

# Model Selection

①

## Tools so far:

- t-test on individual coefficients
- F-test between particular subset models
- Backward / Forward search

Issues → selection bias

Ultimate validation of a model = prediction.  
Does the model generalize to future data?

Usually, parsimonious models / simple models work best for prediction. Why?

Easier to estimate coefficients of a simple model with a moderate sample size. Limit risk of including a spurious relationship that doesn't generalize.

Go for simplicity - also easier to interpret!  
Occam's razor

(2)

Flexible models	Simple models
- many parameters	- few parameters
- data adaptive	- rigid
- little (or no) bias	- possible bias
- high estimation variance	- smaller variance
- risk of overfitting	- risk of underfitting
bias ↓ variance ↑	bias ↑ variance ↓

Always — Need to check the overall fit of the model

A simple model won't cut it if

- patterns in residuals
- no better than an intercept model
- $R^2$  too small
- outliers, leverage, ... 5 basic assumptions

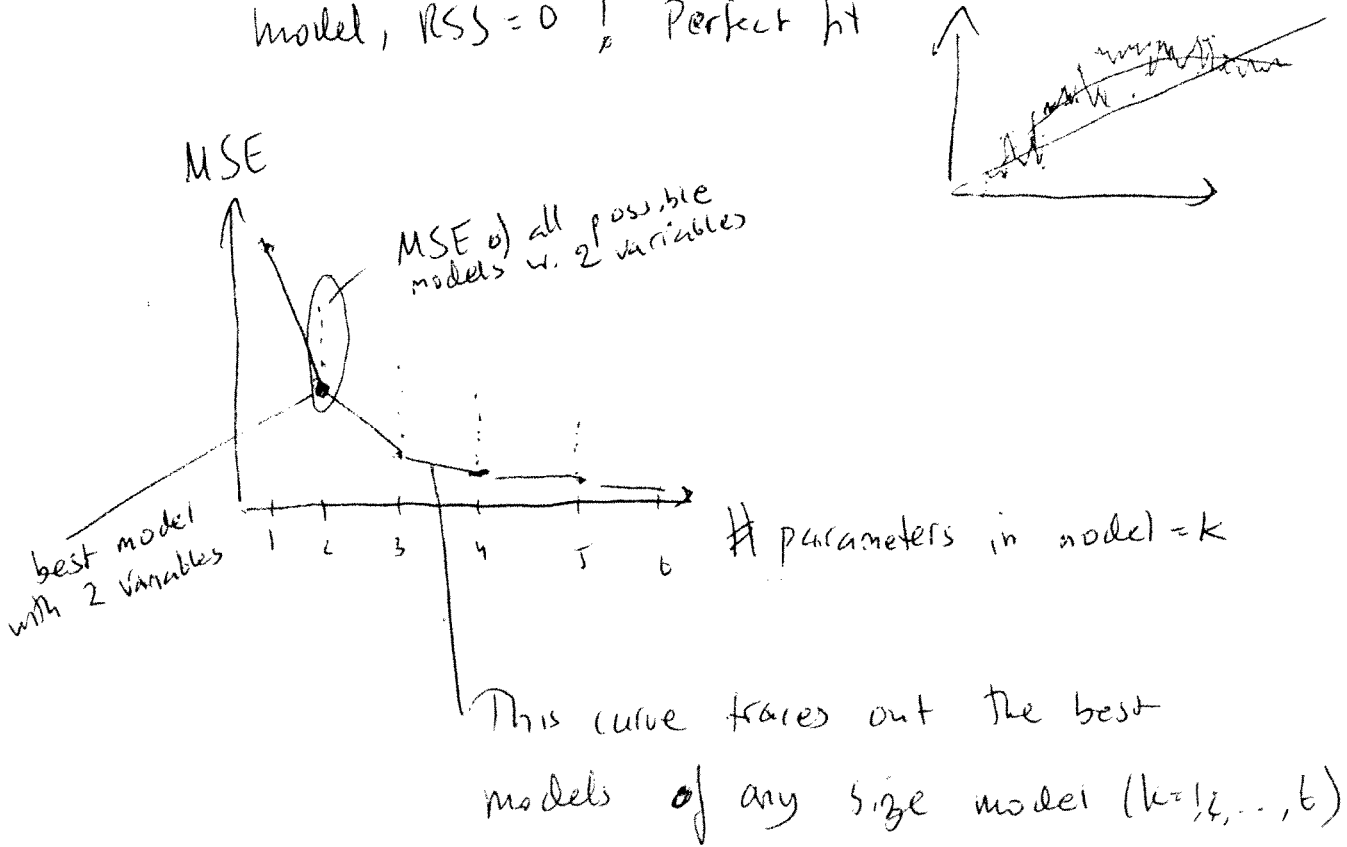
How do we identify a good prediction model?

RSS does not work

Why?  $RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ ,  $MSE = \frac{1}{n} RSS$

always decreases with the complexity of the model.

Note, in extreme case with n variables in the model,  $RSS = 0$  ! Perfect fit



So, if we used RSS to choose a model  $\Rightarrow$  we end up with the full / largest model.

④ If we had a separate test data set we could choose model that minimizes the prediction MSE

④

$$pMSE = \frac{1}{n} \sum_{i=1}^n (y_i^{new} - \hat{y}_i)^2$$

Strategy  $\left\{ \begin{array}{l} \text{Training data } (Tr) \quad (x_i, y_i)_{i=1}^n \\ \text{Test data } (Te) \quad (x_i, y_i^{new})_{i=1}^n \end{array} \right.$

In practice could be different  $x_i$ 's & sample sizes.

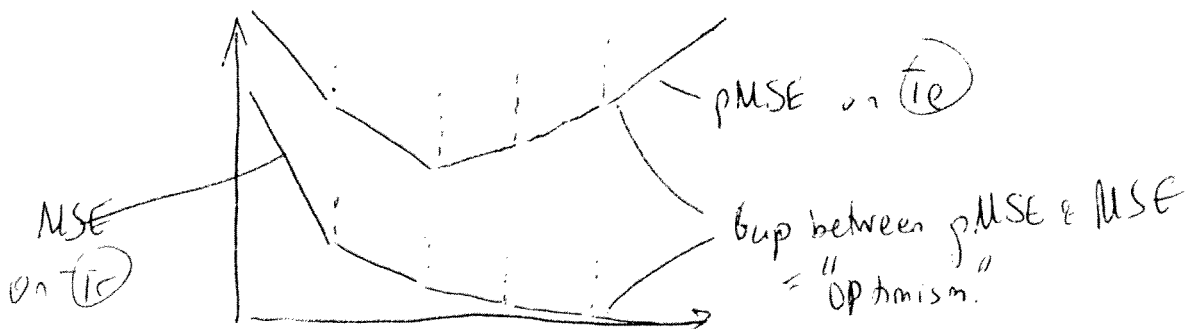
① On  $(Tr)$  - fit all models of interest;  $M_1, M_2, \dots, M_j$  with coefficients  $(\hat{\beta}_1^1, \hat{\beta}_1^2, \dots, \hat{\beta}_1^j)$

② On  $(Te)$  - compute  $pMSE^j$  Model  $M_j$  predictions

$$pMSE^j = \frac{1}{n} \sum_{i=1}^n (y_i^{new} - \hat{y}_i(M_j))^2$$

$$= \frac{1}{n} \sum_{i=1}^n (y_i^{new} - \sum_i \hat{\beta}_i^j)^2$$

③ Choose the model  $j^*$  that minimizes  $pMSE$



⑤

④  $\rho \text{MSE}^j - \text{MSE}^j$  = optimism for model  $j$ .

That is, because we fit the model on  $\mathcal{T}$

Using the LS criterion  $\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \text{RSS}$

So the  $\text{RSS}^j$  or  $\text{MSE}^j$  is an underestimate of how well the model would perform on future data.

Now, in practice we don't have an independent test data

⇒

Alternative 1

We use part of the data for training, part for testing

Cross-validation

Alternative 2

Estimate the expected size of the gap  $\rho \text{MSE}^j - \text{MSE}^j$

- or at least an approximation of the gap.

Model selection criteria

like AIC, BIC

Let's take a closer look at the gap

(6)

$$\text{Prediction error of model } M_j = PE^j = E_{\text{new}} \left( \frac{1}{n} \sum_{i=1}^n (y_i^{\text{new}} - \hat{y}_i(M_j))^2 \right)$$

$$= \frac{1}{n} \sum_{i=1}^n E_{\text{new}} \left( y_i - \hat{y}_i(M_j) \right)^2$$

prediction using model  $M_j$

where coefficients are estimated from training data  $(x_i, y_i)_{i=1}^n$

$$\text{Notation } \hat{y}_i(M_j) = \hat{y}_i^j$$

Expand  $\textcircled{*}$  into 6 terms, adding & subtracting  $E(y_i^{\text{new}})$  &  $E(\hat{y}_i^j)$

$$\textcircled{*} = E \left( y_i^{\text{new}} - E(y_i^{\text{new}}) + E(y_i^{\text{new}}) - E(\hat{y}_i^j) + E(\hat{y}_i^j) - \hat{y}_i^j \right)^2$$

$$= E \left( y_i^{\text{new}} - E(y_i^{\text{new}}) \right)^2 + E \left( E(y_i^{\text{new}}) - E(\hat{y}_i^j) \right)^2 + E \left( E(\hat{y}_i^j) - \hat{y}_i^j \right)^2$$

$\textcircled{1} \qquad \qquad \qquad \textcircled{2} \qquad \qquad \qquad \textcircled{3}$

$$+ 2E \left( (y_i^{\text{new}} - E(y_i^{\text{new}})) (E(y_i^{\text{new}}) - E(\hat{y}_i^j)) \right) + 2E \left( (y_i^{\text{new}} - E(y_i^{\text{new}})) (E(\hat{y}_i^j) - \hat{y}_i^j) \right)$$

$\textcircled{4} \qquad \qquad \qquad \textcircled{5}$

$$+ 2E \left( (E(y_i^{\text{new}}) - E(\hat{y}_i^j)) (E(\hat{y}_i^j) - \hat{y}_i^j) \right)$$

$\textcircled{6}$

# Meaning of terms ①-⑥

⑦

①  $E(y_i^{new} - E(y_i^{new}))^2 = \sigma^2$

- The irreducible error  
- noise around true model

②  $E(E(y_i^{new}) - E(\hat{y}_i^j))^2 = \text{bias}^2$  of model  $j$  at obs.  $i$

true model

- if model  $j$  is adequate the bias is 0

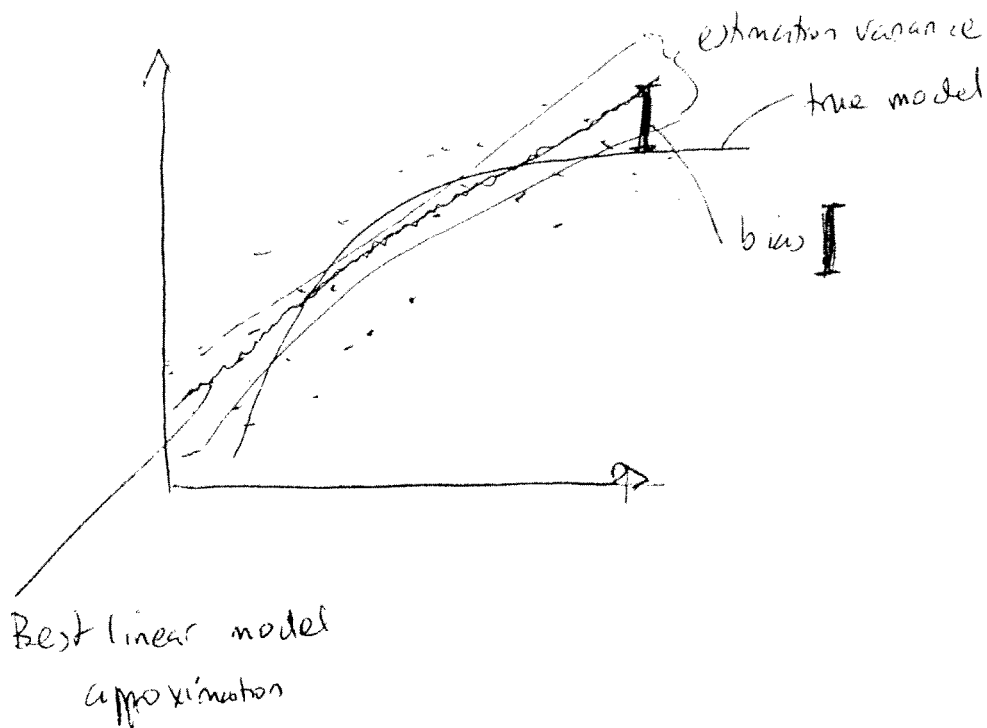
③  $E(E(\hat{y}_i^j) - \hat{y}_i^j)^2 = V(\hat{y}_i^j)$  - estimation variance of model  $j$

- increases with complexity of model.

④ = 0  
⑥ = 0 } can pull constant outside

⑤ = 0 since  $y_i^{new}$  &  $y_i$  are independent data sets

So  $PE^j = \frac{1}{n} \sum_{i=1}^n (\sigma^2 + \text{bias}^2(i, M_j) + V(\hat{y}_i^j))$   
 $= \sigma^2 + \overline{\text{bias}^2(M_j)} + \frac{1}{n} \sum_{i=1}^n V(\hat{y}_i^j)$



Example: LS estimate

$$\hat{\beta} = (X'X)^{-1} X'y$$

-  $X = n \times j$  design matrix

- a model with  $j$  variables

if model  $j$  is adequate  $\rightarrow$  no bias

$$V(\hat{y}) = V(Hy) = \sigma^2 H$$

$$V(\hat{y}_i) = \sigma^2 h_{ii} \rightarrow \frac{1}{n} \sum V(\hat{y}_i) = \frac{\sigma^2}{n} \sum h_{ii}$$

$$\begin{aligned} \text{Now, } H &= X(X'X)^{-1}X' \\ &= \frac{\sigma^2}{n} \text{Tr}\{H\} \end{aligned}$$

$$\text{Tr}\{H\} = \text{Tr}\{X(X'X)^{-1}X'\} = \text{Tr}\{(X'X)(X'X)^{-1}\} = \text{Tr}\{I_j\} = j$$

$$\text{So } PEJ = \sigma^2 \left(1 + \frac{j}{n}\right) \text{ so increases with } j!$$

# parameters in model  $j$



What about The training error?

(9)

$$TE^j = E \left( \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i^j)^2 \right) = \frac{1}{n} \sum_{i=1}^n E (y_i - \hat{y}_i^j)^2$$

Like before, we expand  $E(y_i - \hat{y}_i^j)^2$  into 6 terms,

but now (5) is not 0

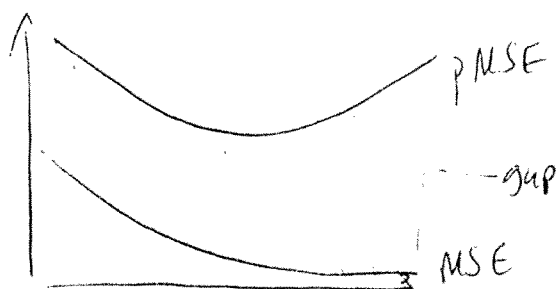
$$\begin{aligned} (5) &= 2E((y_i - E(y_i)) (E(\hat{y}_i^j) - \hat{y}_i^j)) = 2E(y_i (E(\hat{y}_i^j) - \hat{y}_i^j)) \\ &= 2E(y_i)E(\hat{y}_i^j) - 2E(y_i \hat{y}_i^j) = -2 \text{Cov}(y_i, \hat{y}_i^j) \end{aligned}$$

$$(6) \quad TE^j = PE^j - \frac{2}{n} \sum \text{Cov}(y_i, \hat{y}_i^j)$$

What did we learn? Well,  $TE < PE$  always

(and)  $TE$  much smaller than  $PE$  if  $y_i$  &  $\hat{y}_i$  are highly correlated, which happens for complex models.

That is, the gap between the  $MSE$  &  $pMSE$  increases with the size of the model.



Example For LS fits  $\hat{y} = Ay$  &  $\text{cov}(y, \hat{y}_i) = \sigma^2 h_{ii}$  (10)

$$\frac{2}{n} \sum \text{cov}(y, \hat{y}_i) = 2 \frac{\sigma^2}{n} \sum h_{ii} \quad \# \text{ of parameters in model } j.$$

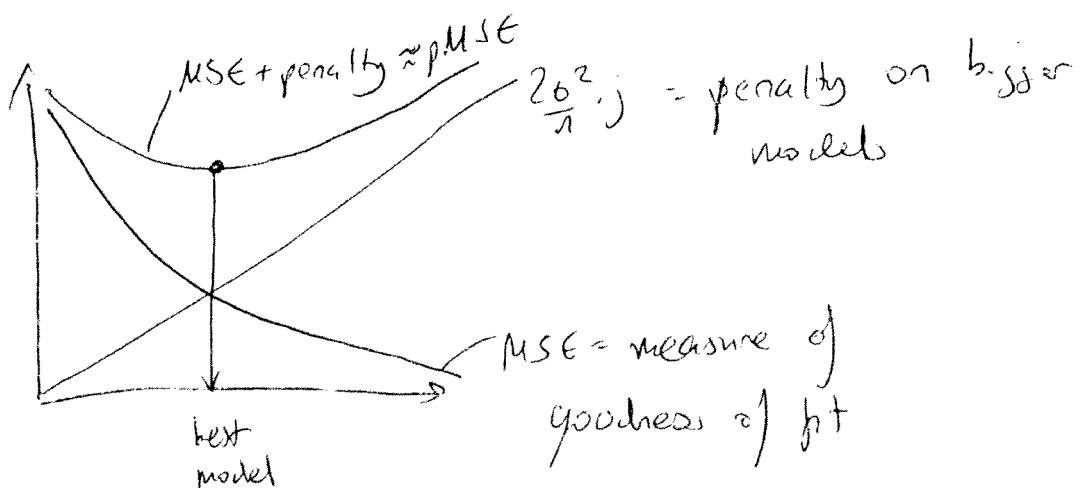
i.e. gap grows linearly with model size.

What's the point?

Alternative 2 • We don't have an independent test set so we can't estimate pMSE, but we need it for selection a good prediction model

• but we have MSE for each model and we know the expected size of the gap between pMSE & MSE  $\Rightarrow 2 \frac{\sigma^2}{n} \cdot j$

• Approximate pMSE<sup>j</sup> by  $\boxed{\text{MSE}^j + 2 \frac{\sigma^2}{n} \cdot j}$



Many model selection criteria have this form

II

Goodness-of-fit measure + penalty on model size

The one based on the expected gap size

$$MSE + \frac{2\hat{\sigma}^2}{n} \times (\text{number of parameters})$$

is called Mallow's  $C_p$ , which you also see

as  $RSS + 2\hat{\sigma}^2 \cdot (\text{number of parameters})$  in the texts

To use we need to plug in  $\hat{\sigma}^2$  in the penalty. We use the same  $\hat{\sigma}^2$  for all models — the  $\hat{\sigma}^2$  from an unbiased model (usually the largest model we fit)

- Strategy
- ① Enumerate all models  $M_1, \dots, M_J$  ( $j=1, \dots, J$ )
  - ② Evaluate  $MSE_j, RSS_j$  for all models
  - ③ Compute  $C_p^j = RSS_j + 2\hat{\sigma}^2 \cdot j$ , where  $\hat{\sigma}^2 = MSE^J$  (largest model)
  - ④ Pick model with smallest  $C_p$ .

Other model selection criteria

$$\begin{cases} AIC(M_j) = n \log(RSS_j) + 2 \cdot j \\ BIC(M_j) = n \cdot \log(RSS_j) + j \cdot \log(n) \end{cases}$$

Goodness  
of fit  
measure

penalty  
on model size

Note, AIC & BIC measure goodness of fit on the likelihood scale (normal error distribution)

Also, BIC penalty is larger than AIC or Cp so BIC tends to pick smaller models.

In practice, AIC & Cp tend to overfit the data.

A note on PE;  $PE = \underbrace{\sigma^2}_{\text{irreducible}} + \overbrace{\text{bias}^2 + \text{Var}(\text{model estimate})}^{\text{trade-off}}$

irreducible

trade-off

Perhaps better to sacrifice some bias if we can lower variance?

→ Ridge regression ...