

S&DS 265 / 565
Introductory Machine Learning

Some Context and Concepts

Thursday, September 9

Logistics

- Recordings posted to Canvas under Media Library
- Assignment 1 posted on Tuesday
- Quiz 0 available on Canvas at noon today, for 24 hours
- Check Canvas / EdD for office hours

Plan for Today

- Continue Python elements
- Basics of classification, regression, overfitting
- Linear regression example

Some Terminology

- supervised vs. unsupervised
- classification vs. regression
- prediction vs. inference

Supervised Learning vs. Unsupervised Learning

Supervised learning:

- Given a set of (x, y) , learn to predict y using x .
- e.g.
 - ▶ Predicting whether a loan will default based on customer characteristics

Supervised Learning vs. Unsupervised Learning

Supervised learning:

- Given a set of (x, y) , learn to predict y using x .
- e.g.
 - ▶ Predicting whether a loan will default based on customer characteristics

Unsupervised learning:

- Given a set of x , learn underlying structure or relationships of x .
- e.g.
 - ▶ Identifying market segments with similar spending patterns.

Classification vs. Regression

The `Income` dataset:

Education	Seniority	Income
21.58621	113.1034	99.91717
18.27586	119.3103	92.57913
12.06897	100.6897	34.67873
17.03448	187.5862	78.70281
19.93103	20.0000	68.00992
18.27586	26.2069	71.50449

Information for 30 *simulated*
individuals.

Classification vs. Regression

The `Income` dataset:

Education	Seniority	Income
21.58621	113.1034	99.91717
18.27586	119.3103	92.57913
12.06897	100.6897	34.67873
17.03448	187.5862	78.70281
19.93103	20.0000	68.00992
18.27586	26.2069	71.50449

Regression: Model `income` based on other characteristics.

Information for 30 *simulated individuals*.

Classification vs. Regression

The `Income` dataset:

Education	Seniority	Income
21.58621	113.1034	99.91717
18.27586	119.3103	92.57913
12.06897	100.6897	34.67873
17.03448	187.5862	78.70281
19.93103	20.0000	68.00992
18.27586	26.2069	71.50449

Information for 30 *simulated individuals*.

Regression: Model **income** based on other characteristics.

Classification: Model **whether someone will earn above the median income** based on other characteristics.

Inference vs. Prediction

The `Income` dataset:

Education	Seniority	Income
21.58621	113.1034	99.91717
18.27586	119.3103	92.57913
12.06897	100.6897	34.67873
17.03448	187.5862	78.70281
19.93103	20.0000	68.00992
18.27586	26.2069	71.50449

Prediction: accurately predict Y for new observations

Information for 30 *simulated* individuals.

Inference vs. Prediction

The `Income` dataset:

Education	Seniority	Income
21.58621	113.1034	99.91717
18.27586	119.3103	92.57913
12.06897	100.6897	34.67873
17.03448	187.5862	78.70281
19.93103	20.0000	68.00992
18.27586	26.2069	71.50449

Prediction: accurately predict Y for new observations

Inference: explain the underlying relationship between Y and X

Information for 30 *simulated* individuals.

Example: Handwritten Digit Recognition

- Data: images of handwritten digits (grayscale pixel values)
- Classify images as digits 0 to 9.



Example: Handwritten Digit Recognition

- Data: images of handwritten digits (grayscale pixel values)
- Classify images as digits 0 to 9.



80322-4129 80206

40004 14310

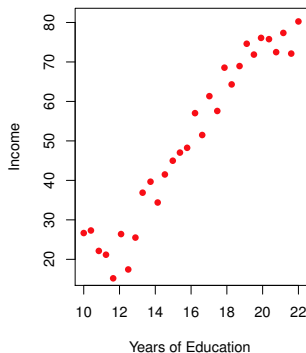
37879 05153

5502 75216

35460 44209

Regression Example

The `Income` dataset:



Quantitative response Y

Predictors $X = (X_1, \dots, X_p)$

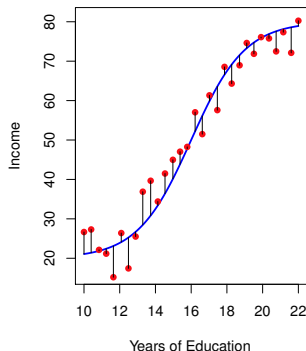
Assume the relationship can be expressed by:

$$Y = f(X) + \epsilon,$$

where f is a fixed, unknown function and ϵ is error term.

Regression Example

The `Income` dataset:



Quantitative response Y

Predictors $X = (X_1, \dots, X_p)$

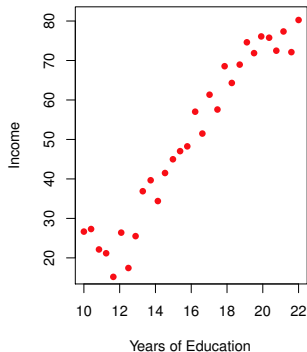
Assume the relationship can be expressed by:

$$Y = f(X) + \epsilon,$$

where f is a fixed, unknown function and ϵ is error term.

Regression Example

Back to regression with $p = 1$:



$$Y = f(X) + \epsilon$$

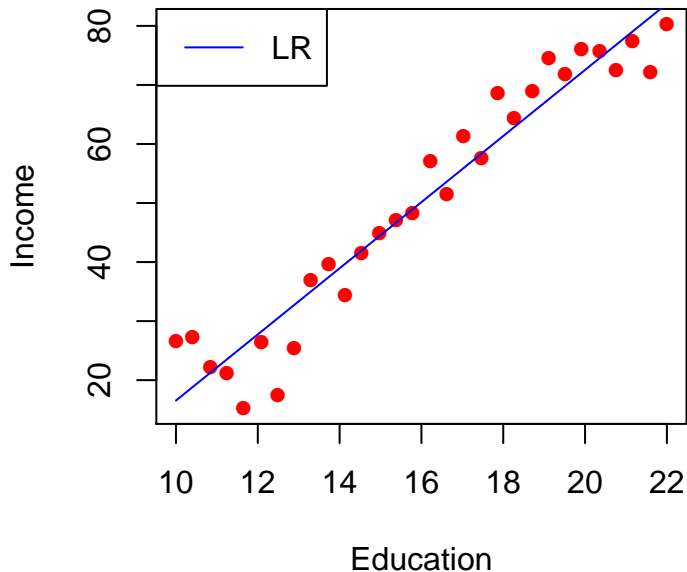
Modeling:

Use a procedure to get \hat{f} . Derive estimates $\hat{Y} = \hat{f}(X)$.

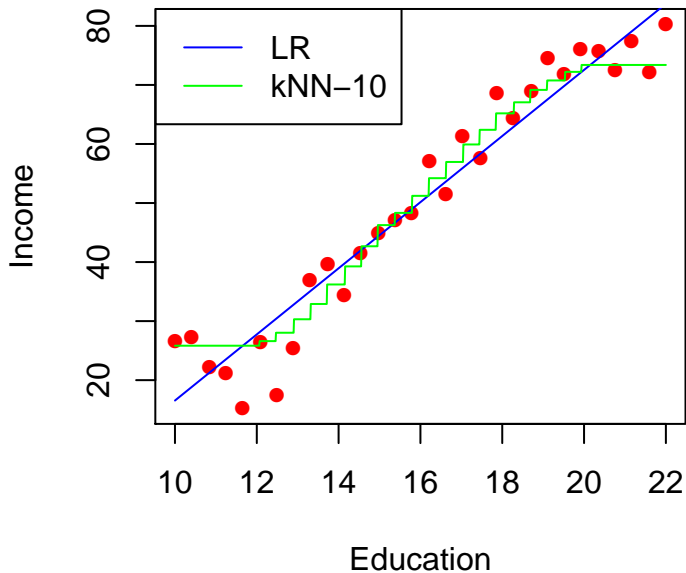
Possible Regression Approaches

- linear regression
 - ▶ Fitting a straight line through the data.
- k -nearest neighbors regression
 - ▶ Average together the y_i for x_i close to x

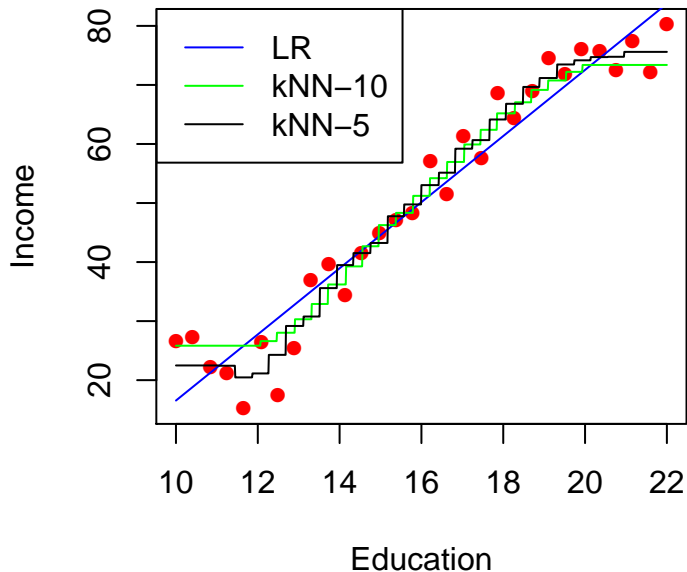
Possible Regression Approaches



Possible Regression Approaches

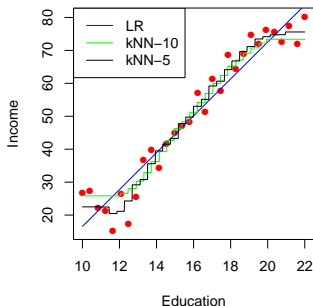


Possible Regression Approaches



Possible Regression Approaches

Measuring performance via **Mean Squared Error**



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

Possible Regression Approaches

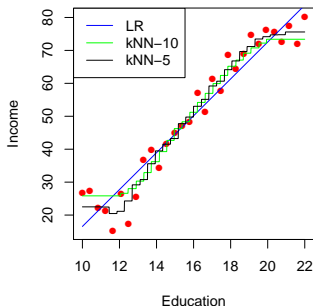
Measuring performance via **Mean Squared Error**

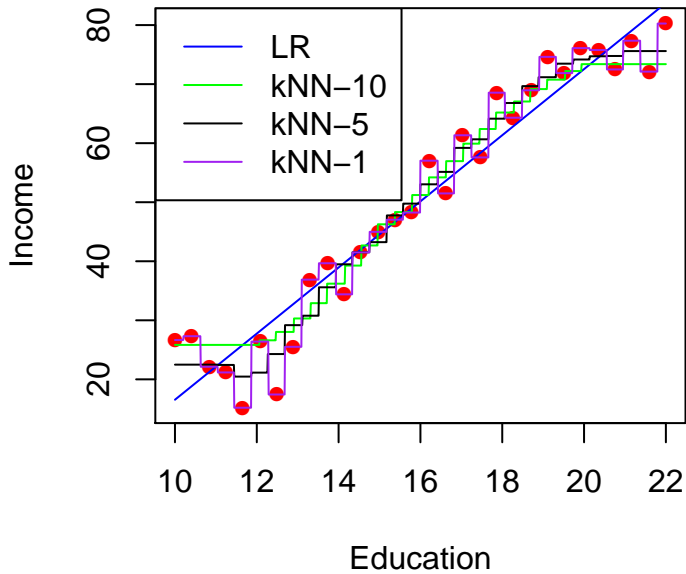
$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

MSEs for three methods:

Linear Regression	29.829
k-Nearest Neighbors (k=10)	23.519
k-Nearest Neighbors (k=5)	16.21

A k -nearest neighbors model with $k = 5$ achieves lowest error. Is it the best?





Training MSE vs. Test MSE

MSE in the previous table, **training MSE**, was computed based on data used in fitting the model.

We are more interested in **test MSE** computed on *unseen data*.

Training MSE vs. Test MSE

MSE in the previous table, **training MSE**, was computed based on data used in fitting the model.

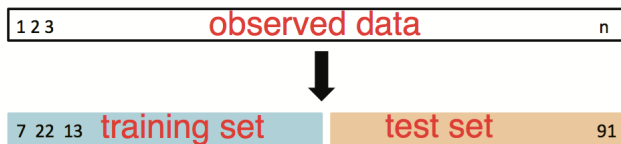
We are more interested in **test MSE** computed on *unseen data*. What if we don't have other data?

Training MSE vs. Test MSE

MSE in the previous table, **training MSE**, was computed based on data used in fitting the model.

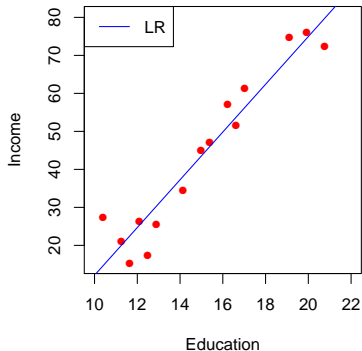
We are more interested in **test MSE** computed on *unseen data*. What if we don't have other data?

We can randomly split our data into a test set and a training set.

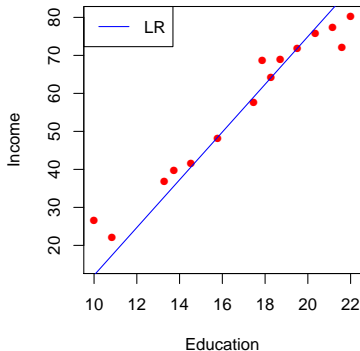


Regression Approaches Revisited

Training Set

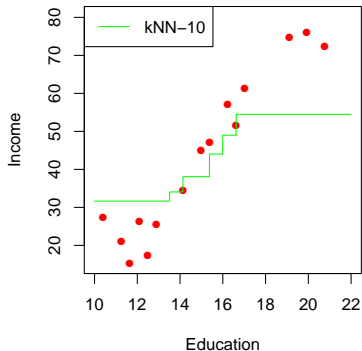


Test Set

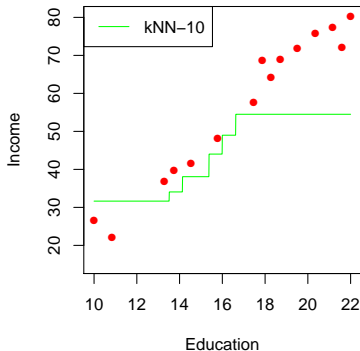


Regression Approaches Revisited

Training Set

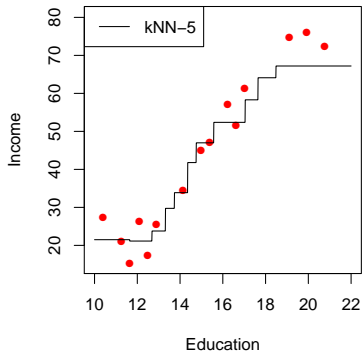


Test Set

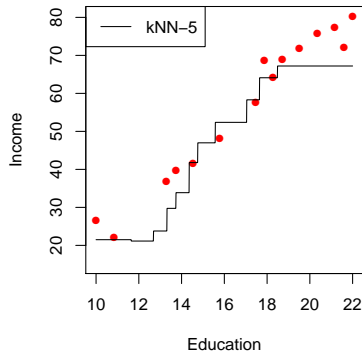


Regression Approaches Revisited

Training Set



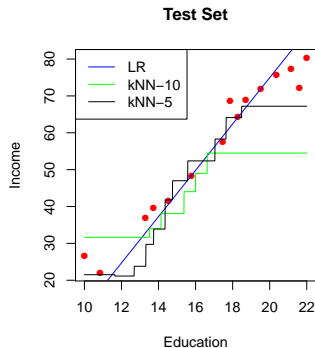
Test Set



Regression Approaches Revisited

Compute MSE on the test set:

$$MSE = \frac{1}{n} \sum (y_i - \hat{f}(x_i))^2$$

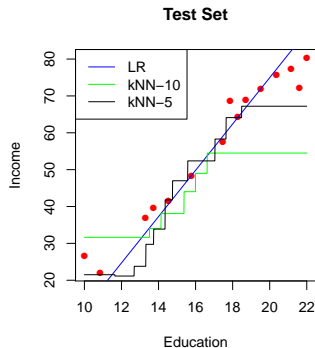


Linear Regression	37.807
k-Nearest Neighbors (k=10)	197.809
k-Nearest Neighbors (k=5)	48.682

Regression Approaches Revisited

Compute MSE on the test set:

$$MSE = \frac{1}{n} \sum (y_i - \hat{f}(x_i))^2$$



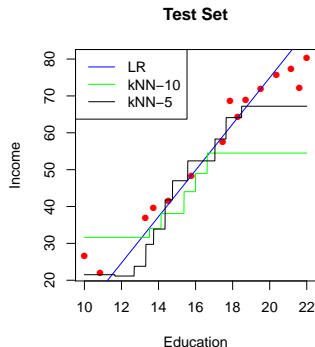
Linear Regression	37.807
k-Nearest Neighbors (k=10)	197.809
k-Nearest Neighbors (k=5)	48.682

So it appears that linear regression wins.

Regression Approaches Revisited

Compute MSE on the test set:

$$MSE = \frac{1}{n} \sum (y_i - \hat{f}(x_i))^2$$



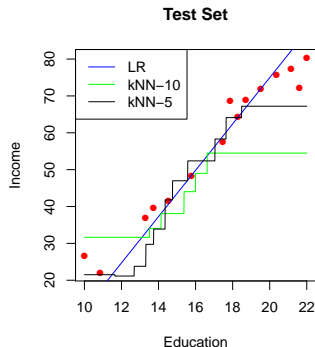
Linear Regression	37.807
k-Nearest Neighbors (k=10)	197.809
k-Nearest Neighbors (k=5)	48.682

So it appears that linear regression wins.
Does it?

Regression Approaches Revisited

Compute MSE on the test set:

$$MSE = \frac{1}{n} \sum (y_i - \hat{f}(x_i))^2$$



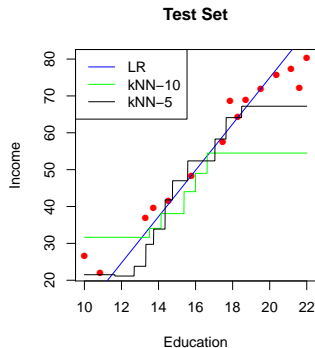
Linear Regression	37.807
k-Nearest Neighbors (k=10)	197.809
k-Nearest Neighbors (k=5)	48.682

So it appears that linear regression wins. Does it? With different random splits of test vs. training, we could have gotten different results.

Regression Approaches Revisited

Compute MSE on the test set:

$$MSE = \frac{1}{n} \sum (y_i - \hat{f}(x_i))^2$$



Linear Regression	37.807
k-Nearest Neighbors (k=10)	197.809
k-Nearest Neighbors (k=5)	48.682

So it appears that linear regression wins. Does it? With different random splits of test vs. training, we could have gotten different results. [We'll talk about ways around this later.](#)

Overfitting

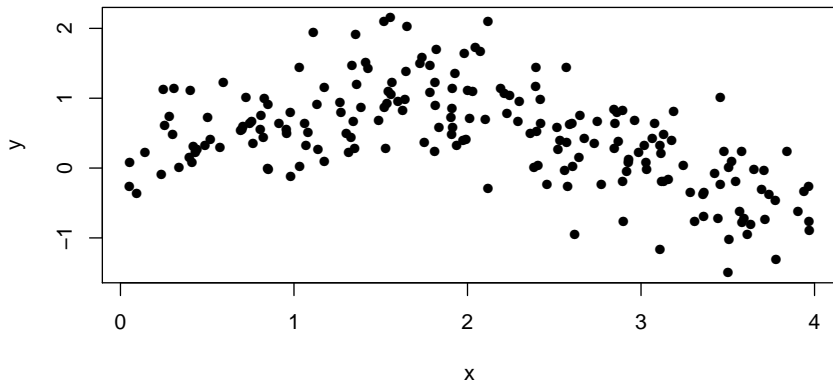
A method is **overfitting** the data when it has a small training MSE but a large test MSE.

Overfitting

A method is **overfitting** the data when it has a small training MSE but a large test MSE.

Let's examine this phenomenon using a bigger dataset:

Simulated Data

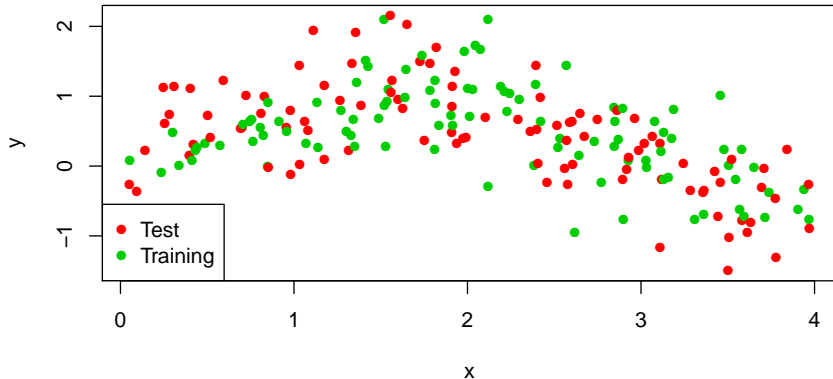


Overfitting

A method is **overfitting** the data when it has a small training MSE but a large test MSE.

Let's examine this phenomenon using a bigger dataset:

Simulated Data

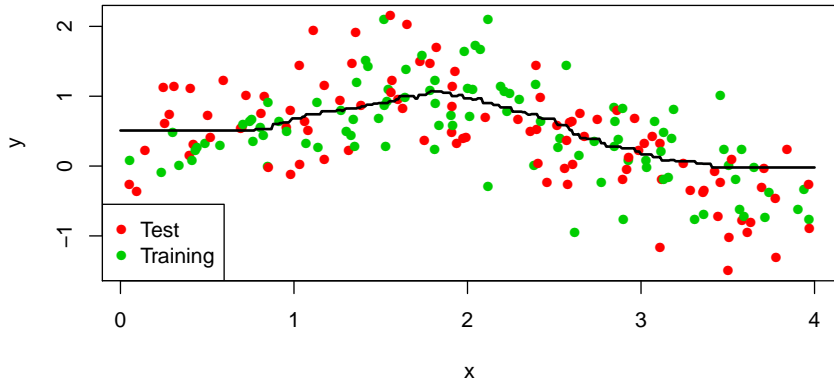


Overfitting

A method is **overfitting** the data when it has a small training MSE but a large test MSE.

Let's examine this phenomenon using a bigger dataset:

kNN fit ($k=30$)

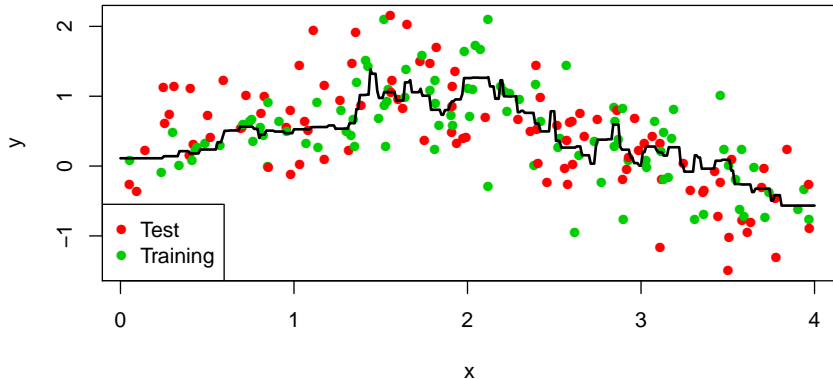


Overfitting

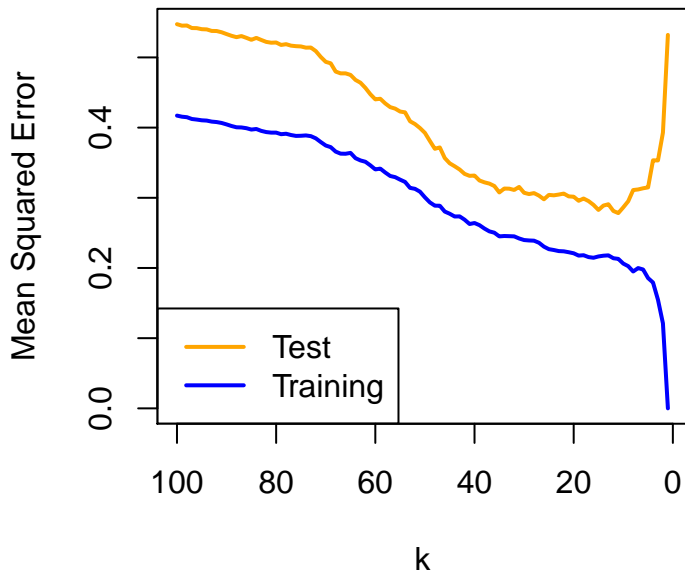
A method is **overfitting** the data when it has a small training MSE but a large test MSE.

Let's examine this phenomenon using a bigger dataset:

kNN fit ($k=5$)



Overfitting via k-Nearest Neighbors



Summary

- Two cultures: model based and prediction based
- Prediction based approaches are sometimes not interpretable
- Overfitting is easy with very flexible models and algorithms

Next week: Linear regression and classification