

Model Selection: Frequentist vs Bayesian

Sargur Srihari
srihari@cedar.buffalo.edu

Topics

- Regression based on basis functions are determined by
 - M and λ
- How to choose them (model complexity)?
- Two approaches
 1. Frequentist setting
 - Cross-validation
 2. Bayesian model selection

Models in Curve Fitting

- Regression using basis functions and MSE:

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

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{ t_n - \mathbf{w}^T \boldsymbol{\phi}(x_n) \}^2$$

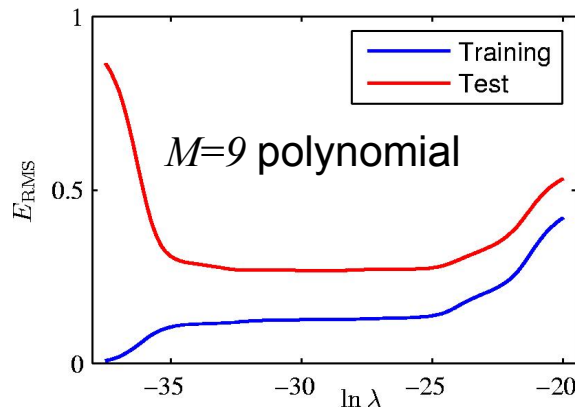
- Need an M that gives best generalization
 - M = No of free parameters in model or *model complexity*
- With regularized least squares

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{ t_n - \mathbf{w}^T \boldsymbol{\phi}(x_n) \}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- λ also controls model complexity (and hence degree of over-fitting)

Impact of Regularization on Error

- λ controls model complexity (similar to choice of M)
- Frequentist Approach:
 - Training set
 - To determine coefficients w for different values of (M or λ)
 - Validation set (holdout)
 - to optimize model complexity (M or λ)



What value of λ minimizes error?

$\lambda = -38$: low error on training, high error on testing set

$\lambda = -30$ is best for testing set

$$E_{RMS} = \sqrt{2E(w^*) / N}$$

Division by N allows different data sizes to be compared since E is a sum over N

Sqrt (of squared error) measures on same scale as t

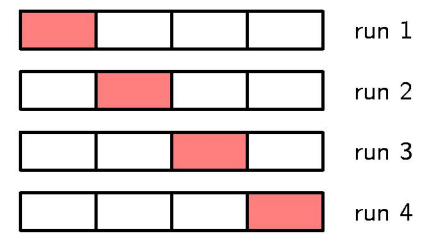
Validation Set to Select Model

- Performance on training set is not a good indicator of predictive performance
- If there is plenty of data,
 - use some of the data to train a range of models Or a given model with a range of values for its parameters
 - Compare them on an independent set, called validation set
 - Select one having best predictive performance
- If data set is small then some over-fitting can occur and it is necessary to keep aside a test set

S -fold Cross Validation

- Supply of data is limited
- All available data is partitioned into S groups
- $S-1$ groups are used to train and evaluated on remaining group
- Repeat for all S choices of held-out group
- Performance scores from S runs are averaged

$S=4$



If $S=N$ this is the leave-one-out method

Disadvantage of Cross-Validation

- No. of training runs is increased by factor of S
 - Problematic if training itself is expensive
- Different data sets can yield different complexity parameters for a single model
 - E.g., for given M several values of λ
- Combinations of parameters is exponential
- Need a better approach
 - Ideally one that depends only on a single run with training data for a given model

Akaike Information Criterion (AIC)

- AIC chooses model that maximizes

$$\ln p(D|w_{ML}) - M$$

- First term is best-fit log-likelihood for given model
 - M is no of adjustable parameters in model
-
- Penalty term M is added for over-fitting using many parameters
 - BIC is a variant of this quantity
 - Disadvantages:
 - need runs with each model, prefers overly simple models

AIC for Linear Regression

- AIC chooses model that maximizes

$$\ln p(D|\mathbf{w}_{\text{ML}}) - M$$

- In linear regression with basis fns

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

- Assuming additive Gaussian noise with precision β

- Log-likelihood term $\ln p(D|\mathbf{w}_{\text{ML}})$ can be written as

$$\ln p(t|\mathbf{w}, \beta) = \sum_{n=1}^N \ln N(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln 2\pi - \beta E_D(\mathbf{w})$$

- where

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{ t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) \}^2 + \lambda E_w(\mathbf{w})$$

$$E_w(\mathbf{w}) = \frac{1}{2} \sum_{j=1}^M |w_j|^q, \quad q=0.5, 1 \text{ (lasso)}, 2 \text{ (quadratic)}, 4$$

- M is the no of basis functions M plus one (for λ)

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

- A model is a probability distribution over data D
 - E.g., a polynomial model is a distribution over target values t when input x is known, i.e., $p(t|x, D)$
- Uncertainty in model itself can be represented by probabilities
- Compare a set of models $M_i, i=1, \dots, L$
- Given training set, wish to evaluate posterior

$$p(M_i|D) \propto \underbrace{p(M_i)}_{\text{Prior expresses preference for different models}} \underbrace{p(D|M_i)}_{\text{Model Evidence (or Marginal Likelihood)}}$$

We assume equal priors

Model Evidence (or Marginal Likelihood)

- $p(D|M_i)$ is preference shown by data for model
- Called *Model evidence* or *marginal likelihood*
 - because parameters have been marginalized

Bayes Factor

- From model evidence $p(D|M_i)$ for each model:
 - Ratio $\frac{p(D|M_i)}{p(D|M_j)}$ is called *Bayes Factor*
 - Preference for Model i over Model j

Predictive Distribution

- Given $p(M_i|D)$, predictive distribution is

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

Prediction under model

– *A mixture*

- predictions under different models are weighted by posterior probabilities of models
- Instead of all models M_i , approximate with single most probable model
 - Known as *Model Selection*

Model Evidence in terms of parameters

- Probability of generating the data from the model
- For model M_i governed by parameters w :
from sum and product rules

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

This is obtained by marginalizing over 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

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

- Insight into model evidence with one parameter

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}

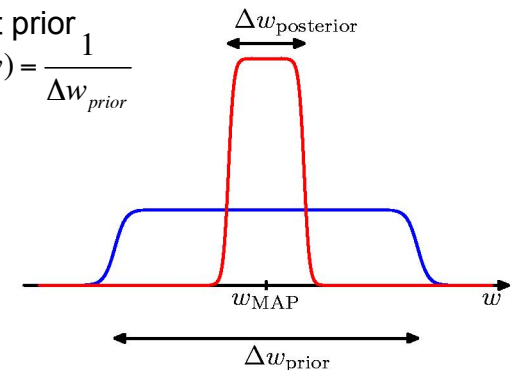
$$p(D) = \int p(D|w)p(w)dw \approx p(D|w_{MAP}) \frac{\Delta w_{posterior}}{\Delta w_{prior}}$$

Integral = maximum times width of peak
Flat prior
 $p(w) = \frac{1}{\Delta w_{prior}}$

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

Fit to data given most probable parameter values

Penalizes model according to its complexity since $\Delta w_{posterior} < \Delta w_{prior}$ and term is negative
If parameters are finely tuned to data this term is large



Optimal Model Complexity

- For a given model \mathcal{M}_i 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

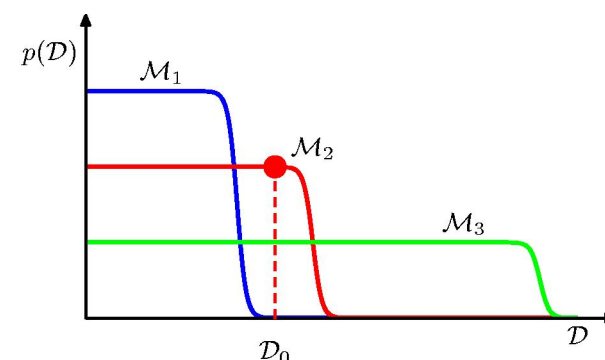
For observed data set D_0

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



x-axis is space of possible data sets

Bayes factor always favors correct model

- We are assuming the true model is contained among the M_i
- If this is so, Bayes model comparison will favor the correct model
- If M_1 is the true model, and we average over the distribution of data sets, Bayes factor has the form

$$\int p(D|M_1) \ln \frac{p(D|M_1)}{p(D|M_2)} dD$$

- This is K-L Divergence which is always positive
 - Thus Bayes factor favors the correct model

Summary of Bayesian Model Comparison

- Avoids problem of over-fitting
- Allows models to be compared using training data alone
- However needs to make assumptions about form of model
 - Model evidence can be sensitive to many aspects of prior such as behavior of tails
- In practice necessary to keep aside independent test data to evaluate performance