

Bayesian Model Comparison

Sargur Srihari

srihari@cedar.buffalo.edu

Model Selection: Bayesian Perspective

- Use of probabilities to represent uncertainties in choice of model
- Compare a set of models $M_i, i=1,..L$
- A model is a probability distribution over observed data D
 - E.g., in polynomial curve-fitting distribution over target values t when input X is known
- Given training set, wish to evaluate posterior distribution

$$P(M_i|D) \propto p(M_i) \underbrace{p(D|M_i)}$$

Model Evidence

Using Model Distribution

- Given Posterior Distribution over models, predictive distribution can be obtained

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

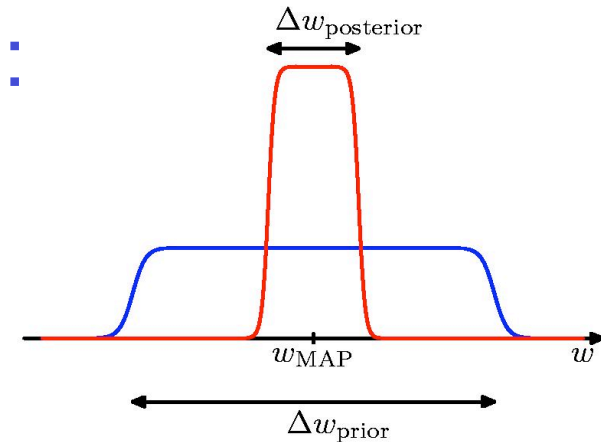
- A mixture distribution where predictions under different models are weighted by model evidence
- Using single most probable model is *Model Selection*
- From model evidence $p(D/M_j)$ for each model we can compute
 - Ratio $p(D/M_j)/p(D/M_j)$ called *Bayes Factor*

Optimal Model Complexity

Model Evidence, by sum rule:

$$p(D | M_i) = \int p(D | w, M_i) p(w | M_i) dw$$

A simple approximation



Consider model with single parameter w

If posterior probability is sharply peaked around most probable value w_{MAP}

$$\ln p(D) \approx \underbrace{\ln p(D | w_{MAP})}_{\text{Fit to data given most probable parameter values}} + M \ln \left(\frac{\Delta w_{posterior}}{\Delta w_{prior}} \right)$$

Fit to data given most
probable parameter values

First term will decrease with model complexity since it better fits the data

Second term will increase due to dependence on M

Trade-off between terms determines optimal model complexity

Favoring Intermediate Complexity

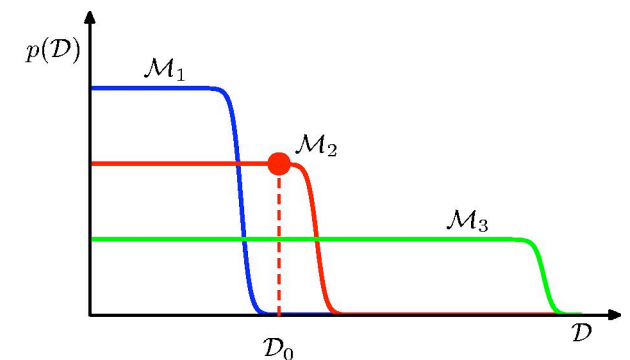
M_1 , M_2 and M_3 are models of increasing complexity

Simple one (e.g., first order polynomial)
can only generate simple data sets

Complex one (ninth order polynomial)
can generate complex sets

Middle one is favored
over the too simple (M_1) and too complex (M_3)

Distribution of Data Sets for 3 Models



Evidence Approximation

Application to Linear Regression

- Fully Bayesian treatment includes
 - Prior distributions over hyper-parameters α and β
 - Prediction by marginalizing wrt hyper-parameters and w
- Analytically intractable
- Instead hyper-parameters are set to values determined by maximizing marginal likelihood functions

Predictive Distribution

Obtained by marginalizing over w , α and β

Dependence on input variable x is omitted

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

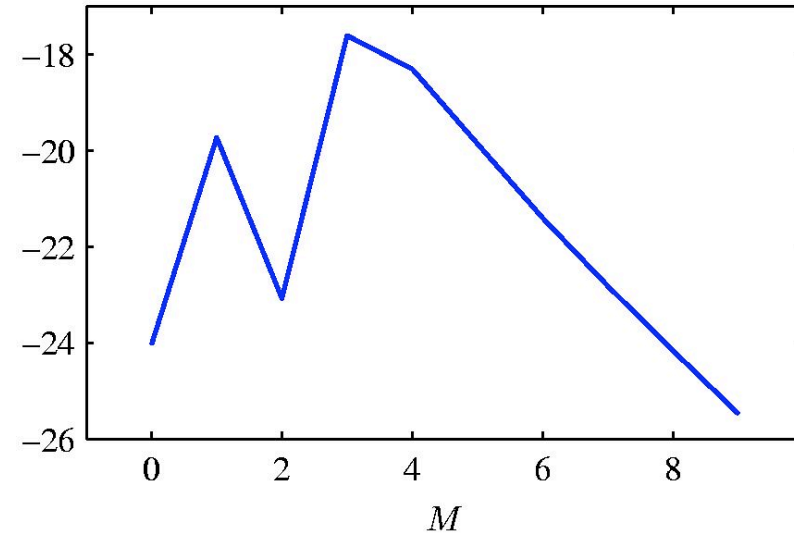
Evaluation of the Evidence Function

$$p(t | \alpha, \beta) = \int p(t | w, \beta) p(w | \alpha) dw$$

Model Evidence vs Order M

Polynomial Regression

Evidence favors model with $M=3$

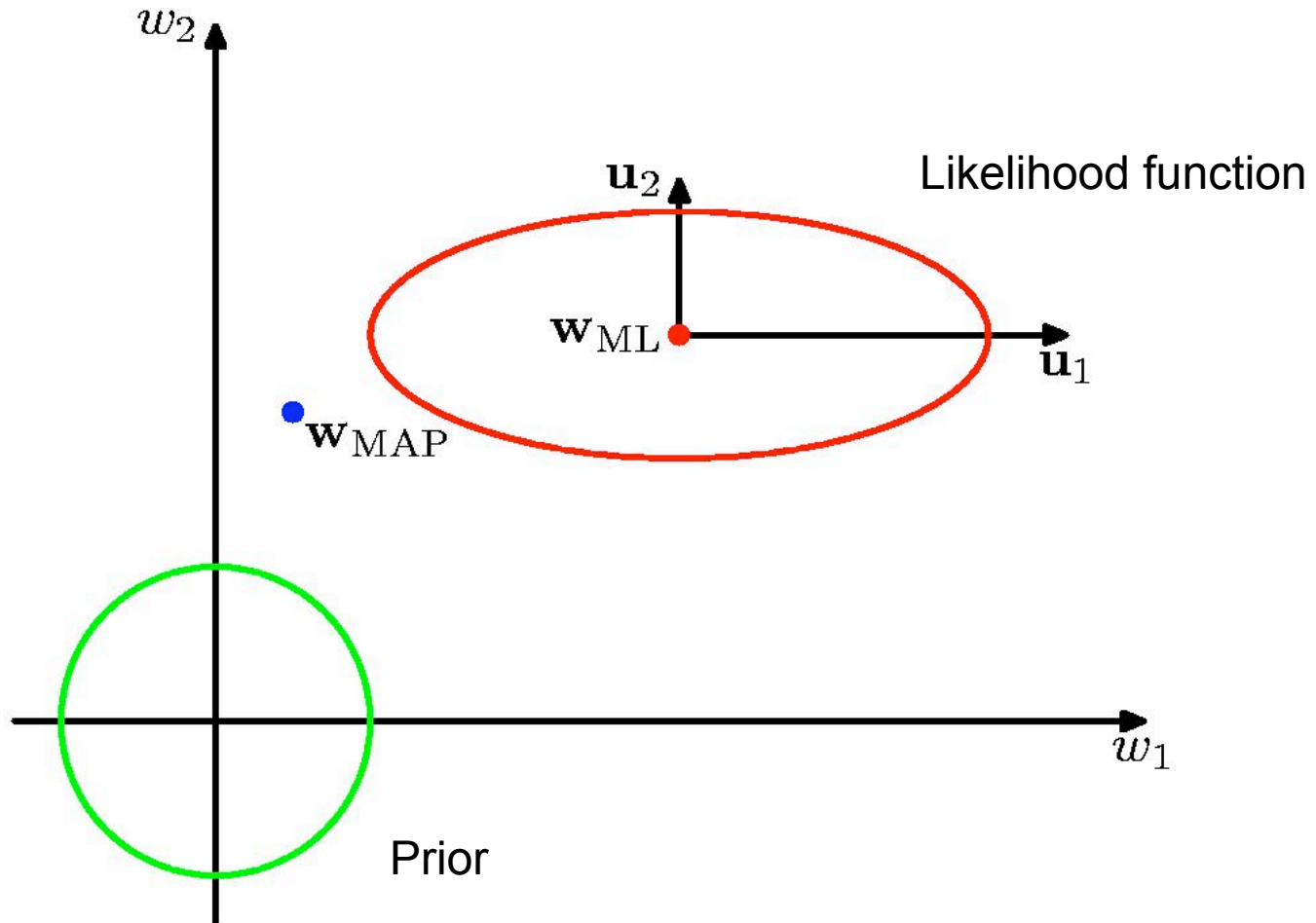


Order of Polynomial

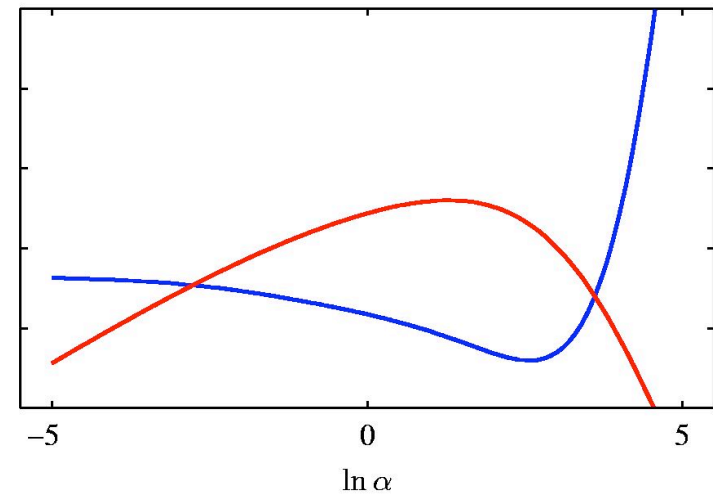
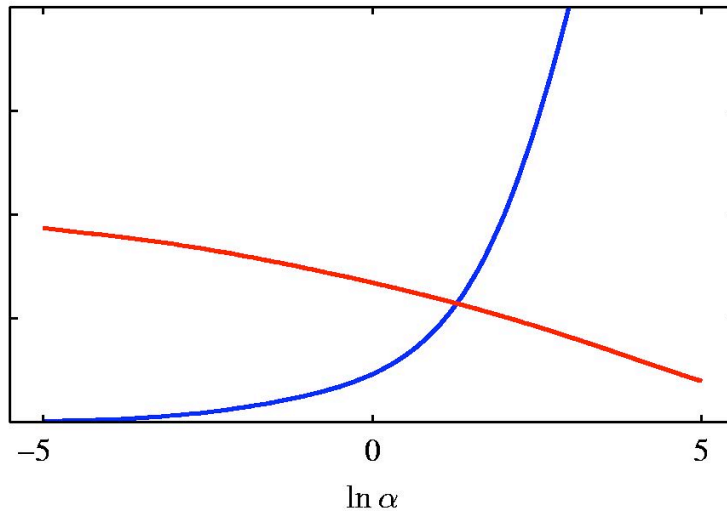
Data has a poor fit for $M=2$ since the polynomial expansion for sinusoid is an odd function and has no even terms

Maximizing the Evidence Function

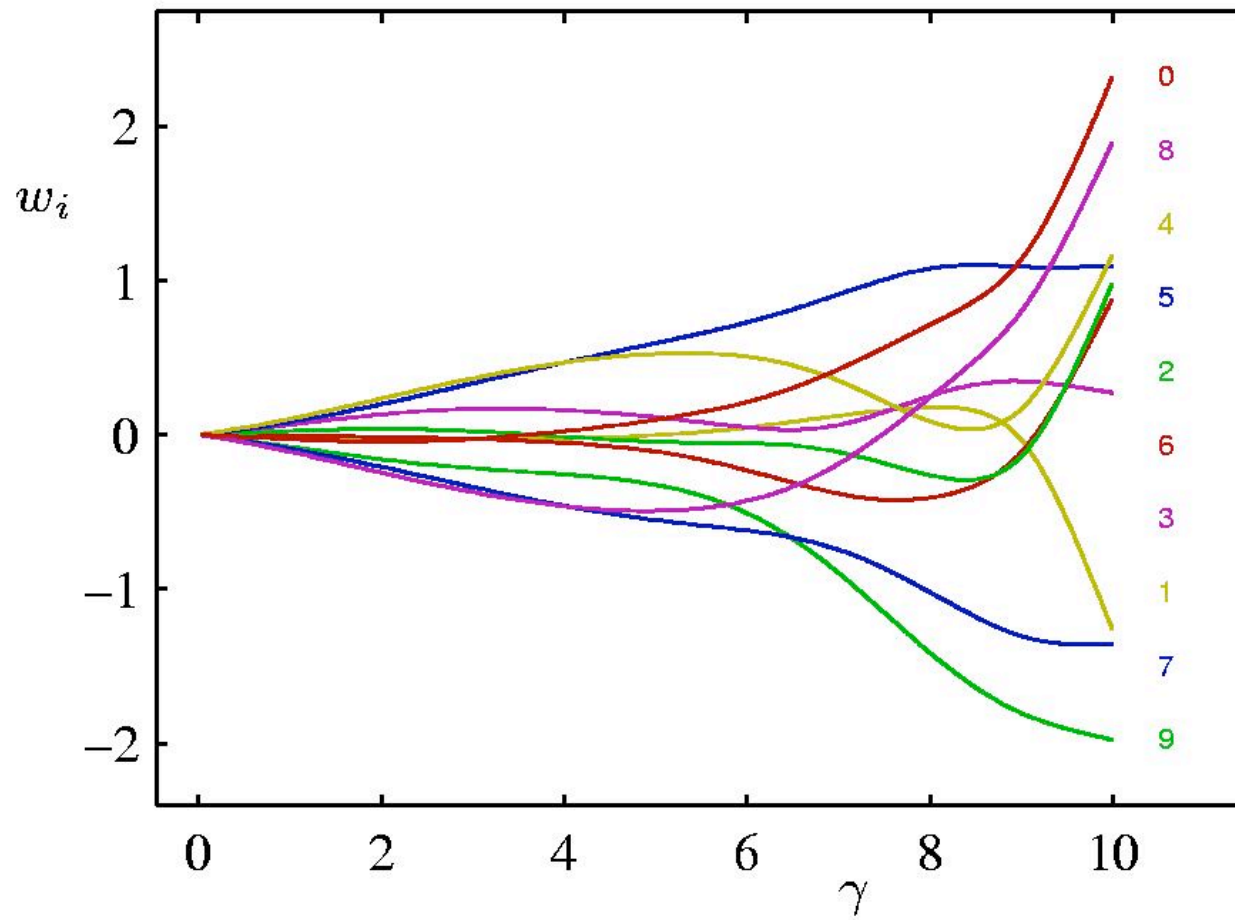
Effective No of Parameters



Optimum value for α



Plot of 10 parameters from Gaussian Basis



Limitations of Fixed Basis Functions

- Although we replaced D parameters with M basis functions, the number of basis functions needs to grow exponentially with D
- Solutions
 - Radial basis functions networks and SVMs exploit intrinsic dimensionality of data
 - Neural networks, with sigmoidal nonlinearity adapt to the data manifold

Summary of Linear Regression

- Goal: predict value of one or more target variables t given value of D -dimensional vector \mathbf{x} of input variables
 - Polynomial is a simple form
- We looked at probabilistic methods
- Linear regression models are linear in parameter values, but use nonlinear functions of input variables (called basis functions)
- Maximum likelihood formulation is equivalent to Least Squares
 - Bias Variance tradeoff exists
 - Partitioning data set is not satisfactory use of data
- Bayesian formulation is most general solution
 - Weighted average over parameter values
- Limitations when the dimensionality is large