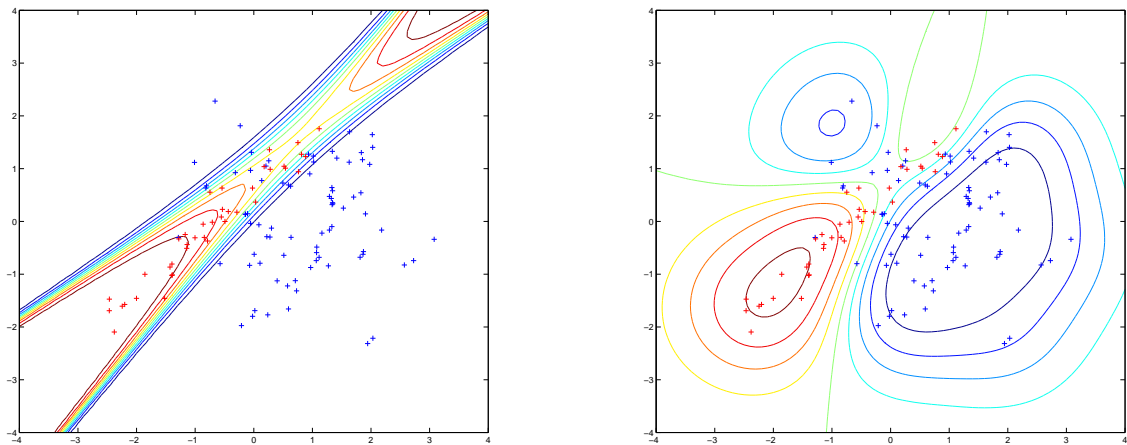
The maximum can be simply found using Newton's method for example.

Classification Example

Data generated from two Gaussian Distributions

$$\mathbf{x}^{(1)} \sim \mathcal{N} \left(\begin{bmatrix} 0.75 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) \quad \mathbf{x}^{(2)} \sim \mathcal{N} \left(\begin{bmatrix} -0.75 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 0.95 \\ 0.95 & 1 \end{bmatrix} \right)$$

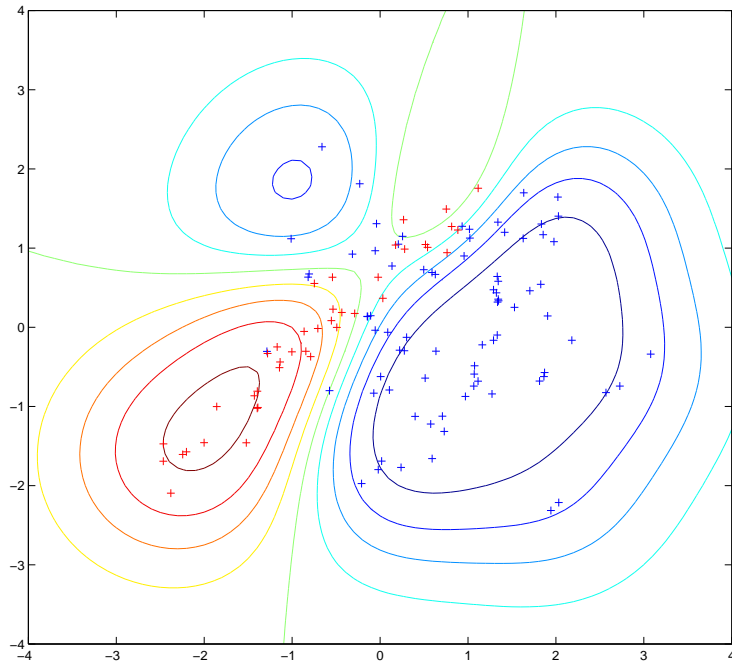
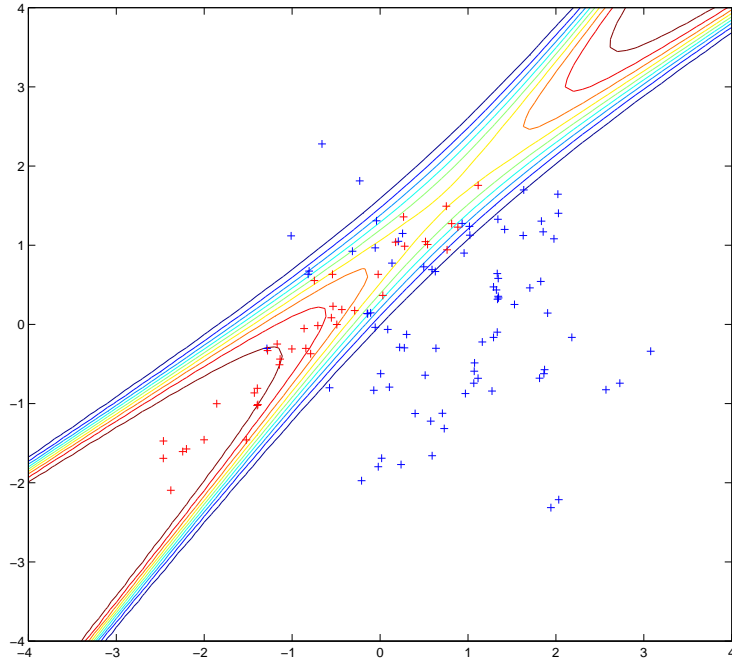
Given the two class-conditional PDFs (and assuming equal priors) it is possible to define decision boundaries/class posteriors (left diagram).



The right diagram shows an application of Laplace's approximation with a squared exponential covariance functions with

$$\sigma_1 = 1, \quad \alpha = 1$$

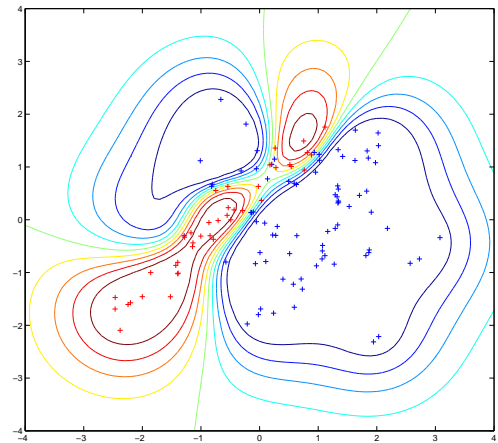
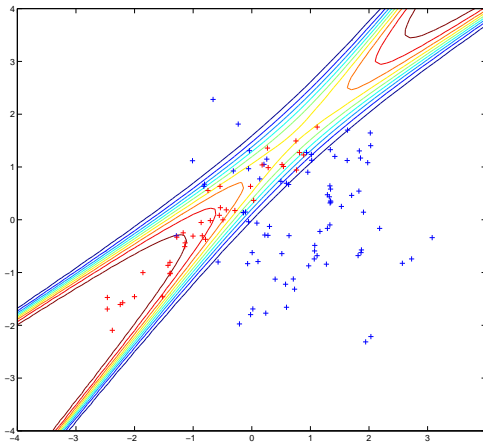
- posteriors differ, but similar in high data density regions
- tend to posteriors of 0.5 as we move away from the data



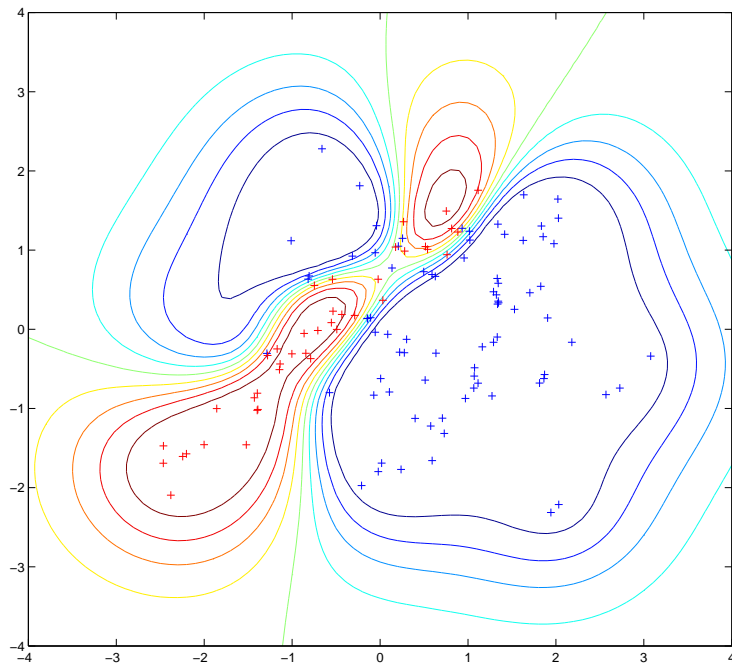
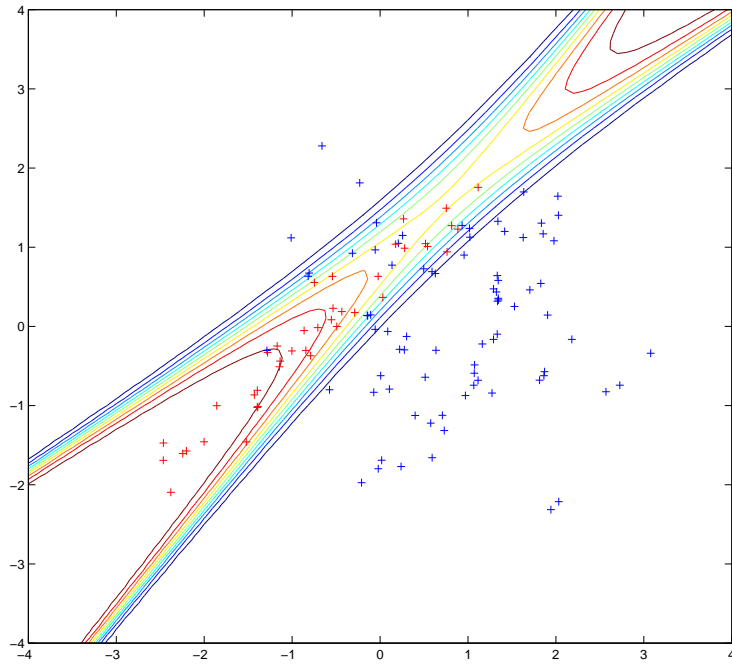
Classification Hyper-parameter Estimation

Rather than using the arbitrary hyper-parameters, these can be trained by maximising the likelihood of the outputs (using the Laplace's approximation). this yields

$$\sigma_1 = 1.280, \quad \alpha = 14.440$$



- length scale moves only a small amount
- signal variance increased - more “extreme” probabilities



Summary

The last two lectures have examined the use of **Gaussian Processes** for prediction/regression and classification. In particular:

- linear/basis function regression
- Bayesian regression
- definition of a Gaussian process
- covariance functions and kernels
- Gaussian processes for regression
- learning hyper-parameters
- Gaussian processes for classification
- Laplace's approximation