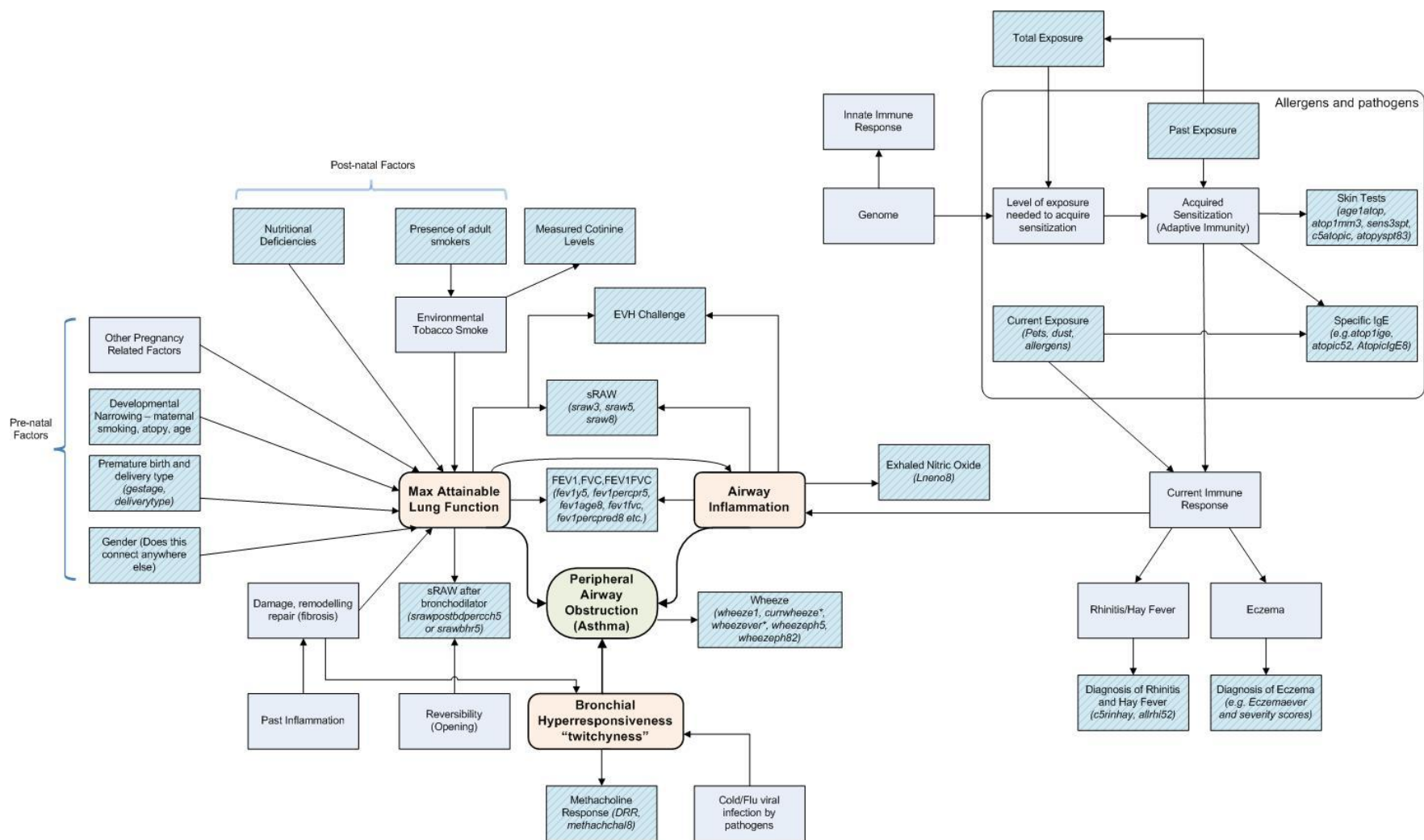


INTRODUCTION TO BAYESIAN INFERENCE – PART 2

CHRIS BISHOP



Personal Healthcare Revolution

Electronic health records (CFH)

Personal genomics

(DeCode, Navigenics, 23andMe)

X-prize: first \$10k human genome technology

NIH: \$1k by 2014

Microsoft Research Cambridge:

PhD Scholarships

Internships: 3 months

Postdoctoral Fellowships

Why Probabilities?

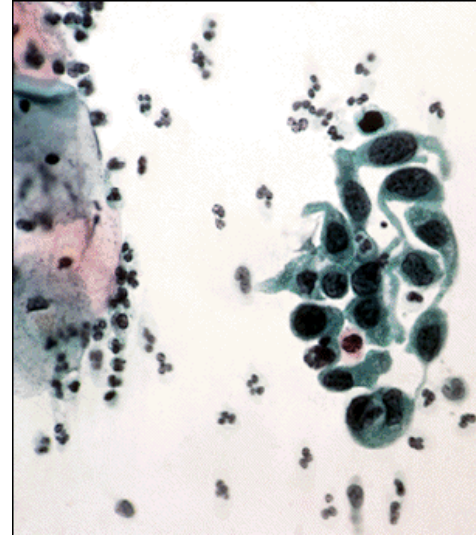
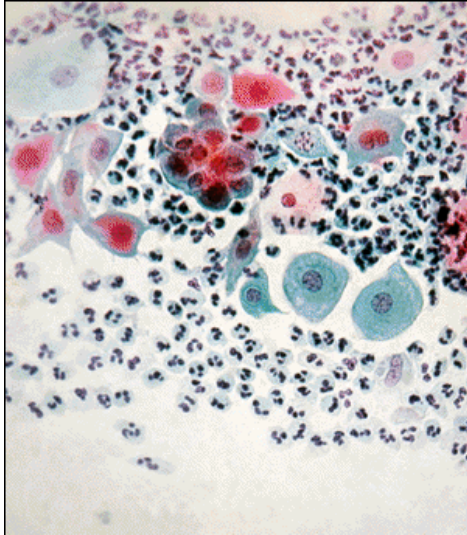


Image vector x

Class C_k “cancer” or “normal”

Decisions

One-step solution

train a function to decide the class

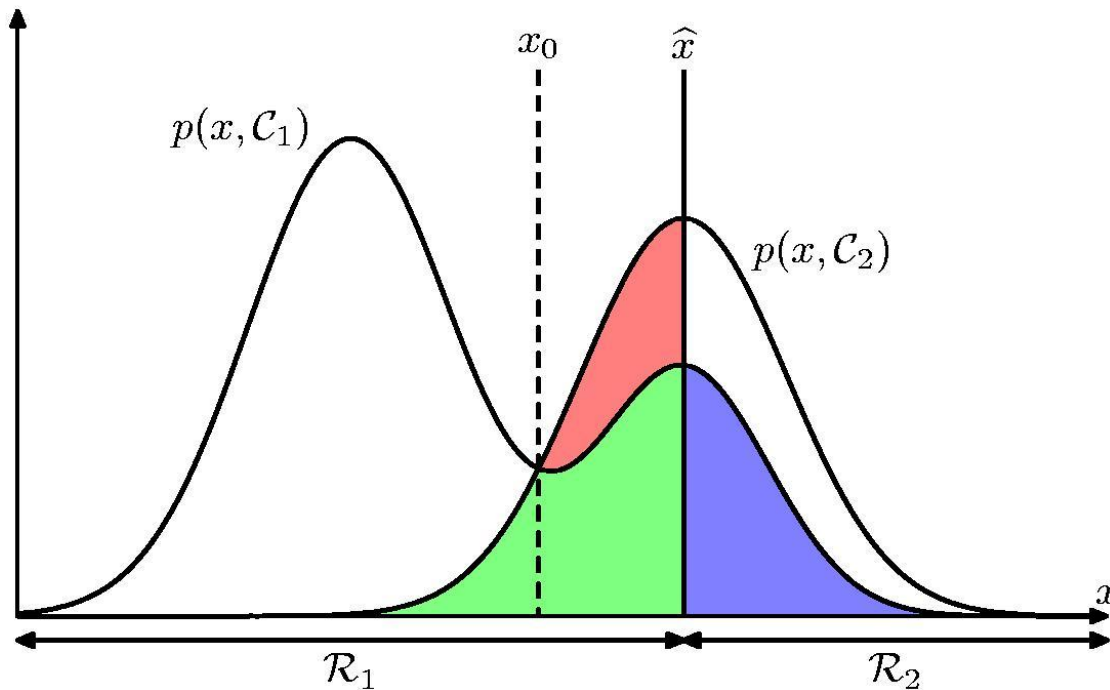
Two-step solution

inference : infer posterior probabilities

$$p(\mathcal{C}_k | \mathbf{x})$$

decision : use probabilities to decide the class

Minimum Misclassification Rate



$$p(\mathbf{x}, C_1) > p(\mathbf{x}, C_2)$$

$$p(C_1|\mathbf{x}) > p(C_2|\mathbf{x})$$

$$\begin{aligned} p(\text{mistake}) &= p(\mathbf{x} \in \mathcal{R}_1, C_2) + p(\mathbf{x} \in \mathcal{R}_2, C_1) \\ &= \int_{\mathcal{R}_1} p(\mathbf{x}, C_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, C_1) d\mathbf{x}. \end{aligned}$$

Why Separate Inference and Decision?

- Minimizing risk (loss matrix may change over time)
- Reject option
- Unbalanced class priors
- Combining models

Loss Matrix

		Decision	
		cancer	normal
True class	cancer	0	1000
	normal	1	0

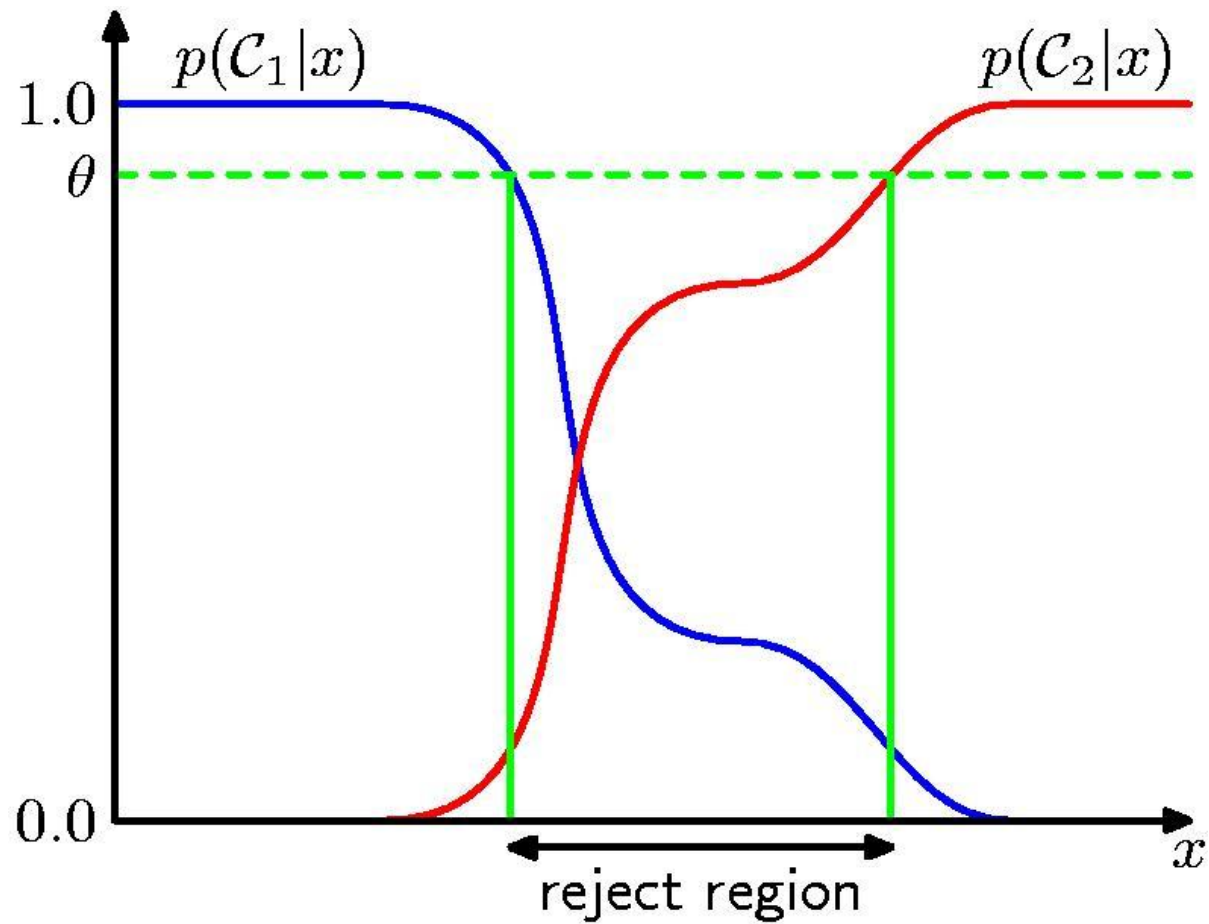
Minimum Expected Loss

$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x}$$

Regions \mathcal{R}_j are chosen, at each \mathbf{x} , to minimize

$$\sum_k L_{kj} p(\mathcal{C}_k | \mathbf{x})$$

Reject Option



Unbalanced class priors

In screening application, cancer is very rare

Use “balanced” data sets to train models, then use Bayes’ theorem to correct the posterior probabilities

Combining models

Image data and blood tests

Assume independent for each class:

$$p(\mathbf{x}_I, \mathbf{x}_B | \mathcal{C}_k) \propto p(\mathbf{x}_I | \mathcal{C}_k) p(\mathbf{x}_B | \mathcal{C}_k)$$

$p(\mathcal{C}_k | \mathbf{x}_I, \mathbf{x}_B) \propto p(\mathbf{x}_I, \mathbf{x}_B | \mathcal{C}_k) p(\mathcal{C}_k)$

$\propto p(\mathbf{x}_I | \mathcal{C}_k) p(\mathbf{x}_B | \mathcal{C}_k) p(\mathcal{C}_k)$

$\propto \frac{p(\mathcal{C}_k | \mathbf{x}_I) p(\mathcal{C}_k | \mathbf{x}_B)}{p(\mathcal{C}_k)} \mathbf{x}_B$

Binary Variables (1)

Coin flipping: heads=1, tails=0

$$p(x = 1|\mu) = \mu \quad \mu \in [0, 1]$$

$$p(x = 0|\mu) = 1 - \mu$$

Bernoulli Distribution

$$\text{Bern}(x|\mu) = \mu^x (1 - \mu)^{1-x}$$

Expectation and Variance

In general

$$\mathbb{E}[f] = \sum_x p(x) f(x) \qquad \mathbb{E}[f] = \int p(x) f(x) \, dx$$

$$\text{var}[f] = \mathbb{E} \left[(f(x) - \mathbb{E}[f(x)])^2 \right] = \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2$$

For Bernoulli

$$\text{Bern}(x|\mu) = \mu^x (1 - \mu)^{1-x}$$

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

$$\text{var}[x] = \mu(1 - \mu)$$

Likelihood function

Data set

$$\mathcal{D} = \{x_1, \dots, x_N\}, \text{ } m \text{ heads } (x = 1), \text{ } N - m \text{ tails } (x = 0)$$

Likelihood function

$$\begin{aligned} p(\mathcal{D}|\mu) &= \prod_{n=1}^N p(x_n|\mu) \\ &= \prod_{n=1}^N \mu^{x_n} (1 - \mu)^{1-x_n} \\ &= \mu^m (1 - \mu)^{N-m} \end{aligned}$$

Prior Distribution

Simplification if prior has same functional form as likelihood function

$$p(\mu) \propto \mu^{a-1}(1 - \mu)^{b-1}$$

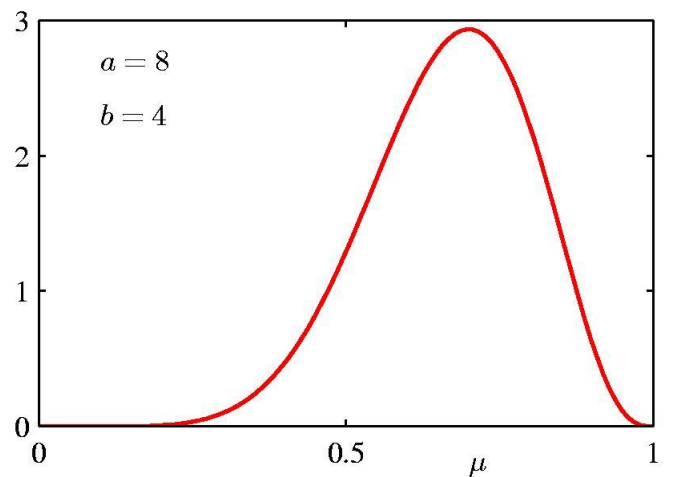
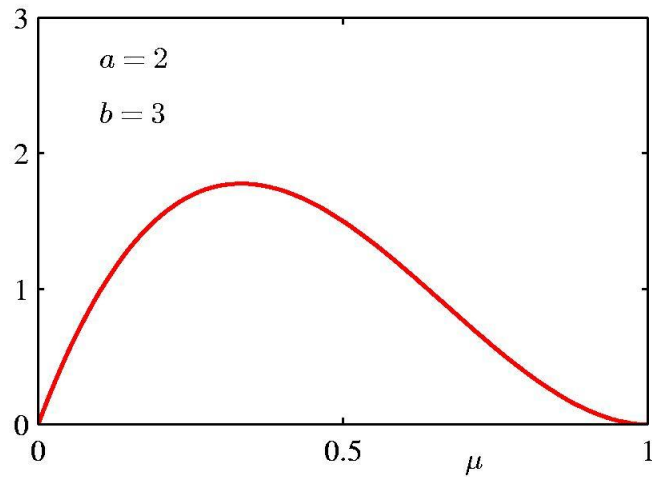
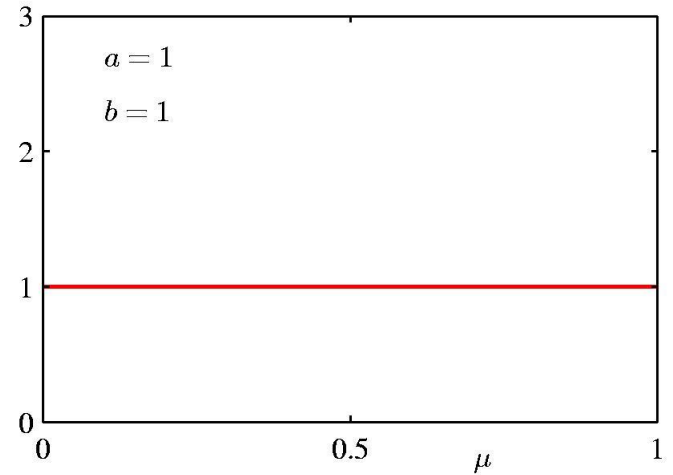
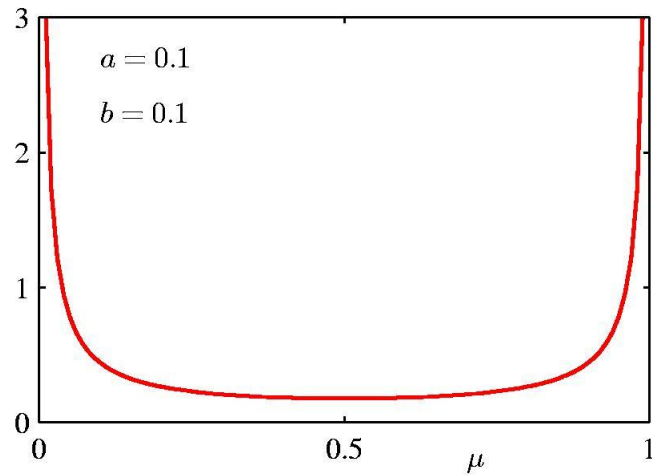
Called *conjugate prior*

$$\text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1}(1 - \mu)^{b-1}$$

$$\mathbb{E}[\mu] = \frac{a}{a+b}$$

$$\text{var}[\mu] = \frac{ab}{(a+b)^2(a+b+1)}$$

Beta Distribution



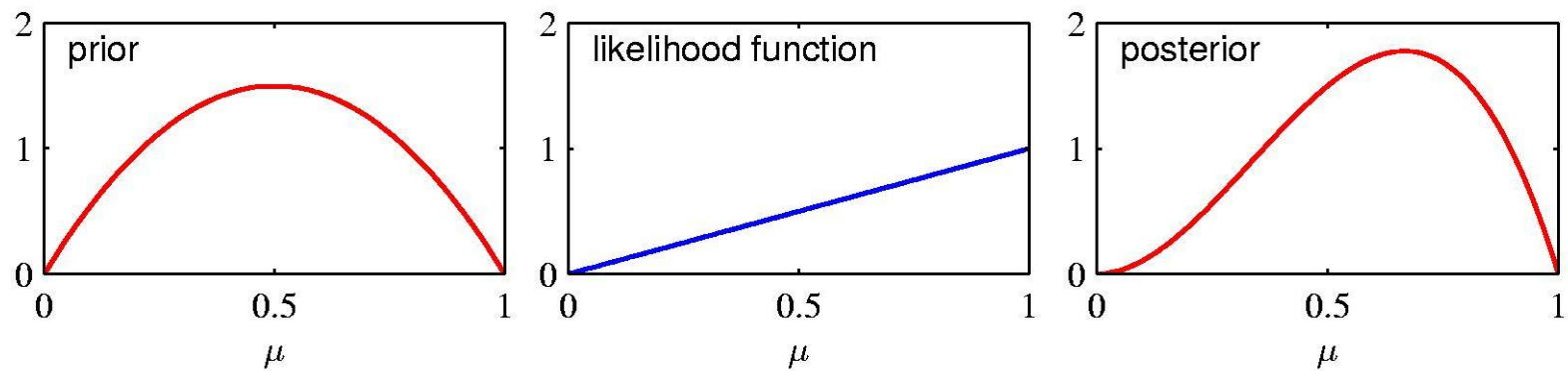
Posterior Distribution

$$\begin{aligned} p(\mu|a_0, b_0, \mathcal{D}) &\propto p(\mathcal{D}|\mu)p(\mu|a_0, b_0) \\ &\propto \left(\prod_{n=1}^N \mu^{x_n} (1 - \mu)^{1-x_n} \right) \text{Beta}(\mu|a_0, b_0) \\ &\propto \mu^{m+a_0-1} (1 - \mu)^{(N-m)+b_0-1} \end{aligned}$$

$$p(\mu|a_0, b_0, \mathcal{D}) = \text{Beta}(\mu|a_N, b_N)$$

$$a_N = a_0 + m \quad b_N = b_0 + (N - m)$$

Posterior Distribution



Properties of the Posterior

As the size N of the data set increases

$$a_N \rightarrow m$$

$$b_N \rightarrow N - m$$

$$\mathbb{E}[\mu] = \frac{a_N}{a_N + b_N} \rightarrow \frac{m}{N}$$

$$\text{var}[\mu] = \frac{a_N b_N}{(a_N + b_N)^2 (a_N + b_N + 1)} \rightarrow 0$$

Predictive Distribution

What is the probability that the next coin flip will be heads?

$$\begin{aligned} p(x = 1 | a_0, b_0, \mathcal{D}) &= \int_0^1 p(x = 1 | \mu) p(\mu | a_0, b_0, \mathcal{D}) \, d\mu \\ &= \int_0^1 \mu p(\mu | a_0, b_0, \mathcal{D}) \, d\mu \\ &= \mathbb{E}[\mu | a_0, b_0, \mathcal{D}] \\ &= \frac{a_N}{a_N + b_N} \end{aligned}$$

The Exponential Family

$$p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x})g(\boldsymbol{\eta}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \}$$

where $\boldsymbol{\eta}$ is the *natural parameter*

$$g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \} \, d\mathbf{x} = 1$$

We can interpret $g(\boldsymbol{\eta})$ as the normalization coefficient

Likelihood Function

Give a data set, $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$

$$p(\mathbf{X}|\boldsymbol{\eta}) = \left(\prod_{n=1}^N h(\mathbf{x}_n) \right) g(\boldsymbol{\eta})^N \exp \left\{ \boldsymbol{\eta}^T \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n) \right\}$$

Depends on data through *sufficient statistics*

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

Expected Sufficient Statistics

$$g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \} \, d\mathbf{x} = 1$$

$$\underbrace{\nabla g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \} \, d\mathbf{x}}_{1/g(\boldsymbol{\eta})} + \underbrace{g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \} \mathbf{u}(\mathbf{x}) \, d\mathbf{x}}_{\mathbb{E}[\mathbf{u}(\mathbf{x})]} = 0$$

$$-\nabla \ln g(\boldsymbol{\eta}) = \mathbb{E}[\mathbf{u}(\mathbf{x})]$$

Conjugate priors

For the exponential family

$$p(\boldsymbol{\eta}|\boldsymbol{\chi}, \nu) = f(\boldsymbol{\chi}, \nu)g(\boldsymbol{\eta})^\nu \exp \{ \nu \boldsymbol{\eta}^\text{T} \boldsymbol{\chi} \}$$

Combining with the likelihood function, we get

$$p(\boldsymbol{\eta}|\mathbf{X}, \boldsymbol{\chi}, \nu) \propto g(\boldsymbol{\eta})^{\nu+N} \exp \left\{ \boldsymbol{\eta}^\text{T} \left(\sum_{n=1}^N \mathbf{u}(\mathbf{x}_n) + \nu \boldsymbol{\chi} \right) \right\}$$

Prior corresponds to ν pseudo-observations with statistic $\boldsymbol{\chi}$

Bernoulli revisited

The Bernoulli distribution

$$\begin{aligned} p(x|\mu) &= \text{Bern}(x|\mu) = \mu^x (1 - \mu)^{1-x} \\ &= \exp \{ x \ln \mu + (1 - x) \ln(1 - \mu) \} \\ &= (1 - \mu) \exp \left\{ \ln \left(\frac{\mu}{1 - \mu} \right) x \right\} \end{aligned}$$

Comparing with the general form we see that

$$\eta = \ln \left(\frac{\mu}{1 - \mu} \right) \quad \text{and so} \quad \mu = \underbrace{\sigma(\eta)}_{\text{Logistic sigmoid}} = \frac{1}{1 + \exp(-\eta)}$$

Bernoulli revisited

The Bernoulli distribution in canonical form

$$p(x|\eta) = h(x)g(\eta) \exp \{ \eta^T u(x) \}$$

where

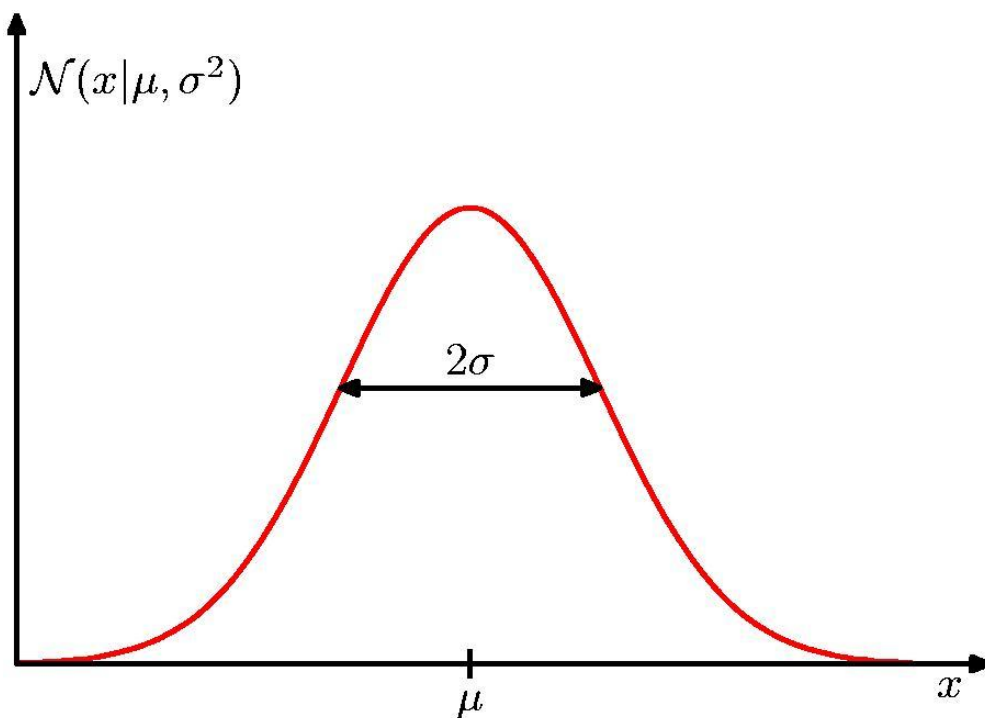
$$u(x) = x$$

$$h(x) = 1$$

$$g(\eta) = 1 - \sigma(\eta) = \sigma(-\eta)$$

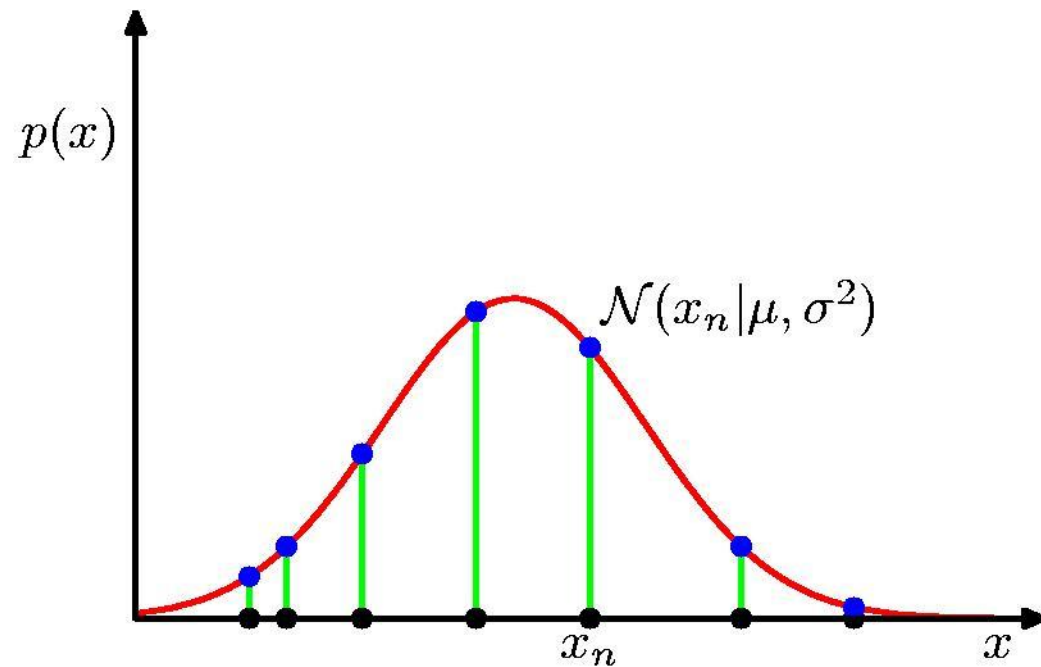
The Gaussian Distribution

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\}$$



Likelihood Function

$$p(\mathbf{x}|\mu, \sigma^2) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \sigma^2)$$



Bayesian Inference – unknown mean

Assume σ^2 is known

Data set

$$\mathbf{x} = \{x_1, \dots, x_N\}$$

Likelihood function for μ

$$p(\mathbf{x}|\mu) = \prod_{n=1}^N p(x_n|\mu) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 \right\}.$$

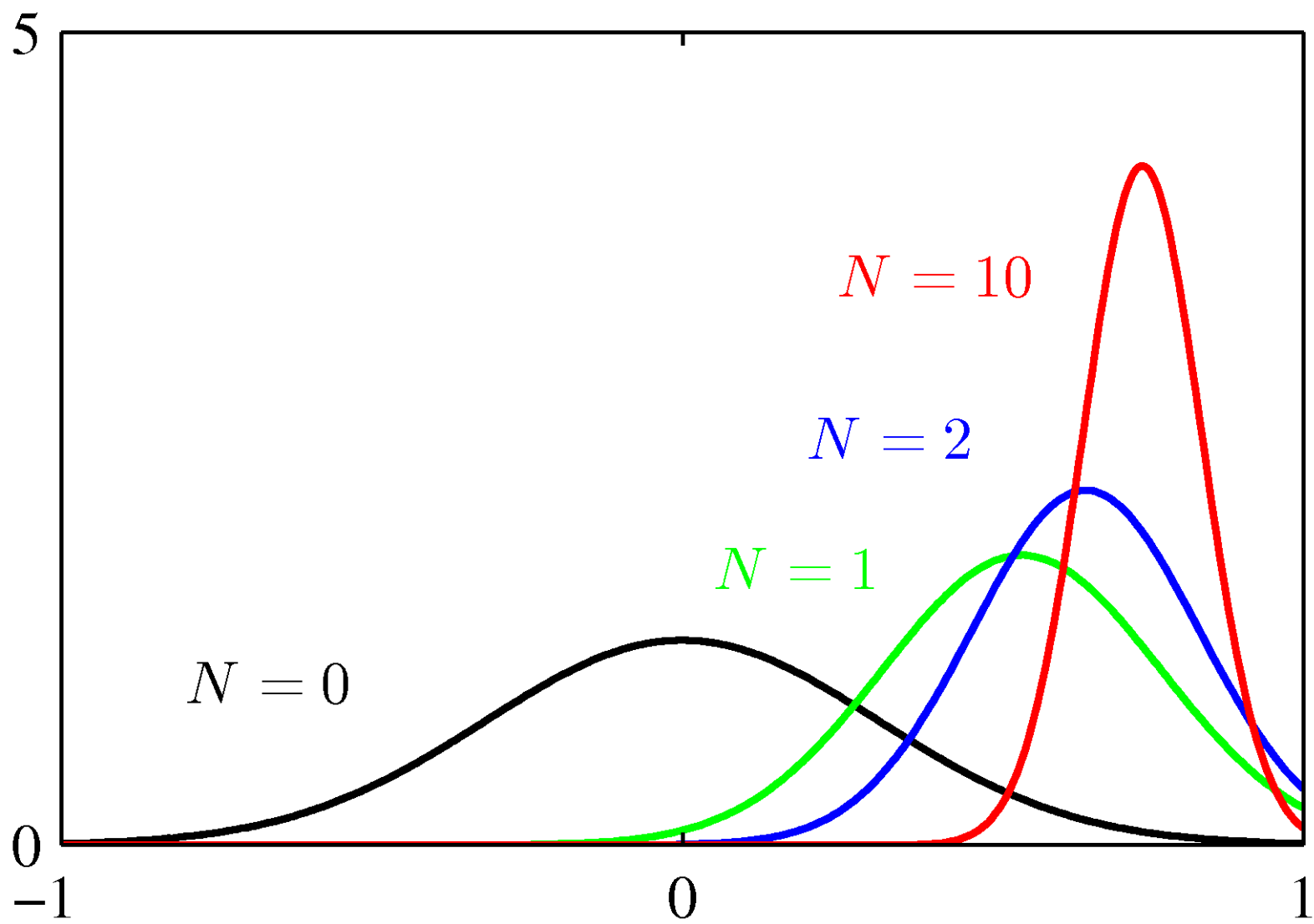
Bayesian Inference – unknown mean

Conjugate prior is a Gaussian

$$p(\mu) = \mathcal{N}(\mu | \mu_0, \sigma_0^2)$$

which gives a Gaussian posterior

$$p(\mu | \mathbf{x}) \propto p(\mathbf{x} | \mu) p(\mu)$$



Bayesian Inference – unknown precision

Now assume μ is known

Likelihood function for precision $\lambda = 1/\sigma^2$

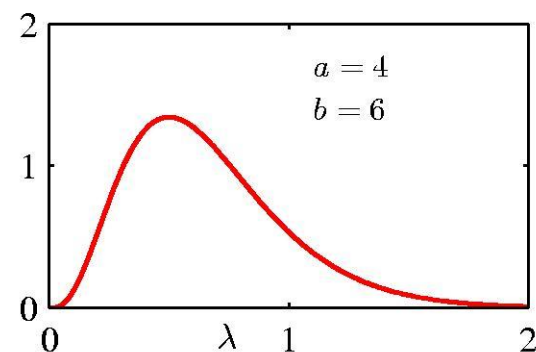
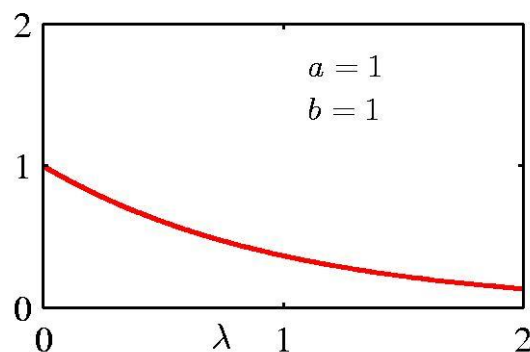
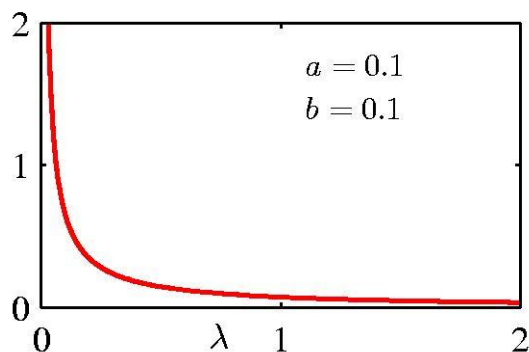
$$p(\mathbf{x}|\lambda) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \lambda^{-1}) \propto \lambda^{N/2} \exp \left\{ -\frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2 \right\}$$

Conjugate prior

Gamma distribution

$$\text{Gam}(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda)$$

$$\mathbb{E}[\lambda] = \frac{a}{b} \qquad \text{var}[\lambda] = \frac{a}{b^2}$$



Unknown Mean and Precision

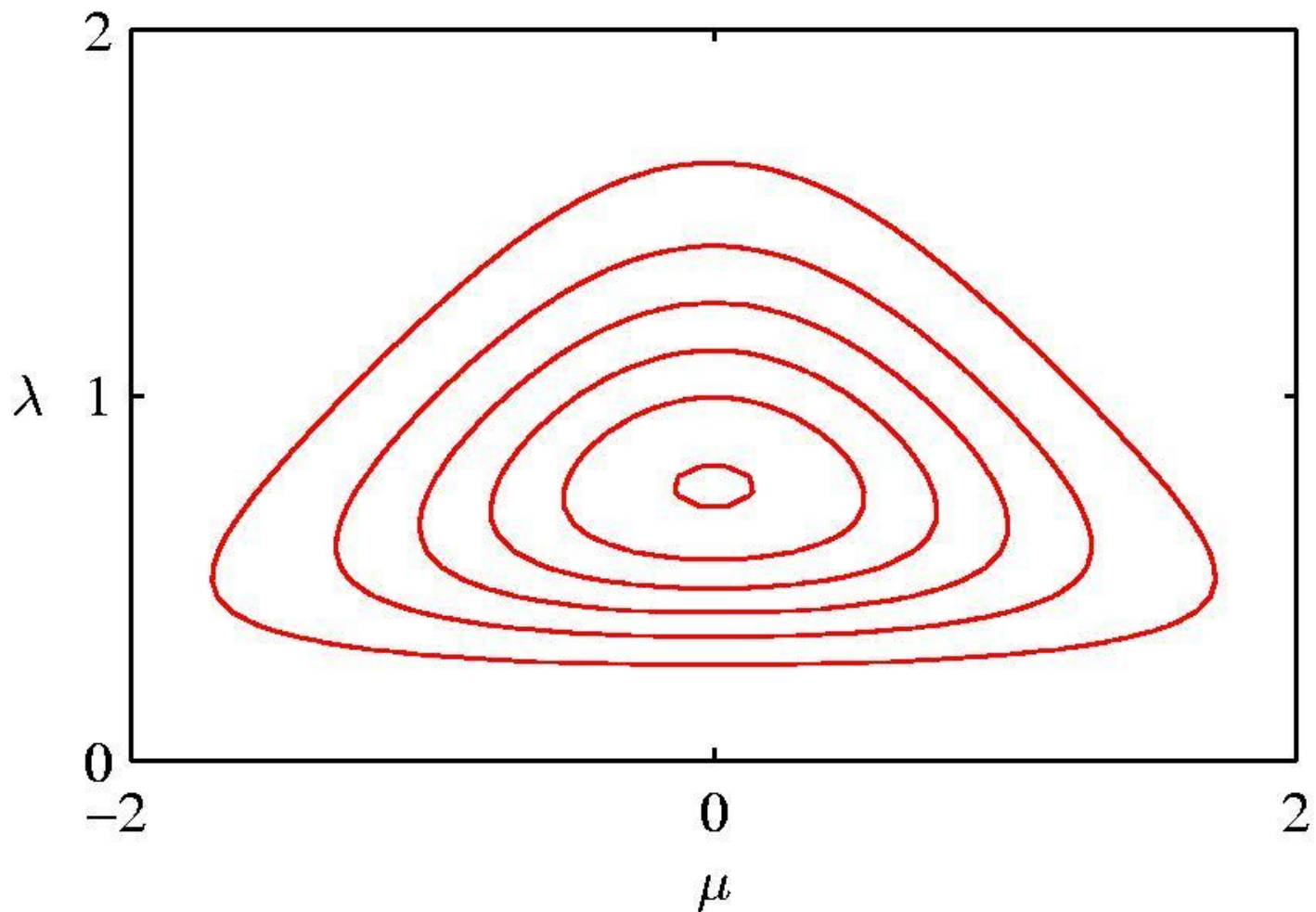
Likelihood function

$$\begin{aligned} p(\mathbf{x}|\mu, \lambda) &= \prod_{n=1}^N \left(\frac{\lambda}{2\pi} \right)^{1/2} \exp \left\{ -\frac{\lambda}{2} (x_n - \mu)^2 \right\} \\ &\propto \left[\lambda^{1/2} \exp \left(-\frac{\lambda \mu^2}{2} \right) \right]^N \exp \left\{ \lambda \mu \sum_{n=1}^N x_n - \frac{\lambda}{2} \sum_{n=1}^N x_n^2 \right\}. \end{aligned}$$

Gaussian-gamma distribution

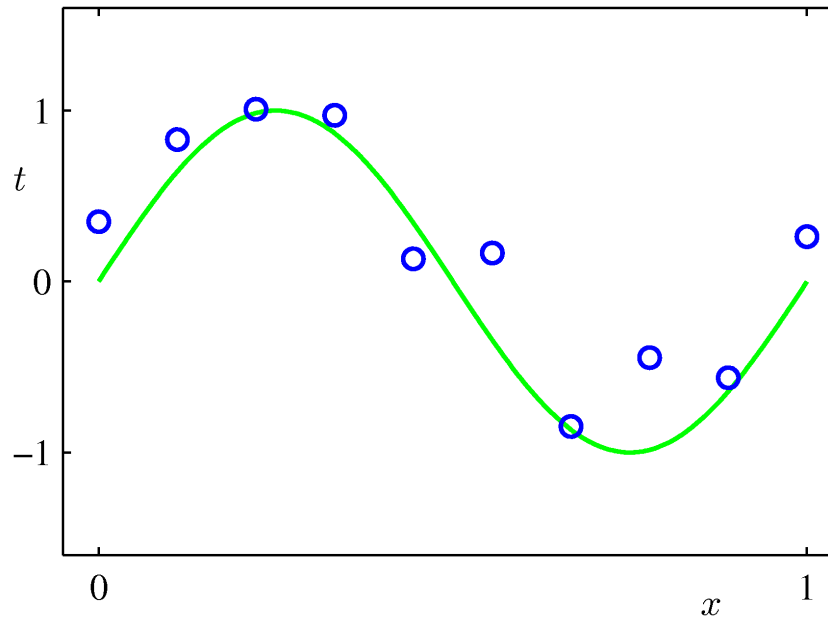
$$\begin{aligned} p(\mu, \lambda) &= p(\mu|\lambda)p(\lambda) = \mathcal{N}(\mu|\mu_0, (\beta\lambda)^{-1}) \text{Gam}(\lambda|a, b) \\ &\propto \exp \left\{ -\frac{\beta\lambda}{2} (\mu - \mu_0)^2 \right\} \lambda^{a-1} \exp \{-b\lambda\} \end{aligned}$$

Gaussian-gamma Distribution



Linear Regression (1)

Noisy sinusoidal data



Linear Regression (2)

Linear combination of basis functions

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

Noise model

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

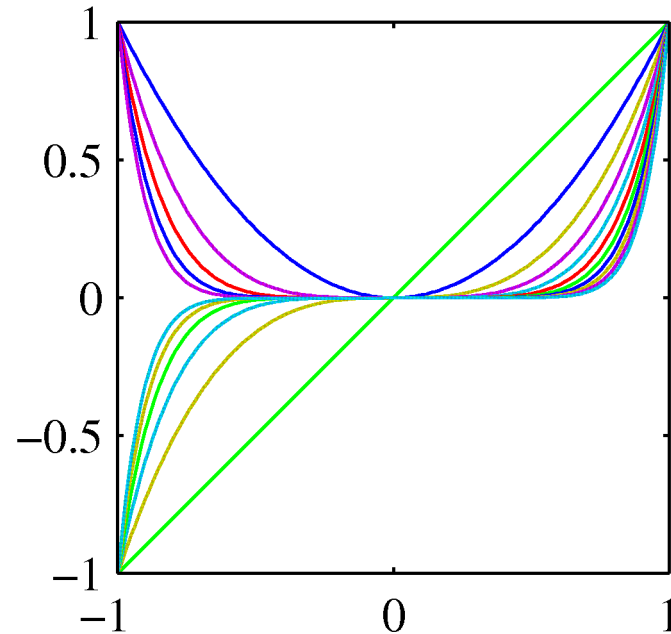
Likelihood function

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

Linear Regression (3)

Polynomial basis functions

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



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

Linear Regression (4)

Define a conjugate prior over \mathbf{w}

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

Combining with likelihood function gives the posterior

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

where

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

$$\Phi_{nj} = \phi_j(x_n)$$

Simple Example (1)

Data from straight line with Gaussian noise

$$t = a + bx + \epsilon$$

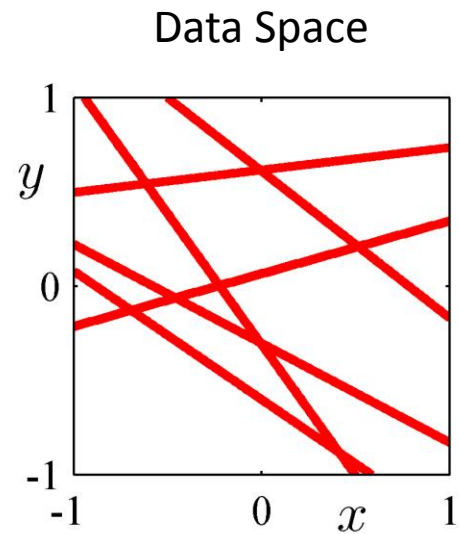
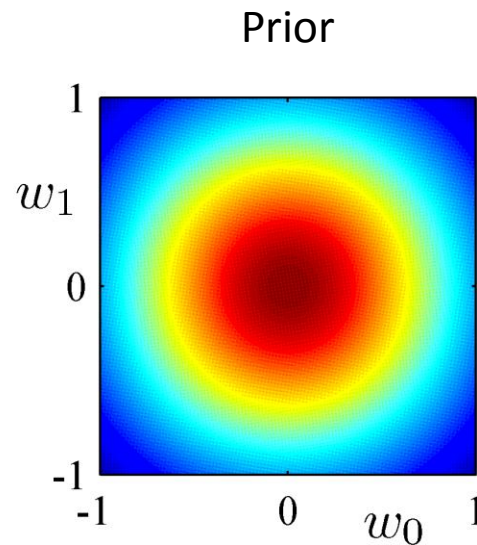
$$\epsilon \sim \mathcal{N}(\cdot|0, 1)$$

First order polynomial model

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

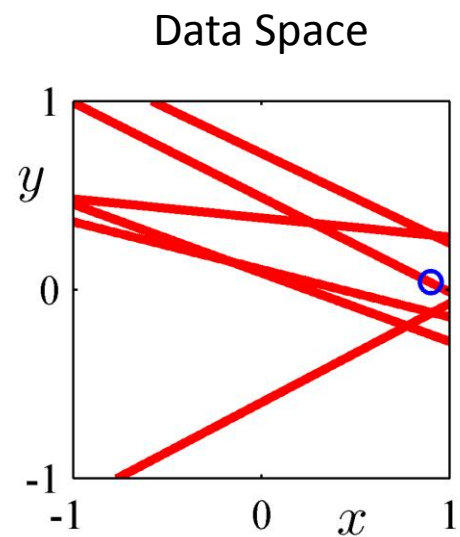
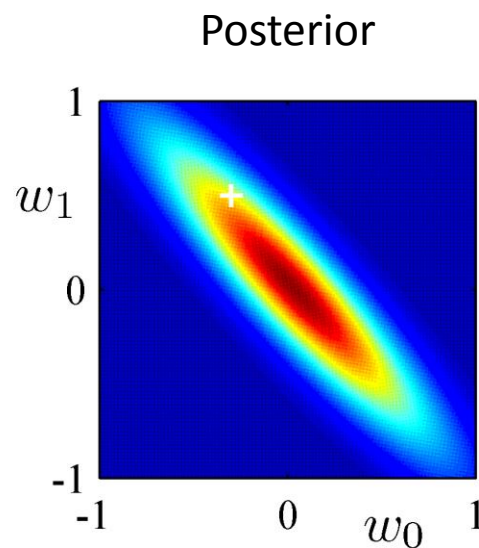
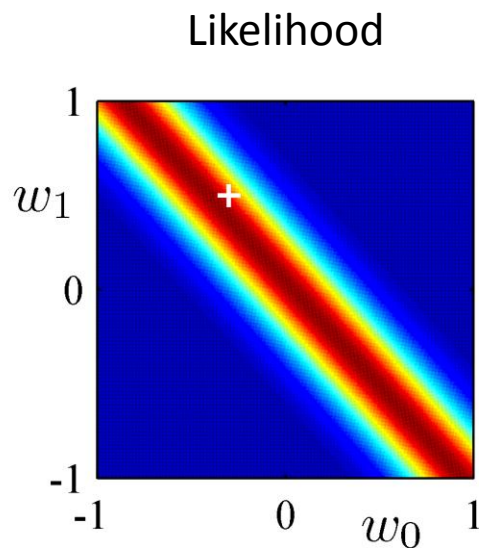
Simple Example (2)

0 data points observed



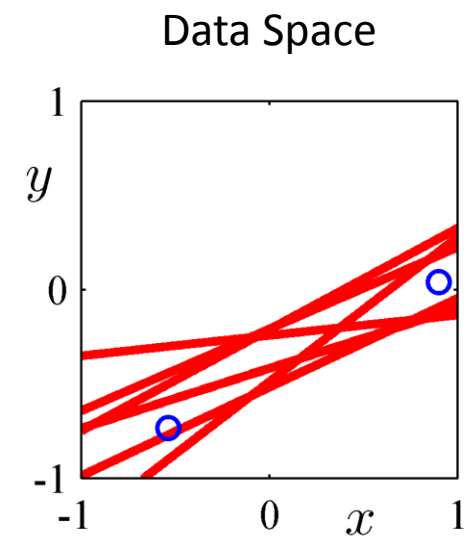
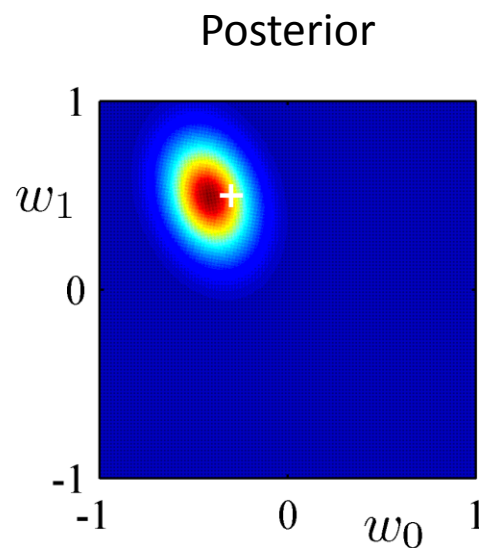
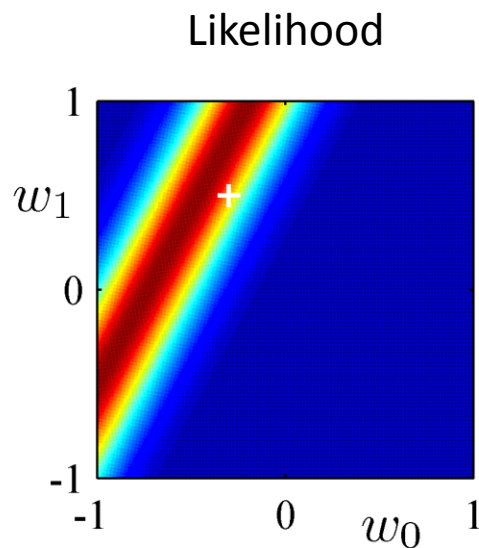
Simple Example (3)

1 data point observed



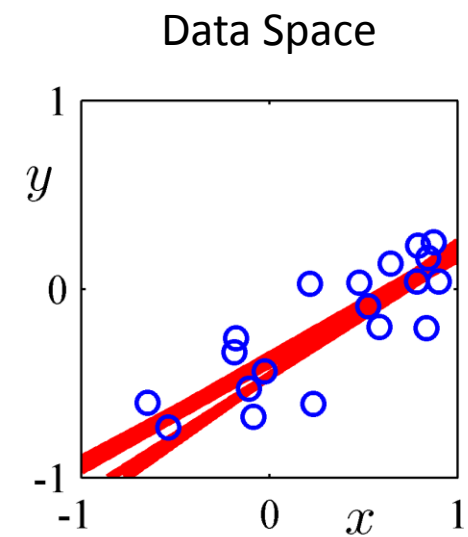
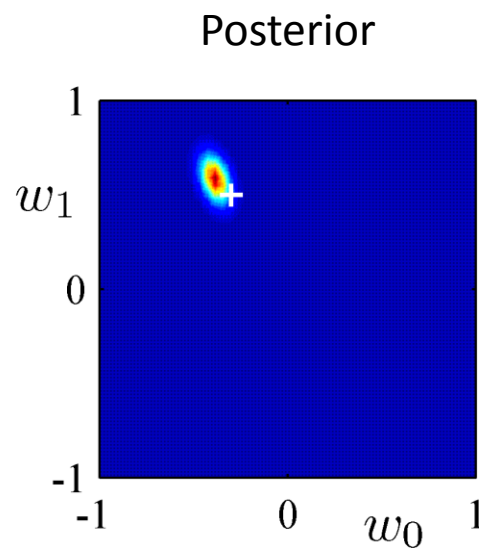
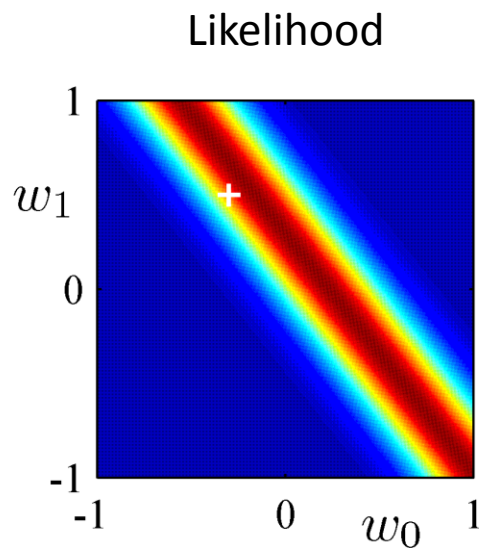
Simple Example (4)

2 data points observed



Simple Example (5)

20 data points observed



Predictive Distribution (1)

Predict t for new values of x by integrating over \mathbf{w} :

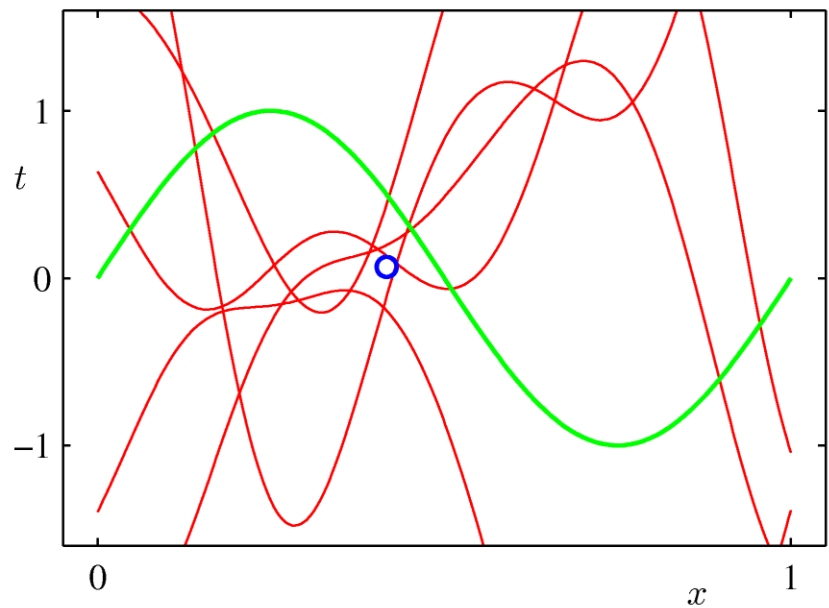
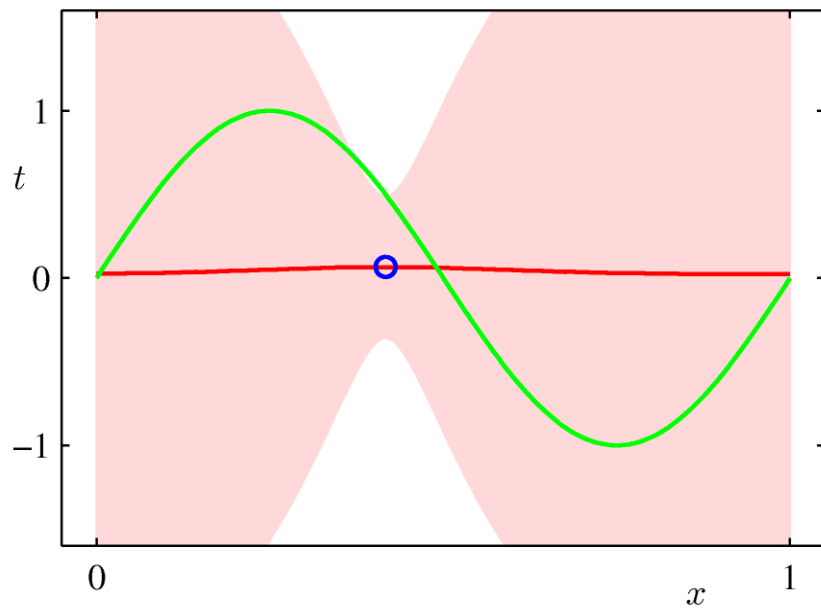
$$\begin{aligned} p(\hat{t}|\mathbf{t}, \alpha, \beta, \hat{x}) &= \int p(\hat{t}|\mathbf{w}, \beta, \hat{x})p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} \\ &= \mathcal{N}(\hat{t}|\mathbf{m}_N^T\boldsymbol{\phi}(\hat{x}), \sigma_N^2(\hat{x})) \end{aligned}$$

where

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

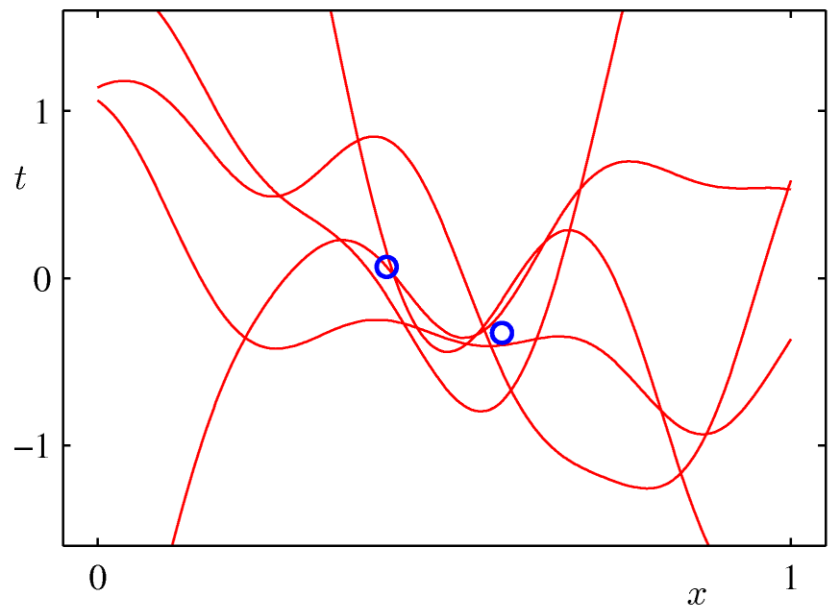
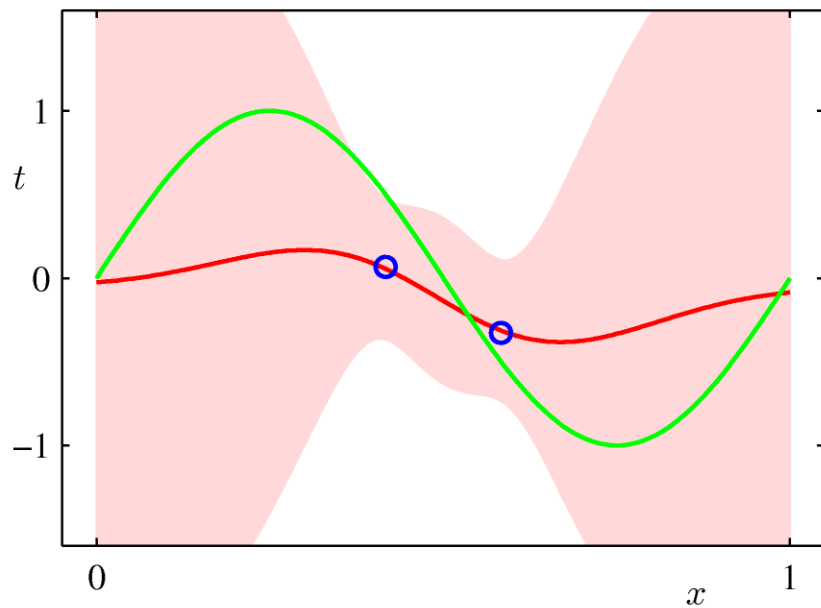
Predictive Distribution (3)

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



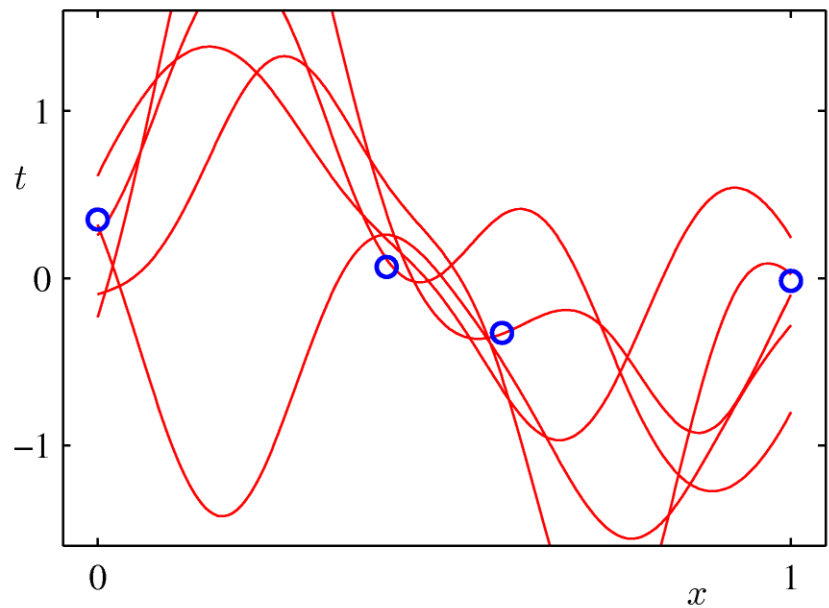
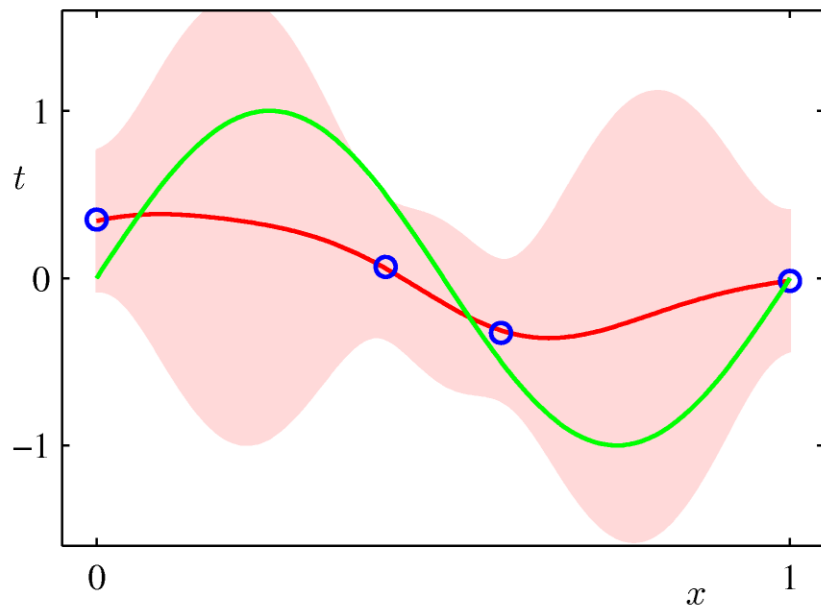
Predictive Distribution (4)

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



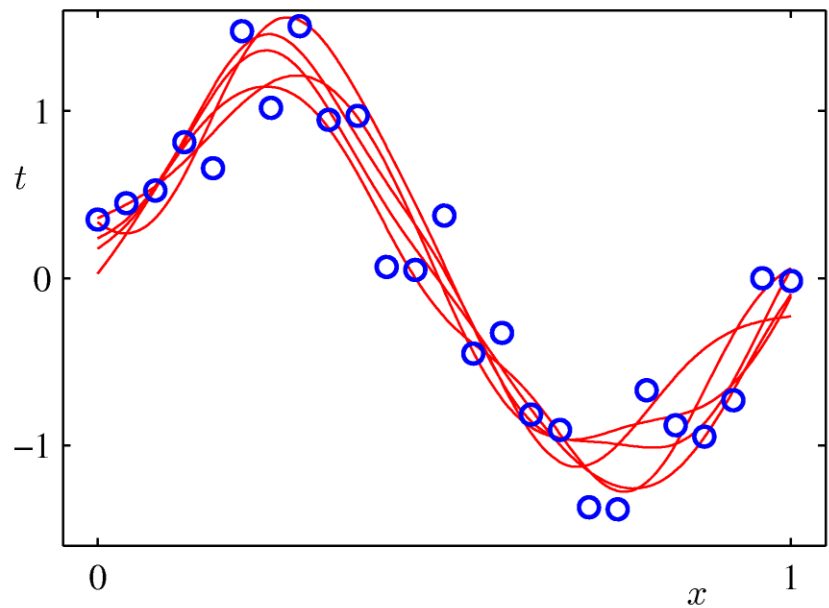
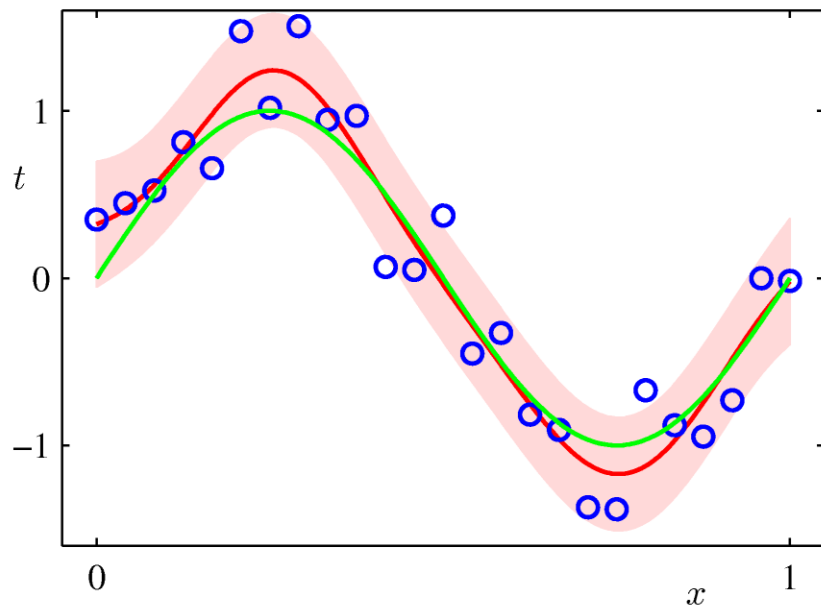
Predictive Distribution (5)

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



Predictive Distribution (6)

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



Bayesian Model Comparison (1)

Alternative models $\mathcal{M}_i, i=1, \dots, L$

Predictive distribution is a mixture

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

Model selection: keep only most probable model

Bayesian Model Comparison (2)

From Bayes' theorem

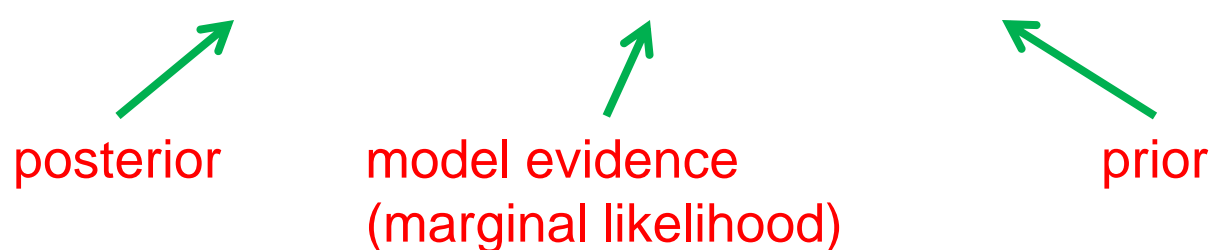
$$p(\mathcal{M}_i|\mathcal{D}) \propto p(\mathcal{D}|\mathcal{M}_i)p(\mathcal{M}_i)$$


Diagram illustrating the components of Bayes' theorem:

- posterior (points to $p(\mathcal{M}_i|\mathcal{D})$)
- model evidence (marginal likelihood) (points to $p(\mathcal{D}|\mathcal{M}_i)$)
- prior (points to $p(\mathcal{M}_i)$)

For equal priors, models ranked by marginal likelihood

Bayesian Model Comparison (4)

For a model with parameters \mathbf{w}

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

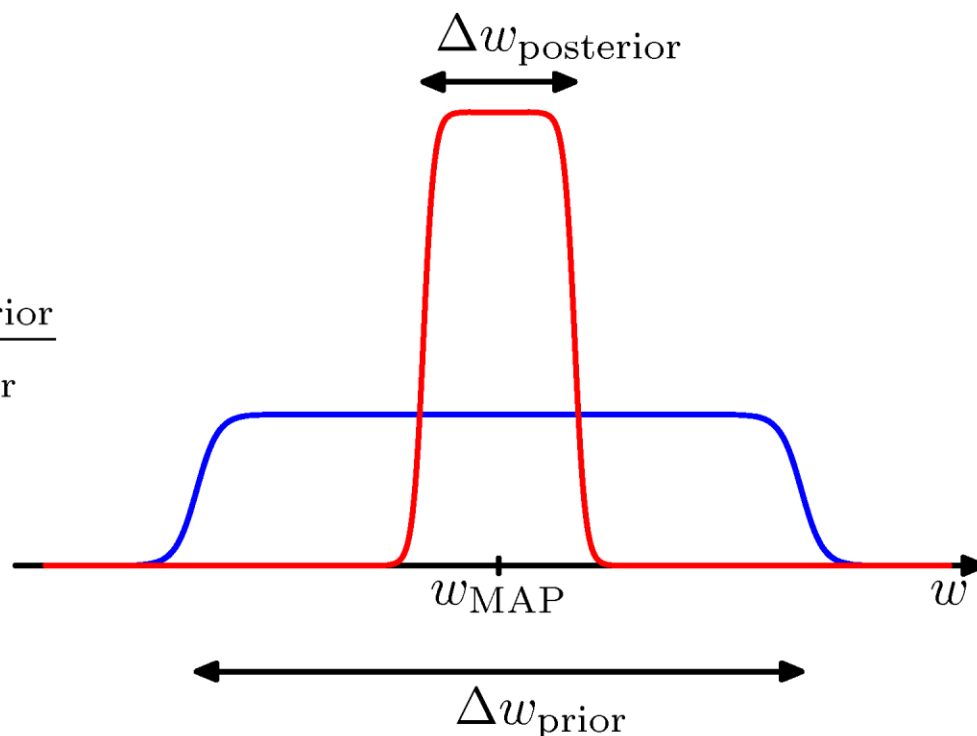
Note that

$$p(\mathbf{w}|\mathcal{D}, \mathcal{M}_i) = \frac{p(\mathcal{D}|\mathbf{w}, \mathcal{M}_i)p(\mathbf{w}|\mathcal{M}_i)}{p(\mathcal{D}|\mathcal{M}_i)}$$

Bayesian Model Comparison (5)

Consider model with a single parameter w

$$p(\mathcal{D}) = \int p(\mathcal{D}|w)p(w) dw$$
$$\simeq p(\mathcal{D}|w_{\text{MAP}}) \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}}$$



Bayesian Model Comparison (6)

Taking logarithms, we obtain

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

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

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\mathbf{w}_{\text{MAP}}) + M \ln \left(\frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \right)$$

Linear Regression revisited

Marginal likelihood

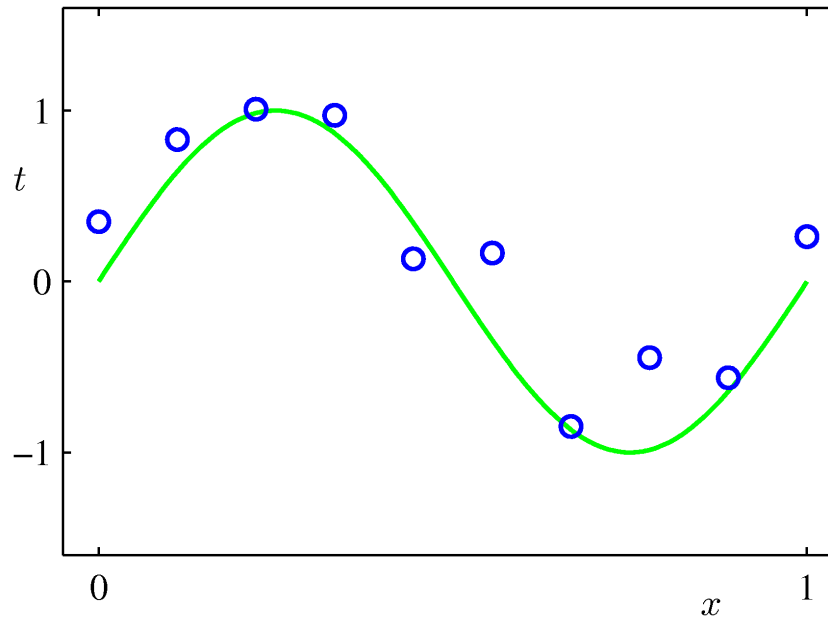
$$p(\mathbf{t}|\alpha, \beta) = \int p(\mathbf{t}|\mathbf{w}, \beta)p(\mathbf{w}|\alpha) d\mathbf{w}$$

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

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

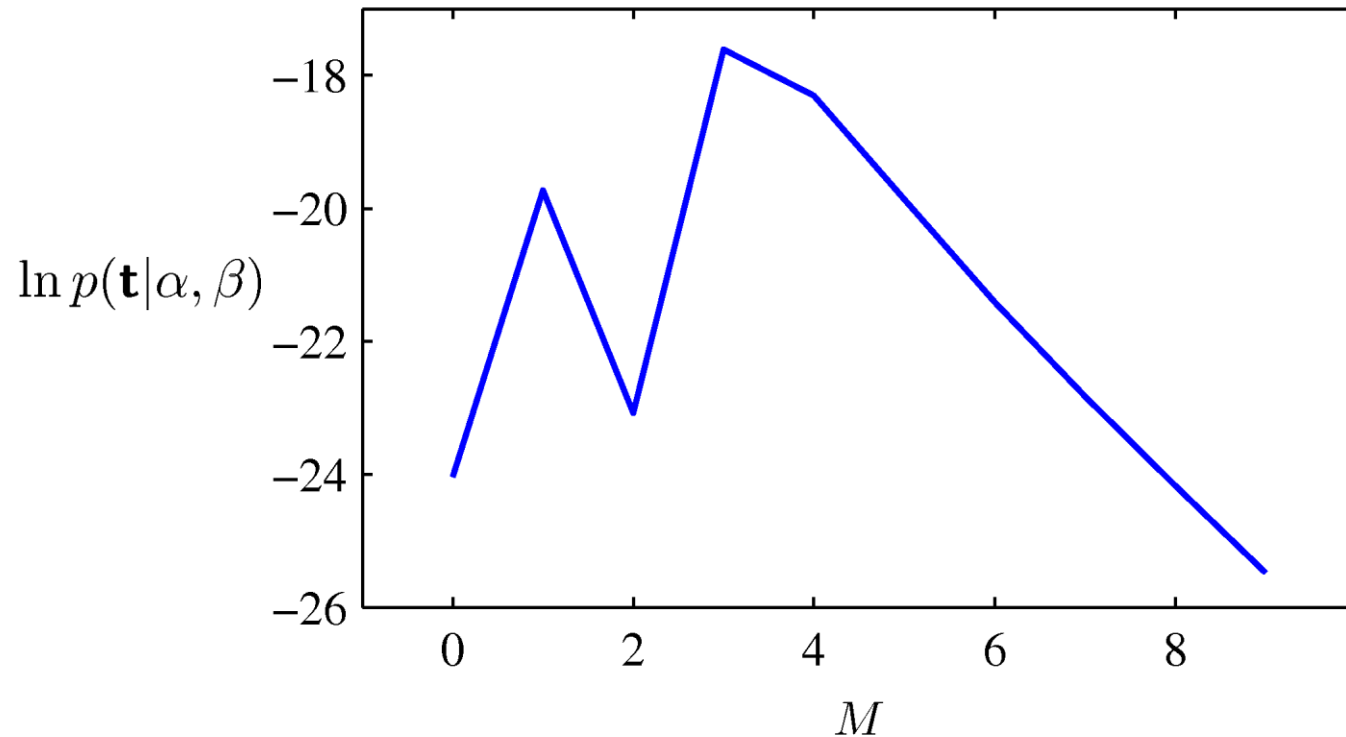
Linear Regression revisited

Noisy sinusoidal data



Linear Regression revisited

Polynomial of order M , $\alpha = 5 \times 10^{-3}$



Bayesian Model Comparison

Matching data and model complexity

