

Bayesian Linear Regression

Sargur Srihari

srihari@cedar.buffalo.edu

Motivation

- In maximum likelihood
 - model complexity (number of basis functions) needs to be controlled according to the size of the data set
- Adding regularization term helps control model complexity
- Number and choice of basis functions is still important
- Bayesian treatment avoids this problem

Parameter Distribution

- Prior probability distribution $p(\mathbf{w})$ over model parameters \mathbf{w}
- Noise precision parameter β assumed known
- Since Likelihood function $p(t|\mathbf{w})$ with Gaussian noise has an exponential form

– Conjugate prior is given by Gaussian

$$p(\mathbf{w}) = N(\mathbf{w} | \mathbf{m}_0, S_0)$$

with mean m_0 and covariance S_0

Posterior Distribution of Parameters

- Given by product of likelihood function and prior
 - $p(\mathbf{w}|D) = p(D|\mathbf{w})p(\mathbf{w})/p(D)$
- Due to choice of conjugate Gaussian prior, posterior is also Gaussian
- Posterior can be written directly in the form

$p(\mathbf{w}|\mathbf{t}) = N(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$ where

$$\mathbf{m}_N = \mathbf{S}_N(\mathbf{S}_0^{-1}\mathbf{m}_0 + \beta\Phi^T\mathbf{t}),$$

and

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta\Phi^T\Phi$$

Gaussian Prior for Linear Regression

- Zero mean isotropic Gaussian

$$p(\mathbf{w} | \alpha) = N(\mathbf{w} | 0, \alpha^{-1} I) \quad \text{Single precision parameter}$$

- Corresponding posterior distribution is

$$p(\mathbf{w} | \mathbf{t}) = N(\mathbf{w} | \mathbf{m}_N, S_N) \quad \text{where}$$

$$\mathbf{m}_N = \beta S_N \Phi^T \mathbf{t} \quad \text{and} \quad S_N^{-1} = \alpha I + \beta \Phi^T \Phi$$

Note:

β is noise precision and

α is distribution of parameter \mathbf{w}

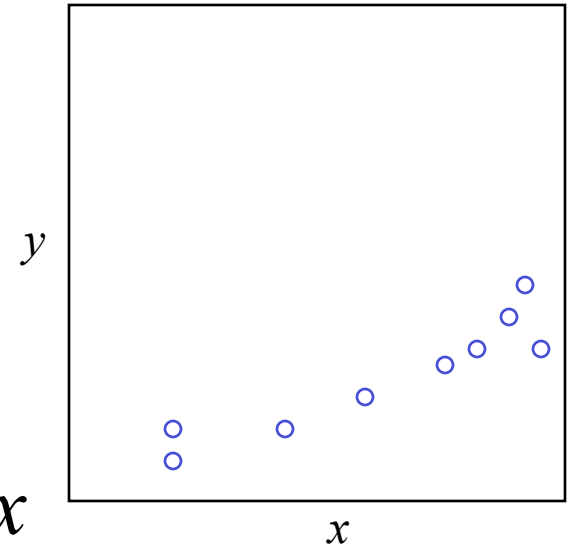
Log Posterior Distribution

$$\ln p(\mathbf{w} | t) = -\frac{\beta}{2} \sum_{n=1}^N \left\{ t_n - \mathbf{w}^T \phi(x_n) \right\}^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const}$$

- Maximization wrt \mathbf{w} is equivalent to minimization of sum-of-squares error function with addition of quadratic regularization term with $\lambda = \alpha/\beta$

Bayesian Linear Regression Example

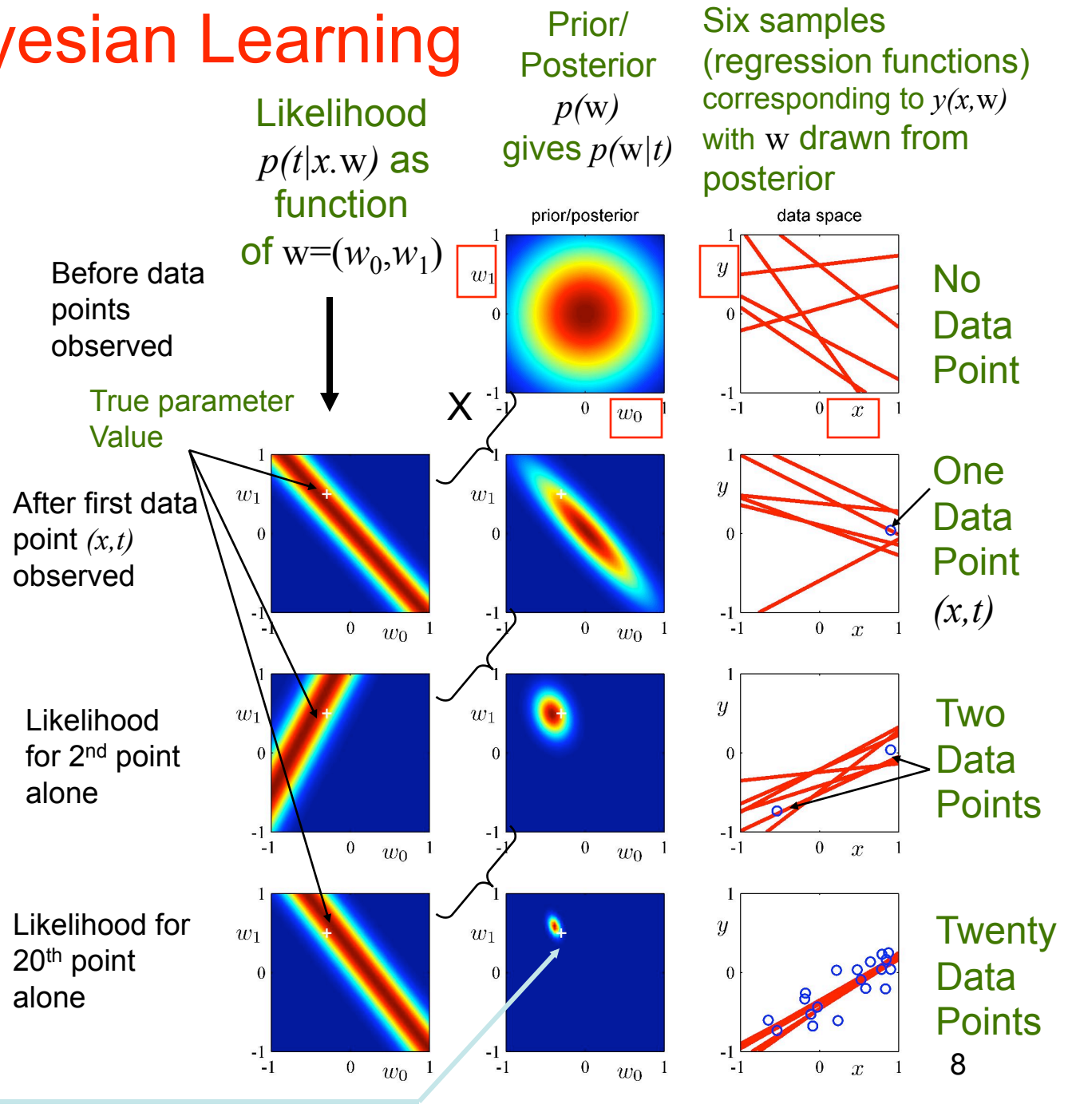
- Straight line fitting
- Single input variable x
- Single target variable t
- Linear model $y(x, \mathbf{w}) = w_0 + w_1 x$
 - t and y used interchangeably here
- Since there are only two parameters
 - We can plot prior and posterior distributions in parameter space



Sequential Bayesian Learning

- Synthetic data generated from $f(x, a) = a_0 + a_1 x$ with parameter values $a_0 = -0.3$ and $a_1 = 0.5$
- By first choosing x_n from $U(x | -1, 1)$, then evaluating $f(x_n, a)$ and then adding Gaussian noise with std dev 0.2 to obtain target values t_n
- Goal is to recover values of a_0 and a_1

With infinite points posterior is a delta function centered at true parameters (white cross)

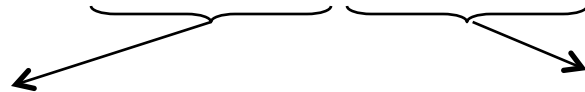


Predictive Distribution

- Usually not interested in the value of \mathbf{w} itself
- But predicting t for new values of \mathbf{x}
- This requires evaluating an integral

$$p(t | \mathbf{x}, \alpha, \beta) = \int p(t | \mathbf{w}, \beta) p(\mathbf{w} | \mathbf{x}, \alpha, \beta) d\mathbf{w} \quad \text{Marginalizing } \mathbf{w}$$

Conditioning variable \mathbf{x} has been left out for convenience.



conditioning distribution of target variable

$$p(t | \mathbf{x}, \mathbf{w}, \beta) = N(t | y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

posterior weight distribution

$$p(\mathbf{w} | \mathbf{x}) = N(\mathbf{w} | \mathbf{m}_N, S_N)$$

RHS is convolution of two Gaussian distributions whose result is a Gaussian

Predictive Distribution

- Predictive distribution takes the form

$$p(t | \mathbf{x}, t, \alpha, \beta) = N(t | m_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$

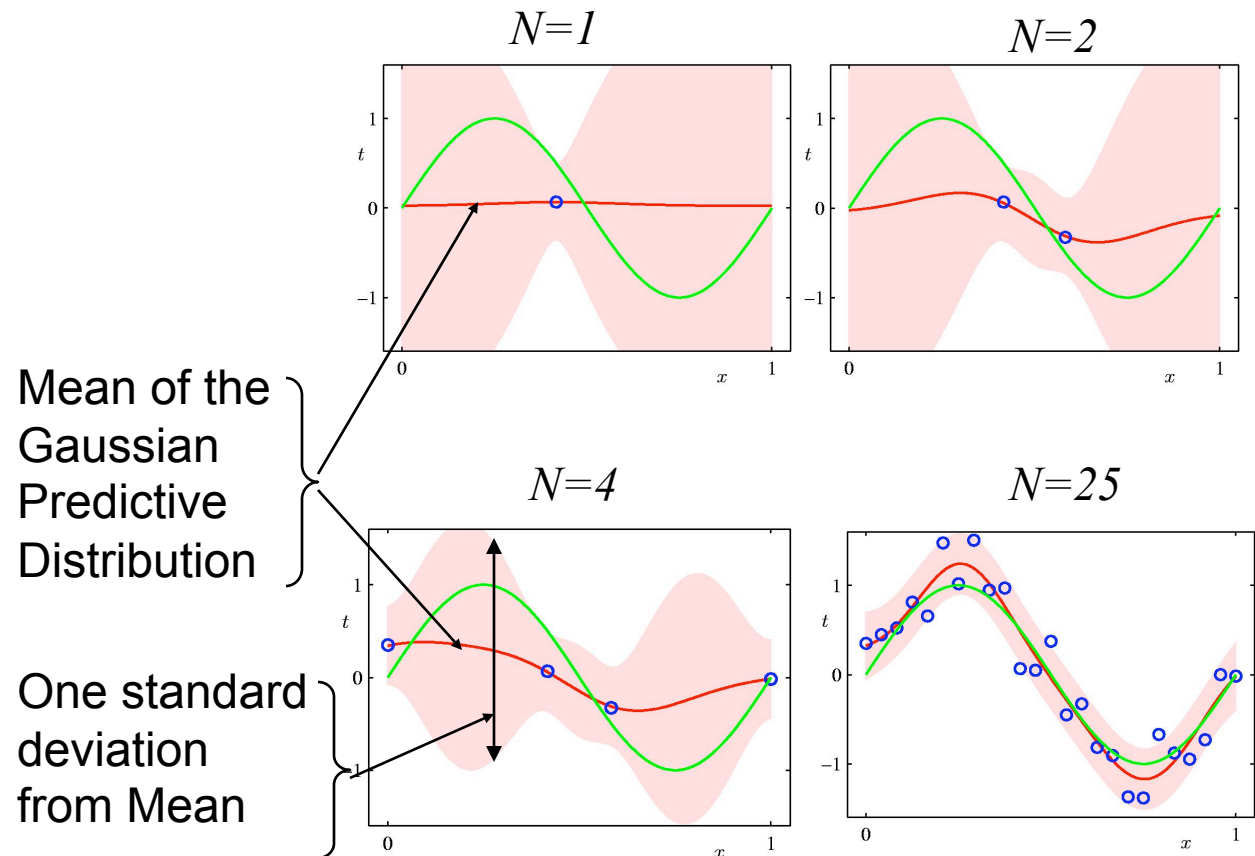
$$\text{where } \sigma_N^2(\mathbf{x}) = \underbrace{\frac{1}{\beta}}_{\text{Noise in data}} + \underbrace{\phi(\mathbf{x})^T S_N \phi(\mathbf{x})}_{\text{Uncertainty associated with parameters } \mathbf{w}}$$

Noise in data

Uncertainty associated
with parameters \mathbf{w}

Example of Predictive Distribution

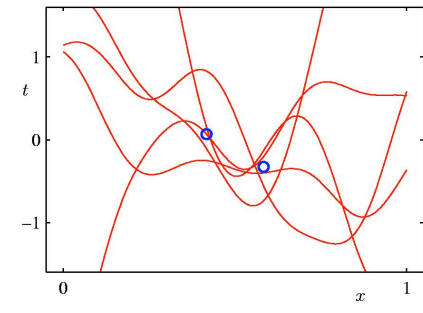
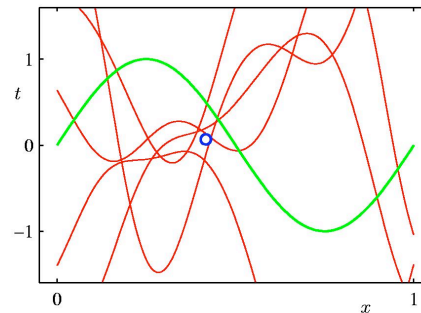
- Predictive uncertainty depends on x and is smallest in neighborhood of data points
- Level of uncertainty decreases as more data points are observed



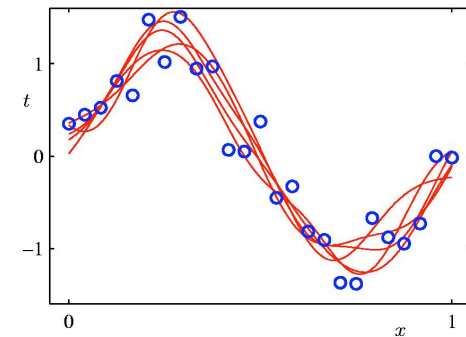
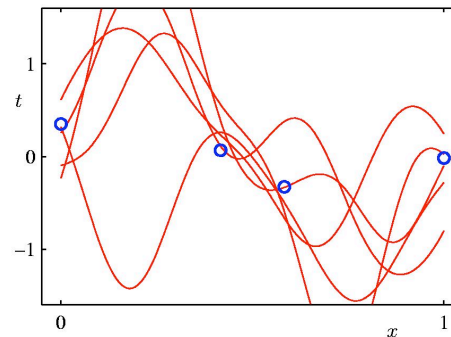
Plot only shows point-wise predictive variance
 To show covariance between predictions at different values of x draw samples from posterior distribution over w

Plots of function $y(x, w)$

- Draw samples from from posterior distribution over w



- Plots of samples from $y(x, w)$ corresponding to predictive distributions in previous slide



Equivalent Kernel

- Posterior mean solution has a kernel interpretation

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

Sets stage for kernel methods and Gaussian processes

- Substituting $\mathbf{m}_N = \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t}$ where \mathbf{S}_N is defined by the design matrix

Mean of predictive distribution at point \mathbf{x}

$$y(\mathbf{x}, \mathbf{m}_N) = \mathbf{m}_N^T \boldsymbol{\phi}(\mathbf{x}) = \beta \boldsymbol{\phi}(\mathbf{x})^T \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t} = \sum_{n=1}^N \beta \boldsymbol{\phi}(\mathbf{x})^T \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}_n) t_n = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n$$

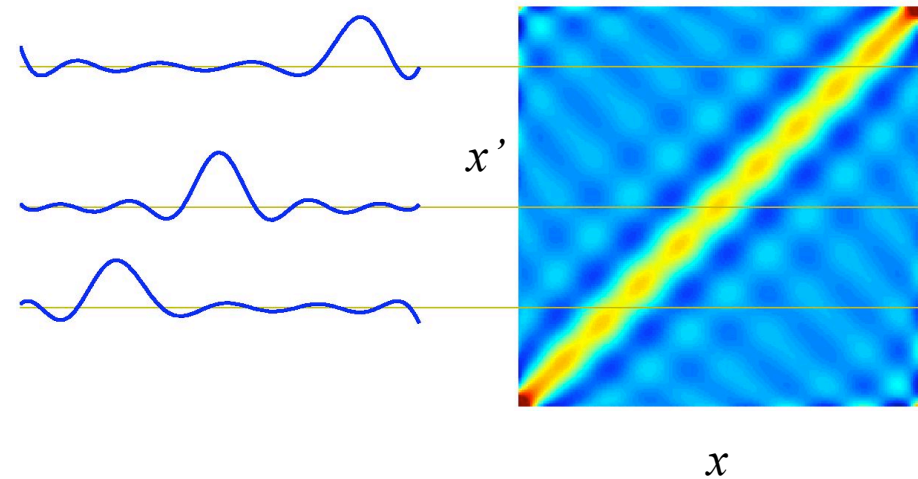
- Where the function $k(\mathbf{x}, \mathbf{x}') = \beta \boldsymbol{\phi}(\mathbf{x})^T \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}')$ is known as the equivalent kernel
 - Equivalent kernel depends on the input values \mathbf{x}_n

Equivalent kernel for Gaussian Basis

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

Plot of $k(x, x')$ shown as a function of x and x'

For three values of x' the behavior of $k(x, x')$ is shown as a slice

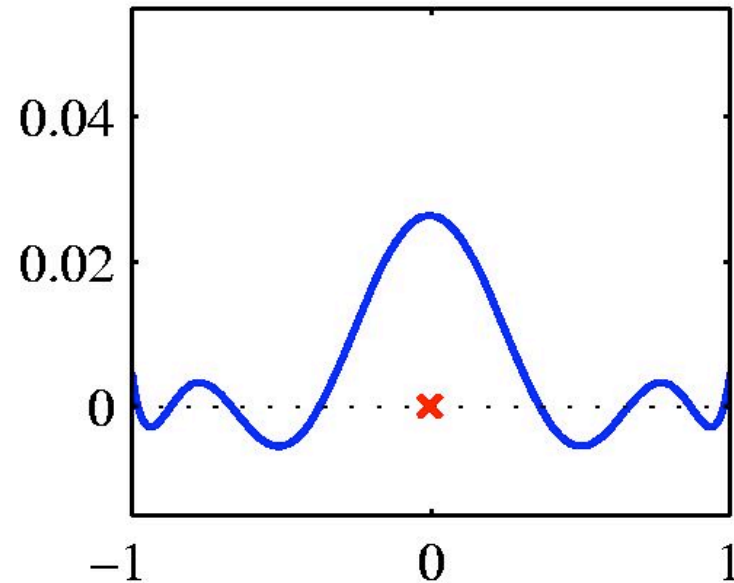


They are localized around x , i.e., peaks when $x' = x$
Local evidence is weighted more than distant evidence

Data set used to generate kernel were 200 values of x equally spaced in $(-1, 1)$

Equivalent Kernel for Polynomial Basis Function

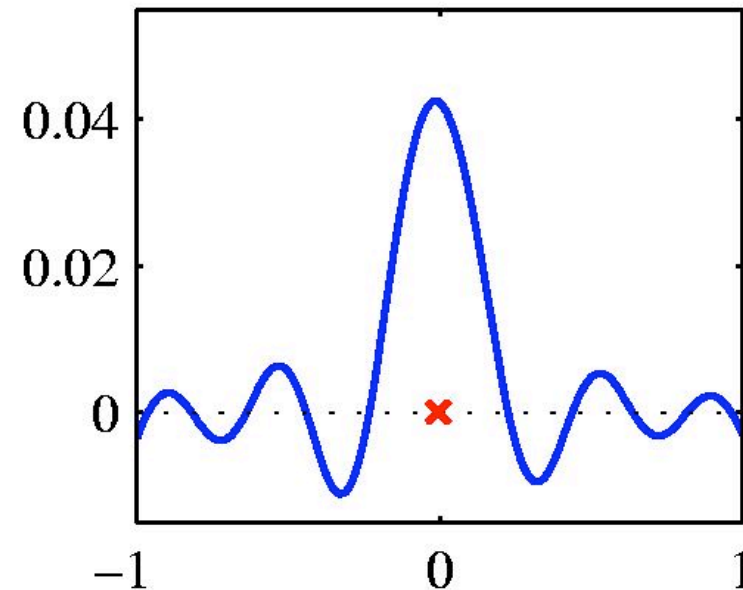
$$k(\mathbf{x}, \mathbf{x}') = \beta \phi(\mathbf{x})^T S_N \phi(\mathbf{x}')$$



Localized function of \mathbf{x}' even though corresponding basis function is nonlocal

Equivalent Kernel for Sigmoidal Basis Function

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



Localized function of x' even though corresponding basis function is nonlocal

Covariance between $y(\mathbf{x})$ and $y(\mathbf{x}')$

$$\begin{aligned}\text{COV} [y(\mathbf{x}'), y(\mathbf{x})] &= \text{COV}[\phi(\mathbf{x})^\top \mathbf{w}, \mathbf{w}^\top \phi(\mathbf{x}')] \\ &= \phi(\mathbf{x})^\top \mathbf{S}_N \phi(\mathbf{x}') \\ &= \beta^{-1} k(\mathbf{x}, \mathbf{x}')\end{aligned}$$

From the form of the equivalent kernel, the predictive mean at nearby points will be highly correlated

For more distant pairs correlation is smaller

Kernel Function formulation

- Formulation of Linear Regression in terms of kernel function suggests:
 - Directly define kernel functions and use to make predictions for input \mathbf{x}
- Leads to Gaussian Processes for Regression and Classification
- Can be shown that

$$\sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) = 1 \quad k(\mathbf{x}, \mathbf{z}) = \psi(\mathbf{x})^T \psi(\mathbf{z})$$

- where $\psi(\mathbf{x}) = \beta^{1/2} S_N^{1/2} \phi(\mathbf{x})$