

Gaussian Processes

MLAI: Week 10

Neil D. Lawrence

Department of Computer Science
Sheffield University

2nd December 2014

Outline

Bayesian Polynomials

Distributions over Functions

Covariance from Basis Functions

Basis Function Representations

Covariance from Basis Functions

Basis Function Representations

GP Limitations

Revisit Olympics Data

- ▶ Use Bayesian approach on olympics data with polynomials.
- ▶ Choose a prior $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \alpha \mathbf{I})$ with $\alpha = 1$.
- ▶ Choose noise variance $\sigma^2 = 0.01$

Sampling the Prior

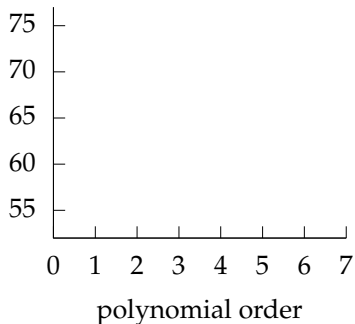
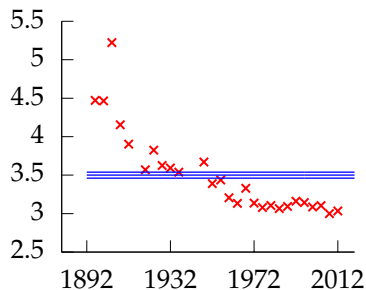
- ▶ Always useful to perform a ‘sanity check’ and sample from the prior before observing the data.
- ▶ Since $\mathbf{y} = \Phi\mathbf{w} + \epsilon$ just need to sample

$$w \sim \mathcal{N}(0, \alpha)$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2)$$

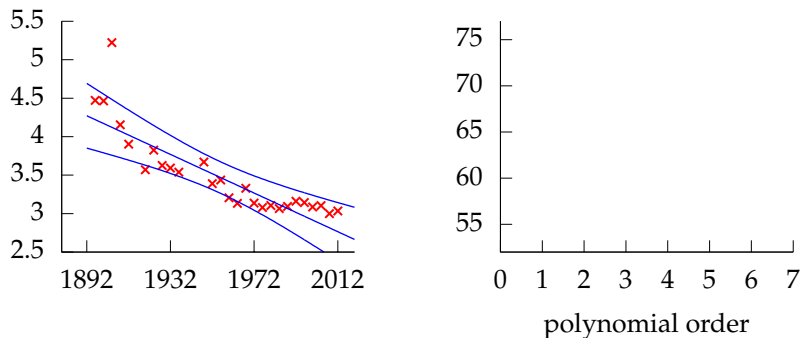
with $\alpha = 1$ and $\epsilon = 0.01$.

Polynomial Fits to Olympics Data



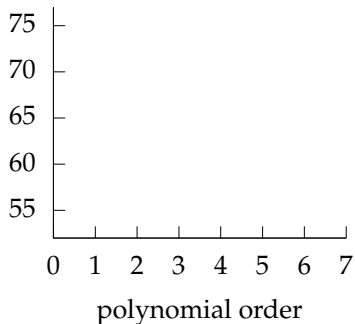
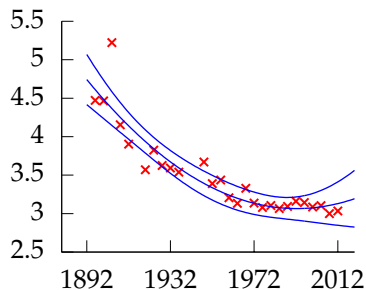
Left: fit to data, Right: marginal log likelihood. Polynomial order 0, model error 29.757, $\sigma^2 = 0.286$, $\sigma = 0.535$.

Polynomial Fits to Olympics Data



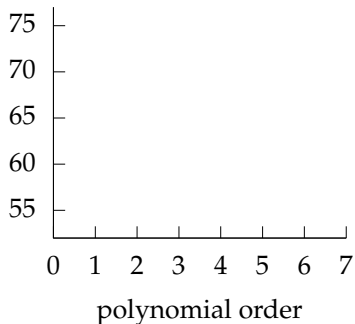
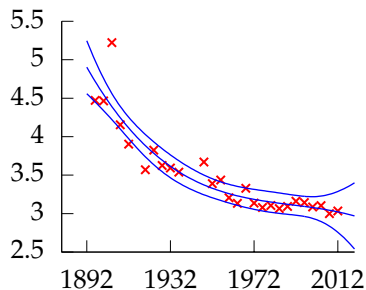
Left: fit to data, Right: marginal log likelihood. Polynomial order 1, model error 14.942, $\sigma^2 = 0.0749$, $\sigma = 0.274$.

Polynomial Fits to Olympics Data



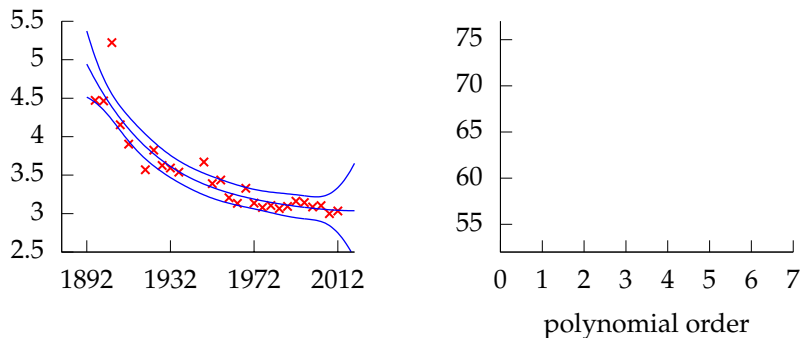
Left: fit to data, Right: marginal log likelihood. Polynomial order 2, model error 9.7206, $\sigma^2 = 0.0427$, $\sigma = 0.207$.

Polynomial Fits to Olympics Data



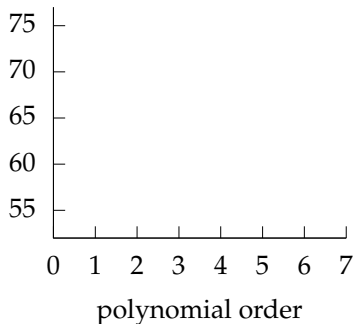
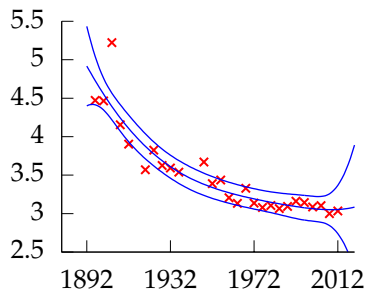
Left: fit to data, Right: marginal log likelihood. Polynomial order 3, model error 10.416, $\sigma^2 = 0.0402$, $\sigma = 0.200$.

Polynomial Fits to Olympics Data



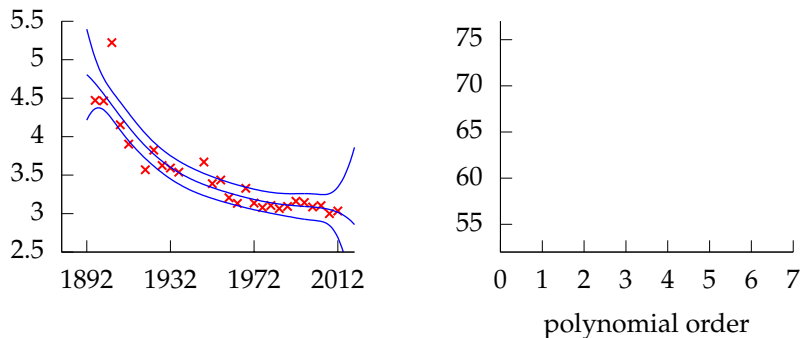
Left: fit to data, Right: marginal log likelihood. Polynomial order 4, model error 11.34, $\sigma^2 = 0.0401$, $\sigma = 0.200$.

Polynomial Fits to Olympics Data



Left: fit to data, Right: marginal log likelihood. Polynomial order 5, model error 11.986, $\sigma^2 = 0.0399$, $\sigma = 0.200$.

Polynomial Fits to Olympics Data

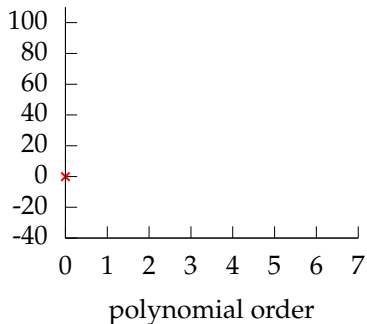
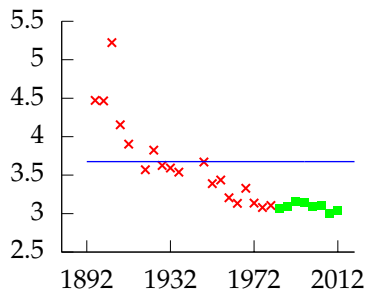


Left: fit to data, Right: marginal log likelihood. Polynomial order 6, model error 12.369, $\sigma^2 = 0.0384$, $\sigma = 0.196$.

Model Fit

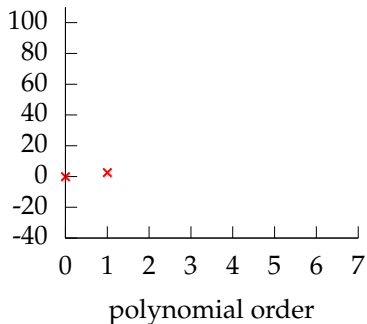
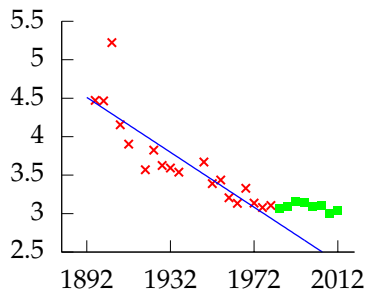
- ▶ Marginal likelihood doesn't always increase as model order increases.
- ▶ Bayesian model always has 2 parameters, regardless of how many basis functions (and here we didn't even fit them).
- ▶ Maximum likelihood model over fits through increasing number of parameters.
- ▶ Revisit maximum likelihood solution with validation set.

Recall: Validation Set for Maximum Likelihood



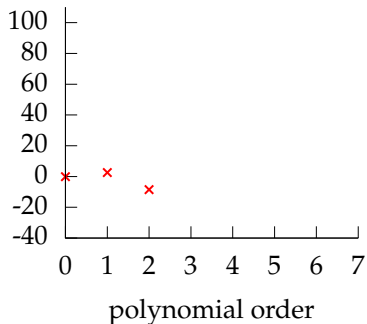
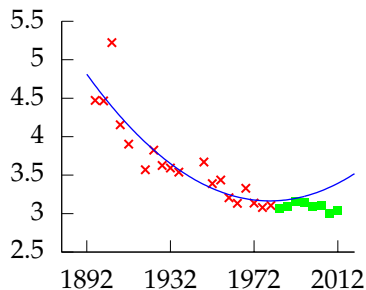
Left: fit to data, Right: model error. Polynomial order 0, training error -1.8774, validation error -0.13132, $\sigma^2 = 0.302$, $\sigma = 0.549$.

Recall: Validation Set for Maximum Likelihood



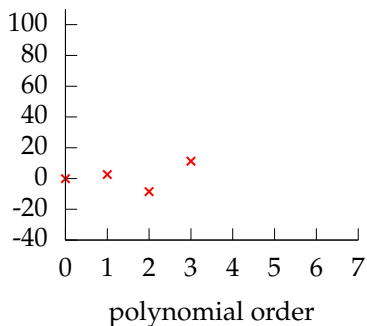
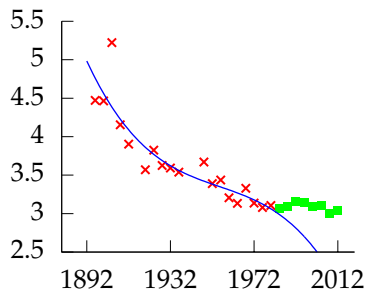
Left: fit to data, Right: model error. Polynomial order 1, training error -15.325, validation error 2.5863, $\sigma^2 = 0.0733$, $\sigma = 0.271$.

Recall: Validation Set for Maximum Likelihood



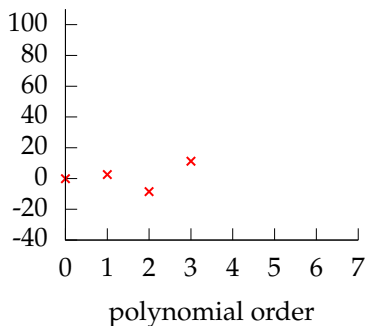
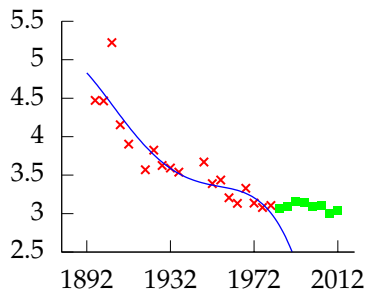
Left: fit to data, Right: model error. Polynomial order 2, training error -17.579, validation error -8.4831, $\sigma^2 = 0.0578$, $\sigma = 0.240$.

Recall: Validation Set for Maximum Likelihood



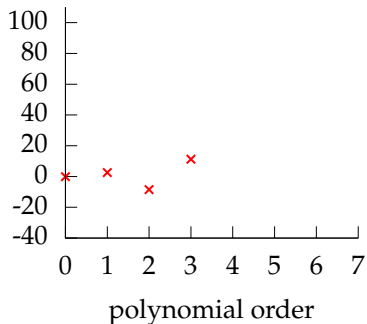
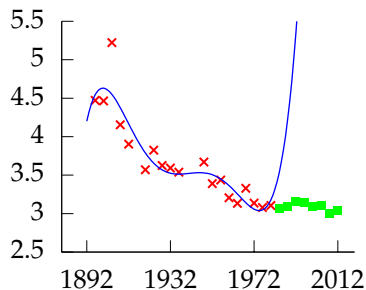
Left: fit to data, Right: model error. Polynomial order 3, training error -18.064, validation error 11.27, $\sigma^2 = 0.0549$, $\sigma = 0.234$.

Recall: Validation Set for Maximum Likelihood



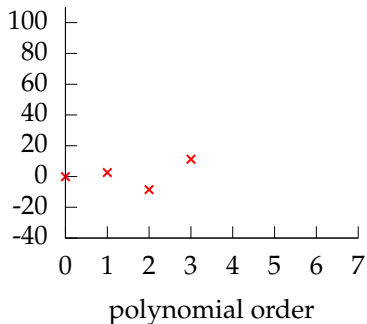
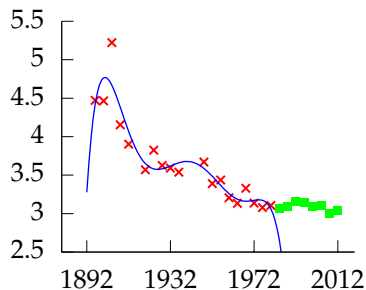
Left: fit to data, Right: model error. Polynomial order 4, training error -18.245, validation error 232.92, $\sigma^2 = 0.0539$, $\sigma = 0.232$.

Recall: Validation Set for Maximum Likelihood



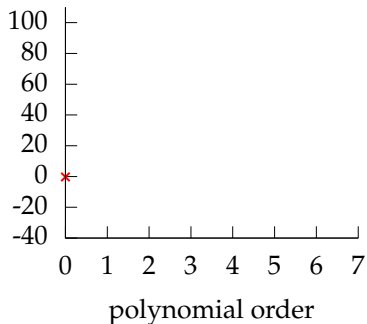
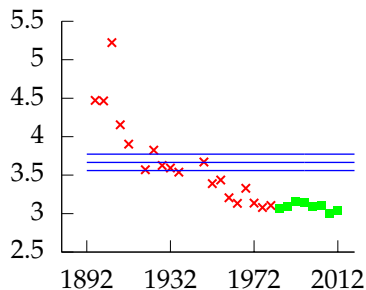
Left: fit to data, Right: model error. Polynomial order 5, training error -20.471, validation error 9898.1, $\sigma^2 = 0.0426$, $\sigma = 0.207$.

Recall: Validation Set for Maximum Likelihood



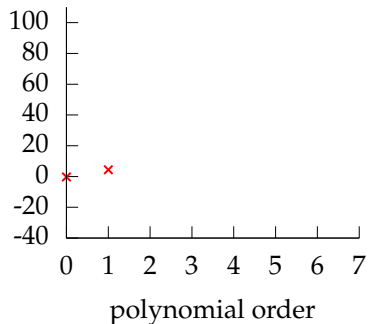
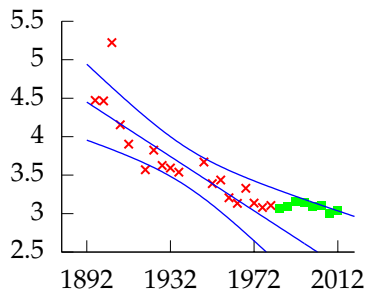
Left: fit to data, Right: model error. Polynomial order 6, training error -22.881, validation error 67775, $\sigma^2 = 0.0331$, $\sigma = 0.182$.

Validation Set



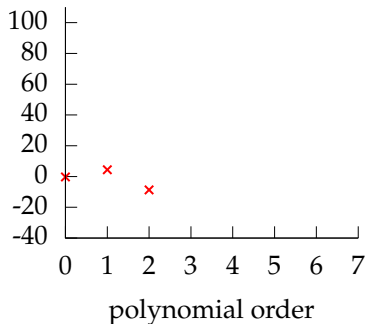
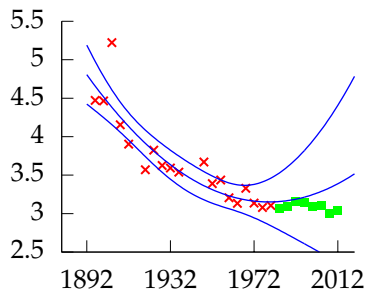
Left: fit to data, Right: model error. Polynomial order 0, training error 29.757, validation error -0.29243, $\sigma^2 = 0.302$, $\sigma = 0.550$.

Validation Set



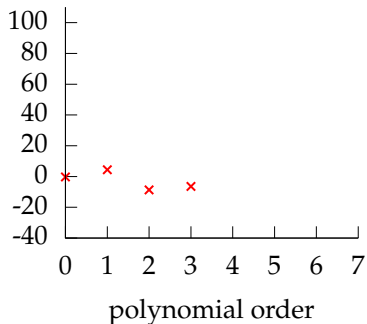
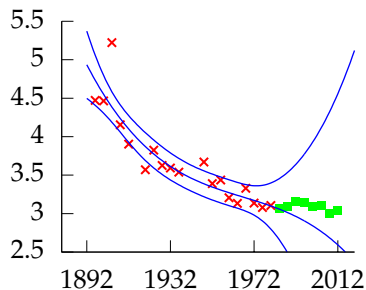
Left: fit to data, Right: model error. Polynomial order 1, training error 14.942, validation error 4.4027, $\sigma^2 = 0.0762$, $\sigma = 0.276$.

Validation Set



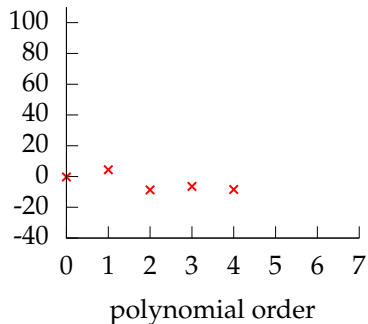
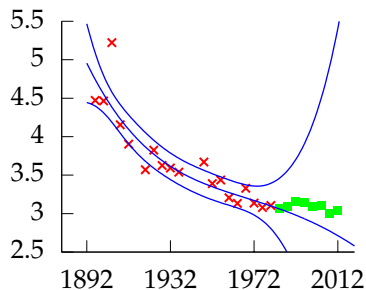
Left: fit to data, Right: model error. Polynomial order 2, training error 9.7206, validation error -8.6623, $\sigma^2 = 0.0580$, $\sigma = 0.241$.

Validation Set



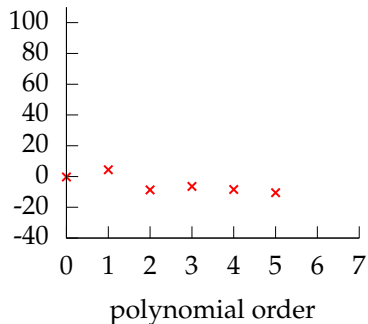
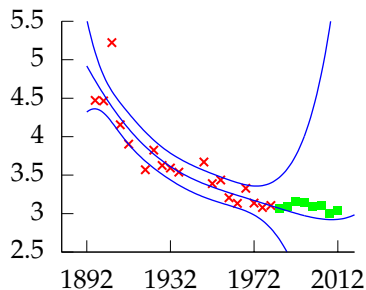
Left: fit to data, Right: model error. Polynomial order 3, training error 10.416, validation error -6.4726, $\sigma^2 = 0.0555$, $\sigma = 0.236$.

Validation Set



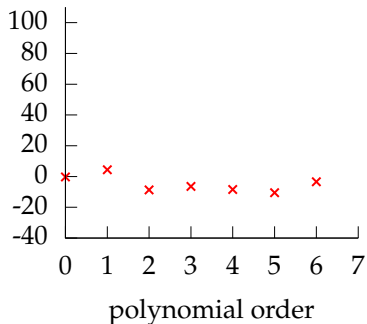
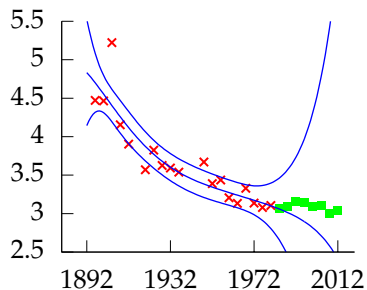
Left: fit to data, Right: model error. Polynomial order 4, training error 11.34, validation error -8.431, $\sigma^2 = 0.0555$, $\sigma = 0.236$.

Validation Set



Left: fit to data, Right: model error. Polynomial order 5, training error 11.986, validation error -10.483, $\sigma^2 = 0.0551$, $\sigma = 0.235$.

Validation Set



Left: fit to data, Right: model error. Polynomial order 6, training error 12.369, validation error -3.3823, $\sigma^2 = 0.0537$, $\sigma = 0.232$.

Regularized Mean

- ▶ Validation fit here based on mean solution for \mathbf{w} only.
- ▶ For Bayesian solution

$$\boldsymbol{\mu}_w = \left[\sigma^{-2} \boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \alpha^{-1} \mathbf{I} \right]^{-1} \sigma^{-2} \boldsymbol{\Phi}^\top \mathbf{y}$$

instead of

$$\mathbf{w}^* = \left[\boldsymbol{\Phi}^\top \boldsymbol{\Phi} \right]^{-1} \boldsymbol{\Phi}^\top \mathbf{y}$$

- ▶ Two are equivalent when $\alpha \rightarrow \infty$.
- ▶ Equivalent to a prior for \mathbf{w} with infinite variance.
- ▶ In other cases $\alpha \mathbf{I}$ *regularizes* the system (keeps parameters smaller).

Sampling the Posterior

- ▶ Now check samples by extracting \mathbf{w} from the *posterior*.
- ▶ Now for $\mathbf{y} = \mathbf{\Phi}\mathbf{w} + \epsilon$ need

$$w \sim \mathcal{N}(\boldsymbol{\mu}_w, \mathbf{C}_w)$$

$$\text{with } \mathbf{C}_w = [\sigma^{-2}\mathbf{\Phi}^\top\mathbf{\Phi} + \alpha^{-1}\mathbf{I}]^{-1} \text{ and } \boldsymbol{\mu}_w = \mathbf{C}_w\sigma^{-2}\mathbf{\Phi}^\top\mathbf{y}$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2)$$

with $\alpha = 1$ and $\epsilon = 0.01$.

Marginal Likelihood

- ▶ The marginal likelihood can also be computed, it has the form:

$$p(\mathbf{y}|\mathbf{X}, \sigma^2, \alpha) = \frac{1}{(2\pi)^{\frac{n}{2}} |\mathbf{K}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}\right)$$

where $\mathbf{K} = \alpha \mathbf{\Phi} \mathbf{\Phi}^\top + \sigma^2 \mathbf{I}$.

- ▶ So it is a zero mean n -dimensional Gaussian with covariance matrix \mathbf{K} .

Computing the Expected Output

- ▶ Given the posterior for the parameters, how can we compute the expected output at a given location?
- ▶ Output of model at location \mathbf{x}_i is given by

$$f(\mathbf{x}_i; \mathbf{w}) = \phi_i^\top \mathbf{w}$$

- ▶ We want the expected output under the posterior density, $p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \sigma^2, \alpha)$.
- ▶ Mean of mapping function will be given by

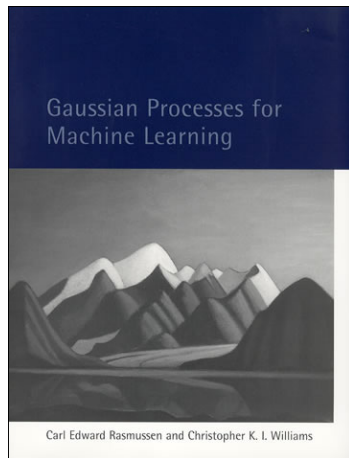
$$\begin{aligned}\langle f(\mathbf{x}_i; \mathbf{w}) \rangle_{p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \sigma^2, \alpha)} &= \phi_i^\top \langle \mathbf{w} \rangle_{p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \sigma^2, \alpha)} \\ &= \phi_i^\top \boldsymbol{\mu}_w\end{aligned}$$

Variance of Expected Output

- ▶ Variance of model at location \mathbf{x}_i is given by

$$\begin{aligned}\text{var}(f(\mathbf{x}_i; \mathbf{w})) &= \langle (f(\mathbf{x}_i; \mathbf{w}))^2 \rangle - \langle f(\mathbf{x}_i; \mathbf{w}) \rangle^2 \\ &= \boldsymbol{\phi}_i^\top \langle \mathbf{w}\mathbf{w}^\top \rangle \boldsymbol{\phi}_i - \boldsymbol{\phi}_i^\top \langle \mathbf{w} \rangle \langle \mathbf{w} \rangle^\top \boldsymbol{\phi}_i \\ &= \boldsymbol{\phi}_i^\top \mathbf{C}_w \boldsymbol{\phi}_i\end{aligned}$$

where all these expectations are taken under the posterior density, $p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \sigma^2, \alpha)$.



Rasmussen and Williams (2006)

Outline

Bayesian Polynomials

Distributions over Functions

Covariance from Basis Functions

Basis Function Representations

Covariance from Basis Functions

Basis Function Representations

GP Limitations

Outline

Bayesian Polynomials

Distributions over Functions

Covariance from Basis Functions

Basis Function Representations

Covariance from Basis Functions

Basis Function Representations

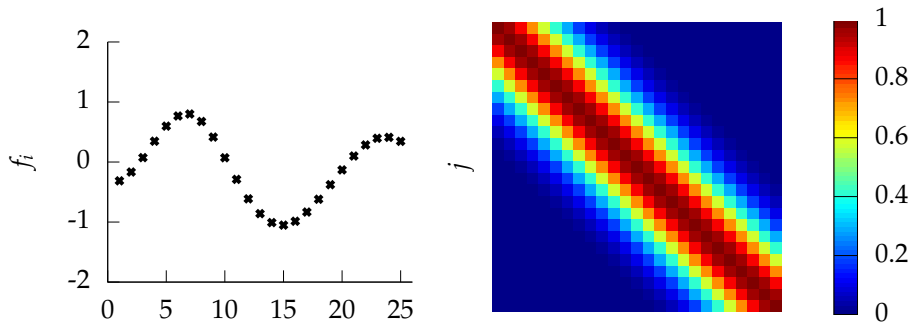
GP Limitations

Sampling a Function

Multi-variate Gaussians

- ▶ We will consider a Gaussian with a particular structure of covariance matrix.
- ▶ Generate a single sample from this 25 dimensional Gaussian distribution, $\mathbf{f} = [f_1, f_2 \dots f_{25}]$.
- ▶ We will plot these points against their index.

Gaussian Distribution Sample

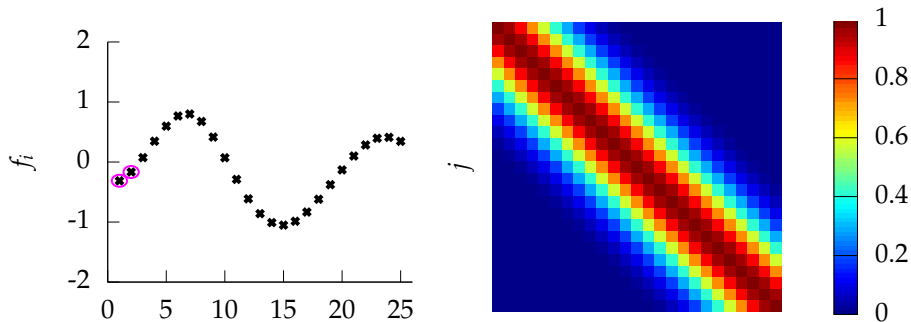


(a) A 25 dimensional correlated random variable (values plotted against index)

(b) colormap showing correlations between dimensions.

Figure : A sample from a 25 dimensional Gaussian distribution.

Gaussian Distribution Sample

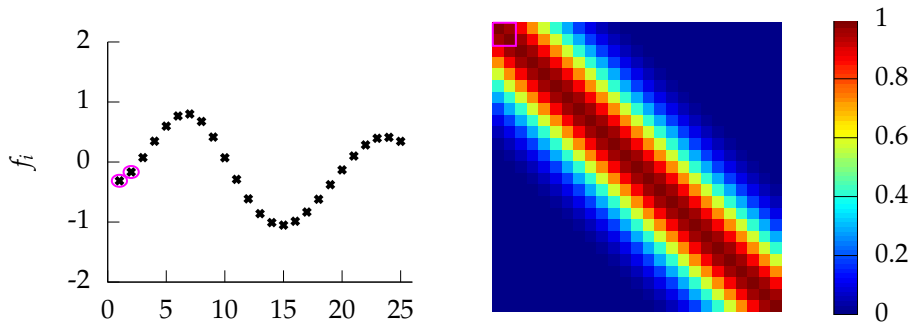


(a) A 25 dimensional correlated random variable (values plotted against index)

(b) colormap showing correlations between dimensions.

Figure : A sample from a 25 dimensional Gaussian distribution.

Gaussian Distribution Sample

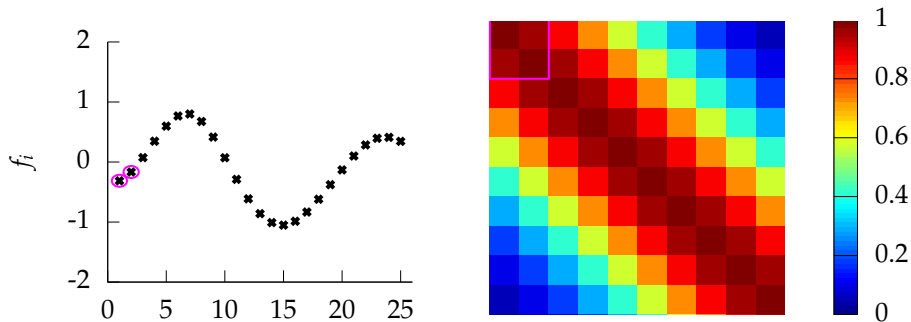


(a) A 25 dimensional correlated random variable (values plotted against index)

(b) colormap showing correlations between dimensions.

Figure : A sample from a 25 dimensional Gaussian distribution.

Gaussian Distribution Sample

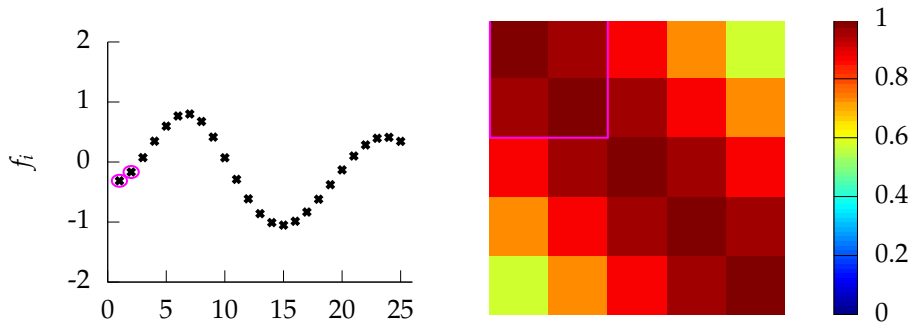


(a) A 25 dimensional correlated random variable (values plotted against index)

(b) colormap showing correlations between dimensions.

Figure : A sample from a 25 dimensional Gaussian distribution.

Gaussian Distribution Sample

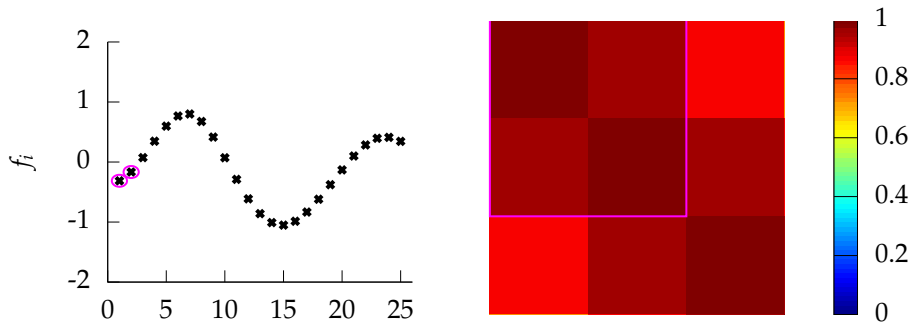


(a) A 25 dimensional correlated random variable (values plotted against index)

(b) colormap showing correlations between dimensions.

Figure : A sample from a 25 dimensional Gaussian distribution.

Gaussian Distribution Sample

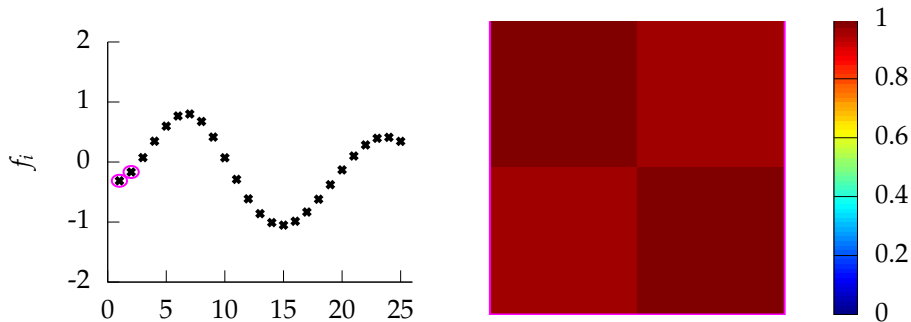


(a) A 25 dimensional correlated random variable (values plotted against index)

(b) colormap showing correlations between dimensions.

Figure : A sample from a 25 dimensional Gaussian distribution.

Gaussian Distribution Sample

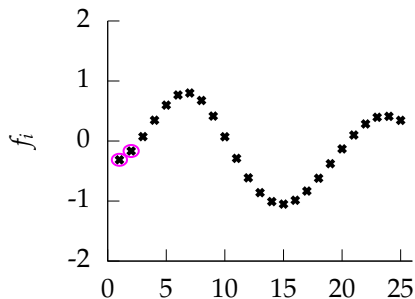


(a) A 25 dimensional correlated random variable (values plotted against index)

(b) colormap showing correlations between dimensions.

Figure : A sample from a 25 dimensional Gaussian distribution.

Gaussian Distribution Sample



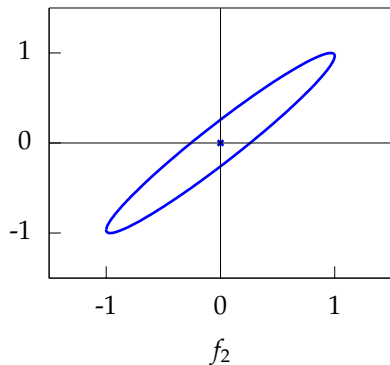
(a) A 25 dimensional correlated random variable (values plotted against index)



(b) correlation between f_1 and f_2 .

Figure : A sample from a 25 dimensional Gaussian distribution.

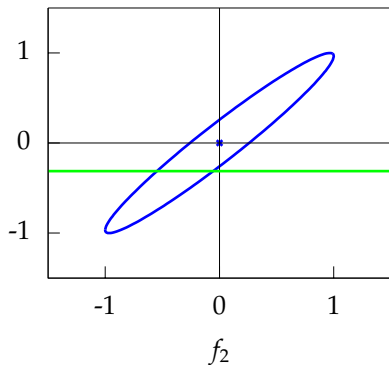
Prediction of f_2 from f_1



$$\begin{bmatrix} 1 & 0.96587 \\ 0.96587 & 1 \end{bmatrix}$$

- ▶ The single contour of the Gaussian density represents the joint distribution, $p(f_1, f_2)$.

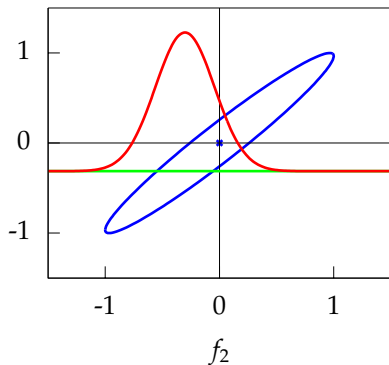
Prediction of f_2 from f_1



$$\begin{bmatrix} 1 & 0.96587 \\ 0.96587 & 1 \end{bmatrix}$$

- ▶ The single contour of the Gaussian density represents the **joint distribution**, $p(f_1, f_2)$.
- ▶ We observe that $f_1 = -0.313$.

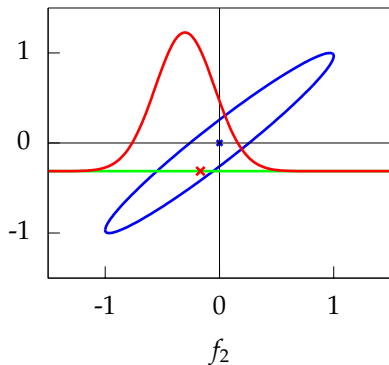
Prediction of f_2 from f_1



$$\begin{bmatrix} 1 & 0.96587 \\ 0.96587 & 1 \end{bmatrix}$$

- ▶ The single contour of the Gaussian density represents the **joint distribution**, $p(f_1, f_2)$.
- ▶ We observe that $f_1 = -0.313$.
- ▶ Conditional density: $p(f_2|f_1 = -0.313)$.

Prediction of f_2 from f_1



$$\begin{bmatrix} 1 & 0.96587 \\ 0.96587 & 1 \end{bmatrix}$$

- ▶ The single contour of the Gaussian density represents the **joint distribution**, $p(f_1, f_2)$.
- ▶ We observe that $f_1 = -0.313$.
- ▶ Conditional density: $p(f_2|f_1 = -0.313)$.

Prediction with Correlated Gaussians

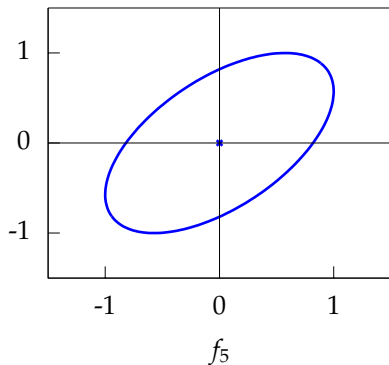
- ▶ Prediction of f_2 from f_1 requires *conditional density*.
- ▶ Conditional density is *also* Gaussian.

$$p(f_2|f_1) = \mathcal{N}\left(f_2 \mid \frac{k_{1,2}}{k_{1,1}} f_1, k_{2,2} - \frac{k_{1,2}^2}{k_{1,1}}\right)$$

where covariance of joint density is given by

$$\mathbf{K} = \begin{bmatrix} k_{1,1} & k_{1,2} \\ k_{2,1} & k_{2,2} \end{bmatrix}$$

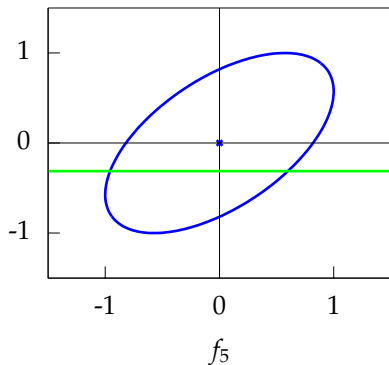
Prediction of f_5 from f_1



$$\begin{bmatrix} 1 & 0.57375 \\ 0.57375 & 1 \end{bmatrix}$$

- ▶ The single contour of the Gaussian density represents the joint distribution, $p(f_1, f_5)$.

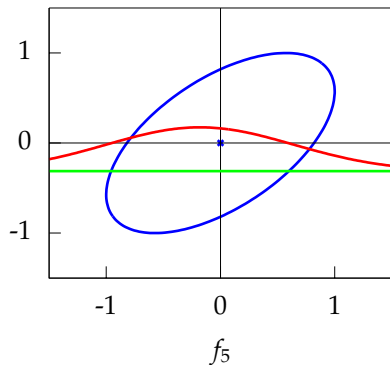
Prediction of f_5 from f_1



$$\begin{bmatrix} 1 & 0.57375 \\ 0.57375 & 1 \end{bmatrix}$$

- ▶ The single contour of the Gaussian density represents the **joint distribution**, $p(f_1, f_5)$.
- ▶ We observe that $f_1 = -0.313$.

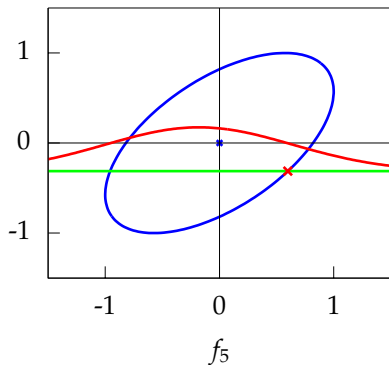
Prediction of f_5 from f_1



$$\begin{bmatrix} 1 & 0.57375 \\ 0.57375 & 1 \end{bmatrix}$$

- ▶ The single contour of the Gaussian density represents the **joint distribution**, $p(f_1, f_5)$.
- ▶ We observe that $f_1 = -0.313$.
- ▶ Conditional density: $p(f_5|f_1 = -0.313)$.

Prediction of f_5 from f_1



$$\begin{bmatrix} 1 & 0.57375 \\ 0.57375 & 1 \end{bmatrix}$$

- ▶ The single contour of the Gaussian density represents the **joint distribution**, $p(f_1, f_5)$.
- ▶ We observe that $f_1 = -0.313$.
- ▶ Conditional density: $p(f_5|f_1 = -0.313)$.

Prediction with Correlated Gaussians

- ▶ Prediction of \mathbf{f}_* from \mathbf{f} requires multivariate *conditional density*.
- ▶ Multivariate conditional density is *also* Gaussian.

$$p(\mathbf{f}_*|\mathbf{f}) = \mathcal{N}\left(\mathbf{f}_*|\mathbf{K}_{*,f}\mathbf{K}_{f,f}^{-1}\mathbf{f}, \mathbf{K}_{*,*} - \mathbf{K}_{*,f}\mathbf{K}_{f,f}^{-1}\mathbf{K}_{f,*}\right)$$

- ▶ Here covariance of joint density is given by

$$\mathbf{K} = \begin{bmatrix} \mathbf{K}_{f,f} & \mathbf{K}_{*,f} \\ \mathbf{K}_{f,*} & \mathbf{K}_{*,*} \end{bmatrix}$$

Prediction with Correlated Gaussians

- ▶ Prediction of \mathbf{f}_* from \mathbf{f} requires multivariate *conditional density*.
- ▶ Multivariate conditional density is *also* Gaussian.

$$p(\mathbf{f}_*|\mathbf{f}) = \mathcal{N}(\mathbf{f}_*|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\boldsymbol{\mu} = \mathbf{K}_{*,\mathbf{f}}\mathbf{K}_{\mathbf{f},\mathbf{f}}^{-1}\mathbf{f}$$

$$\boldsymbol{\Sigma} = \mathbf{K}_{*,*} - \mathbf{K}_{*,\mathbf{f}}\mathbf{K}_{\mathbf{f},\mathbf{f}}^{-1}\mathbf{K}_{\mathbf{f},*}$$

- ▶ Here covariance of joint density is given by

$$\mathbf{K} = \begin{bmatrix} \mathbf{K}_{\mathbf{f},\mathbf{f}} & \mathbf{K}_{*,\mathbf{f}} \\ \mathbf{K}_{\mathbf{f},*} & \mathbf{K}_{*,*} \end{bmatrix}$$

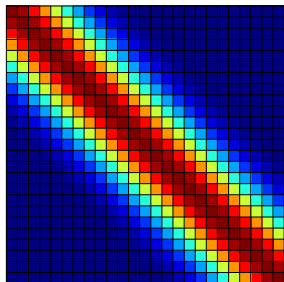
Covariance Functions

Where did this covariance matrix come from?

Exponentiated Quadratic Kernel Function (RBF, Squared Exponential, Gaussian)

$$k(\mathbf{x}, \mathbf{x}') = \alpha \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell^2}\right)$$

- ▶ Covariance matrix is built using the *inputs* to the function \mathbf{x} .
- ▶ For the example above it was based on Euclidean distance.
- ▶ The covariance function is also known as a kernel.



Covariance Functions

Where did this covariance matrix come from?

Exponentiated Quadratic Kernel Function (RBF, Squared Exponential, Gaussian)

$$k(\mathbf{x}, \mathbf{x}') = \alpha \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell^2}\right)$$

- ▶ Covariance matrix is built using the *inputs* to the function \mathbf{x} .
- ▶ For the example above it was based on Euclidean distance.
- ▶ The covariance function is also known as a kernel.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = -3.0, x_1 = -3.0$$

$$k_{1,1} = 1.00 \times \exp\left(-\frac{(-3.0 - -3.0)^2}{2 \times 2.00^2}\right)$$

$x_1 = -3.0, x_2 = 1.20, \text{ and } x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = -3.0, x_1 = -3.0$$

$$k_{1,1} = 1.00 \times \exp\left(-\frac{(-3.0 - -3.0)^2}{2 \times 2.00^2}\right)$$

$$\begin{bmatrix} 1.00 \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20, \text{ and } x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.20, x_1 = -3.0$$

$$k_{2,1} = 1.00 \times \exp\left(-\frac{(1.20 - (-3.0))^2}{2 \times 2.00^2}\right)$$

$$\begin{bmatrix} 1.00 \\ \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20,$ and $x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.20, x_1 = -3.0$$

$$k_{2,1} = 1.00 \times \exp\left(-\frac{(1.20 - (-3.0))^2}{2 \times 2.00^2}\right)$$

$$\begin{bmatrix} 1.00 \\ 0.110 \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20,$ and $x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.20, x_1 = -3.0$$

$$k_{2,1} = 1.00 \times \exp\left(-\frac{(1.20 - (-3.0))^2}{2 \times 2.00^2}\right)$$

$$\begin{bmatrix} 1.00 & 0.110 \\ 0.110 & \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20,$ and $x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = 1.20, x_2 = 1.20$$

$$k_{2,2} = 1.00 \times \exp\left(-\frac{(1.20-1.20)^2}{2 \times 2.00^2}\right)$$

$$\begin{bmatrix} 1.00 & 0.110 \\ 0.110 & 1.00 \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20, \text{ and } x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = 1.20, x_2 = 1.20$$

$$k_{2,2} = 1.00 \times \exp\left(-\frac{(1.20-1.20)^2}{2 \times 2.00^2}\right)$$

$$\begin{bmatrix} 1.00 & 0.110 \\ 0.110 & 1.00 \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20, \text{ and } x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_1 = -3.0$$

$$k_{3,1} = 1.00 \times \exp\left(-\frac{(1.40 - (-3.0))^2}{2 \times 2.00^2}\right)$$

$$\begin{bmatrix} 1.00 & 0.110 \\ 0.110 & 1.00 \end{bmatrix}$$

$x_1 = -3.0$, $x_2 = 1.20$, and $x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_1 = -3.0$$

$$k_{3,1} = 1.00 \times \exp\left(-\frac{(1.40 - (-3.0))^2}{2 \times 2.00^2}\right)$$

$$\begin{bmatrix} 1.00 & 0.110 \\ 0.110 & 1.00 \\ 0.0889 & & \end{bmatrix}$$

$x_1 = -3.0$, $x_2 = 1.20$, and $x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_1 = -3.0$$

$$k_{3,1} = 1.00 \times \exp\left(-\frac{(1.40 - (-3.0))^2}{2 \times 2.00^2}\right)$$

$$\begin{bmatrix} 1.00 & 0.110 & 0.0889 \\ 0.110 & 1.00 & \\ 0.0889 & & \end{bmatrix}$$

$x_1 = -3.0$, $x_2 = 1.20$, and $x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_2 = 1.20$$

$$k_{3,2} = 1.00 \times \exp\left(-\frac{(1.40-1.20)^2}{2 \times 2.00^2}\right)$$

$$\begin{bmatrix} 1.00 & 0.110 & 0.0889 \\ 0.110 & 1.00 & \\ 0.0889 & & \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20, \text{ and } x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_2 = 1.20$$

$$k_{3,2} = 1.00 \times \exp\left(-\frac{(1.40-1.20)^2}{2 \times 2.00^2}\right)$$

$$\begin{bmatrix} 1.00 & 0.110 & 0.0889 \\ 0.110 & 1.00 & \\ 0.0889 & 0.995 & \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20, \text{ and } x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_2 = 1.20$$

$$k_{3,2} = 1.00 \times \exp\left(-\frac{(1.40-1.20)^2}{2 \times 2.00^2}\right)$$

$$\begin{bmatrix} 1.00 & 0.110 & 0.0889 \\ 0.110 & 1.00 & 0.995 \\ 0.0889 & 0.995 & 1.00 \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20, \text{ and } x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = 1.40, x_2 = 1.40$$

$$k_{3,3} = 1.00 \times \exp\left(-\frac{(1.40-1.40)^2}{2 \times 2.00^2}\right)$$

$$\begin{bmatrix} 1.00 & 0.110 & 0.0889 \\ 0.110 & 1.00 & 0.995 \\ 0.0889 & 0.995 & 1.00 \end{bmatrix}$$

$x_1 = -3.0$, $x_2 = 1.20$, and $x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_3 = 1.40$$

$$k_{3,3} = 1.00 \times \exp\left(-\frac{(1.40-1.40)^2}{2 \times 2.00^2}\right)$$

$$\begin{bmatrix} 1.00 & 0.110 & 0.0889 \\ 0.110 & 1.00 & 0.995 \\ 0.0889 & 0.995 & 1.00 \end{bmatrix}$$

$x_1 = -3.0$, $x_2 = 1.20$, and $x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

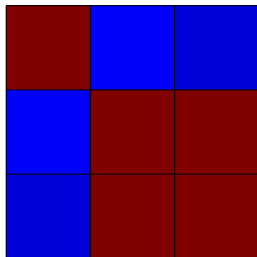
Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_3 = 1.40$$

$$k_{3,3} = 1.00 \times \exp\left(-\frac{(1.40-1.40)^2}{2 \times 2.00^2}\right)$$



$x_1 = -3.0$, $x_2 = 1.20$, and $x_3 = 1.40$ with $\ell = 2.00$ and $\alpha = 1.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = -3, x_1 = -3$$

$$k_{1,1} = 1.0 \times \exp\left(-\frac{(-3--3)^2}{2 \times 2.0^2}\right)$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = -3, x_1 = -3$$

$$k_{1,1} = 1.0 \times \exp\left(-\frac{(-3 - -3)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 \\ \vdots \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.2, x_1 = -3$$

$$k_{2,1} = 1.0 \times \exp\left(-\frac{(1.2 - (-3))^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.2, x_1 = -3$$

$$k_{2,1} = 1.0 \times \exp\left(-\frac{(1.2 - (-3))^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 \\ 0.11 \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.2, x_1 = -3$$

$$k_{2,1} = 1.0 \times \exp\left(-\frac{(1.2 - (-3))^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 \\ 0.11 & \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = 1.2, x_2 = 1.2$$

$$k_{2,2} = 1.0 \times \exp\left(-\frac{(1.2-1.2)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 \\ 0.11 & \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.2, x_2 = 1.2$$

$$k_{2,2} = 1.0 \times \exp\left(-\frac{(1.2-1.2)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 \\ 0.11 & 1.0 \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.4, x_1 = -3$$

$$k_{3,1} = 1.0 \times \exp\left(-\frac{(1.4 - (-3))^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 \\ 0.11 & 1.0 \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.4, x_1 = -3$$

$$k_{3,1} = 1.0 \times \exp\left(-\frac{(1.4 - (-3))^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 \\ 0.11 & 1.0 \\ 0.089 & & \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.4, x_1 = -3$$

$$k_{3,1} = 1.0 \times \exp\left(-\frac{(1.4 - (-3))^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 \\ 0.11 & 1.0 & \\ 0.089 & & \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.4, x_2 = 1.2$$

$$k_{3,2} = 1.0 \times \exp\left(-\frac{(1.4-1.2)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 \\ 0.11 & 1.0 & \\ 0.089 & & \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.4, x_2 = 1.2$$

$$k_{3,2} = 1.0 \times \exp\left(-\frac{(1.4-1.2)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 \\ 0.11 & 1.0 & \\ 0.089 & 1.0 & \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.4, x_2 = 1.2$$

$$k_{3,2} = 1.0 \times \exp\left(-\frac{(1.4-1.2)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 \\ 0.11 & 1.0 & 1.0 \\ 0.089 & 1.0 & \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.4, x_3 = 1.4$$

$$k_{3,3} = 1.0 \times \exp\left(-\frac{(1.4-1.4)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 \\ 0.11 & 1.0 & 1.0 \\ 0.089 & 1.0 & \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.4, x_3 = 1.4$$

$$k_{3,3} = 1.0 \times \exp\left(-\frac{(1.4-1.4)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 \\ 0.11 & 1.0 & 1.0 \\ 0.089 & 1.0 & 1.0 \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_1 = -3$$

$$k_{4,1} = 1.0 \times \exp\left(-\frac{(2.0 - (-3))^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 \\ 0.11 & 1.0 & 1.0 \\ 0.089 & 1.0 & 1.0 \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_1 = -3$$

$$k_{4,1} = 1.0 \times \exp\left(-\frac{(2.0 - (-3))^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 & \\ 0.11 & 1.0 & 1.0 & \\ 0.089 & 1.0 & 1.0 & \\ 0.044 & & & \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_1 = -3$$

$$k_{4,1} = 1.0 \times \exp\left(-\frac{(2.0 - (-3))^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 & 0.044 \\ 0.11 & 1.0 & 1.0 & \\ 0.089 & 1.0 & 1.0 & \\ 0.044 & & & \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_2 = 1.2$$

$$k_{4,2} = 1.0 \times \exp\left(-\frac{(2.0-1.2)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 & 0.044 \\ 0.11 & 1.0 & 1.0 & \\ 0.089 & 1.0 & 1.0 & \\ 0.044 & & & \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_2 = 1.2$$

$$k_{4,2} = 1.0 \times \exp\left(-\frac{(2.0-1.2)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 & 0.044 \\ 0.11 & 1.0 & 1.0 & \\ 0.089 & 1.0 & 1.0 & \\ 0.044 & 0.92 & & \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_2 = 1.2$$

$$k_{4,2} = 1.0 \times \exp\left(-\frac{(2.0-1.2)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 & 0.044 \\ 0.11 & 1.0 & 1.0 & 0.92 \\ 0.089 & 1.0 & 1.0 & \\ 0.044 & 0.92 & & \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_3 = 1.4$$

$$k_{4,3} = 1.0 \times \exp\left(-\frac{(2.0-1.4)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 & 0.044 \\ 0.11 & 1.0 & 1.0 & 0.92 \\ 0.089 & 1.0 & 1.0 & \\ 0.044 & 0.92 & & \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_3 = 1.4$$

$$k_{4,3} = 1.0 \times \exp\left(-\frac{(2.0-1.4)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 & 0.044 \\ 0.11 & 1.0 & 1.0 & 0.92 \\ 0.089 & 1.0 & 1.0 & \\ 0.044 & 0.92 & \boxed{0.96} & \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_3 = 1.4$$

$$k_{4,3} = 1.0 \times \exp\left(-\frac{(2.0-1.4)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 & 0.044 \\ 0.11 & 1.0 & 1.0 & 0.92 \\ 0.089 & 1.0 & 1.0 & 0.96 \\ 0.044 & 0.92 & 0.96 & 1.0 \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_4 = 2.0$$

$$k_{4,4} = 1.0 \times \exp\left(-\frac{(2.0-2.0)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 & 0.044 \\ 0.11 & 1.0 & 1.0 & 0.92 \\ 0.089 & 1.0 & 1.0 & 0.96 \\ 0.044 & 0.92 & 0.96 & 1.0 \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_4 = 2.0$$

$$k_{4,4} = 1.0 \times \exp\left(-\frac{(2.0-2.0)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 & 0.044 \\ 0.11 & 1.0 & 1.0 & 0.92 \\ 0.089 & 1.0 & 1.0 & 0.96 \\ 0.044 & 0.92 & 0.96 & 1.0 \end{bmatrix}$$

$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

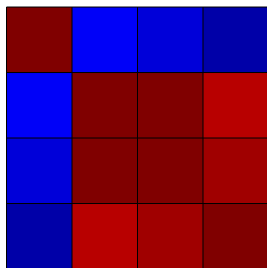
Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_4 = 2.0$$

$$k_{4,4} = 1.0 \times \exp\left(-\frac{(2.0-2.0)^2}{2 \times 2.0^2}\right)$$



$x_1 = -3, x_2 = 1.2, x_3 = 1.4,$ and $x_4 = 2.0$ with $\ell = 2.0$ and $\alpha = 1.0$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = -3.0, x_1 = -3.0$$

$$k_{1,1} = 4.00 \times \exp\left(-\frac{(-3.0 - -3.0)^2}{2 \times 5.00^2}\right)$$

$x_1 = -3.0, x_2 = 1.20, \text{ and } x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = -3.0, x_1 = -3.0$$

$$k_{1,1} = 4.00 \times \exp\left(-\frac{(-3.0 - -3.0)^2}{2 \times 5.00^2}\right)$$

$$\begin{bmatrix} 4.00 \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20, \text{ and } x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.20, x_1 = -3.0$$

$$k_{2,1} = 4.00 \times \exp\left(-\frac{(1.20 - (-3.0))^2}{2 \times 5.00^2}\right)$$

$$\begin{bmatrix} 4.00 \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20$, and $x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.20, x_1 = -3.0$$

$$k_{2,1} = 4.00 \times \exp\left(-\frac{(1.20 - (-3.0))^2}{2 \times 5.00^2}\right)$$

$$\begin{bmatrix} & & \\ & 4.00 & \\ & 2.81 & \\ & & & \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20,$ and $x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.20, x_1 = -3.0$$

$$k_{2,1} = 4.00 \times \exp\left(-\frac{(1.20 - (-3.0))^2}{2 \times 5.00^2}\right)$$

$$\begin{bmatrix} 4.00 & 2.81 \\ 2.81 & \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20,$ and $x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = 1.20, x_2 = 1.20$$

$$k_{2,2} = 4.00 \times \exp\left(-\frac{(1.20-1.20)^2}{2 \times 5.00^2}\right)$$

$$\begin{bmatrix} 4.00 & 2.81 \\ 2.81 & \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20, \text{ and } x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = 1.20, x_2 = 1.20$$

$$k_{2,2} = 4.00 \times \exp\left(-\frac{(1.20-1.20)^2}{2 \times 5.00^2}\right)$$

$$\begin{bmatrix} 4.00 & 2.81 \\ 2.81 & 4.00 \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20, \text{ and } x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_1 = -3.0$$

$$k_{3,1} = 4.00 \times \exp\left(-\frac{(1.40 - (-3.0))^2}{2 \times 5.00^2}\right)$$

$$\begin{bmatrix} 4.00 & 2.81 \\ 2.81 & 4.00 \end{bmatrix}$$

$x_1 = -3.0$, $x_2 = 1.20$, and $x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_1 = -3.0$$

$$k_{3,1} = 4.00 \times \exp\left(-\frac{(1.40 - (-3.0))^2}{2 \times 5.00^2}\right)$$

$$\begin{bmatrix} 4.00 & 2.81 \\ 2.81 & 4.00 \\ 2.72 & \end{bmatrix}$$

$x_1 = -3.0$, $x_2 = 1.20$, and $x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_1 = -3.0$$

$$k_{3,1} = 4.00 \times \exp\left(-\frac{(1.40 - (-3.0))^2}{2 \times 5.00^2}\right)$$

$$\begin{bmatrix} 4.00 & 2.81 & 2.72 \\ 2.81 & 4.00 & \\ 2.72 & & \end{bmatrix}$$

$x_1 = -3.0$, $x_2 = 1.20$, and $x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_2 = 1.20$$

$$k_{3,2} = 4.00 \times \exp\left(-\frac{(1.40-1.20)^2}{2 \times 5.00^2}\right)$$

$$\begin{bmatrix} 4.00 & 2.81 & 2.72 \\ 2.81 & 4.00 & \\ 2.72 & & \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20, \text{ and } x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_2 = 1.20$$

$$k_{3,2} = 4.00 \times \exp\left(-\frac{(1.40-1.20)^2}{2 \times 5.00^2}\right)$$

$$\begin{bmatrix} 4.00 & 2.81 & 2.72 \\ 2.81 & 4.00 & \\ 2.72 & 4.00 & \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20, \text{ and } x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_2 = 1.20$$

$$k_{3,2} = 4.00 \times \exp\left(-\frac{(1.40-1.20)^2}{2 \times 5.00^2}\right)$$

$$\begin{bmatrix} 4.00 & 2.81 & 2.72 \\ 2.81 & 4.00 & 4.00 \\ 2.72 & 4.00 & \end{bmatrix}$$

$x_1 = -3.0, x_2 = 1.20, \text{ and } x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = 1.40, x_2 = 1.40$$

$$k_{3,3} = 4.00 \times \exp\left(-\frac{(1.40-1.40)^2}{2 \times 5.00^2}\right)$$

$$\begin{bmatrix} 4.00 & 2.81 & 2.72 \\ 2.81 & 4.00 & 4.00 \\ 2.72 & 4.00 & 4.00 \end{bmatrix}$$

$x_1 = -3.0$, $x_2 = 1.20$, and $x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_3 = 1.40$$

$$k_{3,3} = 4.00 \times \exp\left(-\frac{(1.40-1.40)^2}{2 \times 5.00^2}\right)$$

$$\begin{bmatrix} 4.00 & 2.81 & 2.72 \\ 2.81 & 4.00 & 4.00 \\ 2.72 & 4.00 & 4.00 \end{bmatrix}$$

$x_1 = -3.0$, $x_2 = 1.20$, and $x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

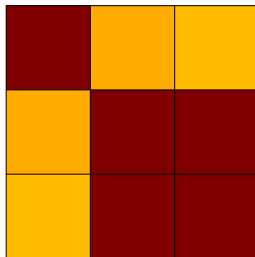
Covariance Functions

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_3 = 1.40$$

$$k_{3,3} = 4.00 \times \exp\left(-\frac{(1.40-1.40)^2}{2 \times 5.00^2}\right)$$



$x_1 = -3.0$, $x_2 = 1.20$, and $x_3 = 1.40$ with $\ell = 5.00$ and $\alpha = 4.00$.

Outline

Bayesian Polynomials

Distributions over Functions

Covariance from Basis Functions

Basis Function Representations

Covariance from Basis Functions

Basis Function Representations

GP Limitations

Basis Function Form

Radial basis functions commonly have the form

$$\phi_k(\mathbf{x}_i) = \exp\left(-\frac{|\mathbf{x}_i - \boldsymbol{\mu}_k|^2}{2\ell^2}\right).$$

- ▶ Basis function maps data into a “feature space” in which a linear sum is a non linear function.

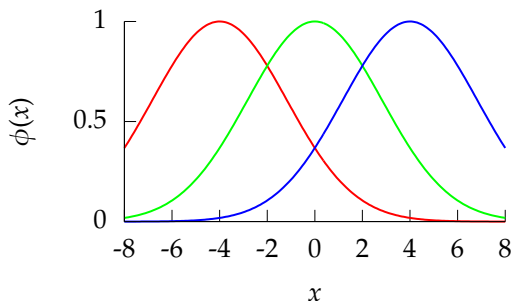


Figure : A set of radial basis functions with width $\ell = 2$ and location parameters $\boldsymbol{\mu} = [-4 \ 0 \ 4]^T$.

Basis Function Representations

- ▶ Represent a function by a linear sum over a basis,

$$f(\mathbf{x}_{i,:}; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_{i,:}), \quad (1)$$

- ▶ Here: m basis functions and $\phi_k(\cdot)$ is k th basis function and

$$\mathbf{w} = [w_1, \dots, w_m]^\top .$$

- ▶ For standard linear model: $\phi_k(\mathbf{x}_{i,:}) = x_{i,k}$.

Random Functions

Functions derived
using:

$$f(x) = \sum_{k=1}^m w_k \phi_k(x),$$

where \mathbf{W} is sampled
from a Gaussian
density,

$$w_k \sim \mathcal{N}(0, \alpha).$$

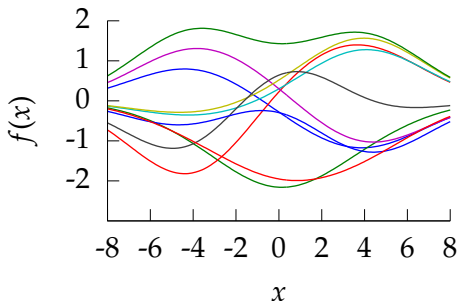


Figure : Functions sampled using the basis set from figure 4. Each line is a separate sample, generated by a weighted sum of the basis set. The weights, \mathbf{w} are sampled from a Gaussian density with variance $\alpha = 1$.

Outline

Bayesian Polynomials

Distributions over Functions

Covariance from Basis Functions

Basis Function Representations

Covariance from Basis Functions

Basis Function Representations

GP Limitations

Basis Function Form

Radial basis functions commonly have the form

$$\phi_k(\mathbf{x}_i) = \exp\left(-\frac{|\mathbf{x}_i - \boldsymbol{\mu}_k|^2}{2\ell^2}\right).$$

- ▶ Basis function maps data into a “feature space” in which a linear sum is a non linear function.

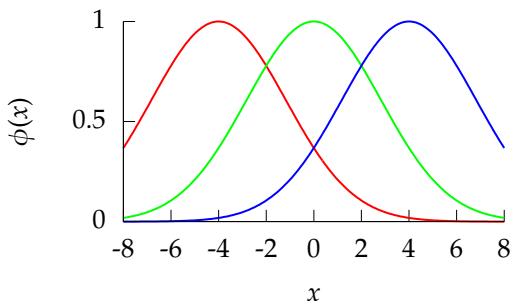


Figure : A set of radial basis functions with width $\ell = 2$ and location parameters $\boldsymbol{\mu} = [-4 \ 0 \ 4]^T$.

Basis Function Representations

- ▶ Represent a function by a linear sum over a basis,

$$f(\mathbf{x}_{i,:}; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_{i,:}), \quad (2)$$

- ▶ Here: m basis functions and $\phi_k(\cdot)$ is k th basis function and

$$\mathbf{w} = [w_1, \dots, w_m]^\top.$$

- ▶ For standard linear model: $\phi_k(\mathbf{x}_{i,:}) = x_{i,k}$.

Random Functions

Functions derived
using:

$$f(x) = \sum_{k=1}^m w_k \phi_k(x),$$

where \mathbf{W} is sampled
from a Gaussian
density,

$$w_k \sim \mathcal{N}(0, \alpha).$$

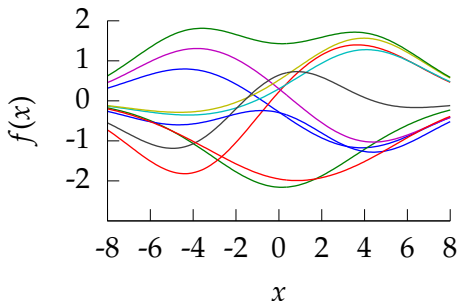


Figure : Functions sampled using the basis set from figure 4. Each line is a separate sample, generated by a weighted sum of the basis set. The weights, \mathbf{w} are sampled from a Gaussian density with variance $\alpha = 1$.

Direct Construction of Covariance Matrix

Use matrix notation to write function,

$$f(\mathbf{x}_i; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_i)$$

Direct Construction of Covariance Matrix

Use matrix notation to write function,

$$f(\mathbf{x}_i; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_i)$$

computed at training data gives a vector

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

Direct Construction of Covariance Matrix

Use matrix notation to write function,

$$f(\mathbf{x}_i; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_i)$$

computed at training data gives a vector

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

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \alpha\mathbf{I})$$

Direct Construction of Covariance Matrix

Use matrix notation to write function,

$$f(\mathbf{x}_i; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_i)$$

computed at training data gives a vector

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

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \alpha \mathbf{I})$$

\mathbf{w} and \mathbf{f} are only related by an *inner product*.

Direct Construction of Covariance Matrix

Use matrix notation to write function,

$$f(\mathbf{x}_i; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_i)$$

computed at training data gives a vector

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

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \alpha \mathbf{I})$$

\mathbf{w} and \mathbf{f} are only related by an *inner product*.

$\mathbf{\Phi} \in \mathcal{R}^{n \times p}$ is a *design matrix*

Direct Construction of Covariance Matrix

Use matrix notation to write function,

$$f(\mathbf{x}_i; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_i)$$

computed at training data gives a vector

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

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \alpha \mathbf{I})$$

\mathbf{w} and \mathbf{f} are only related by an *inner product*.

$\mathbf{\Phi} \in \mathcal{R}^{n \times p}$ is a *design matrix*

$\mathbf{\Phi}$ is fixed and non-stochastic for a given training set.

Direct Construction of Covariance Matrix

Use matrix notation to write function,

$$f(\mathbf{x}_i; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_i)$$

computed at training data gives a vector

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

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \alpha \mathbf{I})$$

\mathbf{w} and \mathbf{f} are only related by an *inner product*.

$\mathbf{\Phi} \in \mathcal{R}^{n \times p}$ is a *design matrix*

$\mathbf{\Phi}$ is fixed and non-stochastic for a given training set.

\mathbf{f} is Gaussian distributed.

Expectations

- ▶ We have

$$\langle \mathbf{f} \rangle = \Phi \langle \mathbf{w} \rangle.$$

We use $\langle \cdot \rangle$ to denote expectations under prior distributions.

Expectations

- ▶ We have

$$\langle \mathbf{f} \rangle = \mathbf{\Phi} \langle \mathbf{w} \rangle.$$

- ▶ Prior mean of \mathbf{w} was zero giving

$$\langle \mathbf{f} \rangle = \mathbf{0}.$$

We use $\langle \cdot \rangle$ to denote expectations under prior distributions.

Expectations

- ▶ We have

$$\langle \mathbf{f} \rangle = \mathbf{\Phi} \langle \mathbf{w} \rangle.$$

- ▶ Prior mean of \mathbf{w} was zero giving

$$\langle \mathbf{f} \rangle = \mathbf{0}.$$

- ▶ Prior covariance of \mathbf{f} is

$$\mathbf{K} = \langle \mathbf{f}\mathbf{f}^\top \rangle - \langle \mathbf{f} \rangle \langle \mathbf{f} \rangle^\top$$

We use $\langle \cdot \rangle$ to denote expectations under prior distributions.

Expectations

- ▶ We have

$$\langle \mathbf{f} \rangle = \mathbf{\Phi} \langle \mathbf{w} \rangle.$$

- ▶ Prior mean of \mathbf{w} was zero giving

$$\langle \mathbf{f} \rangle = \mathbf{0}.$$

- ▶ Prior covariance of \mathbf{f} is

$$\mathbf{K} = \langle \mathbf{f}\mathbf{f}^\top \rangle - \langle \mathbf{f} \rangle \langle \mathbf{f} \rangle^\top$$

$$\langle \mathbf{f}\mathbf{f}^\top \rangle = \mathbf{\Phi} \langle \mathbf{w}\mathbf{w}^\top \rangle \mathbf{\Phi}^\top,$$

giving

$$\mathbf{K} = \alpha \mathbf{\Phi} \mathbf{\Phi}^\top.$$

We use $\langle \cdot \rangle$ to denote expectations under prior distributions.

Covariance between Two Points

- ▶ The prior covariance between two points \mathbf{x}_i and \mathbf{x}_j is

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \phi: (\mathbf{x}_i)^\top \phi: (\mathbf{x}_j),$$

Covariance between Two Points

- ▶ The prior covariance between two points \mathbf{x}_i and \mathbf{x}_j is

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j),$$

or in sum notation

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \sum_{k=1}^m \phi_k(\mathbf{x}_i) \phi_k(\mathbf{x}_j)$$

Covariance between Two Points

- ▶ The prior covariance between two points \mathbf{x}_i and \mathbf{x}_j is

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j),$$

or in sum notation

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \sum_{k=1}^m \phi_k(\mathbf{x}_i) \phi_k(\mathbf{x}_j)$$

- ▶ For the radial basis used this gives

Covariance between Two Points

- ▶ The prior covariance between two points \mathbf{x}_i and \mathbf{x}_j is

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \phi: (\mathbf{x}_i)^\top \phi: (\mathbf{x}_j),$$

or in sum notation

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \sum_{k=1}^m \phi_k(\mathbf{x}_i) \phi_k(\mathbf{x}_j)$$

- ▶ For the radial basis used this gives

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \sum_{k=1}^m \exp\left(-\frac{|\mathbf{x}_i - \boldsymbol{\mu}_k|^2 + |\mathbf{x}_j - \boldsymbol{\mu}_k|^2}{2\ell^2}\right).$$

Selecting Number and Location of Basis

- ▶ Need to choose
 1. location of centers

Selecting Number and Location of Basis

- ▶ Need to choose
 1. location of centers
 2. number of basis functions

Selecting Number and Location of Basis

- ▶ Need to choose
 1. location of centers
 2. number of basis functions
- ▶ Consider uniform spacing over a region:

$$k(x_i, x_j) = \alpha' \Delta\mu \sum_{k=1}^m \exp\left(-\frac{x_i^2 + x_j^2 - 2\mu_k(x_i + x_j) + 2\mu_k^2}{2\ell^2}\right),$$

Restrict analysis to 1-D input, x .

Uniform Basis Functions

- ▶ Set each center location to

$$\mu_k = a + \Delta\mu \cdot (k - 1).$$

Uniform Basis Functions

- ▶ Set each center location to

$$\mu_k = a + \Delta\mu \cdot (k - 1).$$

- ▶ Specify the basis functions in terms of their indices,

$$k(x_i, x_j) = \alpha' \Delta\mu \sum_{k=0}^{m-1} \exp\left(-\frac{x_i^2 + x_j^2}{2\ell^2} - \frac{2(a + \Delta\mu \cdot k)(x_i + x_j) + 2(a + \Delta\mu \cdot k)^2}{2\ell^2}\right).$$

Uniform Basis Functions

- ▶ Set each center location to

$$\mu_k = a + \Delta\mu \cdot (k - 1).$$

- ▶ Specify the basis functions in terms of their indices,

$$k(x_i, x_j) = \alpha' \Delta\mu \sum_{k=0}^{m-1} \exp\left(-\frac{x_i^2 + x_j^2}{2\ell^2} - \frac{2(a + \Delta\mu \cdot k)(x_i + x_j) + 2(a + \Delta\mu \cdot k)^2}{2\ell^2}\right).$$

- ▶ Here we've scaled variance of process by $\Delta\mu$.

Infinite Basis Functions

- ▶ Take $\mu_0 = a$ and $\mu_m = b$ so $b = a + \Delta\mu \cdot (m - 1)$.

Infinite Basis Functions

- ▶ Take $\mu_0 = a$ and $\mu_m = b$ so $b = a + \Delta\mu \cdot (m - 1)$.
- ▶ Take limit as $\Delta\mu \rightarrow 0$ so $m \rightarrow \infty$

Infinite Basis Functions

- ▶ Take $\mu_0 = a$ and $\mu_m = b$ so $b = a + \Delta\mu \cdot (m - 1)$.
- ▶ Take limit as $\Delta\mu \rightarrow 0$ so $m \rightarrow \infty$

$$k(x_i, x_j) = \alpha' \int_a^b \exp\left(-\frac{x_i^2 + x_j^2}{2\ell^2} + \frac{2\left(\mu - \frac{1}{2}(x_i + x_j)\right)^2 - \frac{1}{2}(x_i + x_j)^2}{2\ell^2}\right) d\mu,$$

where we have used $k \cdot \Delta\mu \rightarrow \mu$.

Result

- ▶ Performing the integration leads to

$$k(x_i, x_j) = \alpha' \frac{\sqrt{\pi\ell^2}}{2} \exp\left(-\frac{(x_i - x_j)^2}{4\ell^2}\right) \times \left[\operatorname{erf}\left(\frac{(b - \frac{1}{2}(x_i + x_j))}{\ell}\right) - \operatorname{erf}\left(\frac{(a - \frac{1}{2}(x_i + x_j))}{\ell}\right) \right],$$

Result

- ▶ Performing the integration leads to

$$k(x_i, x_j) = \alpha' \frac{\sqrt{\pi} \ell^2}{2} \exp\left(-\frac{(x_i - x_j)^2}{4\ell^2}\right) \times \left[\operatorname{erf}\left(\frac{(b - \frac{1}{2}(x_i + x_j))}{\ell}\right) - \operatorname{erf}\left(\frac{(a - \frac{1}{2}(x_i + x_j))}{\ell}\right) \right],$$

- ▶ Now take limit as $a \rightarrow -\infty$ and $b \rightarrow \infty$

Result

- ▶ Performing the integration leads to

$$k(x_i, x_j) = \alpha' \frac{\sqrt{\pi\ell^2}}{2} \exp\left(-\frac{(x_i - x_j)^2}{4\ell^2}\right) \times \left[\operatorname{erf}\left(\frac{(b - \frac{1}{2}(x_i + x_j))}{\ell}\right) - \operatorname{erf}\left(\frac{(a - \frac{1}{2}(x_i + x_j))}{\ell}\right) \right],$$

- ▶ Now take limit as $a \rightarrow -\infty$ and $b \rightarrow \infty$

$$k(x_i, x_j) = \alpha \exp\left(-\frac{(x_i - x_j)^2}{4\ell^2}\right).$$

where $\alpha = \alpha' \sqrt{\pi\ell^2}$.

Infinite Feature Space

- ▶ An RBF model with infinite basis functions is a Gaussian process.

Infinite Feature Space

- ▶ An RBF model with infinite basis functions is a Gaussian process.
- ▶ The covariance function is given by the exponentiated quadratic covariance function.

$$k(x_i, x_j) = \alpha \exp\left(-\frac{(x_i - x_j)^2}{4\ell^2}\right).$$

where $\alpha = \alpha' \sqrt{\pi\ell^2}$.

Infinite Feature Space

- ▶ An RBF model with infinite basis functions is a Gaussian process.
- ▶ The covariance function is the exponentiated quadratic.
- ▶ **Note:** The functional form for the covariance function and basis functions are similar.
 - ▶ this is a special case,
 - ▶ in general they are very different

Similar results can obtained for multi-dimensional input models Williams (1998); Neal (1996).

Nonparametric Gaussian Processes

- ▶ We've seen how we go from parametric to non-parametric.
- ▶ The limit implies infinite dimensional \mathbf{w} .
- ▶ Gaussian processes are generally non-parametric: combine data with covariance function to get model.
- ▶ This representation *cannot* be summarized by a parameter vector of a fixed size.

The Parametric Bottleneck

- ▶ Parametric models have a representation that does not respond to increasing training set size.
- ▶ Bayesian posterior distributions over parameters contain the information about the training data.
 - ▶ Use Bayes' rule from training data, $p(\mathbf{w}|\mathbf{y}, \mathbf{X})$,
 - ▶ Make predictions on test data

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

- ▶ \mathbf{w} becomes a bottleneck for information about the training set to pass to the test set.
- ▶ Solution: increase m so that the bottleneck is so large that it no longer presents a problem.
- ▶ How big is big enough for m ? Non-parametrics says $m \rightarrow \infty$.

The Parametric Bottleneck

- ▶ Now no longer possible to manipulate the model through the standard parametric form.

The Parametric Bottleneck

- ▶ Now no longer possible to manipulate the model through the standard parametric form.
- ▶ However, it is possible to express *parametric* as GPs:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi: (\mathbf{x}_i)^\top \phi: (\mathbf{x}_j).$$

The Parametric Bottleneck

- ▶ Now no longer possible to manipulate the model through the standard parametric form.
- ▶ However, it is possible to express *parametric* as GPs:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j).$$

- ▶ These are known as degenerate covariance matrices.

The Parametric Bottleneck

- ▶ Now no longer possible to manipulate the model through the standard parametric form.
- ▶ However, it is possible to express *parametric* as GPs:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi: (\mathbf{x}_i)^\top \phi: (\mathbf{x}_j).$$

- ▶ These are known as degenerate covariance matrices.
- ▶ Their rank is at most m , non-parametric models have full rank covariance matrices.

The Parametric Bottleneck

- ▶ Now no longer possible to manipulate the model through the standard parametric form.
- ▶ However, it is possible to express *parametric* as GPs:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j).$$

- ▶ These are known as degenerate covariance matrices.
- ▶ Their rank is at most m , non-parametric models have full rank covariance matrices.
- ▶ Most well known is the “linear kernel”, $k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$.

Making Predictions

- ▶ For non-parametrics prediction at new points \mathbf{f}_* is made by conditioning on \mathbf{f} in the joint distribution.

Making Predictions

- ▶ For non-parametrics prediction at new points \mathbf{f}_* is made by conditioning on \mathbf{f} in the joint distribution.
- ▶ In GPs this involves combining the training data with the covariance function and the mean function.

Making Predictions

- ▶ For non-parametrics prediction at new points \mathbf{f}_* is made by conditioning on \mathbf{f} in the joint distribution.
- ▶ In GPs this involves combining the training data with the covariance function and the mean function.
- ▶ Parametric is a special case when conditional prediction can be summarized in a *fixed* number of parameters.

Making Predictions

- ▶ For non-parametrics prediction at new points \mathbf{f}_* is made by conditioning on \mathbf{f} in the joint distribution.
- ▶ In GPs this involves combining the training data with the covariance function and the mean function.
- ▶ Parametric is a special case when conditional prediction can be summarized in a *fixed* number of parameters.
- ▶ Complexity of parametric model remains fixed regardless of the size of our training data set.

Making Predictions

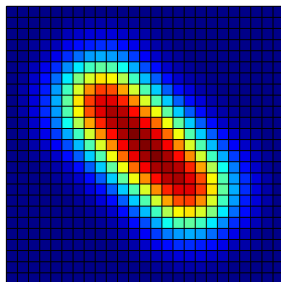
- ▶ For non-parametrics prediction at new points \mathbf{f}_* is made by conditioning on \mathbf{f} in the joint distribution.
- ▶ In GPs this involves combining the training data with the covariance function and the mean function.
- ▶ Parametric is a special case when conditional prediction can be summarized in a *fixed* number of parameters.
- ▶ Complexity of parametric model remains fixed regardless of the size of our training data set.
- ▶ For a non-parametric model the required number of parameters grows with the size of the training data.

RBF Basis Functions

$$k(\mathbf{x}, \mathbf{x}') = \alpha \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

$$\phi_i(x) = \exp\left(-\frac{\|x - \mu_i\|_2^2}{\ell^2}\right)$$

$$\mu = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}$$



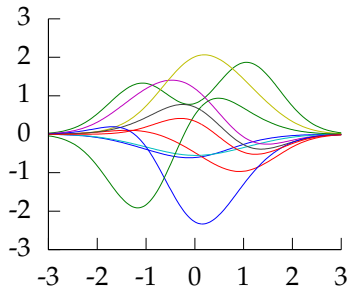
Covariance Functions

RBF Basis Functions

$$k(\mathbf{x}, \mathbf{x}') = \alpha \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

$$\phi_i(x) = \exp\left(-\frac{\|x - \mu_i\|_2^2}{\ell^2}\right)$$

$$\mu = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}$$



Covariance Functions and Mercer Kernels

- ▶ Mercer Kernels and Covariance Functions are similar.

Covariance Functions and Mercer Kernels

- ▶ Mercer Kernels and Covariance Functions are similar.
- ▶ the kernel perspective does not make a probabilistic interpretation of the covariance function.

Covariance Functions and Mercer Kernels

- ▶ Mercer Kernels and Covariance Functions are similar.
- ▶ the kernel perspective does not make a probabilistic interpretation of the covariance function.
- ▶ Algorithms can be simpler, but probabilistic interpretation is crucial for kernel parameter optimization.

Gaussian Process Interpolation

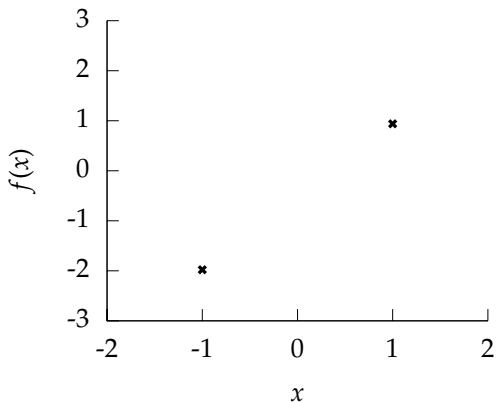


Figure : Real example: BACCO (see *e.g.* (Oakley and O'Hagan, 2002)). Interpolation through outputs from slow computer simulations (*e.g.* atmospheric carbon levels).

Gaussian Process Interpolation

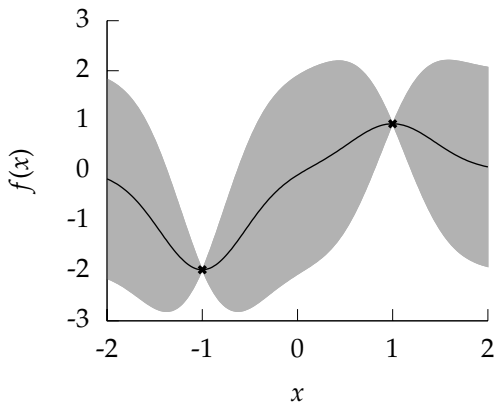


Figure : Real example: BACCO (see *e.g.* (Oakley and O'Hagan, 2002)). Interpolation through outputs from slow computer simulations (*e.g.* atmospheric carbon levels).

Gaussian Process Interpolation

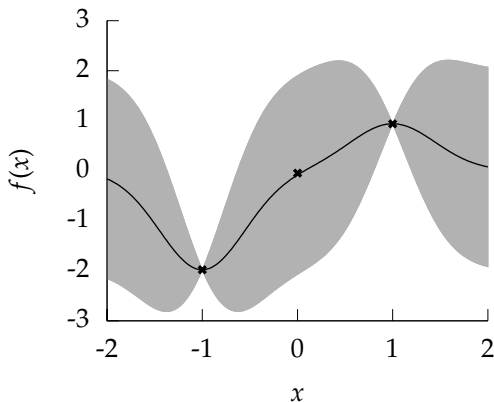


Figure : Real example: BACCO (see *e.g.* (Oakley and O'Hagan, 2002)). Interpolation through outputs from slow computer simulations (*e.g.* atmospheric carbon levels).

Gaussian Process Interpolation

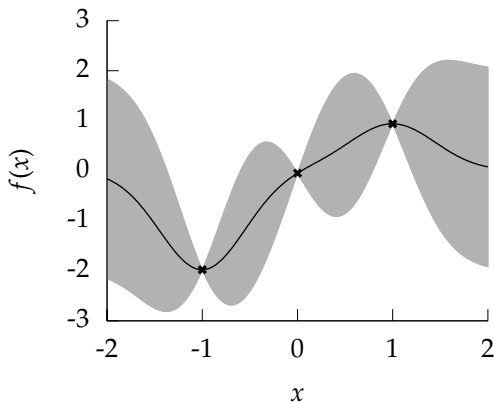


Figure : Real example: BACCO (see *e.g.* (Oakley and O'Hagan, 2002)). Interpolation through outputs from slow computer simulations (*e.g.* atmospheric carbon levels).

Gaussian Process Interpolation

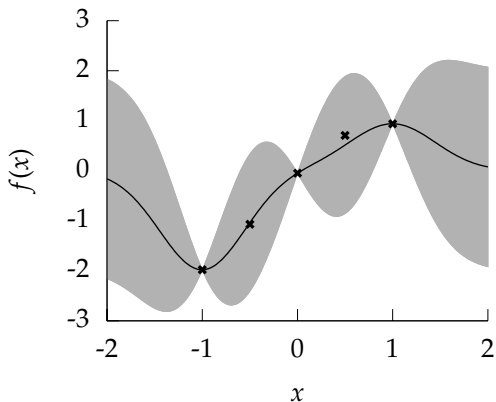


Figure : Real example: BACCO (see *e.g.* (Oakley and O'Hagan, 2002)). Interpolation through outputs from slow computer simulations (*e.g.* atmospheric carbon levels).

Gaussian Process Interpolation

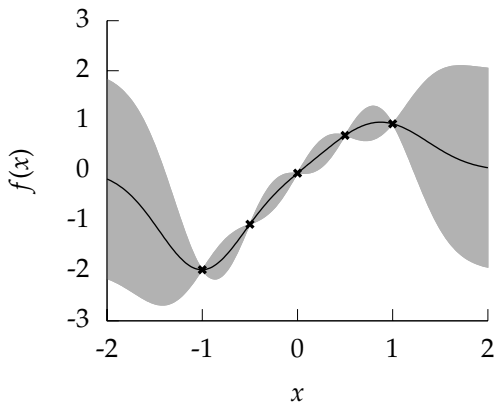


Figure : Real example: BACCO (see *e.g.* (Oakley and O'Hagan, 2002)). Interpolation through outputs from slow computer simulations (*e.g.* atmospheric carbon levels).

Gaussian Process Interpolation

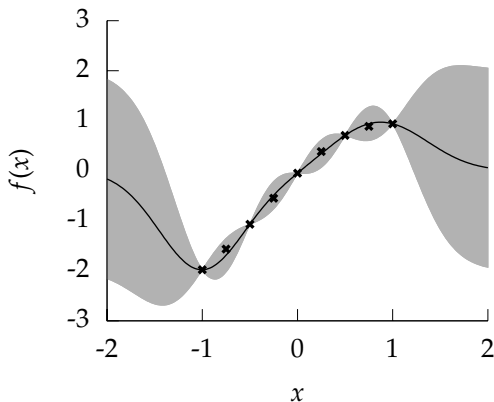


Figure : Real example: BACCO (see *e.g.* (Oakley and O'Hagan, 2002)). Interpolation through outputs from slow computer simulations (*e.g.* atmospheric carbon levels).

Gaussian Process Interpolation

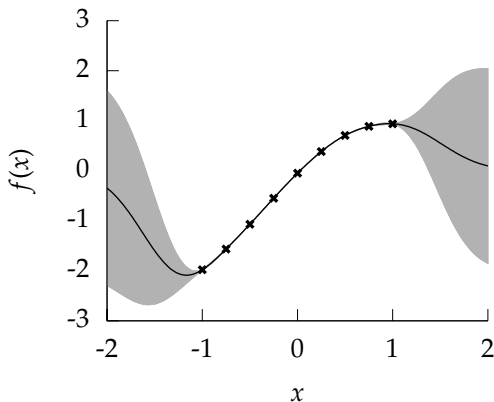


Figure : Real example: BACCO (see *e.g.* (Oakley and O'Hagan, 2002)). Interpolation through outputs from slow computer simulations (*e.g.* atmospheric carbon levels).

Gaussian Noise

- ▶ Gaussian noise model,

$$p(y_i|f_i) = \mathcal{N}(y_i|f_i, \sigma^2)$$

where σ^2 is the variance of the noise.

- ▶ Equivalent to a covariance function of the form

$$k(\mathbf{x}_i, \mathbf{x}_j) = \delta_{i,j}\sigma^2$$

where $\delta_{i,j}$ is the Kronecker delta function.

- ▶ Additive nature of Gaussians means we can simply add this term to existing covariance matrices.

Gaussian Process Regression

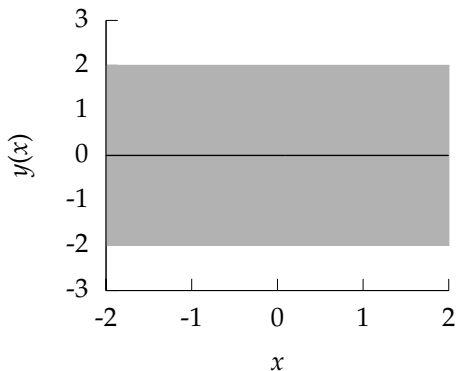


Figure : Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

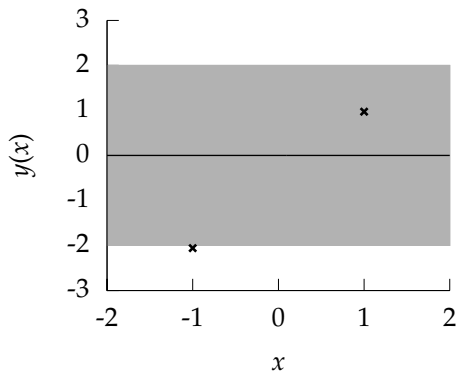


Figure : Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

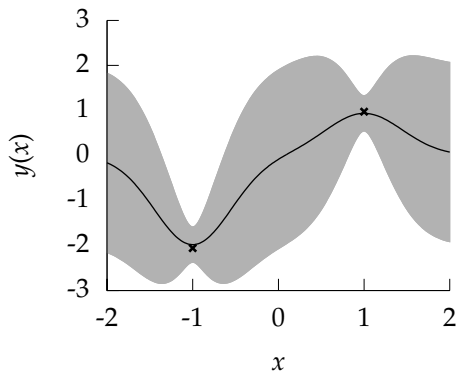


Figure : Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

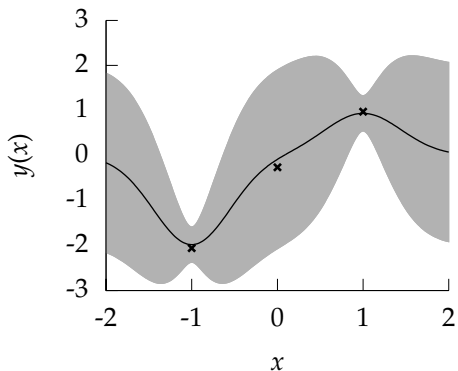


Figure : Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

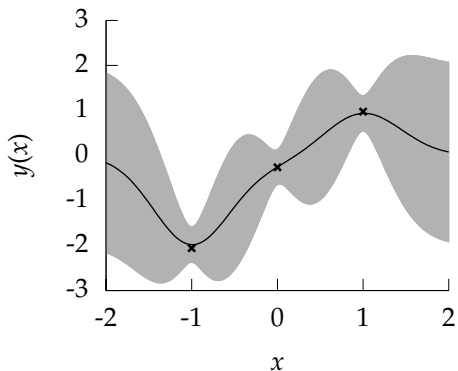


Figure : Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

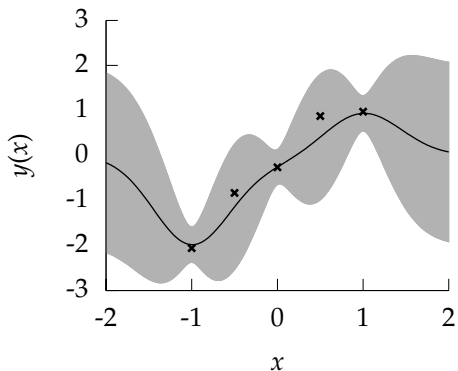


Figure : Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

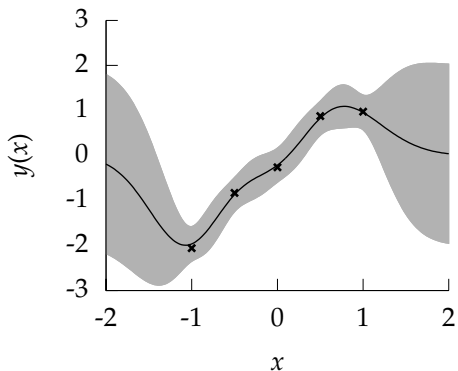


Figure : Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

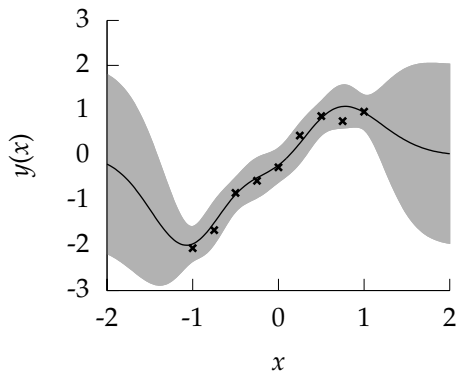


Figure : Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

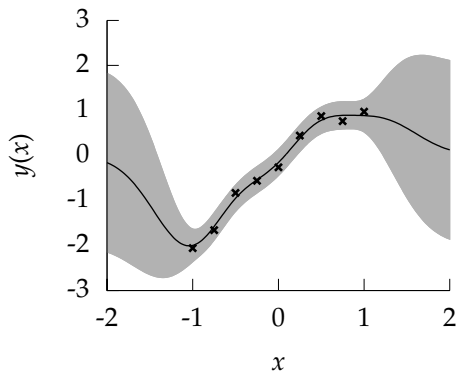


Figure : Examples include WiFi localization, C14 calibration curve.

Learning Covariance Parameters

Can we determine covariance parameters from the data?

$$\mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}) = \frac{1}{(2\pi)^{\frac{n}{2}} |\mathbf{K}|^{\frac{1}{2}}} \exp\left(-\frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2}\right)$$

The parameters are *inside* the covariance function (matrix).

$$k_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j; \theta)$$

Learning Covariance Parameters

Can we determine covariance parameters from the data?

$$\mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}) = \frac{1}{(2\pi)^{\frac{n}{2}} |\mathbf{K}|^{\frac{1}{2}}} \exp\left(-\frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2}\right)$$

The parameters are *inside* the covariance function (matrix).

$$k_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j; \theta)$$

Learning Covariance Parameters

Can we determine covariance parameters from the data?

$$\log \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}) = -\frac{1}{2} \log |\mathbf{K}| - \frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2} - \frac{n}{2} \log 2\pi$$

The parameters are *inside* the covariance function (matrix).

$$k_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j; \theta)$$

Learning Covariance Parameters

Can we determine covariance parameters from the data?

$$E(\boldsymbol{\theta}) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2}$$

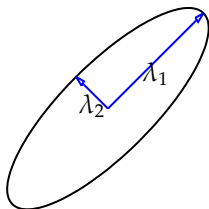
The parameters are *inside* the covariance function (matrix).

$$k_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j; \boldsymbol{\theta})$$

Eigendecomposition of Covariance

A useful decomposition for understanding the objective function.

$$\mathbf{K} = \mathbf{R}\mathbf{\Lambda}^2\mathbf{R}^\top$$



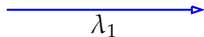
Diagonal of $\mathbf{\Lambda}$ represents distance along axes.

\mathbf{R} gives a rotation of these axes.

where $\mathbf{\Lambda}$ is a *diagonal* matrix and $\mathbf{R}^\top\mathbf{R} = \mathbf{I}$.

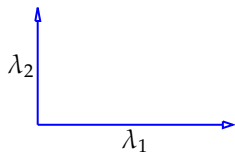
Capacity control: $\log |\mathbf{K}|$

$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$



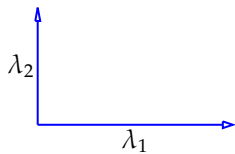
Capacity control: $\log |\mathbf{K}|$

$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$



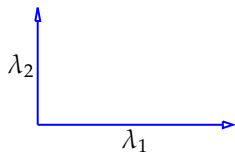
Capacity control: $\log |\mathbf{K}|$

$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$



Capacity control: $\log |\mathbf{K}|$

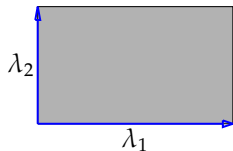
$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$



$$|\mathbf{\Lambda}| = \lambda_1 \lambda_2$$

Capacity control: $\log |\mathbf{K}|$

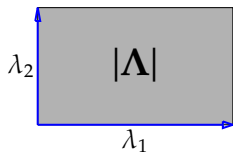
$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$



$$|\mathbf{\Lambda}| = \lambda_1 \lambda_2$$

Capacity control: $\log |\mathbf{K}|$

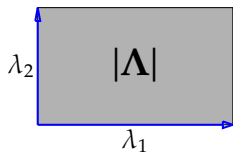
$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$



$$|\mathbf{\Lambda}| = \lambda_1 \lambda_2$$

Capacity control: $\log |\mathbf{K}|$

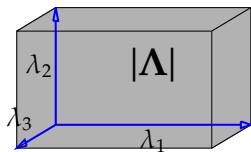
$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{bmatrix}$$



$$|\mathbf{\Lambda}| = \lambda_1 \lambda_2$$

Capacity control: $\log |\mathbf{K}|$

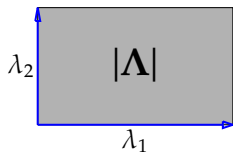
$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{bmatrix}$$



$$|\mathbf{\Lambda}| = \lambda_1 \lambda_2 \lambda_3$$

Capacity control: $\log |\mathbf{K}|$

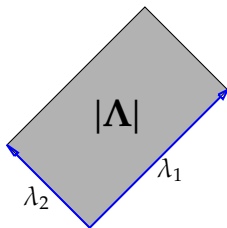
$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$



$$|\mathbf{\Lambda}| = \lambda_1 \lambda_2$$

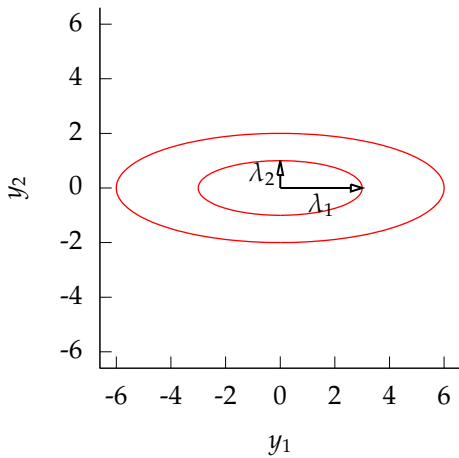
Capacity control: $\log |\mathbf{K}|$

$$\mathbf{R}\mathbf{\Lambda} = \begin{bmatrix} w_{1,1} & w_{1,2} \\ w_{2,1} & w_{2,2} \end{bmatrix}$$

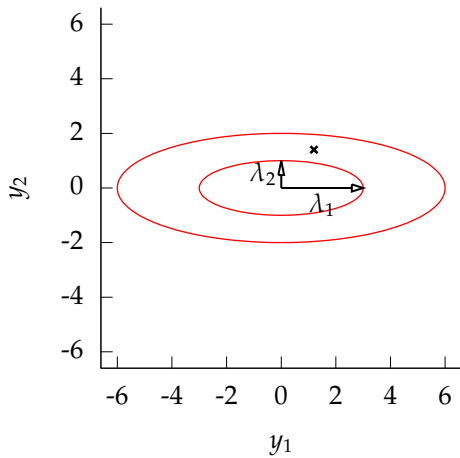


$$|\mathbf{R}\mathbf{\Lambda}| = \lambda_1 \lambda_2$$

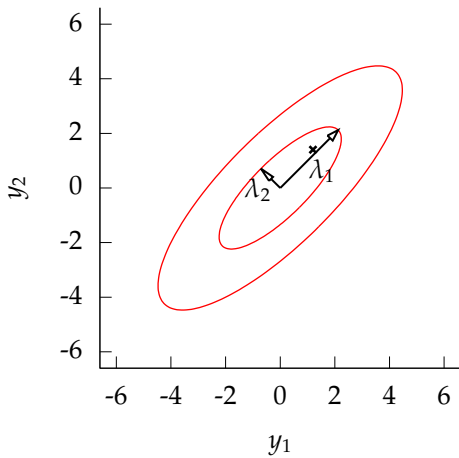
Data Fit: $\frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2}$



Data Fit: $\frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2}$

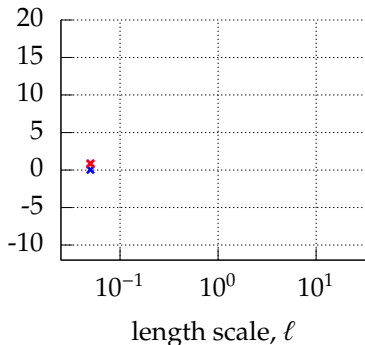
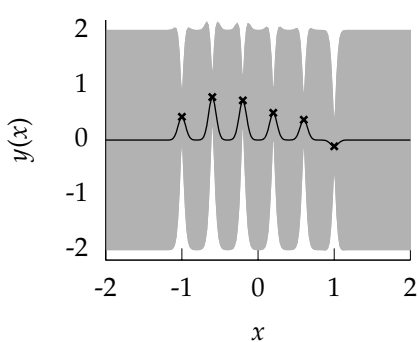


Data Fit: $\frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2}$



Learning Covariance Parameters

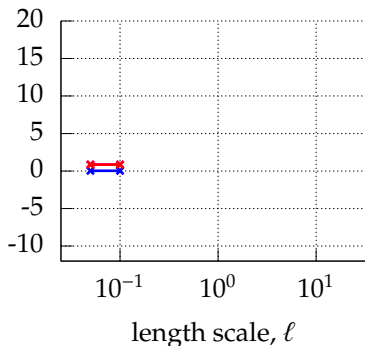
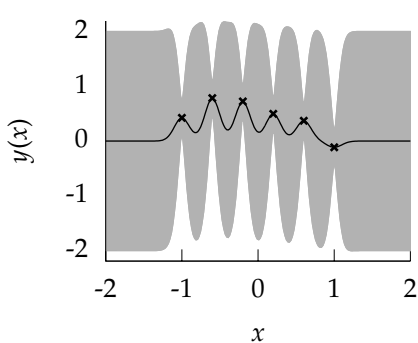
Can we determine length scales and noise levels from the data?



$$E(\theta) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2}$$

Learning Covariance Parameters

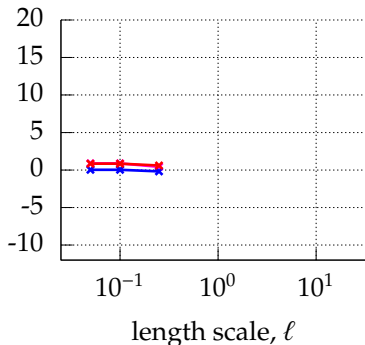
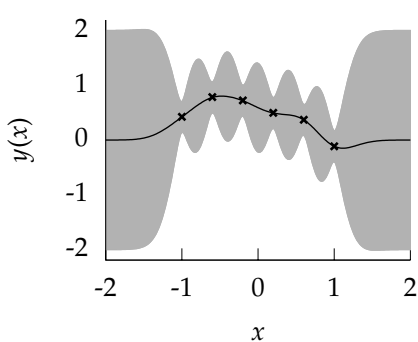
Can we determine length scales and noise levels from the data?



$$E(\theta) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2}$$

Learning Covariance Parameters

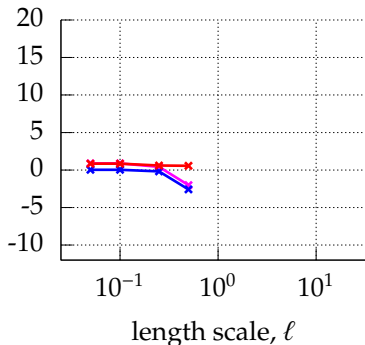
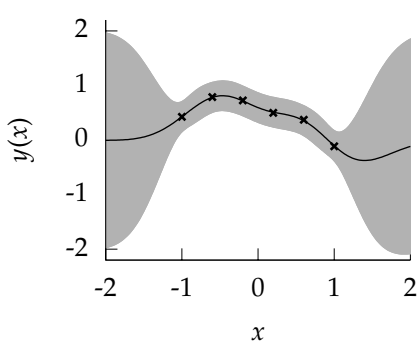
Can we determine length scales and noise levels from the data?



$$E(\theta) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2}$$

Learning Covariance Parameters

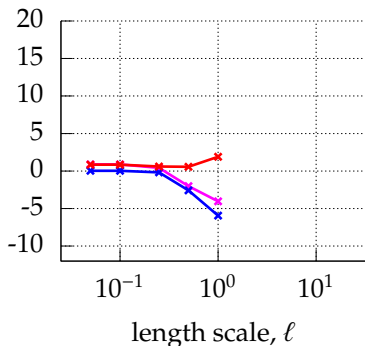
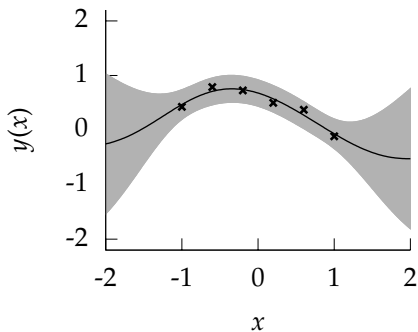
Can we determine length scales and noise levels from the data?



$$E(\theta) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2}$$

Learning Covariance Parameters

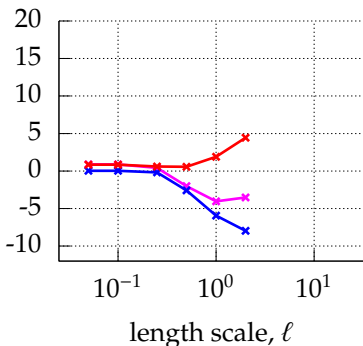
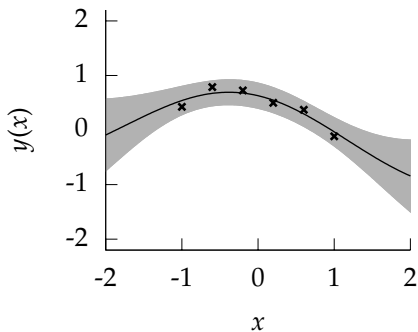
Can we determine length scales and noise levels from the data?



$$E(\theta) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2}$$

Learning Covariance Parameters

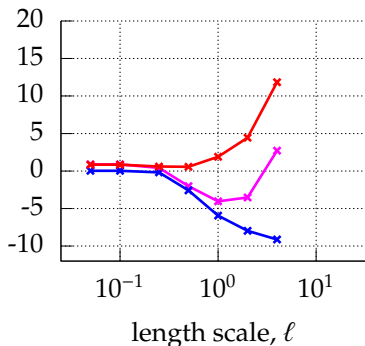
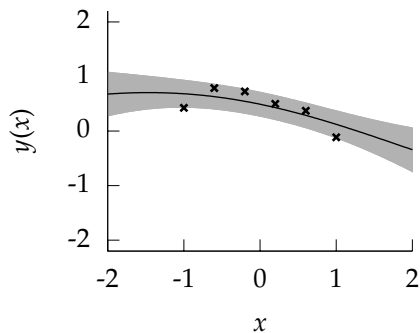
Can we determine length scales and noise levels from the data?



$$E(\theta) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2}$$

Learning Covariance Parameters

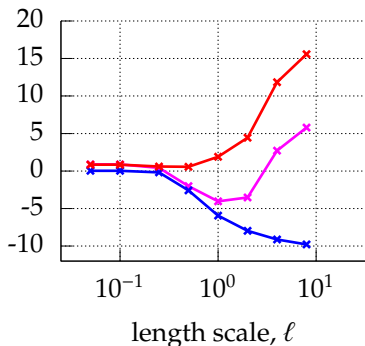
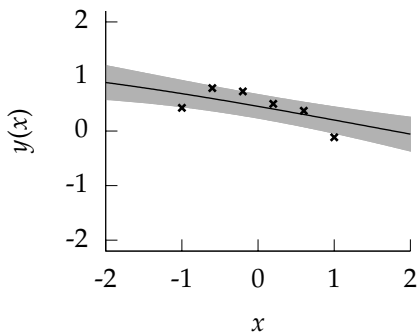
Can we determine length scales and noise levels from the data?



$$E(\theta) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2}$$

Learning Covariance Parameters

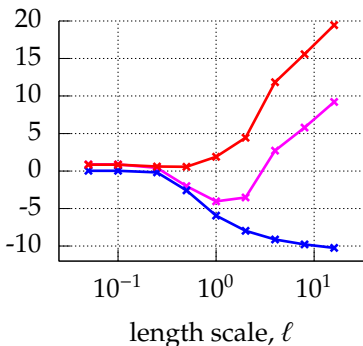
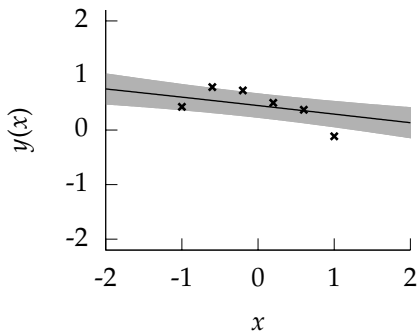
Can we determine length scales and noise levels from the data?



$$E(\theta) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2}$$

Learning Covariance Parameters

Can we determine length scales and noise levels from the data?



$$E(\theta) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}}{2}$$

Gene Expression Example

- ▶ Given given expression levels in the form of a time series from Della Gatta et al. (2008).
- ▶ Want to detect if a gene is expressed or not, fit a GP to each gene (Kalaitzis and Lawrence, 2011).

RESEARCH ARTICLE

Open Access

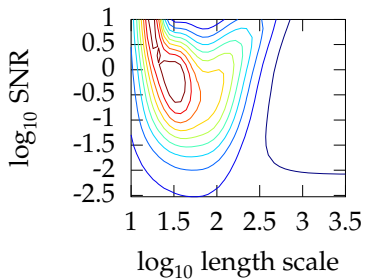
A Simple Approach to Ranking Differentially Expressed Gene Expression Time Courses through Gaussian Process Regression

Alfredo A Kalaitzis^{*} and Neil D Lawrence^{*}

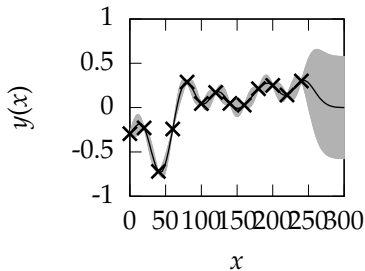
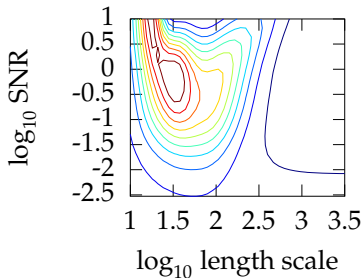
Abstract

Background: The analysis of gene expression from time series underpins many biological studies. Two basic forms of analysis recur for data of this type: removing inactive (quiet) genes from the study and determining which genes are differentially expressed. Often these analysis stages are applied disregarding the fact that the data is drawn from a time series. In this paper we propose a simple model for accounting for the underlying temporal nature of the data based on a Gaussian process.

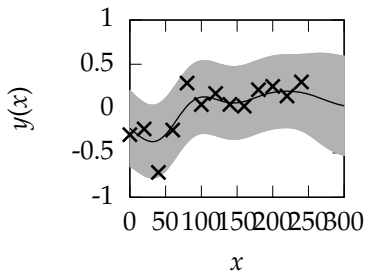
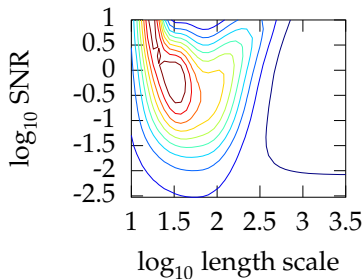
Results: We review Gaussian process (GP) regression for estimating the continuous trajectories underlying in gene expression time-series. We present a simple approach which can be used to filter quiet genes, or for the case of time series in the form of expression ratios, quantify differential expression. We assess via ROC curves the rankings produced by our regression framework and compare them to a recently proposed hierarchical Bayesian model for the analysis of gene expression time-series (BATS). We compare on both simulated and experimental data showing that the proposed approach considerably outperforms the current state of the art.



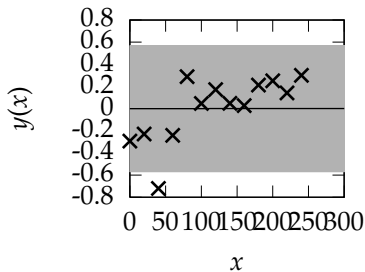
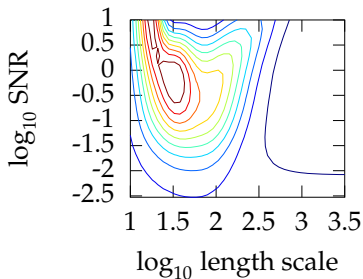
Contour plot of Gaussian process likelihood.



Optima: length scale of 1.2221 and \log_{10} SNR of 1.9654
 log likelihood is -0.22317.



Optima: length scale of 1.5162 and \log_{10} SNR of 0.21306
log likelihood is -0.23604.



Optima: length scale of 2.9886 and \log_{10} SNR of -4.506
 log likelihood is -2.1056.

Outline

Bayesian Polynomials

Distributions over Functions

Covariance from Basis Functions

Basis Function Representations

Covariance from Basis Functions

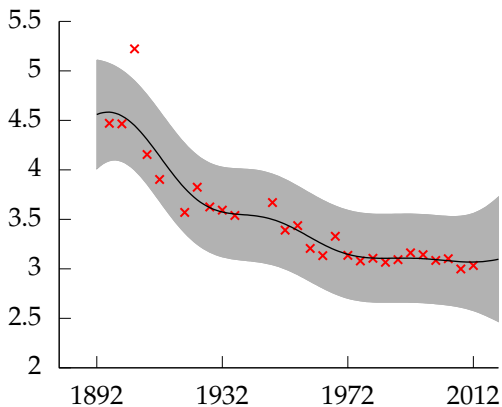
Basis Function Representations

GP Limitations

Limitations of Gaussian Processes

- ▶ Inference is $O(n^3)$ due to matrix inverse (in practice use Cholesky).
- ▶ Gaussian processes don't deal well with discontinuities (financial crises, phosphorylation, collisions, edges in images).
- ▶ Widely used exponentiated quadratic covariance (RBF) can be too smooth in practice (but there are many alternatives!!).

Gaussian Process Fit to Olympic Marathon Data



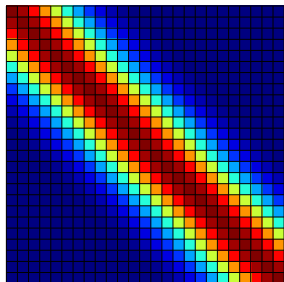
Covariance Functions

Where did this covariance matrix come from?

Exponentiated Quadratic Kernel Function (RBF, Squared Exponential, Gaussian)

$$k(\mathbf{x}, \mathbf{x}') = \alpha \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell^2}\right)$$

- ▶ Covariance matrix is built using the *inputs* to the function \mathbf{x} .
- ▶ For the example above it was based on Euclidean distance.
- ▶ The covariance function is also known as a kernel.



Covariance Functions

Where did this covariance matrix come from?

Exponentiated Quadratic Kernel Function (RBF, Squared Exponential, Gaussian)

$$k(\mathbf{x}, \mathbf{x}') = \alpha \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell^2}\right)$$

- ▶ Covariance matrix is built using the *inputs* to the function \mathbf{x} .
- ▶ For the example above it was based on Euclidean distance.
- ▶ The covariance function is also known as a kernel.

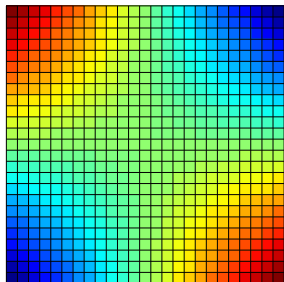
Covariance Functions

Linear Covariance Function

$$k(\mathbf{x}, \mathbf{x}') = \alpha \mathbf{x}^\top \mathbf{x}'$$

- ▶ Bayesian linear regression.

$$\alpha = 1$$



Covariance Functions

Linear Covariance Function

$$k(\mathbf{x}, \mathbf{x}') = \alpha \mathbf{x}^\top \mathbf{x}'$$

- ▶ Bayesian linear regression.

$$\alpha = 1$$

Covariance Functions

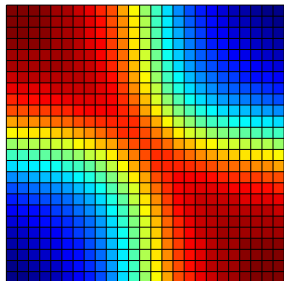
MLP Covariance Function

$$k(\mathbf{x}, \mathbf{x}') = \alpha \sin\left(\frac{w\mathbf{x}^\top \mathbf{x}' + b}{\sqrt{w\mathbf{x}^\top \mathbf{x} + b + 1} \sqrt{w\mathbf{x}'^\top \mathbf{x}' + b + 1}}\right)$$

- ▶ Based on infinite neural network model.

$$w = 40$$

$$b = 4$$



Covariance Functions

MLP Covariance Function

$$k(\mathbf{x}, \mathbf{x}') = \alpha \sin\left(\frac{w\mathbf{x}^\top \mathbf{x}' + b}{\sqrt{w\mathbf{x}^\top \mathbf{x} + b + 1} \sqrt{w\mathbf{x}'^\top \mathbf{x}' + b + 1}}\right)$$

- ▶ Based on infinite neural network model.

$$w = 40$$

$$b = 4$$

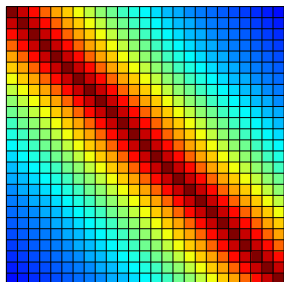
Covariance Functions

Where did this covariance matrix come from?

Ornstein-Uhlenbeck (stationary Gauss-Markov) covariance function

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

- ▶ In one dimension arises from a stochastic differential equation. Brownian motion in a parabolic tube.
- ▶ In higher dimension a Fourier filter of the form $\frac{1}{\pi(1+x^2)}$.



Covariance Functions

Where did this covariance matrix come from?

Ornstein-Uhlenbeck (stationary Gauss-Markov) covariance function

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

- ▶ In one dimension arises from a stochastic differential equation. Brownian motion in a parabolic tube.
- ▶ In higher dimension a Fourier filter of the form $\frac{1}{\pi(1+x^2)}$.

Summary

- ▶ Broad introduction to Gaussian processes.
 - ▶ Started with Gaussian distribution.
 - ▶ Motivated Gaussian processes through the multivariate density.
- ▶ Emphasized the role of the covariance (not the mean).
- ▶ Performs nonlinear regression with error bars.
- ▶ Parameters of the covariance function (kernel) are easily optimized with maximum likelihood.

Reading

- ▶ Section 3.7–3.8 of Rogers and Girolami (pg 122–133).
- ▶ Section 3.4 of Bishop (pg 161–165).
- ▶ Chapter 1 & 2 of Rasmussen and Williams.

References I

- C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer-Verlag, 2006. [\[Google Books\]](#) .
- G. Della Gatta, M. Bansal, A. Ambesi-Impiombato, D. Antonini, C. Missero, and D. di Bernardo. Direct targets of the trp63 transcription factor revealed by a combination of gene expression profiling and reverse engineering. *Genome Research*, 18(6):939–948, Jun 2008. [\[URL\]](#). [\[DOI\]](#).
- A. A. Kalaitzis and N. D. Lawrence. A simple approach to ranking differentially expressed gene expression time courses through Gaussian process regression. *BMC Bioinformatics*, 12(180), 2011. [\[DOI\]](#).
- R. M. Neal. *Bayesian Learning for Neural Networks*. Springer, 1996. Lecture Notes in Statistics 118.
- J. Oakley and A. O’Hagan. Bayesian inference for the uncertainty distribution of computer model outputs. *Biometrika*, 89(4):769–784, 2002.

References II

- C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, Cambridge, MA, 2006. [\[Google Books\]](#) .
- S. Rogers and M. Girolami. *A First Course in Machine Learning*. CRC Press, 2011. [\[Google Books\]](#) .
- C. K. I. Williams. Computation with infinite neural networks. *Neural Computation*, 10(5):1203–1216, 1998.