

THE UNIVERSITY OF TEXAS AT AUSTIN



**BUSINESS
ANALYTICS**

Introduction to Machine Learning

Book Chapters 1, 2 and 5.1.

Carlos M. Carvalho

The University of Texas McCombs School of Business

1. Introduction
2. Measuring Accuracy
3. Out-of Sample Predictions
4. Bias-Variance Trade-Off
5. Cross-Validation

1. Introduction to Predictive Models

Simply put, the goal is to predict a target variable Y with input variables X !

In Data Mining terminology this is known as **supervised learning** (also called *Predictive Analytics*).

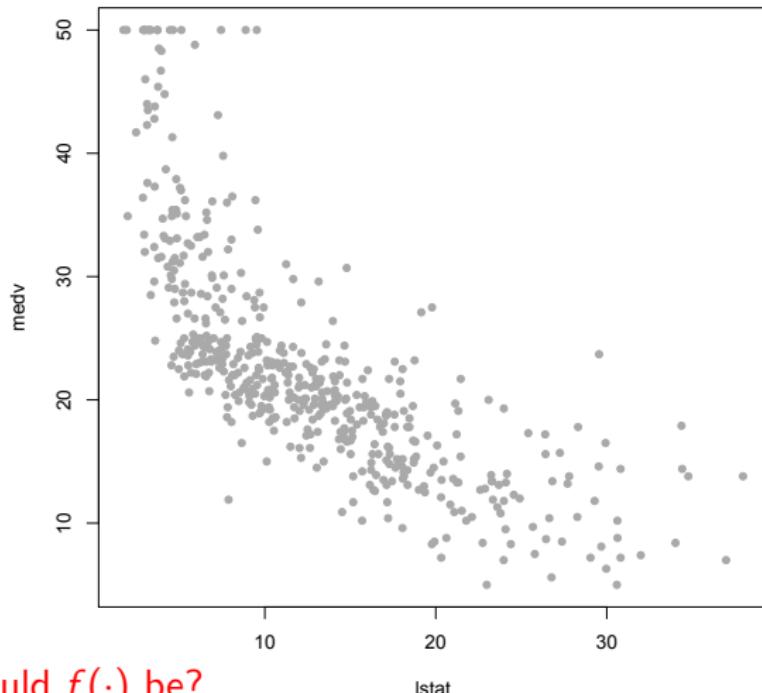
In general, a useful way to think about it is that Y and X are related in the following way:

$$Y_i = f(X_i) + \epsilon_i$$

The main purpose of this part of the course is to *learn or estimate* $f(\cdot)$ from data

Example: Boston Housing

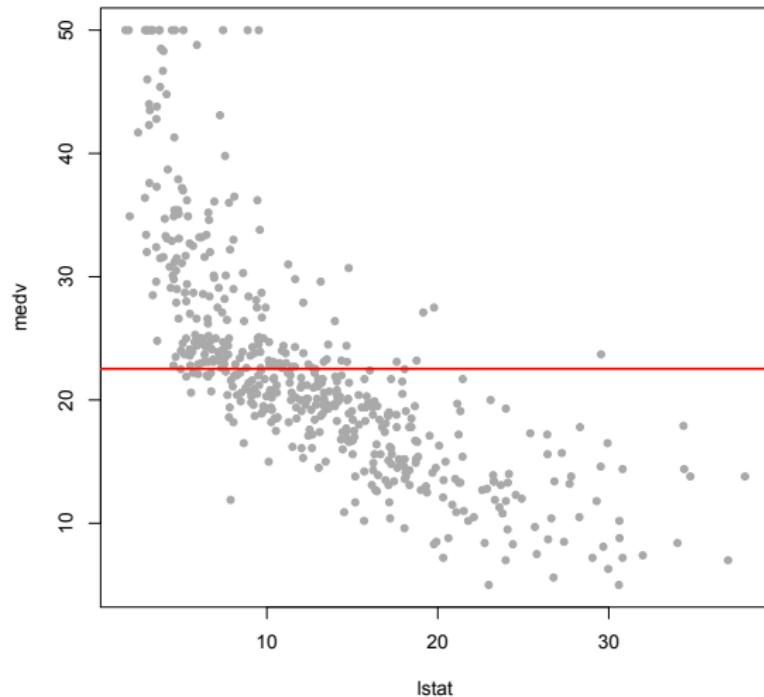
We might be interested in predicting the median house value as a function of some measure of social economic level... here's some data:



What should $f(\cdot)$ be?

Example: Boston Housing

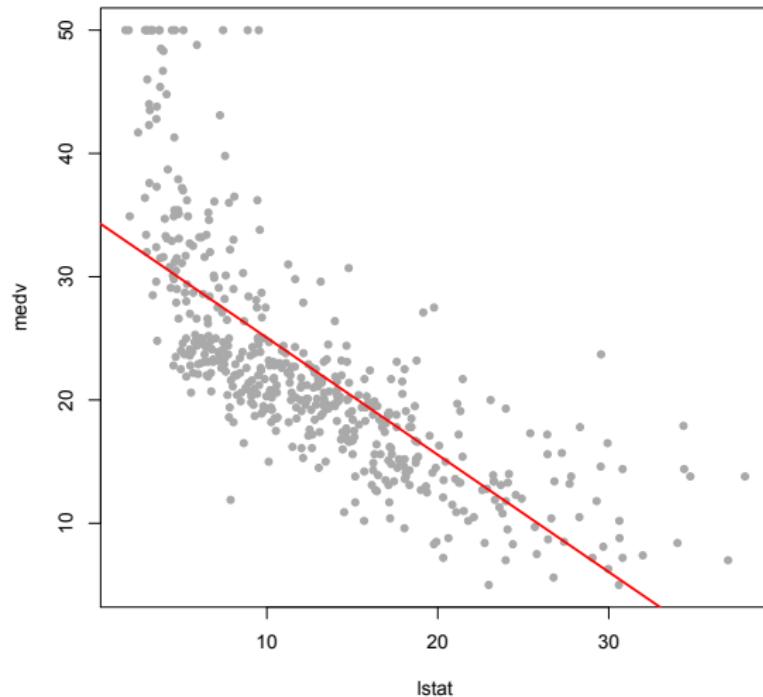
How about this...



If $lstat = 30$ what is the prediction for $medv$?

Example: Boston Housing

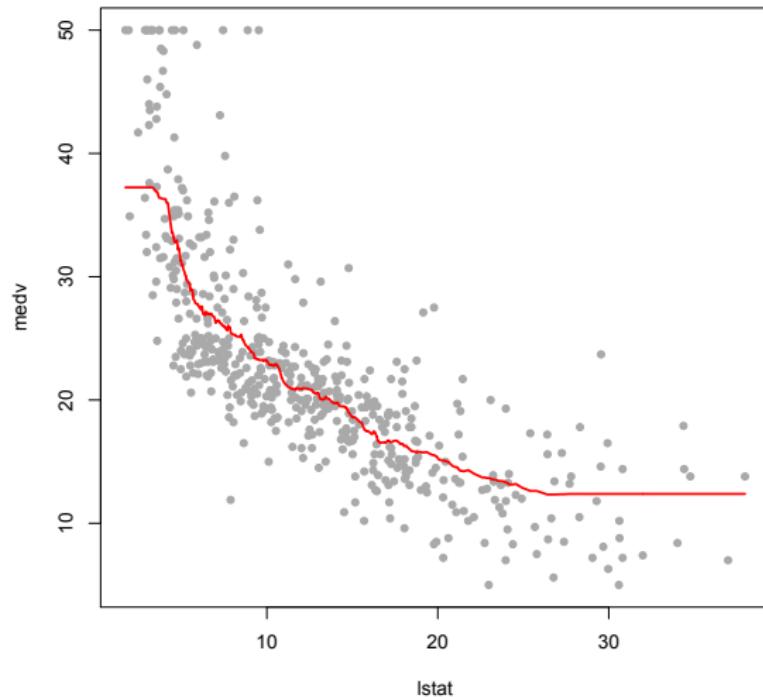
or this...



If $lstat = 30$ what is the prediction for $medv$?

Example: Boston Housing

or even this?



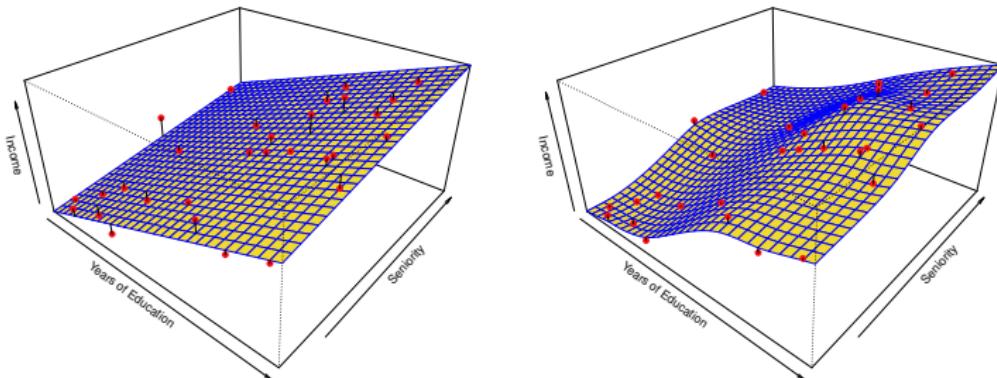
If $lstat = 30$ what is the prediction for $medv$?

How do we estimate $f(\cdot)$?

- ▶ Using *training data*:

$$\{(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)\}$$

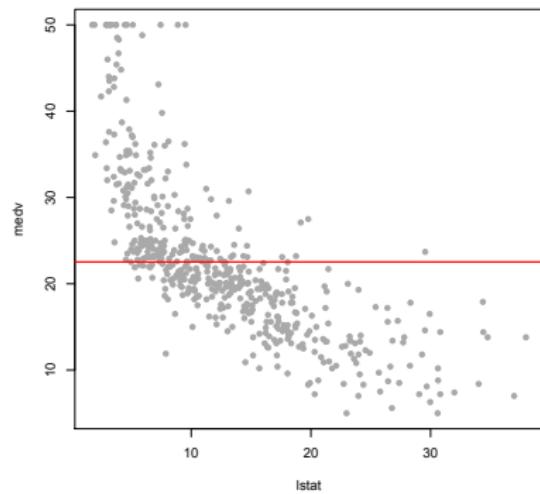
- ▶ We use a statistical method to *estimate* the function $f(\cdot)$
- ▶ Two general methodological strategies:
 1. parametric models (restricted assumptions about $f(\cdot)$)
 2. non-parametric models (flexibility in defining $f(\cdot)$)



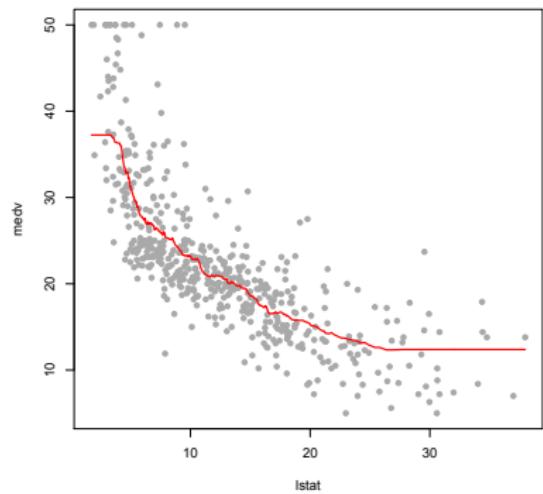
Back to Boston Housing

Parametric Model

$$(Y = \mu + \epsilon)$$



Non-Parametric Model
(k-nearest neighbors)



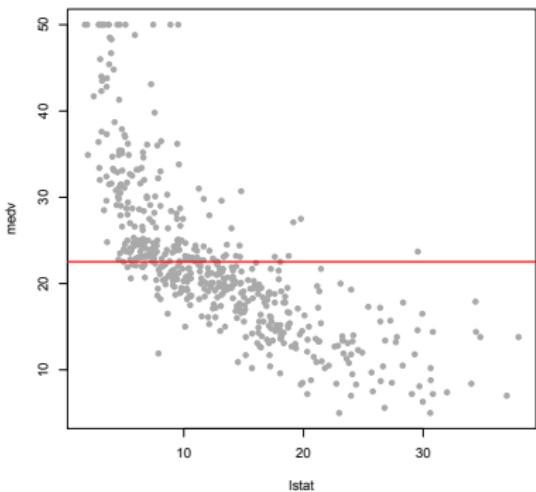
Back to Boston Housing

Simplest parametric model:

$$Y_i = \mu + \epsilon_i$$

Using the training data, we estimate
 $f(\cdot)$ as

$$\widehat{f(\cdot)} = \hat{\mu} = \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$



Back to Boston Housing

The above strategy averages all points in the training set... maybe points that are “closer” to the place I am trying to predict should be more relevant...

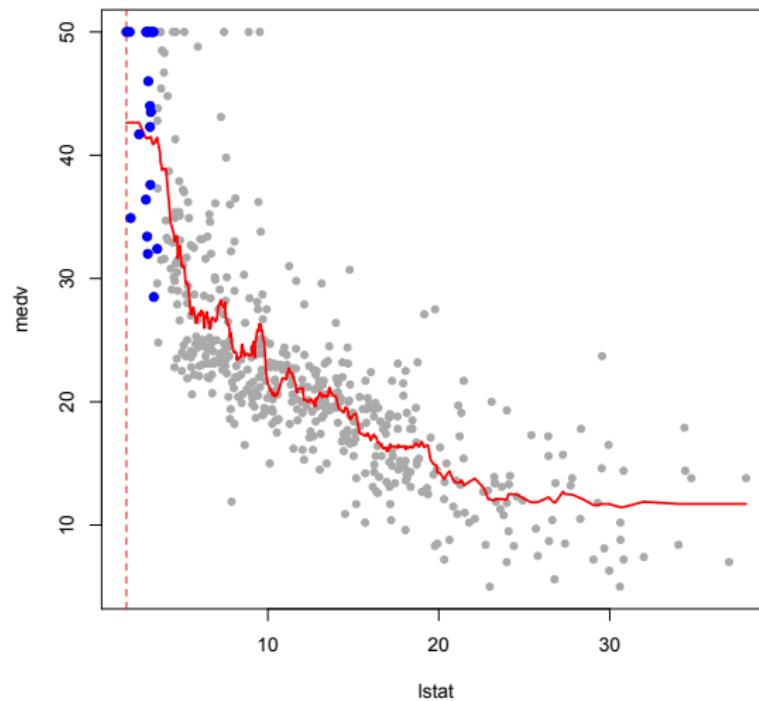
How about averaging the closest 20 neighbors?

What do I mean by closest? We will choose the 20 points that are closest to the X value we are trying to predict.

This is what is called the *k*-nearest neighbors algorithm

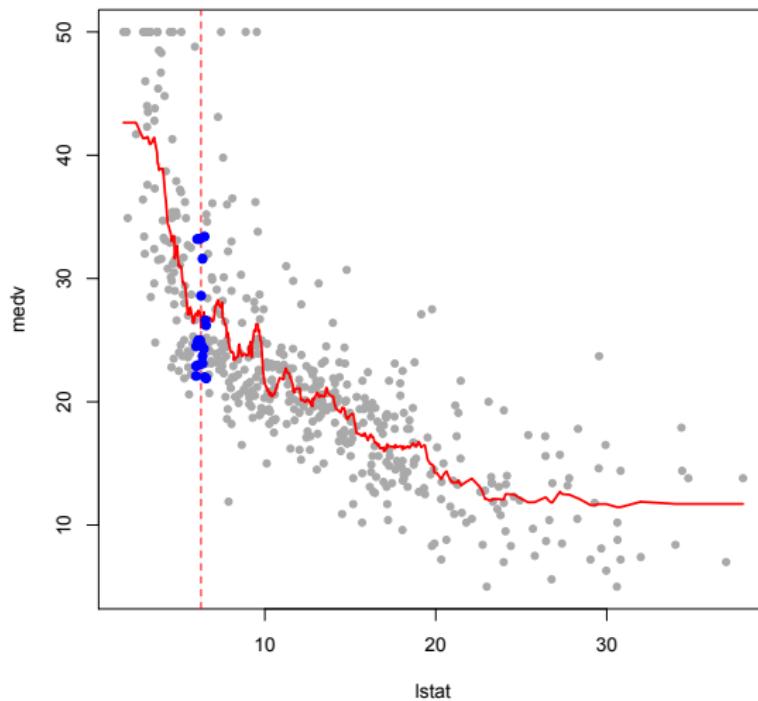
Back to Boston Housing

k= 20



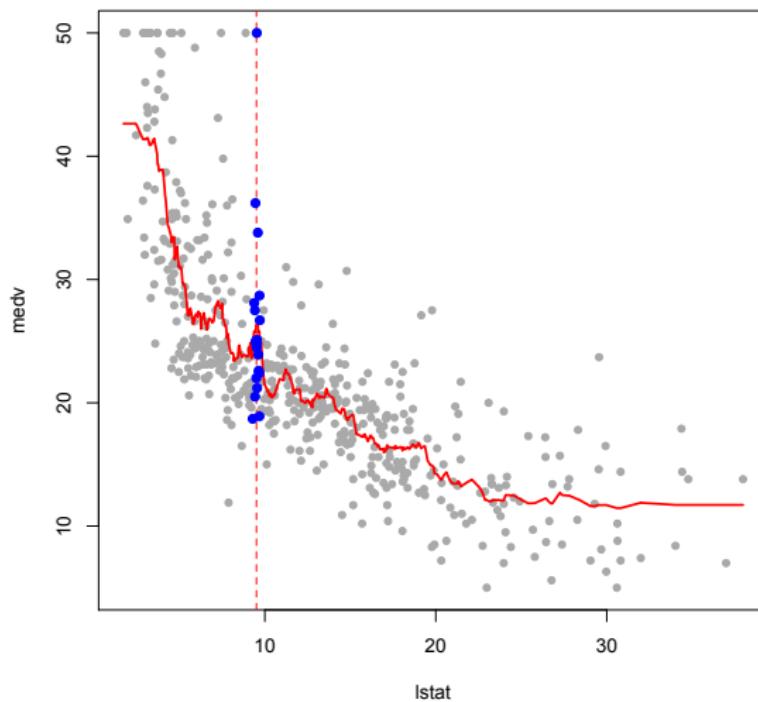
Back to Boston Housing

k= 20



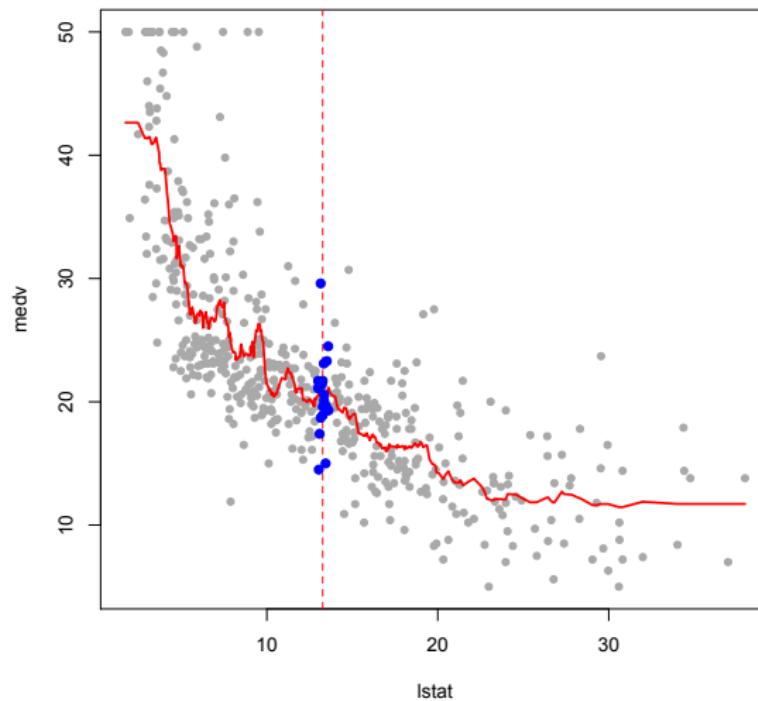
Back to Boston Housing

k= 20



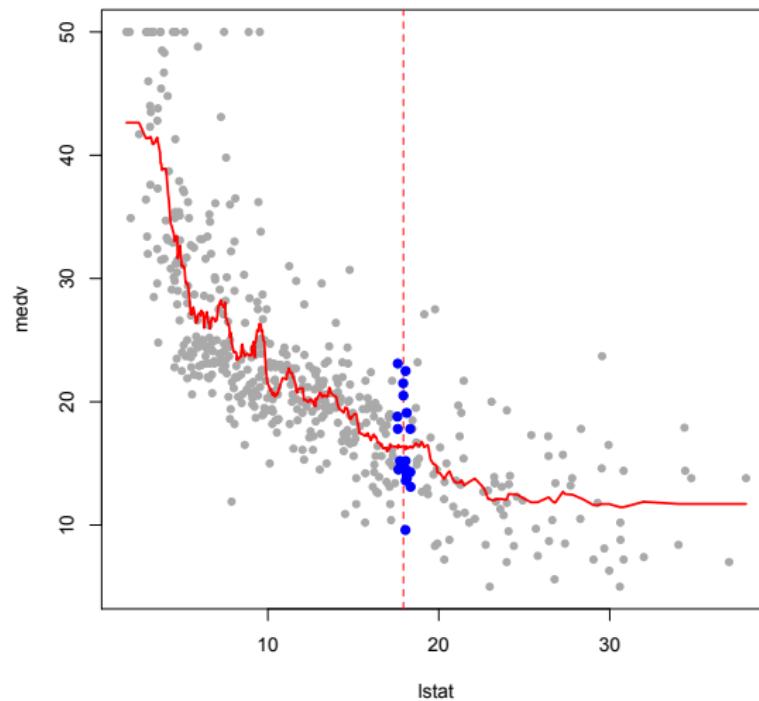
Back to Boston Housing

k = 20



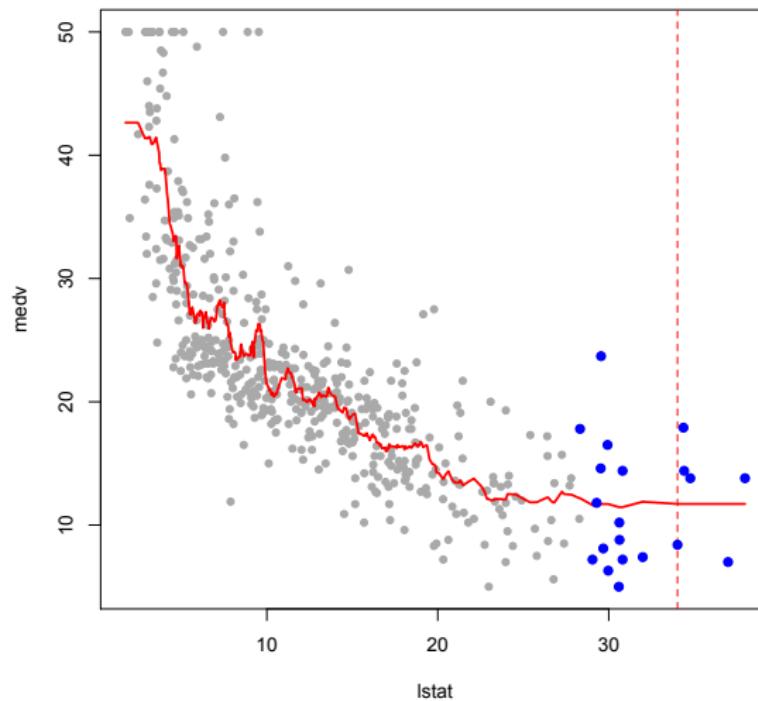
Back to Boston Housing

k = 20



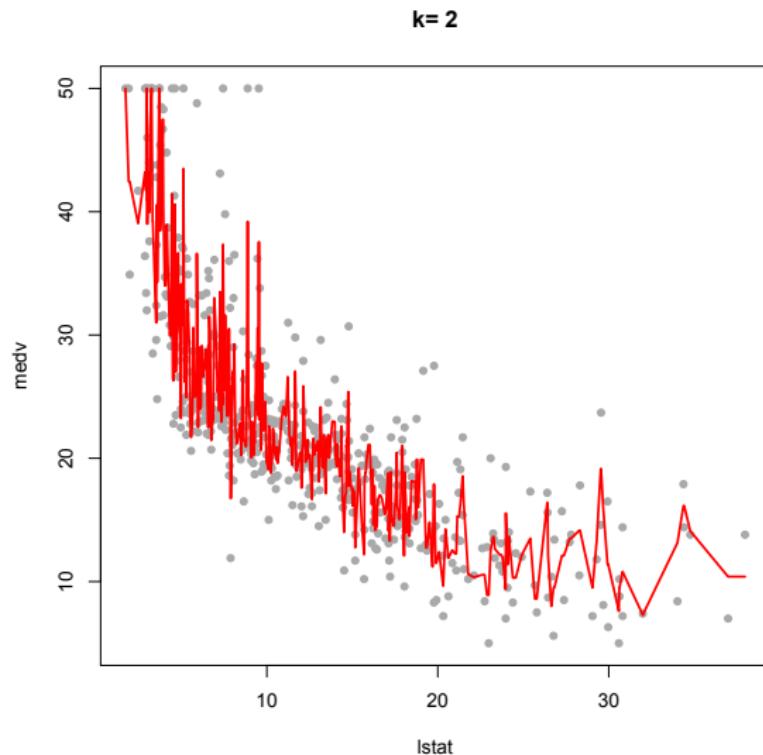
Back to Boston Housing

k= 20



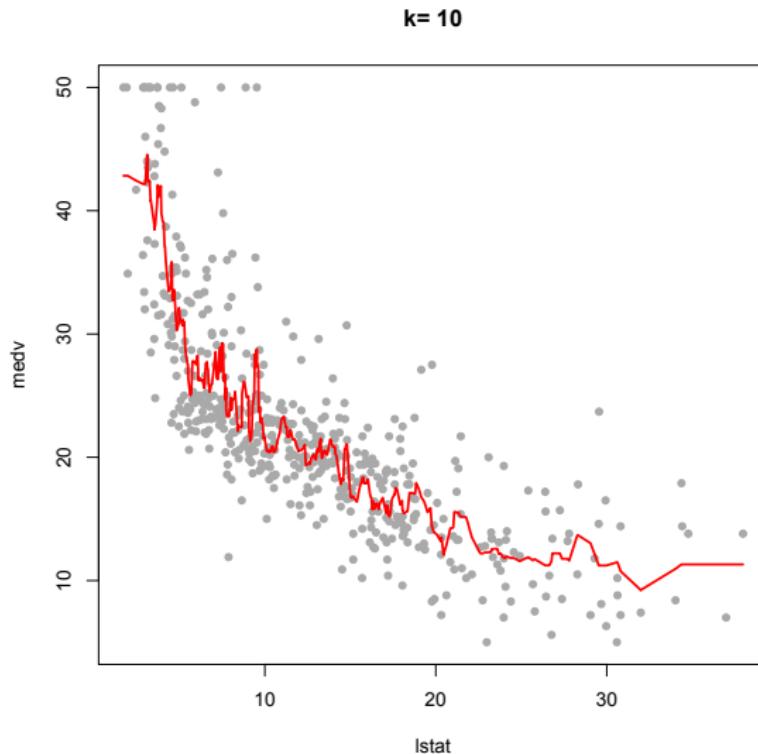
Back to Boston Housing

Okay, that seems sensible but why not use 2 neighbors or 200 neighbors?



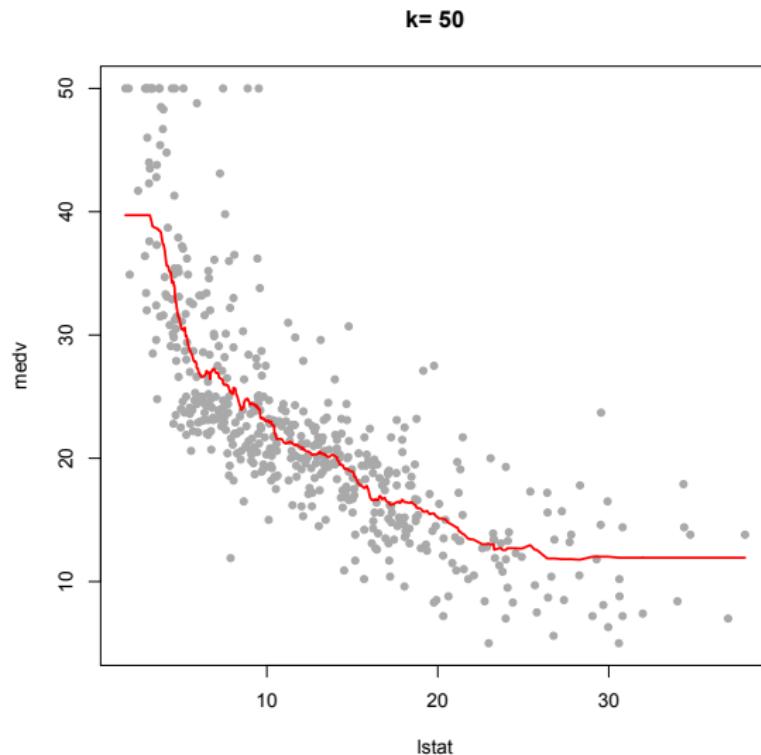
Back to Boston Housing

Okay, that seems sensible but why not use 5 neighbors or 200 neighbors?



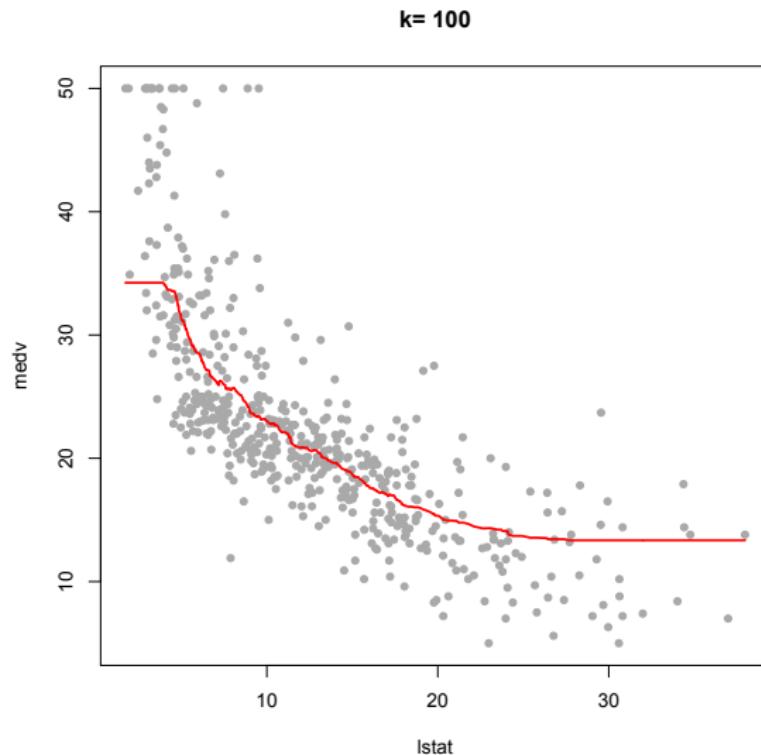
Back to Boston Housing

Okay, that seems sensible but why not use 5 neighbors or 200 neighbors?



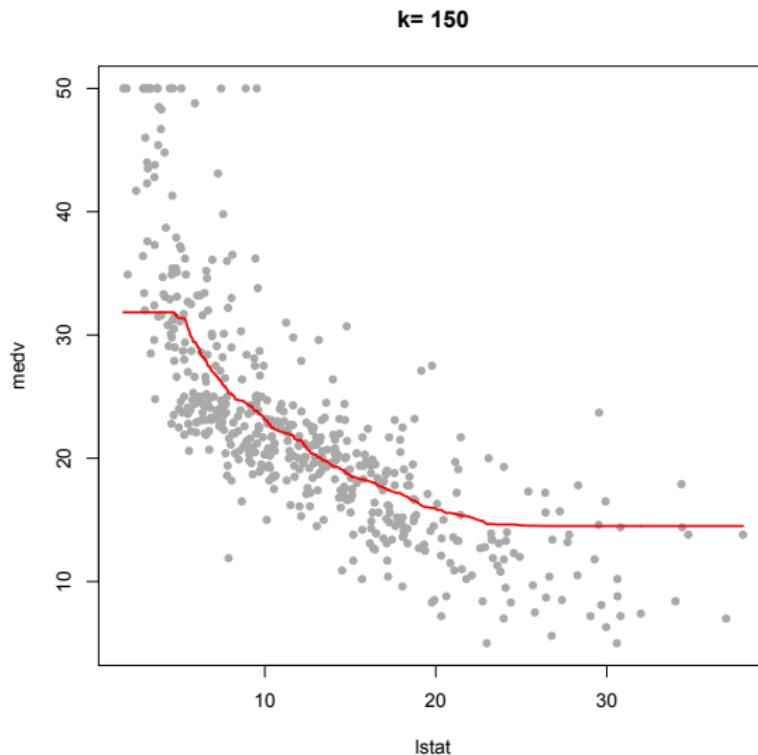
Back to Boston Housing

Okay, that seems sensible but why not use 5 neighbors or 200 neighbors?



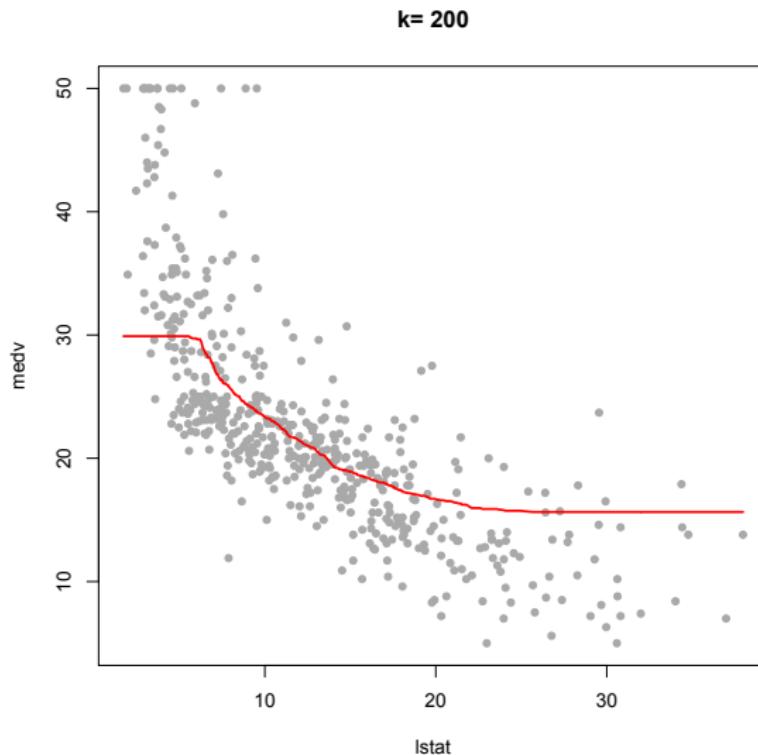
Back to Boston Housing

Okay, that seems sensible but why not use 5 neighbors or 200 neighbors?



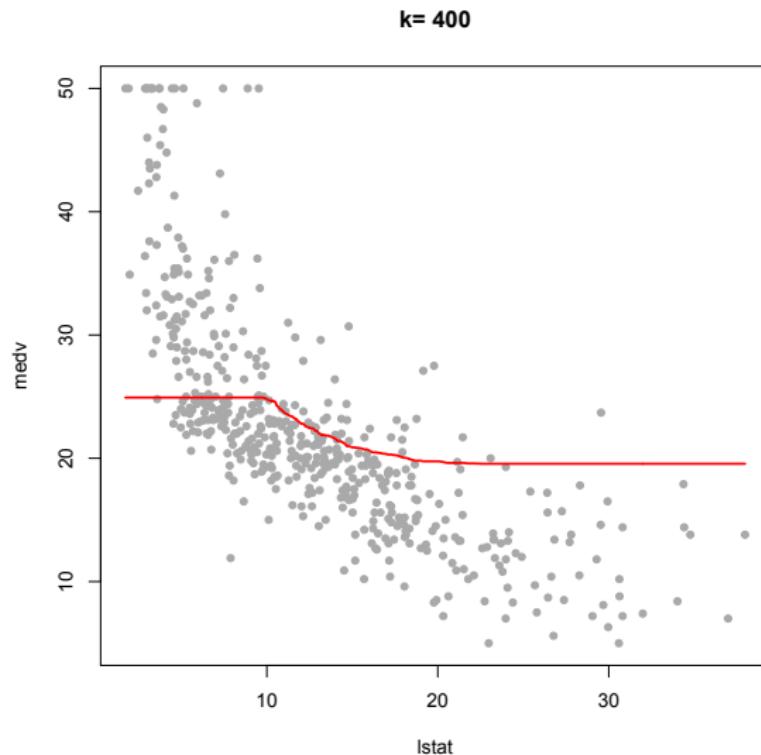
Back to Boston Housing

Okay, that seems sensible but why not use 5 neighbors or 200 neighbors?



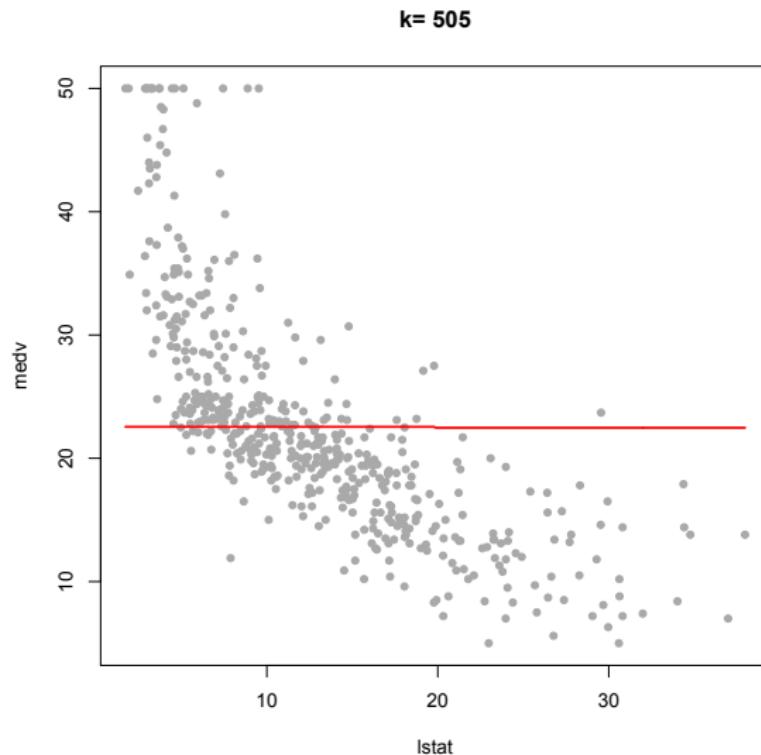
Back to Boston Housing

Okay, that seems sensible but why not use 5 neighbors or 200 neighbors?



Back to Boston Housing

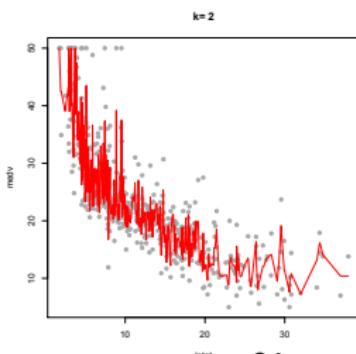
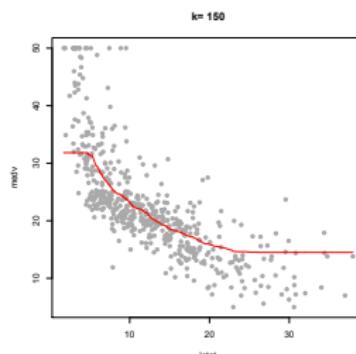
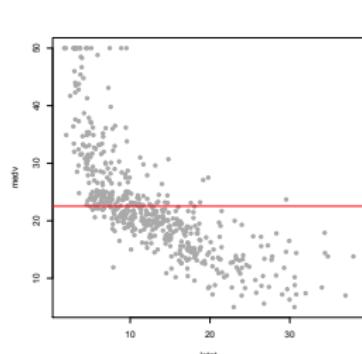
Okay, that seems sensible but why not use 5 neighbors or 200 neighbors?



Complexity, Generalization and Interpretation

- ▶ As we have seen in the examples above, there are lots of options in estimating $f(X)$.
- ▶ Some methods are very flexible some are not... *why would we ever choose a less flexible model?*
 1. Simple, more restrictive methods are usually easier to interpret
 2. More importantly, it is often the case that simpler models are **more accurate** in making future predictions.

Not too simple, but not too complex!



2. Measuring Accuracy

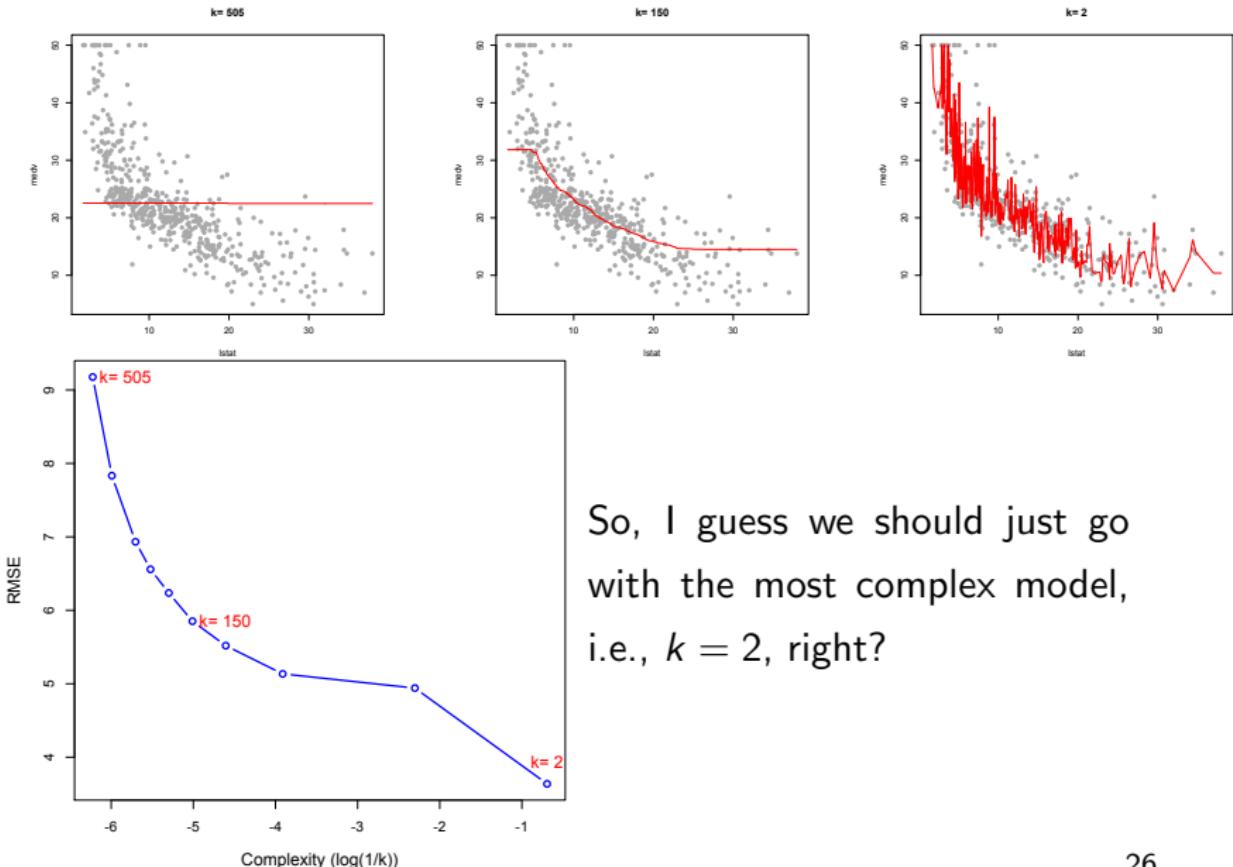
How accurate are each of these models?

Using the training data a standard measure of accuracy is the *root mean-squared error*

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n [Y_i - \widehat{f}(X_i)]^2}$$

This measure, on average, how large the “mistakes” (errors) made by the model are...

Measuring Accuracy (Boston housing, again)



3. Out-of Sample Predictions

But, do we really care about explaining what we have already seen?

Key Idea: what really matters is our prediction accuracy
out-of-sample!!!

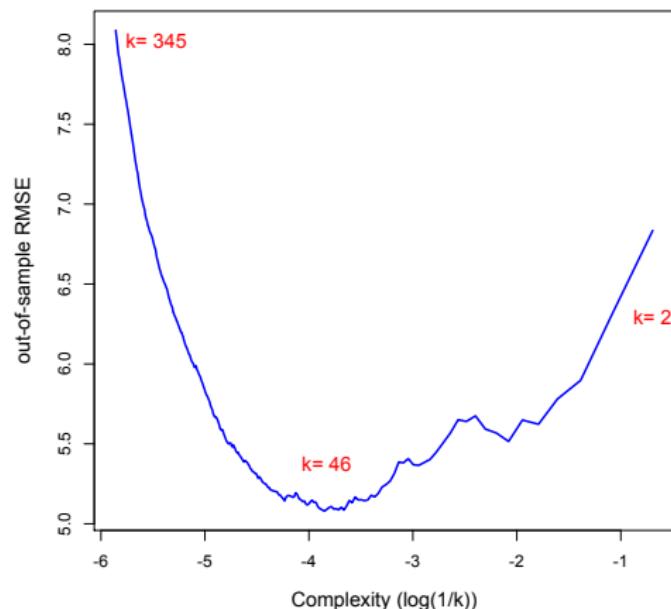
Suppose we have m additional observations (X_i^o, Y_i^o) , for $i = 1, \dots, m$, that we did not use to fit the model. Let's call this dataset the **validation set** (also known as *hold-out set* or *test set*)

Let's look at the out-of-sample RMSE:

$$RMSE^o = \sqrt{\frac{1}{m} \sum_{i=1}^m [Y_i^o - \widehat{f}(X_i^o)]^2}$$

Out-of Sample Predictions

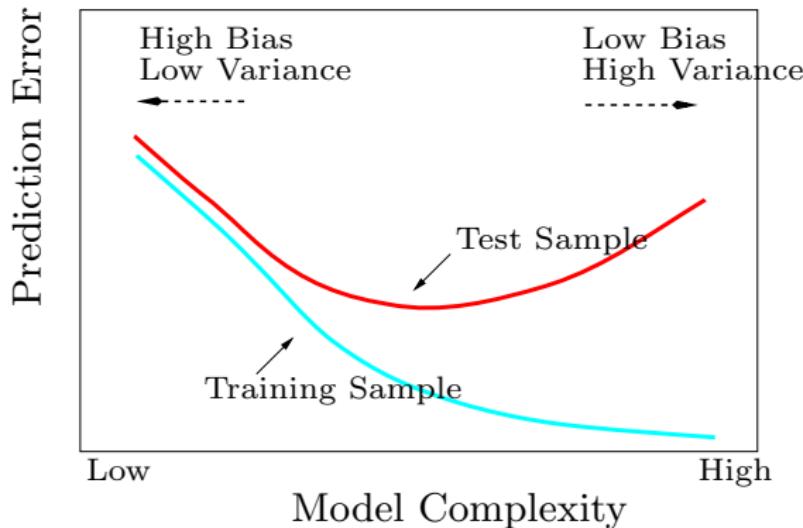
In our Boston housing example, I randomly chose a training set of size 400. I re-estimate the models using only this set and use the models to predict the remaining 106 observations (validation set)...



*Now, the model where
 $k = 46$ looks like the
most accurate choice!!*

**Not too simple but not
too complex!!!**

The Key Idea of the Course!!



This shows the typical behavior of the in and out-of-sample prediction error as a function of the complexity of the model... too flexible models will adapt itself too closely to the training set and will not generalize well, i.e., not be very good for the test data.

4. Bias-Variance Trade-Off

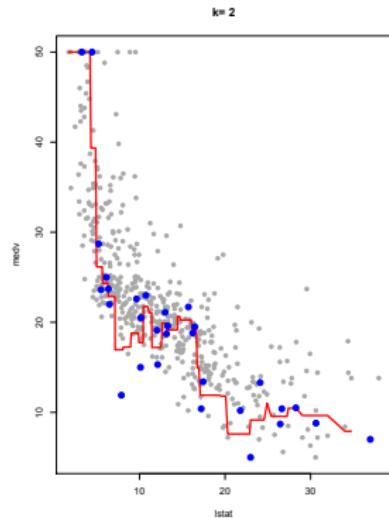
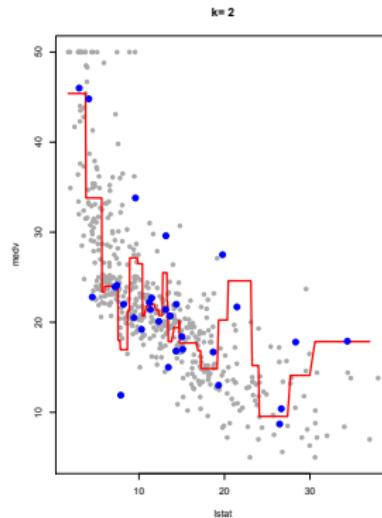
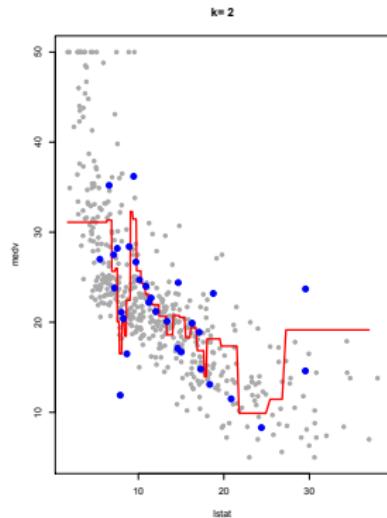
Why do complex models behave poorly in making predictions?

Let's start with an example...

- ▶ In the Boston housing example, I will randomly choose 30 observations to be in the training set 3 different times...
- ▶ for each training set I will estimate $f(\cdot)$ using the k -nearest neighbors idea... first with $k = 2$ and then with $k = 20$

Bias-Variance Trade-Off

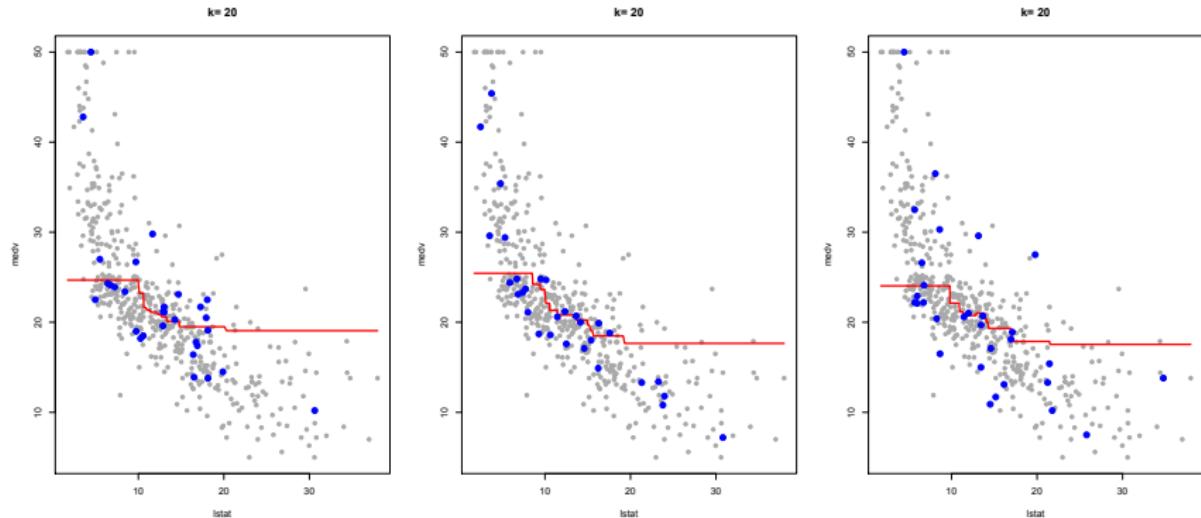
k=2 High variability...



(blue points are the training data used)

Bias-Variance Trade-Off

k=20 Low variability ... but BIAS!!



(blue points are the training data used)

Bias-Variance Trade-Off

What did we see here?

- ▶ When $k = 2$, it seems that the estimate of $f(\cdot)$ varies a lot between training sets...
- ▶ When $k = 20$ the estimates look a lot more stable...

Now, imagine that you are trying to predict *medv* when
Istat = 20... compare the changes in the predictions made by the
different training sets under $k = 2$ and $k = 20$... what do you see?

Bias-Variance Trade-Off

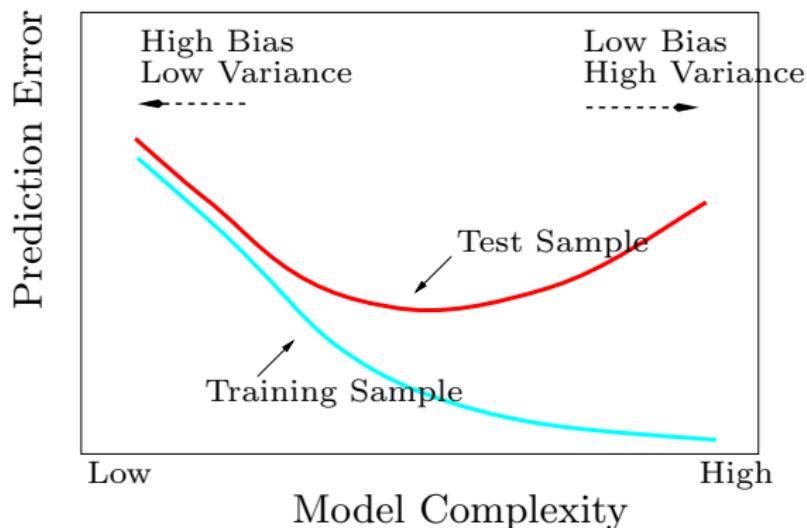
- ▶ This is an illustration of what is called the *bias-variance trade-off*.
- ▶ In general, simple models are trying to explain a complex, real problem with not a lot of flexibility so it introduces **bias**... on the other hand, by being simple the estimates tend to have low **variance**
- ▶ On the other hand, complex models are able to quickly adapt to the real situation and hence lead to small **bias**... however, by being too adaptable, it tends to vary a lot, i.e., high **variance**.

Bias-Variance Trade-Off

- ▶ In other words, we are trying to capture important patterns of the data that **generalize** to the future observations we are trying to predict. Models that are too simple are not able to capture relevant patterns and might have too big of a bias in predictions...
- ▶ Models that are too complex will “chase” irrelevant patterns in the training data that are not likely to exist in future data... so, it will lead to predictions that will vary a lot as things could change a lot depending on what sample we happen to see.
- ▶ Our goal is to find the sweet spot in the bias-variance trade-off!!

Bias-Variance Trade-Off

Once again, this is the key idea of the course!!



Bias-Variance Trade-Off

Let's get back to our original representation of the problem... it helps us understand what is going on...

$$Y_f = f(X_f) + \epsilon$$

- ▶ We need flexible enough models to find $f(\cdot)$ without imposing bias...
- ▶ ... but, too flexible models will “chase” non-existing patterns in ϵ leading to unwanted variability

Bias-Variance Trade-Off

A more detailed look at this idea... assume we are trying to make a prediction for Y_f using X_f and our inferred $\widehat{f(\cdot)}$. We hope to make small mistake measured by squared distance... Let's explore how our mistakes will behave *on average*.

$$\begin{aligned} \mathbb{E} \left[(Y_f - \widehat{f(X_f)})^2 \right] &= \left(f(X_f) - \mathbb{E} \left[\widehat{f(X_f)} \right] \right)^2 + \text{Var} \left[\widehat{f(X_f)} \right] + \text{Var}(\epsilon) \\ &= \text{Bias} \left[\widehat{f(X_f)} \right]^2 + \text{Var} \left[\widehat{f(X_f)} \right] + \text{Var}(\epsilon) \end{aligned}$$

hence, the *Bias-Variance Trade-Off!!*

5. Cross-Validation

- ▶ Using a validation-set to evaluate the performance of competing models has two potential drawbacks:
 1. the results can be highly dependent on the choice of the validation set... what samples? how many?
 2. by leaving aside a subset of data for validation we end up estimating the models with less information. It is harder to *learn* with fewer samples and this might lead to an overestimation of errors.
- ▶ Cross-Validation is a refinement of the validation strategy that help address both of these issues.

Leave-One-Out Cross-Validation (loocv)

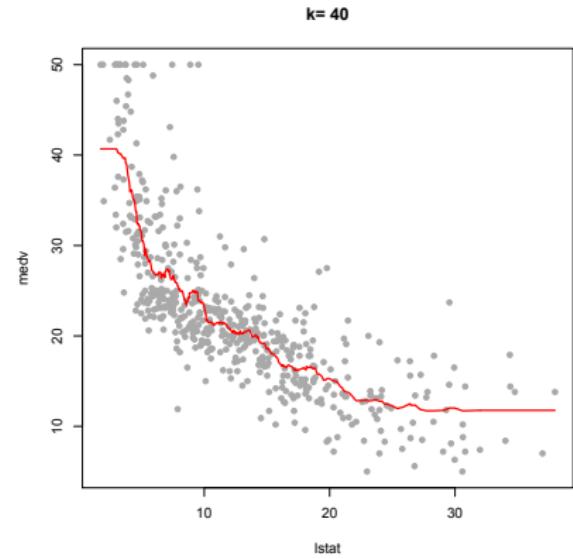
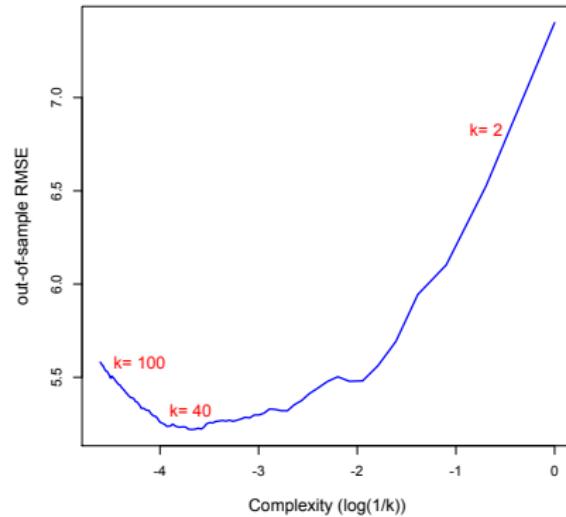
The name says it all!

- ▶ Assume we have n samples in our dataset. Define the validation set by choosing only **one sample**. Call it the i^{th} observation...
- ▶ The model is then trained in the remaining $n - 1$ samples and the results are used to predict the left-out sample. Compute the squared-error $MSE_i = (Y_i - \hat{Y}_i)^2$
- ▶ Repeat the procedure for every sample in the dataset (n times) and compute the average cross-validation MSE:

$$MSE^{loocv} = \frac{1}{n} \sum_{i=1}^n MSE_i$$

Leave-One-Out Cross-Validation

In the Boston data....



k-fold Cross Validation

LOOCV can be computationally expensive as each model being considered has to be estimated n times! A popular alternative is what is called *k*-fold Cross Validation.

- ▶ This approach randomly divides the original dataset into k groups of approximately the same size
- ▶ Choose one of the groups as a validation set. Estimate the models with the remaining $k - 1$ groups and predict the samples in the validation set. Compute the average squared-error for the samples in the validation set

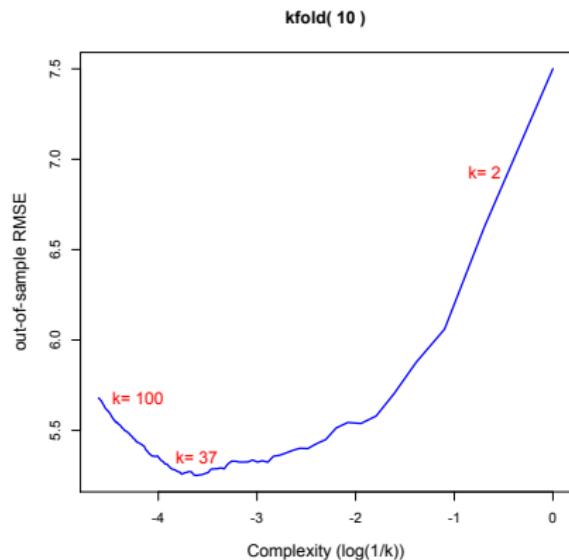
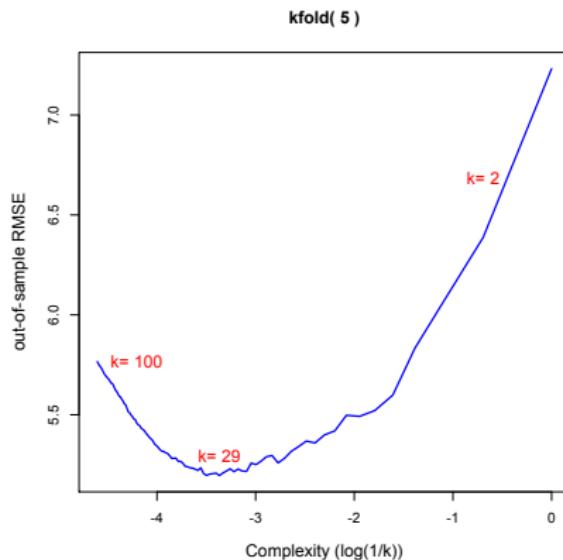
$$MSE_i = \text{Average} \left[(Y_i - \hat{Y}_i)^2 \right]$$

- ▶ Repeat the procedure for every *fold* in the dataset (k times) and compute the average cross-validation MSE:

$$MSE^{kcv} = \frac{1}{k} \sum_{i=1}^k MSE_i$$

k -fold Cross Validation

The usual choices are between $k = 5$ and $k = 10$...



knn: California Housing

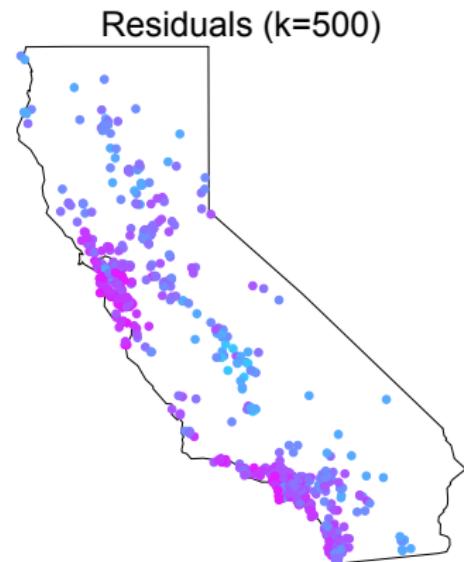
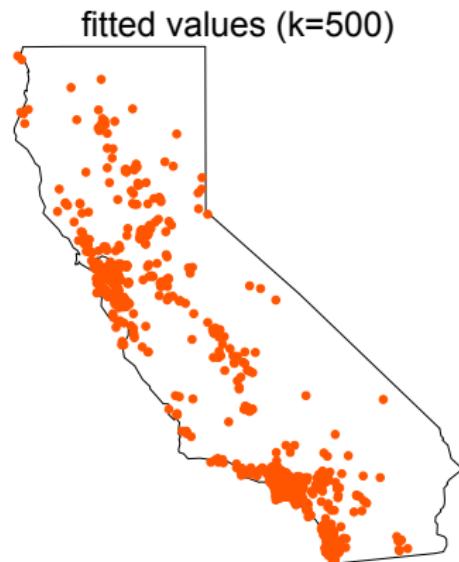
Data: Medium home values in census tract plus the following information:

- ▶ Location (latitude, longitude)
- ▶ Demographic information: population, income, etc...
- ▶ Average room/bedroom number, home age
- ▶ Let's start using just location as our X 's... euclidian distance is quite natural here, right?

Goal: Predict $\log(\text{MediumValue})$ (why logs? more on this later)

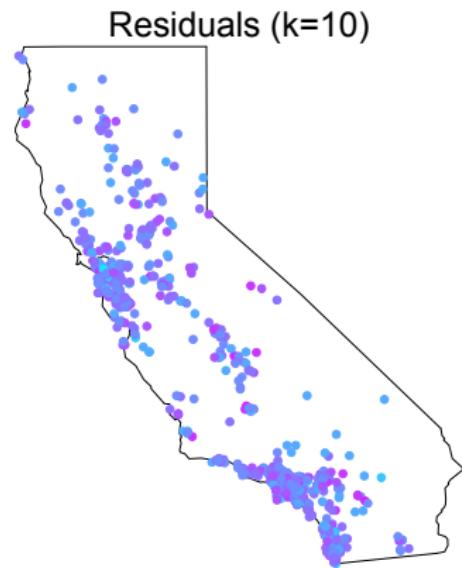
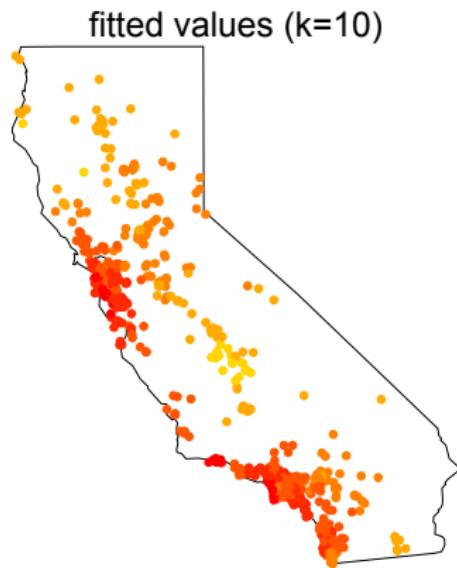
knn: California Housing

Using a training data of $n = 1000$ samples, here's a picture of the results for $k = 500$... left \hat{Y}_i , right $(Y_i - \hat{Y}_i)$



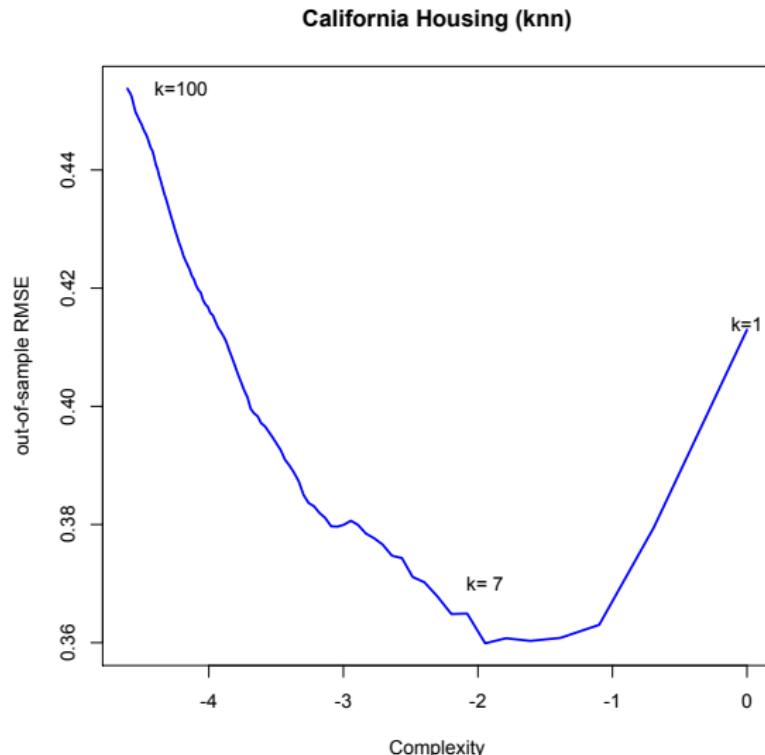
knn: California Housing

Now, $k = 10$... does this make sense?



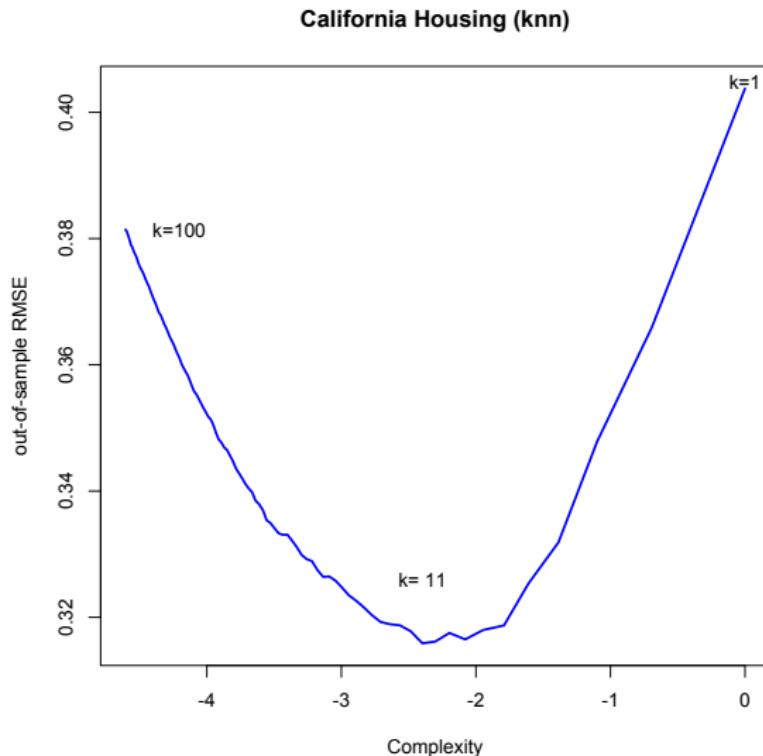
knn: California Housing

k -fold cross validation... $kfold = 10$



knn: California Housing

Adding Income as a predictor... think about the scale of X ...



knn: California Housing

Adding Income as a predictor... think about the scale of X ...

