

# Model Selection

## General Strategies:

- Penalize
- Cross validation
- Model posterior
- Model Averaging.

## Per Issue:

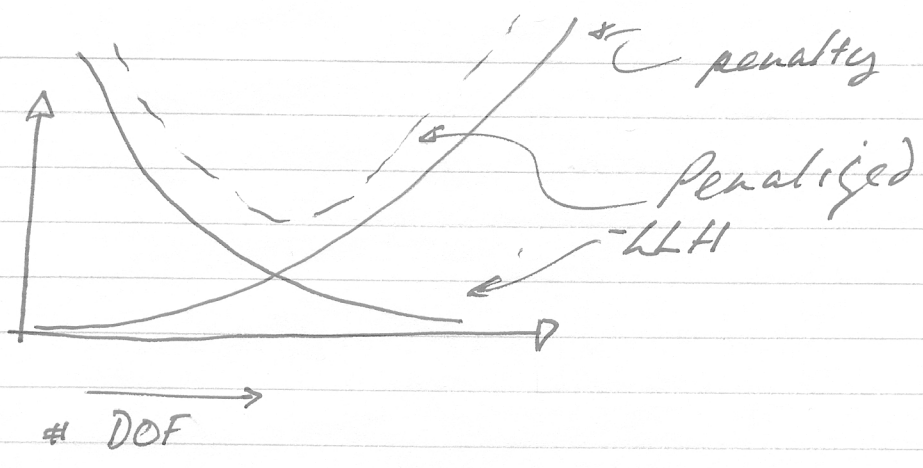
- Models with more DCF appear to fit training data better (but predict test data worse)
- Can't select on fitting LLH alone

## Penalization methods

- wish to select from  $M_1 \dots M_K$
- compare  $L_1(\hat{\theta}) + P(M_1)$

$$L_K(\hat{\theta}) + P(M_K)$$

Typically



Popular penalties:

AIC

$$-2 \log(L(\hat{\theta}^1(y))) + 2K$$

↑ # of parameters in model

~~BIC~~<sub>c</sub>

$$-2 \log(L(\hat{\theta}^1(y))) + 2K \left( \frac{n}{n-K-1} \right)$$

↘ # of samples

~~MDL~~  
BIC  
|||  
MDL

$$\neq -2 \log(L(\hat{\theta}^1(y))) + K \log(n)$$

Which one to use?

→ never AIC, unless  $n \gg K$

→ otherwise, opinion is divided

mechanics

- fit each of  $M_1, \dots, M_K$
- evaluate favored criterion
- smallest value wins.

Cross validation

• AIC, AIC<sub>c</sub> are attempts to estimate

$$KL(f, g(\theta)) = \int f(x) \log \left( \frac{f(x)}{g(x|\hat{\theta})} \right) dx$$

• But we might do this directly.

Assume we fit model to  $x_1 \dots x_n$ ,

now evaluate  $\frac{-1}{N-n} \sum_{i=n+1}^N \log p(x_i | \hat{\theta})$

$\approx - \int f \log p(x | \hat{\theta}) dx$

$= KL(f, p(x | \hat{\theta})) + H_f$

we don't know this, but it is shared for all models

∴ take model with smallest held out averaged LLH

→ waste of data:

Solu: Average over multiple VAR splits of  $x_i$

→ expensive computationally:

- buy a faster computer

# Model posteriors and Bayesian m.s.

typical Bayesian model, we will have

$M_1 \dots M_K \leftarrow$  Models

each will have  $\theta \leftarrow$  hyperparameters

which parametrize the priors of the parameters. E.g. GP, we had length scale, ~~ratio of~~ model var, ~~to~~ noise and parameters  $w \leftarrow$  eg. pars of linear regression.

We will select a model by

$$p(M_i | y, X) = \frac{p(y | X, M_i) p(M_i)}{p(y | X)}$$

Now  $p(y/x, M_i)$

$$= \int p(y/x, M_i, \theta) p(\theta/H_i) d\theta$$

~~and~~

~~$p(\theta)$~~

hyper priors

Marginal likelihood, or evidence

$$p(y/x, M_i, \theta) = \int p(y/x, M_i, \theta, \omega) p(\omega/\theta, H_i) d\omega$$

Problem: • all these integrals

- Strategies:
- get lucky, and have analytical
  - sampled ests
  - fix  $\theta = \hat{\theta}$  and work with  $p(y/x, M_i | y, x, \hat{\theta})$

## Model Averaging

e.g. • 
$$p(y_* / y, X, x_*) = \int p(y_* / y, X, x_*, w) p(w / y, X) dw$$

- In ~~a~~ some strong sense, best thing to do

$$p(y_* / y, X, x_*) = \sum_i p(y_* / y, X, x_*, M_i) p(M_i / y, X).$$

### BUT

- can be violently impractical
- can obscure comprehension

### PROS

- occasional, important practical examples where not averaging ~~obs~~ leads to false sense of security

