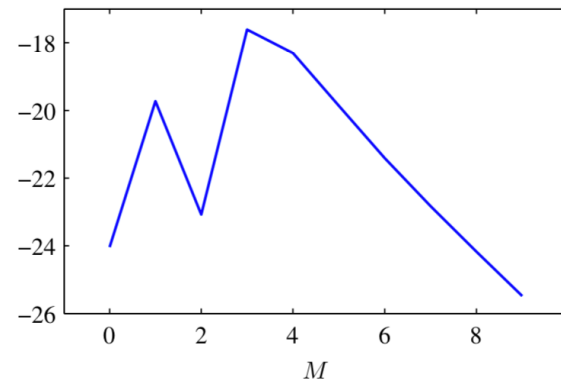
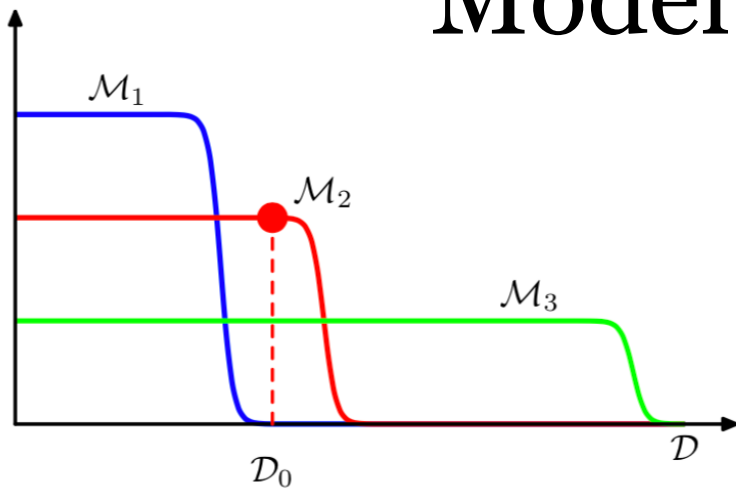


Bayesian Linear Models: Model Selection



SPR Day 09, Spring 2020

Prof. Mike Hughes

<https://www.cs.tufts.edu/comp/136/2020s/>

Recap: Bayesian Linear Regression

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$$

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

Recap: Bayesian Linear Regression

General prior

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$$

*Simpler prior: assume zero mean,
Just need to define precision*

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

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

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

$$\mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t})$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi.$$

Problem:

Hyperparameter Selection

How do we pick the prior hyperparameters?

alpha

How do we pick the likelihood hyperparameters?

beta

Hyperparameter Selection

	Fixed valid. set (fraction f)	K-fold cross-validation	Bayesian evidence	
Fraction data used for training run	$(1.0 - f)$	$(K-1) / K$	1.0	Higher is better Better use of training data
Total runs/ examples seen for training	1 run $(1 - f) N$	K runs $(K-1) * N$	1 run N	Lower is better Faster training
Total runs/ examples seen for evaluation of fitness	1 run fN	K runs N	1 run N	Lower is better Faster evaluation
Fitness function	Heldout likelihood	Heldout likelihood	Evidence	

Why use Model Evidence?

3.4. Bayesian Model Comparison

As we shall see, the over-fitting associated with maximum likelihood can be avoided by marginalizing (summing or integrating) over the model parameters instead of making point estimates of their values. Models can then be compared directly on the training data, without the need for a validation set. This allows all available data to be used for training and avoids the multiple training runs for each model associated with cross-validation.

Related Problem: Model Selection

How do we pick the feature transform?

$$\phi(x_n) = [\phi_1(x_n) \quad \phi_2(x_n) \quad \dots \quad \phi_M(x_N)]$$

What size M?

Which functions to include?

Challenge: Could have unbounded number of choices!

Need to enumerate **small set** of possible models (size L) if want to average over them properly

Model Selection for Regression

 \mathcal{M}_i

Model family (size M, feature functions, hyperparameters)

 w

Specific parameters for chosen model

 $t_n | x_n$

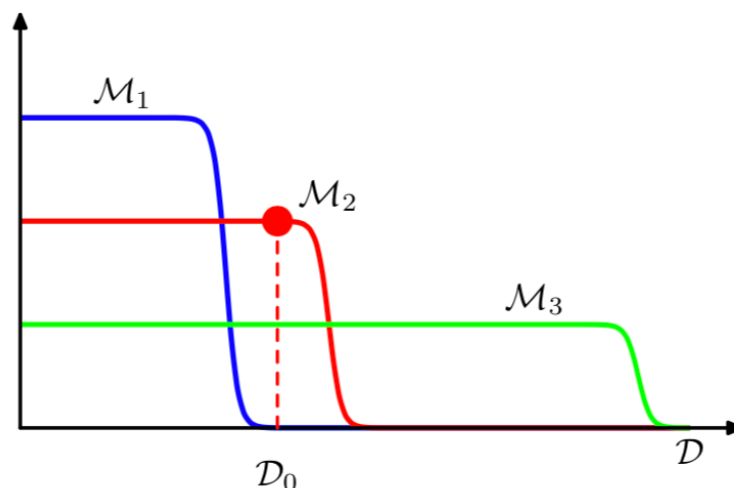
Observed outputs (and corresponding inputs)

$$t_n | x_n \sim \mathcal{N}(w^T \phi(x_n), \beta^{-1})$$

Model Selection via Evidence

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

Figure 3.13 Schematic illustration of the distribution of data sets for three models of different complexity, in which \mathcal{M}_1 is the simplest and \mathcal{M}_3 is the most complex. Note that the distributions are normalized. In this example, for the particular observed data set \mathcal{D}_0 , the model \mathcal{M}_2 with intermediate complexity has the largest evidence.



Key idea: “Goldilocks” principle

If model **too simple**, it puts high mass only few datasets

If model **too complex**, it puts mass on too many datasets

Predictive distribution

Average over predictive distribution for each of L possible models, weighted by posterior probability

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

Predictive distrib. for t given model i Posterior of model i

Key idea: We can use all L models, don't need to pick one

Ideal predictive posterior

If we want to predict new data given old data, ideally we would **average over** all parameters w , α , β , weighted by posterior prob.

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

But, this integral is hard!

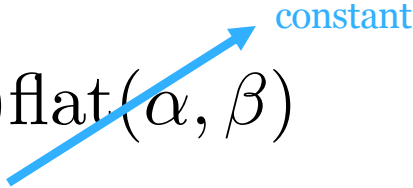
Tractable predictive posterior

Assume:

- We have enough data that the posterior $p(\alpha, \beta | \mathbf{t})$ is peaked at MAP point estimate

$$\hat{\alpha}, \hat{\beta} = \arg \max_{\alpha, \beta} p(\alpha, \beta | \mathbf{t})$$

- Prior on alpha, beta is relatively uniform (“flat”), so these estimates might as well be ML estimates

$$\hat{\alpha}, \hat{\beta} = \arg \max_{\alpha, \beta} p(\mathbf{t} | \alpha, \beta) \text{flat}(\alpha, \beta)$$


Then the **tractable estimate** of predictive posterior becomes:

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

Now, we wish to solve this hyperparameter estimation

$$\hat{\alpha}, \hat{\beta} = \arg \max_{\alpha, \beta} p(\mathbf{t} | \alpha, \beta)$$

“Evidence”

But first, what is this “evidence” anyway?

Evidence for Linear Regression

- Probability of training data \mathbf{t} given alpha, beta
- Marginalizes over the weights \mathbf{w}

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

Simplifying the Evidence

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

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

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

Key ideas:

- Bring constants outside integral
- Recognize inside integral as Gaussian, “complete the square”

Closed-form Log Evidence

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

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

$$\mathbf{A} = \alpha \mathbf{I} + \beta \Phi^T \Phi$$

Precision matrix for posterior $p(\mathbf{w} | \text{data})$

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

Mean vector for posterior $p(\mathbf{w} | \text{data})$

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

How to estimate?

$$\hat{\alpha}, \hat{\beta} = \arg \max_{\alpha, \beta} p(\mathbf{t} | \alpha, \beta)$$

- Can do gradient descent
- Can do coordinate descent (EM, later in course)
- Can get estimates analytically
 - See textbook!

Cycle these until convergence!

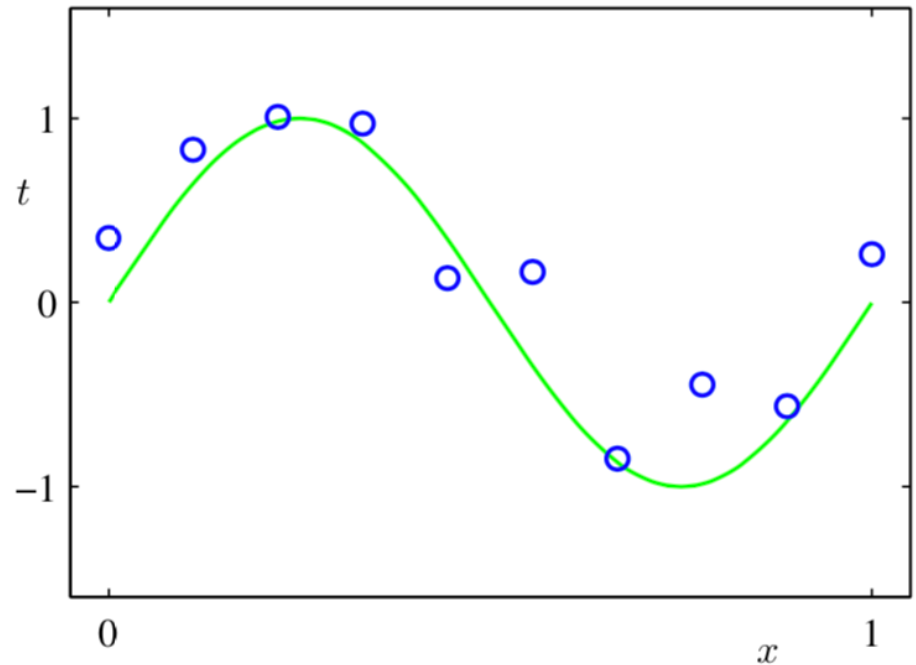
$$\alpha = \frac{\gamma}{\mathbf{m}_N^T \mathbf{m}_N} \quad \frac{1}{\beta} = \frac{1}{N - \gamma} \sum_{n=1}^N \{t_n - \mathbf{m}_N^T \phi(\mathbf{x}_n)\}^2$$

$$\gamma = \sum_i \frac{\lambda_i}{\alpha + \lambda_i} \quad (\beta \Phi^T \Phi) \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

Eigendecomposition!

Example: 1D sinusoid data

Figure 1.2 Plot of a training data set of $N = 10$ points, shown as blue circles, each comprising an observation of the input variable x along with the corresponding target variable t . The green curve shows the function $\sin(2\pi x)$ used to generate the data. Our goal is to predict the value of t for some new value of x , without knowledge of the green curve.



Model Selection for Linear Regr. (using polynomial features of order M)

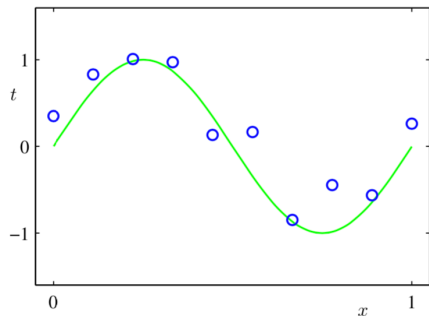
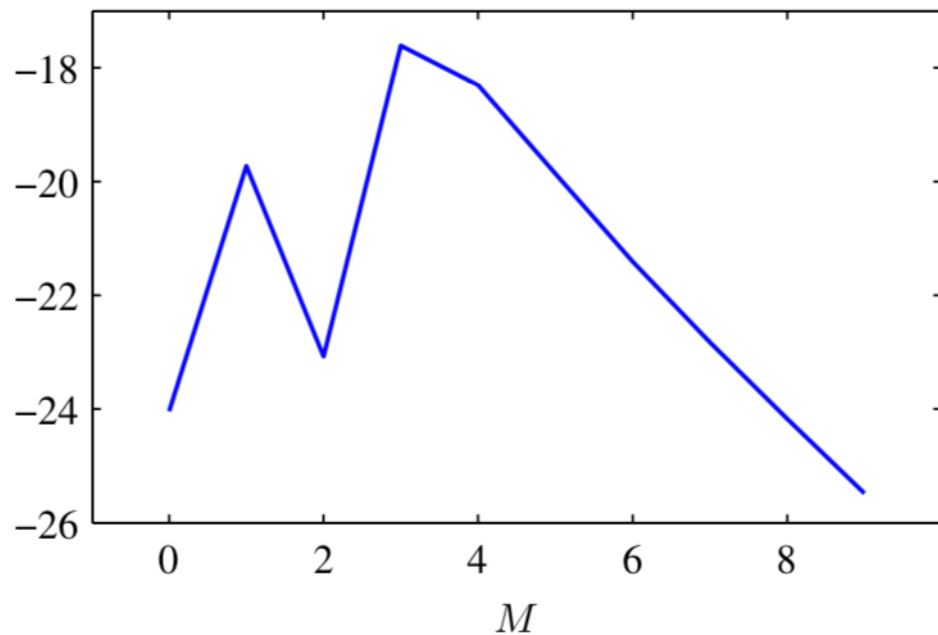


Figure 3.14 Plot of the model evidence versus the order M , for the polynomial regression model, showing that the evidence favours the model with $M = 3$.



Why does $M=2$ have low evidence?

Can set quadratic term to 0, but then we have a model that is “too complex”

Can set quadratic term non zero, but sinusoid (odd function) not a good fit by quadratic