

Model Selection

10

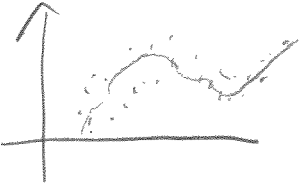

Tools so far:

- t-test on individual coefficients β_0
- F-test Comparing 2 models
- Which models? \rightarrow Backward search

Ultimate validation tool = PREDICTION

Does the model generalize to future data?

Usually it is safer to use as simple models as possible for prediction. Why? Complex models are more difficult to estimate (especially if you have a small sample size).

Flexible models	Simple models
<ul style="list-style-type: none">- Many parameters <p>\Downarrow</p> <ul style="list-style-type: none">- adapts to data	<ul style="list-style-type: none">- few parameters <p>\Downarrow</p> <ul style="list-style-type: none">- rigid
	
<ul style="list-style-type: none">- little or no bias (since flexible)- high estimation variance (small changes to data \rightarrow large impact)	<ul style="list-style-type: none">- possible bias (since rigid)- smaller variance (since rigid, can't change much)

} bias = average across many data sets, how much model deviates from true relationship

To summarize:

(2)

Bias - Variance trade-off

$$E(\hat{y}) - E(y)$$

↑
Average fitted value
across many data sets

↑
True model value

Average fitted value

across many data sets

If the ^{fitted} model is incorrect, can't attain true value

flexible models - bias ↓

simple models - bias ↑

$$V(\hat{y})$$

↑
a function of $\hat{\beta}$
since $\hat{y} = X\hat{\beta}$

The more complex model is, the more uncertainty in $\hat{\beta}$ → and therefore \hat{y}

flexible models - variance ↑

simple models - variance ↓

We want to control both bias and variance.

Just controlling bias is not enough. Yes, on average (across many data sets) we do well, but if variance is high we can be far from this average in one particular data set and we don't know when.

(combining bias & variance ⇒ $MSE = BIAS^2 + VARIANCE$)

So, we should choose a model that minimizes

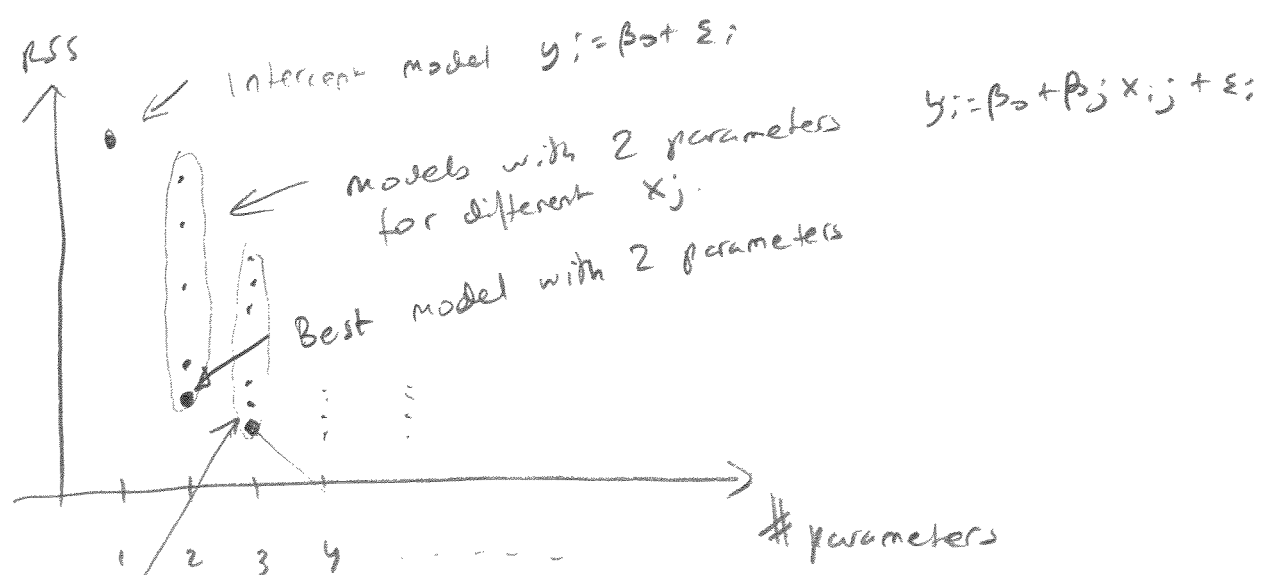
$$MSE = BIAS^2 + VARIANCE.$$

How? To compute the bias we need to know the true model !!!

Can we just use the residual sum of squares (RSS)?

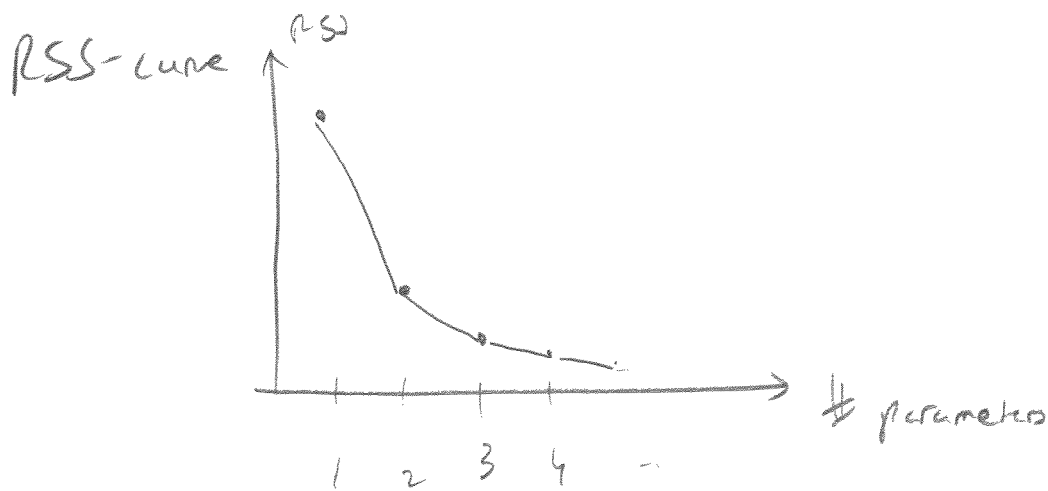
NO! $RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ always decreases the more complex the model is (more parameters to help match the model to the data).

In fact, if you use n parameters you get a perfect fit, $RSS=0$!



models with 3 parameters
 $y_i = \beta_0 + \beta_j x_j + \beta_k x_k + \epsilon_i$
for all combinations of x_j and x_k

Best model with 3 parameters

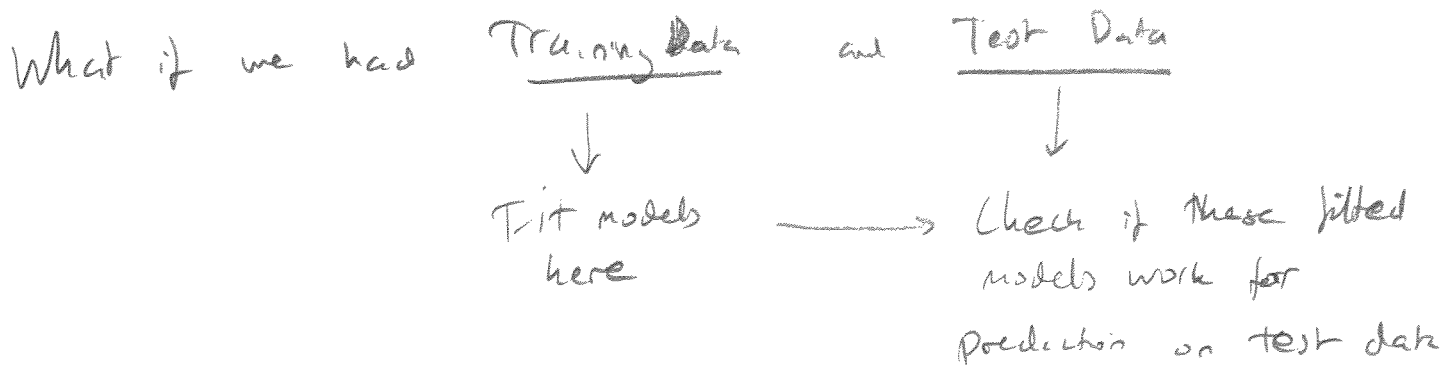


RSS-curve - curve that connects the best model of each size (i.e. best one-variable model, best two-variable model, etc.)

So RSS is minimized for the largest model.

RSS is not a good substitute for the prediction MSE.

A hypothetical situation



Setup Training Data $(\tilde{x}_i, y_i)_{i=1}^n$, $\tilde{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip-1})$

Test Data $(\tilde{x}_i, y_i^{new})_{i=1}^n$, $\tilde{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip-1})$

That is, assume we get to observed new y -values
at same X -locations

(5)

True model $\beta^* = (\beta_0^*, \beta_1^*, \beta_2^*, \dots, \beta_{p-1}^*)$

Some of these are 0

$$\left\{ \begin{array}{l} \text{Training data: } y_i = x_i \beta^* + \varepsilon_i, \quad \varepsilon_i \text{ drawn from } N(0, \sigma^2) \\ \text{Test data: } y_i^{\text{new}} = x_i \beta^* + \varepsilon_i^{\text{new}}, \quad \varepsilon_i^{\text{new}} \dots \dots N(0, \sigma^2) \end{array} \right.$$

Definition

$$\left\{ \begin{array}{l} \text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ \text{MSE}_{\text{train}} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (\text{training mean-squared-error}) \end{array} \right.$$

$\hat{y}_i = x_i \hat{\beta}$ ← LS estimate from training data

$$\text{MSE}_{\text{test}} = \frac{1}{n} \sum_{i=1}^n (y_i^{\text{new}} - \hat{y}_i)^2 \quad (\text{test or prediction MSE})$$

$\hat{y}_i = x_i \hat{\beta}$ ← LS estimate from training data

If this is small, the model generalizes to test data.

So no refitting of model using test data.

Strategy

(6)

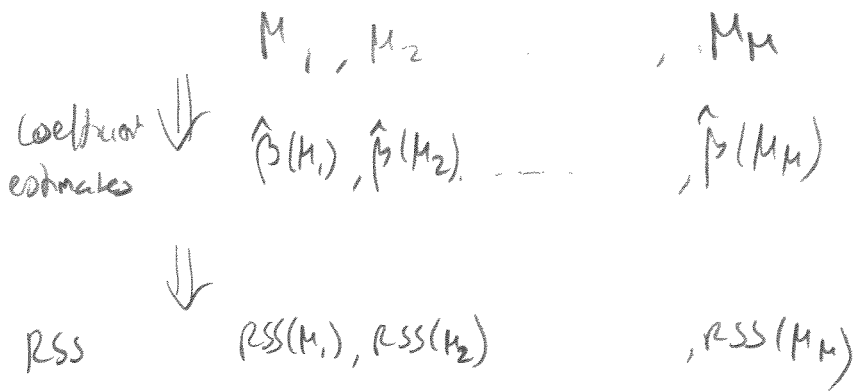
① Enumerate all models of interest (maybe all possible models)

$$M_1, M_2, \dots, M_M$$

M_j (model j) is a particular combination of variables

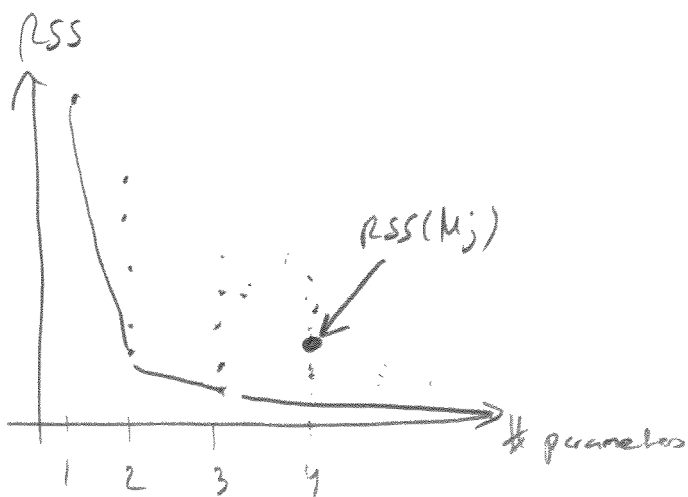
(Example $M_j = \{X_1, X_5, X_{10}\}$ - a 3-variable model involving variable 1, 5 and 10. (4 parameters))

② On the training data, fit all M models to the data using least squares



Example

$$\hat{\beta}(M_j) = \begin{pmatrix} \hat{\beta}_1(M_j) \\ 0 \\ 0 \\ 0 \\ \hat{\beta}_5(M_j) \\ 0 \\ 0 \\ 0 \\ 0 \\ \hat{\beta}_{10}(M_j) \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$



③ On the test data, evaluate the prediction MSE for all models

⑦

M_1, M_2

M_M

$\hat{\beta}(M_1), \hat{\beta}(M_2)$

$\hat{\beta}(M_M)$

← from training data



$PMSE(M_1), PMSE(M_2)$

$PMSE(M_M)$

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

↑ fitted values for model M_j

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

↑ coefficient estimates from training data

④

