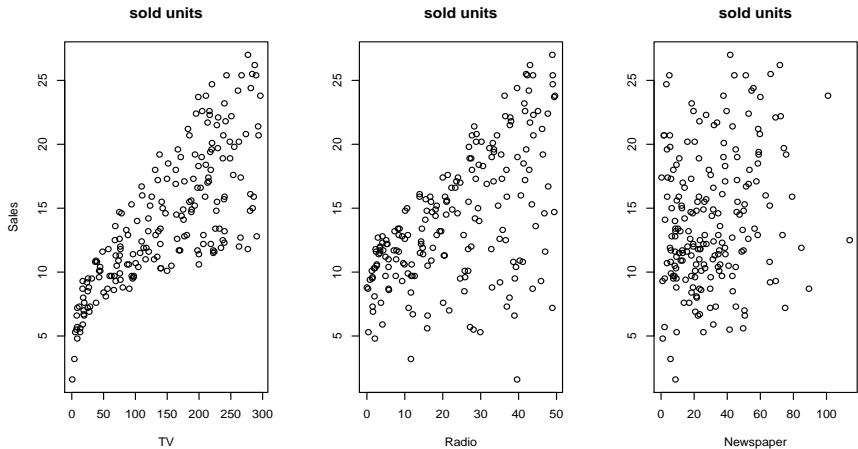


STA314, Lecture 1

September 6, 2019
Stanislav Volgushev

Motivating example (regression)

You work for a marketing firm. One of your clients collected data in 200 market locations. For each market, they have the amount of money spent on TV, Radio and newspaper advertisement (in 1000 \$) of a certain product and the sales (in thousands of units) of that product during a given time frame. They want you to help them predict sales and help decide how to advertise.



Some terminology

Typically, we will have a *data set (sample)* consisting of *data (observations)*
 $(x_1, Y_1), \dots, (x_n, Y_n)$

In *regression* Y_i are usually real values, x_i can be real values, vectors or other variables.
In *classification* (later in lectures) Y_i will take on a small number of different values.

x_i are called *predictors (covariates, regressors, features, input variables)*.

Y_i are called *response (outcome, target, dependent variable)*.

What are the aims in Statistical Learning?

Prediction: Given a new predictor value x_0 , predict response Y_0 .

- ▶ Example: assume we spend 151 thousand \$ on TV advertisement for a certain market. How many units do we expect to sell?
- ▶ Example: assume we spend 120 thousand \$ on TV advertisement, 15 thousand \$ on radio advertisement and 5 thousand \$ on newspaper advertisement. How many units do we expect to sell?

Inference: Understand the relationship between predictor and response in more detail and use that to make decisions. Quantify uncertainty.

- ▶ Example: does spending money of TV advertisement have any effect on sales?
- ▶ Example: assume we increase TV advertisement by 10 thousand \$. What is the effect on sales? How precisely can we quantify this effect?
- ▶ Example: assume we can spend a given budget on advertising on TV, radio and newspaper. How should we allocate the budget to maximise sales?

Regression models

Given observations $(x_1, Y_1), \dots, (x_n, Y_n)$ try to find relationship between x_i and Y_i
(example above: x_i is money spent on TV, Y_i are sales)

Machine Learning/Statistics: build a *regression model* to describe influence of x_i on Y_i .

Such models take the form (for the experts: this is the fixed design formulation)

$$Y_i = f(x_i) + \varepsilon_i, \quad \mathbb{E}[\varepsilon_i] = 0$$

where f is called *regression function* and ε_i are called *errors*.

- ▶ f describes the 'systematic' influence of x_i on Y_i .
- ▶ ε_i captures everything that f does not describe.
- ▶ $\mathbb{E}[\varepsilon_i] = 0$ implies: $f(x_i)$ is the average value of Y_i at predictor value x_i .

Remarks on ε_i

- ▶ In advertisement example: variation due unobserved factors such as size of market, local population, 'randomness' (ex.: weather for umbrella sales) etc.
- ▶ In Physics experiments: measurement error
- ▶ In Biology: genetic variation, environmental influence etc.

How do we use this model for prediction?

$$Y_i = f(x_i) + \varepsilon_i, \quad \mathbb{E}[\varepsilon_i] = 0$$

A typical way of measuring the quality of a prediction \hat{Y}_0 at point x_0 is the *mean squared error* (short: **MSE**)

$$MSE(\hat{Y}_0) = \mathbb{E}[(Y_0 - \hat{Y}_0)^2].$$

Which \hat{Y}_0 will minimise the MSE if f is known?

Some math (see blackboard in lectures for details): if $Y_0 = f(x_0) + \varepsilon_0$, x_0 fixed number and ε_0 independent of \hat{Y}_0 then

$$\mathbb{E}[(\hat{Y}_0 - Y_0)^2] = \mathbb{E}[(\hat{Y}_0 - f(x_0))^2] + \text{Var}(\varepsilon_0).$$

- ▶ $\text{Var}(\varepsilon_0)$ irreducible part. Even if we know f_0 this part can not be improved.
- ▶ $\mathbb{E}[(\hat{Y}_0 - f(x_0))^2] \geq 0$ depends on \hat{Y}_0 .
- ▶ MSE minimized at $\hat{Y}_0 = f(x_0)$
- ▶ If we know f our best prediction is $f(x_0)$. But we usually don't know f .

A first estimator for f : the K-nn method

The best possible value for \hat{Y}_0 is $\hat{Y}_0 = \mathbb{E}[Y_0] = f(x_0)$. Given data $(x_1, Y_1), \dots, (x_n, Y_n)$, how do we *learn/estimate* $f(x_0)$?

First: assume x_1, \dots, x_n take only values 0 or 1 and $x_0 = 0$. A natural approach:

$$\hat{f}(0) = \text{Average}(Y_i : x_i = 0) = \frac{\sum_{i=1}^n Y_i I\{x_i = 0\}}{\sum_{i=1}^n I\{x_i = 0\}}$$

Here

$$I\{x_i = 0\} = \begin{cases} 1, & \text{if } x_i = 0 \\ 0, & \text{otherwise} \end{cases}$$

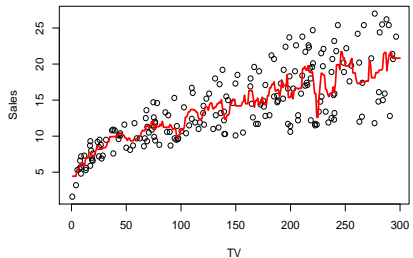
Problem: this works in data sets where x_1, \dots, x_n take only few distinct values and we are interested in predicting outcome for one of those values. Example: advertisement data set has no point with $x = 152$.

Idea: instead of requiring $x_i = x_0$ take K of the 'closest' x_i . **K-nn** (K nearest neighbours) method.

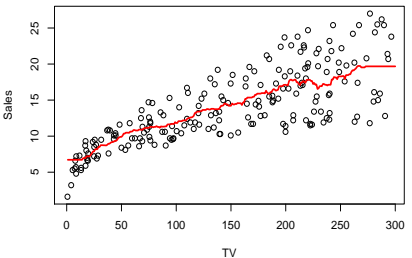
$$\hat{f}(x_0) = \frac{\sum_{i=1}^n Y_i I\{x_i \text{ among closest } K \text{ to } x_0\}}{K}$$

Examples of K-nn regression for regressing Sales on TV advertisement

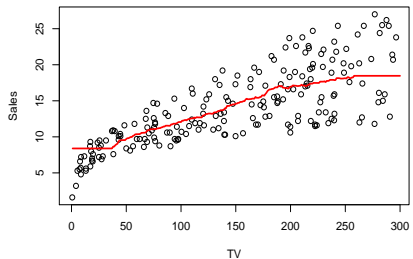
K = 5



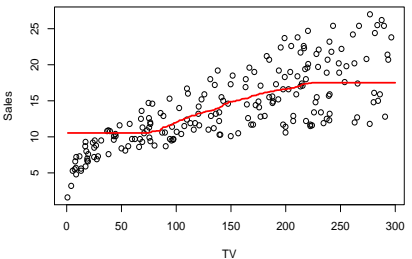
K = 25



K = 50



K = 100



Let's look at some **R** code: create plots on previous slide in **R**.

Some useful R functions you should be familiar with after today's lecture (and tutorials next week):

`read.csv(file = 'file name')`: read csv file into **R**

arrays, lists, data frames: working in **R**, details in tutorials.

`install.packages('package name')`: install a library from the internet

`library('library name')`: load an installed library

`knnreg`: function in library `caret` for doing k-nn regression

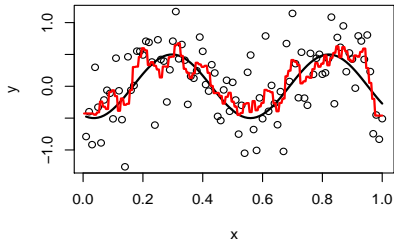
`predict`: function for predicting, we will encounter this a lot

`set.seed`: make results involving 'random' numbers reproducible

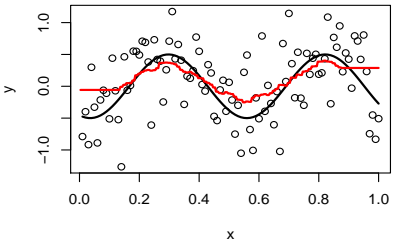
`plot`, `lines`: for making plots

Examples of K-nn with simulated data

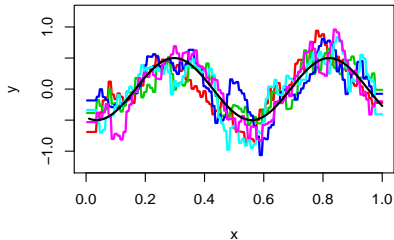
$f(x) = \sin(-2+12*x)$, $K = 5$



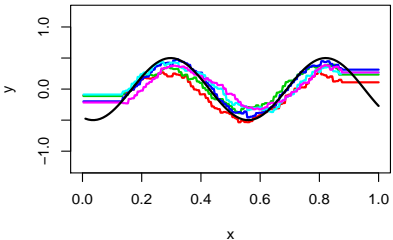
$f(x) = \sin(-2+12*x)$, $K = 25$



$K=5$



$K=25$



Bias-variance decomposition

- ▶ For smaller K , regression function is very 'wiggly'. A lot of variation between samples generated from the same model.
- ▶ If K is very large, data are not described well. Less variation between samples generated from the same model.
- ▶ Intermediate values of K seem to be sensible.

This can be formalized through the concepts of *bias* and *variance*. Recall that the irreducible part of the MSE takes the form $\mathbb{E}[(\hat{Y}_0 - f(x_0))^2]$. After some calculations (see blackboard in lectures for details)

$$\mathbb{E}[(\hat{Y}_0 - f(x_0))^2] = \text{Var}(\hat{Y}_0) + \{\mathbb{E}[\hat{Y}_0] - f(x_0)\}^2.$$

- ▶ $\mathbb{E}[\hat{Y}_0] - f(x_0)$ is called **bias**. It describes how far from the truth the prediction \hat{Y}_0 is *on average*.
- ▶ $\text{Var}(\hat{Y}_0)$ describes how much variation the estimator \hat{Y}_0 has.
- ▶ Ideal estimator would have small bias and small variance, but usually that is impossible.

The bias and variance of K-nn

$Y_i = f(x_i) + \varepsilon_i$, ε_i i.i.d. with $\mathbb{E}[\varepsilon_i] = 0$, $\text{Var}(\varepsilon_i) = \sigma^2$, x_i fixed, $\hat{f}(x_0)$ K-nn estimator. Then, if all x_i take distinct values, (see derivations in lectures for details)

$$\text{Var}[\hat{f}(x_0)] = \sigma^2/K,$$

Bias: more complicated. For simplicity: $x_i = i/n$, $i = 1, \dots, n$, $f : [0, 1] \rightarrow \mathbb{R}$ two times continuously differentiable, $K = 2\ell + 1$, $x_0 = j/n$ with $\ell < j < n - \ell$. Then (see derivations in lectures for details)

$$\mathbb{E}[\hat{f}(x_0)] = \frac{1}{K} \sum_{u=-\ell}^{\ell} f\left(\frac{j+u}{n}\right) \approx f(x_0) + \frac{1}{24}f''(x_0)(K/n)^2 + r_{K,n}$$

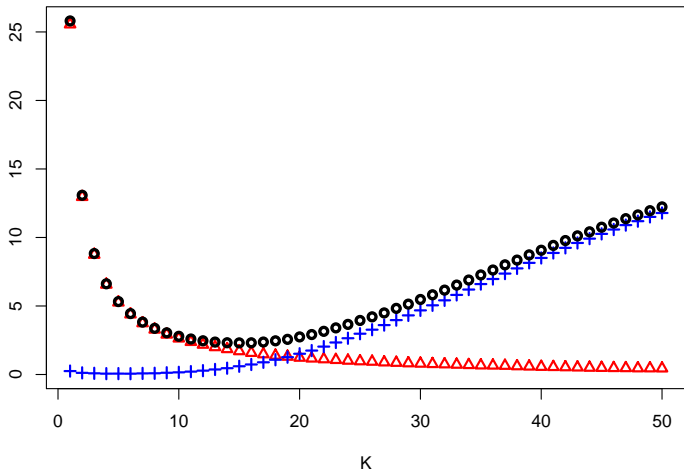
where $r_{K,n}$ remainder term, 'small' under suitable conditions (more details: homework).

- ▶ Variance decreases as K increases
- ▶ Squared bias increases as K increases
- ▶ To get a good MSE we have to balance bias and variance. This is called *bias-variance trade-off*

Remark: if we let $\ell = \ell_n$ with $\ell_n/n \rightarrow 0$ in the analysis above this would be 'Asymptotic Statistics'. There will be some advanced homework related to that.

Let's look at some R code for running simulations

Bias and variance of K-nn: a simulated example



Average Bias, Variance and MSE (which curve is which?) of K-nn for $f(x) = 0.5 \sin(-2 + 12x)$ as a function of K .

Main message from previous slides: need to select a 'good' value of K to obtain good predictions. How do we go about that?

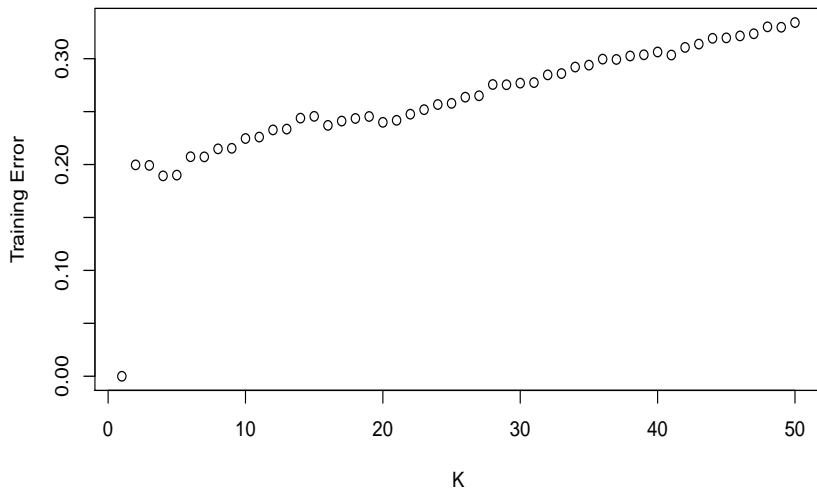
More general question: given any method for prediction, how do we evaluate the performance of our method?

First try: utilize available data set $(x_1, Y_1), \dots, (x_n, Y_n)$ to compute *training error* (*in-sample prediction error*) :

$$MSE_{train} = \frac{1}{n} \sum_{i=1}^n (\hat{f}(x_i) - Y_i)^2$$

and use this to measure quality of our prediction. Let's see how it works with K-nn.

The training error of K-nn regression... is not very useful!



Problem: same data are used to compute \hat{f} and to compute error. This leads to *overfitting* (following the data too closely).