

Evidence Approximation: Determining hyper-parameters

Sargur Srihari

srihari@cedar.buffalo.edu

Topics

- Linear Regression with Basis Functions
- Fully Bayesian Treatment
 - Hyper-parameters for noise and weights
 - Predictive distribution
 - Marginalize over hyper-parameters and weights
- Need for Approximation
 - Called *evidence approximation* or *empirical Bayes*
 - Evaluation of evidence function
 - Maximizing the evidence function
- Interpretation: Effective no. of parameters

Linear Regression with Basis Functions

- Polynomial regression extended by considering nonlinear functions of input variables

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x})$$

- where $\phi_j(\mathbf{x})$ are called Basis functions
- There are M parameters instead of d parameters
- Can be written as

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

- where $\mathbf{w} = (w_0, w_1, \dots, w_{M-1})$ and $\boldsymbol{\phi} = (\phi_0, \phi_1, \dots, \phi_{M-1})^T$

Noise Model

- 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

- Thus distribution of t is univariate normal:

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

Fully 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 $S_N^{-1} = S_0^{-1} + \beta\Phi^T\Phi$

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

Design matrix

Hyper-parameters

- There are two hyper-parameters in this Bayesian treatment
 - β is precision of noise
 - α is precision of weights
- Zero mean isotropic Gaussian

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

Single precision parameter

Predictive Distribution

- Usually not interested in the value of w itself
- But predicting t for new values of x
- For known α and β

$$p(t | t, \alpha, \beta) = \int p(t | w, \beta) p(w | t, \alpha, \beta) dw \quad \text{Marginalizing } w$$

Conditioning variable x has been left out for convenience.

- If we have prior distributions for α and β

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

marginalizing over w, α and β

Need for Approximation

- Fully Bayesian treatment of linear basis function model includes
 - Prior distributions over hyper-parameters α and β
 - Prediction by marginalizing wrt α , β and w
- Analytically intractable
- Instead α and β are set to values determined by maximizing *marginal likelihood* functions
 - obtained by integrating over parameters w

Several Names for Approach

- Set hyper-parameters to specific values
 - by maximizing marginal likelihood function
 - obtained by integrating over parameters w
- This framework has several names:
 - *Empirical Bayes* (Statistics literature)
 - *Type 2 Maximum Likelihood*
 - *Generalized Maximum Likelihood*
 - *Evidence Approximation* (ML literature)

Predictive Distribution

Predictive distribution obtained by marginalizing over w , α and β

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

Dependence on input variable \mathbf{x} is omitted

Need posterior distribution of hyper-parameters

If $p(\alpha, \beta | \mathbf{t})$ is sharply peaked around $\hat{\alpha}$, $\hat{\beta}$ then

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

Posterior of Hyperparameters

- From Bayes theorem

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

- If prior is flat, values $\hat{\alpha}$ and $\hat{\beta}$ obtained by maximizing marginal likelihood function

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

- Can analytically evaluate $p(t | \alpha, \beta)$ for linear basis function model and then find its maximum

Evaluation of the Evidence Function

Also called the *Marginal likelihood* function

It is obtained by integrating over weight parameter w

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

Can derive the following log-marginal likelihood

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

We need α and β to maximize this

where M is the dimensionality of w , N is no of samples, and

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

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

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

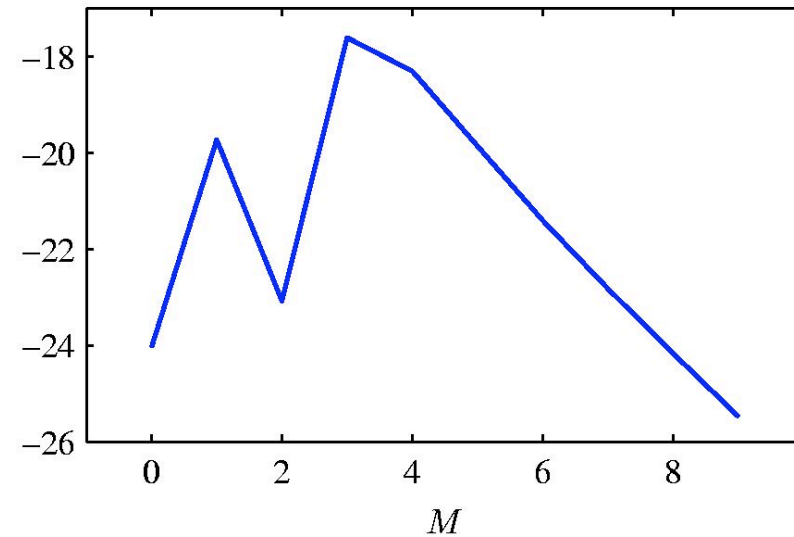
$$m_N = \beta A^{-1} \Phi^T t$$

Model Evidence vs Order M

Polynomial Regression

Plot of model evidence against
Order of polynomial

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

Maximize Evidence Function for α

- Maximization of $p(\mathbf{t}|\alpha, \beta)$ wrt weight precision 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 eigen values of $\Phi^T \Phi$
- An implicit solution
 - since γ and mode \mathbf{m}_N of posterior depend on α
- Iterative procedure is used:
 - Initial choice for α , find \mathbf{m}_N etc until convergence
- α determined purely from training data
 - No independent data set required to optimize model complexity

Maximize Evidence Function for β

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

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

- Again, an implicit solution
- Solved by choosing initial value of β , using this to calculate \mathbf{m}_N and γ and re-estimate β until convergence

Effective No of Parameters

Interpretation provides insight into Bayesian solution for α

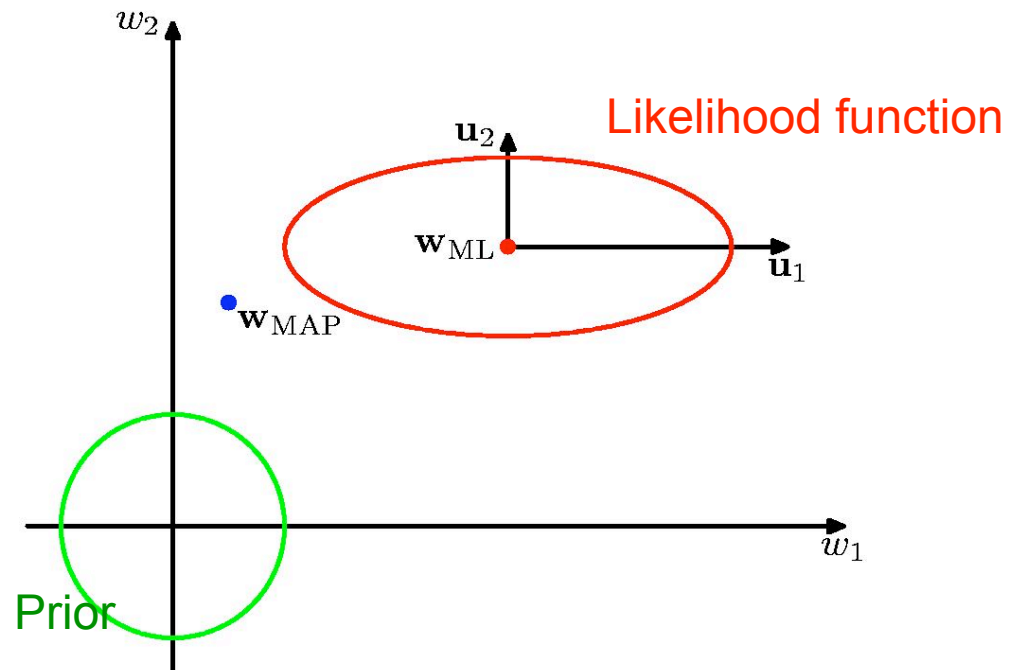
$$\alpha = \frac{\gamma}{\mathbf{m}_N^T \mathbf{m}_N}$$

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

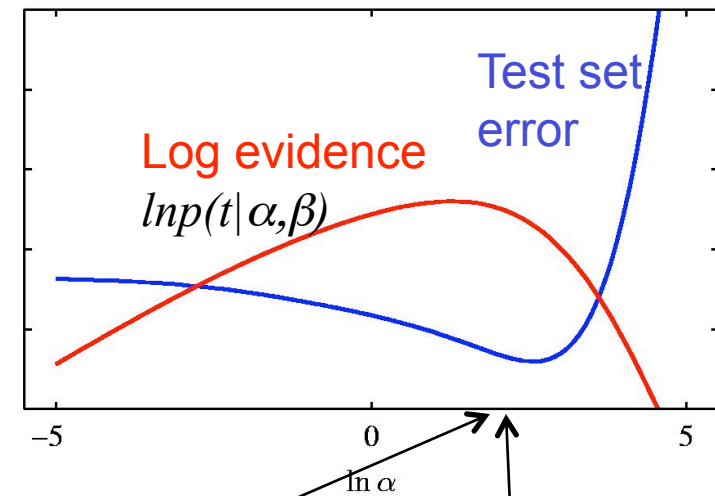
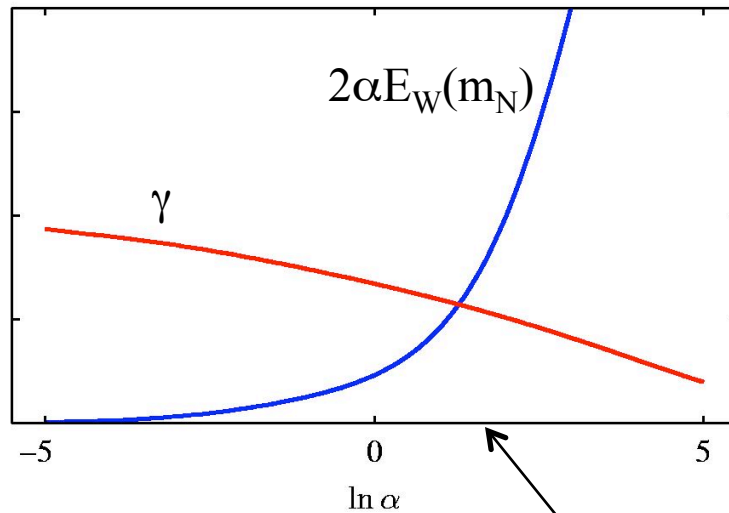
where λ_i are eigen values of Hessian $\Phi^T \Phi$

γ measures the effective Total no of parameters

Axes in parameter space rotated
To align with Eigen vectors \mathbf{u}_i of Hessian



Optimum for hyper-parameter α

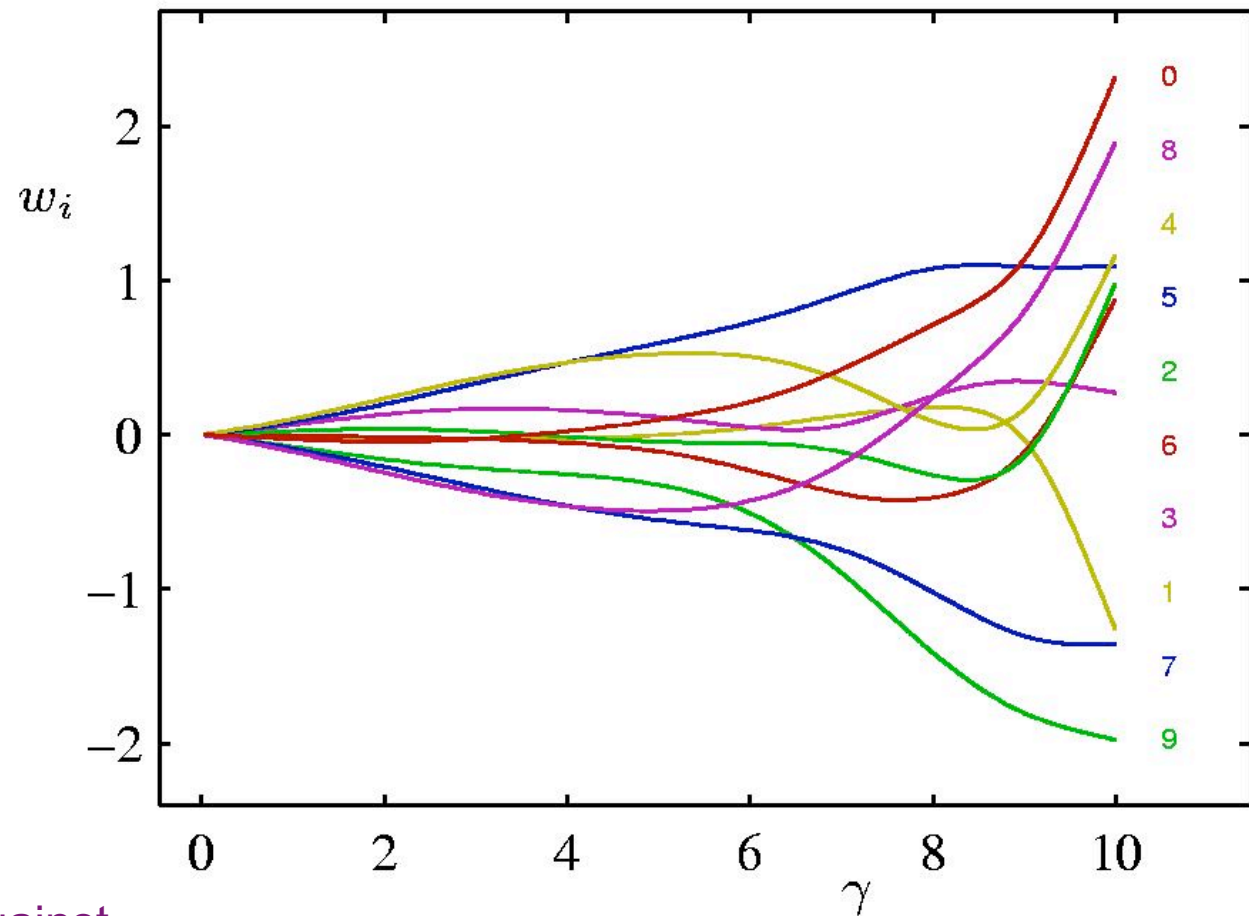


For synthetic
sinusoidal data set
with 9 basis functions
Noise precision β
set to true value of 11.1

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