

Bayesian Model Comparison

Sargur Srihari

srihari@cedar.buffalo.edu

Motivation

- In frequentist setting
 - Over-fitting is a problem
 - Cross-validation used for
 - Setting values for regularization parameters
 - Choosing between alternative models
- Bayesian model selection
 - Used to determine regularization parameters

Bayesian Perspective

- Avoids over-fitting
 - By marginalizing over model parameters
 - sum over model parameters instead of point estimates
- Models compared directly over training data
 - No need for validation set
 - Allows use of all available data in training set
 - Avoids multiple training runs for each model associated with cross-validation
 - Allows multiple complexity parameters to be simultaneously determined during training
 - Relevance vector m/c is Bayesian model with one complexity parameter for every data point

Model Selection based on Model Evidence

- Probabilities used to represent uncertainty 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

Model Evidence and Bayes Factor

- Posterior over models is

$$p(M_i|D) \propto \underbrace{p(M_i)}_{\text{Prior}} \underbrace{p(D|M_i)}_{\text{Model Evidence}}$$

- Prior expresses preference for different models
 - We assume models have equal prior probability
- $p(D|M_i)$ is preference shown by data for model
 - Called *Model evidence* or *marginal likelihood* (because parameters have been marginalized)
- From model evidence $p(D/M_i)$ for each model:
 - Ratio $p(D/M_i)/p(D/M_j)$ called *Bayes Factor*
 - Preference for one model over other

Predictive Distribution

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

$$p(t | \mathbf{x}, D) = \sum_{i=1}^L \underbrace{p(t | \mathbf{x}, M_i, D)}_{\text{Prediction under model}} p(M_i | D)$$

- This is a *mixture distribution* where predictions under different models are weighted by posterior probabilities of models
- Rather than use all models, one approximation is to use single most probable model
 - Known as *Model Selection*

Model Evidence (or Marginal Likelihood)

We are interested in best model as given by

$$p(M_i|D) \propto p(M_i) p(D|M_i)$$

How to get $p(D|M_i)$: marginal likelihood

For model M_i governed by parameters w :
from sum and product rules

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

Probability of generating the data from the model

This is obtained by marginalizing over parameters

Note: model specifies
no of basis functions
not values of parameters

Sampling perspective:

Marginal likelihood is probability of generating the data D from a model M_i whose parameters are sampled from the prior distribution of w

It is also the denominator in posterior distribution of w

$$p(w | D, M_i) = \frac{p(D | w, M_i) p(w | M_i)}{p(D | M_i)}$$

Approximation to Model Evidence $p(D|M_i)$

- Simple approx. to integral over parameters

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

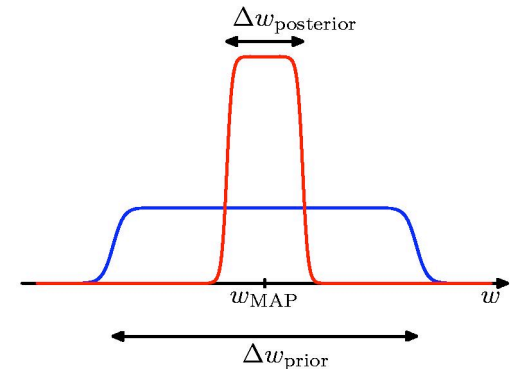
Consider model M_i with single parameter w

$$p(w|D) \propto p(D|w)p(w)$$

Omit dependence on M_i
to keep notation uncluttered

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

$$p(D) = \int p(D|w)p(w)dw \approx p(D|w_{MAP}) \frac{\Delta w_{posterior}}{\Delta w_{prior}}$$



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

Fit to data given most probable parameter values

Penalizes model according to its complexity

If parameters are finely tuned to data this term is large

Optimal Model Complexity

- For a model with M parameters

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

Fit to data given most probable parameter values

- Size of complexity penalty increases linearly with M

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

Second term will increase due to dependence on M

Trade-off between terms determines optimal model complexity

Favoring Intermediate Complexity

Marginal likelihood $p(D|M_i)$ favors models of 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

For observed data set D_0

Middle one has the highest value of the evidence for the model

Simpler model (M_1) cannot fit the data well

More complex model (M_3) assigns too small a probability for any given data set

Distribution of Data Sets for 3 Models

