

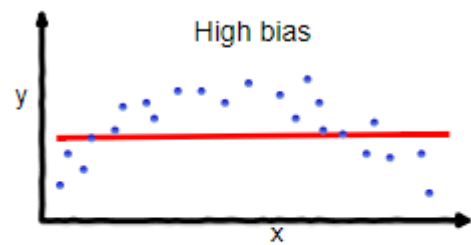
# Class 4 – Polynomial Regression and Resampling methods

---

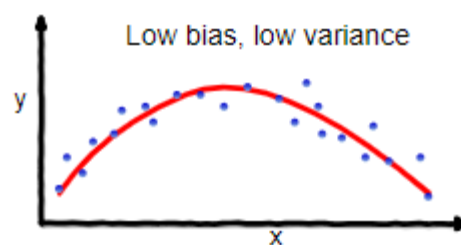
Pedram Jahangiry

Fall 2019

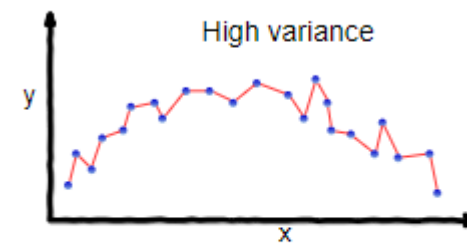




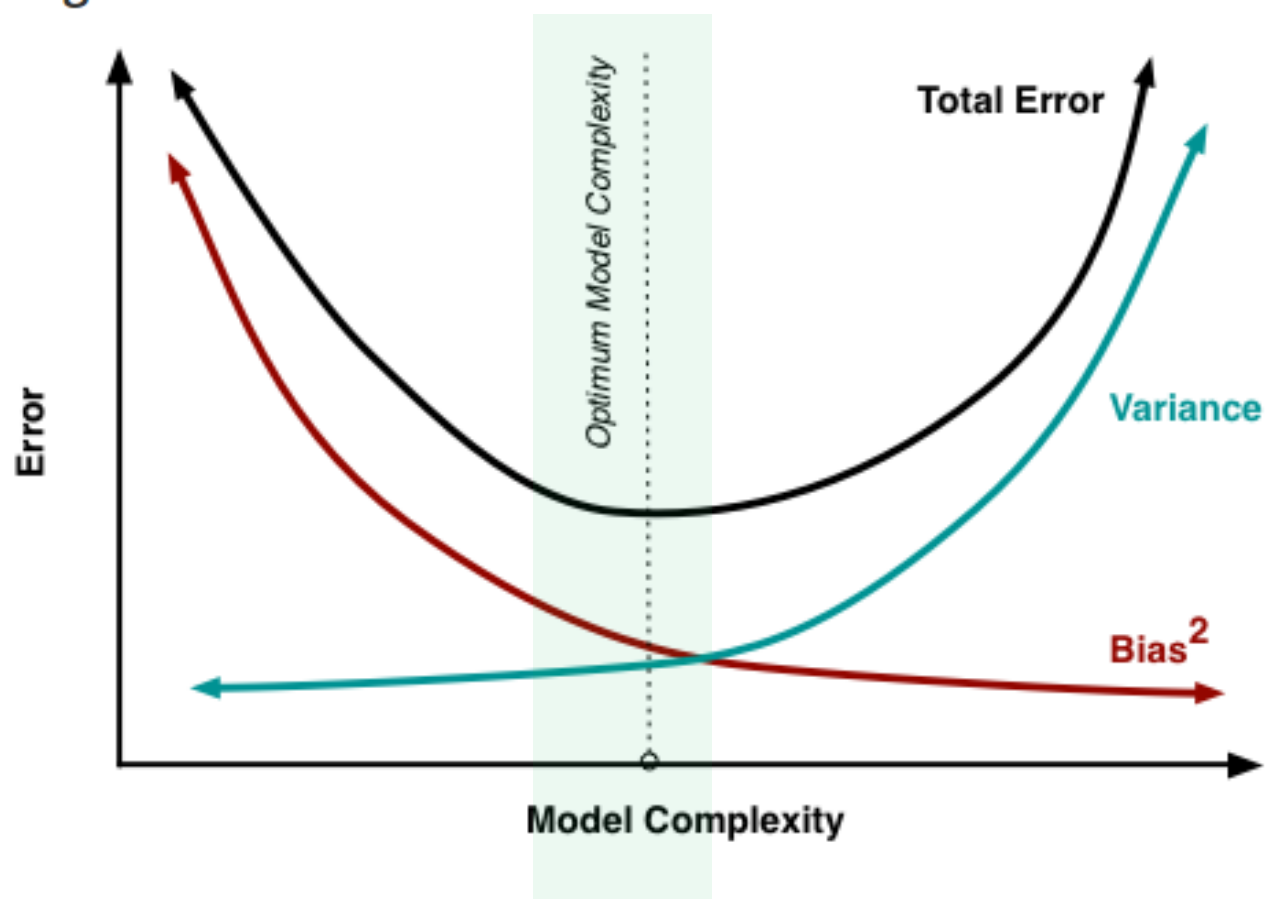
underfitting



Good balance



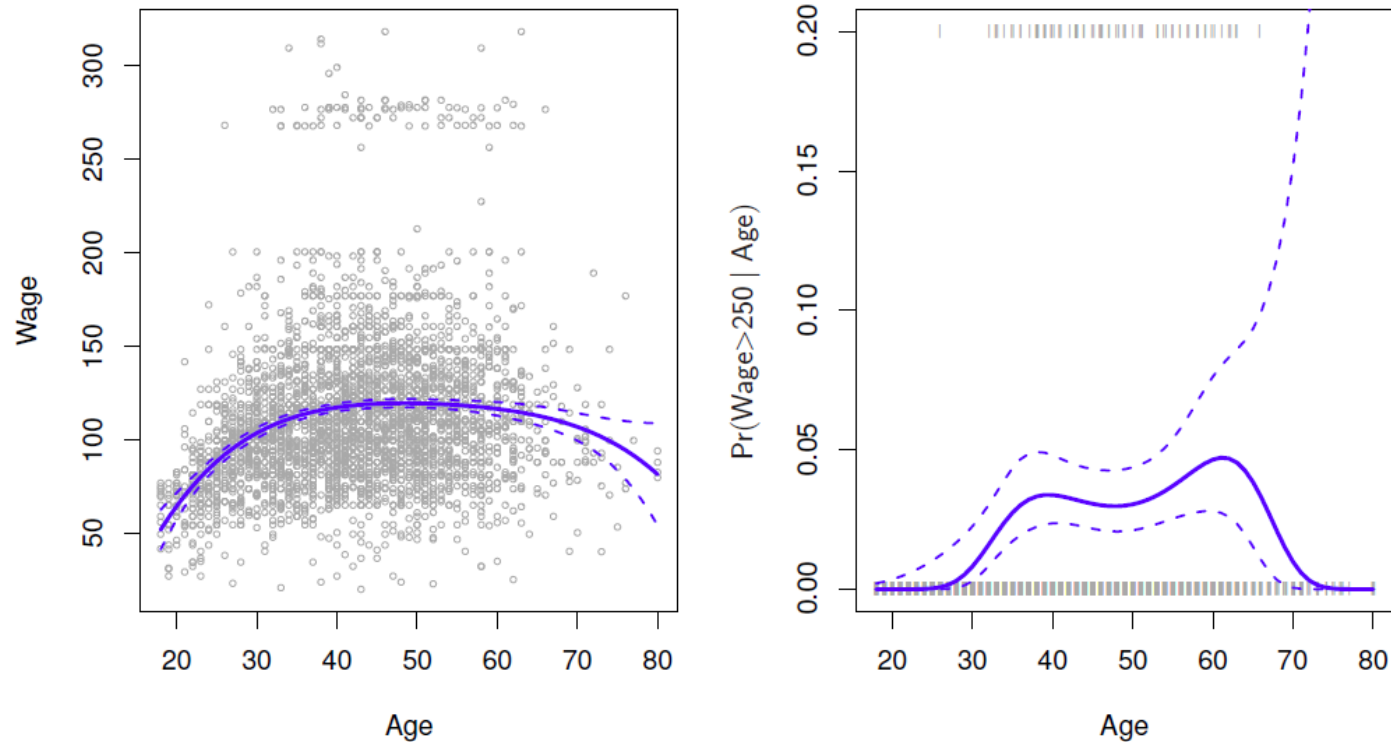
overfitting



# Polynomial Regression

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \dots + \beta_d x_i^d + \epsilon_i$$

**Degree-4 Polynomial**



# Polynomial Regression

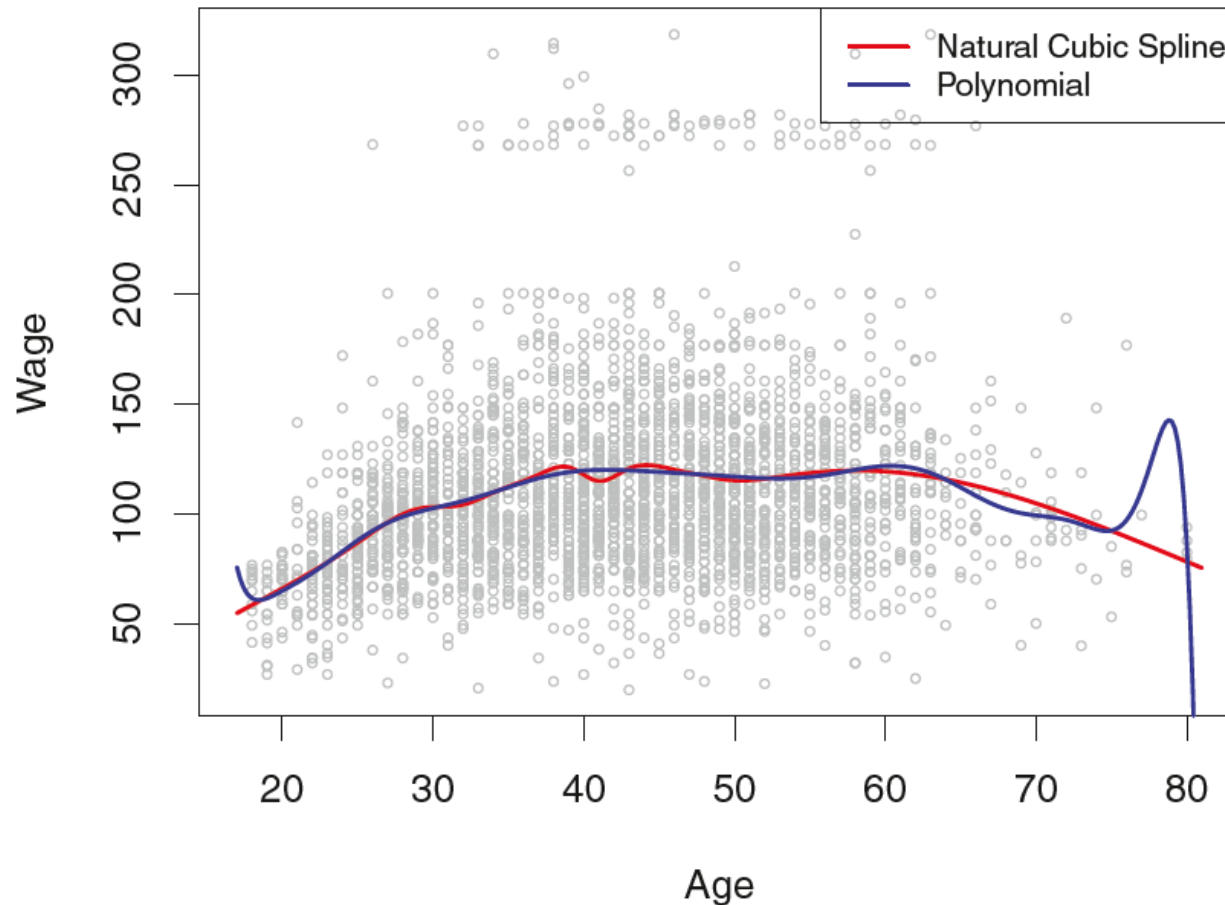
- Create new variables  $X_1 = X$ ,  $X_2 = X^2$ , etc and then treat as multiple linear regression.
- Not really interested in the coefficients; more interested in the fitted function values at any value  $x_0$ :

$$\hat{f}(x_0) = \hat{\beta}_0 + \hat{\beta}_1 x_0 + \hat{\beta}_2 x_0^2 + \hat{\beta}_3 x_0^3 + \hat{\beta}_4 x_0^4.$$

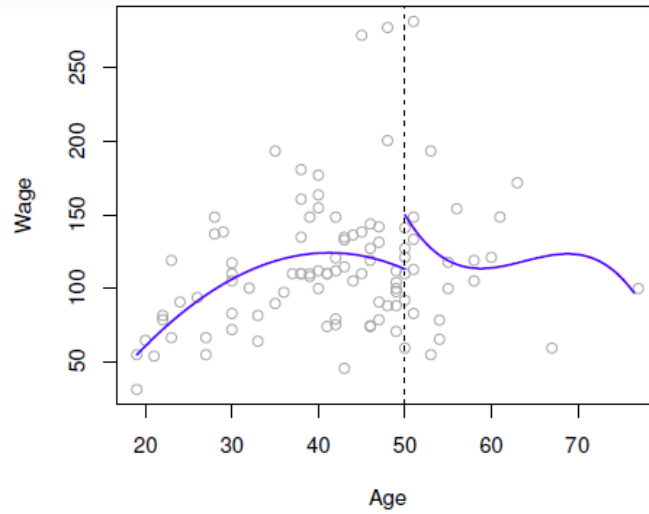
- Since  $\hat{f}(x_0)$  is a linear function of the  $\hat{\beta}_\ell$ , can get a simple expression for *pointwise-variances*  $\text{Var}[\hat{f}(x_0)]$  at any value  $x_0$ . In the figure we have computed the fit and pointwise standard errors on a grid of values for  $x_0$ . We show  $\hat{f}(x_0) \pm 2 \cdot \text{se}[\hat{f}(x_0)]$ .
- We either fix the degree  $d$  at some reasonably low value, else use cross-validation to choose  $d$ .

# Polynomial Regression (Caveat!)

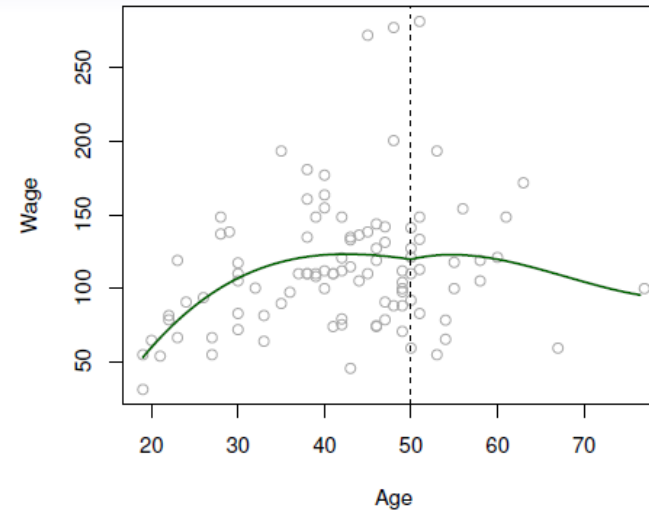
- Polynomials have notorious tail behavior – Very bad for extrapolation.
- Polynomials are global fit! Solution: Piecewise polynomial, splines and local regressions.



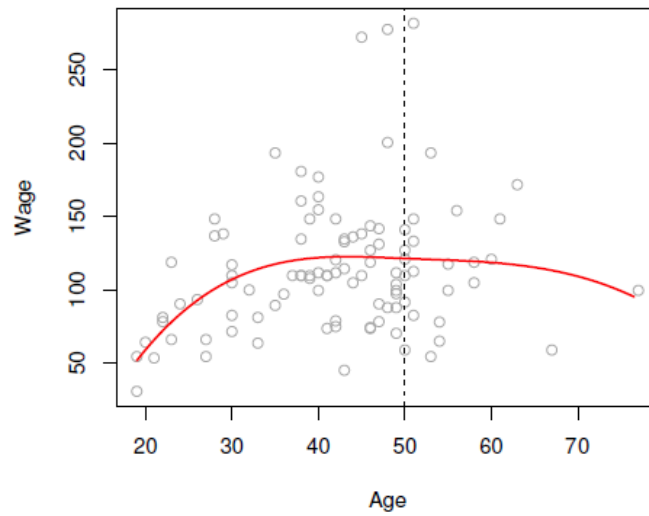
**Piecewise Cubic**



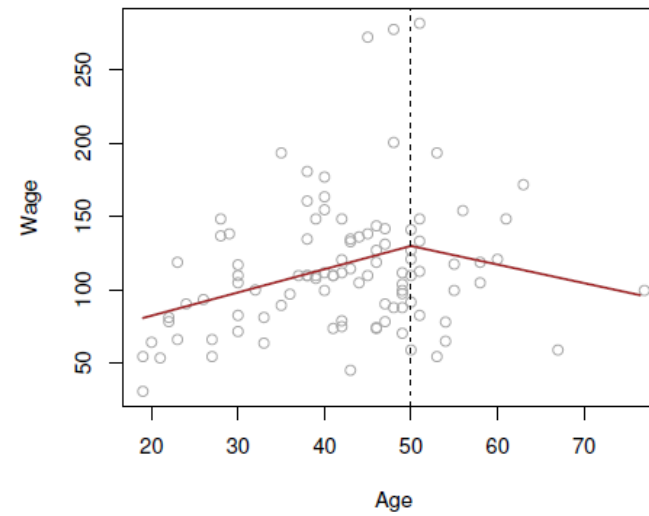
**Continuous Piecewise Cubic**



**Cubic Spline**



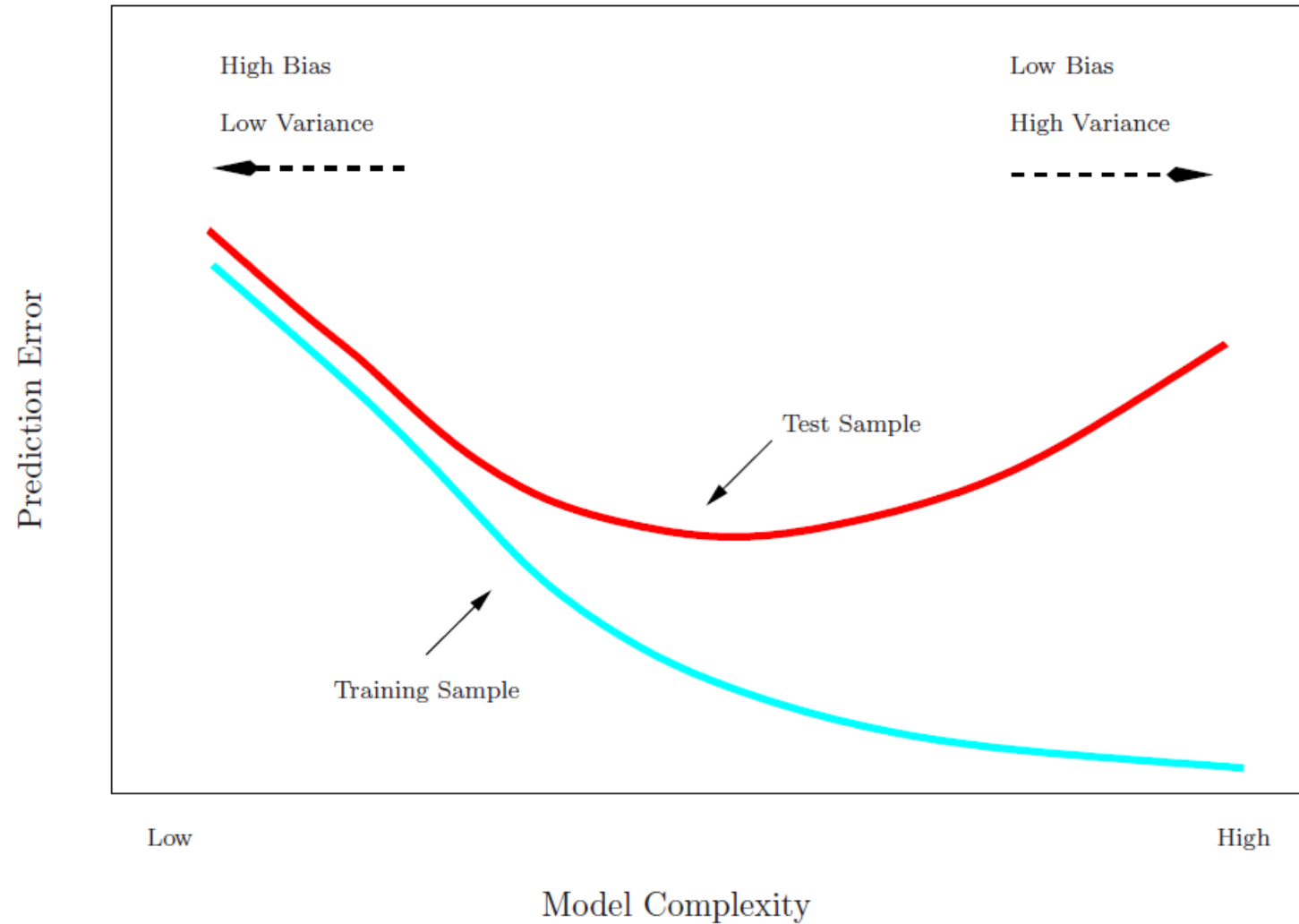
**Linear Spline**



# Resampling methods

- In the section we discuss two *resampling* methods: cross-validation and the bootstrap.
- These methods refit a model of interest to samples formed from the training set, in order to obtain additional information about the fitted model.
- For example, they provide estimates of test-set prediction error, and the standard deviation and bias of our parameter estimates

# Training vs Test set performance





# Validation Set Approach

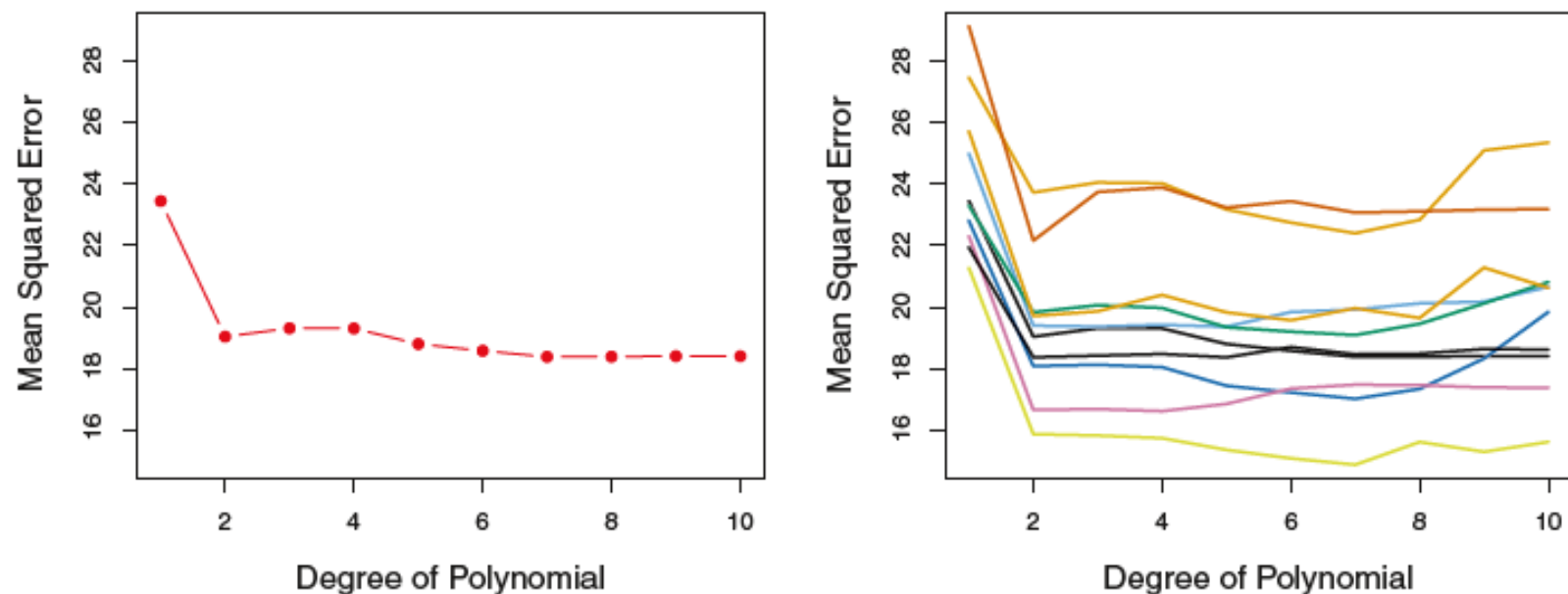
- Here we randomly divide the available set of samples into two parts: a *training set* and a *validation* or *hold-out set*.
- The model is fit on the training set, and the fitted model is used to predict the responses for the observations in the validation set.
- The resulting validation-set error provides an estimate of the test error. This is typically assessed using MSE in the case of a quantitative response and misclassification rate in the case of a qualitative (discrete) response.

# The validation process



A random splitting into two halves: left part is training set, right part is validation set

# Example

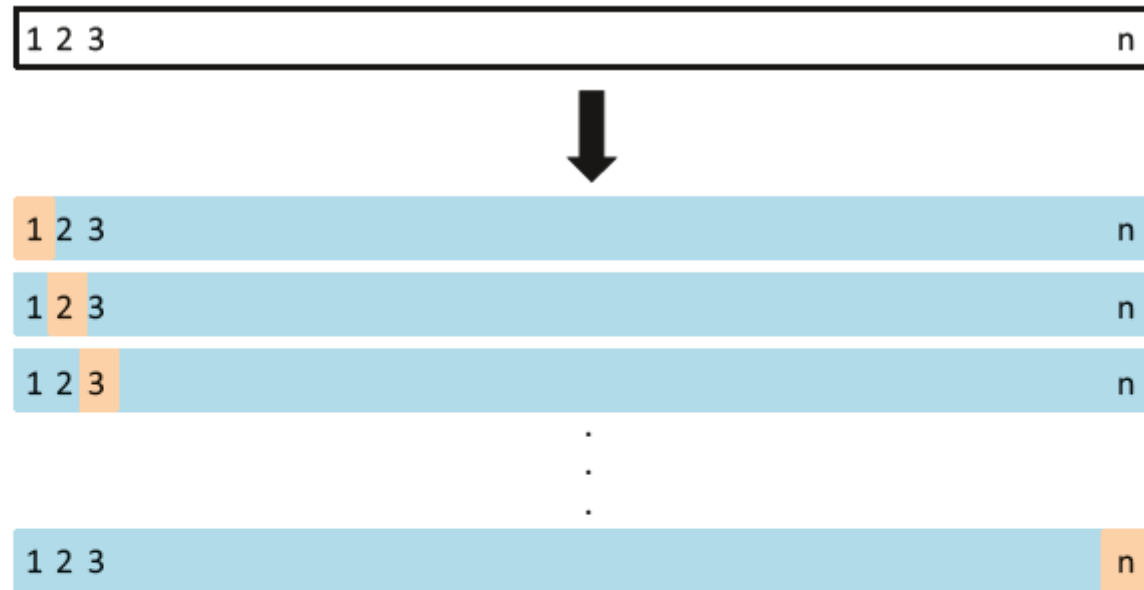


**FIGURE 5.2.** The validation set approach was used on the **Auto** data set in order to estimate the test error that results from predicting **mpg** using polynomial functions of **horsepower**. Left: Validation error estimates for a single split into training and validation data sets. Right: The validation method was repeated ten times, each time using a different random split of the observations into a training set and a validation set. This illustrates the variability in the estimated test MSE that results from this approach.

# Drawbacks of validation set approach

- 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.
- In the validation approach, only a subset of the observations — those that are included in the training set rather than in the validation set — are used to fit the model.
- This suggests that the validation set error may tend to *overestimate* the test error for the model fit on the entire data set. *Why?*

# Leave-One-Out Cross-Validation (LOOCV)

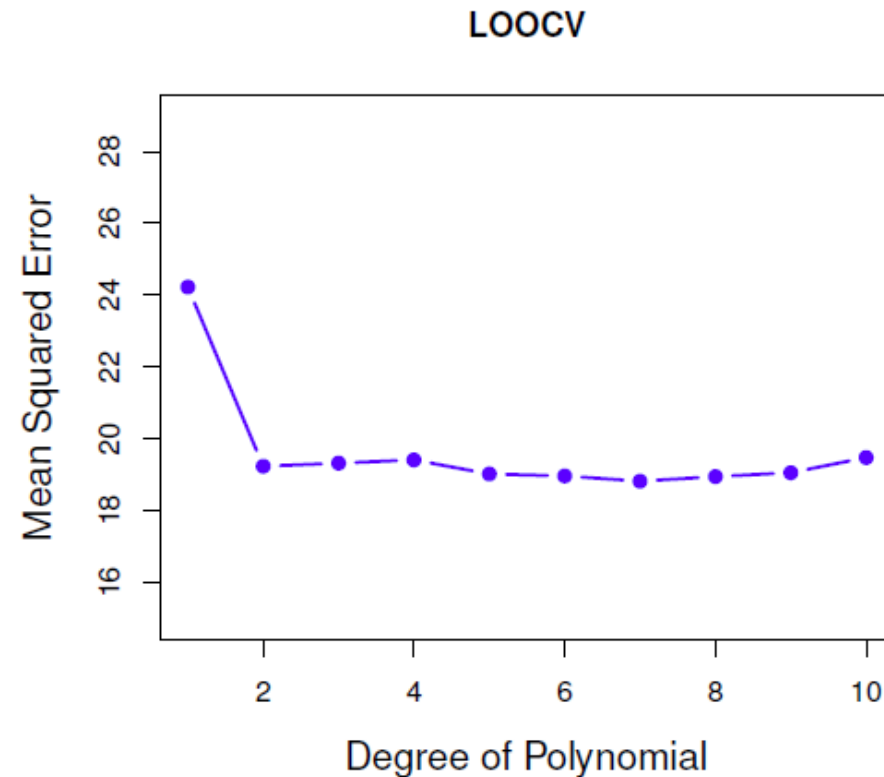


$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n MSE_i$$

**FIGURE 5.3.** A schematic display of LOOCV. A set of  $n$  data points is repeatedly split into a training set (shown in blue) containing all but one observation, and a validation set that contains only that observation (shown in beige). The test error is then estimated by averaging the  $n$  resulting MSE's. The first training set contains all but observation 1, the second training set contains all but observation 2, and so forth.

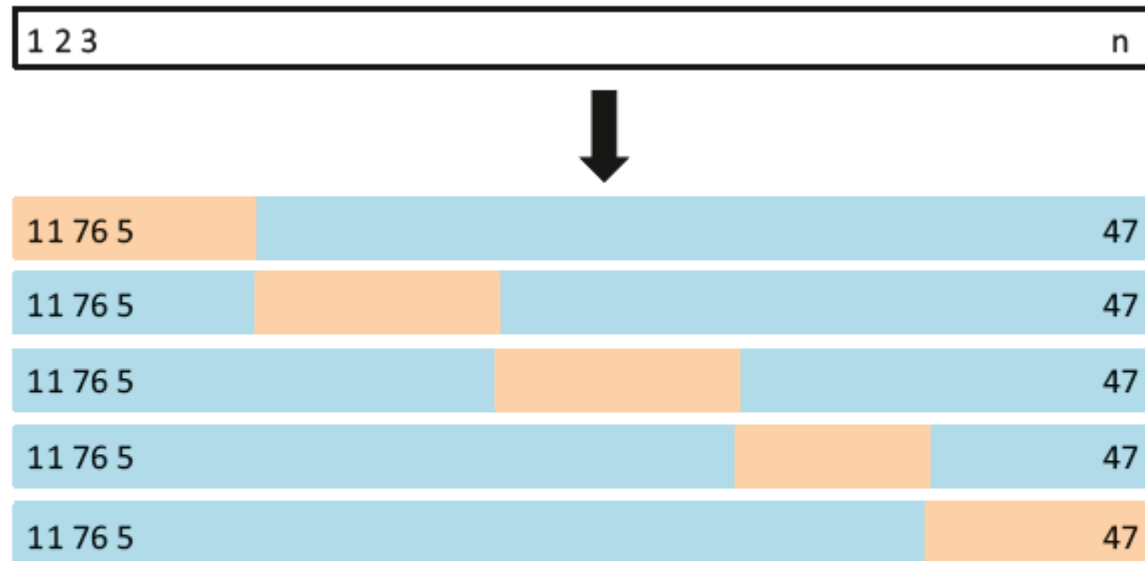
# LOOCV advantages and disadvantages

- 1) **Very small bias:**  
contain  $(n - 1)$  observations, almost as many as are in the entire data set
- 2) **Always yield the same results:**  
there is no randomness in the training/validation set splits.
- 3) **Expensive to implement**
- 4) **High variance!**



# K-fold Cross-Validation

Divide data into  $K$  roughly equal-sized parts ( $K = 5$  here)



$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k MSE_i$$

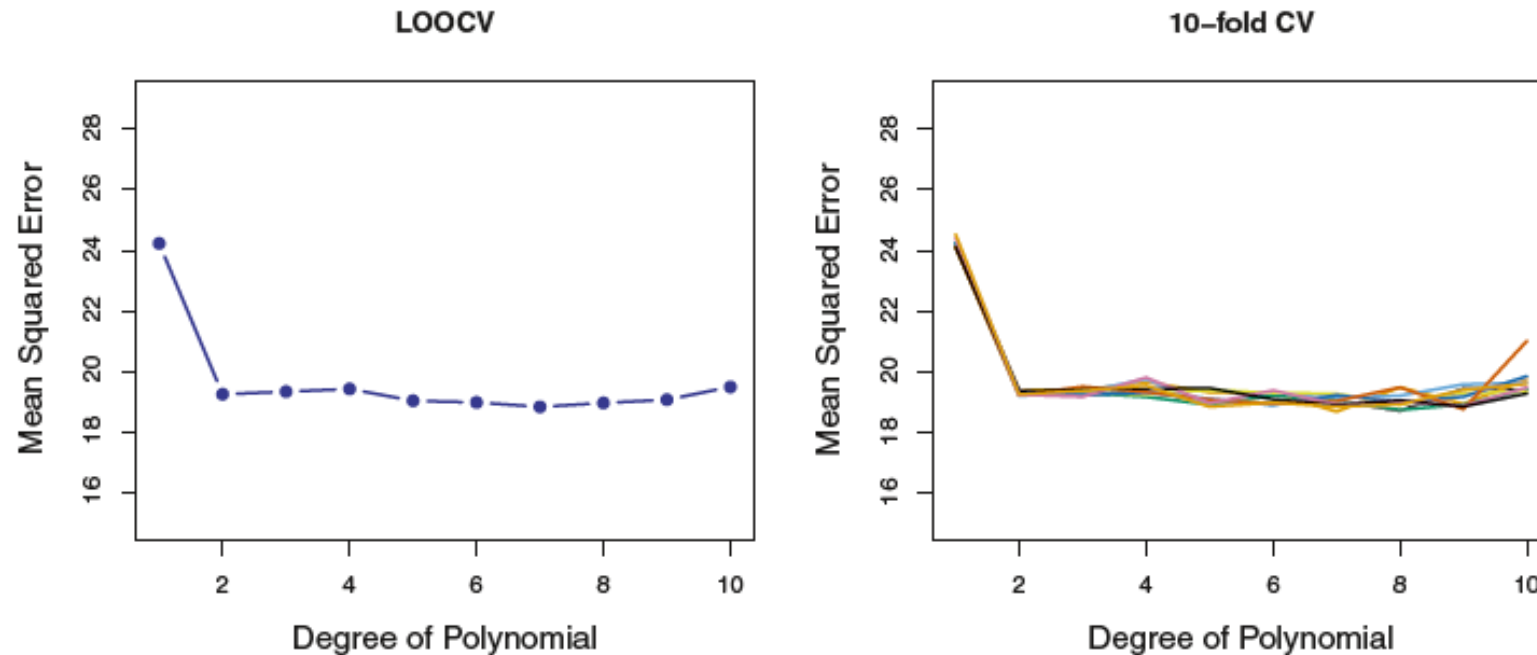
**FIGURE 5.5.** A schematic display of 5-fold CV. A set of  $n$  observations is randomly split into five non-overlapping groups. Each of these fifths acts as a validation set (shown in beige), and the remainder as a training set (shown in blue). The test error is estimated by averaging the five resulting MSE estimates.

# K-fold Cross-Validation

- *Widely used approach* for estimating test error.
- Estimates can be used to select best model, and to give an idea of the test error of the final chosen model.
- Idea is to randomly divide the data into  $K$  equal-sized parts. We leave out part  $k$ , fit the model to the other  $K - 1$  parts (combined), and then obtain predictions for the left-out  $k$ th part.
- This is done in turn for each part  $k = 1, 2, \dots, K$ , and then the results are combined.



# K-fold Cross-validation vs LOOCV



**FIGURE 5.4.** Cross-validation was used on the **Auto** data set in order to estimate the test error that results from predicting **mpg** using polynomial functions of **horsepower**. Left: The LOOCV error curve. Right: 10-fold CV was run nine separate times, each with a different random split of the data into ten parts. The figure shows the nine slightly different CV error curves.