
Bayesian Linear Regression: Model and Variable Selection, Equivalent Kernel Representation

*Prof. Nicholas Zabaras
University of Notre Dame
Notre Dame, IN, USA*

*Email: nzabaras@gmail.com
URL: <https://www.zabaras.com/>*

September 20, 2017



Contents

- Bayesian model comparison, Model Averaging and Model Selection, Model Complexity
- The evidence approximation for our regression example
- Another example of computing model evidence
- Limitations of fixed basis functions
- Laplace approximation, BIC criterion, Another Regression example and MatLab implementation of model selection
- Equivalent Kernel Representation
- Introduction to Variable Selection

Following closely Chris Bishop's PRML book, Chapter 3



Bayesian Model Comparison

- The Bayesian view of model comparison involves the use of probabilities to represent uncertainty in the choice of model.
- Suppose we wish to compare L models $\{\mathcal{M}_i\}$ where $i = 1, \dots, L$. Here a model refers to a probability distribution over the observed data \mathcal{D} . Our uncertainty is expressed through a prior $p(\mathcal{M}_i)$. Given a training set \mathcal{D} , then the posterior distribution is:

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

Posterior

Prior

*Model evidence or
marginal likelihood*

- We have already defined the **Bayes Factor** as the ratio of two model evidences

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



Model Averaging and Model Selection

- Once we know the posterior distribution over models, the predictive distribution is given, by

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

↑
Test data ↑
Training data

- This has the form of a mixture distribution in which the overall predictive distribution is obtained by averaging the predictive distributions $p(t | \mathbf{x}, \mathcal{M}_i, \mathcal{D})$ of individual models, weighted by the posterior probabilities $p(\mathcal{M}_i | \mathcal{D})$ of those models.
- A simple approximation to model averaging is to use the single most probable model alone to make predictions. This is known as model selection.

Model Evidence

- For a model governed by a set of parameters w , the model evidence is given, from the sum and product rules of probability, by

$$p(\mathcal{D} | \mathcal{M}_i) = \int p(\mathcal{D} | w_i, \mathcal{M}_i) p(w_i | \mathcal{M}_i) dw_i$$

- From a sampling perspective, the marginal likelihood can be viewed as **the probability of generating the data set \mathcal{D} from a model whose parameters are sampled at random from the prior.**
- Also we can see the model evidence as normalizing factor:

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

Occam's Razor and Model Selection

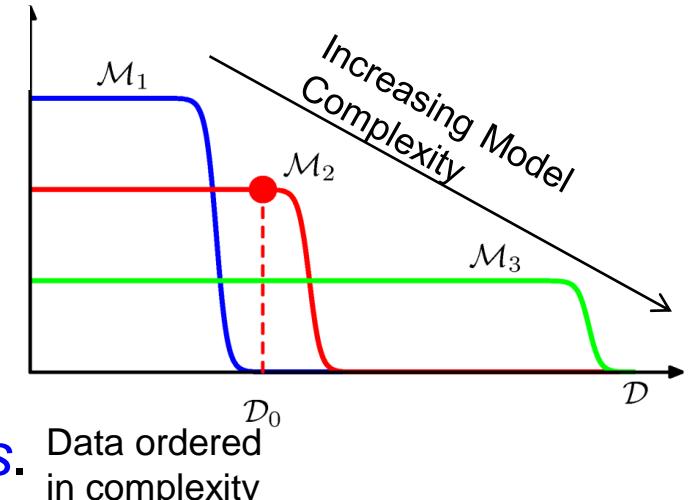
- Compare model classes \mathcal{M}_i using their posterior probability given the data \mathcal{D} :

$$p(\mathcal{M}_i | \mathcal{D}) \propto p(\mathcal{M}_i) p(\mathcal{D} | \mathcal{M}_i) \quad p(\mathcal{D} | \mathcal{M}_i) = \int_{\Theta_i} p(\mathcal{D} | w_i, \mathcal{M}_i) p(w_i | \mathcal{M}_i) dw_i$$

- The marginal likelihood (Bayesian evidence) $p(\mathcal{D} | \mathcal{M}_i)$ is viewed as **the probability that randomly selected parameter values from the model class would generate the data set \mathcal{D} .**

- Simple model classes are unlikely to generate \mathcal{D} .
- Too complex model classes $p(\mathcal{D} | \mathcal{M}_i)$ can generate many data sets so it is unlikely to generate the particular data set \mathcal{D} .

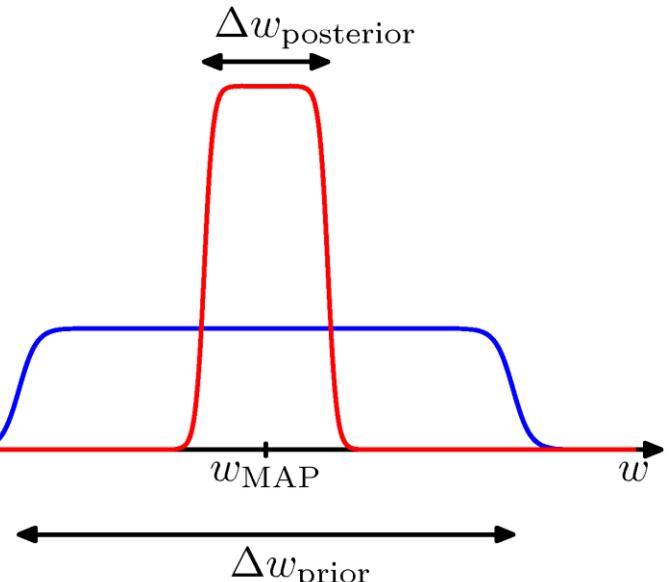
- Bayesian inference automatically implements Occam's Razor Principle:
Prefer simple than complex explanations.



Bayesian Model Comparison

- For a given model (omit $| \mathcal{M}_i$) with a single parameter, w , consider the approximation (use Bayes rule)

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



where the posterior is assumed to be sharply peaked around w_{MAP} (we already have seen the more accurate Laplace approximation)

- Taking logs we obtain

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

Negative

Note: the evidence is not defined if the prior is improper.

Optimal Model Complexity

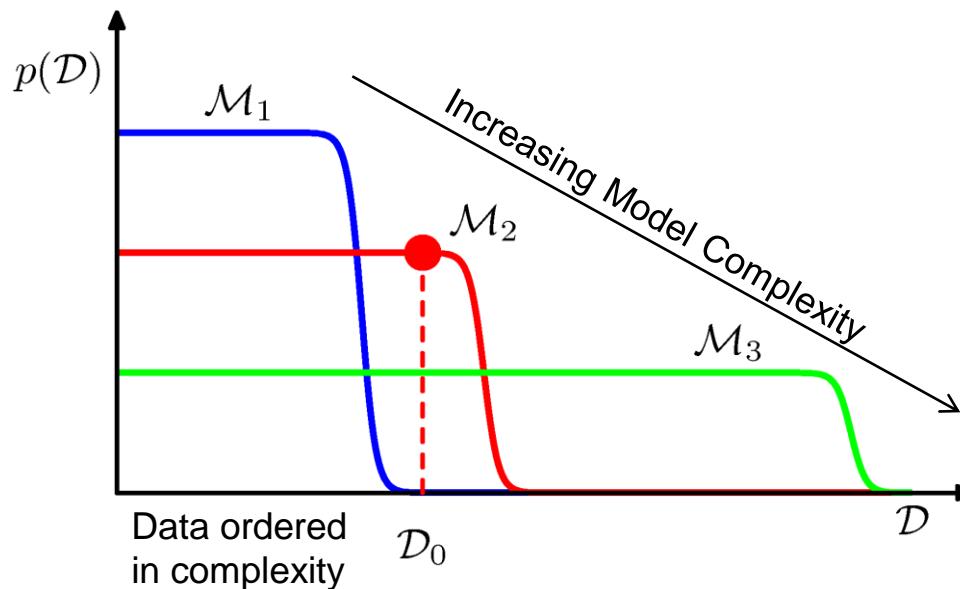
- For a model with M parameters, we can make a similar approximation for each parameter in turn. Assuming that all parameters have the same ratio of $\Delta w_{posterior}/\Delta w_{prior}$,

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

- The size of the complexity penalty increases linearly with M . As we increase the complexity of the model
 - the 1st term increases, because a more complex model is better able to fit the data,
 - whereas the 2nd term decreases due to the dependence on M .
- The optimal model complexity determined by maximum evidence will be given by a trade-off between these two terms.

Matching Data and Model Complexity

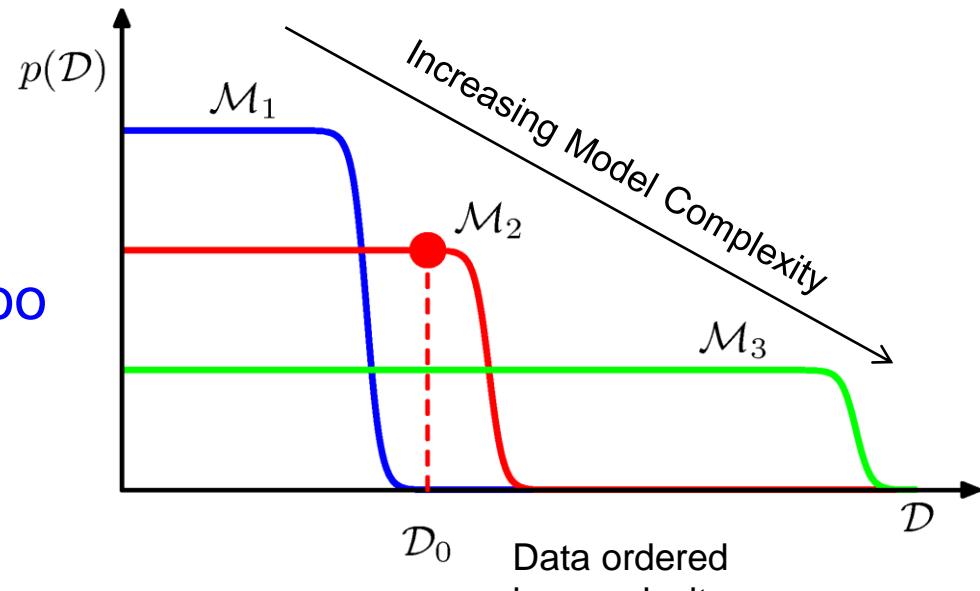
- The marginal likelihood favors models of intermediate complexity.
- Let us think of the regression model and consider the models \mathcal{M}_1 , \mathcal{M}_2 and \mathcal{M}_3 represent linear, quadratic and cubic fitting.
- The data \mathcal{D} are ordered in complexity – for a given model, we choose w from the prior $p(w)$, then sample the data from $p(\mathcal{D}|w)$.



Matching Data and Model Complexity

- A 1st order polynomial has little variability, generates data that are similar, $p(\mathcal{D})$ is confined to a small region in the \mathcal{D} axis.
- A 9th order polynomial generates a variety of different data, and so its $p(\mathcal{D})$ is spread over a large region in the \mathcal{D} axis.
- Because $p(\mathcal{D}|\mathcal{M}_i)$ are normalized, a particular \mathcal{D}_0 can have the highest evidence for the model of intermediate complexity.

- The simpler model cannot fit the data well, whereas the more complex model spreads its predictive probability over too broad a range of data sets.



Matching Data and Model Complexity

- A Bayesian model comparison in an average (over the data \mathcal{D}) sense will favor the correct model.
- Let \mathcal{M}_1 be the correct model and \mathcal{M}_2 another model. We can show that the evidence for model \mathcal{M}_1 is higher. Using the definition and properties of the Kullback-Leibler distance:

$$KL\left(p(\mathcal{D} | \mathcal{M}_1) \| p(\mathcal{D} | \mathcal{M}_2)\right) = \int \underbrace{p(\mathcal{D} | \mathcal{M}_1)}_{\substack{\text{Averaged with the} \\ \text{exact probability}}} \ln \frac{p(\mathcal{D} | \mathcal{M}_1)}{\underbrace{p(\mathcal{D} | \mathcal{M}_2)}_{\text{Bayes factor}}} d\mathcal{D} \geq 0$$

- This analysis assumes that the true distribution from which the data are generated is contained in our class of models.

The Evidence Approximation

- The fully Bayesian predictive distribution for our regression model is given by

The hyperparameters α and β are now random variables

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

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

Dependence on x and \mathbf{x} not shown to simplify the notation

but this integral is intractable. Approximate with

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

where $(\hat{\alpha}, \hat{\beta})$ is the mode of $p(\alpha, \beta|t)$, which is assumed to be sharply peaked.

a.k.a. empirical Bayes, type II or generalized maximum likelihood, or evidence approximation.



The Evidence Approximation

- From Bayes' theorem, the posterior distribution for α and β is given by

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

- If the prior is relatively flat, then in the evidence framework the values of $(\hat{\alpha}, \hat{\beta})$ are obtained by maximizing the marginal likelihood function $p(\mathbf{t} | \alpha, \beta)$.
- The marginal likelihood function $p(\mathbf{t} | \alpha, \beta)$ is obtained by integrating over the parameters \mathbf{w} , so that

$$p(\alpha, \beta | \mathbf{t}) \propto p(\mathbf{t} | \alpha, \beta) = \underbrace{\int p(\mathbf{t} | \mathbf{w}, \beta)}_{\frac{1}{(2\pi)^{N/2}} \beta^{N/2} e^{-\beta E_D}} \underbrace{p(\mathbf{w} | \alpha)}_{\mathcal{N}(\mathbf{w}, \alpha^{-1} \mathbf{I}_{M \times M})} d\mathbf{w}$$

Evidence/
Marginal Likelihood

- One can evaluate this integral using the completion of the square procedure typical of Gaussian marginalizations.



The Evidence Approximation

- We can write the evidence function in the form

$$p(t/\alpha, \beta) = \left(\frac{\beta}{2\pi}\right)^{N/2} \left(\frac{\alpha}{2\pi}\right)^{M/2} \int \exp\{-E(w)\} dw$$

where M is the dimensionality of w , and we have defined

$$\begin{aligned} E(w) &= \beta E_D(w) + \alpha E_W(w) = \frac{\beta}{2} \|t - \Phi w\|^2 + \frac{\alpha}{2} w^T w \\ E(w) &= E(m_N) + \frac{1}{2} (w - m_N)^T A (w - m_N) \end{aligned}$$

- We have introduced here:

$$A = \alpha I + \beta \Phi^T \Phi, \quad E(m_N) = \frac{\beta}{2} \|t - \Phi m_N\|^2 + \frac{\alpha}{2} m_N^T m_N, \quad m_N = \beta A^{-1} \Phi^T t$$

- Note that the Hessian matrix A corresponds to the matrix of 2nd derivatives of the error function:

$$A = \nabla \nabla E(w)$$



The Evidence Approximation

- The integral over \mathbf{w} can now be evaluated simply by appealing to the standard result for the normalization coefficient of a multivariate Gaussian, giving

$$\begin{aligned}\int \exp\{-E(\mathbf{w})\} d\mathbf{w} &= \exp\{-E(\mathbf{m}_N)\} \int \exp\left\{-\frac{1}{2}(\mathbf{w} - \mathbf{m}_N)^T \mathbf{A}(\mathbf{w} - \mathbf{m}_N)\right\} d\mathbf{w} \\ &= \exp\{-E(\mathbf{m}_N)\} (2\pi)^{M/2} |\mathbf{A}|^{-1/2}\end{aligned}$$

- We can then write the log of the marginal likelihood in the form

$$p(\mathbf{t} | \alpha, \beta) = \left(\frac{\beta}{2\pi}\right)^{N/2} \left(\frac{\alpha}{2\pi}\right)^{M/2} \int \exp\{-E(\mathbf{w})\} d\mathbf{w} = \left(\frac{\beta}{2\pi}\right)^{N/2} \left(\frac{\alpha}{2\pi}\right)^{M/2} e^{-E(\mathbf{m}_N)} (2\pi)^{M/2} |\mathbf{A}|^{-1/2} \Rightarrow$$

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

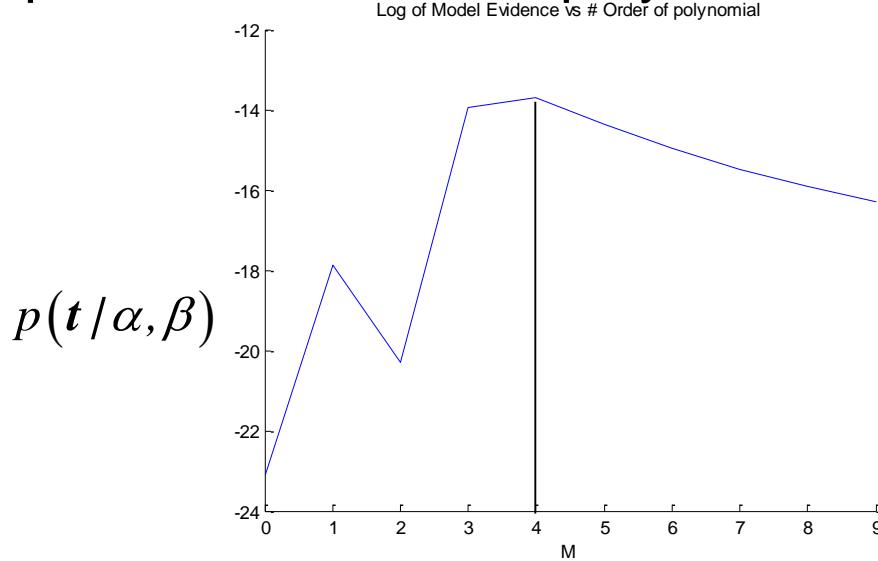
M : number of parameters in the model

N : number of data



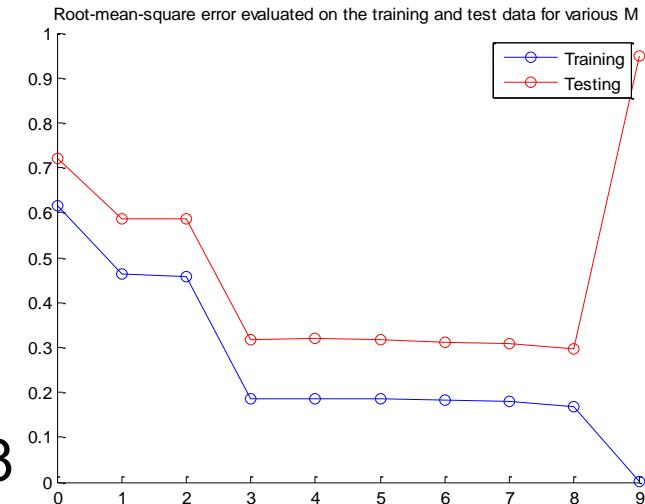
The Evidence Approximation

- Example: sinusoidal data, polynomial regression, $\alpha = 5 \times 10^{-3}$, $\beta = 11.1$



[MatLab Code](#)
and [data](#)

- From the plot of the model evidence for given α and β , we see that the evidence favors the model with $M = 5$ (4th degree polynomial)
- Looking at the non-Bayesian approach, one cannot distinguish the performance of polynomials of order 3....8



Maximizing the Evidence Function

- Let us first consider the maximization of $p(\mathbf{t}|\alpha, \beta)$ with respect to α . This can be done by first defining the following eigenvector equation

$$(\beta \Phi^T \Phi) \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

- Thus, $A = \alpha I + \beta \Phi^T \Phi$ has eigenvalues $\alpha + \lambda_i$.
- Now consider the derivative of the term involving $\ln|A|$ with respect to α . We have

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

$$\frac{d}{d\alpha} \ln |A| = \frac{d}{d\alpha} \ln \prod_i (\lambda_i + \alpha) = \frac{d}{d\alpha} \sum_i \ln(\lambda_i + \alpha) = \sum_i \frac{1}{\lambda_i + \alpha}$$



Maximizing the Evidence Function

$$\frac{d}{d\alpha} \ln |A| = \frac{d}{d\alpha} \ln \prod_i (\lambda_i + \alpha) = \frac{d}{d\alpha} \sum_i \ln(\lambda_i + \alpha) = \sum_i \frac{1}{\lambda_i + \alpha}$$

□ Thus the stationary points of

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

$$E(\mathbf{m}_N) = \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^T \mathbf{m}_N, \quad \mathbf{m}_N = \beta A^{-1} \Phi^T \mathbf{t}, \quad A = \alpha I + \beta \Phi^T \Phi$$

with respect to α satisfy $0 = \frac{M}{2\alpha} - \frac{1}{2} \mathbf{m}_N^T \mathbf{m}_N - \frac{1}{2} \sum_i \frac{1}{\lambda_i + \alpha}$

□ Multiplying through by 2α and rearranging, we obtain

$$\alpha \mathbf{m}_N^T \mathbf{m}_N = M - \alpha \sum_i \frac{1}{\lambda_i + \alpha} = \sum_i \left(1 - \frac{\alpha}{\lambda_i + \alpha} \right) = \sum_i \frac{\lambda_i}{\lambda_i + \alpha} \equiv \gamma$$

Implicit solution
for α

$$\alpha = \frac{\gamma}{\mathbf{m}_N^T \mathbf{m}_N}$$

1. Choose α
2. Calculate \mathbf{m}_N, γ
3. Re-estimate α



Maximizing the Evidence Function

Implicit Solution for Computing α

1. Choose α
2. Calculate \mathbf{m}_N, γ :

$$\mathbf{m}_N = \beta \mathbf{A}^{-1} \Phi^T \mathbf{t}, \quad \mathbf{A} = a\mathbf{I} + \beta \Phi^T \Phi, \quad (\beta \Phi^T \Phi) \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

$$\gamma = \sum_i \frac{\lambda_i}{\lambda_i + \alpha}$$

3. Re-estimate α

$$\alpha = \frac{\gamma}{\mathbf{m}_N^T \mathbf{m}_N}$$



Maximizing the Evidence Function

- We can similarly maximize the log marginal likelihood with respect to β .

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

- To do this, we note that the eigenvalues λ_i defined by

$$(\beta \Phi^T \Phi) \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

are proportional to β , and hence $d\lambda_i / d\beta = \lambda_i / \beta$ giving

$$\frac{d}{d\beta} \ln |\mathbf{A}| = \frac{d}{d\beta} \ln \prod_i (\lambda_i + \alpha) = \frac{1}{\beta} \sum_i \frac{\lambda_i}{\lambda_i + \alpha} = \frac{\gamma}{\beta}$$

Maximizing the Evidence Function

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

$$E(\mathbf{m}_N) = \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^T \mathbf{m}_N \quad \frac{d}{d\beta} \ln |\mathbf{A}| = \frac{1}{\beta} \sum_i \frac{\lambda_i}{\lambda_i + \alpha} = \frac{\gamma}{\beta}$$

$$\mathbf{m}_N = \beta \mathbf{A}^{-1} \Phi^T \mathbf{t}, \quad \mathbf{A} = \alpha \mathbf{I} + \beta \Phi^T \Phi, \quad (\beta \Phi^T \Phi) \mathbf{u}_i = \lambda_i \mathbf{u}_i, \quad \gamma = \sum_i \frac{\lambda_i}{\lambda_i + \alpha}$$

- Setting the derivative wrt β equal to zero, the stationary point of the marginal likelihood therefore satisfies

$$0 = \frac{N}{2\beta} - \frac{1}{2} \sum_{n=1}^N \left\{ \mathbf{t}_n - \mathbf{m}_N^T \boldsymbol{\phi}(x_n) \right\}^2 - \frac{\gamma}{2\beta}$$

- Rearranging we obtain

$$\frac{1}{\beta} = \frac{1}{N - \gamma} \sum_{n=1}^N \left\{ \mathbf{t}_n - \mathbf{m}_N^T \boldsymbol{\phi}(x_n) \right\}^2$$

Implicit solution
for β

1. Choose β
2. Calculate \mathbf{m}_N, γ
3. Re-estimate β



Maximizing the Evidence Function

- It is interesting to note that in the evidence framework (using the optimal computed values of α and β), the following is true:

$$E(\mathbf{m}_N) = \frac{N}{2}$$

- This can be easily shown using the earliest derived results:

$$E(\mathbf{m}_N) = \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^T \mathbf{m}_N$$

with

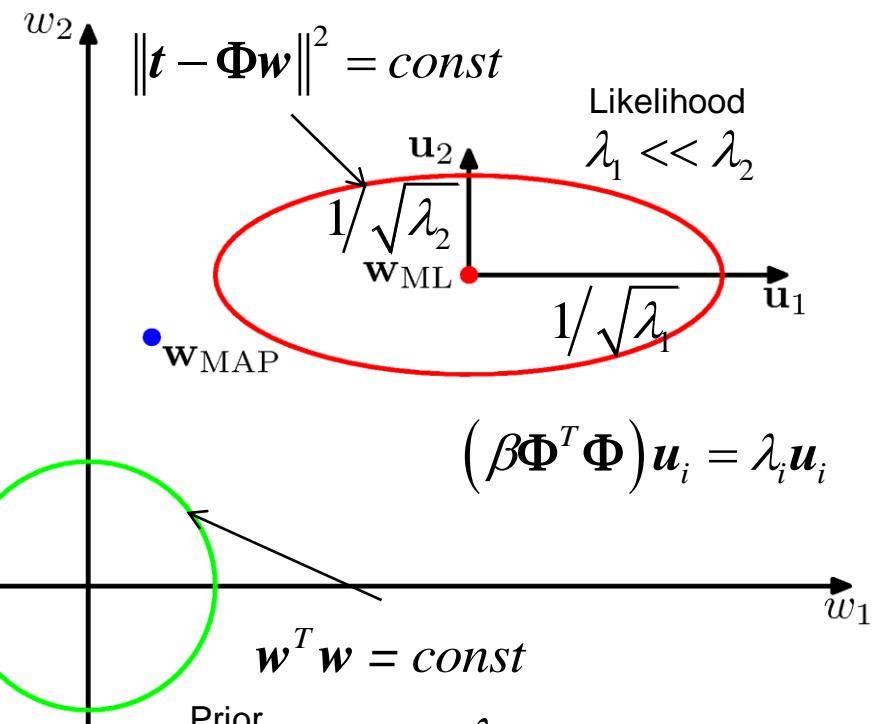
$$\alpha = \frac{\gamma}{\mathbf{m}_N^T \mathbf{m}_N}$$

and

$$\frac{1}{\beta} = \frac{1}{N - \gamma} \sum_{n=1}^N \left\{ \mathbf{t}_n - \mathbf{m}_N^T \boldsymbol{\phi}(x_n) \right\}^2 = \frac{1}{N - \gamma} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2$$

Effective Number of Parameters

- Consider the contours of the likelihood & prior in which the axes in parameter space have been rotated to align with the eigenvectors \mathbf{u}_i .
- For $\alpha = 0$, the mode of the posterior is given by the MLE solution \mathbf{w}_{ML} , whereas for nonzero α the mode is at $\mathbf{w}_{\text{MAP}} = \mathbf{m}_N$.
- In the direction w_1 , λ_1 is small compared with α and so $\lambda_1/(\lambda_1 + \alpha)$ is close to zero, and the corresponding MAP value of w_1 , $w_{1\text{MAP}} = \frac{\lambda_1}{\lambda_1 + \alpha} w_{1\text{MLE}}$ is also close to zero.
- In the direction w_2 , λ_2 is large compared with α and so $\lambda_2/(\lambda_2 + \alpha)$ is close to unity, and the MAP value of w_2 is close to its MLE value.



$$0 < \frac{\lambda_i}{\lambda_i + \alpha} \leq 1,$$

$$0 \leq \gamma = \sum_i \frac{\lambda_i}{\lambda_i + \alpha} \leq M$$

Effective Number of Parameters

- In directions w_i , $\lambda_i \ll \alpha$, so $\lambda_i/(\lambda_i + \alpha)$ is close to zero, and the corresponding MAP value of w_i is also close to zero.

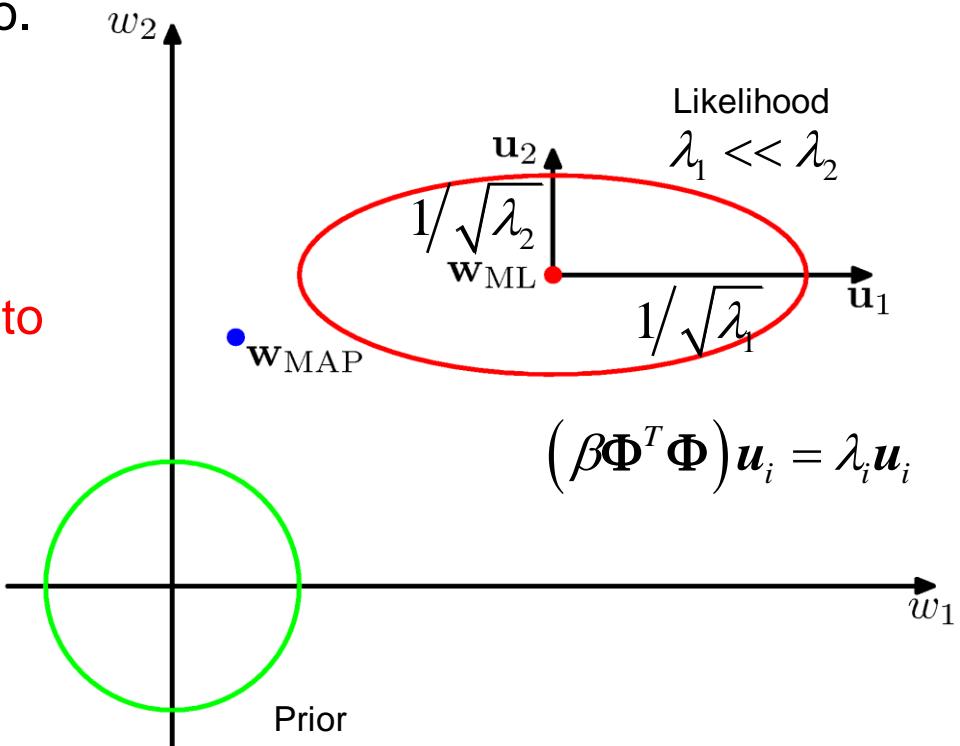
These are directions in which the likelihood function is relatively insensitive to the parameter value and so the parameter has been set to a small value by the prior.

- The quantity γ

$$0 < \frac{\lambda_i}{\lambda_i + \alpha} \leq 1,$$

$$0 \leq \gamma = \sum_i \frac{\lambda_i}{\lambda_i + \alpha} \leq M$$

therefore measures the effective total number of well determined parameters.



Effective Number of Parameters

- We can obtain some insight into the equation for β

$$\frac{1}{\beta} = \frac{1}{N - \gamma} \sum_{n=1}^N \{t_n - \mathbf{m}_N^T \phi(x_n)\}^2$$

by comparing it with the MLE result derived earlier:

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \{t_n - \mathbf{m}_N^T \phi(x_n)\}^2$$

- These formulas express the variance as an average of the squared differences between the targets and model predictions.
- They differ in that the # of data points
 - N in the MLE result is replaced by $N - \gamma$ in the Bayesian result.

Effective Number of Parameters

$$\frac{1}{\beta} = \frac{1}{N - \gamma} \sum_{n=1}^N \left\{ t_n - \mathbf{m}_N^T \boldsymbol{\phi}(x_n) \right\}^2$$

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \left\{ t_n - \mathbf{m}_N^T \boldsymbol{\phi}(x_n) \right\}^2$$

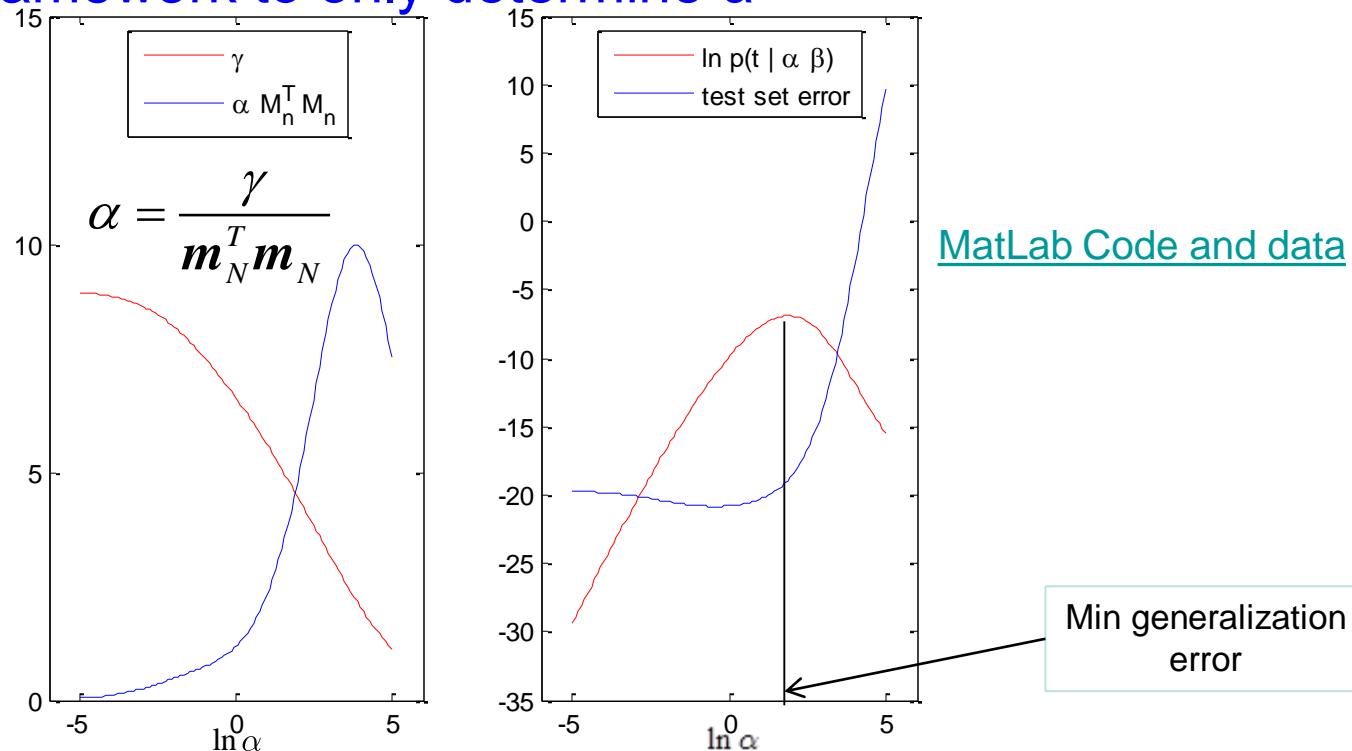
- The effective number of parameters determined by the data is γ .
- The remaining $M - \gamma$ parameters are set to small values by the prior.
- This is reflected in the Bayesian result for the variance that has a factor $N - \gamma$ in the denominator correcting for the bias of the MLE.
- These results are analogous to the estimation of the variance of a Gaussian:

$$\sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N \left\{ x_n - \mu_{ML} \right\}^2 \quad vs. \quad \sigma_{MAP}^2 = \frac{1}{N-1} \sum_{n=1}^N \left\{ x_n - \mu_{ML} \right\}^2$$

1 degree of freedom has been used to fit the mean and the MAP estimate for the variance accounts for that.

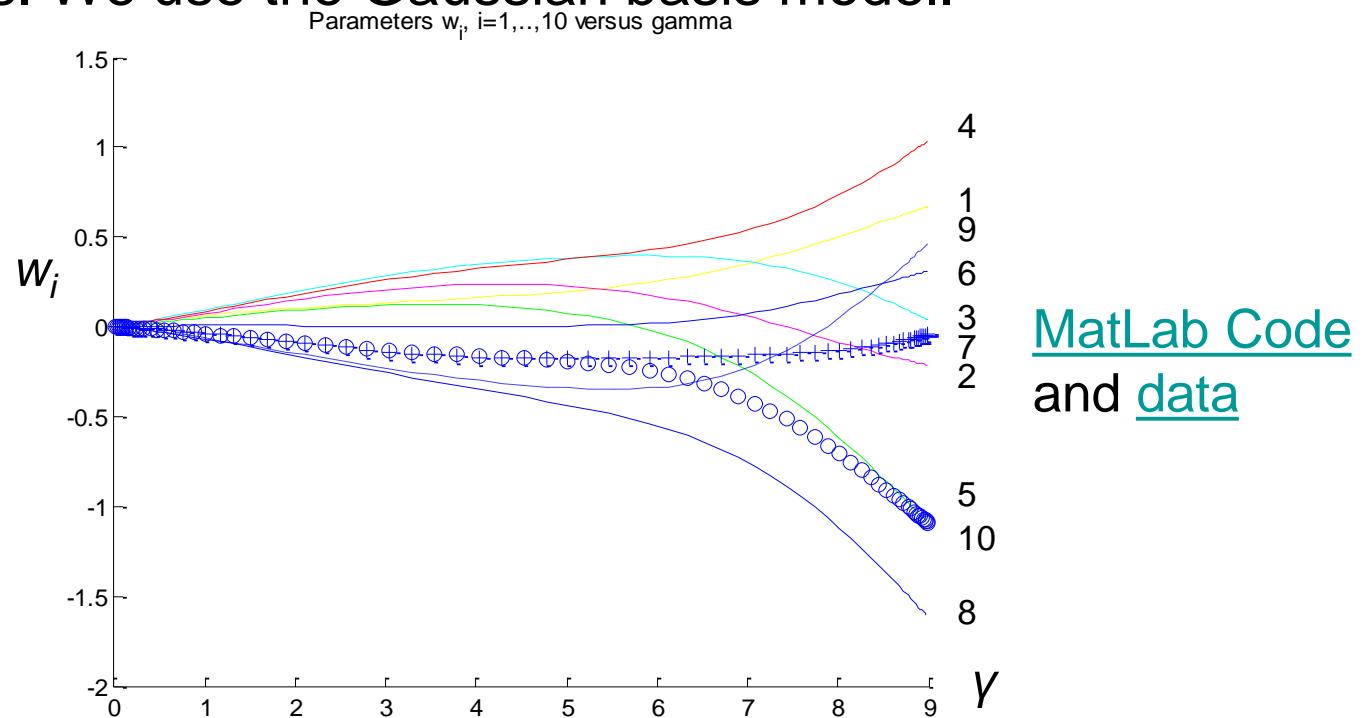
Effective Number of Parameters

- We illustrate the evidence framework for setting hyperparameters using the sinusoidal synthetic data, together with the 9 Gaussian basis functions. The total # of parameters is thus $M = 10$ including the bias.
- For simplicity, we set $\beta = 11.1$ (true value) and **use the evidence framework to only determine α**



Effective Number of Parameters

- We can also see how α controls the magnitude of the parameters $\{w_i\}$, by plotting the individual parameters (posterior means) versus the effective number γ of parameters. We use the Gaussian basis model.



- For the simulation, α is varied $0 \leq \alpha \leq \infty$ causing γ to vary in the range $0 \leq \gamma \leq M$.

Case of $N \gg M$

- For $N \gg M$, all of the parameters are well determined by the data because $\Phi^T\Phi$ involves an implicit sum over data points, and so the eigenvalues λ_i increase with the size of the data set.
- In this case, $\gamma = M$, and the re-estimation equations for α and β become

$$\alpha = \frac{\gamma}{\mathbf{m}_N^T \mathbf{m}_N} = \frac{M}{2E_W(\mathbf{m}_N)} = \frac{M}{\mathbf{m}_N^T \mathbf{m}_N} \quad (\gamma = M)$$

$$\frac{1}{\beta} = \frac{1}{N} \sum_{n=1}^N \left\{ t_n - \mathbf{m}_N^T \boldsymbol{\phi}(x_n) \right\}^2 \quad (N \gg M)$$

- These results are useful as they do not require computing the eigenspectrum of the Hessian.

Another Example of Model Evidence

- We have seen in an earlier lecture that the conjugate prior for a Gaussian with unknown mean and unknown precision is a normal-gamma distribution.
- We can apply the same for the case of our likelihood

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

for which the conjugate prior for w and β is:

$$p(w, \beta) = \mathcal{N}(w | m_0, \beta^{-1} S_0) \mathcal{G}am(\beta | a_0, b_0)$$

- It can be shown that the corresponding posterior for this takes the form:

$$p(w, \beta | t) = \mathcal{N}(w | m_N, \beta^{-1} S_N) \mathcal{G}am(\beta | a_N, b_N)$$

Posterior Distribution

□ The posterior takes the form:

$$\begin{aligned} p(\mathbf{w}, \beta | \mathbf{t}) &\propto \beta^{N/2} \exp \left\{ -\frac{1}{2} \mathbf{w}^T \beta \Phi^T \Phi \mathbf{w} - \beta \mathbf{w}^T \Phi^T \mathbf{t} - \frac{1}{2} \beta \sum_{n=1}^N t_n^2 \right\} \\ &\quad \beta^{M/2} \exp \left\{ -\frac{1}{2} (\mathbf{w} - \mathbf{m}_0)^T \beta S_0^{-1} (\mathbf{w} - \mathbf{m}_0) \right\} \beta^{a_0-1} e^{-b_0 \beta} \\ &\propto \beta^{N/2} e^{-\beta(\mathbf{w}-\mathbf{m}_N)^T S_N^{-1}(\mathbf{w}-\mathbf{m}_N)} e^{\frac{1}{2} \mathbf{m}_N^T \beta S_N^{-1} \mathbf{m}_N} e^{-\frac{1}{2} \beta \sum_{n=1}^N t_n^2} \beta^{M/2} e^{-\frac{1}{2} \mathbf{m}_0^T \beta S_0^{-1} \mathbf{m}_0} \beta^{a_0-1} e^{-b_0 \beta} \\ &\propto \beta^{M/2} e^{-\beta(\mathbf{w}-\mathbf{m}_N)^T S_N^{-1}(\mathbf{w}-\mathbf{m}_N)} \frac{1}{\Gamma(a_N)} b_N^{a_N} \beta^{a_N-1} e^{-b_N \beta} \\ &\boxed{\propto \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \beta^{-1} S_N) \mathcal{G}am(\beta | a_N, b_N)} \end{aligned}$$

where:

$$\begin{aligned} \mathbf{m}_N &= S_N \left[S_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t} \right], \quad S_N^{-1} = \left[S_0^{-1} + \Phi^T \Phi \right] \\ a_N &= a_0 + \frac{N}{2}, \quad b_N = b_0 + \frac{1}{2} \left(\mathbf{m}_0^{-1} S_0^{-1} \mathbf{m}_0 - \mathbf{m}_N^{-1} S_N^{-1} \mathbf{m}_N + \sum_{n=1}^N t_n^2 \right) \end{aligned}$$



Model Evidence

- The model evidence for our example is given as:

$$\begin{aligned} p(\mathbf{t}) &= \iint p(\mathbf{t} / \mathbf{w}, \beta) p(\mathbf{w} / \beta) d\mathbf{w} p(\beta) d\beta \\ &= \iint \left(\frac{\beta}{2\pi} \right)^{N/2} \exp \left\{ -\frac{\beta}{2} (\mathbf{t} - \Phi\mathbf{w})^T (\mathbf{t} - \Phi\mathbf{w}) \right\} \\ &\quad \left(\frac{\beta}{2\pi} \right)^{M/2} |\mathbf{S}_0|^{-1/2} \exp \left\{ -\frac{\beta}{2} (\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{t} - \mathbf{m}_0) \right\} d\mathbf{w} \\ &\quad \Gamma(a_0)^{-1} b_0^{a_0} \beta^{a_0-1} e^{-b_0\beta} d\beta = \\ &= \frac{b_0^{a_0}}{\left((2\pi)^{M+N} |\mathbf{S}_0| \right)^{1/2}} \iint \exp \left\{ -\frac{\beta}{2} (\mathbf{t} - \Phi\mathbf{w})^T (\mathbf{t} - \Phi\mathbf{w}) \right\} \\ &\quad \exp \left\{ -\frac{\beta}{2} (\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{t} - \mathbf{m}_0) \right\} d\mathbf{w} \\ &\quad \Gamma(a_0)^{-1} \beta^{a_0-1} \beta^{N/2} \beta^{M/2} e^{-b_0\beta} d\beta \end{aligned}$$



Model Evidence

- Using some of our earlier results in deriving the posterior:

$$\begin{aligned} p(\mathbf{t}) &= \frac{b_0^{a_0}}{\left((2\pi)^{M+N} |\mathbf{S}_0|\right)^{1/2}} \iint \exp \left\{ -\frac{\beta}{2} (\mathbf{t} - \mathbf{m}_N)^T \mathbf{S}_N^{-1} (\mathbf{t} - \mathbf{m}_N) \right\} d\mathbf{w} \\ &\quad \exp \left\{ -\frac{\beta}{2} (\mathbf{t}^T \mathbf{t} + \mathbf{m}_0^T \mathbf{S}_0^{-1} \mathbf{m}_0 - \mathbf{m}_N^T \mathbf{S}_N^{-1} \mathbf{m}_N) \right\} \\ &\quad \Gamma(a_0)^{-1} \beta^{a_N-1} \beta^{M/2} \exp(-b_0 \beta) d\beta \end{aligned}$$

- Performing the integration in \mathbf{w} and using the normalization factor for the Gamma distribution:

$$\begin{aligned} p(\mathbf{t}) &= \frac{b_0^{a_0}}{\left((2\pi)^{M+N} |\mathbf{S}_0|\right)^{1/2}} (2\pi)^{M/2} |\mathbf{S}_N|^{1/2} \underbrace{\Gamma(a_0)^{-1} \int \beta^{a_N-1} \exp(-b_N \beta) d\beta}_{\Gamma(a_N)/b_N^{a_N}} \\ &= \frac{1}{(2\pi)^{N/2}} \frac{|\mathbf{S}_N|^{1/2}}{|\mathbf{S}_0|^{1/2}} \frac{b_0^{a_0}}{b_N^{a_N}} \frac{\Gamma(a_N)}{\Gamma(a_0)} \end{aligned}$$

Limitations of Fixed Basis Functions

- Up to now the basis functions $\phi_j(\mathbf{x})$ are fixed before the training data set is observed.
- As a consequence, **the number of basis functions grows exponentially with the dimensionality D of the input space.**
- There are two properties of real data sets that we can exploit to alleviate the curse of dimensionality:
 - the data vectors $\{\mathbf{x}_n\}$ typically lie close to a nonlinear manifold whose intrinsic dimensionality is smaller than that of the input space as a result of strong correlations between the input variables.
 - If we are using localized basis functions, we can arrange that they are scattered in input space only in regions containing data (radial basis functions, support vector and relevance vector machines).



Adaptive Basis Functions

- Neural network models using adaptive basis functions having sigmoidal nonlinearities, can adapt the parameters so that the regions of input space over which the basis functions vary correspond to the data manifold.
- The target variables may have significant dependence on only a small number of possible directions within the data manifold.
- Neural networks can exploit this property by choosing the directions in input space to which the basis functions respond.

Laplace Approximation

- As we have seen earlier, the Laplace approximation allows a Gaussian approximation of the parameter posterior about the maximum a posteriori (MAP) parameter estimate.
- Consider a data set \mathcal{D} and M models \mathcal{M}_i , $i=1,\dots,M$ with corresponding parameters θ_i , $i=1,\dots,M$. We compare models using the posteriors:

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

- For large sets of data \mathcal{D} (relative to the model parameters), the parameter posterior is approximately Gaussian around the MAP estimate $\boldsymbol{\theta}_m^{MAP}$ (can also use 2nd order Taylor expansion of the log-posterior):

$$p(\boldsymbol{\theta}_m | \mathcal{D}, M_m) \simeq (2\pi)^{-d/2} |\mathbf{A}|^{1/2} \exp\left(-\frac{1}{2} (\boldsymbol{\theta}_m - \boldsymbol{\theta}_m^{MAP})^T \mathbf{A} (\boldsymbol{\theta}_m - \boldsymbol{\theta}_m^{MAP})\right),$$

$$A_{ij} = -\frac{\partial^2 \log P(\boldsymbol{\theta}_m | \mathcal{D}, M_m)}{\partial \boldsymbol{\theta}_{mi} \partial \boldsymbol{\theta}_{mj}} \Big|_{\boldsymbol{\theta}_m^{MAP}}$$



Laplace Approximation

- We can write the model evidence as

$$p(\mathcal{D} | \mathcal{M}_m) = \frac{p(\theta_m, \mathcal{D} | \mathcal{M}_m)}{p(\theta_m | \mathcal{D}, \mathcal{M}_m)}$$

- Using the Laplace approximation for the posterior of the parameters and evaluating the equation above at θ_m^{MAP} :

$$\begin{aligned}\log p(D | M_m) &\simeq \log p(\theta_m^{MAP}, \mathcal{D} | M_m) - \log p(\theta_m^{MAP} | \mathcal{D}, M_m) \\ &\simeq \log p(\mathcal{D} | \theta_m^{MAP}, M_m) + \log p(\theta_m^{MAP} | M_m) + \frac{d}{2} \log(2\pi) - \frac{1}{2} \log|A|\end{aligned}$$

- This Laplace approximation is used often for model comparison.
- Other approximations are also very useful:
 - Bayesian Information Criterion (BIC) (on the limit of $N \rightarrow \infty$)
 - MCMC (Sampling approach)
 - Variational Methods

Bayesian Information Criterion

- We start with the Laplace approximation on the limit of large data sets $N \rightarrow \infty$,

$$\log p(\mathcal{D}|M_m) \simeq \log p(\mathcal{D}|\boldsymbol{\theta}_m^{MAP}, M_m) + \log p(\boldsymbol{\theta}_m^{MAP}|M_m) + \frac{d}{2} \log(2\pi) - \frac{1}{2} \log|A|$$

- As N grows, A grows as $N\mathbf{A}_0$ for some fixed matrix \mathbf{A}_0 , thus

$$\log|A| \rightarrow \log|N\mathbf{A}_0| = \log(N^d |\mathbf{A}_0|) = d \log N + \log(|\mathbf{A}_0|) \xrightarrow{N \rightarrow \infty} d \log N$$

- Then the Laplace approximation is simplified as:

$$\log p(\mathcal{D}|M_m) \simeq \log p(\mathcal{D}|\boldsymbol{\theta}_m^{MAP}, M_m) - \frac{d}{2} \log N \quad (\text{limit } N \rightarrow \infty)$$

- Note interesting properties of (the easy to compute) BIC:

- No dependence on the prior
- One can use the MLE rather than the MAP estimate of $\boldsymbol{\theta}_m$
- If not all parameters are well determined from the data,
 $d=\text{number of effective parameters}$.



Another Example of Model Selection

- Let us consider a regression model with the following particulars ($d=k+1$ dimensional data):^a

$$y | \mathbf{w}, \sigma^2, \Phi \sim \mathcal{N}_n(\Phi\mathbf{w}, \sigma^2 \mathbf{I}_n)$$

$$\begin{aligned} \mathbf{w} | \sigma^2, \Phi &\sim \mathcal{N}_{k+1}(\mathbf{w}_0, \sigma^2 \mathbf{M}^{-1}), \mathbf{M} \text{ a } (k+1) \times (k+1) \text{ pos. def. symm. matrix} \\ \sigma^2 / \Phi &\sim \text{InvGamma}(a, b), a, b > 0 \end{aligned}$$

$$\mathbf{M} = \mathbf{I}_{k+1} / c, c > 0 \text{ and } \mathbf{w}_0 = \mathbf{0}_{k+1}$$

- Our data are in a matrix form (dimension $N \times (k + 1)$):

$$\Phi = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ 1 & x_{31} & x_{32} & \dots & x_{3k} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix}$$

We only slightly change our notation here to conform with [the MatLab program](#) implementing this example (from Zoubin Ghahramani)



Example: Likelihood Calculation

- We will derive the model evidence analytically. At first the likelihood can be written as:

$$\ell(\mathbf{w}, \sigma^2 | \mathbf{y}, \Phi) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} (\mathbf{y} - \Phi\mathbf{w})^T (\mathbf{y} - \Phi\mathbf{w})\right)$$

- With simple algebra, we can rewrite the likelihood introducing the MLE estimate of the parameters as follows:

$$\begin{aligned}\ell(\mathbf{w}, \sigma^2 | \mathbf{y}, \Phi) &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} (\mathbf{y} - \Phi\mathbf{w}_{ML})^T (\mathbf{y} - \Phi\mathbf{w}_{ML}) - \frac{1}{2\sigma^2} (\mathbf{w} - \mathbf{w}_{ML})^T \Phi^T \Phi (\mathbf{w} - \mathbf{w}_{ML})\right) = \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} s^2 - \frac{1}{2\sigma^2} (\mathbf{w} - \mathbf{w}_{ML})^T \Phi^T \Phi (\mathbf{w} - \mathbf{w}_{ML})\right)\end{aligned}$$

where

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

$$s^2 \triangleq (\mathbf{y} - \Phi\mathbf{w}_{ML})^T (\mathbf{y} - \Phi\mathbf{w}_{ML})$$

Computing Model Evidence

- Our posterior is then of the following form:

$$p(\mathbf{w}, \sigma^2 | \mathbf{w}_{ML}, s^2, \Phi) = \frac{1}{(2\pi)^{N/2}} \frac{b^a}{\Gamma(a)} \sigma^{-k-n-2a-3} \exp\left(-\frac{1}{2\sigma^2} \left\{ s^2 + 2b + \mathbf{w}^T \left(\frac{1}{c} \mathbf{I} + \Phi^T \Phi \right) \mathbf{w} - 2\mathbf{w}^T \Phi^T \Phi \mathbf{w}_{ML} + \mathbf{w}_{ML}^T (\Phi^T \Phi) \mathbf{w}_{ML} \right\}\right)$$

- The evidence is now computed as (use first the normalization of the *Inverse Gamma* distribution):

$$p(\mathbf{y} / \mathcal{M}) = \int \int \ell(\mathbf{w}, \sigma^2 | \mathbf{y}, \Phi) p(\mathbf{w}, \sigma^2 | \mathbf{w}_{ML}, s^2, \Phi) d\mathbf{w} d\sigma^2 =$$

$$\frac{1}{(2\pi)^{N/2}} \frac{b^a}{\Gamma(a)} \frac{1}{(2\pi)^{d/2}} c^{d/2} \int \int \sigma^{-k-n-2a-3} \exp\left(-\frac{1}{2\sigma^2} \underbrace{\left\{ s^2 + 2b + \mathbf{w}^T \left(\frac{1}{c} \mathbf{I} + \Phi^T \Phi \right) \mathbf{w} - 2\mathbf{w}^T \Phi^T \Phi \mathbf{w}_{ML} + \mathbf{w}_{ML}^T (\Phi^T \Phi) \mathbf{w}_{ML} \right\}}_A\right) d\sigma^2 d\mathbf{w} =$$

$$\frac{1}{(2\pi)^{N/2}} \frac{b^a}{\Gamma(a)} \frac{1}{(2\pi)^{d/2}} c^{d/2} \int \int \sigma^{-k-n-2a-3} \exp\left(-\frac{1}{2\sigma^2} A\right) d\sigma^2 d\mathbf{w} = \frac{1}{(2\pi)^{N/2}} \frac{b^a}{\Gamma(a)} \frac{1}{(2\pi)^{d/2}} c^{d/2} \int \int (\sigma^2)^{(-k-n-2a-3)/2} \exp\left(-\frac{A/2}{\sigma^2}\right) d\sigma^2 d\mathbf{w} =$$

$$\frac{1}{(2\pi)^{N/2}} \frac{b^a}{\Gamma(a)} \frac{1}{(2\pi)^{d/2}} c^{d/2} \int \int (\sigma^2)^{\frac{-(k+n+2a+1)}{\alpha}/2-1} \exp\left(-\underbrace{(A/2)}_{\beta} \frac{1}{\sigma^2}\right) d\sigma^2 d\mathbf{w} = \frac{1}{(2\pi)^{N/2}} \frac{b^a}{\Gamma(a)} \frac{1}{(2\pi)^{d/2}} c^{d/2} \int \frac{\Gamma((k+n+2a+1)/2)}{(A/2)^{(k+n+2a+1)/2}} d\mathbf{w} =$$

$$\frac{1}{(2\pi)^{N/2}} \frac{b^a}{\Gamma(a)} \frac{1}{(2\pi)^{d/2}} c^{d/2} \Gamma((k+n+2a+1)/2) 2^{(k+n+2a+1)/2} \int \left[s^2 + 2b + \mathbf{w}^T \left(\frac{1}{c} \mathbf{I} + \Phi^T \Phi \right) \mathbf{w} - 2\mathbf{w}^T \Phi^T \Phi \mathbf{w}_{ML} + \mathbf{w}_{ML}^T (\Phi^T \Phi) \mathbf{w}_{ML} \right]^{-(k+n+2a+1)/2} d\mathbf{w}$$

Useful formulas for the *Inverse Gamma*: $p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{-(\alpha+1)} e^{-\beta/\theta}$, $\theta > 0$

Computing Model Evidence

- The evidence can be further simplified (we will use now the normalization of the multivariate *Student-t* distribution)

$$p(\mathbf{y} / \mathcal{M}) = \frac{1}{(2\pi)^{N/2}} \frac{b^a}{\Gamma(a)} \frac{1}{(2\pi)^{d/2}} c^{d/2} \Gamma\left(\frac{d+n}{2} + a\right) 2^{\frac{d+n}{2} + a}$$

$$\int \underbrace{s^2 + 2b + \mathbf{w}_{ML}^T (\Phi^T \Phi) \mathbf{w}_{ML} - \boldsymbol{\mu}^T \Sigma^{-1} \boldsymbol{\mu}}_g \left[\underbrace{\left(\frac{1}{c} \mathbf{I} + \Phi^T \Phi \right)^{-1}}_{\Sigma} \right]^{-1} (\mathbf{w} - \boldsymbol{\mu})^{-(d+n)/2-a} d\mathbf{w}$$

$$\begin{aligned} & \frac{1}{(2\pi)^{N/2}} \frac{b^a}{\Gamma(a)} \frac{1}{(2\pi)^{d/2}} c^{d/2} \Gamma\left(\frac{d+n}{2} + a\right) 2^{\frac{d+n}{2} + a} g^{-(d+n)/2-a} \int \left[1 + \frac{1}{g} \{ (\mathbf{w} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{w} - \boldsymbol{\mu}) \} \right]^{-(d+n+2a)/2} d\mathbf{w} \\ & \frac{1}{(2\pi)^{N/2}} \frac{b^a}{\Gamma(a)} \frac{1}{(2\pi)^{d/2}} c^{d/2} \Gamma\left(\frac{d+n}{2} + a\right) 2^{\frac{d+n}{2} + a} g^{-(d+n)/2-a} \int \left[1 + \frac{1}{v} \left\{ (\mathbf{w} - \boldsymbol{\mu})^T \left(\frac{g}{n+2a} \Sigma \right)^{-1} (\mathbf{w} - \boldsymbol{\mu}) \right\} \right]^{-(d+v)/2} d\mathbf{w} \end{aligned}$$

Useful formulas for the *Student-t*

$$p(\theta) = \frac{\Gamma((\nu+d)/2)}{\Gamma(\nu/2)\nu^{d/2}\pi^{d/2}} |\Sigma|^{-1/2} \times (1 + \frac{1}{\nu}(\theta - \boldsymbol{\mu})^T \Sigma^{-1} (\theta - \boldsymbol{\mu}))^{-(\nu+d)/2}$$



Computing Model Evidence

- Performing the integration in β in the last slide:

$$\begin{aligned}
 p(\mathbf{y} / \mathcal{M}) &= \frac{1}{(2\pi)^{N/2}} \frac{b^a}{\Gamma(a)} \frac{1}{(2\pi)^{d/2}} c^{d/2} 2^{\frac{d+n}{2}+a} g^{-\frac{(d+n)/2-a}{2}} \Gamma\left(\frac{n}{2}+a\right) (n+2a)^{d/2} \pi^{d/2} \left(\frac{g}{n+2a}\right)^{d/2} \left|\left(\frac{1}{c} \mathbf{I} + \Phi^T \Phi\right)\right|^{-1/2} = \\
 &\frac{1}{(2\pi)^{N/2}} \frac{b^a}{\Gamma(a)} c^{d/2} \Gamma\left(\frac{n}{2}+a\right) \left(\frac{1}{2} s^2 + b + \frac{1}{2} \mathbf{w}_{ML}^T (\Phi^T \Phi) \mathbf{w}_{ML} - \frac{1}{2} \boldsymbol{\mu}^T \Sigma^{-1} \boldsymbol{\mu}\right)^{-n/2-a} \left|\left(\frac{1}{c} \mathbf{I} + \Phi^T \Phi\right)\right|^{-1/2} = \\
 &\frac{1}{(2\pi)^{N/2}} \frac{b^a}{\Gamma(a)} c^{d/2} \Gamma\left(\frac{n}{2}+a\right) \left(\underbrace{\frac{1}{2} (\mathbf{y} - \Phi \mathbf{w}_{ML})^T (\mathbf{y} - \Phi \mathbf{w}_{ML}) + b + \frac{1}{2} \mathbf{w}_{ML}^T (\Phi^T \Phi) \mathbf{w}_{ML} - \frac{1}{2} \boldsymbol{\mu}^T \Sigma^{-1} \boldsymbol{\mu}}_{\mathcal{E}} \right)^{-n/2-a} \left|\left(\frac{1}{c} \mathbf{I} + \Phi^T \Phi\right)\right|^{-1/2}
 \end{aligned}$$

- Using some of the earlier definitions,

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \quad s^2 \triangleq (\mathbf{y} - \Phi \mathbf{w}_{ML})^T (\mathbf{y} - \Phi \mathbf{w}_{ML}) \quad \boldsymbol{\mu} = \left(\frac{1}{c} \mathbf{I} + \Phi^T \Phi\right)^{-1} (\Phi^T \Phi) \mathbf{w}_{ML}$$

we can simplify \mathcal{E} as:

$$\begin{aligned}
 \mathcal{E} &= \frac{1}{2} (\mathbf{y} - \Phi \mathbf{w}_{ML})^T (\mathbf{y} - \Phi \mathbf{w}_{ML}) + b + \frac{1}{2} \mathbf{w}_{ML}^T (\Phi^T \Phi) \mathbf{w}_{ML} - \frac{1}{2} \boldsymbol{\mu}^T \Sigma^{-1} \boldsymbol{\mu} = b + \frac{1}{2} \mathbf{y}^T \mathbf{y} + \frac{1}{2} \mathbf{w}_{ML}^T (\Phi^T \Phi) \mathbf{w}_{ML} - \mathbf{y}^T \Phi \mathbf{w}_{ML} + \frac{1}{2} \mathbf{w}_{ML}^T (\Phi^T \Phi) \mathbf{w}_{ML} \\
 &- \frac{1}{2} \mathbf{w}_{ML}^T (\Phi^T \Phi) \left(\frac{1}{c} \mathbf{I} + \Phi^T \Phi\right)^{-1} \left(\frac{1}{c} \mathbf{I} + \Phi^T \Phi\right) \left(\frac{1}{c} \mathbf{I} + \Phi^T \Phi\right)^{-1} (\Phi^T \Phi) \mathbf{w}_{ML} = b + \frac{1}{2} \mathbf{y}^T \mathbf{y} + \mathbf{y}^T \Phi (\Phi^T \Phi)^{-1} (\Phi^T \Phi) (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \\
 &- \mathbf{y}^T \Phi (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} - \frac{1}{2} \mathbf{w}_{ML}^T (\Phi^T \Phi) \left(\frac{1}{c} \mathbf{I} + \Phi^T \Phi\right)^{-1} (\Phi^T \Phi) (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} = b + \frac{1}{2} \mathbf{y}^T \mathbf{y} - \frac{1}{2} \mathbf{y}^T \Phi \left(\frac{1}{c} \mathbf{I} + \Phi^T \Phi\right)^{-1} \Phi^T \mathbf{y}
 \end{aligned}$$

Model Evidence

- The final evidence in analytical form is given as:

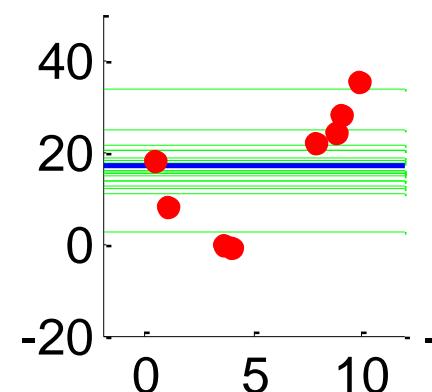
$$p(\mathbf{y} / \mathcal{M}) = \frac{1}{(2\pi)^{N/2}} \frac{b^a}{\Gamma(a)} c^{d/2} \Gamma\left(\frac{n}{2} + a\right) \left(b + \frac{1}{2} \mathbf{y}^T \mathbf{y} - \frac{1}{2} \mathbf{y}^T \Phi \left(\frac{1}{c} \mathbf{I} + \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{y} \right)^{-n/2-a} \left| \left(\frac{1}{c} \mathbf{I} + \Phi^T \Phi \right) \right|^{-1/2}$$

- Compare this with what is given in this [MatLab Implementation](#).
- The model evidence and samples of different order (M) regression models are given below. The specific data of the problem can be found in the MatLab file.
- We are looking for the order of the polynomial that maximizes the evidence. Note that the MatLab implementation utilized random input/output data.

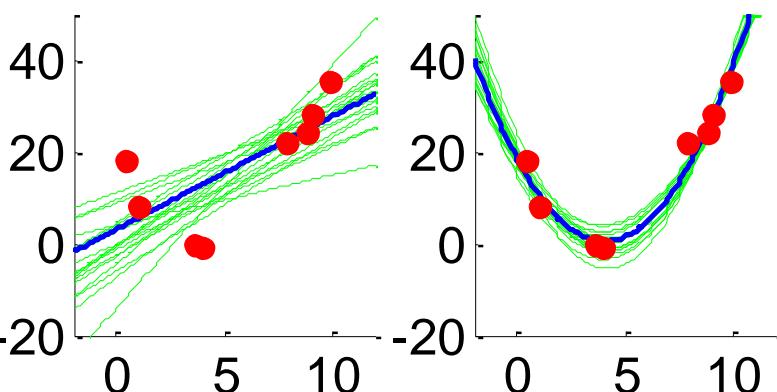


Bayesian Model Comparison

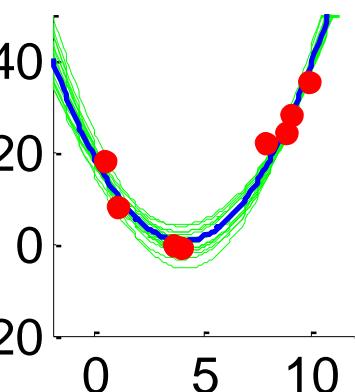
$M = 0$



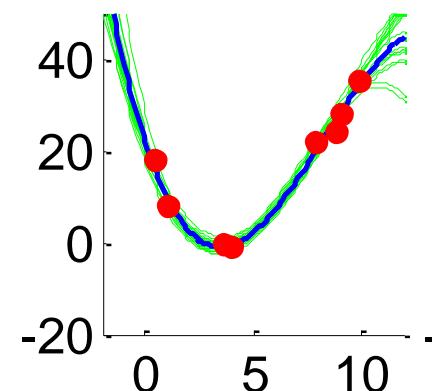
$M = 1$



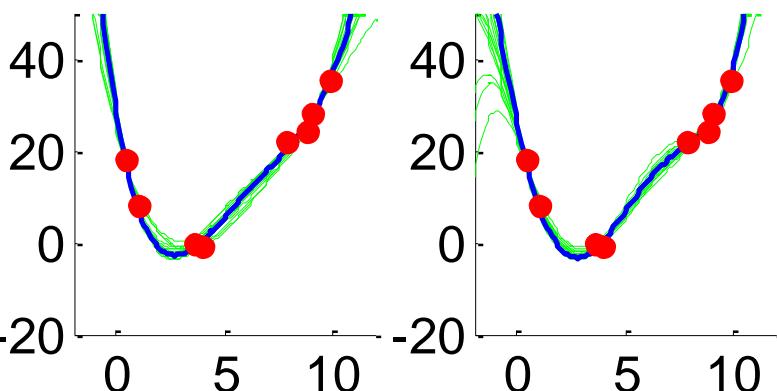
$M = 2$



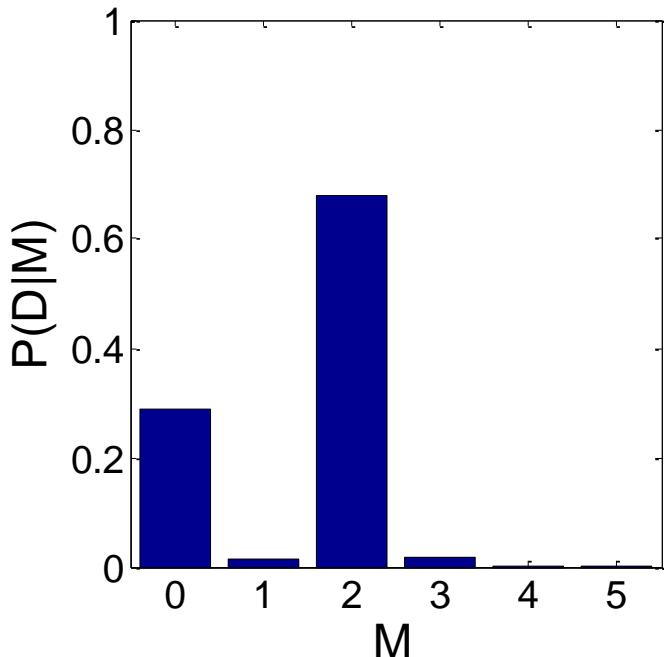
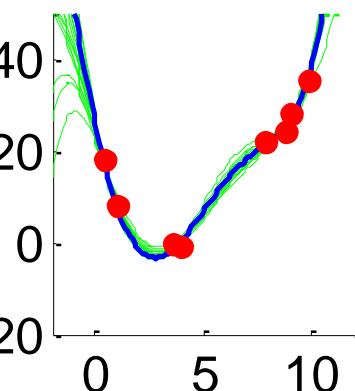
$M = 3$



$M = 4$



$M = 5$



[MatLab implementation](#) of Bayesian Model Selection (from [Zoubin Ghahramani](#))

Statistical Computing, University of Notre Dame, Notre Dame, IN, USA (Fall 2017, N. Zabaras)



Linear Models Of Regression: Equivalent Kernel Representation



Equivalent Kernel

- The predictive mean in our regression model can be written as:

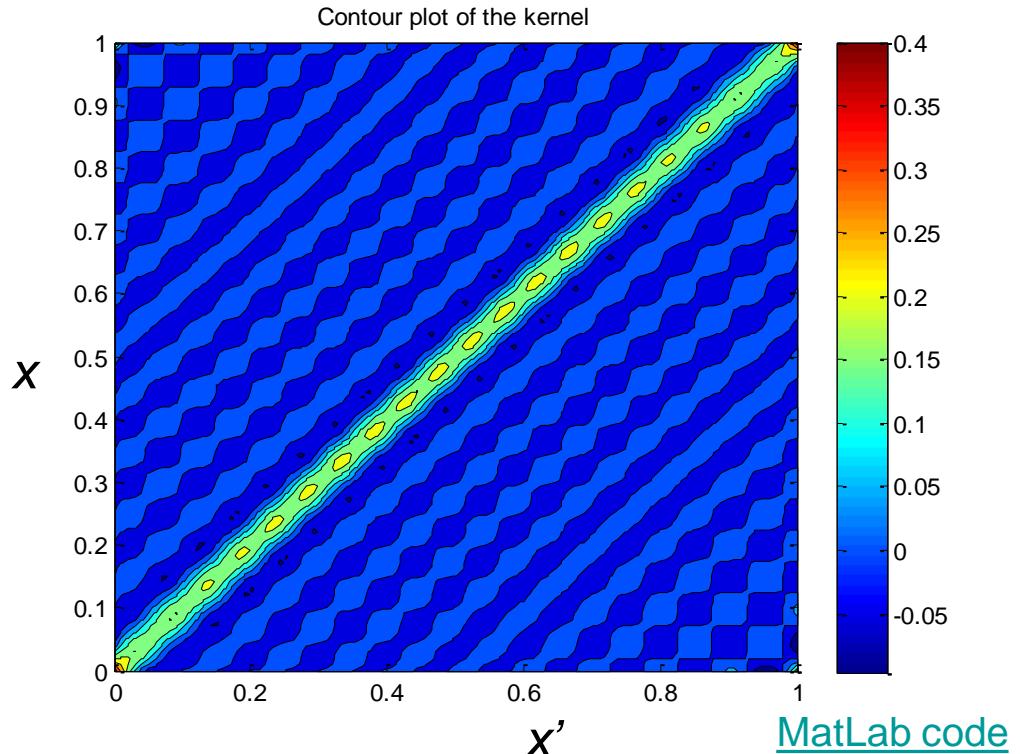
$$y(\mathbf{x}, \mathbf{m}_N) = \boldsymbol{\phi}(\mathbf{x})^T \mathbf{m}_N = \beta \boldsymbol{\phi}(\mathbf{x})^T S_N \Phi^T \mathbf{t} = \sum_{n=1}^N \underbrace{\beta \boldsymbol{\phi}(\mathbf{x})^T S_N \boldsymbol{\phi}(\mathbf{x}_n)}_{\text{Equivalent kernel } k(\mathbf{x}, \mathbf{x}_n)} t_n \Rightarrow$$
$$y(\mathbf{x}, \mathbf{m}_N) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n$$
$$k(\mathbf{x}, \mathbf{x}_n) = \beta \boldsymbol{\phi}(\mathbf{x})^T S_N \boldsymbol{\phi}(\mathbf{x}_n)$$

The kernel is shown here as a plot of \mathbf{x} vs. \mathbf{x}' .

20 samples \mathbf{x}_n equally spaced in $[0, 1]$

Gaussian kernels
($s=0.05$), $\alpha=5 \times 10^{-3}$,
 $\beta=11.1$.

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



Equivalent Kernel

- The predictive mean can be written as follows:

$$y(\mathbf{x}, \mathbf{m}_N) = \boldsymbol{\phi}(\mathbf{x})^T \mathbf{m}_N = \beta \boldsymbol{\phi}(\mathbf{x})^T S_N \Phi^T \mathbf{t} = \sum_{n=1}^N \underbrace{\beta \boldsymbol{\phi}(\mathbf{x})^T S_N \boldsymbol{\phi}(\mathbf{x}_n)}_{\text{Equivalent kernel } k(\mathbf{x}, \mathbf{x}_n)} t_n \Rightarrow$$

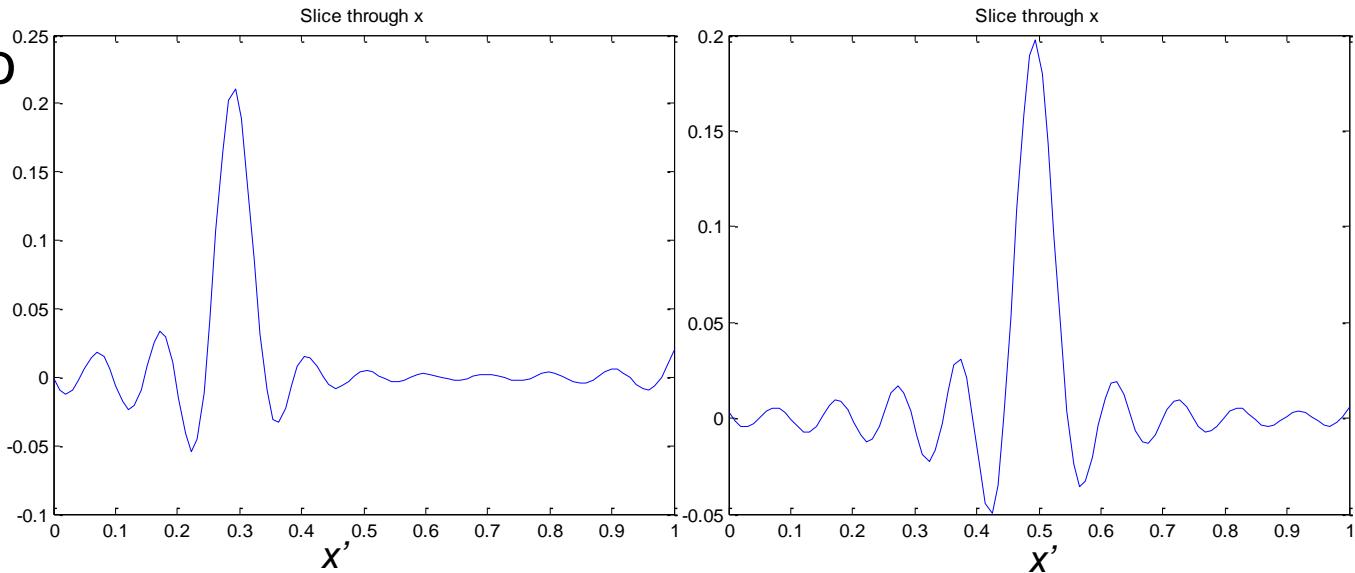
$$y(\mathbf{x}, \mathbf{m}_N) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n$$

$$k(\mathbf{x}, \mathbf{x}_n) = \beta \boldsymbol{\phi}(\mathbf{x})^T S_N \boldsymbol{\phi}(\mathbf{x}_n)$$

Two slices
corresponding to
 $\mathbf{x}=0.3, \mathbf{x}=0.5.$

20 Gaussians
($s=0.05$)

Plots for 100
points \mathbf{x}
uniformly
in $(0,1)$



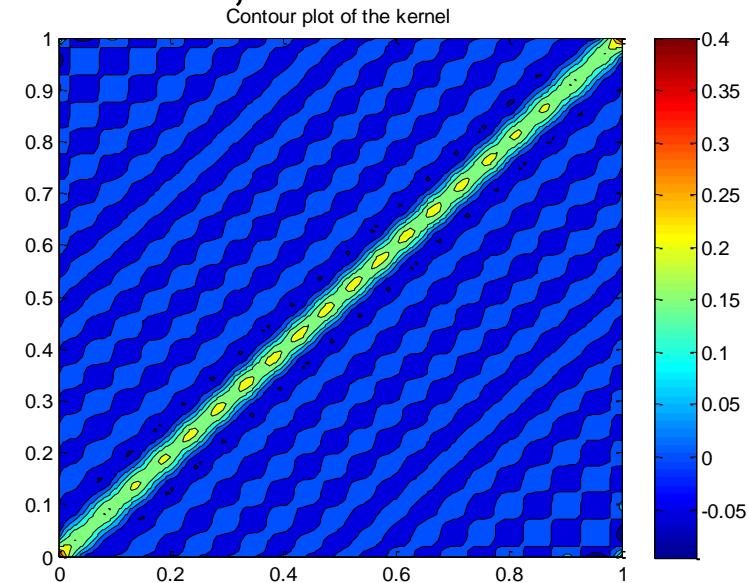
[MatLab code](#)

Equivalent Kernel

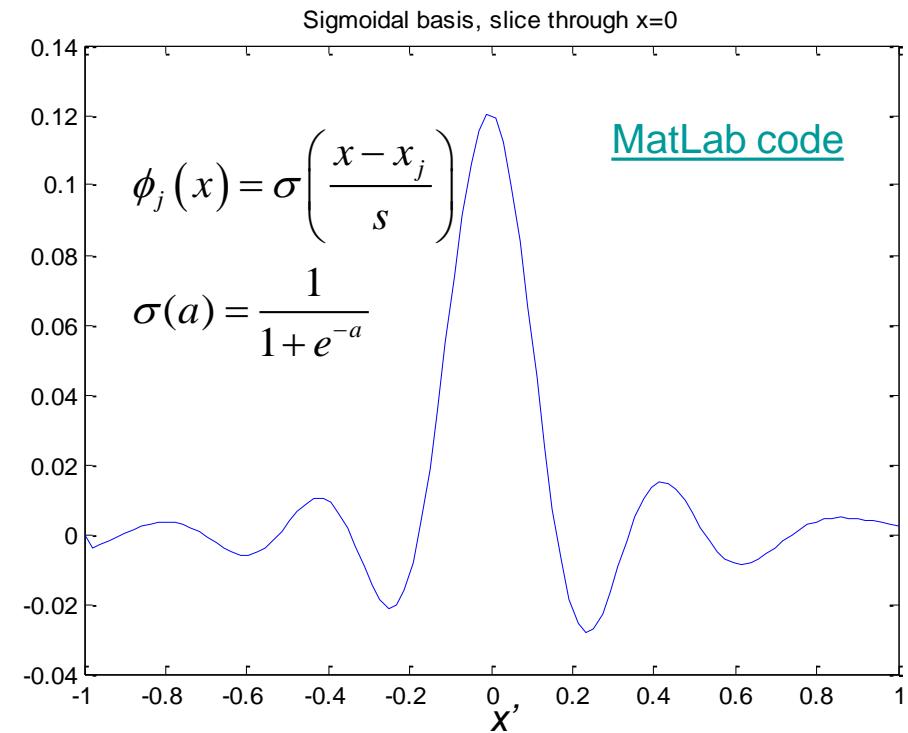
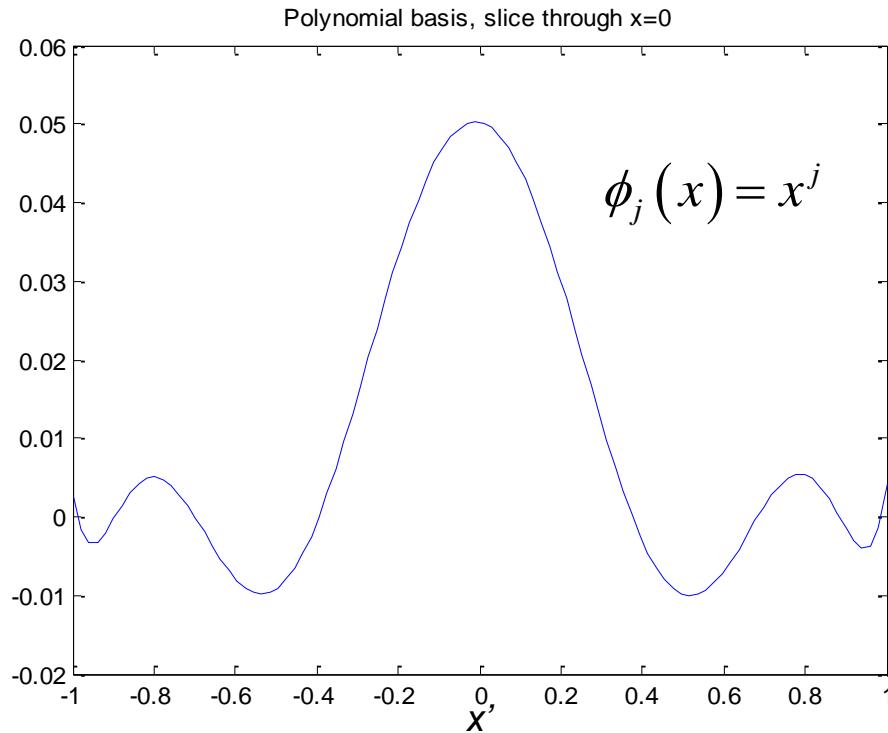
$$y(\mathbf{x}, \mathbf{m}_N) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n$$

$$k(\mathbf{x}, \mathbf{x}_n) = \beta \phi(\mathbf{x})^T S_N \phi(\mathbf{x}_n)$$

- Note that we make predictions by taking linear combinations of the training set target values (*linear smoothers*).
- The equivalent kernel depends on the input values \mathbf{x}_n since these appear in S_N .
- See from the figure that its localized around \mathbf{x} , and so the mean of the predictive distribution $y(\mathbf{x}, \mathbf{m}_N)$, is obtained by forming a weighted combination of t_n in which data points close to \mathbf{x} are given higher weight than points (evidence) further removed from \mathbf{x} .



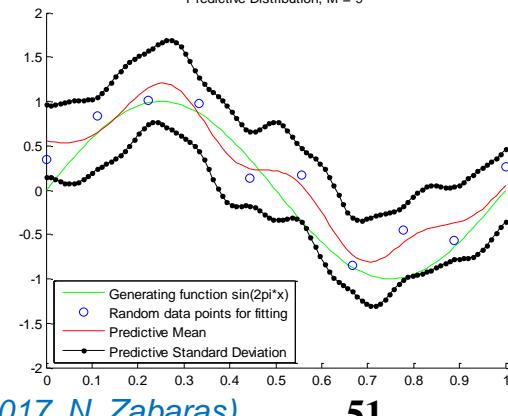
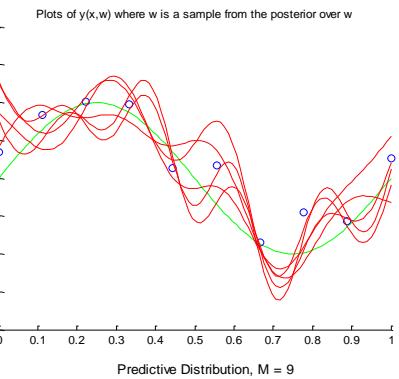
Equivalent Kernel



- Examples of equivalent kernels $k(x, x')$ for $x = 0$ plotted as a function of x' . 10 basis functions were used in each case. The sigmoidal basis is centered at 10 equally spaced x_n in $[-1,1]$.
- These are localized functions of x' even though the corresponding basis functions are nonlocal.

Equivalent Kernel

- Consider the covariance between $y(\mathbf{x})$ and $y(\mathbf{x}')$, which is given by (recall that $p(\mathbf{w} | \mathbf{x}, t, \alpha, \beta) = \mathcal{N}\left(\mathbf{w}, \beta S_N \sum_{n=1}^N t_n \phi(x_n), S_N\right)$):
$$\text{cov}[y(\mathbf{x}), y(\mathbf{x}')] = \text{cov}[\phi(\mathbf{x})^T \mathbf{w}, \mathbf{w}^T \phi(\mathbf{x}')] = \phi(\mathbf{x})^T S_N \phi(\mathbf{x}') = \beta^{-1} k(\mathbf{x}, \mathbf{x}')$$
- From the earlier discussion on the equivalent kernel, we see that the predictive mean at nearby points will be highly correlated, whereas for more distant points the correlation will be smaller.
- By drawing samples from the posterior distribution over \mathbf{w} , and plotting the corresponding model functions $y(\mathbf{x}, \mathbf{w})$, we are visualizing the joint uncertainty in the posterior distribution between the y values at two (or more) x values, as governed by the equivalent kernel.
- This is contrary to the Figure here where we visualized pointwise uncertainty in the predictions.



Equivalent Kernel

- We can avoid the use of basis functions and define the kernel function directly, leading e.g. to *Gaussian Processes* (more on Kernel methods and GPs later on in this course).
- Note that

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

for all values of \mathbf{x} . However, the equivalent kernel

$$k(\mathbf{x}, \mathbf{x}_n) = \beta \phi(\mathbf{x})^T S_N \phi(\mathbf{x}_n)$$

may be negative for some values of \mathbf{x} .

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

$$k(\mathbf{x}, \mathbf{z}) = \beta \phi(\mathbf{x})^T S_N \phi(\mathbf{z}) \Rightarrow$$

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

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

- This property can be used to “design” interesting kernels.





Introduction to Variable Selection



Variable Selection

- Let us return to our regression model with one dependent random variable y and a set of k $\{x_1, x_2, \dots, x_k\}$ explanatory variables.
- Are all the x_i 's needed in the regression?
- We assume that every q -subset $\{i_1, i_2, \dots, i_q\}$, $0 \leq q \leq k$, of the explanatory variables,
$$\{1, x_{i_1}, x_{i_2}, \dots, x_{i_q}\}$$
is a proper set of explanatory variables for the regression of y (as before, the intercept is included in all models).
- We have a total of 2^k models to select from!

Variable Selection

- Following earlier notation, we denote: $X = [\mathbf{1}_n \ x_1 \ x_2 \dots x_k]$ as the matrix that contains $\mathbf{1}_n$ and the k potential predictor variables.
- Each model M_γ is associated with binary indicator vector

$$\gamma \in \Gamma = \{0,1\}^k$$

where $\gamma_i=1$ means that the variable x_i is included in the model M_γ and $\gamma_i=0$ that it is not.

- The number of variables included in the model M_γ is:
- The indices of the variables included in the model and not included in the model are denoted, respectively, as: $t_1(\gamma), t_0(\gamma)$



Variable Selection - Models in Competition

- For $\beta \in \mathbb{R}^{k+1}$ and X , we define β_γ as the sub-vector

$$\beta_\gamma = \left(\beta_0, (\beta_i)_{i \in t_1(\gamma)} \right)$$

- Let X_γ be the submatrix of X where only the column 1_n and the columns in $t_1(\gamma)$ have been left.
- The model M_γ is then defined as:

$$y | \gamma, \beta_\gamma, \sigma^2, X \sim \mathcal{N}_n(X_\gamma \beta_\gamma, \sigma^2 I_n)$$

where $\beta_\gamma \in \mathbb{R}^{q_\gamma+1}$, $\sigma^2 \in \mathbb{R}_+^*$ are the unknown parameters.

- The σ^2 is common to all models and we use the same prior for all models.



Variable Selection - Models in Competition

- We have a high number 2^k of models in competition.
- We cannot specify a prior on every M_γ in a completely subjective and autonomous manner.
- We derive all priors from a single global prior **associated with the full model** that corresponds to $\gamma = (1, \dots, 1)$.



Zellner's Informative Prior: Variable Selection

- For the full model that corresponds to $\gamma = (1, \dots, 1)$, we use the Zellner's informative G-prior:

$$\begin{aligned}\boldsymbol{\beta} | \sigma^2, \mathbf{X} &\sim \mathcal{N}_{k+1}(\tilde{\boldsymbol{\beta}}, c\sigma^2(\mathbf{X}^T \mathbf{X})^{-1}) \\ \sigma^2 &\sim \pi(\sigma^2 | \mathbf{X}) \propto \sigma^{-2} \text{ improper Jeffreys prior}\end{aligned}$$

- For each model M_γ , the prior distribution of β_γ conditional on σ^2 is fixed as:

$$\boldsymbol{\beta}_\gamma | \gamma, \sigma^2 \sim \mathcal{N}_{q_\gamma+1}(\tilde{\boldsymbol{\beta}}_\gamma, c\sigma^2(\mathbf{X}_\gamma^T \mathbf{X}_\gamma)^{-1})$$

where $\tilde{\boldsymbol{\beta}}_\gamma = (\mathbf{X}_\gamma^T \mathbf{X}_\gamma)^{-1} \mathbf{X}_\gamma^T \tilde{\boldsymbol{\beta}}$ and same prior on σ^2 .



Zellner's Informative Prior: Variable Selection - Prior

- The joint prior for model M_γ is the improper prior

$$\pi(\boldsymbol{\beta}_\gamma, \sigma^2 | \mathbf{y}) \propto (\sigma^2)^{-(q_\gamma+1)/2-1} \exp \left[-\frac{1}{2(c\sigma^2)} (\boldsymbol{\beta}_\gamma - \tilde{\boldsymbol{\beta}}_\gamma)^T (\mathbf{X}_\gamma^T \mathbf{X}_\gamma) (\boldsymbol{\beta}_\gamma - \tilde{\boldsymbol{\beta}}_\gamma) \right]$$

- Infinitely many ways of defining a prior on the model index γ :

Our choice is a uniform prior $p(\gamma|\mathbf{X}) = 2^{-k}$.

- Posterior distribution of γ is central to variable selection since it is proportional to marginal density of \mathbf{y} on M_γ (or evidence of M_γ)

Zellner's Informative Prior: Variable Selection - Prior

- Posterior distribution of γ is proportional to the marginal density of γ on M_γ (so it can also be used to compute Bayes factors)

$$\begin{aligned}\pi(\gamma | \mathbf{y}, \mathbf{X}) &\propto f(\mathbf{y} | \gamma, \mathbf{X}) \pi(\gamma | \mathbf{X}) \propto f(\mathbf{y} | \gamma, \mathbf{X}) \\ &= \int \left(\int f(\mathbf{y} | \gamma, \boldsymbol{\beta}, \sigma^2, \mathbf{X}) \pi(\boldsymbol{\beta} | \gamma, \sigma^2, \mathbf{X}) d\boldsymbol{\beta} \right) \pi(\sigma^2 | \mathbf{X}) d\sigma^2\end{aligned}$$

where it can be shown

$$\begin{aligned}f(\mathbf{y} | \boldsymbol{\gamma}, \sigma^2, \mathbf{X}) &= \int f(\mathbf{y} | \boldsymbol{\gamma}, \boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta} | \boldsymbol{\gamma}, \sigma^2) d\boldsymbol{\beta} = \\ &= (c + 1)^{-(q_\gamma + 1)/2} (2\pi)^{-n/2} (\sigma^2)^{-n/2} \times \\ \exp \left(-\frac{1}{2\sigma^2} \mathbf{y}^T \mathbf{y} + \frac{1}{2\sigma^2(c + 1)} \left\{ c \mathbf{y}^T \mathbf{X}_\gamma (\mathbf{X}_\gamma^T \mathbf{X}_\gamma)^{-1} \mathbf{X}_\gamma^T \mathbf{y} - \tilde{\boldsymbol{\beta}}_\gamma^T \mathbf{X}_\gamma^T \mathbf{X}_\gamma \tilde{\boldsymbol{\beta}}_\gamma + 2 \mathbf{y}^T \mathbf{X}_\gamma \tilde{\boldsymbol{\beta}}_\gamma \right\} \right)\end{aligned}$$

Zellners Informative Prior: Variable Selection - Prior

$$\pi(\gamma | \mathbf{y}, \mathbf{X}) \propto \int f(\mathbf{y} | \gamma, \sigma^2, \mathbf{X}) \pi(\sigma^2 | \mathbf{X}) d\sigma^2 = \int f(\mathbf{y} | \gamma, \sigma^2, \mathbf{X}) \frac{1}{\sigma^2} d\sigma^2$$

□ Posterior distribution of γ is then given as:

$$f(\gamma | \mathbf{y}, \mathbf{X}) \propto (c + 1)^{-(q_\gamma + 1)/2} \times \\ \left(\mathbf{y}^T \mathbf{y} - \frac{c}{(c + 1)} \mathbf{y}^T \mathbf{X}_\gamma (\mathbf{X}_\gamma^T \mathbf{X}_\gamma)^{-1} \mathbf{X}_\gamma^T \mathbf{y} + \frac{1}{c + 1} \tilde{\boldsymbol{\beta}}_\gamma^T \mathbf{X}_\gamma^T \mathbf{X}_\gamma \tilde{\boldsymbol{\beta}}_\gamma - \frac{2}{c + 1} \mathbf{y}^T \mathbf{X}_\gamma \tilde{\boldsymbol{\beta}}_\gamma \right)^{-n/2}$$

□ This is of similar form to a familiar result for a fixed c :

$$f(\mathbf{y} | \mathbf{X}, c) = (c + 1)^{-(k+1)/2} \pi^{-n/2} \Gamma\left(\frac{n}{2}\right) \left[\mathbf{y}^T \mathbf{y} - \frac{c}{c + 1} \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} + \frac{1}{c + 1} \tilde{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} \tilde{\boldsymbol{\beta}} - \frac{2}{c + 1} \mathbf{y}^T \mathbf{X} \tilde{\boldsymbol{\beta}} \right]^{-n/2}$$

Model Selection

- Most likely models ordered by decreasing posterior probabilities using Zellner's informative G-prior with $c=100$.

$t_1(\gamma)$	$\pi(\gamma y, X)$
t1_gamma	$\text{pi}(\text{gamma} y, X)$
0 1 2 4 5	0.231543
0 1 2 4 5 9	0.037358
0 1 9	0.034435
0 1 2 4 5 10	0.032975
0 1 4 5	0.030606
0 1 2 9	0.025016
0 1 2 4 5 7	0.024144
0 1 2 4 5 8	0.023784
0 1 2 4 5 6	0.023735
0 1 2 3 4 5	0.023207
0 1 6 9	0.014587
0 1 2 3 9	0.014491
0 9	0.014281
0 1 2 6 9	0.013551
0 1 4 5 9	0.012761
0 1 3 9	0.011712
0 1 2 8	0.011477
0 1 8	0.009519
0 1 2 3 4 5 9	0.009036
0 1 2 4 5 6 9	0.009031

The data set corresponds to the caterpillar regression problem (details in [Bayesian Core](#)) with 10 explanatory variables.



Model Selection

- Model M_γ with the highest posterior probability is $t_1(\gamma) = (1, 2, 4, 5)$, which corresponds to the variables
 - altitude,
 - slope,
 - height of the tree sampled in the center of the area, and
 - diameter of the tree sampled in the center of the area.

Model Selection

- For the Zellner's non-informative prior with $\pi(c)=1/c$, we have ($\tilde{\beta} = \mathbf{0}_{k+1}$) :

$$\pi(\gamma / \mathbf{y}, \mathbf{X}) = \sum_{c=1}^{\infty} c^{-1} (c+1)^{-(q_\gamma+1)/2} \left[\mathbf{y}^T \mathbf{y} - \frac{c}{c+1} \mathbf{y}^T \mathbf{X}_\gamma (\mathbf{X}_\gamma^T \mathbf{X}_\gamma)^{-1} \mathbf{X}_\gamma^T \mathbf{y} \right]^{-n/2}$$

- Again this is of the same form as (see [Bayesian Core](#))

$$f(\mathbf{y} / \mathbf{X}) = \sum_{c=1}^{\infty} f(\mathbf{y} / \mathbf{X}, c) c^{-1} \propto \\ \sum_{c=1}^{\infty} c^{-1} (c+1)^{-(k+1)/2} \left[\mathbf{y}^T \mathbf{y} - \frac{c}{c+1} \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \right]^{-n/2}$$

Model Selection

- Most likely models ordered by decreasing posterior probabilities using Zellner's non-informative G-prior.

$t_1(\gamma)$	$\pi(\gamma y, X)$
t1_gamma	$\pi(\text{gamma} y, X)$
0 1 2 4 5	0.092914
0 1 2 4 5 9	0.032553
0 1 2 4 5 10	0.029512
0 1 2 4 5 7	0.023114
0 1 2 4 5 8	0.022843
0 1 2 4 5 6	0.022807
0 1 2 3 4 5	0.022409
0 1 2 3 4 5 9	0.016733
0 1 2 4 5 6 9	0.016725
0 1 2 4 5 8 9	0.013726
0 1 4 5	0.011031
0 1 2 4 5 9 10	0.009933
0 1 2 3 9	0.009698
0 1 2 9	0.009316
0 1 2 4 5 7 9	0.009253
0 1 2 6 9	0.009189
0 1 4 5 9	0.008756
0 1 2 3 4 5 10	0.007933
0 1 2 4 5 8 10	0.007901
0 1 2 4 5 7 10	0.007896

Stochastic Search for the Most Likely Model

- When k is large, it becomes computationally intractable to compute the posterior probabilities of the 2^k models.
- Need of a tailored algorithm that samples from $\pi(\gamma|y, X)$ and selects the most likely models.
- Can be done by Gibbs sampling*, given the availability of the **full conditional posterior probabilities of the γ_i 's**. If

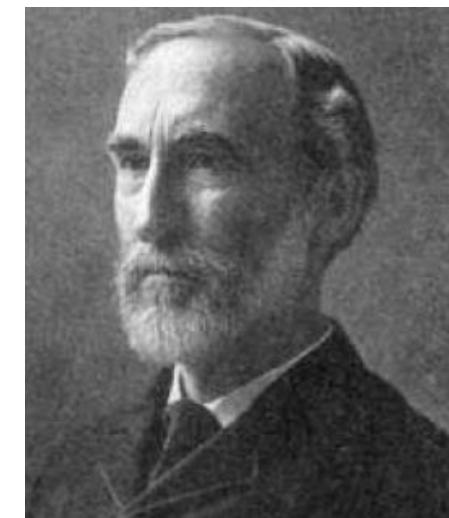
$$\gamma_{-i} = (\gamma_1, \dots, \gamma_{i-1}, \gamma_{i+1}, \dots, \gamma_k) \quad (1 \leq i \leq k)$$

then

$$\pi(\gamma_i | y, \gamma_{-i}, X) \propto \pi(\gamma | y, X)$$

(to be evaluated in both $\gamma_i = 0$ and $\gamma_i = 1$)

* Gibbs and other sampling algorithms will be introduced and discussed in detail in



forthcoming lectures.

Gibbs Sampling for Variable Selection

Initialization: Draw γ^0 from the uniform distribution on Γ

Iteration t: Given $(\gamma_1^{(t-1)}, \dots, \gamma_k^{(t-1)})$, generate

- $\gamma_1^{(t)}$ according to $\pi(\gamma_1 | \mathbf{y}, \gamma_2^{(t-1)}, \dots, \gamma_k^{(t-1)}, \mathbf{X})$
- $\gamma_2^{(t)}$ according to $\pi(\gamma_2 | \mathbf{y}, \gamma_1^{(t)}, \gamma_3^{(t-1)}, \dots, \gamma_k^{(t-1)}, \mathbf{X})$
- ..
- ..
- $\gamma_k^{(t)}$ according to $\pi(\gamma_k | \mathbf{y}, \gamma_1^{(t)}, \gamma_2^{(t)}, \dots, \gamma_{k-1}^{(t)}, \mathbf{X})$



Gibbs Sampling for Variable Selection

Question: How to sample γ_1^t according to $\pi(\gamma_1 | \mathbf{y}, \gamma_2^{(t-1)}, \dots, \gamma_k^{(t-1)}, \mathbf{X})$

1. The conditional distribution $\pi(\gamma_1 | \mathbf{y}, \gamma_2^{(t-1)}, \dots, \gamma_k^{(t-1)}, \mathbf{X})$ is proportional to $\pi(\gamma | \mathbf{y}, \mathbf{X})$
2. Since γ_1^t only has two possible values which are 0 and 1, we get

$$p_0 \propto \pi(\gamma_1 = 0, \gamma_2^{(t-1)}, \dots, \gamma_k^{(t-1)} | \mathbf{y}, \mathbf{X})$$

$$p_1 \propto \pi(\gamma_1 = 1, \gamma_2^{(t-1)}, \dots, \gamma_k^{(t-1)} | \mathbf{y}, \mathbf{X})$$

So the probability that $\gamma_1^t = 1$ is p_1 , then we can use Gibbs Sampling to approximate the distribution of $\{\gamma_i^t\}$

Gibbs Sampling: Posterior Probabilities

- After $T \gg 1$ MCMC iterations, we approximate the posterior probabilities $p(\gamma|y, X)$ by empirical averages

$$\hat{\pi}(\gamma|y, X) = \frac{1}{T - T_0 + 1} \sum_{t=T_0}^T \mathbb{I}_{\gamma^{(t)}=\gamma}$$

The T_0 first values (burn in) in the MCMC chain are eliminated.



Model Choice Comparison: Gibbs Estimates

First level Informative G-prior model with ($\tilde{\beta} = 0_{11}$, c=100) compared with the Gibbs estimates of the top ten posterior probabilities

$t_1(\gamma)$	$\pi(\gamma y, X)$	$\hat{\pi}(\gamma y, X)$
t1_gamma	pi(gamma y, X)	pi_hat(gamma y, X)
0 1 2 4 5	0.231543	0.239276
0 1 2 4 5 9	0.037358	0.034397
0 1 9	0.034435	0.032397
0 1 2 4 5 10	0.032975	0.030097
0 1 4 5	0.030606	0.029397
0 1 2 9	0.025016	0.025297
0 1 2 4 5 7	0.024144	0.022498
0 1 2 4 5 8	0.023784	0.024898
0 1 2 4 5 6	0.023735	0.023598
0 1 2 3 4 5	0.023207	0.022998



Model Choice Comparison: Gibbs Estimates

Non-informative G-prior variable model choice compared with the Gibbs estimates of the top ten posterior probabilities

$t_1(\gamma)$	$\pi(\gamma y, X)$	$\hat{\pi}(\gamma y, X)$
t1_gamma	$\pi(\text{gamma} y, X)$	$\hat{\pi}(\text{gamma} y, X)$
0 1 2 4 5	0.092914	0.093391
0 1 2 4 5 9	0.032553	0.033097
0 1 2 4 5 10	0.029512	0.032597
0 1 2 4 5 7	0.023114	0.025097
0 1 2 4 5 8	0.022843	0.023098
0 1 2 4 5 6	0.022807	0.022498
0 1 2 3 4 5	0.022409	0.021698
0 1 2 3 4 5 9	0.016733	0.015998
0 1 2 4 5 6 9	0.016725	0.014899
0 1 2 4 5 8 9	0.013726	0.013399



Gibbs Sampling: Probabilities of Inclusion

- An approximation of the probability to include the i-th variable:

$$\hat{P}^\pi(\gamma_i = 1 | \mathbf{y}, \mathbf{X}) = \frac{1}{T - T_0 + 1} \sum_{t=T_0}^T \mathbb{I}_{\gamma_i^{(t)}=1}$$



Probabilities of Inclusion Estimates

Informative ($\tilde{\beta} = 0_{11}$, $c=100$) and non-informative G-prior variable inclusion estimates (based on the same Gibbs output as in the earlier two tables)

γ_i	$\hat{P}^\pi(\gamma_i = 1 \mathbf{y}, \mathbf{X})$	$\hat{P}^\pi(\gamma_i = 1 \mathbf{y}, \mathbf{X})$
gamma_i	P(gamma_i y, X)	P(gamma_i y, X)
gamma_1	0.8733	0.8806
gamma_2	0.7100	0.7789
gamma_3	0.1515	0.2958
gamma_4	0.6842	0.7422
gamma_5	0.6635	0.7234
gamma_6	0.1659	0.2992
gamma_7	0.1343	0.2812
gamma_8	0.1478	0.2740
gamma_9	0.3942	0.5015
gamma_10	0.1135	0.2556