

Evidence Approximation

Sargur Srihari

srihari@cedar.buffalo.edu

Background

- Bayesian methods need priors with hyper-parameters
- Analytical treatment with distributions over hyper-parameters is intractable
- Approximate methods are used
- Also called as:
 - *Empirical Bayes* (Statistics literature)
 - *Type 2 Maximum Likelihood*
 - *Generalized Maximum Likelihood*
 - *Evidence Approximation* (ML literature)

Topics

- Linear Regression: Bayesian Treatment
 - Hyper-parameters for noise and weights
 - Predictive distribution
 - Marginalize over hyper-parameters and weights
- Need for Approximation
 - Evaluation of evidence function
 - Maximizing the evidence function
- Interpretation: Effective no. of parameters

Linear Regression with Basis Functions

- Target variable is a scalar t given by
 - deterministic function $y(\mathbf{x}, \mathbf{w})$ with additive Gaussian noise

$$t = y(\mathbf{x}, \mathbf{w}) + \varepsilon$$

ε is *zero-mean* Gaussian with *precision* β

- Linear regression

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

- where $\mathbf{w} = (w_0, w_1, \dots, w_{M-1})$ and $\boldsymbol{\phi} = (\phi_0, \phi_1, \dots, \phi_{M-1})^T$
- $\phi_j(\mathbf{x})$ are called Basis functions
- There are M parameters

Bayesian treatment

- Prior distribution of parameter $p(\mathbf{w})$
- Since likelihood $p(t|\mathbf{w})$ with Gaussian noise has an exponential form

– Conjugate prior is chosen to be Gaussian

$$p(\mathbf{w}) = N(\mathbf{w} | \mathbf{m}_0, S_0) \text{ with mean } \mathbf{m}_0 \text{ and covariance } S_0$$

- Posterior is a Gaussian $p(\mathbf{w} | \mathbf{t}) = N(\mathbf{w} | \mathbf{m}_N, S_N)$
 where $\mathbf{m}_N = S_N(S_0^{-1}\mathbf{m}_0 + \beta\Phi^T\mathbf{t})$ and $S_N^{-1} = S_0^{-1} + \beta\Phi^T\Phi$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & & & \\ & & & \\ \phi_0(\mathbf{x}_N) & & & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

$N \times M$ Design matrix

Note that $\Phi^T\Phi$ is a $M \times M$ square matrix

Parameters and Hyper-parameters

- Parameters \mathbf{w} , a vector of M elements
 - Where M is model complexity

- Two hyper-parameters:

- α is weight precision

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

- Zero-mean isotropic Gaussian
- Has single parameter α

- β is noise precision

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

Predictive Distribution

- Usually not interested in the value of w itself
 - But predicting t for new values of x
- Introducing hyper-priors over α and β and
 - marginalizing over w , α and β

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

Conditioning variable x has been left out for convenience.

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

$$p(w | t) = N(w | m_N, S_N)$$

- Marginalizing wrt α , β , w is analytically intractable
- If $p(\alpha, \beta | t)$ is sharply peaked around $\hat{\alpha}$, $\hat{\beta}$ then

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

Strategy of Approximation

- Goal: determine $\hat{\alpha}, \hat{\beta}$ where $p(\alpha, \beta | t)$ is maximum
- From Bayes theorem

$$p(\alpha, \beta | t) \propto p(t | \alpha, \beta) p(\alpha, \beta)$$

- If prior $p(\alpha, \beta)$ is flat, values $\hat{\alpha}$ and $\hat{\beta}$ obtained by maximizing marginal likelihood $p(t | \alpha, \beta)$
- Can analytically determine maximum of $p(t | \alpha, \beta)$ for linear basis model

Evidence Function Evaluation

Obtained by integrating over weight parameter \mathbf{w}

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

Where we substitute

$$\begin{aligned} p(\mathbf{t} | \mathbf{w}, \beta) &= \prod_{i=1}^N N(\mathbf{t}_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \\ p(\mathbf{w} | \alpha) &= N(\mathbf{w} | \mathbf{0}, \alpha^{-1} I) \end{aligned}$$

And get

$$\begin{aligned} p(\mathbf{t} | \alpha, \beta) &= \left(\frac{\beta}{2\pi} \right)^{N/2} \left(\frac{\alpha}{2\pi} \right)^{M/2} \int \exp(-E(\mathbf{w})) d\mathbf{w} \\ \text{where} \\ E(\mathbf{w}) &= \beta E_D(\mathbf{w}) + \alpha E_w(\mathbf{w}) \\ &= \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{w}\|^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \end{aligned}$$

$E(\mathbf{w})$ is the regularized sum of squares error
Upto a constant of proportionality

Final Expression for Model Evidence

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

- Where M is the dimensionality of \mathbf{w} , N is no of samples, and

$A = \alpha I + \beta \Phi^T \Phi$ which implies that $A = \nabla \nabla E(\mathbf{w})$ is the Hessian matrix of Error

$$E(\mathbf{w}) = E(\mathbf{m}_N) + \frac{1}{2} (\mathbf{w} - \mathbf{m}_N)^T A (\mathbf{w} - \mathbf{m}_N)$$

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

$$\mathbf{m}_N = \beta A^{-1} \Phi^T \mathbf{t}$$

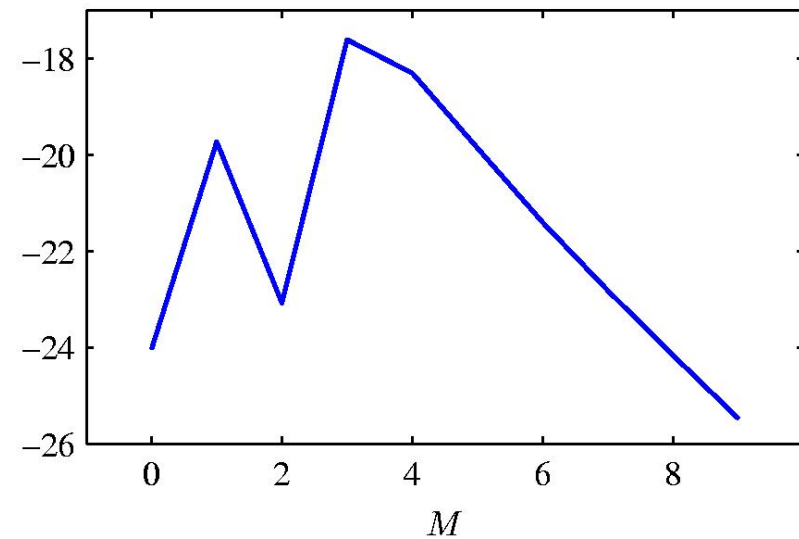
- We wish to maximize wrt M, α, β

Maximize Model Evidence over M

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

We need α and β to maximize this

Polynomial Regression
Model Evidence vs M



Order of Polynomial

Evidence favors model with $M=3$

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

Maximize Evidence over α

- Maximization of $p(\mathbf{t}|\alpha, \beta)$ wrt α yields

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

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

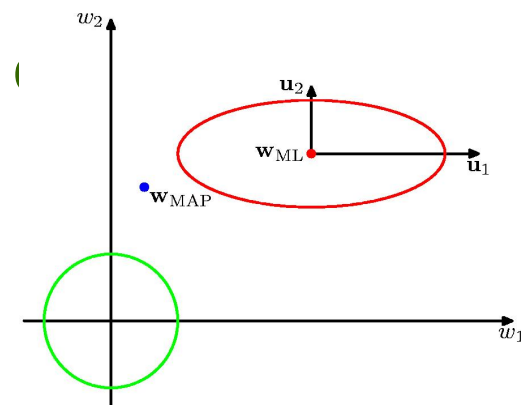
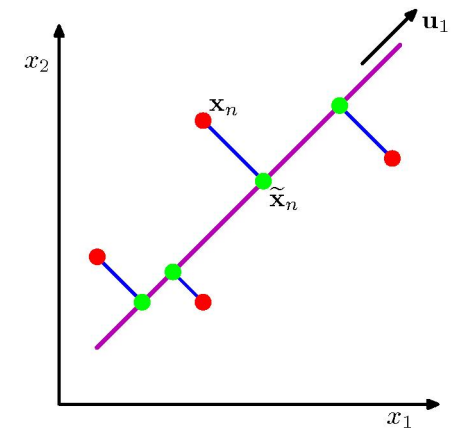
- where λ_i are β times the eigen values of $\Phi^T \Phi$
 - Defined by eigenvector equation $(\beta \Phi^T \Phi) \mathbf{u}_i = \lambda_i \mathbf{u}_i$

Eigen Value tutorial

- For a square matrix A of size $M \times M$
 - *Eigen vector equation* is $A\mathbf{u}_i = \lambda_i \mathbf{u}_i$ for $i=1, \dots, M$ where \mathbf{u}_i is an eigenvector and λ_i the corresponding eigenvalue
 - Can be viewed as M simultaneous equations which are solutions to *characteristic equation* $|A - \lambda_i I| = 0$

Eigen Value Interpretation

- Eigen values are real-valued for symmetric matrices
- For a covariance matrix Σ
 - Eigen vector corresponding to largest eigen value indicates direction of maximum variance
 - Principal Components
- For $\Phi^T \Phi$
 - large eigen values indicate direction of high curvature of likelihood function



Iterative Solution for α

- Maximization of $p(\mathbf{t}|\alpha, \beta)$ wrt α yields

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

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

- An implicit solution

– since γ and mode \mathbf{m}_N of posterior depend on α

- Iterative procedure:

– Initial choice for α , find \mathbf{m}_N etc until convergence

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

- α determined purely from training data

– No independent data set required to optimize model

Maximize Evidence over β

- Maximization of $\ln p(t|\alpha,\beta)$ wrt noise precision β yields

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

- Again, an implicit solution

$$\begin{aligned} m_N &= \beta A^{-1} \Phi^T t \\ A &= \alpha I + \beta \Phi^T \Phi \end{aligned}$$

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

- Iterative Solution:
 - Choose initial value of β , calculate m_N , γ and re-estimate β until convergence

Effective No of Parameters

Interpretation provides insight into Bayesian solution for α

$$\alpha = \frac{\gamma}{\mathbf{m}_N^T \mathbf{m}_N} \quad \gamma = \sum_i \frac{\lambda_i}{\alpha + \lambda_i}$$

where λ_i are eigen values of Hessian $\Phi^T \Phi$

$\lambda_i / (\lambda_i + \alpha)$ lies between 0 and 1.

When λ_i is small, corresponding weight is close to zero.

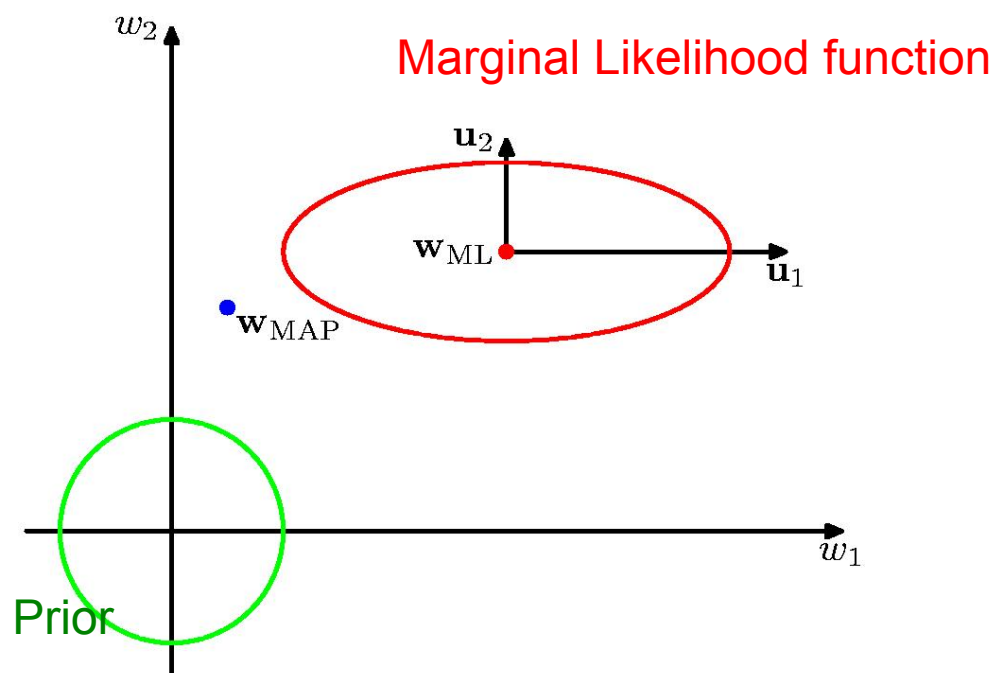
Here w_1 is small and w_2 is larger.

Thus

γ measures the effective total no of parameters

Axes in parameter space rotated

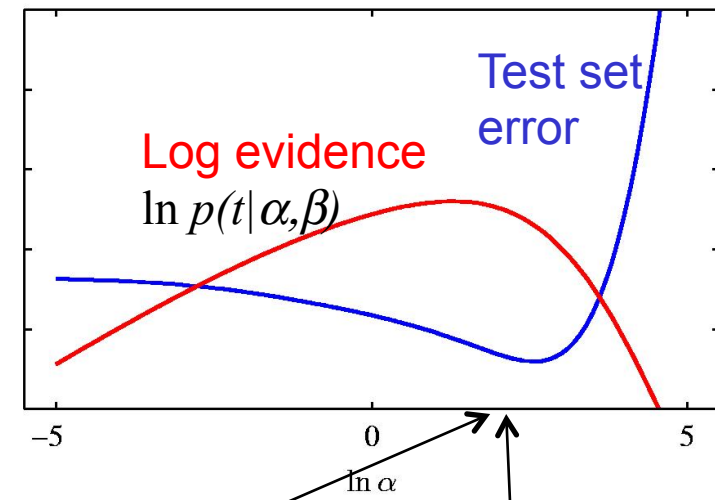
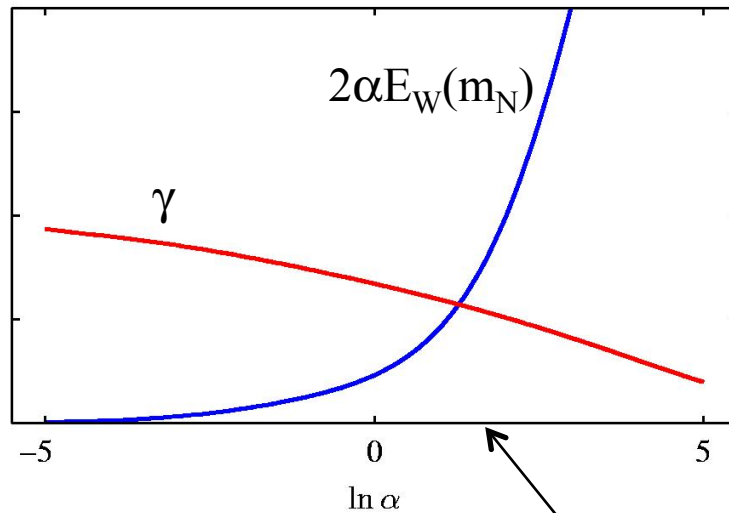
To align with Eigen vectors \mathbf{u}_i of Hessian $A = \alpha I + \beta \Phi^T \Phi$



For $\alpha = 0$, \mathbf{w}_{ML} is the mode (max likelihood solution)

For nonzero α , $\mathbf{w}_{MAP} = \mathbf{m}_N$ is the mode

Determining optimum value of α

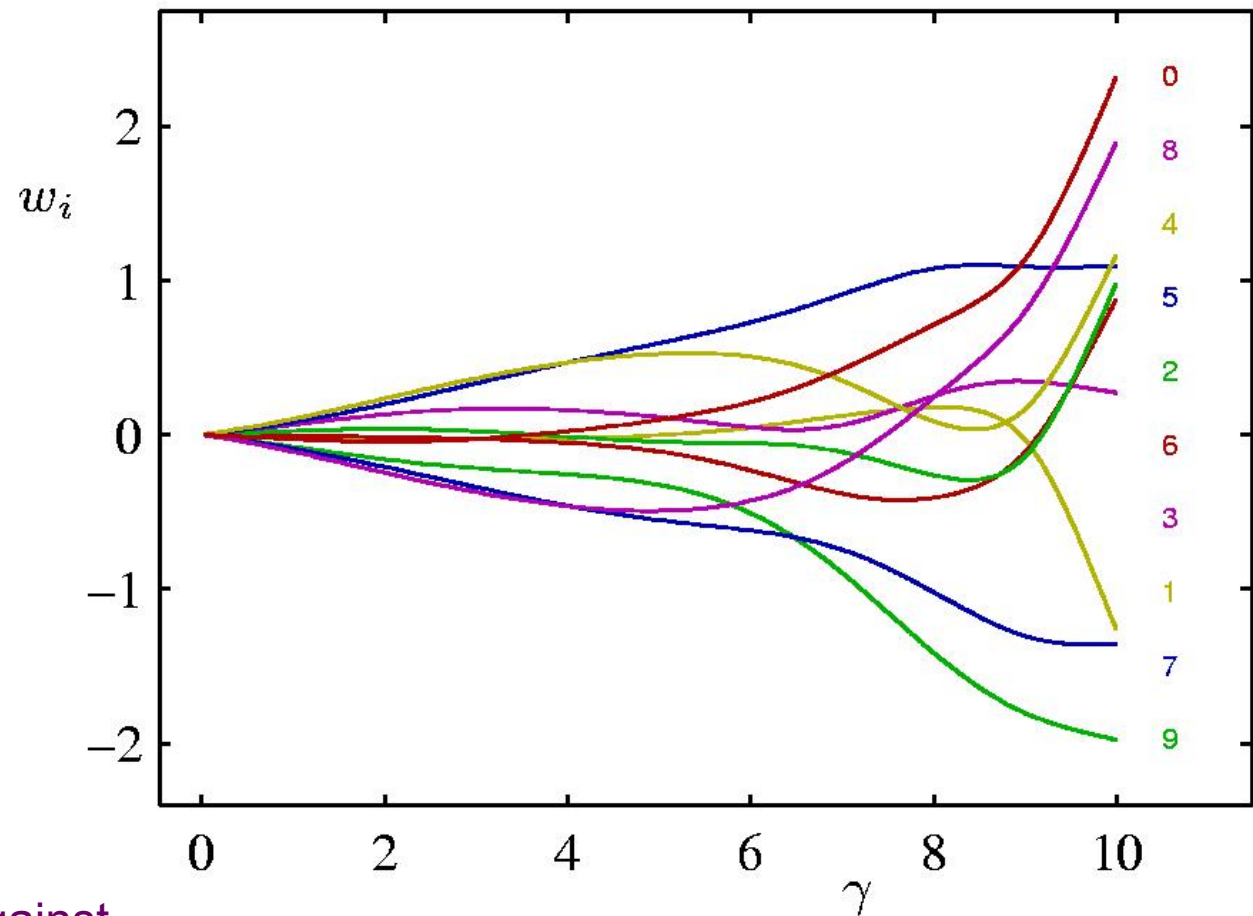


For synthetic
sinusoidal data set
with 9 basis functions
Noise precision β
set to true value of 11.1

Crossing point and peak
occur at same point

Evidence maximum
occurs close to point
of best generalization

How α controls parameters w_i



Plot of 10 parameters from Gaussian basis function model:
Individual parameters against effective no of parameters

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
 - Hyper-parameters are estimated using evidence approximation
- Limitations when the dimensionality is large