

# Prediction Error: The Bias-Variance Trade-Off

Aldo Solari



# Mean squared error

- The mean squared error for the training data

$$\text{MSE}_{\text{Tr}} = \frac{1}{n} \sum_{i=1}^n [y_i - \hat{f}(x_i)]^2$$

is not a good measure of performance

- We would like to have a good performance

$$\text{MSE}_{\text{Te}} = \frac{1}{m} \sum_{i=1}^m [y_i^* - \hat{f}(x_i^*)]^2$$

on the test data

$$(x_1^*, y_1^*), (x_2^*, y_2^*), \dots, (x_m^*, y_m^*)$$



```
load("poly.Rdata")

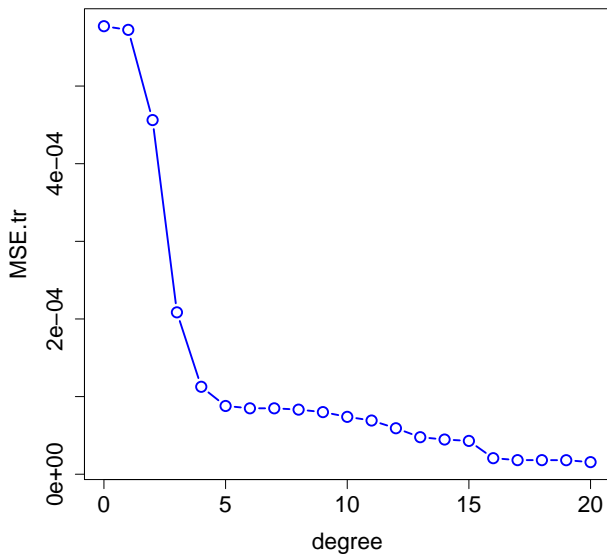
n <- nrow(train)
ds = 0:20

fun <- function(d) if (d==0) lm(y~1, train)
                        else lm(y~poly(x,degree=d), train)
fits <- sapply(ds, fun)

MSEs.tr <- unlist( lapply(fits, deviance) )/n
plot(ds, MSEs.tr, type="b")
```



# Training data MSE

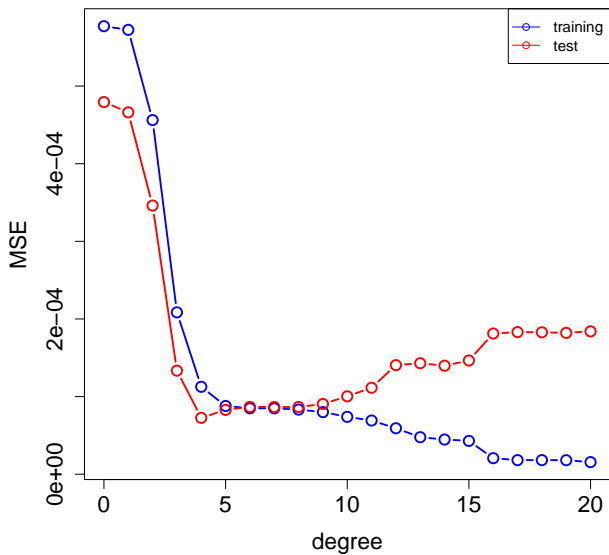


```
yhats <- lapply(fits, predict)
MSEs.te <- unlist( lapply(yhats,
                          function(yhat) mean((test$y - yhat)^2)) )

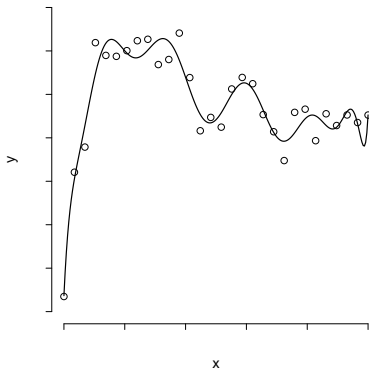
plot(ds, MSEs.tr, type="b", col=4)
lines(ds, MSEs.te, type="b", col=2)
```



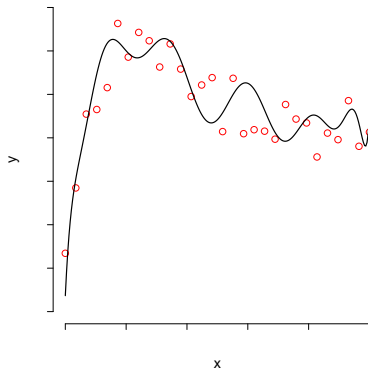
# Test data MSE



# Overfitting



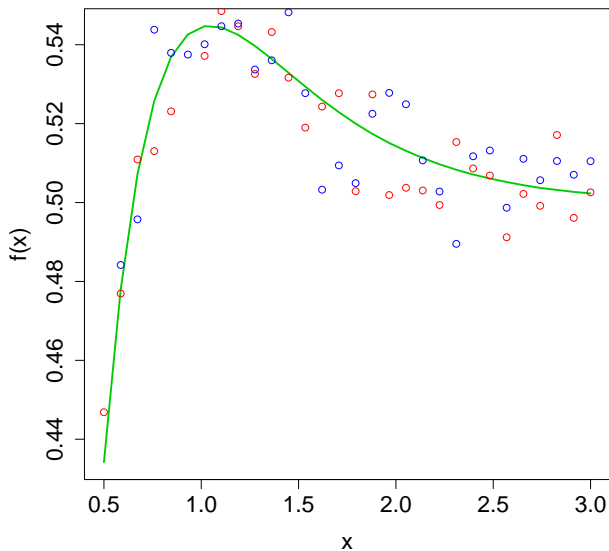
Yesterday data



Tomorrow data

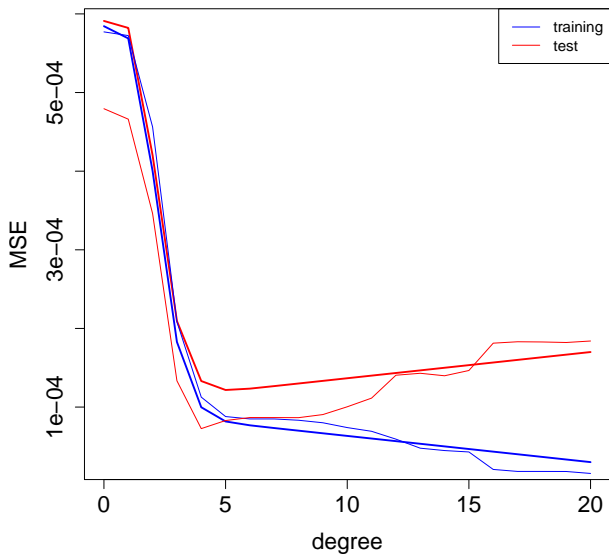


# True function

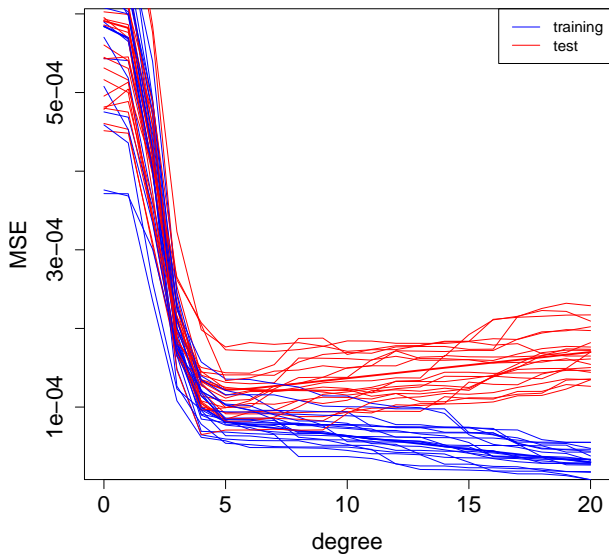




# Expected MSE



# Expected MSE



# Outline

## ① The Bias-Variance Decomposition

## ② How to Avoid Overfitting



# Regression function

- $(X, Y)$  have some unknown joint distribution
- We want to predict  $Y$  from  $X$
- Over all functions  $f$ , the expected test error, measured in terms of squared error loss,

$$\mathbb{E}[(Y - f(X))^2] \tag{1}$$

is minimized at

$$f(x) = \mathbb{E}(Y|X = x) \tag{2}$$

- $f(x)$  is called the (true) regression function
- **Homework:** prove that (1) is minimized at (2)



# Example of regression function

- Joint distribution:

$$\begin{pmatrix} Y \\ X \end{pmatrix} \sim \mathcal{N} \left[ \begin{pmatrix} \mu_y \\ \mu_x \end{pmatrix}, \begin{pmatrix} \sigma_y^2 & \rho\sigma_x\sigma_y \\ \rho\sigma_x\sigma_y & \sigma_x^2 \end{pmatrix} \right]$$

- Conditional distribution:

$$(Y|X = x) \sim \mathcal{N} \left( \mu_y + \rho \frac{\sigma_y}{\sigma_x} (x - \mu_x), \sigma_y^2 (1 - \rho^2) \right)$$

- Regression function:

$$f(x) = E(Y|X = x) = \left( \mu_y - \rho \frac{\sigma_y}{\sigma_x} \mu_x \right) + \left( \rho \frac{\sigma_y}{\sigma_x} \right) x$$



# The general model

$$Y = f(X) + \varepsilon \quad (3)$$

- $Y$  : continuous response
- $X = (X_1, \dots, X_p)^T$  : predictors  
 $p \times 1$
- $f(X) = \mathbb{E}(Y|X)$  : regression function
- $\varepsilon$  : error term, assumed with

$$\mathbb{E}(\varepsilon) = 0, \quad \text{Var}(\varepsilon) = \sigma^2$$

and independent of  $X$



# Expected test error

- The true  $f(x) = \mathbb{E}(Y|X = x)$  is unknown
- We can use the training data

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

to get an estimate  $\hat{f}(x)$  of  $f(x)$

- The expected test error of  $\hat{f}$  is

$$\mathbb{E}[(Y - \hat{f}(X))^2] \tag{4}$$

where the expectation is over all that is random, namely, the training data  $(X_i, Y_i), i = 1, \dots, n$ , and the test point  $(X, Y)$



# Sources of error

## Irreducible error

Can we ever predict  $Y$  from  $X$  with zero error? No. Even the true regression function  $f$  cannot do this:  $\mathbb{E}[(Y - f(X))^2] = \sigma^2$

## Estimation bias

What happens if our fitted function  $\hat{f}$  belongs to a model class that is far from the true  $f$ ? E.g. we choose to fit a linear model in a setting where the true relationship is far from linear?

## Estimation variance

What happens if our fitted (random) function  $\hat{f}$  is itself quite variable? In other words, over different copies of the training data, we end up constructing substantially different functions  $\hat{f}$ ?

Source: Tibshirani R. (2015) Statistical Machine Learning





# The bias-variance decomposition

The expected test error, conditional on  $X = x$ , decomposes into

$$\mathbb{E}[(Y - \hat{f}(x))^2 | X = x] = \sigma^2 + \mathbb{E}[(f(x) - \hat{f}(x))^2] \quad (5)$$

where  $\sigma^2$  is the irreducible error. The reducible error

$$\begin{aligned} \mathbb{E}[(f(x) - \hat{f}(x))^2] &= (f(x) - \mathbb{E}(\hat{f}(x)))^2 + \mathbb{E}[(\hat{f}(x) - \mathbb{E}(\hat{f}(x)))^2] \\ &= \left\{ \text{Bias}(\hat{f}(x)) \right\}^2 + \text{Var}(\hat{f}(x)) \end{aligned} \quad (6)$$

decomposes into the squared (estimation) bias and the (estimation) variance

**Homework:** prove (5) and (6)



# The bias-variance decomposition

- Unconditionally over  $X$ , the bias-variance decomposition is

$$\mathbb{E}[(Y - \hat{f}(X))^2] = \sigma^2 + \int \text{Bias}^2(\hat{f}(x)) P_X(dx) + \int \text{Var}(\hat{f}(x)) P_X(dx)$$

where  $P_X$  is the distribution of  $X$

- Note that

$$\begin{aligned} \mathbb{E}(\text{MSE}_{\text{Te}}) &= \sigma^2 + \frac{1}{m} \sum_{i=1}^m \text{Bias}^2(\hat{f}(x_i^*)) + \frac{1}{m} \sum_{i=1}^m \text{Var}(\hat{f}(x_i^*)) \\ &\approx \mathbb{E}[(Y - \hat{f}(X))^2] \end{aligned}$$

with equality if  $X_i^* = X_i$  are constants ( $X$  is not random)

- In words, the expected test error is the irreducible error + average (squared) bias + average variance



# The bias-variance trade-off

$$\text{Reducible Error} = \text{Bias}^2 + \text{Variance}$$

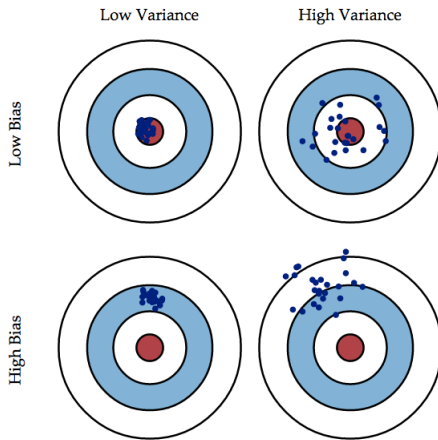
- Models  $\hat{f}$  with low bias tend to have high variance
- Models  $\hat{f}$  with low variance tend to have high bias

We can see that even if our prediction is unbiased, i.e.  $\mathbb{E}(\hat{f}(x)) = f(x)$ , we can still incur a large error if it is highly variable. On the other hand,  $\hat{f}(x) \equiv 0$  has 0 variance but will be terribly biased

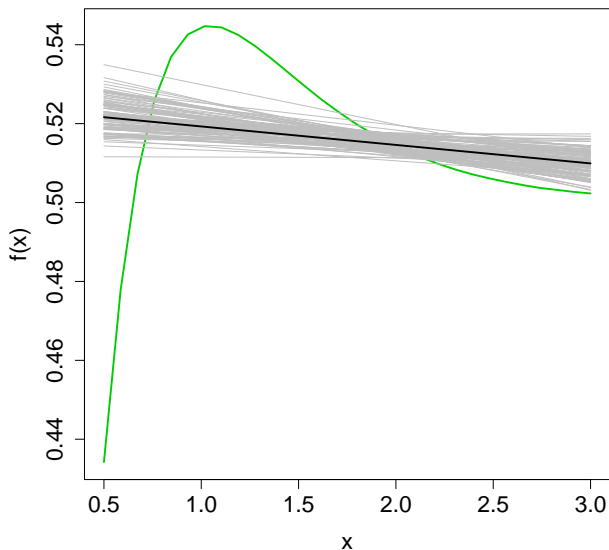
To predict well, we need to balance the bias and the variance



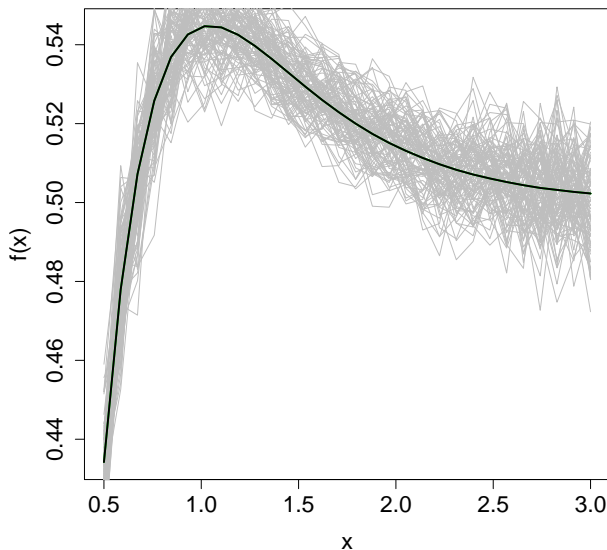
# The bias-variance trade-off



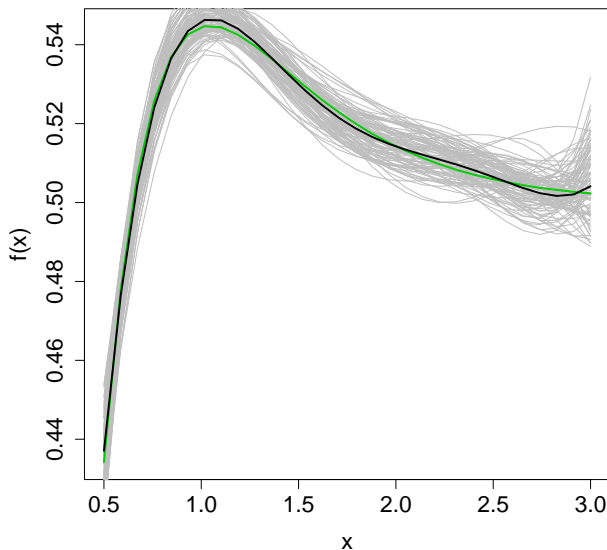
$\hat{f}(x)$  = regression line



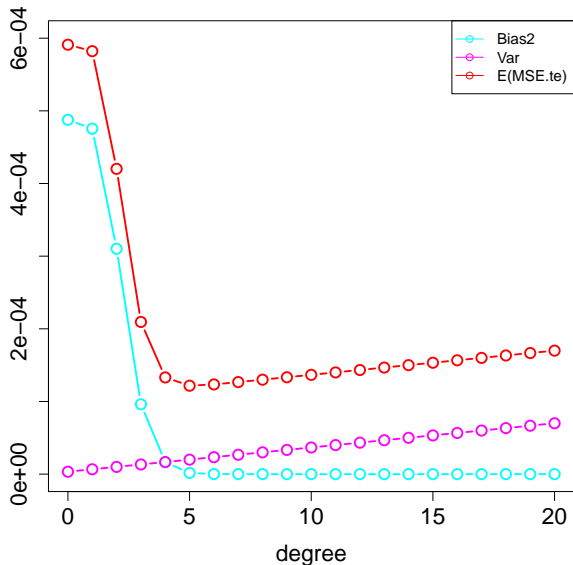
$\hat{f}(x)$  = high-degree polynomial



$\hat{f}(x)$  = **optimal-degree polynomial**

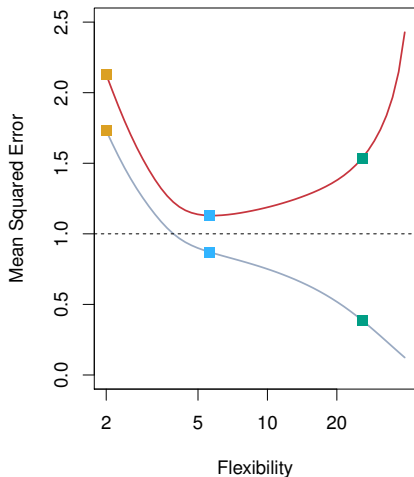
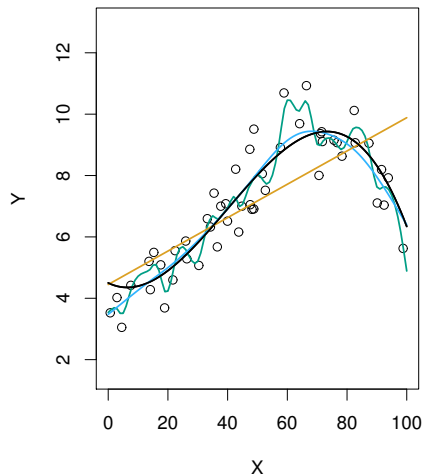


# The bias-variance decomposition





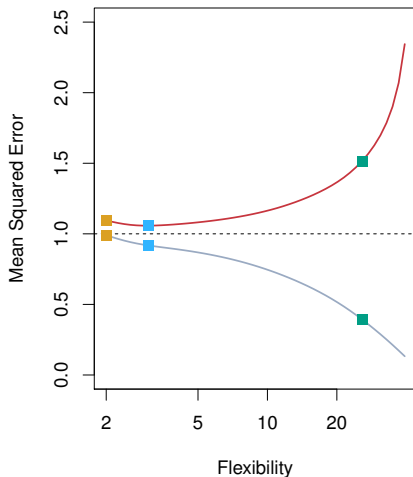
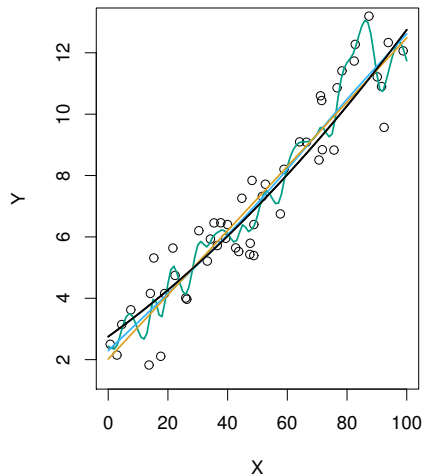
# The bias-variance decomposition



Source: Gareth et al. (2013) p. 31



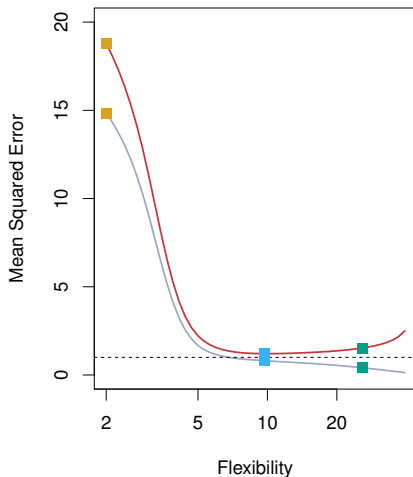
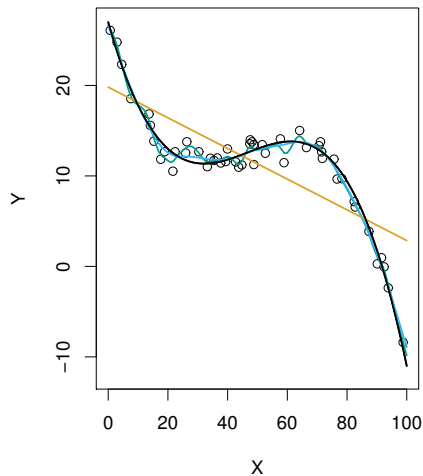
# The bias-variance decomposition



Source: Gareth et al. (2013) p. 33



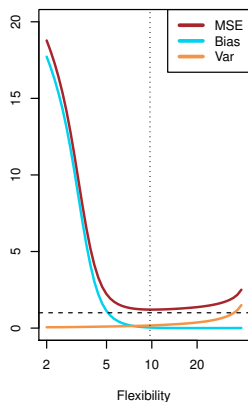
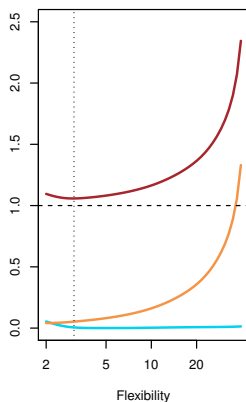
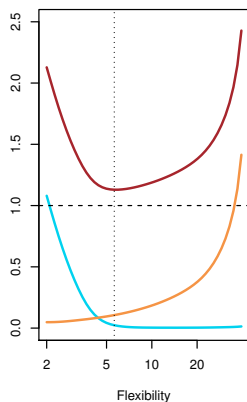
# The bias-variance decomposition



Source: Gareth et al. (2013) p. 34)



# The bias-variance decomposition



Source: Gareth et al. (2013) p. 36



# Outline

① The Bias-Variance Decomposition

② How to Avoid Overfitting



# Expected optimism

- For simplicity, assume that  $x_i = x_i^*$ ,  $i = 1, \dots, n$  are fixed constants ( $X$  is not random)
- The expected optimism for the linear model is

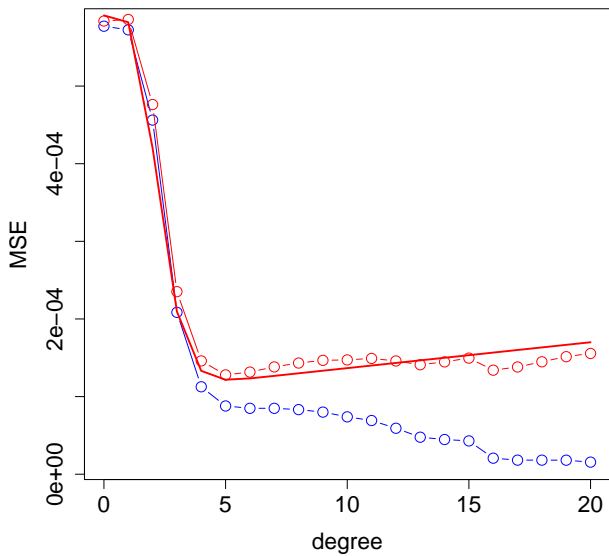
$$\mathbb{E}(\text{Optimism}) = \mathbb{E}(\text{MSE}_{\text{Te}}) - \mathbb{E}(\text{MSE}_{\text{Tr}}) = \frac{2\sigma^2 p}{n}$$

- We can estimate  $\mathbb{E}(\text{MSE}_{\text{Te}})$  as

$$\begin{aligned}\mathbb{E}(\widehat{\text{MSE}}_{\text{Te}}) &= \mathbb{E}(\widehat{\text{MSE}}_{\text{Tr}}) + \frac{2\sigma^2 p}{n} \\ &= \text{MSE}_{\text{Tr}} + \frac{2\sigma^2 p}{n}\end{aligned}$$



$\sigma^2$  known



# $\sigma^2$ unknown: AIC and BIC

- $\sigma^2$  is usually estimated by  $\hat{\sigma}^2 = \frac{n\text{MSE}_{\text{Tr}}}{n-p}$
- Finding the minimum of  $\text{MSE}_{\text{Tr}} + \frac{2\hat{\sigma}^2 p}{n}$  is equivalent to finding the minimum value of the Akaike Information Criterion

$$\text{AIC} = -2\ell(\hat{\beta}, \hat{\sigma}^2) + 2p$$

where for the linear model  $-2\ell(\hat{\beta}, \hat{\sigma}^2) = n \log(\text{MSE}_{\text{Tr}})$

- A different penalty for complexity:

$$\text{BIC} = -2\ell(\hat{\beta}, \hat{\sigma}^2) + \log(n)p$$

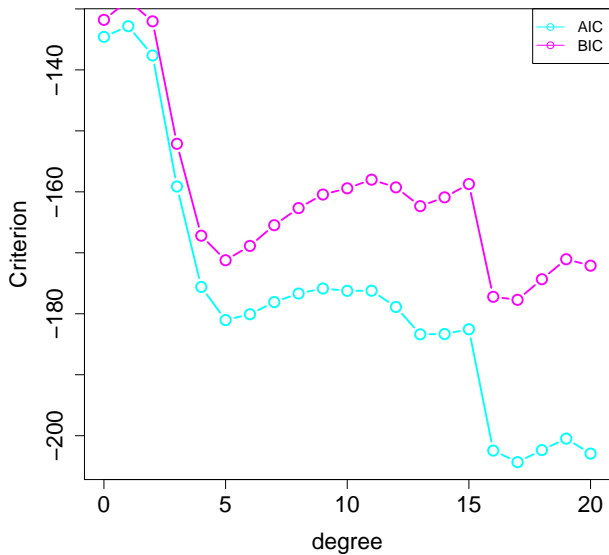




```
AICs <- unlist( lapply(fits, AIC) )  
BICs <- unlist( lapply(fits, BIC) )  
  
plot(ds, AICs, type="b", col=5)  
lines(ds, BICs, type="b", col=6)
```



# AIC and BIC



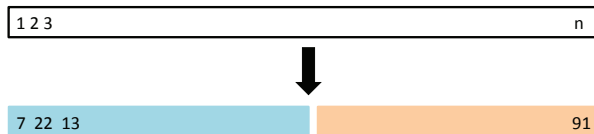
# Cross-validation

- Cross-validation is a method for estimating  $\mathbb{E}(\text{MSE}_{\text{Te}})$
- The idea is to hold out a subset of the training observations from the fitting process, and then applying the model to those held out observations
- Cross-validation is a non-parametric method, i.e. it applies to any model



# Validation-set approach

- A simple approach is to randomly divide the  $n$  observations into two parts: a training set and a validation or hold-out set



- The model is fitted on the training set  $\mathcal{T} \subset \{1, \dots, n\}$ , and the fitted model  $\hat{f}^{\mathcal{T}}$  is used to predict the responses for the observations in the validation set  $\mathcal{V} = \{1, \dots, n\} \setminus \mathcal{T}$ , providing an estimate of the test MSE

$$\mathbb{E}(\widehat{\text{MSE}}_{\text{Te}}) = \frac{1}{\#\mathcal{V}} \sum_{i \in \mathcal{V}} \left( y_i - \hat{f}^{\mathcal{T}}(x_i) \right)^2$$



# Cross-validation

- This scheme reduces the sample size used for fitting the model, but this is not a problem when  $n$  is very large
- If  $n$  is not very large, however, the validation estimate of the test error can be highly variable, depending on precisely which observations are included in the training set and which observations are included in the validation set
- One way of partially overcoming this arbitrariness is to split the data into equal parts  $\mathcal{F}_1, \dots, \mathcal{F}_K$  and then use  $i \in \mathcal{F}_k$  for training the model and  $i \notin \mathcal{F}_k$  for evaluating it:

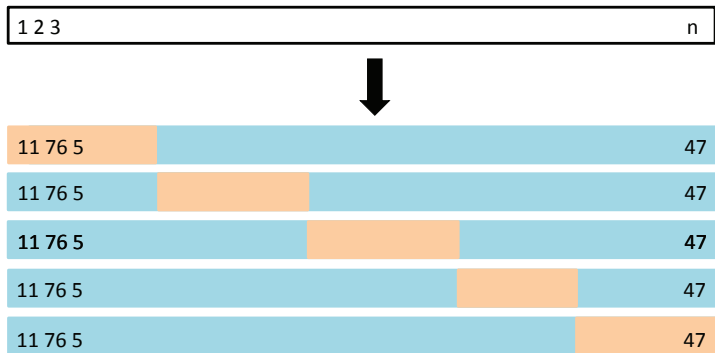
$$\frac{1}{\#\mathcal{F}_k} \sum_{i \notin \mathcal{F}_k} \left( y_i - \hat{f}^{\mathcal{F}_k}(x_i) \right)^2, \quad k = 1, \dots, K$$

- Then take the average

$$\mathbb{E}(\widehat{\text{MSE}}_{\text{Te}}) = \frac{1}{K} \sum_{k=1}^K \left[ \frac{1}{\#\mathcal{F}_k} \sum_{i \notin \mathcal{F}_k} \left( y_i - \hat{f}^{\mathcal{F}_k}(x_i) \right)^2 \right]$$



# K-fold cross-validation



Source: Gareth et al. (2013), p. 179

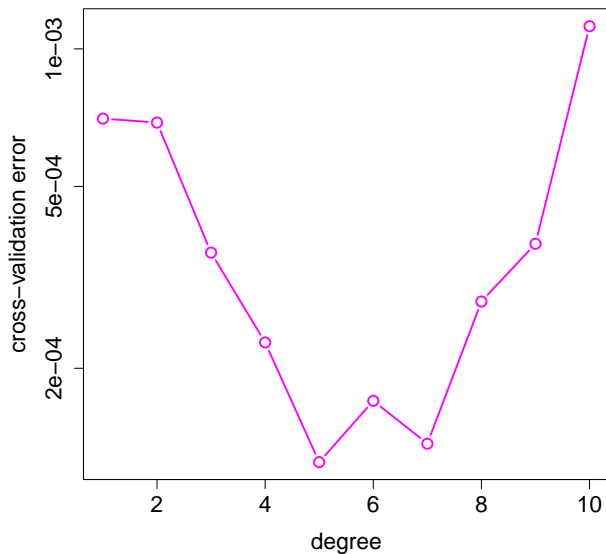


```
set.seed(123)
KCV = sapply(1:10, function(d)
  cv.glm(train, glm(y~poly(x,degree=d),
    train, family = gaussian), K=5 )$delta[1] )

plot(1:10, KCV, type="b", log="y")
```



# 5-fold cross-validation





# A single hold-out test point

For  $i = 1, 2, \dots, n$

- Hold out the  $i$ th training observation  $(x_i, y_i)$
- Use  $n - 1$  training observations

$$(x_1, y_1), \dots, (x_{i-1}, y_{i-1}), (x_{i+1}, y_{i+1}), \dots, (x_n, y_n)$$

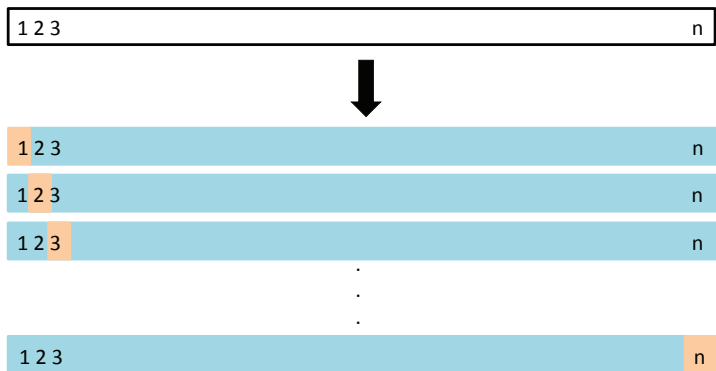
to estimate  $\hat{f}^{-i}(x)$ , and the test observation  $(x_i, y_i)$  to evaluate it:  $(y_i - \hat{f}^{-i}(x_i))^2$

Then take the average

$$\mathbb{E}(\widehat{\text{MSE}}_{\text{Te}}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}^{-i}(x_i))^2$$



# Leave-One-Out Cross-Validation



Source: Gareth et al. (2013), p. 179



# Shortcut for LOOCV

For the linear model

$$\frac{1}{n} \sum_{i=1}^n \left( y_i - \hat{f}^{-i}(x_i) \right)^2 = \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - \hat{f}(x_i)}{1 - h_i} \right)^2$$

where  $h_i$  is the  $i$ th diagonal element of the projection (hat) matrix

$$\mathbf{H}_{n \times n} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

where  $\mathbf{X}_{n \times p}$  is the design matrix



# Generalized Cross-Validation

$$\frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - \hat{f}(x_i)}{1 - h_i} \right)^2 \approx \frac{\text{MSE}_{\text{Tr}}}{\left(1 - \frac{p}{n}\right)^2}$$

where we approximate each  $h_i$  with their average  $\frac{1}{n} \sum_{i=1}^n h_i = \frac{p}{n}$



```
library(boot)

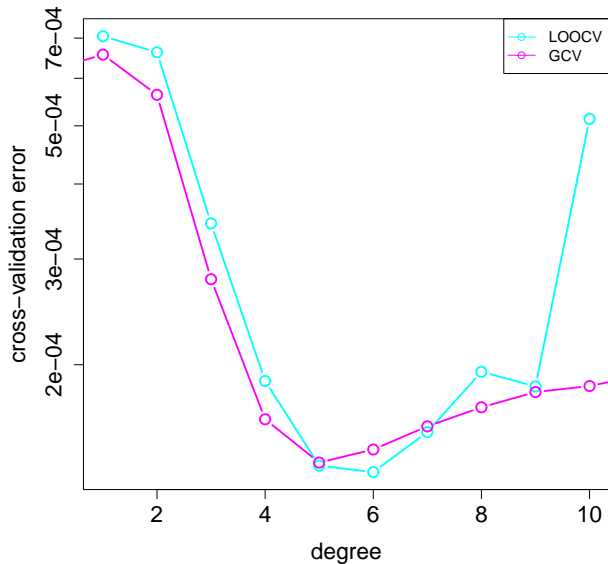
LOOCV = sapply(1:10, function(d)
  cv.glm(train, glm(y~poly(x,degree=d),
    train, family = gaussian) )$delta[1] )

GCV = MSEs.tr/(1-(ds+1)/n )^2

plot(1:10, LOOCV, type="b", log="y")
lines(ds, GCV)
```



# LOOCV and GCV



# Cross-validation bias-variance trade-off

## Bias

- $K$ -fold CV with  $K = 5$  or  $10$  gives a biased (upward) estimate of  $\mathbb{E}(\text{MSE}_{\text{Te}})$  because it uses less information ( $4/5$  or  $9/10$  of the observations)
- LOOCV has very low bias (it uses  $n - 1$  observations)

## Variance

- LOOCV has high variance because it is an average of  $n$  extremely correlated quantities ( $\hat{f}^{-i}$  and  $\hat{f}^{-l}$  are fitted on  $n - 2$  common observations)<sup>1</sup>
- $K$ -fold CV with  $K = 5$  or  $10$  has less variance because it is an average of quantities that are less correlated

---

<sup>1</sup>remember that  $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$

