Slides modified from:

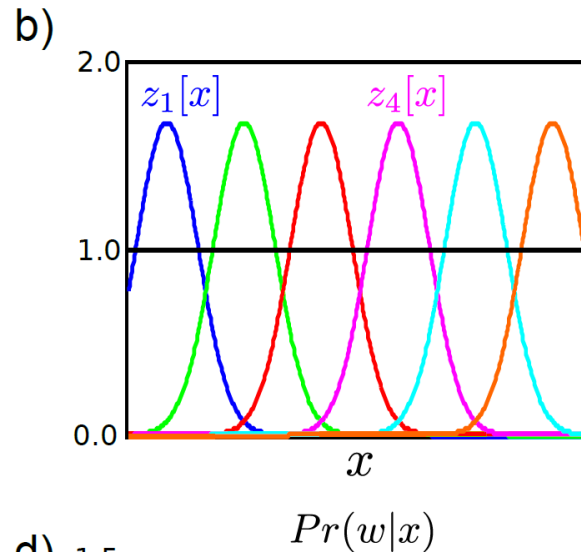PATTERN RECOGNITION

AND MACHINE LEARNING

CHRISTOPHER M. BISHOP
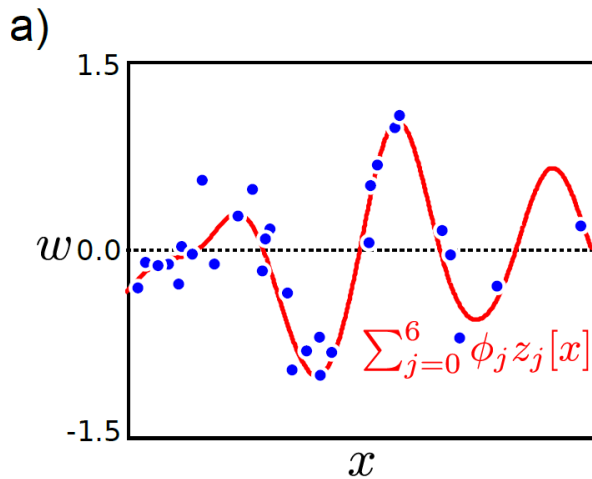
and:
Computer vision: models,
learning and inference.
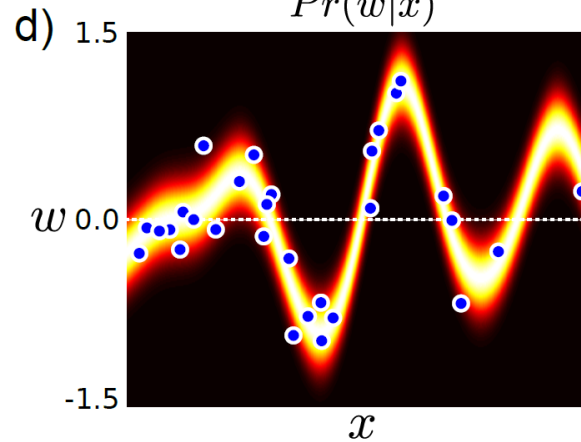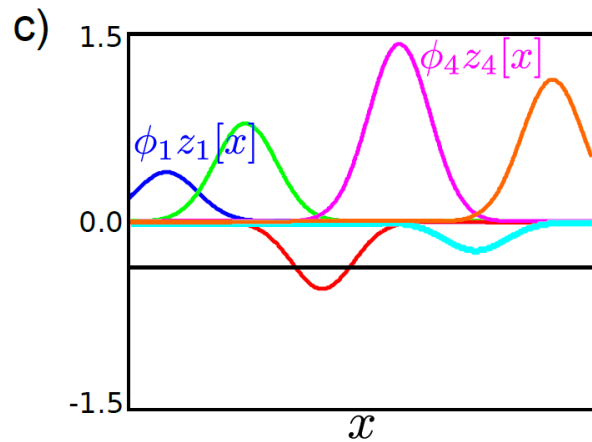©2011 Simon J.D. Prince

# Radial basis functions

a)



$$\sum_{j=0}^{6} \phi_j z_j[x]$$

b)



$z_1[x]$  $z_4[x]$

$$\mathbf{z}_i = \begin{bmatrix} 1 \\ \exp\left[-(x_i - \alpha_1)^2/\lambda\right] \\ \exp\left[-(x_i - \alpha_2)^2/\lambda\right] \\ \exp\left[-(x_i - \alpha_3)^2/\lambda\right] \\ \exp\left[-(x_i - \alpha_4)^2/\lambda\right] \\ \exp\left[-(x_i - \alpha_5)^2/\lambda\right] \\ \exp\left[-(x_i - \alpha_6)^2/\lambda\right] \end{bmatrix}$$

c)



$\phi_1 z_1[x]$  $\phi_4 z_4[x]$

d)

$Pr(w|x)$

# Bayesian regression

# Predictive Distribution (1)

Predict t for new values of x by integrating over w:

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

where

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}).$$

# Bayesian Linear Regression

Likelihood

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

A common choice for the prior is

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

for which the posterior is

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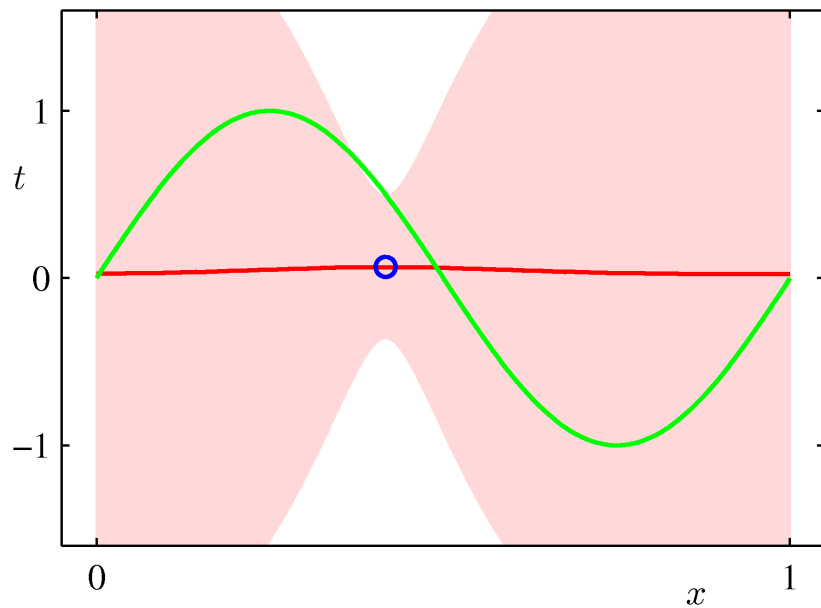$$\mathbf{m}_N = \beta\mathbf{S}_N\boldsymbol{\Phi}^{\mathrm{T}}\mathbf{t}$$
$$\mathbf{S}_N^{-1} = \alpha\mathbf{I} + \beta\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}.$$

# Predictive Distribution (2)

Example: Sinusoidal data, 9 Gaussian basis functions, 1 data point



$$\mathcal{N}(t|\mathbf{m}_N^{\mathrm{T}}\boldsymbol{\phi}(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$
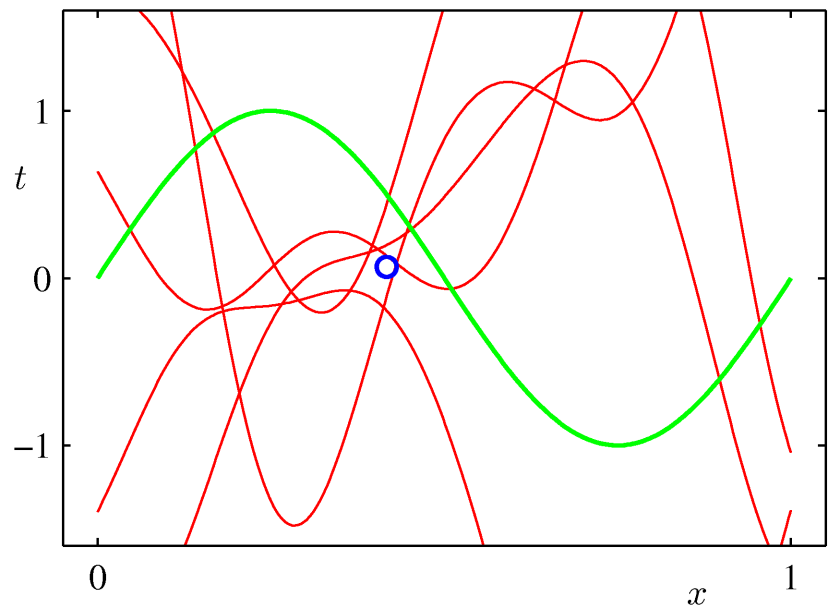
$$y(x, \mathbf{w})$$

# Predictive Distribution (3)

Example: Sinusoidal data, 9 Gaussian basis functions,
2 data points



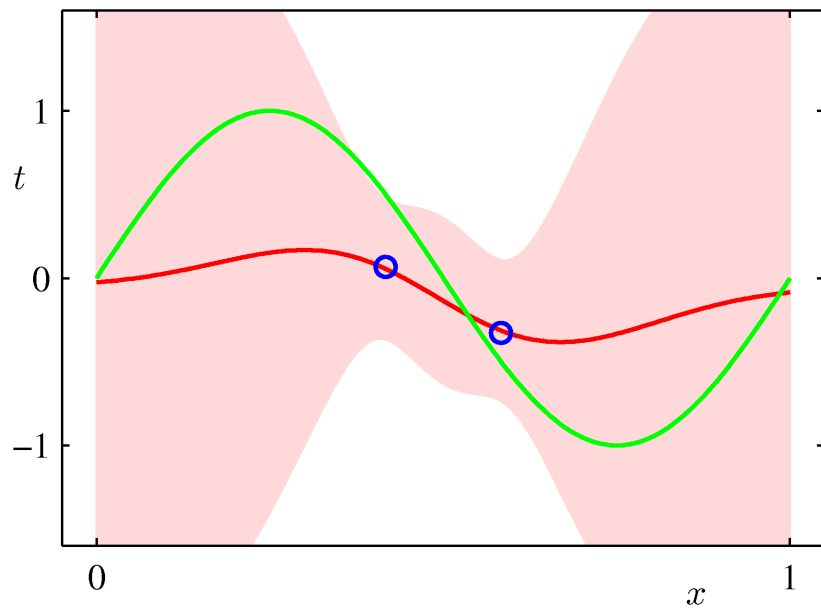$$\mathcal{N}(t|\mathbf{m}_N^{\mathrm{T}}\boldsymbol{\phi}(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$

$$y(x, \mathbf{w})$$
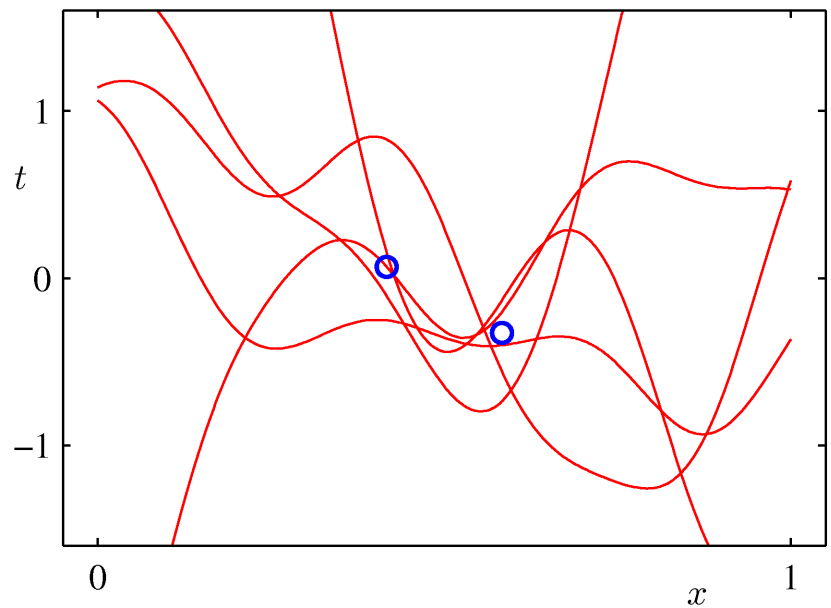
# Predictive Distribution (4)

Example: Sinusoidal data, 9 Gaussian basis functions, 4 data points



$$\mathcal{N}(t|\mathbf{m}_N^{\mathrm{T}}\boldsymbol{\phi}(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$
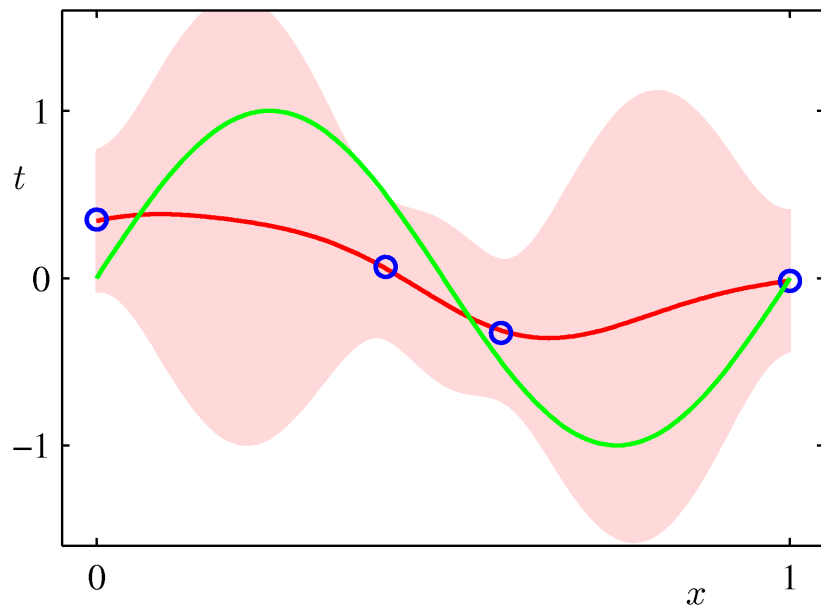
$$y(x, \mathbf{w})$$
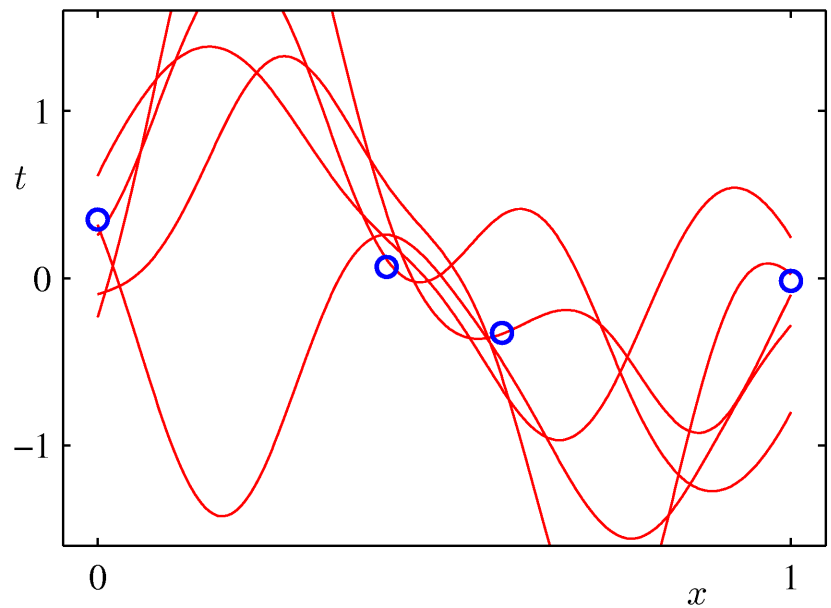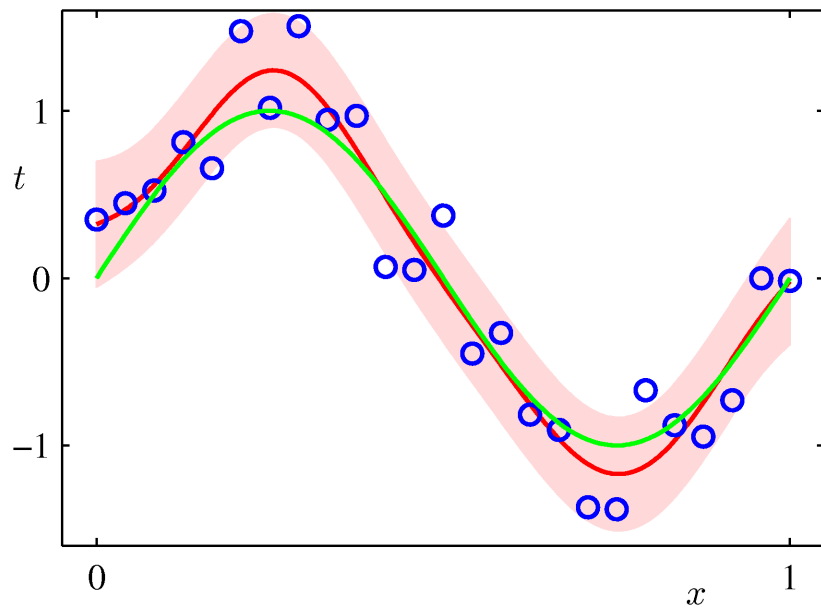
# Predictive Distribution (5)

Example: Sinusoidal data, 9 Gaussian basis functions, 25 data points



$$\mathcal{N}(t|\mathbf{m}_N^{\mathrm{T}}\boldsymbol{\phi}(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$

$$y(x, \mathbf{w})$$

# Predictive Distribution

$$p(t|x, \mathbf{w}_{\mathrm{ML}}, \beta_{\mathrm{ML}}) = \mathcal{N}\left(t|y(x, \mathbf{w}_{\mathrm{ML}}), \beta_{\mathrm{ML}}^{-1}\right)$$

# Bayesian Predictive Distribution

$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}\left(t|m(x), s^2(x)\right)$$

# Equivalent Kernel (1)

The predictive mean can be written

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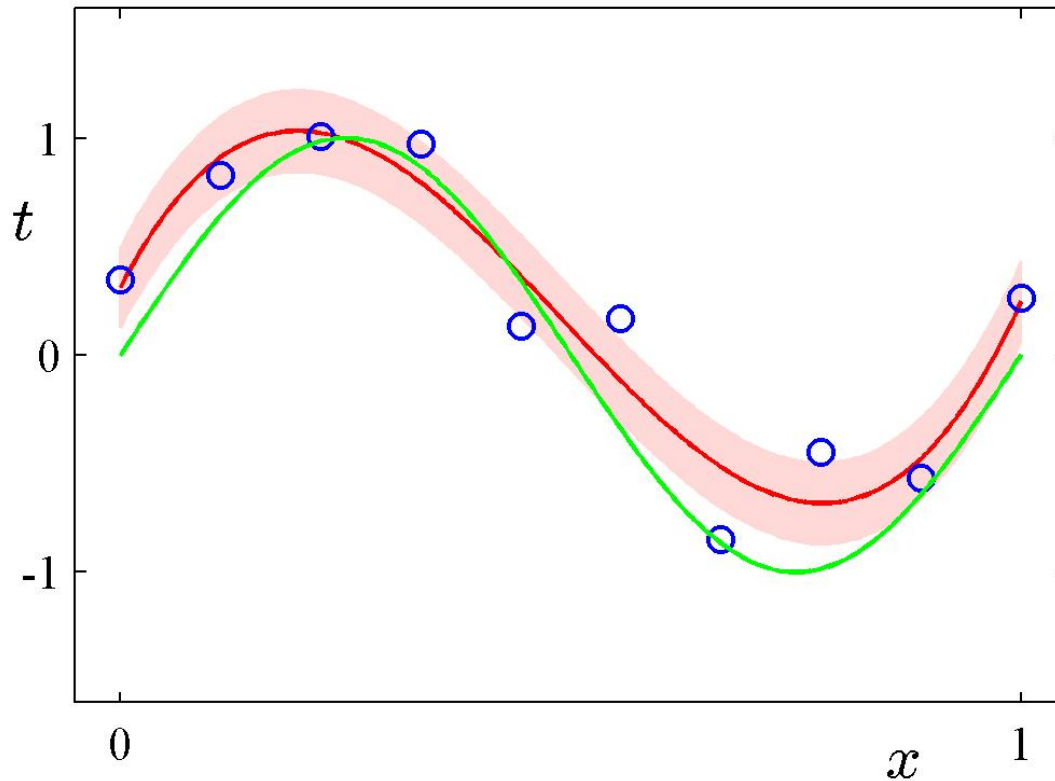y(\mathbf{x}, \mathbf{m}_N) &= \mathbf{m}_N^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}) = \beta \boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} \mathbf{S}_N \boldsymbol{\Phi}^{\mathrm{T}} \mathbf{t} \\
&= \sum_{n=1}^{N} \beta \boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}_n)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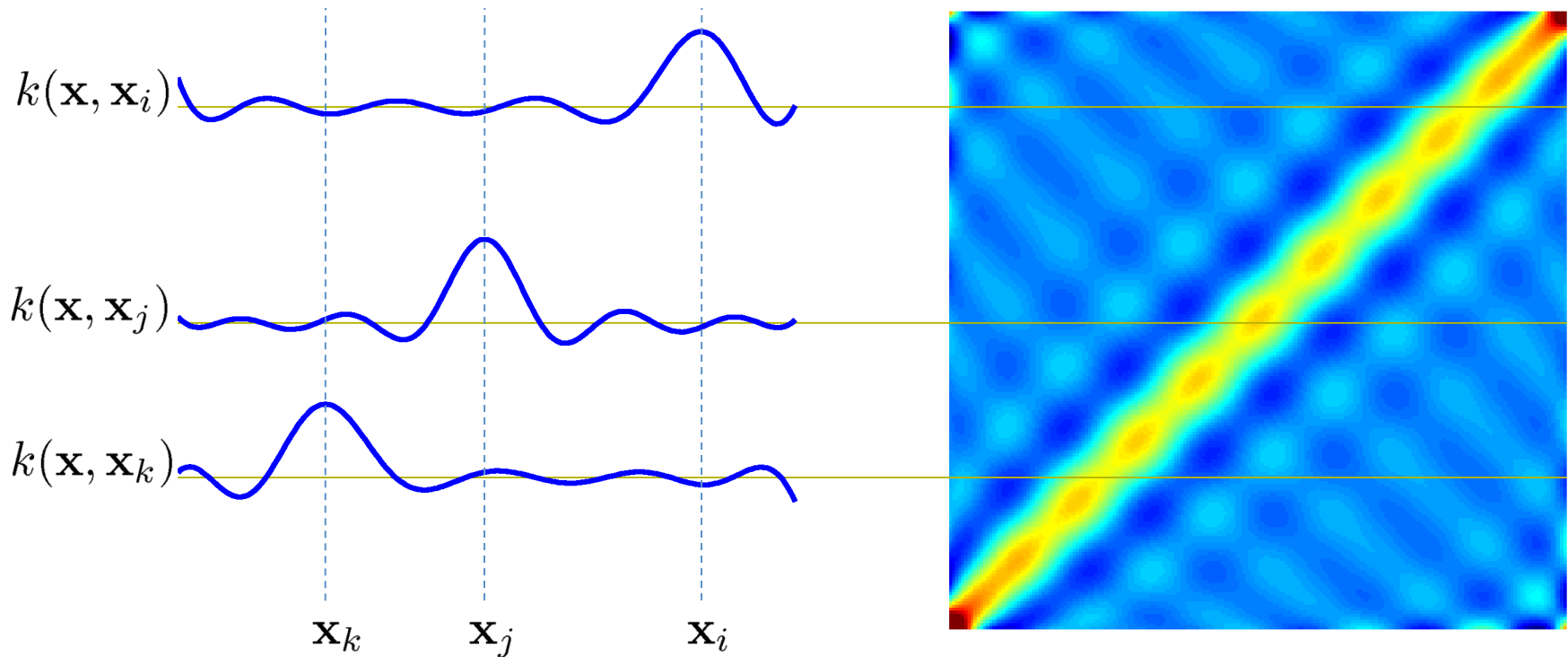&= \sum_{n=1}^{N} k(\mathbf{x}, \mathbf{x}_n) t_n.
\end{aligned}
$$

*Equivalent kernel* or *smoother matrix.*

This is a weighted sum of the training data target values, $t_n$.

# Equivalent Kernel (2)



Weight of $t_n$ depends on distance between x and $x_n$; nearby $x_n$ carry more weight.

# Equivalent Kernel (3)

Non-local basis functions have local equivalent kernels:



Polynomial                    Sigmoidal

# Equivalent Kernel (4)

The kernel as a covariance function: consider

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

We can avoid the use of basis functions and define the kernel function directly, leading to  *Gaussian Processes*.

# Equivalent Kernel (5)

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

for all values of x; however, the equivalent kernel may be negative for some values of x.

Like all kernel functions, the equivalent kernel can be expressed as an inner product:

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

where $\boldsymbol{\psi}(\mathbf{x}) = \beta^{1/2} \mathbf{S}_N^{1/2} \boldsymbol{\phi}(\mathbf{x})$.

# Bayesian Model Comparison (1)

How do we choose the 'right' model?

Assume we want to compare models $M_i$, i=1, …,L, using data $D$; this requires computing

$$p(\mathcal{M}_i|\mathcal{D}) \propto p(\mathcal{M}_i)p(\mathcal{D}|\mathcal{M}_i).$$

Posterior    Prior    *Model evidence* or *marginal likelihood*

*Bayes Factor*: ratio of evidence for two models

$$\frac{p(\mathcal{D}|\mathcal{M}_i)}{p(\mathcal{D}|\mathcal{M}_j)}$$

# Bayesian Model Comparison (2)

Having computed p(M$_i$|D), we can compute the predictive (mixture) distribution

$$p(t|\mathbf{x}, \mathcal{D}) = \sum_{i=1}^{L} p(t|\mathbf{x}, \mathcal{M}_i, \mathcal{D}) p(\mathcal{M}_i|\mathcal{D}).$$

A simpler approximation, known as *model selection*, is to use the model with the highest evidence.

# Bayesian Model Comparison (3)
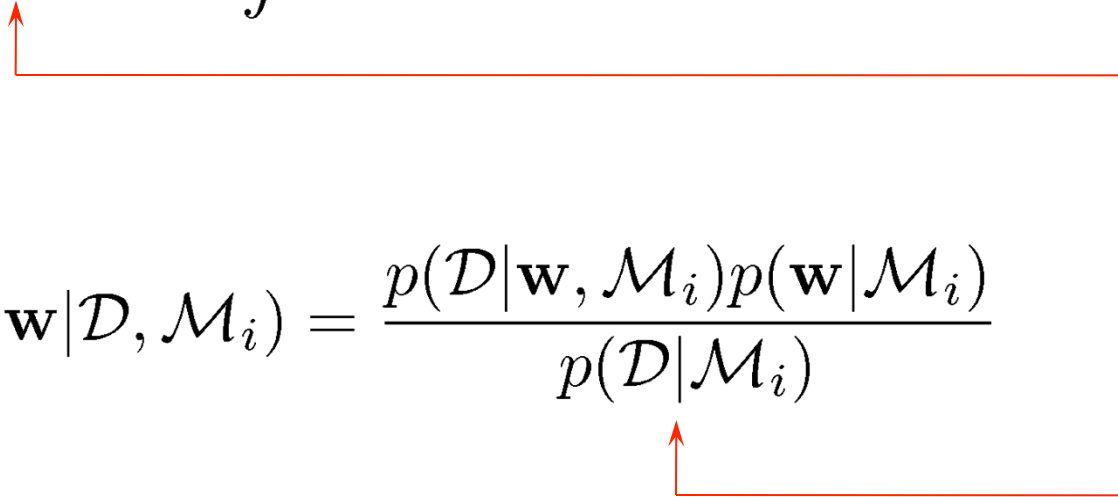
For a model with parameters w, we get the model evidence by marginalizing over w

$$p(\mathcal{D}|\mathcal{M}_i) = \int p(\mathcal{D}|\mathbf{w}, \mathcal{M}_i)p(\mathbf{w}|\mathcal{M}_i)\,\mathrm{d}\mathbf{w}.$$

Note that

$$p(\mathbf{w}|\mathcal{D}, \mathcal{M}_i) = \frac{p(\mathcal{D}|\mathbf{w}, \mathcal{M}_i)p(\mathbf{w}|\mathcal{M}_i)}{p(\mathcal{D}|\mathcal{M}_i)}$$
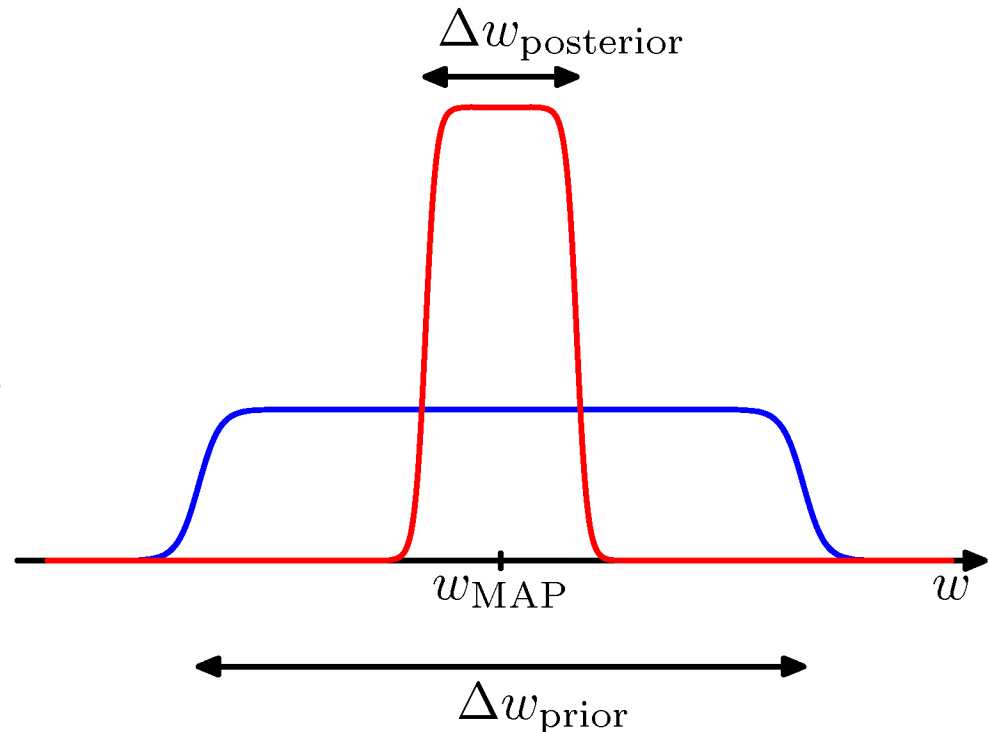
# Bayesian Model Comparison (4)

For a given model with a single parameter, w, consider the approximation

$$p(\mathcal{D}) = \int p(\mathcal{D}|w)p(w)\,\mathrm{d}w$$

$$\simeq \quad p(\mathcal{D}|w_{\mathrm{MAP}})\frac{\Delta w_{\mathrm{posterior}}}{\Delta w_{\mathrm{prior}}}$$

where the posterior is assumed to be sharply peaked.

# Bayesian Model Comparison (5)

Taking logarithms, we obtain

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|w_{\mathrm{MAP}}) + \underbrace{\ln\left(\frac{\Delta w_{\mathrm{posterior}}}{\Delta w_{\mathrm{prior}}}\right)}_{\text{Negative}}.$$

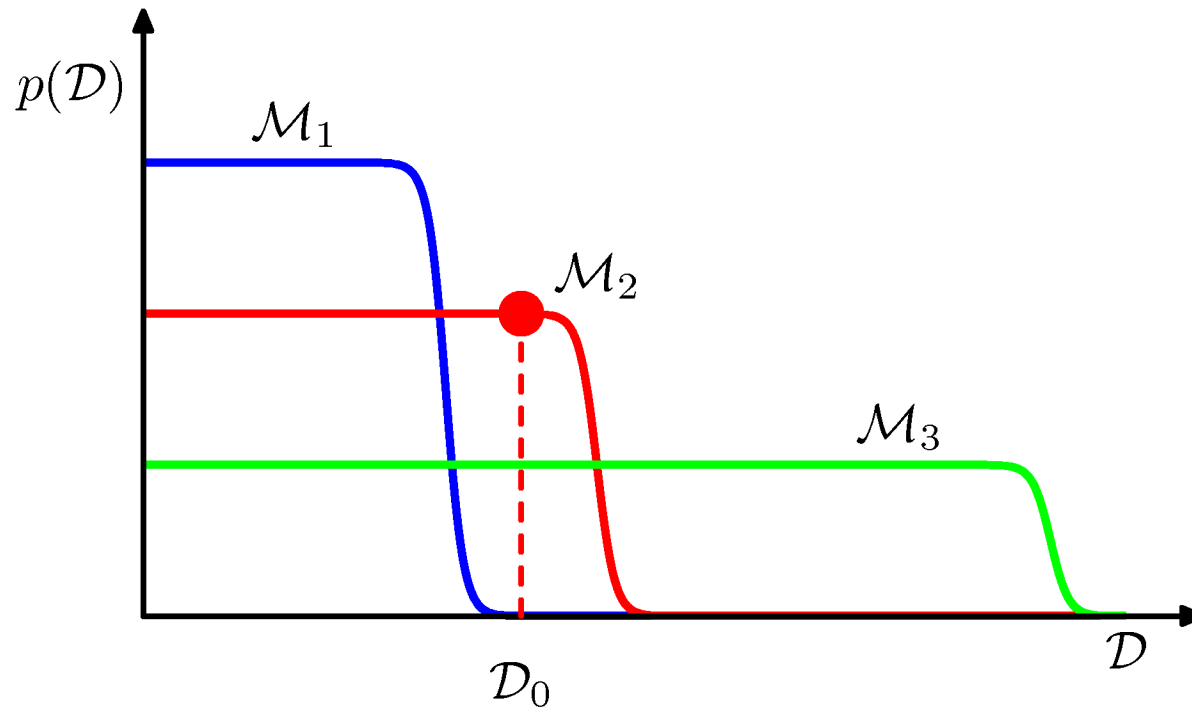With M parameters, all assumed to have the same ratio $\Delta w_{\mathrm{posterior}}/\Delta w_{\mathrm{prior}}$, we get

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\mathbf{w}_{\mathrm{MAP}}) + \underbrace{M \ln\left(\frac{\Delta w_{\mathrm{posterior}}}{\Delta w_{\mathrm{prior}}}\right)}_{\text{Negative and linear in M.}}.$$

# Bayesian Model Comparison (6)

Matching data and model complexity

# The Evidence Approximation (1)

The fully Bayesian predictive distribution is given by

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

but this integral is intractable. Approximate with

$$p(t|\mathbf{t}) \simeq p\left(t|\mathbf{t}, \widehat{\alpha}, \widehat{\beta}\right) = \int p\left(t|\mathbf{w}, \widehat{\beta}\right) p\left(\mathbf{w}|\mathbf{t}, \widehat{\alpha}, \widehat{\beta}\right) \, \mathrm{d}\mathbf{w}$$

where $\left(\widehat{\alpha}, \widehat{\beta}\right)$ is the mode of $p(\alpha, \beta|\mathbf{t})$, which is assumed to be sharply peaked; a.k.a. *empirical Bayes, type II* or *gene-ralized maximum likelihood,* or *evidence approximation*.

# The Evidence Approximation (2)

From Bayes' theorem we have

$$p(\alpha, \beta | \mathbf{t}) \propto p(\mathbf{t} | \alpha, \beta) p(\alpha, \beta)$$

and if we assume p($\alpha$,$\beta$) to be flat we see that

$$
\begin{aligned}
p(\alpha, \beta | \mathbf{t}) &\propto p(\mathbf{t} | \alpha, \beta) \\
&= \int p(\mathbf{t} | \mathbf{w}, \beta) p(\mathbf{w} | \alpha) \, \mathrm{d}\mathbf{w}.
\end{aligned}
$$

General results for Gaussian integrals give

$$\ln p(\mathbf{t} | \alpha, \beta) = \frac{M}{2} \ln \alpha + \frac{N}{2} \ln \beta - E(\mathbf{m}_N) + \frac{1}{2} \ln |\mathbf{S}_N| - \frac{N}{2} \ln(2\pi).$$

# The Evidence Approximation (3)

Example: sinusoidal data, M $^{\text{th}}$ degree polynomial, $\alpha = 5 \times 10^{-3}$