#### **EECS 545: Machine Learning**

# Lecture 9. Kernel methods: Gaussian Processes

Honglak Lee 2/7/2011



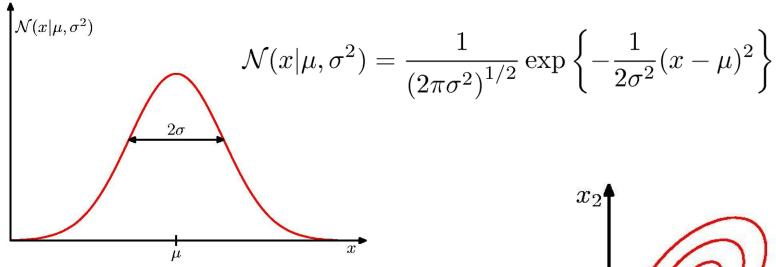


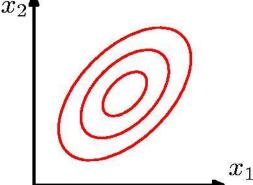
#### Outline

- More about Multivariate Gaussians
  - Marginal and conditional distributions
- Bayesian Linear Regression
- Gaussian Processes
  - GP for Regression

# Multivariate Gaussians – marginal and conditional distributions

#### The Gaussian Distribution





$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right\}$$

# Multivariate Gaussian (mean)

$$\mathbb{E}[\mathbf{x}] = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\mathbf{\Sigma}|^{1/2}} \int \exp\left\{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right\} \mathbf{x} \, d\mathbf{x}$$
$$= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\mathbf{\Sigma}|^{1/2}} \int \exp\left\{-\frac{1}{2} \mathbf{z}^{\mathrm{T}} \mathbf{\Sigma}^{-1} \mathbf{z}\right\} (\mathbf{z} + \boldsymbol{\mu}) \, d\mathbf{z}$$

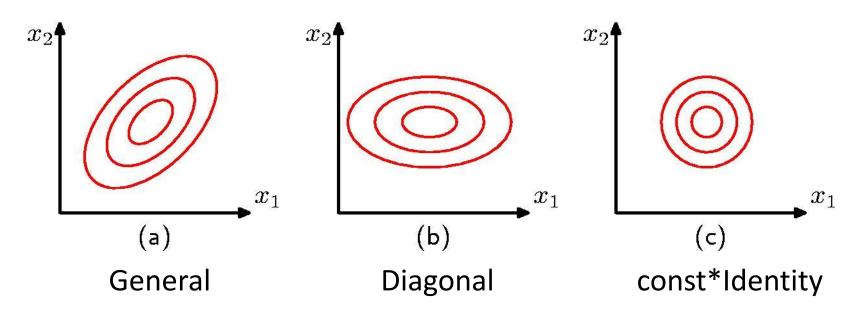
Thanks to the anti-symmetry of **z**, we get:

$$\mathbb{E}[\mathbf{x}] = oldsymbol{\mu}$$

# Multivariate Gaussian (covariance)

$$\mathbb{E}[\mathbf{x}\mathbf{x}^{\mathrm{T}}] = \boldsymbol{\mu}\boldsymbol{\mu}^{\mathrm{T}} + \boldsymbol{\Sigma}$$
 $\operatorname{cov}[\mathbf{x}] = \mathbb{E}\left[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^{\mathrm{T}}\right] = \boldsymbol{\Sigma}$ 

#### Types of covariance matrices:



#### Partitioned Gaussian Distributions

Multivariate Gaussian distribution for x

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

- Partitioning x into x<sub>a</sub> and x<sub>b</sub>.
  - Mean and covariance

$$\mathbf{x} = egin{pmatrix} \mathbf{x}_a \ \mathbf{x}_b \end{pmatrix} \qquad \qquad oldsymbol{\mu} = egin{pmatrix} oldsymbol{\mu}_a \ oldsymbol{\mu}_b \end{pmatrix} \qquad \qquad oldsymbol{\Sigma} = egin{pmatrix} oldsymbol{\Sigma}_{aa} & oldsymbol{\Sigma}_{ab} \ oldsymbol{\Sigma}_{ba} & oldsymbol{\Sigma}_{bb} \end{pmatrix}$$

Precision matrix

$$oldsymbol{\Lambda} \equiv oldsymbol{\Sigma}^{-1} \qquad \qquad oldsymbol{\Lambda} = egin{pmatrix} oldsymbol{\Lambda}_{aa} & oldsymbol{\Lambda}_{ab} \ oldsymbol{\Lambda}_{ba} & oldsymbol{\Lambda}_{bb} \end{pmatrix}$$

#### Partitioned Conditionals and Marginals

Conditional distribution

$$p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b})$$
 $\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Lambda}_{aa}^{-1} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba}$ 
 $\boldsymbol{\mu}_{a|b} = \boldsymbol{\Sigma}_{a|b} \left\{ \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \right\}$ 
 $= \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b)$ 
 $= \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} (\mathbf{x}_b - \boldsymbol{\mu}_b)$ 

Marginal distribution

$$p(\mathbf{x}_a) = \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b$$
$$= \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})$$

#### Derivation of conditional Gaussian

From the exponent of the Gaussian

$$\begin{split} -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) &= \\ -\frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^{\mathrm{T}} \boldsymbol{\Lambda}_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^{\mathrm{T}} \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ -\frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^{\mathrm{T}} \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^{\mathrm{T}} \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b). \end{split}$$

Note that general Gaussians have the following form:

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2} \mathbf{x}^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} \mathbf{x} + \mathbf{x}^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \mathrm{const}$$

• Treat  $x_b$ 's as constant and rearrange terms for  $x_a$ 's.

#### Derivation of conditional Gaussian

Second order term:

$$-\frac{1}{2}\mathbf{x}_{a}^{\mathrm{T}}\boldsymbol{\Lambda}_{aa}\mathbf{x}_{a} \qquad \qquad \boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Lambda}_{aa}^{-1}.$$

First order term:

$$\mathbf{x}_a^{\mathrm{T}} \left\{ \mathbf{\Lambda}_{aa} \boldsymbol{\mu}_a - \mathbf{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \right\}$$

$$-\operatorname{From} \Sigma_{a|b}^{-1} \mu_{a|b} = \Lambda_{aa} \mu_a - \Lambda_{ab} (x_b - \mu_b)$$

$$\mu_{a|b} = \Sigma_{a|b} \{ \Lambda_{aa} \mu_a - \Lambda_{ab} (\mathbf{x}_b - \mu_b) \}$$

$$= \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (\mathbf{x}_b - \mu_b)$$

#### Derivation of conditional Gaussian

Expressing in terms of covariance matrix

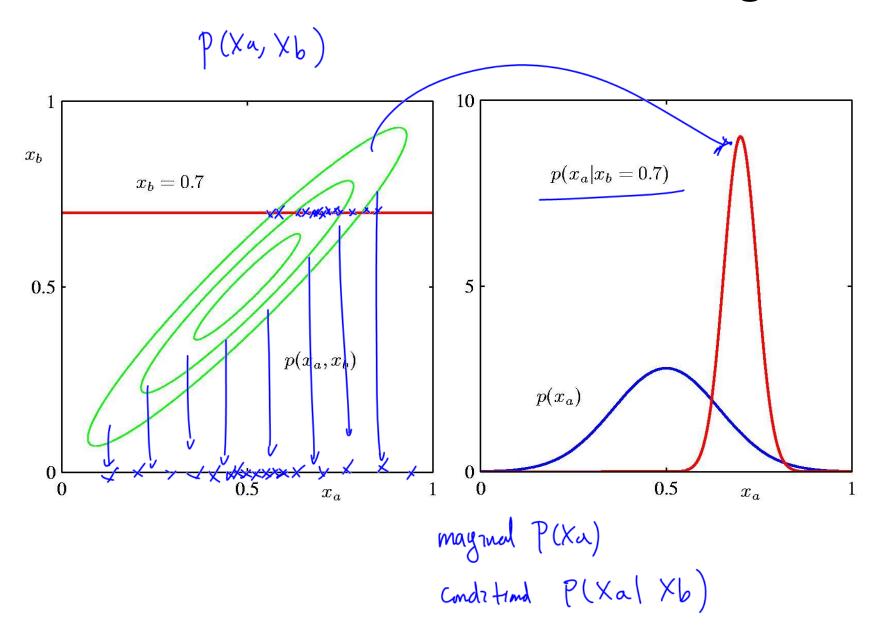
$$egin{array}{lll} \mu_{a|b} &=& \mu_a + \Sigma_{ab}\Sigma_{bb}^{-1}(\mathbf{x}_b - \mu_b) \ \Sigma_{a|b} &=& \Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba}. \end{array} \ - \, ext{where} & \left( egin{array}{lll} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{array} 
ight)^{-1} = \left( egin{array}{lll} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{array} 
ight) \ \Lambda_{aa} &=& \left( \Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba} \right)^{-1} \ \Lambda_{ab} &=& -(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}\Sigma_{ab}\Sigma_{bb}^{-1}. \end{array}$$

Here, we used matrix inversion lemma

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{M} & -\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{M} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \end{pmatrix}$$

- where 
$$M = (A - BD^{-1}C)^{-1}$$
.

#### Partitioned Conditionals and Marginals



#### Linear Gaussian Distributions

- Linear combination of Gaussians is also a Gaussian distribution.
  - Its marginal distribution is also a Gaussian
  - Its conditional distribution is also a Gaussian

# Bayes' Theorem for Gaussian Variables

Given

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$
 $p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$ 
 $p(\mathbf{y}, \boldsymbol{\mu})$ 

we have

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}})$$
$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\mathbf{\Sigma}\{\mathbf{A}^{\mathrm{T}}\mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \mathbf{\Sigma})$$

where

$$\mathbf{\Sigma} = (\mathbf{\Lambda} + \mathbf{A}^{\mathrm{T}} \mathbf{L} \mathbf{A})^{-1}$$

# **Bayesian Linear Regression**

#### Regression

- Given a set of observations:  $x = \{x_1 \dots x_N\}$ 
  - And corresponding target values:  $\mathbf{t} = \{ t_1 \dots t_N \}$

- We want to learn a function y(x)=t to predict future values.
  - We have just learned to find the maximum likelihood weights  $\mathbf{w}_{ML}$ , to predict  $y(x, \mathbf{w}_{ML})$ .

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

$$P(w|D) \ll P(w|D) = P(D|w)P(w)$$
Gaussian Prior on w

• With a likelihood:

$$t_{n} \sim \mathcal{N}(w^{\intercal}\phi(x_{n}), \beta^{\intercal})$$

• and a prior:

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

we get a posterior

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

where

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

and

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

# Simplify the Prior

• Zero-mean isotropic Gaussian: we R

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|0, \underline{\alpha^{-1}\mathbf{I}}) \quad \text{if } \mathbf{w} \in \mathbb{R}^{\mathsf{MAN}}$$

The corresponding posterior is:

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

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

# Derivation of posterior distribution

Use linear Gaussian distributions

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

$$p(t|w,x) = N(t|\Phi w, \beta^{-1}I)$$
• Conditional distribution

$$p(w|t,x) = N(w|\Sigma(\beta\Phi^T t), \Sigma^{-1})$$

where 
$$\Sigma = (\alpha I + \beta \Phi^T \Phi)^{-1}$$

 $= | \phi(x_1)' w |$   $| \phi(x_0)^T \omega |$ 

Recall:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$
 $p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$ 
 $p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}})$ 
 $p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^{\mathrm{T}}\mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma})$ 
 $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{A})^{-1}$ 

# Maximize the (log) Posterior

Same as minimizing sum-of-squared error, wmap
 with regularization term:

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

- This is the same as  $w_{ML}$ , with  $\lambda = \frac{\alpha}{2}$
- So the solution is (same as  $w_{ML}$ )

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

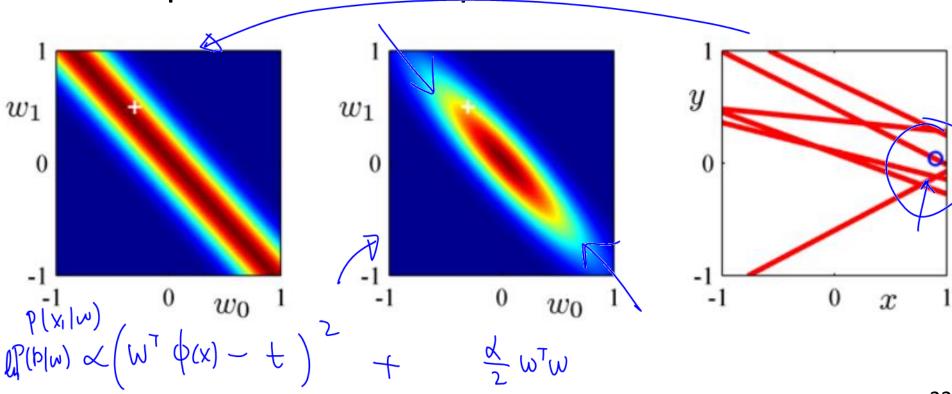
But now we have the variance on w, as well!

• Simple model:  $y(x,w) = w_0 + w_1x$ .

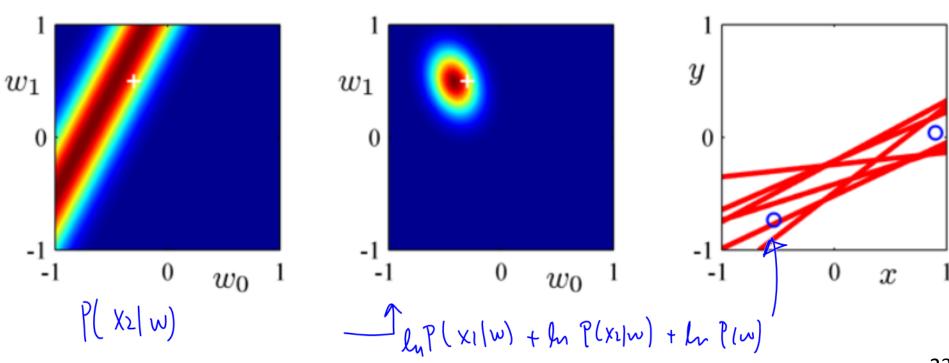
– Prior \* Likelihood = Posterior

P(w) ~ exp(-dw/2) prior/posterior likelihood ~1 data space y $w_1$ 

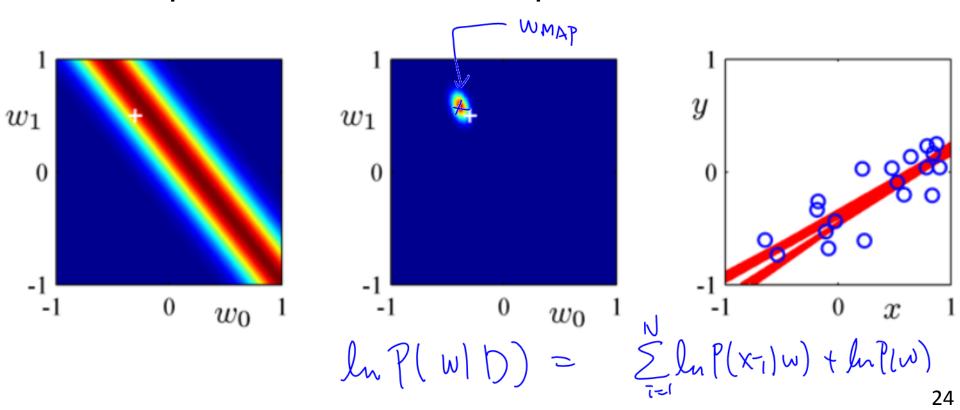
- Likelihood of most recent data point.
  - Posterior = likelihood \* prior
- Samples drawn from posterior.



- Likelihood of second data point (only).
  - Posterior = Likelihood \* Prior
- Sample lines drawn from posterior.



- Likelihood of third data point (only).
  - Posterior = Likelihood \* Prior
- Sample lines drawn from posterior



Our real goal is to predict t given new x, so we evaluate the predictive distribution:

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

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

and we just derived

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$$
  
 $\mathbf{m}_N = \beta \mathbf{S}_N \mathbf{\Phi}^T \mathbf{t}$   $\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \mathbf{\Phi}^T \mathbf{\Phi}$ 

#### **Predictive Distribution**

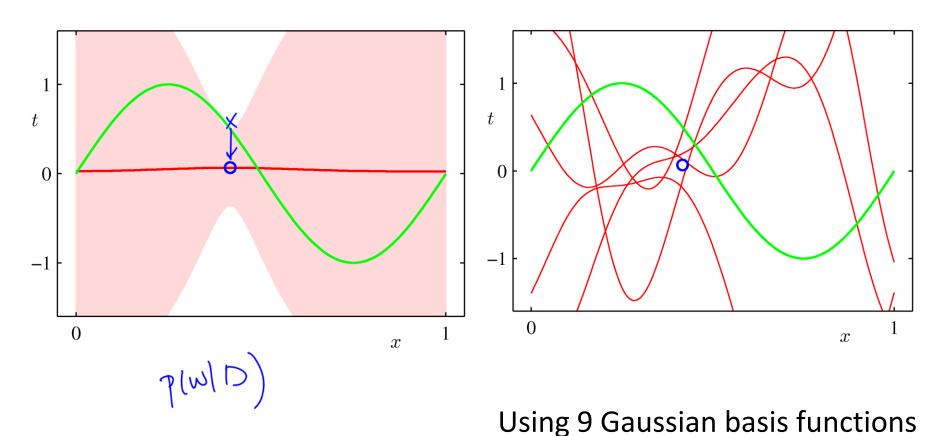
 The integral is a convolution of Gaussians, so the result is:

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

where

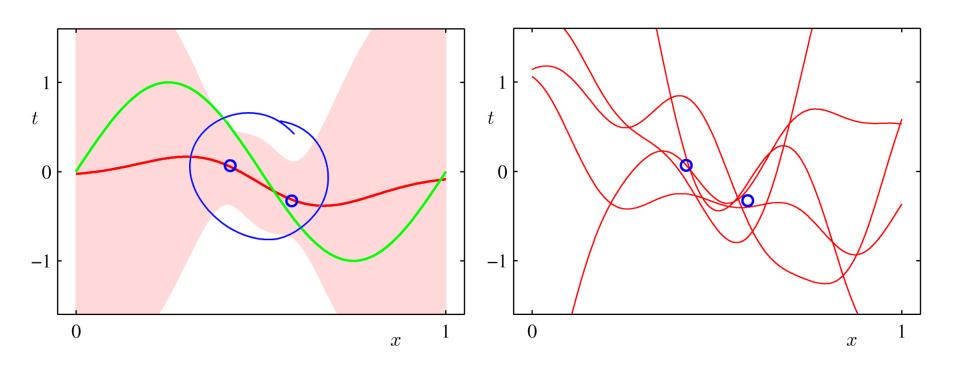
$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_n \phi(\mathbf{x})$$

= noise in the data + uncertainty in w

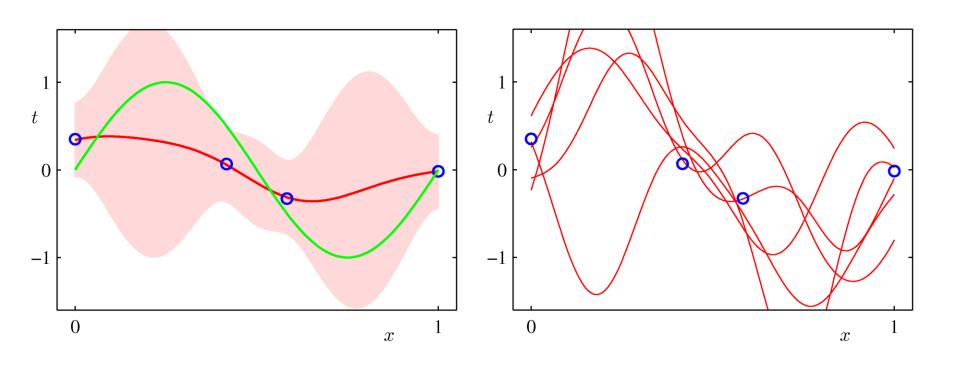


• N=1 observed point

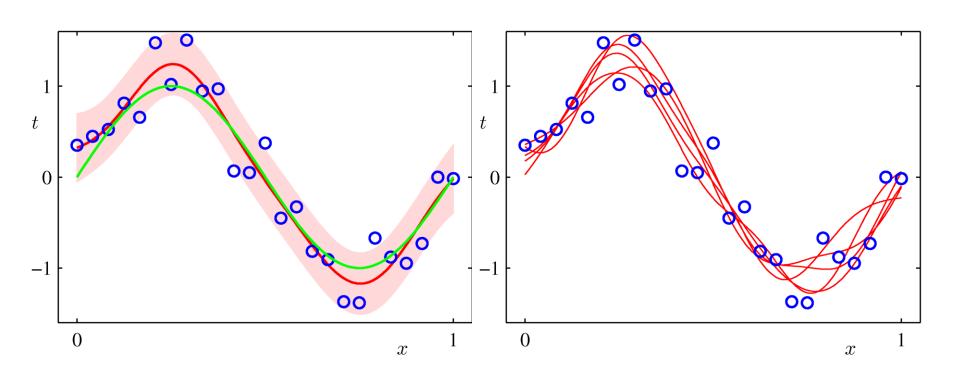
$$\phi_j(x) = \exp\left\{-\frac{(x-\mu_j)^2}{2s^2}\right\}$$



• N=2 observed points



• N=4 observed points



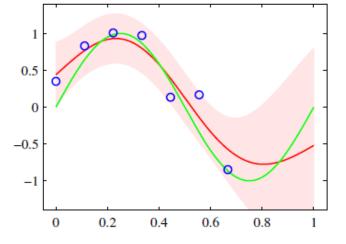
N=25 observed points

#### **Gaussian Processes**

# Why GPs?

Here are some data points. What function did they

come from?



- GPs are a nice way of expressing this "prior on functions" idea.
- Applications:
  - Regression
  - Classification

#### Linear Regression Revisited

• Linear regression model: Combination of M fixed basis functions given by  $\phi(x)$ , so that

$$y(x) = w^T \phi(x)$$

- Prior distribution  $p(w) = N(w | 0, \alpha^{-1}I)$
- Given training data points  $x_1,...,x_n$ , what is the joint distribution of  $y(x_1),...,y(x_n)$ ?
  - y is the vector with elements  $y_n = y(x_n)$ , this vector is given by  $y = \Phi w$
  - where  $\Phi$  is the design matrix with elements  $\Phi_{nk} = \Phi_k(x_n)$

#### Linear Regression Revisited

- $y = \Phi w$ : y is a linear combination of Gaussian distributed variables w, hence itself is Gaussian.
- Mean and covariance

$$E[y] = \Phi E[w] = 0$$

$$E[y] = \Phi E[w] = 0$$

$$= \Phi E[w] = 0$$

$$cov[y] = E[yy^{T}] = \Phi E[ww^{T}] \Phi^{T} = \frac{1}{\alpha} \Phi \Phi^{T} = K$$

where K is the Gram matrix with elements

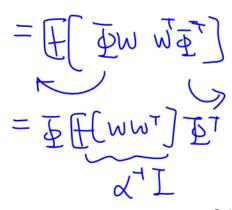
$$(\mathbf{K}_{nm}) = k(x_n, x_m) = \frac{1}{\alpha} \Phi(x_n)^T \Phi(x_m)$$

$$= \mathbf{E} \left[ \mathbf{\Phi} \mathbf{W} \mathbf{W} \mathbf{\Phi}^T \right]$$

$$= \mathbf{E} \left[ \mathbf{W} \mathbf{W}^T \right] \mathbf{E}^T$$

$$= \mathbf{E} \left[ \mathbf{W} \mathbf{W}^T \right] \mathbf{E}^T$$

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

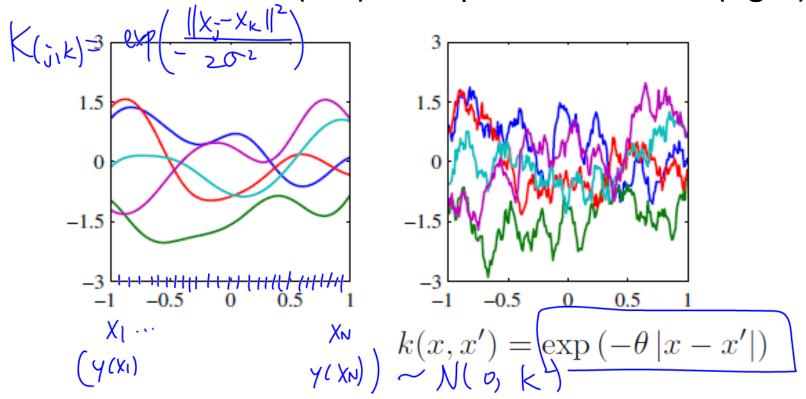


#### Definition of GP

- A Gaussian process is defined as a probability distribution over functions y(x), such that the set of values of y(x) evaluated at an arbitrary set of points x1,.. Xn jointly have a Gaussian  $(x_1, x_2, x_3)$  distribution.  $(y(x_1), y(x_2)) \sim (y(x_1), y(x_2)$ 
  - Gaussian distribution
- What determines the GP is
  - The mean function  $\mu(x) = E(y(x))$
  - The covariance function (kernel) k(x,x')=E(y(x)y(x'))
  - In most applications, we take  $\mu(x)=0$ . Hence the prior is represented by the kernel.

#### Covariance function of GP defines prior

- The figure show samples of functions drawn from Gaussian processes for two different choices of kernel functions
  - Gaussian kernel (left) vs. exponential kernel (right)



#### Linear regression updated by GP

- Specific case of a Gaussian Process
- It is defined by the linear regression model

$$y(x) = w^T \phi(x)$$

with a weight prior

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

the kernel function is given by

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

$$(x_{l...}, x_{l...}) \rightarrow (y(x_{l...}, y(x_{l...})) \sim \mathcal{N}(o_{l...} k_{l...})$$
GP for regression  $\rightarrow (t_{l...}, t_{l...})$ 

 Take into account of the noise on the observed target values, which are given by

$$\mathbf{t}_{\mathbf{n}} = \mathbf{y}_{\mathbf{n}} + \mathbf{\varepsilon}_{\mathbf{n}} \qquad \qquad \mathbf{\varepsilon}_{\mathbf{n}} \sim \mathbf{N}(\mathbf{o}_{\mathbf{n}} \mathbf{b}^{-1})$$

where  $y_n = y(x_n)$ , and  $\varepsilon_n$  is a random noise variable

Here we consider noise processes that have a Gaussian distribution, so that

$$p(t_n | y_n) = N(t_n | y_n, \beta^{-1})$$

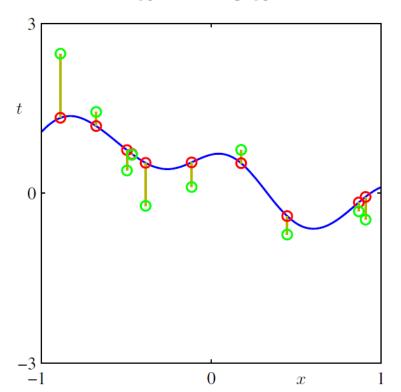
where  $\beta$  is a hyperparameter representing the precision of the noise.

Because the noise is independent, the joint distribution of  $t = (t_1, ..., t_n)^T$  conditioned on  $y = (y_1, ..., y_n)^T$  is given by

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

#### Example: sampling data points

- Sample function from GP (blue):  $y(x) \sim GP(\mu, K)$
- Sample points from GP (red):  $(y(x_1), \dots, y(x_N)) \sim M(\mu, K)$  $(x_1, y(x_1)), (x_2, y(x_2)), (x_N, y(x_N)) \sim GP(\mu, K)$
- Add noise (green):  $t_n(x) \sim y_n(x) + N(0, \beta^{-1})$



#### GP for regression

 From the definition of GP, the marginal distribution p(y) is given by

$$p(y) = N(y | 0, K)$$

The marginal distribution of t is given by

$$p(t) = \int p(t|y)p(y)dy = N(t|0,C) \qquad \text{maginal of}$$

$$\text{the covariance matrix C base elements}$$

Where the covariance matrix C has elements

$$C(x_n, x_m) = k(x_n, x_m) + \beta^{-1} \delta_{nm}$$

$$S_{nm} = \begin{cases} 1 & \text{if } n = m \\ 0 & \text{or } w \end{cases}$$

# **GP** for Regression

- We have used GP to build a model of the joint distribution over sets of data points
- Goal: Given training points  $t_n = (t_1,...,t_n)^T$ , input values  $x_1,...,x_n$ , predict  $t_{n+1}$  for a new input  $x_{n+1}$
- To find  $p(t_{n+1} | t)$ , we begin by writing down the joint distribution

$$p(t_{n+1}) = N(t_{n+1} | 0, C_{n+1})$$
  
where  $C_{n+1}$  is  $(n+1) \times (n+1)$  matrix

joint distribution
$$p(t_{n+1}) = N(t_{n+1} \mid 0, C_{n+1})$$
where  $C_{n+1}$  is  $(n+1) \times (n+1)$  matrix,
$$C_{n+1} = \begin{pmatrix} C_n & k \\ k^T & c \end{pmatrix}, \text{ where } C_n \text{ is } n \times n \text{ matrix, and } c = k(x_{n+1}, x_{n+1}) + \beta^{-1}$$

$$C_n = \begin{pmatrix} C_n & k \\ k^T & c \end{pmatrix}, \text{ where } C_n \text{ is } n \times n \text{ matrix, and } c = k(x_{n+1}, x_{n+1}) + \beta^{-1}$$

$$C_n = \begin{pmatrix} C_n & k \\ k^T & c \end{pmatrix}, \text{ where } C_n \text{ is } n \times n \text{ matrix, and } c = k(x_{n+1}, x_{n+1}) + \beta^{-1}$$

#### **GP** for Regression

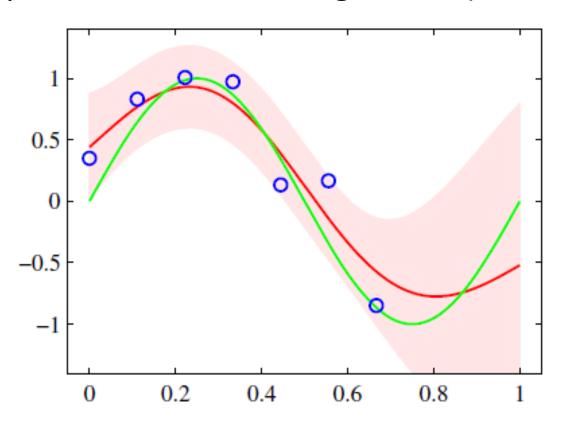
• The conditional distribution  $p(t_{n+1} | t)$  is a Gaussian distribution with mean and covariance given by

$$m(x_{n+1}) = k^T C_n^{-1} t$$
  
 $\sigma^2(x_{n+1}) = c - k^T C_n^{-1} k$ 

- These are the key results that define Gaussian process regression.
- The predictive distribution is a Gaussian whose mean and variance both depend on  $\mathcal{X}_{n+1}$

#### An Example of GP regression

- Green: underlying function (sine function)
- Blue: samples from GP (with noise)
- Red: prediction from GP regression ) with "error bars"



### **GP** for Regression

 The only restriction on the kernel is that the covariance matrix given by

$$C(x_n, x_m) = k(x_n, x_m) + \beta^{-1} \delta_{nm}$$

must be positive definite.

• GP will involve a matrix of size N\*N, for which require  $O(N^3)$  computations.

#### Learning Hyperparameters

Log likelihood

$$\ln p(\mathbf{t}|\boldsymbol{\theta}) = -\frac{1}{2} \ln |\mathbf{C}_N| - \frac{1}{2} \mathbf{t}^{\mathrm{T}} \mathbf{C}_N^{-1} \mathbf{t} - \frac{N}{2} \ln(2\pi).$$

• Gradient Ascent for parameter heta

$$\frac{\partial}{\partial \theta_i} \ln p(\mathbf{t}|\boldsymbol{\theta}) = -\frac{1}{2} \mathrm{Tr} \left( \mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_i} \right) + \frac{1}{2} \mathbf{t}^{\mathrm{T}} \mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_i} \mathbf{C}_N^{-1} \mathbf{t}.$$

- where we used the following:

$$\frac{\partial}{\partial x} (\mathbf{A}^{-1}) = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1}$$
$$\frac{\partial}{\partial x} \ln |\mathbf{A}| = \text{Tr} \left( \mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \right)$$

#### Conclusion

- Distribution over functions
- GP generates a jointly have a Gaussian distribution where the input data can be:
  - scalar
  - real vectors
  - graphs
  - strings
  - **–** ...
- Most interesting structure is in k(x,x'), the 'kernel.'
- GP can be used for regression to predict the target for a new input