

Bayesian Linear Regression

Sargur Srihari

srihari@cedar.buffalo.edu

Topics in Bayesian Regression

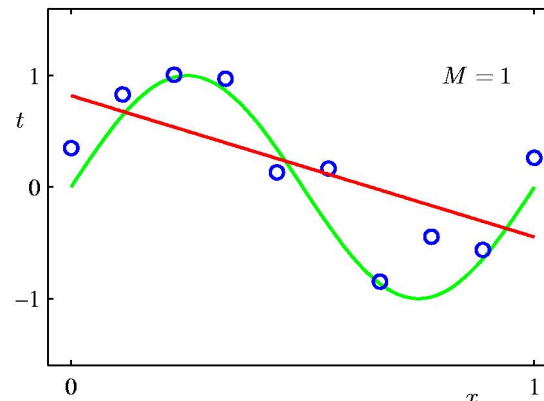
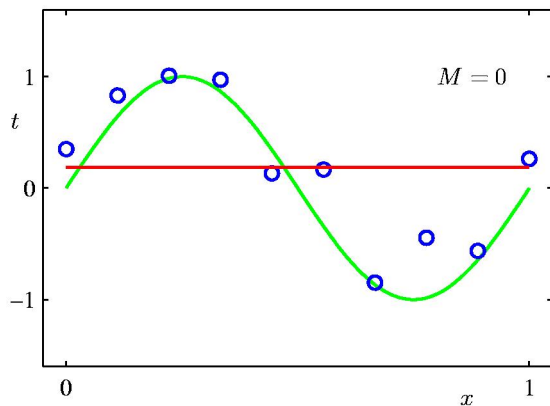
- Recall Max Likelihood Linear Regression
- Parameter Distribution
- Predictive Distribution
- Equivalent Kernel

Linear Regression: model complexity M

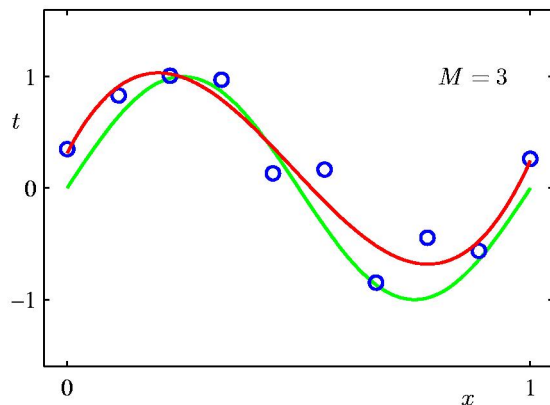
- Polynomial regression

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

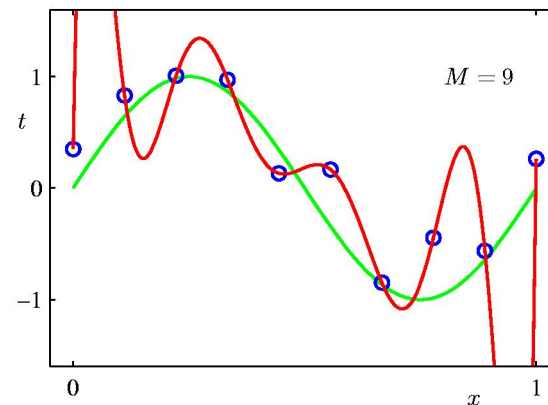
- Red lines are best fits with $M = 0, 1, 3, 9$ and $N=10$



← Poor representations of $\sin(2\pi x)$



← Best Fit to $\sin(2\pi x)$



Over Fit
Poor representation of $\sin(2\pi x)$

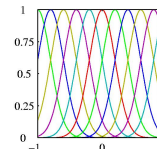
Max Likelihood Regression

- Input vector \mathbf{x} , basis functions $\{\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})\}$:

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

Radial basis fns:

$$\phi_j(\mathbf{x}) = \exp\left[-\frac{1}{2}(\mathbf{x} - \mu_j)^T \Sigma^{-1}(\mathbf{x} - \mu_j)\right]$$



- Objective Function:

Max Likelihood objective with N examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$:
(equivalent to Mean Squared Error Objective)

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - \mathbf{w}^T \phi(\mathbf{x}_n) \right\}^2$$

Regularized MSE with N examples:
(λ is the regularization coefficient)

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

- Closed-form ML solution is:

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

where Φ is the design matrix:
($\Phi^T \Phi$)⁻¹ is Moore-Penrose inverse

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & & & \\ \vdots & & & \\ \phi_0(\mathbf{x}_N) & & & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

Regularized solution is:

$$\mathbf{w}_{ML} = (\lambda I + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

- Gradient Descent: $\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E$

$$\nabla E = - \sum_{n=1}^N \left\{ t_n - \mathbf{w}^{(\tau)T} \phi(\mathbf{x}_n) \right\} \phi(\mathbf{x}_n)$$

Regularized version:

$$\nabla E = \left[- \sum_{n=1}^N \left\{ t_n - \mathbf{w}^{(\tau)T} \phi(\mathbf{x}_n) \right\} \phi(\mathbf{x}_n) \right] - \lambda \mathbf{w}^{(\tau)}$$

Shortcomings of MLE

- M.L.E. of parameters \mathbf{w} does not address
 - M (Model complexity: how many basis functions?)
 - It is controlled by data size N
 - More data allows better fit without overfitting
- Regularization also controls overfit (λ controls effect)

$$E(\mathbf{w}) = E_D(\mathbf{w}) + \lambda E_W(\mathbf{w}) \text{ where } E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - \mathbf{w}^T \phi(\mathbf{x}_n) \right\}^2 \quad E_W(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

 - But M and choice of ϕ_j are still important
 - M can be determined by holdout, but wasteful of data
- Model complexity and over-fitting are better handled using Bayesian approach

Bayesian Linear Regression

- Using Bayes rule, posterior is proportional to Likelihood \times Prior:

$$p(\mathbf{w} | \mathbf{t}) = \frac{p(\mathbf{t} | \mathbf{w})p(\mathbf{w})}{p(\mathbf{t})}$$

- where $p(\mathbf{t} | \mathbf{w})$ is the likelihood of observed data
- $p(\mathbf{w})$ is prior distribution over the parameters
- We will look at:
 - A normal distribution for prior $p(\mathbf{w})$
 - Likelihood $p(\mathbf{t} | \mathbf{w})$ is a product of Gaussians based on the noise model
 - And conclude that posterior is also Gaussian

Gaussian Prior Parameters

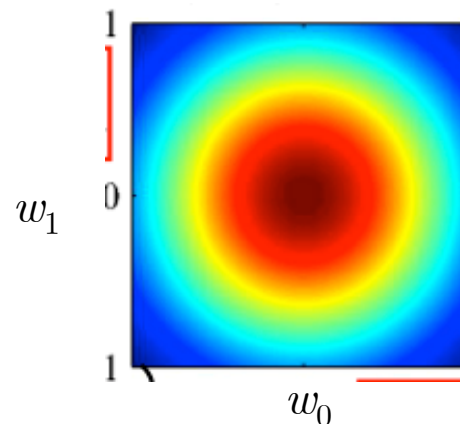
Assume multivariate Gaussian prior for \mathbf{w} (which has components w_0, \dots, w_{M-1})

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

with mean \mathbf{m}_0 and covariance matrix \mathbf{S}_0

If we choose $\mathbf{S}_0 = \alpha^{-1} \mathbf{I}$ it means that the variances of the weights are all equal to α^{-1} and covariances are zero

$p(\mathbf{w})$ with zero mean ($\mathbf{m}_0 = \mathbf{0}$)
and isotropic over weights (*same variances*)



Likelihood of Data is Gaussian

- Assume noise precision parameter β

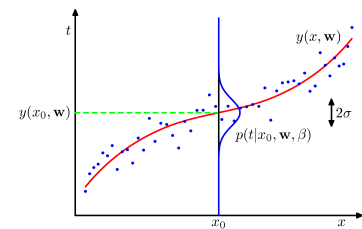
$t = y(\mathbf{x}, \mathbf{w}) + \varepsilon$ where ε is defined probabilistically as Gaussian noise

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

Note that output t is a scalar

- Likelihood of $\mathbf{t} = \{t_1, \dots, t_N\}$ is then

$$p(\mathbf{t} | X, \mathbf{w}, \beta) = \prod_{n=1}^N N(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})$$



- This is the probability of target data \mathbf{t} given the parameters \mathbf{w} and input $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- Due to Gaussian noise, likelihood $p(\mathbf{t} | \mathbf{w})$ is also a Gaussian

Posterior Distribution is also Gaussian

- **Prior:** $p(\mathbf{w}) \sim N(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$ i.e., it is Gaussian
- **Likelihood comes from Gaussian noise**

$$p(\mathbf{t} | X, \mathbf{w}, \beta) = \prod_{n=1}^N N(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})$$

- It follows that posterior $p(\mathbf{w} | \mathbf{t})$ is also Gaussian
- *Proof:* use standard result from Gaussians:
 - If marginal $p(\mathbf{w})$ & conditional $p(\mathbf{t} | \mathbf{w})$ have Gaussian forms then the marginals $p(\mathbf{t})$ and $p(\mathbf{w} | \mathbf{t})$ are also Gaussian:
 - Let $p(\mathbf{w}) = N(\mathbf{w} | \boldsymbol{\mu}, \Lambda^{-1})$ and $p(\mathbf{t} | \mathbf{w}) = N(\mathbf{t} | \mathbf{A}\mathbf{w} + \mathbf{b}, \mathbf{L}^{-1})$
 - Then marginal $p(\mathbf{t}) = N(\mathbf{t} | \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\Lambda^{-1}\mathbf{A}^T)$ and conditional $p(\mathbf{w} | \mathbf{t}) = N(\mathbf{w} | \Sigma \{ \mathbf{A}^T \mathbf{L}(\mathbf{t} - \mathbf{b}) + \Lambda \boldsymbol{\mu} \}, \Sigma)$ where $\Sigma = (\Lambda + \mathbf{A}^T \mathbf{L} \mathbf{A})^{-1}$

Exact form of Posterior Distribution

- We have $p(\mathbf{w}) = N(\mathbf{w} | \mathbf{m}_0, S_0)$ & $p(t | X, \mathbf{w}, \beta) = \prod_{n=1}^N N(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})$
- Posterior is also Gaussian, written directly as

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

– where \mathbf{m}_N is the mean of the posterior

given by $\mathbf{m}_N = S_N (S_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t})$

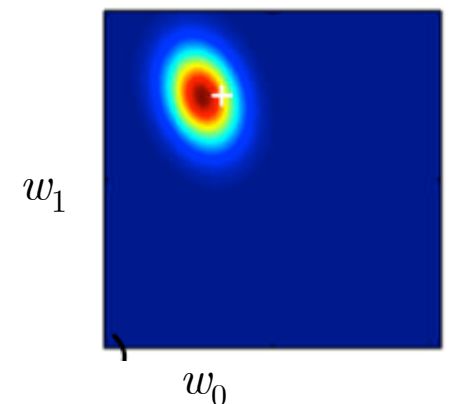
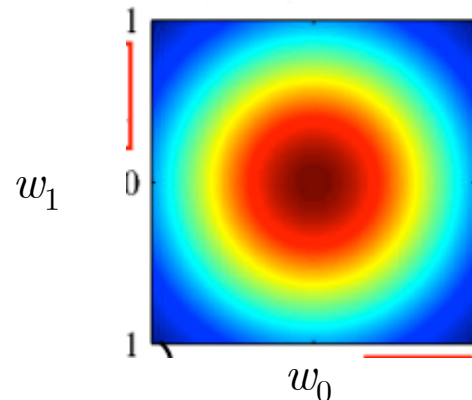
– and S_N is the covariance matrix of posterior

given by $S_N^{-1} = S_0^{-1} + \beta \Phi^T \Phi$

Φ is the design matrix

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & & & \\ \vdots & & & \\ \phi_0(\mathbf{x}_N) & & & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

Prior $p(\mathbf{w} | \alpha) = N(\mathbf{w} | \mathbf{0}, \alpha^{-1} I)$
and Posterior in weight space
for scalar input x and
 $y(x, \mathbf{w}) = w_0 + w_1 x$



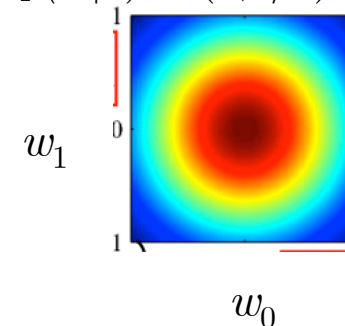
Properties of Posterior

1. Since posterior $p(\mathbf{w}|\mathbf{t}) = N(\mathbf{w}|\mathbf{m}_N, S_N)$ is Gaussian its mode coincides with its mean
 - Thus maximum posterior weight is $\mathbf{w}_{\text{MAP}} = \mathbf{m}_N$
2. Infinitely broad prior $S_0 = \alpha^{-1} \mathbf{I}$, i.e., precision $\alpha \rightarrow 0$
 - Then mean \mathbf{m}_N reduces to the maximum likelihood value, i.e., mean is the solution vector
$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$
3. If $N = 0$, posterior reverts to the prior
4. If data points arrive sequentially, then posterior to any stage acts as prior distribution for subsequent data points

Choose a simple Gaussian prior $p(\mathbf{w})$

$$y(x, \mathbf{w}) = w_0 + w_1 x$$

$$p(\mathbf{w} | \alpha) \sim N(0, 1/\alpha)$$



- Zero mean ($\mathbf{m}_0 = \mathbf{0}$) isotropic
- (*same variances*) Gaussian

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

Single precision
parameter α

- Corresponding posterior distribution is

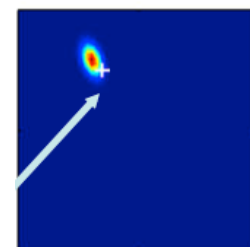
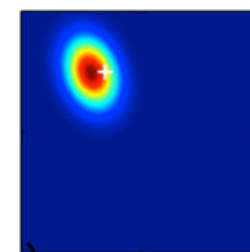
$$p(\mathbf{w} | \mathbf{t}) = N(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$$

where

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

Note:

β is noise precision and
 α is variance of parameter \mathbf{w} in prior



Point
Estimate
with
infinite
samples

Equivalence to MLE with Regularization

- Since $p(t | X, \mathbf{w}, \beta) = \prod_{n=1}^N N(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})$ and $p(\mathbf{w} | \alpha) = N(\mathbf{w} | \mathbf{0}, \alpha^{-1} I)$

- we have

$$p(\mathbf{w} | t) = \prod_{n=1}^N N(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) N(\mathbf{w} | \mathbf{0}, \alpha^{-1} I)$$

- Log of Posterior is

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

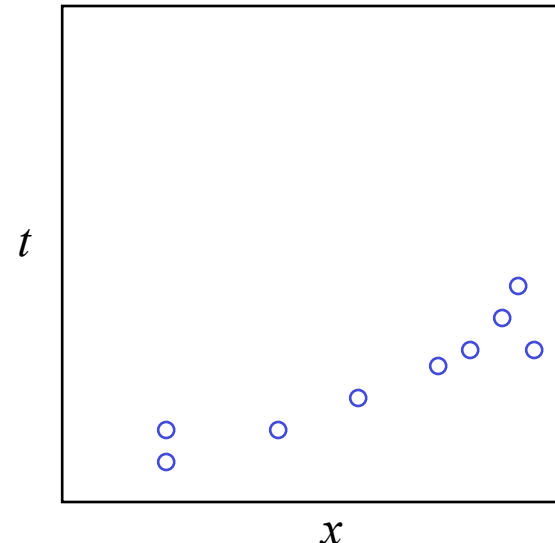
- Thus Maximization of posterior is equivalent to minimization of sum-of-squares error

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

with addition of quadratic regularization term $\mathbf{w}^T \mathbf{w}$
with $\lambda = \alpha / \beta$

Bayesian Linear Regression Example (Straight Line Fit)

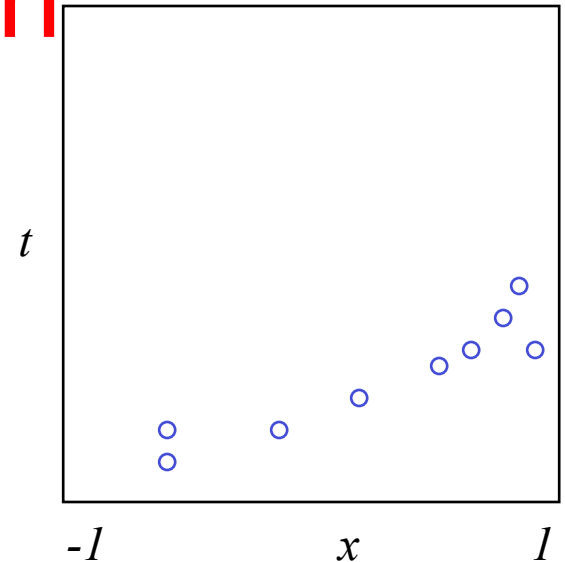
- Single input variable x
- Single target variable t
- Goal is to fit
 - Linear model $y(x, \mathbf{w}) = w_0 + w_1 x$
- Goal of Linear Regression is to recover $\mathbf{w} = [w_0, w_1]$ given the samples



Data Generation

- Synthetic data generated from $f(x, \mathbf{w}) = w_0 + w_1 x$ with parameter values

$$w_0 = -0.3 \text{ and } w_1 = 0.5$$

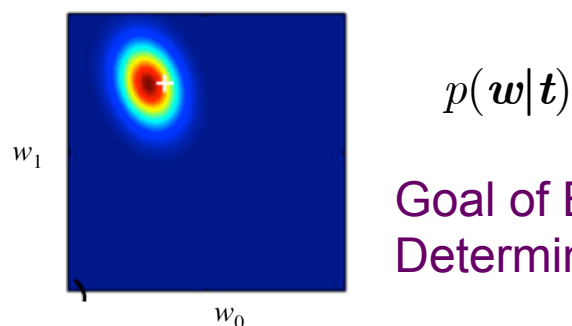
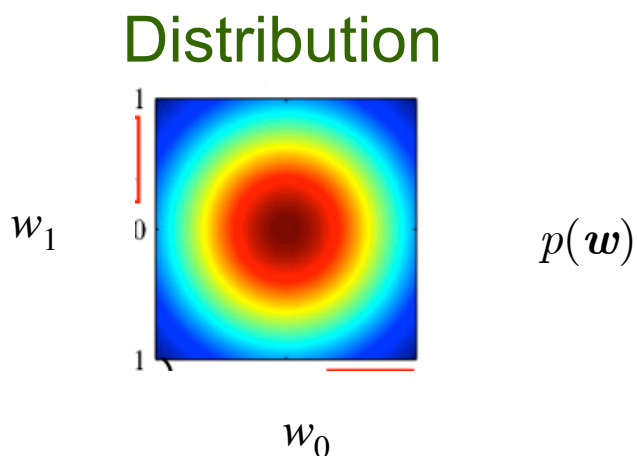


- First choose x_n from $U(x|-1,1)$, then evaluate $f(x_n, \mathbf{w})$
- Add Gaussian noise with st dev 0.2 to get target t_n
 - Precision parameter $\beta = (1/0.2)^2 = 25$
- For prior over \mathbf{w} we choose $\alpha = 2$

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

Sampling $p(\mathbf{w})$ and $p(\mathbf{w}|\mathbf{t})$

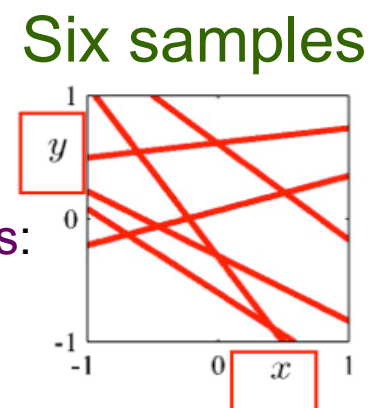
- Each sample represents a straight line in data space (modified by examples)



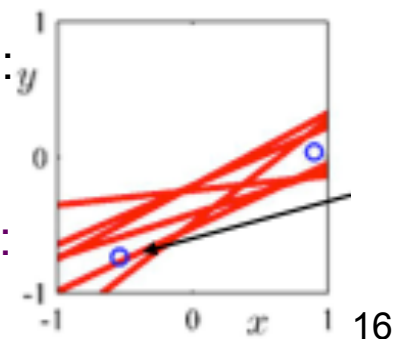
Goal of Bayesian Linear Regression:
Determine $p(\mathbf{w}|\mathbf{t})$

$$y(x, \mathbf{w}) = w_0 + w_1 x$$

With no examples:



With two examples:



Sequential Bayesian Learning

- Since there are only two parameters
 - We can plot prior and posterior distributions in parameter space
- We look at sequential update of posterior

Prior/
Posterior
 $p(\mathbf{w})$
gives $p(\mathbf{w}|t)$

Likelihood $p(t|x, \mathbf{w})$
as function of \mathbf{w}

We are plotting $p(\mathbf{w}|t)$
for a single data point

Six samples
(regression functions)
corresponding to $y(x, \mathbf{w})$
with \mathbf{w} drawn from
posterior

Before data
points
observed

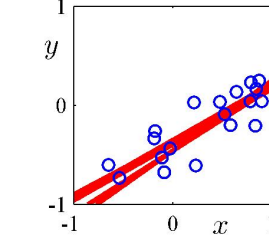
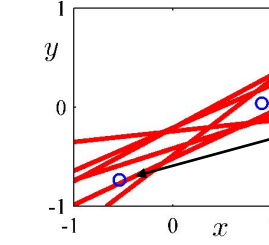
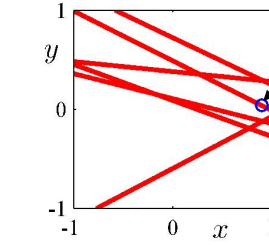
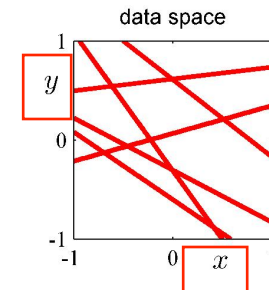
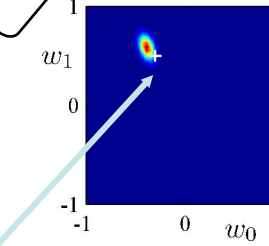
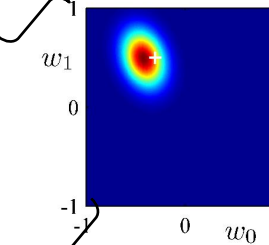
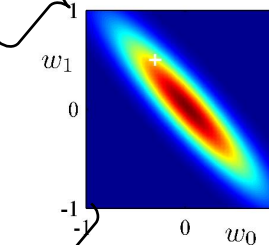
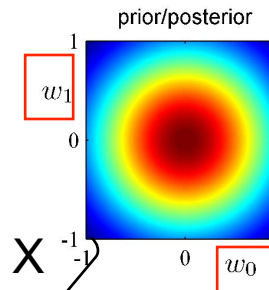
True parameter
Value

After first data point
 (x, t) observed
Band represents values
of w_0, w_1 representing
st lines going near data
point x

Likelihood
for 2nd point
alone

Likelihood for
20th point
alone

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



No
Data
Point

First
Data
Point
 (x_1, t_1)

Second
Data
Point

Twenty
Data
Points

Generalization of Gaussian prior

- The Gaussian prior over parameters is

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

Maximization of posterior $\ln p(\mathbf{w} \mid \mathbf{t})$ is equivalent to minimization of sum of squares error

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

- Other prior yields Lasso and variations:

$$p(\mathbf{w} \mid \alpha) = \left[\frac{q}{2} \left(\frac{\alpha}{2} \right)^{1/q} \frac{1}{\Gamma(1/q)} \right]^M \exp \left(-\frac{\alpha}{2} \sum_{j=1}^M |w_j|^q \right)$$

- $q=2$ corresponds to Gaussian
- Corresponds to minimization of regularized error function

$$\frac{1}{2} \sum_{n=1}^N \left\{ t_n - \mathbf{w}^T \phi(\mathbf{x}_n) \right\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$

Predictive Distribution

- Usually not interested in the value of w itself
- But predicting t for a new value of x

$$p(t|\mathbf{t}, X, \mathbf{x}) \text{ or}$$

$$p(t|\mathbf{t})$$

- Leaving out conditioning variables X and \mathbf{x} for convenience
- Marginalizing over parameter variable w , is the standard Bayesian approach
 - Sum rule of probability
 - We can now write

$$p(t) = \int p(t, w) dw = \int p(t|w)p(w) dw$$

$$p(t | \mathbf{t}) = \int p(t|w)p(w|\mathbf{t}) dw$$

Predictive Distribution with $\alpha, \beta, \mathbf{x}, \mathbf{t}$

- We can predict t for a new value of \mathbf{x} using

$$p(t | \mathbf{t}) = \int p(t | \mathbf{w}) p(\mathbf{w} | \mathbf{t}) d\mathbf{w}$$

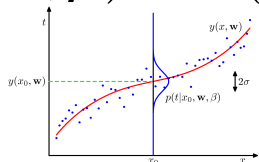
We have left out conditioning variables \mathbf{X} and \mathbf{x} for convenience.
Also we have applied sum rule of probability $p(\mathbf{t}) = \sum_{\mathbf{w}} p(\mathbf{t} | \mathbf{w}) p(\mathbf{w})$

- With explicit dependence on prior parameter α , noise parameter β , & targets in training set \mathbf{t}

$$p(t | \mathbf{t}, \alpha, \beta) = \int \underbrace{p(t | \mathbf{w}, \beta)}_{\text{Conditional of target } t \text{ given weight } \mathbf{w}} \cdot \underbrace{p(\mathbf{w} | \mathbf{t}, \alpha, \beta)}_{\text{posterior of weight } \mathbf{w}} d\mathbf{w}$$

Conditional of target t given weight \mathbf{w}

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



posterior of weight \mathbf{w}

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

where

$$\mathbf{m}_N = \beta S_N \Phi^T \mathbf{t}$$

$$S_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$$

- RHS is a convolution of two Gaussian distributions
 - whose result is the Gaussian:

$$p(t | \mathbf{x}, \mathbf{t}, \alpha, \beta) = N(t | \mathbf{m}_N^T \Phi(\mathbf{x}), \sigma_N^2(\mathbf{x})) \quad \text{where} \quad \sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \Phi(\mathbf{x})^T S_N \Phi(\mathbf{x})$$

Variance of Predictive Distribution

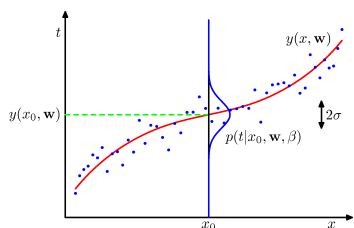
- Predictive distribution is a Gaussian:

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

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

Since noise process and distribution of \mathbf{w} are independent Gaussians their variances are additive

Noise in data



Uncertainty associated with parameters \mathbf{w} :

where $S_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$ is the covariance of $p(\mathbf{w} | \alpha)$

Since $\sigma_{N+1}^2(x) \leq \sigma_N^2(x)$ as no. of samples increases it becomes narrower

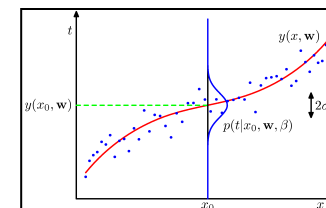
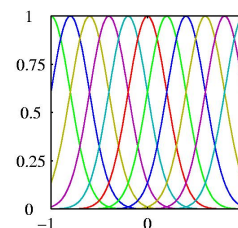
As $N \rightarrow \infty$, second term of variance goes to zero and variance of predictive distribution arises solely from the additive noise parameter β

Example of Predictive Distribution

- Data generated from $\sin(2\pi x)$
- Model: nine Gaussian basis functions

$$y(x, \mathbf{w}) = \sum_{j=0}^8 w_j \phi_j(x) = \mathbf{w}^T \boldsymbol{\phi}(x)$$

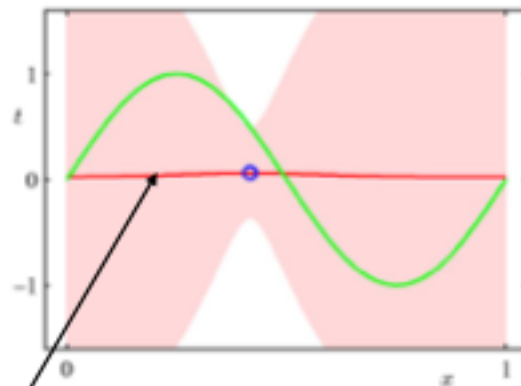
$$\phi_j(x) = \exp\left(-\frac{(x - \mu_j)^2}{2\sigma^2}\right)$$



- Predictive distribution

$$p(t \mid x, \mathbf{t}, \alpha, \beta) = N(t \mid m_N^T \boldsymbol{\phi}(x), \sigma_N^2(x)) \quad \text{where} \quad \sigma_N^2(x) = \frac{1}{\beta} + \boldsymbol{\phi}(x)^T S_N \boldsymbol{\phi}(x)$$

Plot of $p(t|x)$
for one data point
showing mean (red)
and one std dev (pink)



Mean of Predictive Distribution

where $m_N = \beta S_N \Phi^T \mathbf{t}$, $S_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$
and α and β come from assumptions
 $p(\mathbf{w}|\alpha) = N(\mathbf{w}|0, \alpha^{-1} \mathbf{I})$
 $p(t \mid x, \mathbf{w}, \beta) = N(t \mid y(x, \mathbf{w}), \beta^{-1})$

Predictive Distribution Variance

Bayesian prediction:

where we have assumed

Gaussian prior over parameters:

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

Noise model assumed Gaussian:

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

and use design matrix as:

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

$$\Phi = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & & & \\ \vdots & & & \\ \phi_0(x_N) & & & \phi_{M-1}(x_N) \end{bmatrix}$$

Using data from

$\sin(2\pi x)$:

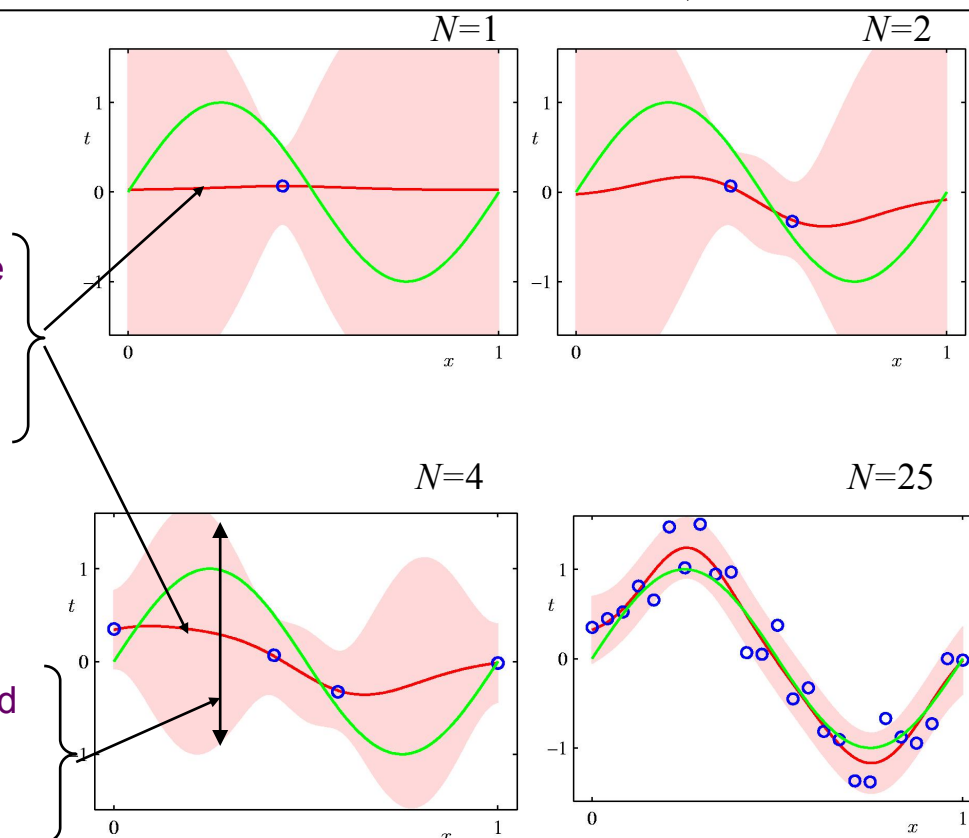
$\sigma_N^2(x)$, std dev of t , is smallest in neighborhood of data points

Uncertainty decreases as more data points are observed

$$p(t | \mathbf{x}, \mathbf{t}, \alpha, \beta) = N(t | m_N^T \phi(\mathbf{x}), \sigma_N^2(x)) \text{ where } \sigma_N^2(x) = \frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x})$$

Mean of the Gaussian Predictive Distribution

One standard deviation from Mean



Plot only shows point-wise predictive variance
To show covariance between predictions at different values of x draw samples from posterior distribution over \mathbf{w} $p(\mathbf{w} | \mathbf{t})$ and plot corresponding functions $y(\mathbf{x}, \mathbf{w})$

Plots of function $y(x, \mathbf{w})$

Draw samples \mathbf{w} from
from posterior
distribution $p(\mathbf{w}|\mathbf{t})$

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

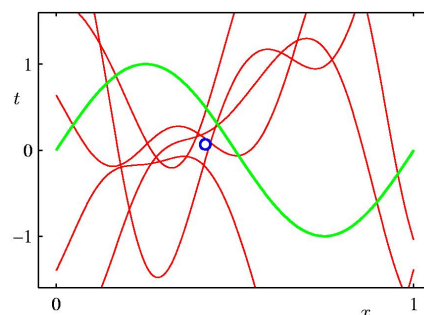
and plot samples

from $y(x, \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(x)$

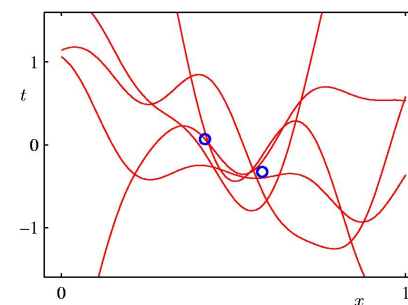
Shows covariance
between predictions at
different values of x

For a given function, for a
pair of x, x' , the values of
 y, y' are determined by $k(x, x')$
which in turn is determined by
the samples

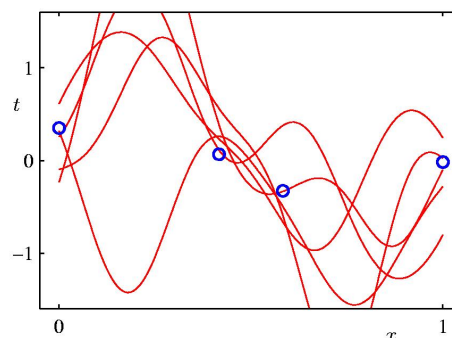
$N=1$



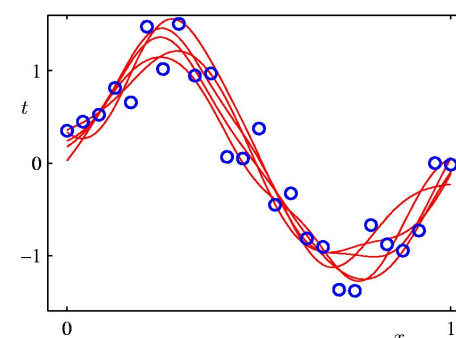
$N=2$



$N=4$



$N=25$



Disadvantage of Local Basis

- Predictive distribution, assuming Gaussian prior

- $p(\mathbf{w} | \alpha) = N(\mathbf{w} | \mathbf{0}, \alpha^{-1}\mathbf{I})$ and Gaussian noise $t = y(\mathbf{x}, \mathbf{w}) + \varepsilon$
- where noise is defined probabilistically as $p(t | \mathbf{x}, \mathbf{w}, \beta) = N(t | y(\mathbf{x}, \mathbf{w}), \beta^{-1})$

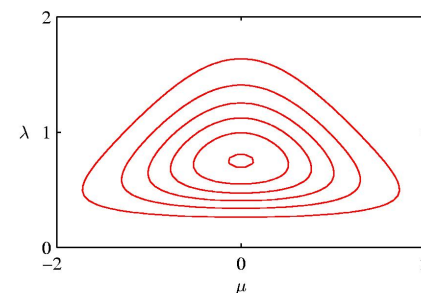
$$p(t | \mathbf{x}, \mathbf{t}, \alpha, \beta) = N(t | m_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x})) \quad \text{where} \quad \sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \phi(\mathbf{x})^T S_N \phi(\mathbf{x}) \quad S_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$$

- With localized basis functions, e.g., Gaussian
 - at regions away from basis function centers, contribution of second term of variance σ_n^2 in will go to zero leaving only noise contribution β^{-1}
 - Model becomes very confident outside of region occupied by basis functions
 - Problem avoided by alternative Bayesian approach of Gaussian Processes

Dealing with unknown β

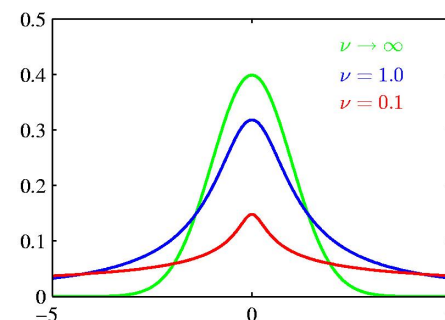
- If both w and β are treated as unknown then we can introduce a conjugate prior distribution $p(w, \beta)$ which is given by a *Gaussian-gamma* distribution

$$p(\mu, \lambda) = N\left(\mu \mid \mu_0, (\beta\lambda)^{-1}\right) \text{Gam}(\lambda \mid a, b)$$



- In this case the predictive distribution is a *Student's t-distribution*

$$St(x \mid \mu, \lambda, \nu) = \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left(\frac{\lambda}{\pi \nu} \right)^{1/2} \left[1 + \frac{\lambda(x - \mu)^2}{\nu} \right]^{-\nu/2 - 1/2}$$



Mean of $p(\mathbf{w}|\mathbf{t})$ has Kernel Interpretation

- Regression function is:

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

- If we take a Bayesian approach with Gaussian prior $p(\mathbf{w}) = N(\mathbf{w}|\mathbf{m}_0, S_0)$ then we have:

– **Posterior** $p(\mathbf{w}|\mathbf{t}) = N(\mathbf{w}|\mathbf{m}_N, S_N)$ where

$$\mathbf{m}_N = S_N (S_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t})$$

$$S_N^{-1} = S_0^{-1} + \beta \Phi^T \Phi$$

- With zero mean isotropic $p(\mathbf{w}|\alpha) = N(\mathbf{w}|0, \alpha^{-1}I)$

$$\begin{aligned} \mathbf{m}_N &= \beta S_N \Phi^T \mathbf{t}, \\ S_N^{-1} &= \alpha I + \beta \Phi^T \Phi \end{aligned}$$

- Posterior mean $\beta S_N \Phi^T \mathbf{t}$ has a kernel interpretation
 - Sets stage for kernel methods and Gaussian processes

Equivalent Kernel

- Posterior mean of \mathbf{w} is $\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}$
 - where $\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi$,
 - \mathbf{S}_0 is the covariance matrix of the prior $p(\mathbf{w})$, β is the noise parameter and Φ is the design matrix that depends on the samples

- Substitute mean value into Regression function

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

- Mean of predictive distribution at point \mathbf{x} is

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

- where $k(\mathbf{x}, \mathbf{x}') = \beta \boldsymbol{\phi}(\mathbf{x})^T \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}')$ is the *equivalent kernel*
- Thus mean of predictive distribution is a linear combination of training set target variables t_n
 - Note: the equivalent kernel depends on input values \mathbf{x}_n from the dataset because they appear in \mathbf{S}_N

Kernel Function

- Regression functions such as

$$y(\mathbf{x}, \mathbf{m}_N) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n$$

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

$$S_N^{-1} = S_0^{-1} + \beta \Phi^T \Phi$$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & & & \\ \vdots & & & \\ \phi_0(\mathbf{x}_N) & & & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

- That take a linear combination of the training set target values are known as *linear smoothers*
- They depend on the input values \mathbf{x}_n from the data set since they appear in the definition of S_N

Example of kernel for Gaussian Basis

Equivalent Kernel

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

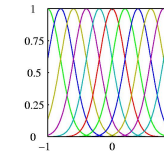
$$S_N^{-1} = S_0^{-1} + \beta \Phi^T \Phi$$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & & & \\ \vdots & & & \\ \phi_0(\mathbf{x}_N) & & & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

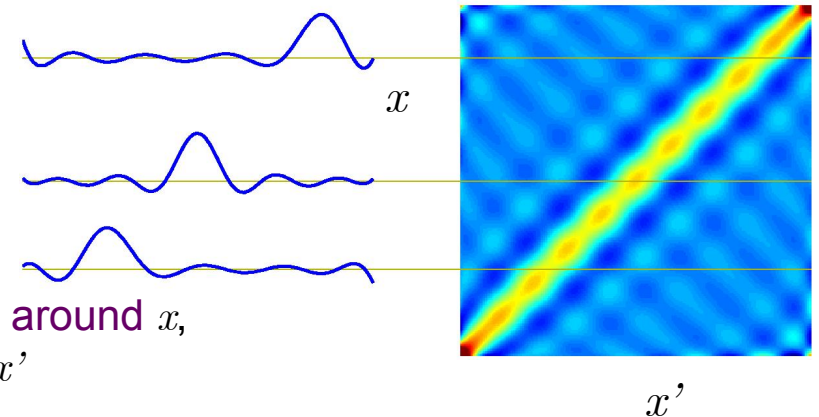
For three values of x the behavior of $k(x, x')$ is shown as a slice

Kernels are localized around x , i.e., peaks when $x = x'$

Gaussian Basis $\phi(x)$



$$\phi_j(x) = \exp\left(-\frac{(x - \mu_j)^2}{2s^2}\right)$$



Kernel used directly in regression.

Mean of the predictive distribution is

$$y(\mathbf{x}, \mathbf{m}_N) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n$$

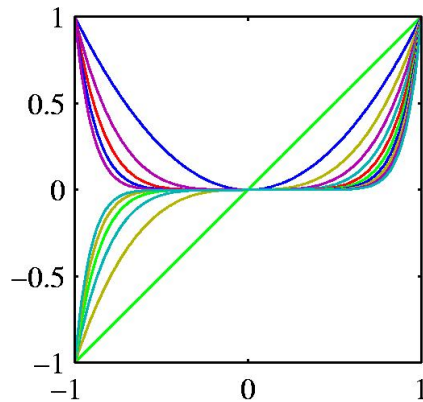
Obtained by forming a weighted combination of target values:

Data points close to x are given higher weight than points further removed from x

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

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

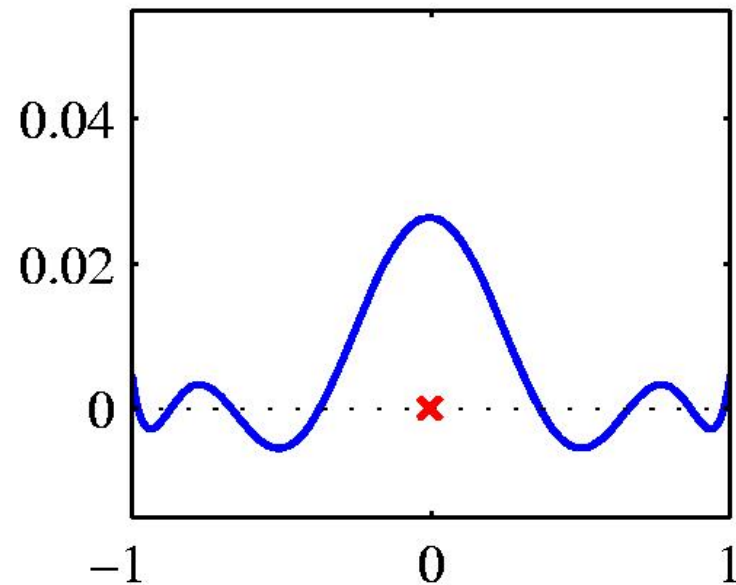
Equivalent Kernel for Polynomial Basis Function



$$\phi_j(x) = x^j$$

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

$$S_N^{-1} = S_0^{-1} + \beta \Phi^T \Phi$$

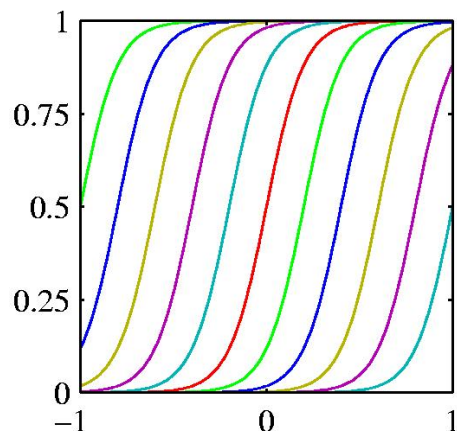


Plotted as a function of x' for $x=0$

Data points close to x are given higher weight than points further removed from x

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

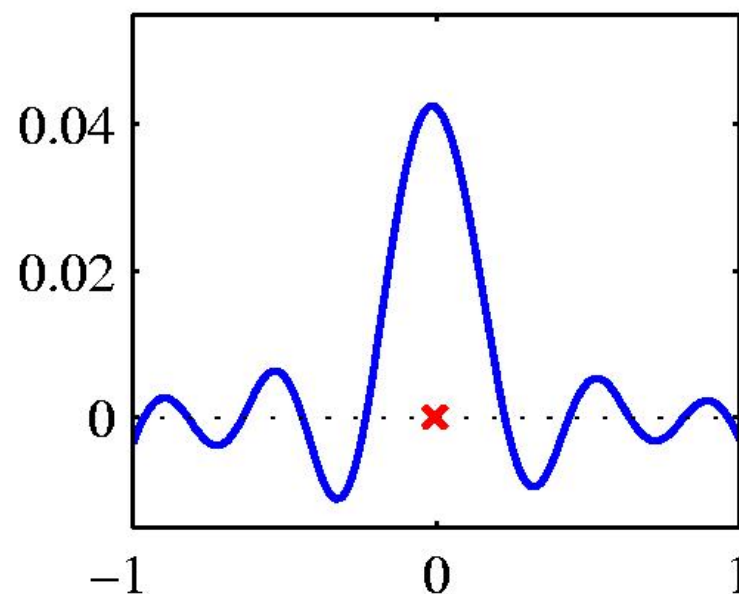
Equivalent Kernel for Sigmoidal Basis Function



$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right) \quad \text{where} \quad \sigma(a) = \frac{1}{1 + \exp(-a)}$$

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

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



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

An important insight:

The value of the kernel function between two points is directly related to the covariance between their target values

$$\begin{aligned}\text{cov} [y(\mathbf{x}), y(\mathbf{x}')] &= \text{cov}[\phi(\mathbf{x})^T \mathbf{w}, \mathbf{w}^T \phi(\mathbf{x}')] \\ &= \phi(\mathbf{x})^T S_N \phi(\mathbf{x}') \\ &= \beta^{-1} k(\mathbf{x}, \mathbf{x}')\end{aligned}$$

where we have used:

$$\begin{aligned}p(\mathbf{w}|\mathbf{t}) &\sim N(\mathbf{w}|\mathbf{m}_N, S_N) \\ k(\mathbf{x}, \mathbf{x}') &= \beta \phi(\mathbf{x})^T S_N \phi(\mathbf{x}')\end{aligned}$$

From the form of the equivalent kernel $k(\mathbf{x}, \mathbf{x}')$

the predictive mean at nearby points $y(\mathbf{x})$, $y(\mathbf{x}')$ will be highly correlated

For more distant pairs correlation is smaller

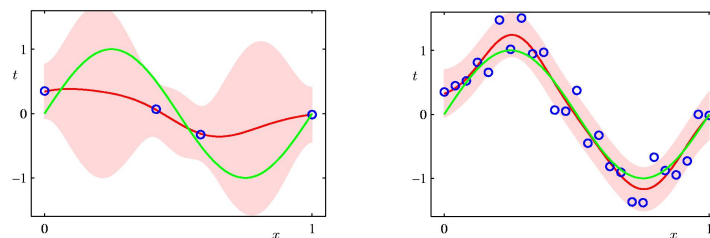
The kernel captures the covariance

Predictive plot vs. Posterior plots

- Predictive distribution

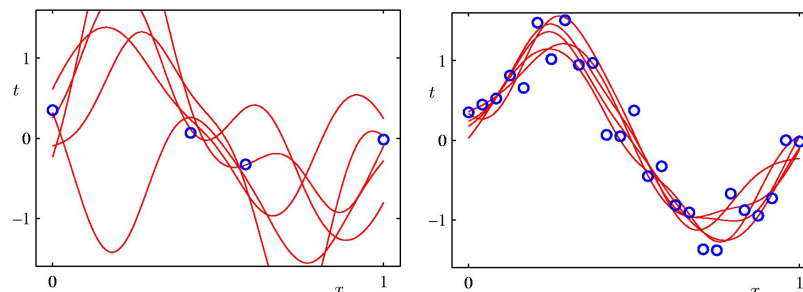
- allows us to visualize pointwise uncertainty in the predictions governed by

$$p(t | x, \mathbf{t}, \alpha, \beta) = N(t | m_N^T \phi(x), \sigma_N^2(x)) \quad \text{where} \quad \sigma_N^2(x) = \frac{1}{\beta} + \phi(x)^T S_N \phi(x)$$



- Drawing samples from posterior $p(\mathbf{w} | \mathbf{t})$

- Plotting corresponding functions $y(x, \mathbf{w})$ we visualize joint uncertainty in the posterior distribution between the y values at two or more x values as governed by the kernel



Directly Specifying Kernel Function

- Formulation of Linear Regression in terms of kernel function suggests an alternative approach to regression:
 - Instead of introducing a set of basis functions, which implicitly determines an equivalent kernel:
 - Directly define kernel functions and use to make predictions for new input \mathbf{x} , given the observation set
- This leads to a practical framework for regression (and classification) called Gaussian Processes

Summing Kernel Values Over samples

- Effective kernel defines weights by which
 - target values combined to make a prediction at \mathbf{x}
- It can be shown that weights sum to one, i.e.,

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

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

$$S_N^{-1} = S_0^{-1} + \beta \Phi^T \Phi$$

- For all values of \mathbf{x}

– This result can be proven intuitively:

- Since $y(\mathbf{x}, \mathbf{m}_N) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n$ summation is equivalent to considering predictive mean $\hat{y}(\mathbf{x})$ for a set of integer data in which $t_n = 1$ for all n
- Provided basis functions are linearly independent, that $N > M$, one of the basis functions is constant (corresponding to the bias parameter), then we can fit training data exactly, and hence $\hat{y}(\mathbf{x}) = 1$

Kernel Function Properties

- Equivalent kernel can be positive or negative

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

- Although it satisfies a summation constraint, the corresponding predictions are not necessarily a convex combination of the training set target variables
- Equivalent kernel satisfies important property shared by kernel functions in general.
 - It can be expressed in the form of an inner product wrt a vector $\psi(\mathbf{x})$ of nonlinear functions:

$$k(\mathbf{x}, \mathbf{z}) = \psi(\mathbf{x})^T \psi(\mathbf{z}) \quad \text{where} \quad \psi(\mathbf{x}) = \beta^{1/2} S_N^{1/2} \phi(\mathbf{x})$$