

Bayesian Linear Regression

Sargur Srihari
srihari@cedar.buffalo.edu

Motivation of Bayesian Approach

- We have seen Linear Regression with M basis functions:

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

- The maximum likelihood solution is

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

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

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 its effect)

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

where

$$E(\mathbf{w}) = E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

$$E_W(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

 - 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 better handled using Bayesian approach

Bayesian Approach

- Posterior is proportional to Likelihood x 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

- 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

Prior Distribution for Parameters

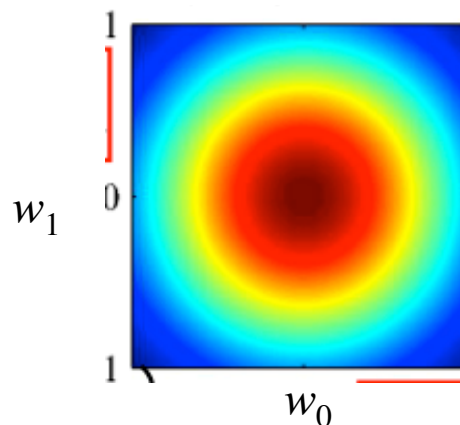
Assume multivariate Gaussian prior for \mathbf{w} (which has components w_0, \dots, w_{M-1})

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

with mean \mathbf{m}_0 and covariance matrix S_0

If we choose $S_0 = \alpha^{-1} \mathbf{I}$ 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*)

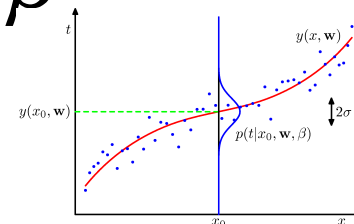


Likelihood of Data

- Assume noise precision parameter β

$$t = y(\mathbf{x}, \mathbf{w}) + \varepsilon$$

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



- Likelihood $p(\mathbf{t}|\mathbf{w})$ with Gaussian noise has an exponential form

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

This is the probability of data \mathbf{t} given the parameters \mathbf{w} and input X

Posterior Distribution of Parameters

- Posterior is proportional to Likelihood x Prior

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

– where likelihood function is

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

Product of Gaussians

- Multiplying by Gaussian prior, $p(\mathbf{w}) = N(\mathbf{w} | \mathbf{m}_0, S_0)$ posterior is also Gaussian, written directly as

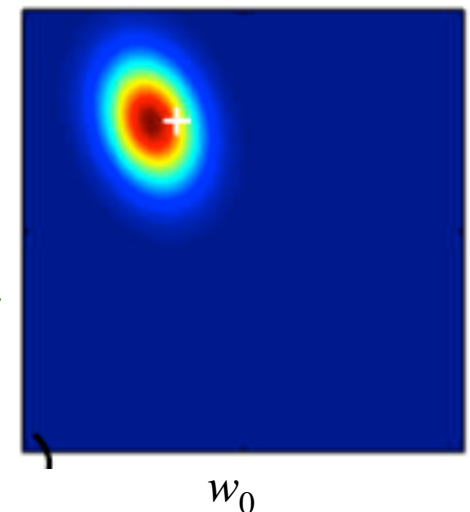
$$p(\mathbf{w} | \mathbf{t}) = N(\mathbf{w} | \mathbf{m}_N, S_N)$$

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



Properties of Posterior

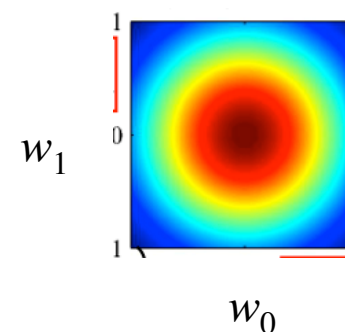
1. Since posterior $p(\mathbf{w}|\mathbf{t}) = N(\mathbf{w}|\mathbf{m}_N, S_N)$ is Gaussian its mode coincides with its mean
 - Thus maximum posterior weight vector is $\mathbf{w}_{\text{MAP}} = \mathbf{m}_N$
2. For infinitely broad prior $S_0 = \alpha^{-1} \mathbf{I}$ with $\alpha \rightarrow 0$
 - Then mean \mathbf{m}_N reduces to the maximum likelihood value
$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$
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 $p(\mathbf{w})$

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

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

Single precision parameter



- Corresponding posterior distribution is

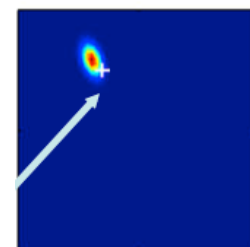
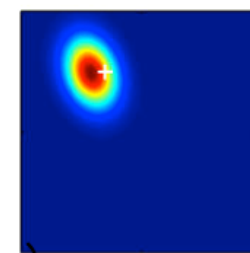
$$p(\mathbf{w} | \mathbf{t}) = N(\mathbf{w} | \mathbf{m}_N, S_N)$$

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:

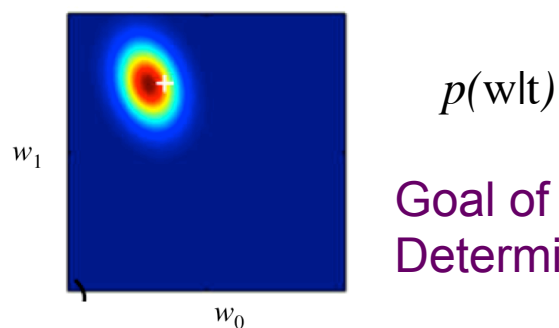
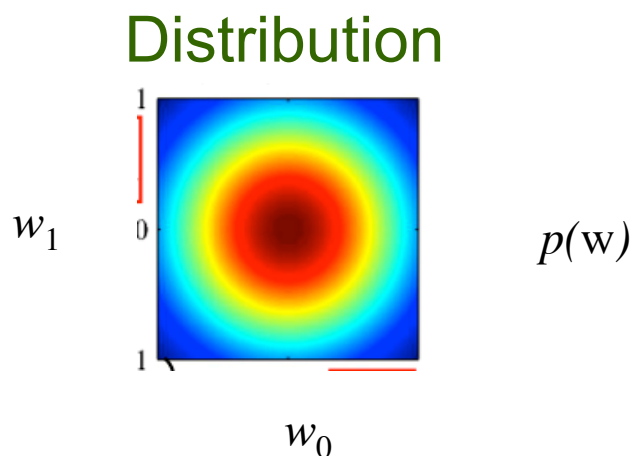
β is noise precision and
 α is variance of parameter \mathbf{w} in prior



Point
Estimate
with
infinite
samples

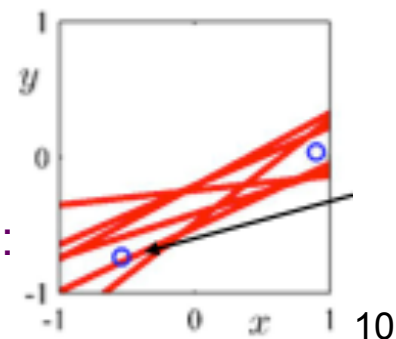
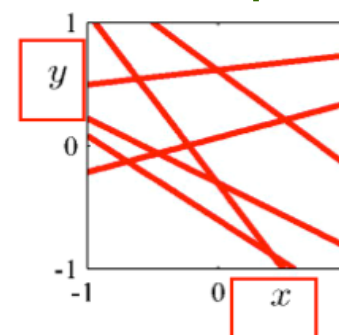
Sampling $p(w)$ and $p(w/t)$

- Each sample represents a straight line in data space



Goal of Bayesian Linear Regression:
Determine $p(w|t)$

Six samples



Interesting Relationship to MLE

- Since

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

- Log of Posterior is

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

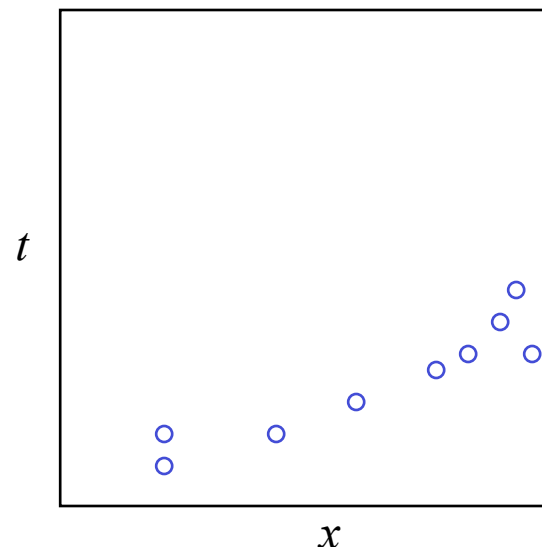
- Thus Maximization of posterior distribution wrt \mathbf{w} is equivalent to minimization of sum-of-squares error function

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

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

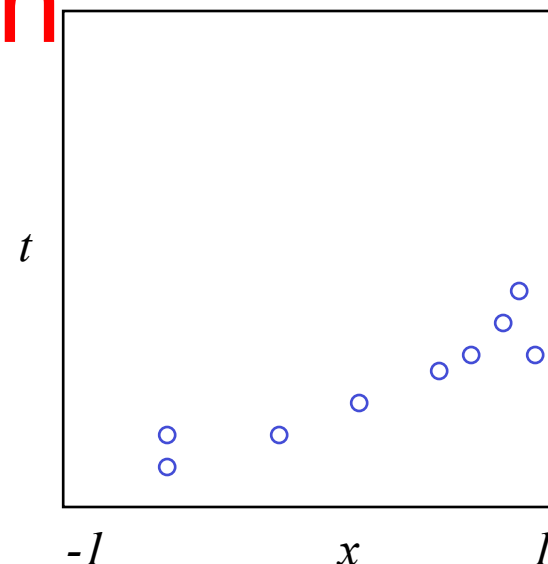


Data Generation

- Synthetic data generated from

$f(x, \mathbf{w}) = w_0 + w_1 x$ with
parameter values

$w_0 = -0.3$ 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 obtain target t_n
 - Precision parameter $\beta = (1/0.2)^2 = 25$

- For prior over \mathbf{w} we choose $\alpha = 2$

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

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

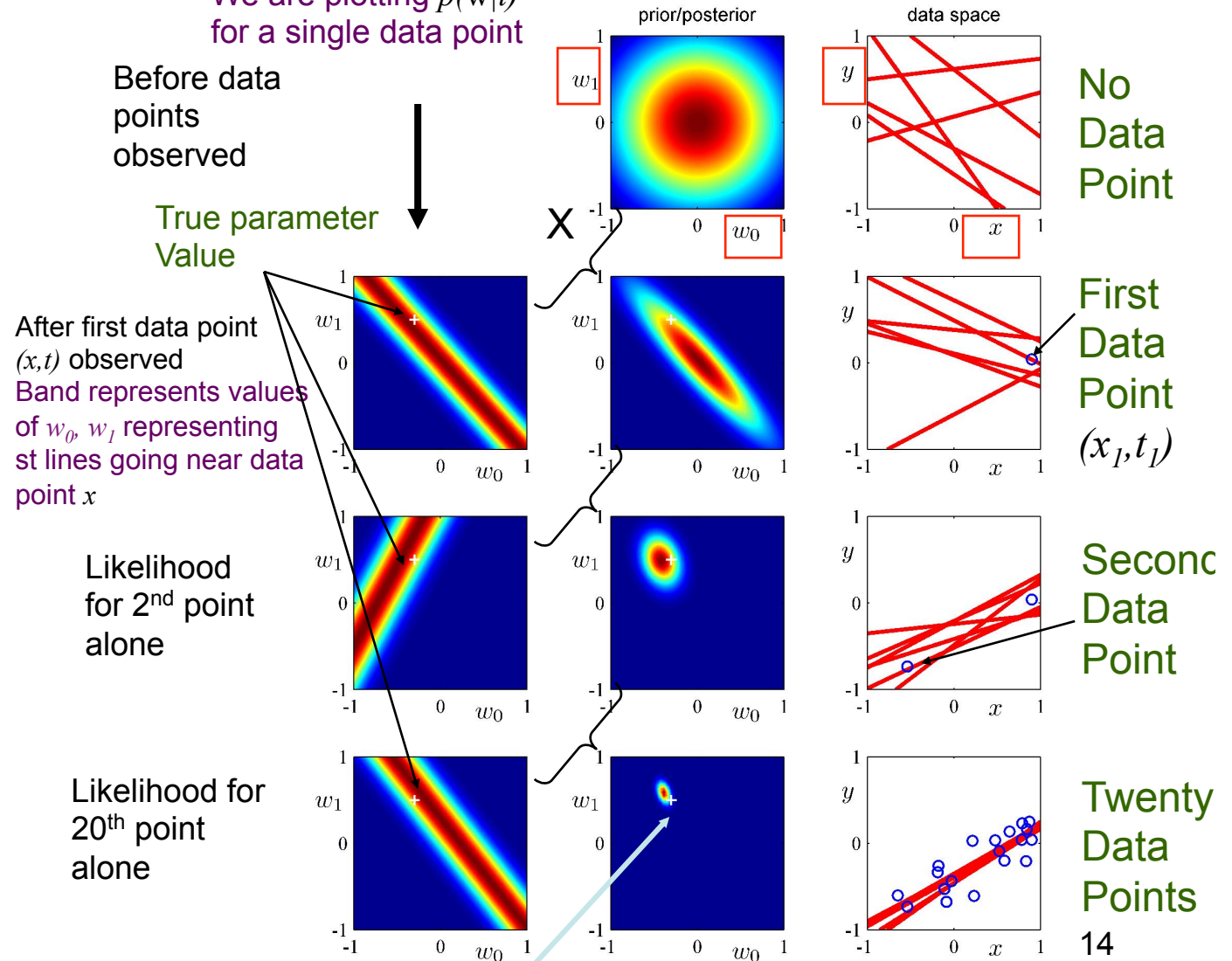
With infinite points posterior is a delta function centered at true parameters (white cross)

Likelihood $p(t|x, w)$ as function of w

We are plotting $p(w|t)$ for a single data point

Prior/
Posterior
 $p(w)$
gives $p(w|t)$

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



Generalization of Gaussian prior

- The Gaussian prior over parameters is

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

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

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

- Generalization of Gaussian prior

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

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

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

Predictive Distribution

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

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

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

- Leaving out conditioning variables \mathbf{X} and \mathbf{x} for convenience
- Marginalizing over parameter variable \mathbf{w} , is the standard Bayesian approach
 - From sum rule

$$\begin{aligned} p(t | \mathbf{t}) &= \int p(t, \mathbf{w}) d\mathbf{w} \\ &= \int p(t|\mathbf{w})p(\mathbf{w}|\mathbf{t}) d\mathbf{w} \end{aligned}$$

Predictive Distribution with α, β

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

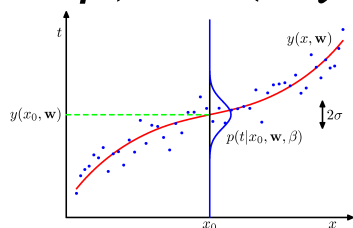
$$p(t | \mathbf{t}) = \int p(t | \mathbf{w}) p(\mathbf{w} | \mathbf{t}) d\mathbf{w}$$

- Introducing α characterizing prior, and β characterizing noise

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

conditional distribution of target variable

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



posterior weight distribution

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

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

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

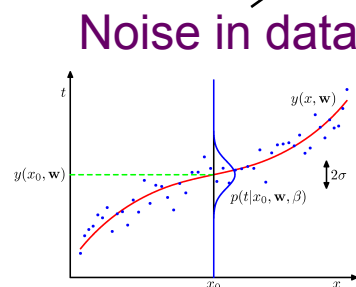
- RHS is convolution of two Gaussian distributions
 - whose result is a Gaussian

Variance of Predictive Distribution

- Predictive distribution takes the form

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

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



Uncertainty associated with parameters \mathbf{w}

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

Properties of Variance of Predictive Distribution

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \varphi(\mathbf{x})^T S_N \varphi(\mathbf{x})$$

- It can be shown that $\sigma_{N+1}^2(\mathbf{x}) \leq \sigma_N^2(\mathbf{x})$
- In the limit as $N \rightarrow \infty$,
 - The second term goes to zero
 - Variance is solely due to additive noise

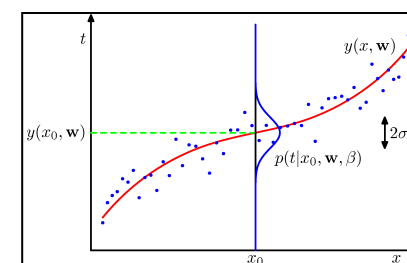
Example of Predictive Distribution

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

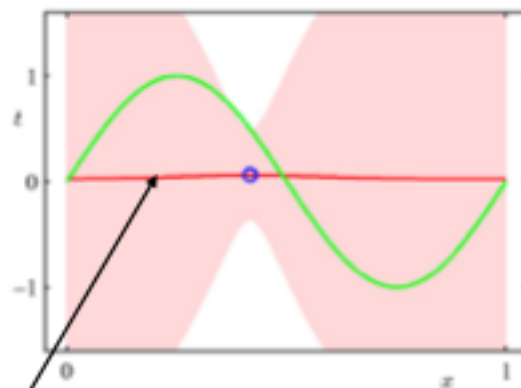
$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^8 w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

- Predictive distribution

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



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



Mean of Predictive Distribution

$$m_N = \beta S_N \boldsymbol{\Phi}^T \mathbf{t}$$

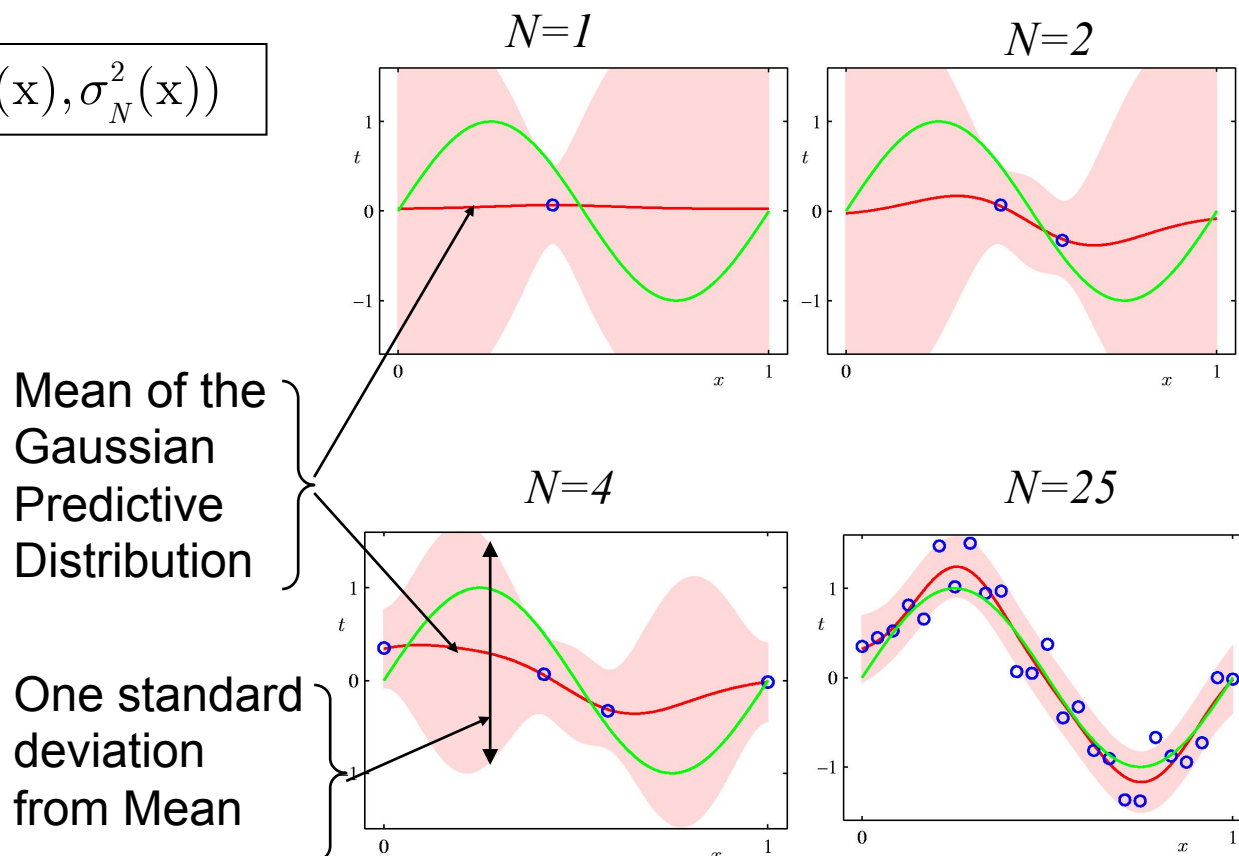
$$S_N^{-1} = \alpha I + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

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

Predictive Distribution (Sinusoidal Data)

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

- Std dev of t is smallest in neighborhood of data points
- Uncertainty decreases as more data points are observed



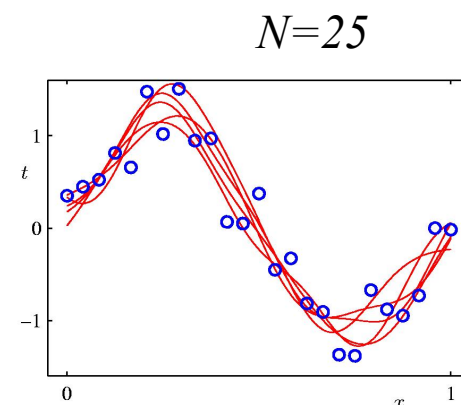
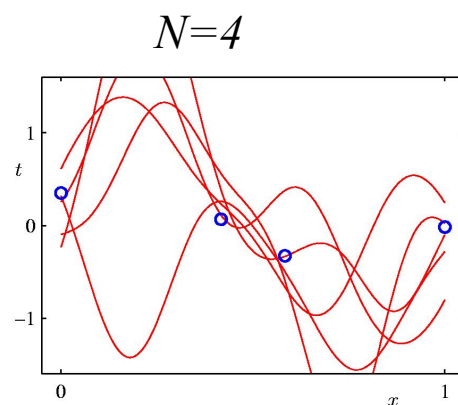
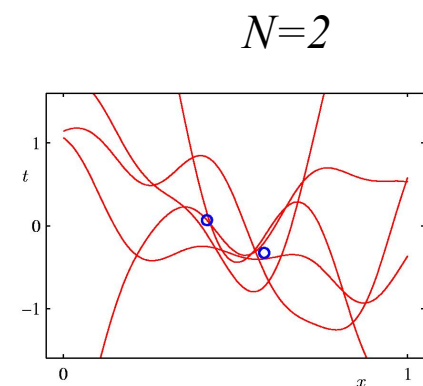
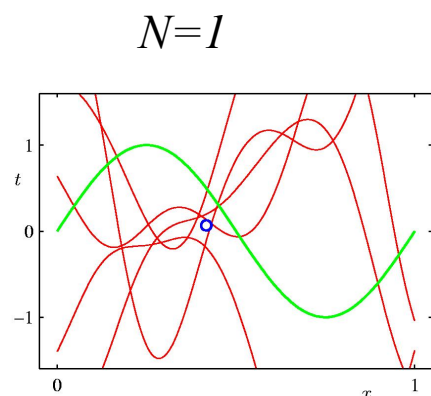
Plot only shows point-wise predictive variance
 To show covariance between predictions at different values of x draw samples from posterior distribution over \mathbf{w} $p(\mathbf{w}|\mathbf{t})$ and draw samples of $y(\mathbf{x}, \mathbf{w})$

Plots of function $y(x, w)$

- Draw samples from from posterior distribution $p(w/t)$

$$p(w|t) = N(w|m_N, S_N)$$

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



These curves represent distribution of the regression function

Kernel Interpretation

- Regression function is:

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

- Bayesian Soln. with Gaussian prior $p(\mathbf{w}) = N(\mathbf{w} | \mathbf{m}_0, S_0)$

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

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

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

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

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

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

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

Equivalent Kernel

- Posterior mean of w is $m_N = \beta S_N \Phi^T t$
 – where S_N is the design matrix
- Substitute mean value into Regression function

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

- Mean of predictive distribution at point x

$$\begin{aligned} y(x, m_N) &= m_N^T \phi(x) = \beta \phi(x)^T S_N \Phi^T t \\ &= \sum_{n=1}^N \beta \phi(x)^T S_N \phi(x_n) t_n \\ &= \sum_{n=1}^N k(x, x_n) t_n \end{aligned}$$

- Where $k(x, x') = \beta \phi(x)^T S_N \phi(x')$ is the *equivalent kernel*
- Mean of predictive distribution is a linear combination of training set target variables

Kernel Function

- Regression functions such as

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

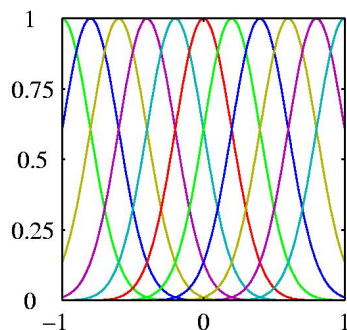
$$k(\mathbf{x}, \mathbf{x}') = \beta \phi(\mathbf{x})^T S_N \phi(\mathbf{x}')$$

- 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

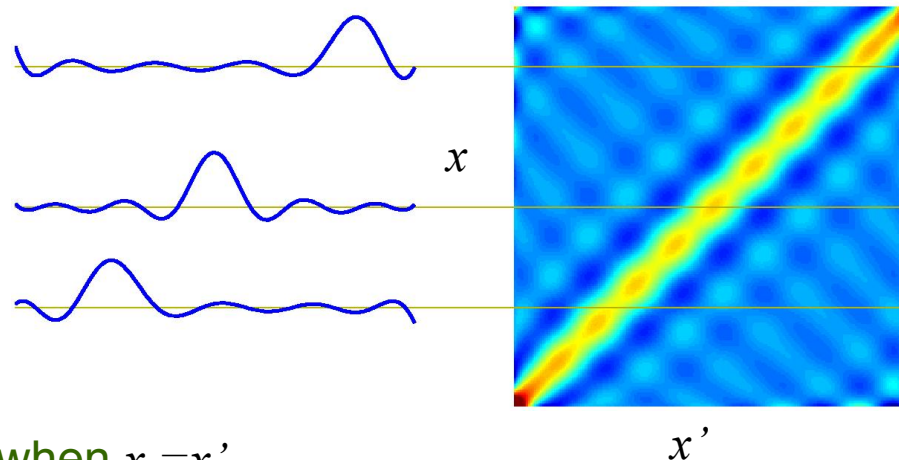
Equivalent kernel for Gaussian Basis

Gaussian Basis $\phi(x)$

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



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



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

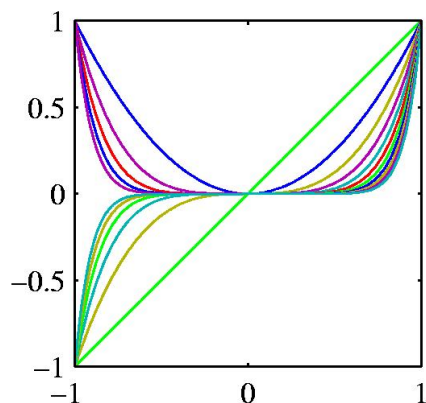
They are localized around x , i.e., peaks when $x = x'$
Local evidence is weighted more than distant evidence

Kernel used directly in regression

$$y(x, m_N) = \sum_{n=1}^N k(x, x_n) t_n$$

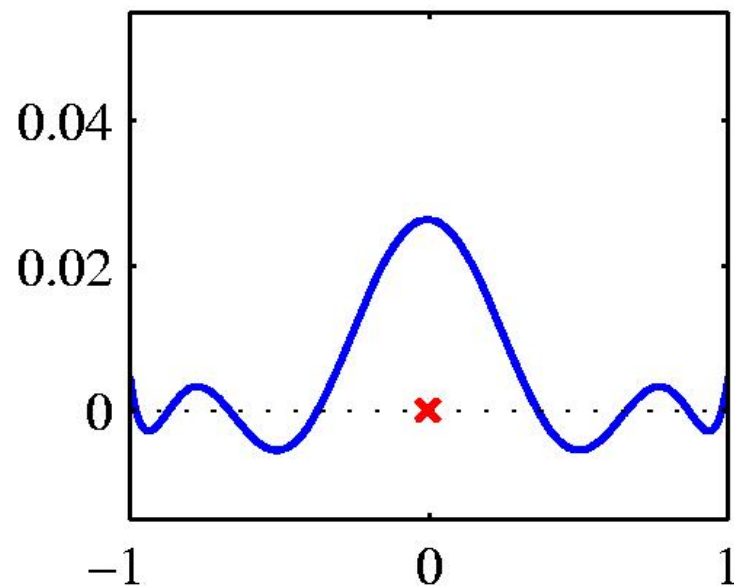
Data set used to
generate kernel were
200 values of
 x equally spaced
in $(-1, 1)$

Equivalent Kernel for Polynomial Basis Function



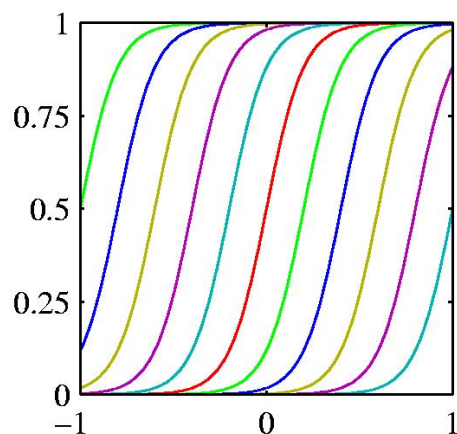
$$\phi_j(x) = x^j$$

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



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

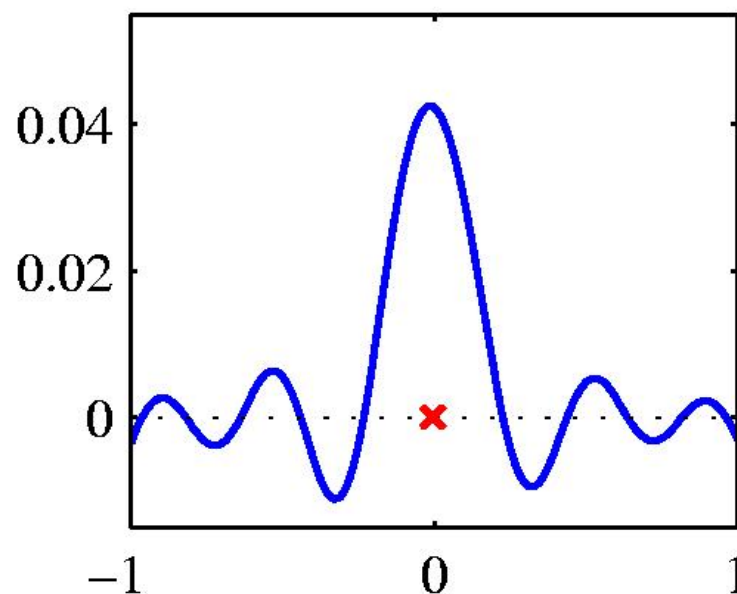
Equivalent Kernel for Sigmoidal Basis Function



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

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

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



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

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

The kernel captures the covariance

From the form of the equivalent kernel,
the predictive mean at nearby points $y(\mathbf{x}) , y(\mathbf{x}')$ 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 $\sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) = 1$
- Equivalent kernel satisfies important property:

$$k(\mathbf{x}, \mathbf{z}) = \psi(\mathbf{x})^T \psi(\mathbf{z})$$

– where

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

Probabilistic Linear Regression Revisited

- Re-derive the predictive distribution by working in terms of distribution over functions $y(\mathbf{x}, \mathbf{w})$
 - It will provide a specific example of a Gaussian Process
- Consider model with M fixed basis functions

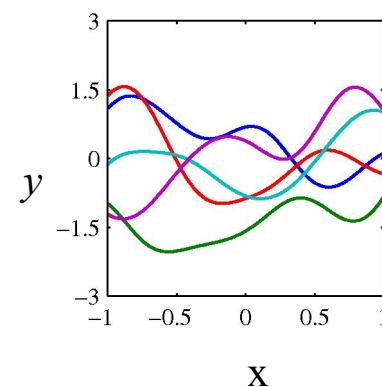
$$y(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

where \mathbf{x} is the input vector and \mathbf{w} is the M -dimensional weight vector

- Assume a Gaussian distribution of weight \mathbf{w}

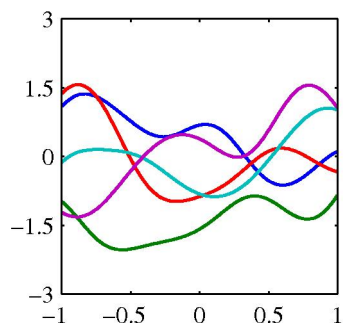
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | 0, \alpha^{-1} \mathbf{I})$$

- Probability distribution over \mathbf{w} induces a probability distribution over functions $y(\mathbf{x})$



Probabilistic Linear Regression is a GP

- We wish to evaluate $y(x)$ at training points x_1, \dots, x_N
- Our interest is the joint distribution of values $y = [y(x_1), \dots, y(x_N)]$ Since $y(x) = w^T \phi(x)$ we can write $y = \Phi w$
 - where Φ is the design matrix with elements $\Phi_{nk} = \phi_k(x_n)$
 $N \times M$ times $M \times 1$ yields a $N \times 1$
- Since y is a linear combination of elements of w which are Gaussian distributed as $p(w) = \mathcal{N}(w|0, \alpha^{-1}I)$ Covariance expressed by Kernel function
 y is itself Gaussian with



mean:

$$E[y] = \Phi E[w] = 0 \quad \text{and}$$

variance:

$$\text{Cov}[y] = E[yy^T] = \Phi E[ww^T] \Phi^T = (1/\alpha) \Phi \Phi^T = K$$

where K is the $N \times N$ Gram Matrix with elements

$$K_{nm} = k(x_n, x_m) = (1/\alpha) \phi(x_n)^T \phi(x_m)$$

and $k(x, x')$ is the kernel function

Thus Gaussian distributed weight vector induces a Gaussian joint distribution over training samples

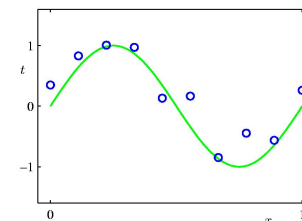
General definition of GP

- We saw a particular example of a Gaussian process
 - Linear regression using the parametric model $y(x)=w^T\phi(x)$
 - With samples $y=[y(x_1),\dots,y(x_N)]$
 - Assume $p(w) = \mathcal{N}(w|0,\alpha^{-1}I)$
 - Then $E[y]=0$ and $\text{Cov}[y]=K$, the Gram matrix which is equivalent to pairwise kernel values
- More generally, a Gaussian process is a probability distribution over functions $y(x)$
 - Such that the set of values of $y(x)$ evaluated at arbitrary points x_1,\dots,x_N jointly have a multivariate Gaussian distribution
 - For a single input x_1 , output y_1 is univariate Gaussian. For two inputs x_1, x_2 the output y_1, y_2 is bivariate Gaussian, etc

Stochastic Process

- A *stochastic process* $y(\mathbf{x})$ is specified by giving the joint probability distribution for any finite set of values $y(\mathbf{x}_1), \dots, y(\mathbf{x}_N)$ in a consistent manner
 - The random variables typically develop over time
- A Gaussian process is a stochastic process which is Gaussian
- When input \mathbf{x} is 2-D it is also known as a *Gaussian Random Field*

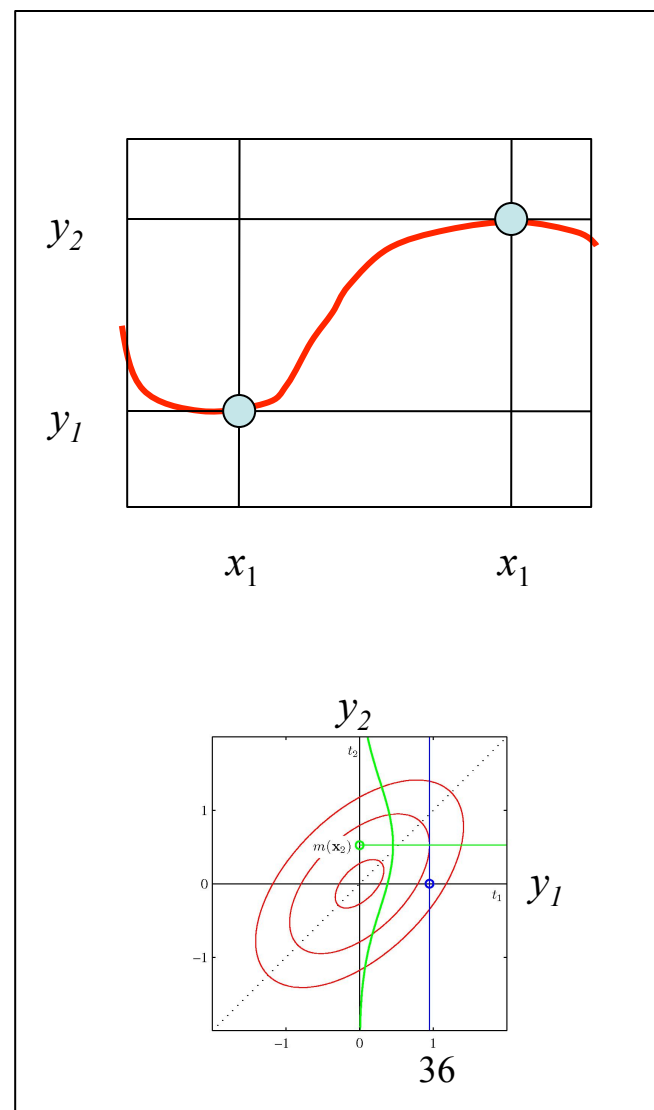
Parametric Regression vs GP Regression



- In parametric regression we have several samples from which we learn the parameters
 - In Gaussian processes we view the samples as one huge input that has a Gaussian distribution
 - with a mean and a covariance matrix
 - A **Gaussian process** is a stochastic process X_t , $t \in T$, for which any finite linear combination of samples has a joint Gaussian Distribution
 - any linear functional applied to the sample function X_t will give a normally distributed result. Notation-wise, one can write $X \sim \text{GP}(m, K)$, meaning the random function X is distributed as a GP with mean function m and
- 35 covariance function K .

Gaussian Process with Two Samples

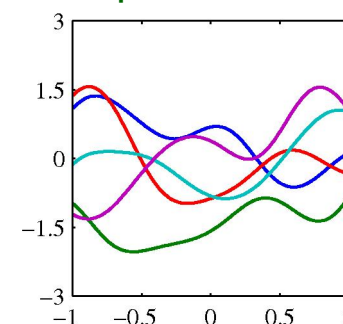
- Let y be a function (curve)
 - of a one-dimensional variable x
- We take two samples y_1 and y_2 corresponding to x_1 and x_2
- Assume they have a bivariate Gaussian distribution $p(y_1, y_2)$
- Each point from this distribution
 - has an associated probability
 - It also defines a function $y(x)$
 - Assuming that two points are enough to define a curve
- More than two points will be needed to define a curve
 - Which leads to a higher dimensional probability distribution



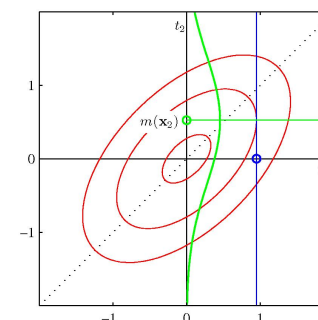
Gaussian Process with N samples

- A process that generates samples over time $\{y_1, \dots, y_N\}$ is a GP iff every set of samples $\mathbf{Y} = [y_1, \dots, y_N]$ is a vector-valued Gaussian random variable
- We define a distribution over all functions
 - Constrained by the samples
- The samples have a joint Gaussian distribution in N -dimensional space

Regression:
Possible functions
constrained by N
samples



Case of two samples
 t_1 and t_2 that are
bivariate Gaussian



Specifying a Gaussian Process

- Key point about Gaussian Stochastic Processes
 - Joint distribution over N variables y_1, \dots, y_N is completely specified by the second-order statistics,
 - i.e., mean and covariance
- With mean zero, it is completely specified by covariance of $y(\mathbf{x})$ evaluated at any two values of \mathbf{x} which is given by a kernel function

$$E[y(\mathbf{x}_n) y(\mathbf{x}_m)] = k(\mathbf{x}_n, \mathbf{x}_m)$$

- For the Gaussian Process defined by the linear regression model $y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$ with prior $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | 0, \alpha^{-1} \mathbf{I})$ the kernel function is

$$K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m) = (1/\alpha) \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$$

Defining a Kernel Directly

- GP can also be specified directly by choice of kernel function (instead of indirectly by basis function)
- Samples of functions drawn from Gaussian processes for two different choices of kernel functions are shown next

Samples from Gaussian Processes for two kernels

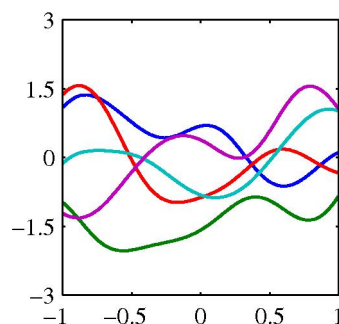
Functions are drawn from

Gaussian processes

Note that each sample is a function

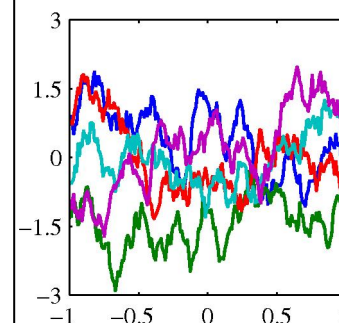
Gaussian Kernel

$$k(x, x') = \exp(-\|x - x'\|^2 / 2\sigma^2)$$



Exponential Kernel

$$k(x, x') = \exp(-\theta \|x - x'\|)$$



Ornstein-Uhlenbeck process for Brownian motion

Prob Distributions and Stochastic Processes

- A probability distribution describes the distributions of scalars (x) and vectors (\mathbf{x})
- A stochastic process describes distribution of functions $f(x)$
 - Can think of a function as a very long vector
 - Each entry in the vector is $f(x)$ which is the value of function at x
- A Gaussian process is a generalization of Gaussian distributions where it is describing the distribution of functions
- Since the vector is infinite-dimensional we constrain it to only those points x corresponding to training and test points

GP for Regression (Direct Approach)

- We specify Gaussian Process directly over functions
 - Abandon approach of defining a distribution over weights w
- Take into account noise on observed target values as
 - $t_n = y_n + \epsilon_n$ where $y_n = y(x_n)$
 - Noise process has a Gaussian distribution $p(t_n|y_n) = N(t_n|y_n, \beta^{-1})$
- Note that target t_n is output y_n corrupted by noise
- Defining $\mathbf{t} = (t_1, \dots, t_N)^T$ our goal is to determine a distribution $p(\mathbf{t})$
 - Which is a distribution over functions

GP for Regression (Direct Approach)

- Assuming noise is independent for each data point

joint distribution of $\mathbf{t} = (t_1, \dots, t_N)^T$ on values $\mathbf{y} = (y_1, \dots, y_N)^T$ is

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

- From definition of GP, marginal distribution of \mathbf{y} is given by a Gaussian with zero mean, covariance matrix given by Gram matrix \mathbf{K}

$$p(\mathbf{y}) = N(\mathbf{y}|0, \mathbf{K})$$

- Where kernel function that determines \mathbf{K} is chosen to express:
property that for points $\mathbf{x}_n, \mathbf{x}_m$ that are similar corresponding values $y(\mathbf{x}_n), y(\mathbf{x}_m)$ will be more strongly correlated than for dissimilar points
 \mathbf{K} depends on application

Regression: Marginal Distribution

- From distributions $p(t|y)$ and $p(y)$ we can get marginal $p(t)$ conditioned on input values x_1, \dots, x_N
- Applying Sum rule and product rule

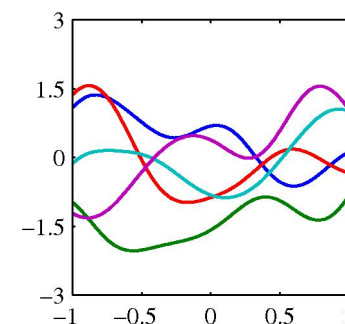
$$p(t) = \int p(t|y)p(y)dy$$

$$= N(t|0, C)$$

From the result that
when $p(y)$ and $p(t|y)$ are Gaussian
 $p(t)$ is also Gaussian

where covariance matrix C has the elements

$$C(x_n, x_m) = \underbrace{k(x_n, x_m)}_{\text{Due to } y(x)} + \underbrace{\beta^{-1} \delta_{nm}}_{\text{Due to } \varepsilon}$$



- The two Gaussian sources of randomness, $y(x)$ and ε are independent, so their covariances simply add

δ specified by ε

Widely used kernel function for GP

- Exponential of a quadratic form
with addition of constant and linear terms

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left\{-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2\right\} + \theta_2 + \theta_3 \underbrace{\mathbf{x}_n^T \mathbf{x}_m}$$

Corresponds to a parametric model that is
a linear function of input variables

- Samples from this prior are plotted for various values of the parameters θ_0 , θ_1 , θ_2 , and θ_3