

## Machine Learning for Official Statistics and SDGs

Statistical learning:  
*You've seen this before*



# WHO IS THIS COURSE FOR?

# WHO IS THIS COURSE FOR?

- ▶ This course is for (present or future) Data Scientists

# WHO IS THIS COURSE FOR?

- ▶ This course is for (present or future) Data Scientists
- ▶ This course uses both statistical and computational concepts

# WHO IS THIS COURSE FOR?

- ▶ This course is for (present or future) Data Scientists
- ▶ This course uses both statistical and computational concepts
- ▶ This course uses many applied examples and a very progressive approach

# WHO IS THIS COURSE FOR?

- ▶ This course is for (present or future) Data Scientists
- ▶ This course uses both statistical and computational concepts
- ▶ This course uses many applied examples and a very progressive approach

*Data Scientist: "Person who is better at statistics than any software engineer and better at software than any statistician."*

J. Wills (2012)

# WHAT IS STATISTICAL LEARNING?

# WHAT IS STATISTICAL LEARNING?

*“ Statistical learning refers to a vast set of tools for understanding data”*

Gareth James, Daniela Witten, Trevor Hastie , Robert Tibshirani (2021)



# WHAT IS STATISTICAL LEARNING?

*“ Statistical learning refers to a vast set of tools for understanding data”*

Gareth James, Daniela Witten, Trevor Hastie , Robert Tibshirani (2021)

- Involves building statistical models

# WHAT IS STATISTICAL LEARNING?

*“ Statistical learning refers to a vast set of tools for understanding data”*

Gareth James, Daniela Witten, Trevor Hastie , Robert Tibshirani (2021)

- ▶ Involves building statistical models
- ▶ Goals are estimation or prediction

# WHAT IS STATISTICAL LEARNING?

Two main learning problems:

# WHAT IS STATISTICAL LEARNING?

Two main learning problems:

- ▶ We observe **both** the outcome  $y$  and regressors (also called features)  $x$ s

# WHAT IS STATISTICAL LEARNING?

Two main learning problems:

- We observe **both** the outcome  $y$  and regressors (also called features)  $x$ s

↪ **Supervised** learning

# WHAT IS STATISTICAL LEARNING?

Two main learning problems:

- We observe **both** the outcome  $y$  and regressors (also called features)  $x$ s

↪ **Supervised** learning

*Most of the examples and applications are supervised learning*

# WHAT IS STATISTICAL LEARNING?

Two main learning problems:

- ▶ We observe **both** the outcome  $y$  and regressors (also called features)  $x$ s

↪ **Supervised** learning

*Most of the examples and applications are supervised learning*

- ▶ We **do not** observe the outcome  $y$  but **only** several  $x$ s

# WHAT IS STATISTICAL LEARNING?

Two main learning problems:

- ▶ We observe **both** the outcome  $y$  and regressors (also called features)  $x$ s

↪ **Supervised** learning

*Most of the examples and applications are supervised learning*

- ▶ We **do not** observe the outcome  $y$  but **only** several  $x$ s

↪ **Unsupervised** learning (or *cluster analysis*)



# WHAT IS STATISTICAL LEARNING?

Two main learning problems:

- ▶ We observe **both** the outcome  $y$  and regressors (also called features)  $x$ s

↪ **Supervised** learning

*Most of the examples and applications are supervised learning*

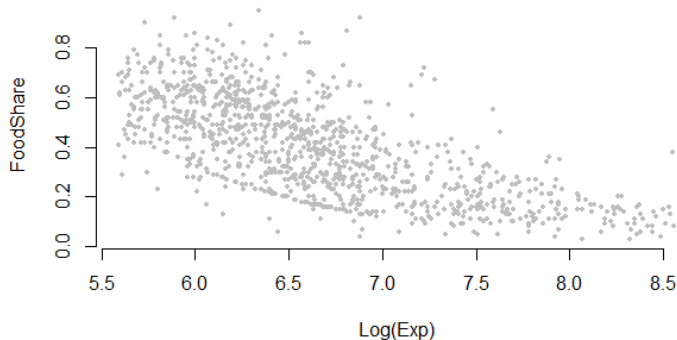
- ▶ We **do not** observe the outcome  $y$  but **only** several  $x$ s

↪ **Unsupervised** learning (or *cluster analysis*)

*More complex models we'll see at the end*

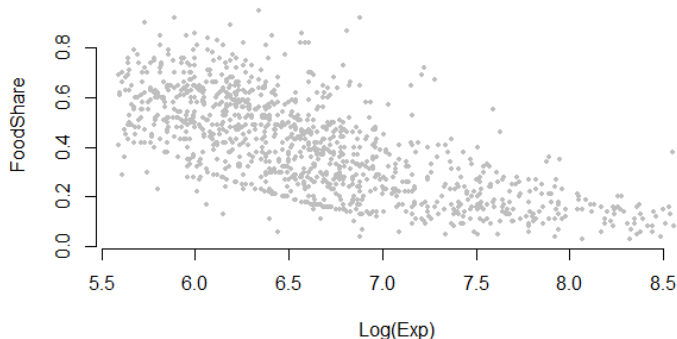
# STATISTICAL LEARNING ON AN EXAMPLE

**Scatter plot of Food Share vs Log(exp)**



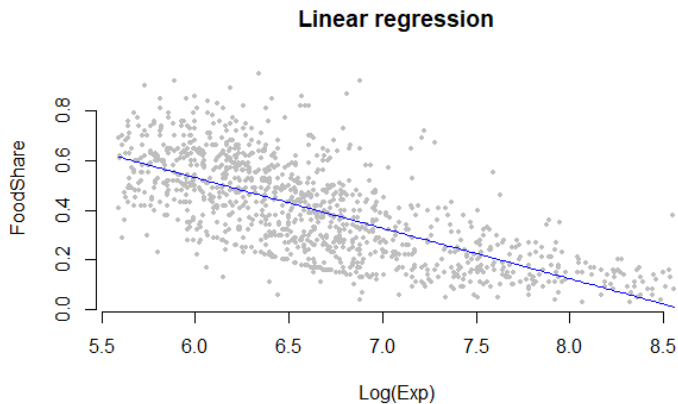
# STATISTICAL LEARNING ON AN EXAMPLE

**Scatter plot of Food Share vs Log(exp)**

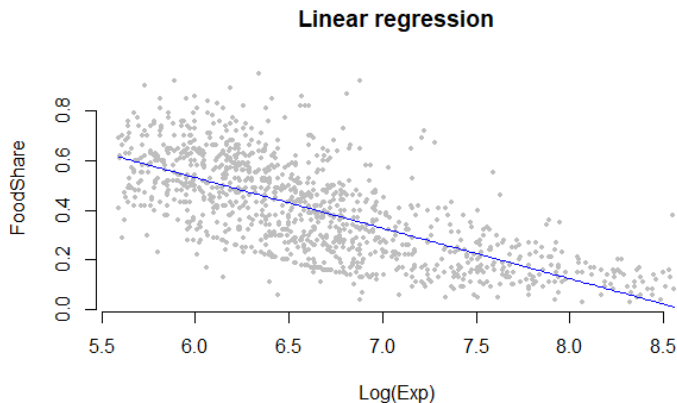


We may be interested in the relationship between the two variables

UNDERSTANDING = ESTIMATE  $f(\cdot)$



UNDERSTANDING = ESTIMATE  $f(\cdot)$



$f(\cdot)$  is the regression line

# WHY ESTIMATING $f(\cdot)$ ?

## ► Inference

# WHY ESTIMATING $f(\cdot)$ ?

## ► Inference

Understand the nature of the relationship between  $X$  and  $Y$

# WHY ESTIMATING $f(\cdot)$ ?

## ► Inference

Understand the nature of the relationship between  $X$  and  $Y$   
Identify "important" variables to understand  $Y$



# WHY ESTIMATING $f(\cdot)$ ?

## ► Inference

Understand the nature of the relationship between  $X$  and  $Y$

Identify "important" variables to understand  $Y$

## ► Prediction

# WHY ESTIMATING $f(\cdot)$ ?

## ► Inference

Understand the nature of the relationship between  $X$  and  $Y$

Identify "important" variables to understand  $Y$

## ► Prediction

Predict  $y$  for any new  $x$  using  $f(\cdot)$

# WHY ESTIMATING $f(\cdot)$ ?

- ▶ Inference

  - Understand the nature of the relationship between  $X$  and  $Y$

  - Identify "important" variables to understand  $Y$

- ▶ Prediction

  - Predict  $y$  for any new  $x$  using  $f(\cdot)$

- ▶ In practice we must estimate  $f(\cdot)$  using a model:

# WHY ESTIMATING $f(\cdot)$ ?

## ► Inference

Understand the nature of the relationship between  $X$  and  $Y$

Identify "important" variables to understand  $Y$

## ► Prediction

Predict  $y$  for any new  $x$  using  $f(\cdot)$

## ► In practice we must estimate $f(\cdot)$ using a model:

$$y = f(x) + \varepsilon$$

# WHY ESTIMATING $f(\cdot)$ ?

## ► Inference

Understand the nature of the relationship between  $X$  and  $Y$

Identify "important" variables to understand  $Y$

## ► Prediction

Predict  $y$  for any new  $x$  using  $f(\cdot)$

## ► In practice we must estimate $f(\cdot)$ using a model:

$$y = f(x) + \varepsilon$$

We denote by  $\widehat{f(\cdot)}$  the estimate of  $f(\cdot)$

# HOW TO ESTIMATE $f(\cdot)$ ?

- ▶ Parametric methods

# HOW TO ESTIMATE $f(\cdot)$ ?

- ▶ Parametric methods

Specify a form for  $f(\cdot)$ , for example linear:

# HOW TO ESTIMATE $f(\cdot)$ ?

## ► Parametric methods

Specify a form for  $f(\cdot)$ , for example linear:

$$y = \beta_0 + \beta_1 x + \varepsilon$$



# HOW TO ESTIMATE $f(\cdot)$ ?

## ► Parametric methods

Specify a form for  $f(\cdot)$ , for example linear:

$$y = \beta_0 + \beta_1 x + \varepsilon$$

- The goal is to find the line that is **minimizing** the distance to the observed points  $(x_i, y_i)$ . The distance is computed as the Mean Square Error (MSE):

$$MSE(\beta_0, \beta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

# HOW TO ESTIMATE $f(\cdot)$ ?

- ▶ Parametric methods

Specify a form for  $f(\cdot)$ , for example linear:

$$y = \beta_0 + \beta_1 x + \varepsilon$$

- ▶ The goal is to find the line that is **minimizing** the distance to the observed points  $(x_i, y_i)$ . The distance is computed as the Mean Square Error (MSE):

$$MSE(\beta_0, \beta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

- ▶ The regression line, defined by  $\beta_0$  and  $\beta_1$ , is simply the solution of:

$$\text{Min}_{(\beta_0, \beta_1)} MSE(\beta_0, \beta_1)$$

# HOW TO ESTIMATE $f(\cdot)$ ?

- ▶ Parametric methods

Specify a form for  $f(\cdot)$ , for example linear:

$$y = \beta_0 + \beta_1 x + \varepsilon$$

- ▶ The goal is to find the line that is **minimizing** the distance to the observed points  $(x_i, y_i)$ . The distance is computed as the Mean Square Error (MSE):

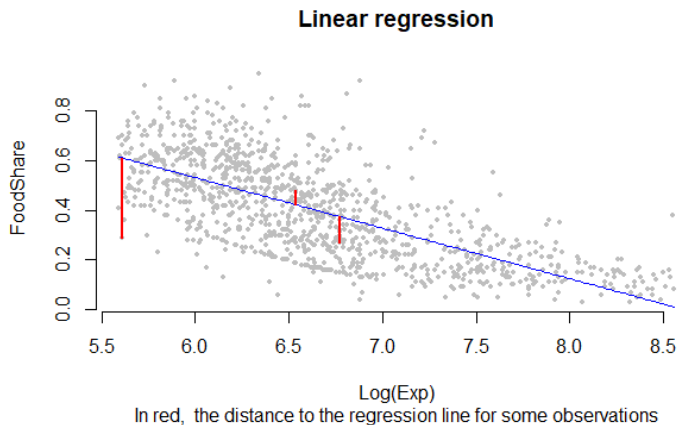
$$MSE(\beta_0, \beta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

- ▶ The regression line, defined by  $\beta_0$  and  $\beta_1$ , is simply the solution of:

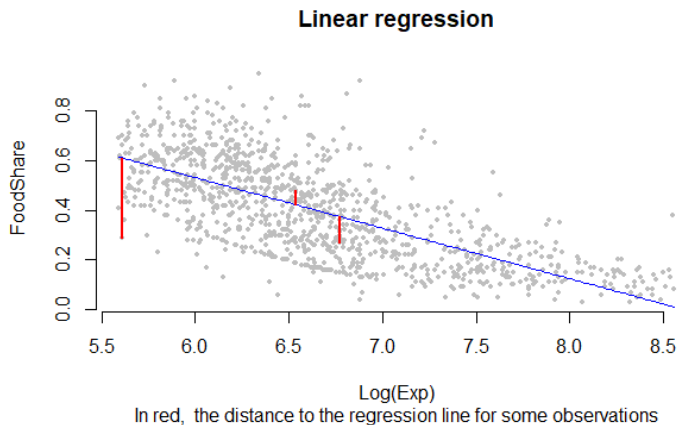
$$\text{Min}_{(\beta_0, \beta_1)} MSE(\beta_0, \beta_1)$$

- ▶ The MSE it is the *cost function* used to estimate  $(\hat{\beta}_0, \hat{\beta}_1)$

# HOW TO ESTIMATE $f(\cdot)$ : IN PRACTICE



# HOW TO ESTIMATE $f(\cdot)$ : IN PRACTICE



The regression line is found by minimizing the sum of all distances or **MSE**

RESULTS:  $\widehat{f(\cdot)}$ 

From the result and the estimated parameters  $(\widehat{\beta}_0, \widehat{\beta}_1)$ , we see that there is a relation, and that it is decreasing.

	Estimate	Std. Error	t value	$Pr(>  t )$
(Intercept)	1.75	0.04	41.09	0
ltexp	-0.20***	0.01	-31.84	0

The quality of the adjustment may be measured by the  $R^2 = 0.478$

# WHAT DID WE LEARN?

- ▶ Is the line fitting the data?

# WHAT DID WE LEARN?

- ▶ Is the line fitting the data?

One may always find a relation between two variables



# WHAT DID WE LEARN?

- Is the line fitting the data?

One may always find a relation between two variables

Statistical test help asses the significance of a variable

# WHAT DID WE LEARN?

- ▶ Is the line fitting the data?

One may always find a relation between two variables

Statistical test help asses the significance of a variable

- ▶ How good is the linear adjustment:  $R^2$

# WHAT DID WE LEARN?

- ▶ Is the line fitting the data?

One may always find a relation between two variables  
Statistical test help asses the significance of a variable

- ▶ How good is the linear adjustment:  $R^2$

$$R^2 = \frac{TSS - RSS}{TSS}$$

# WHAT DID WE LEARN?

- Is the line fitting the data?

One may always find a relation between two variables  
Statistical test help asses the significance of a variable

- How good is the linear adjustment:  $R^2$

$$R^2 = \frac{TSS - RSS}{TSS}$$

with:  $TSS = \sum_{i=1}^n (y_i - \bar{y})^2$  and  $RSS = \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$

# WHAT DID WE LEARN?

- Is the line fitting the data?

One may always find a relation between two variables  
Statistical test help asses the significance of a variable

- How good is the linear adjustment:  $R^2$

$$R^2 = \frac{TSS - RSS}{TSS}$$

with:  $TSS = \sum_{i=1}^n (y_i - \bar{y})^2$  and  $RSS = \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$

- $R^2$  is a very popular measure. The closer to 1, the better.

# WHAT DID WE LEARN?

- ▶ Is the line fitting the data?

One may always find a relation between two variables  
Statistical test help asses the significance of a variable

- ▶ How good is the linear adjustment:  $R^2$

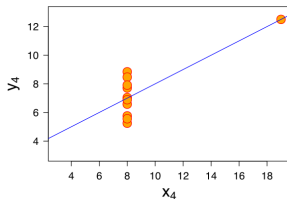
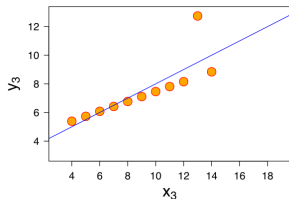
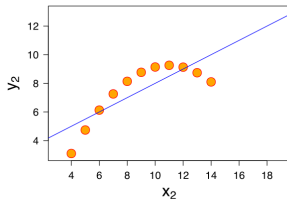
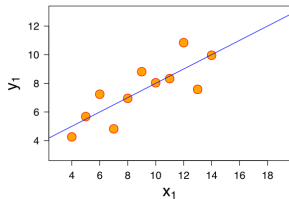
$$R^2 = \frac{TSS - RSS}{TSS}$$

with:  $TSS = \sum_{i=1}^n (y_i - \bar{y})^2$  and  $RSS = \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$

- ▶  $R^2$  is a very popular measure. The closer to 1, the better.
- ▶  $R^2$  can be very misleading

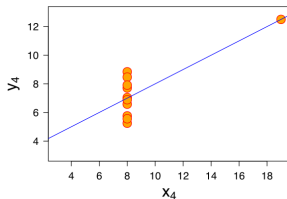
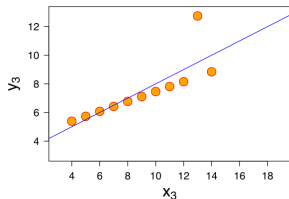
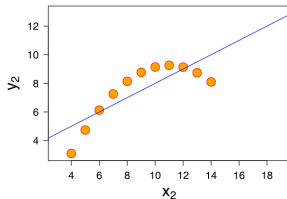
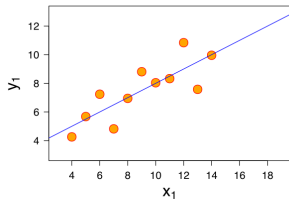
# BEWARE OF $R^2$ : ANSCOMBE QUARTET (1973)

# BEWARE OF $R^2$ : ANSCOMBE QUARTET (1973)





# BEWARE OF $R^2$ : ANSCOMBE QUARTET (1973)



In all these data sets the  $R^2$  is 0.67

# PRACTICE OF STATISTICAL LEARNING

*“ ...make both calculations and graphs. Both sorts of output should be studied; each will contribute to understanding.”*

F. J. Anscombe (1973)

# PRACTICE OF STATISTICAL LEARNING

*“ ...make both calculations and graphs. Both sorts of output should be studied; each will contribute to understanding.”*

F. J. Anscombe (1973)

See also the datasaurus

# PRACTICE OF STATISTICAL LEARNING

*“ ...make both calculations and graphs. Both sorts of output should be studied; each will contribute to understanding.”*

F. J. Anscombe (1973)

See also the datasaurus

- Practice of Statistical learning can be challenging

# PRACTICE OF STATISTICAL LEARNING

*“ ...make both calculations and graphs. Both sorts of output should be studied; each will contribute to understanding.”*

F. J. Anscombe (1973)

See also the datasaurus

► Practice of Statistical learning can be challenging

↪ Need to compute good indicators

# PRACTICE OF STATISTICAL LEARNING

*“ ...make both calculations and graphs. Both sorts of output should be studied; each will contribute to understanding.”*

F. J. Anscombe (1973)

See also the datasaurus

- Practice of Statistical learning can be challenging
  - ↪ Need to compute good indicators
  - ↪ Need to understand the indicators computed

# PRACTICE OF STATISTICAL LEARNING

*“ ...make both calculations and graphs. Both sorts of output should be studied; each will contribute to understanding.”*

F. J. Anscombe (1973)

See also the datasaurus

- Practice of Statistical learning can be challenging
  - ↪ Need to compute good indicators
  - ↪ Need to understand the indicators computed
  - ↪ Need to go beyond linearity

# Machine Learning for Official Statistics and SDGs

## Statistical learning: *Beyond linearity*



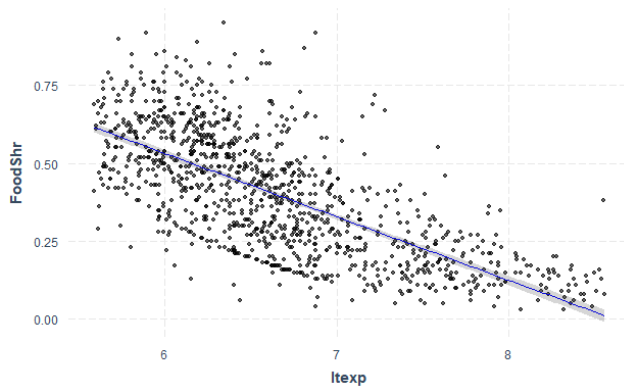


# BEYOND LINEARITY

# BEYOND LINEARITY

- A linear model may be unadapted or too simple

$$y = \beta_0 + \beta_1 x + \varepsilon$$

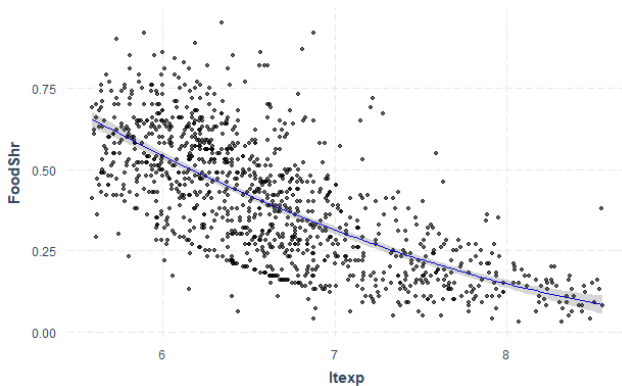


The fit (measured by  $R^2$ ) is:  $R^2 = 0.478$

## BEYOND LINEARITY

- A Polynomial model may be better adapted: **Quadratic** model

$$y = \beta_0 + \beta_1x + \beta_2x^2 + \varepsilon$$

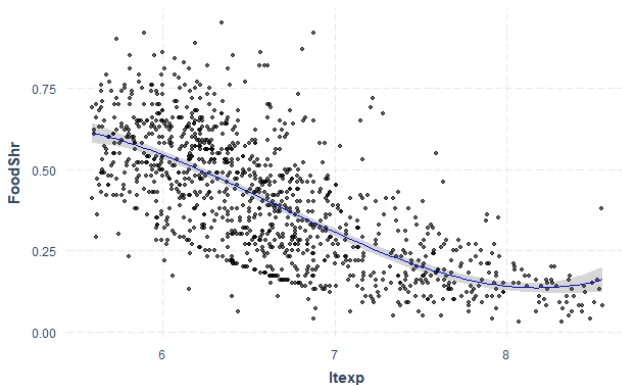


Do we have a better fit?  $R^2 = 0.484$

# BEYOND LINEARITY

- Polynomial may be better adapted: **Cubic** model

$$y = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3 + \varepsilon$$



Do we have a better fit?  $R^2 = 0.490$

# IN PRACTICE

- ▶ Polynomial models may be useful

# IN PRACTICE

- Polynomial models may be useful

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \cdots + \beta_p x^p + \varepsilon$$

# IN PRACTICE

- Polynomial models may be useful

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \cdots + \beta_p x^p + \varepsilon$$

How to choose the degree  $p$  ?

# IN PRACTICE

- Polynomial models may be useful

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \cdots + \beta_p x^p + \varepsilon$$

How to choose the degree  $p$  ?

Collinearity of  $x^p$  and  $x^q$  for  $p \neq q$ ?



# PARAMETRIC VS NONPARAMETRIC MODELS

- ▶ Linear and polynomial models are determined by parameters  $(\beta_0, \beta_2, \dots, \beta_p)$

# PARAMETRIC VS NONPARAMETRIC MODELS

- ▶ Linear and polynomial models are determined by parameters  $(\beta_0, \beta_2, \dots, \beta_p)$
- ▶ Other methods more flexible

# PARAMETRIC VS NONPARAMETRIC MODELS

- ▶ Linear and polynomial models are determined by parameters  $(\beta_0, \beta_2, \dots, \beta_p)$
- ▶ Other methods more flexible
- ▶ Nearest neighbors (or k-NN)

# PARAMETRIC VS NONPARAMETRIC MODELS

- ▶ Linear and polynomial models are determined by parameters  $(\beta_0, \beta_2, \dots, \beta_p)$
- ▶ Other methods more flexible
- ▶ Nearest neighbors (or k-NN)

The goal is to estimate  $f(\cdot)$  not  $\beta_s$ !

# PARAMETRIC VS NONPARAMETRIC MODELS

- ▶ Linear and polynomial models are determined by parameters  $(\beta_0, \beta_2, \dots, \beta_p)$
- ▶ Other methods more flexible
- ▶ Nearest neighbors (or k-NN)

The goal is to estimate  $f(\cdot)$  not  $\beta_s$ !

Similar in spirit to moving average

# PARAMETRIC VS NONPARAMETRIC MODELS

- ▶ Linear and polynomial models are determined by parameters  $(\beta_0, \beta_2, \dots, \beta_p)$
- ▶ Other methods more flexible
- ▶ Nearest neighbors (or k-NN)

The goal is to estimate  $f(\cdot)$  not  $\beta_s$ !

Similar in spirit to moving average

$$\hat{f}(x_i) = \frac{1}{k} \sum_{j \in \{k\text{-nearest neighbours of } x_i\}} y_j$$

# PARAMETRIC VS NONPARAMETRIC MODELS

- ▶ Linear and polynomial models are determined by parameters  $(\beta_0, \beta_2, \dots, \beta_p)$
- ▶ Other methods more flexible
- ▶ Nearest neighbors (or k-NN)

The goal is to estimate  $f(\cdot)$  not  $\beta_s$ !

Similar in spirit to moving average

$$\hat{f}(x_i) = \frac{1}{k} \sum_{j \in \{k\text{-nearest neighbours of } x_i\}} y_j$$

$k$  is the number of neighbors of  $x_i$  taken into account in the estimation.

# PARAMETRIC VS NONPARAMETRIC MODELS

- ▶ Linear and polynomial models are determined by parameters  $(\beta_0, \beta_2, \dots, \beta_p)$
- ▶ Other methods more flexible
- ▶ Nearest neighbors (or k-NN)

The goal is to estimate  $f(\cdot)$  not  $\beta_s$ !

Similar in spirit to moving average

$$\hat{f}(x_i) = \frac{1}{k} \sum_{j \in \{k\text{-nearest neighbours of } x_i\}} y_j$$

$k$  is the number of neighbors of  $x_i$  taken into account in the estimation.

- ▶ The method follows a very general idea:



# PARAMETRIC VS NONPARAMETRIC MODELS

- ▶ Linear and polynomial models are determined by parameters  $(\beta_0, \beta_2, \dots, \beta_p)$
- ▶ Other methods more flexible
- ▶ Nearest neighbors (or k-NN)

The goal is to estimate  $f(\cdot)$  not  $\beta_s$ !

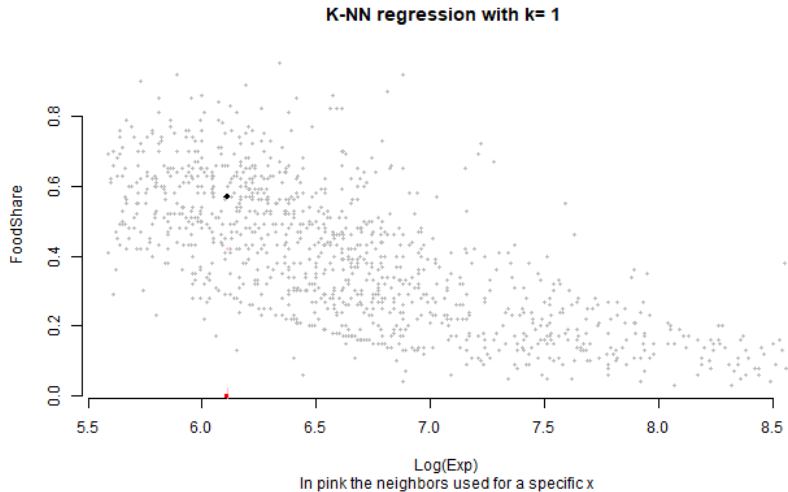
Similar in spirit to moving average

$$\hat{f}(x_i) = \frac{1}{k} \sum_{j \in \{k\text{-nearest neighbours of } x_i\}} y_j$$

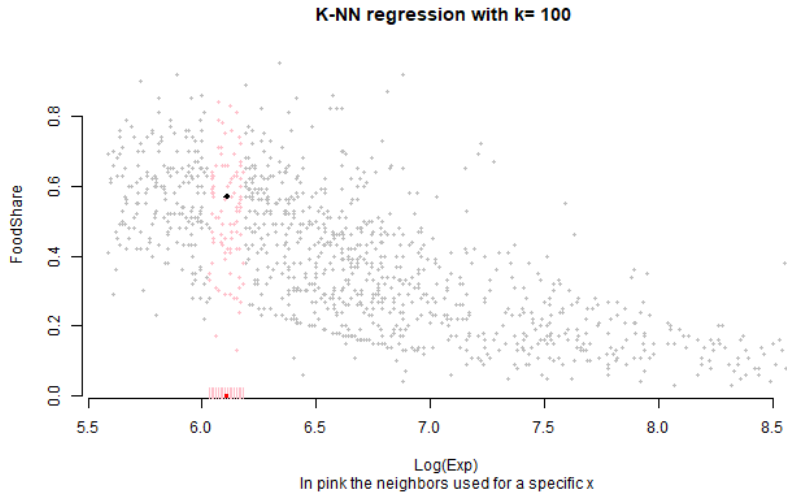
$k$  is the number of neighbors of  $x_i$  taken into account in the estimation.

- ▶ The method follows a very general idea:  
*"Observations close in the  $x$  dimension should be close in the  $y$  dimension"*

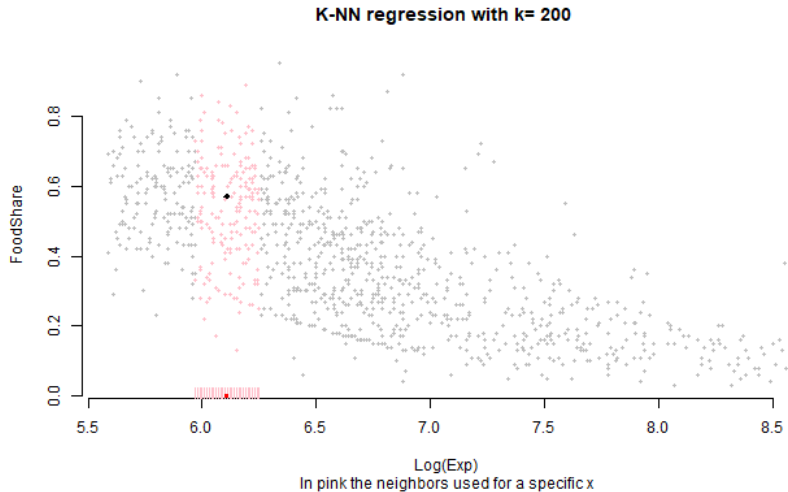
# K-NN IN PRACTICE



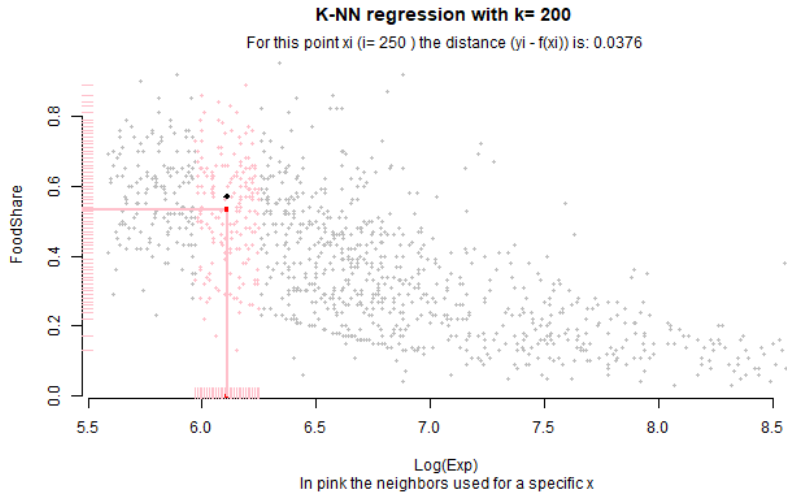
# K-NN IN PRACTICE



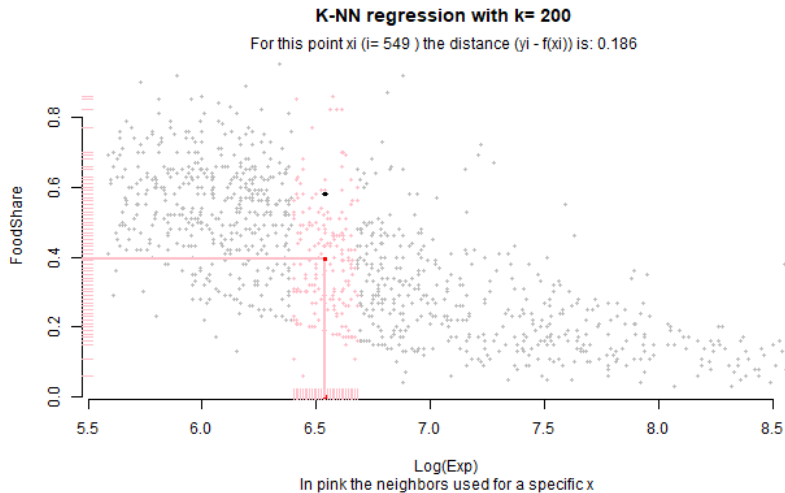
# K-NN IN PRACTICE



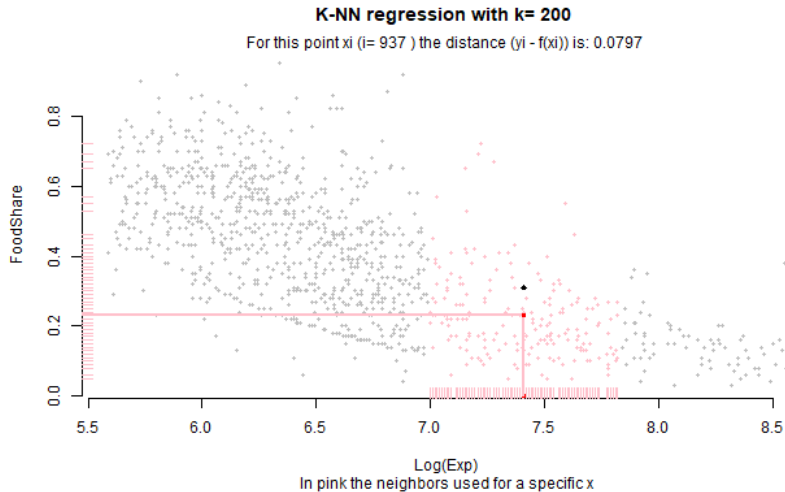
# K-NN IN PRACTICE



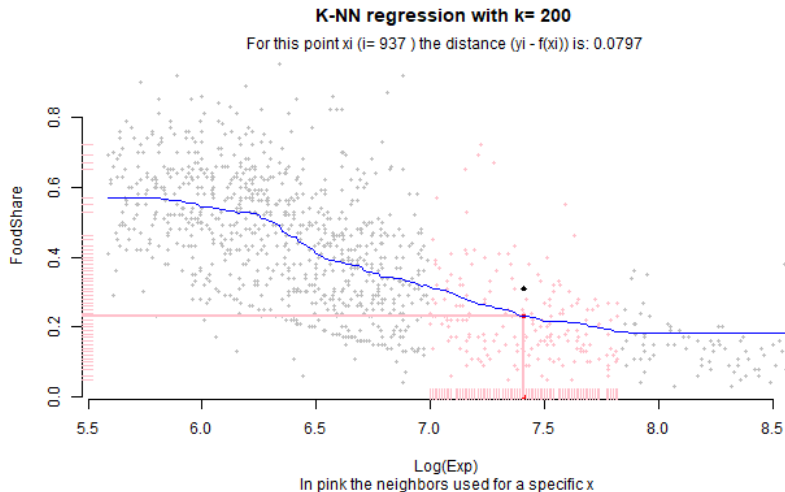
# K-NN IN PRACTICE



# K-NN IN PRACTICE

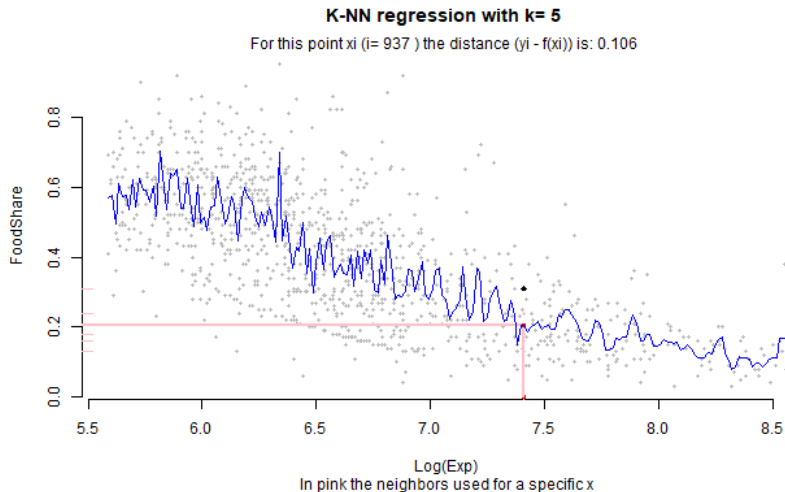


# K-NN IN PRACTICE

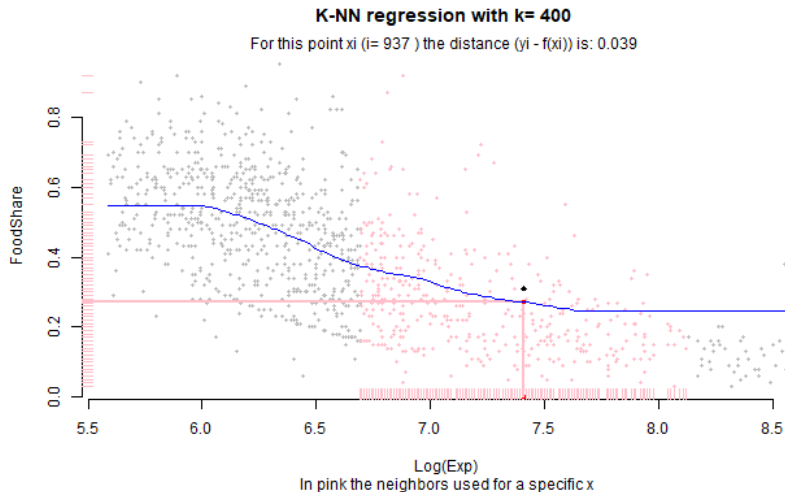




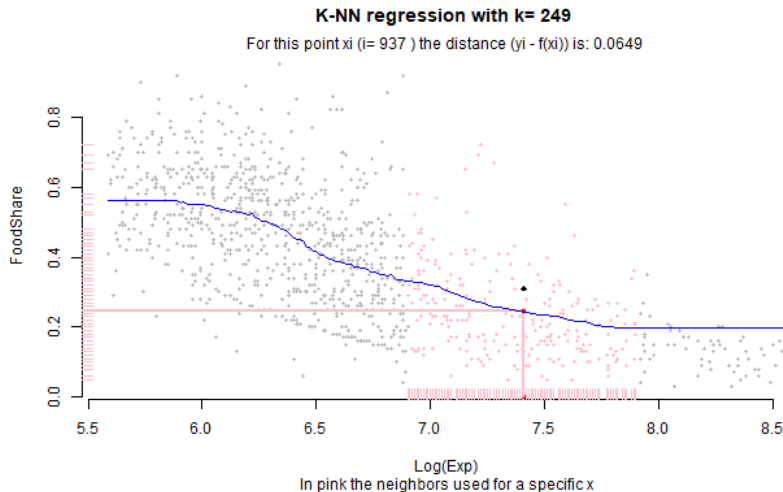
# K-NN IN PRACTICE: CHOOSING K



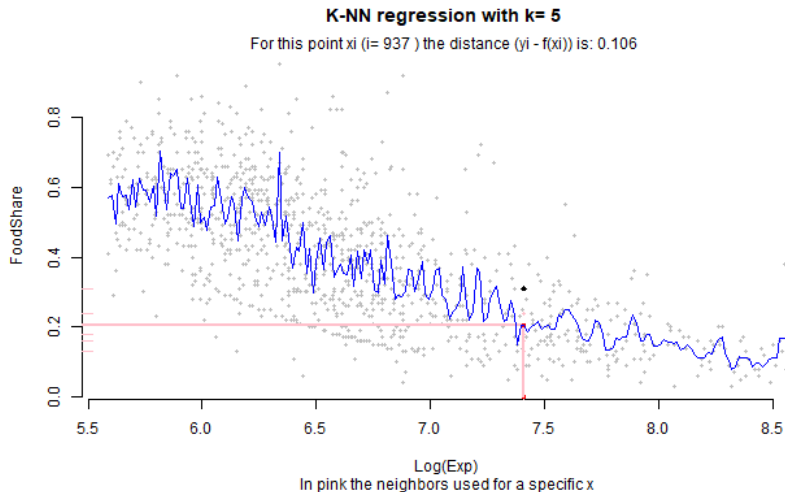
# K-NN IN PRACTICE: CHOOSING K



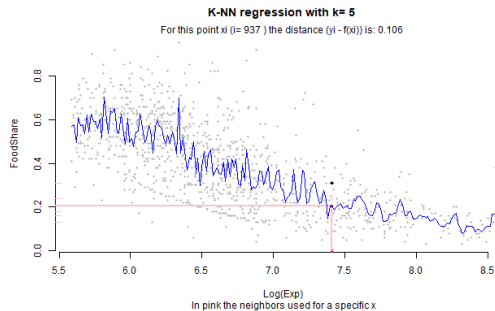
# K-NN IN PRACTICE: CHOOSING K



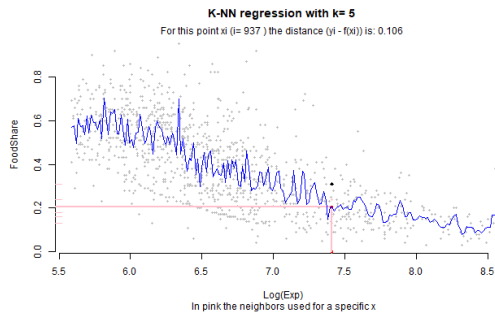
# K-NN IN PRACTICE: CHOOSING K



# K-NN IN PRACTICE: OVERFITTING

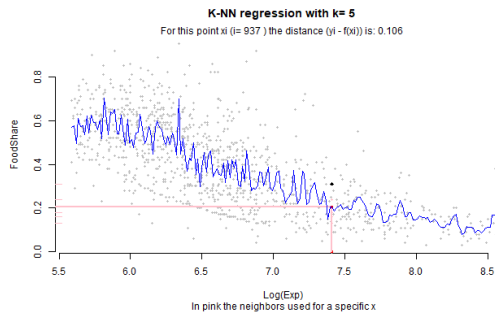


# K-NN IN PRACTICE: OVERFITTING



Overfitting has many consequences

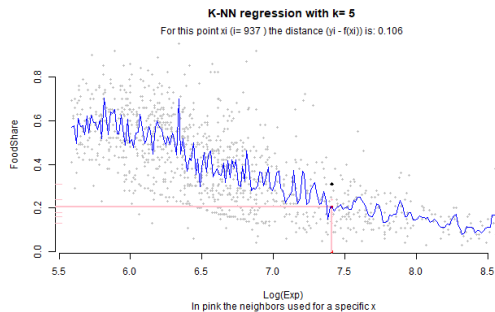
# K-NN IN PRACTICE: OVERFITTING



Overfitting has many consequences

- The estimated curve follows the data too closely

# K-NN IN PRACTICE: OVERFITTING

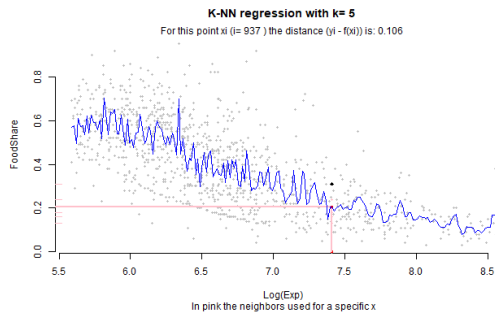


Overfitting has many consequences

- ▶ The estimated curve follows the data too closely
- ▶ The estimated curve follows the **errors** too closely



# K-NN IN PRACTICE: OVERFITTING



Overfitting has many consequences

- ▶ The estimated curve follows the data too closely
- ▶ The estimated curve follows the **errors** too closely
- ▶ The estimated function will not provide good estimates on **new observations**

## Machine Learning for Official Statistics and SDGs

# Statistical learning: *vs* Machine Learning



# WHAT LEARNING MEANS?

The goal is to estimate  $f(\cdot)$

- Many statistical learning methods are relevant and useful to estimate  $f(\cdot)$

## WHAT LEARNING MEANS?

The goal is to estimate  $f(\cdot)$

- ▶ Many statistical learning methods are relevant and useful to estimate  $f(\cdot)$
- ▶  $f(\cdot)$  can take continuous (regression) or discrete values (classification)

## WHAT LEARNING MEANS?

The goal is to estimate  $f(\cdot)$

- ▶ Many statistical learning methods are relevant and useful to estimate  $f(\cdot)$
- ▶  $f(\cdot)$  can take continuous (regression) or discrete values (classification)
- ▶ If predicting is the goal, one may focus on prediction accuracy

# WHAT LEARNING MEANS?

The goal is to estimate  $f(\cdot)$

- ▶ Many statistical learning methods are relevant and useful to estimate  $f(\cdot)$
- ▶  $f(\cdot)$  can take continuous (regression) or discrete values (classification)
- ▶ If predicting is the goal, one may focus on prediction accuracy  
     $\hookrightarrow$  this is the goal of **Machine Learning**

# WHAT LEARNING MEANS?

The goal is to estimate  $f(\cdot)$

- ▶ Many statistical learning methods are relevant and useful to estimate  $f(\cdot)$
- ▶  $f(\cdot)$  can take continuous (regression) or discrete values (classification)
- ▶ If predicting is the goal, one may focus on prediction accuracy  
     $\hookrightarrow$  this is the goal of **Machine Learning**
- ▶ If the goal is to formalize a model, one may focus on testing statistical properties

## WHAT LEARNING MEANS?

The goal is to estimate  $f(\cdot)$

- ▶ Many statistical learning methods are relevant and useful to estimate  $f(\cdot)$
- ▶  $f(\cdot)$  can take continuous (regression) or discrete values (classification)
- ▶ If predicting is the goal, one may focus on prediction accuracy  
     $\hookrightarrow$  this is the goal of **Machine Learning**
- ▶ If the goal is to formalize a model, one may focus on testing statistical properties  
     $\hookrightarrow$  this is the goal of **Statistical Learning**



## WHAT LEARNING MEANS?

The goal is to estimate  $f(\cdot)$

- ▶ Many statistical learning methods are relevant and useful to estimate  $f(\cdot)$
- ▶  $f(\cdot)$  can take continuous (regression) or discrete values (classification)
- ▶ If predicting is the goal, one may focus on prediction accuracy  
     $\hookrightarrow$  this is the goal of **Machine Learning**
- ▶ If the goal is to formalize a model, one may focus on testing statistical properties  
     $\hookrightarrow$  this is the goal of **Statistical Learning**
- ▶ In practice we'll use both tools to "*understand the data*"

## WHAT LEARNING MEANS?

The classical approach

- So far, we have estimated  $f(\cdot)$  on the whole data set

# WHAT LEARNING MEANS?

## The classical approach

- So far, we have estimated  $f(\cdot)$  on the whole data set



# WHAT LEARNING MEANS?

## The classical approach

- So far, we have estimated  $f(\cdot)$  on the whole data set



- We have estimated  $f(\cdot)$  by  $\hat{f}(\cdot)$  and minimized some cost function such as the 
$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

# WHAT LEARNING MEANS?

## The classical approach

- So far, we have estimated  $f(\cdot)$  on the whole data set



- We have estimated  $f(\cdot)$  by  $\hat{f}(\cdot)$  and minimized some cost function such as the  $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$
- The data serve both for estimating  $\hat{f}(\cdot)$  and computing the prediction error

# WHAT LEARNING MEANS?

A different approach: *resampling*

- Our goal is evaluate the prediction accuracy of  $\hat{f}(\cdot)$  on a new, **unseen**, data set

# WHAT LEARNING MEANS?

A different approach: *resampling*

- ▶ Our goal is evaluate the prediction accuracy of  $\hat{f}(\cdot)$  on a new, **unseen**, data set
- ▶ Since we may not have **unseen** data, we will construct one

# WHAT LEARNING MEANS?

A different approach: *resampling*

- ▶ Our goal is evaluate the prediction accuracy of  $\hat{f}(\cdot)$  on a new, **unseen**, data set
- ▶ Since we may not have **unseen** data, we will construct one





# WHAT LEARNING MEANS?

A different approach: *resampling*

- ▶ Our goal is evaluate the prediction accuracy of  $\hat{f}(\cdot)$  on a new, **unseen**, data set
- ▶ Since we may not have **unseen** data, we will construct one



- ▶ And the true value of  $y_i$  will be available to compute MSE of the prediction

## WHY DIFFERENT SETS?

Predicting using predictions capability

- ▶ When estimating  $f(\cdot)$  on the whole data set, over-fitting may occur

## WHY DIFFERENT SETS?

### Predicting using predictions capability

- ▶ When estimating  $f(\cdot)$  on the whole data set, over-fitting may occur
- ▶ The validation set provides a good way to evaluate the prediction capabilities of a model and the prediction error on a new data set

## WHY DIFFERENT SETS?

### Predicting using predictions capability

- ▶ When estimating  $f(\cdot)$  on the whole data set, over-fitting may occur
- ▶ The validation set provides a good way to evaluate the prediction capabilities of a model and the prediction error on a new data set



## WHY DIFFERENT SETS?

### Predicting using predictions capability

- ▶ When estimating  $f(\cdot)$  on the whole data set, over-fitting may occur
- ▶ The validation set provides a good way to evaluate the prediction capabilities of a model and the prediction error on a new data set



- ▶ Prediction accuracy (using  $\hat{f}(\cdot)$ ) is then evaluated on the validation set **only**

# CONSTRUCTING TRAINING & VALIDATION SETS

In practice, the validation set is not a block



# CONSTRUCTING TRAINING & VALIDATION SETS

In practice, the validation set is not a block



- The validation set is constructed from a randomly drawn observations.

# CONSTRUCTING TRAINING & VALIDATION SETS

In practice, the validation set is not a block



- The validation set is constructed from a randomly drawn observations.

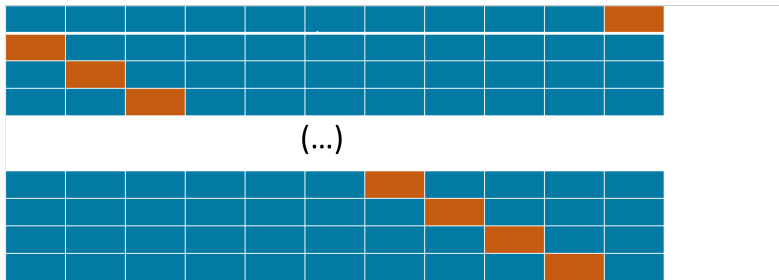
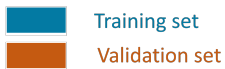




## Many DIFFERENT SETS!

Using resampling methods to estimate the error on the prediction

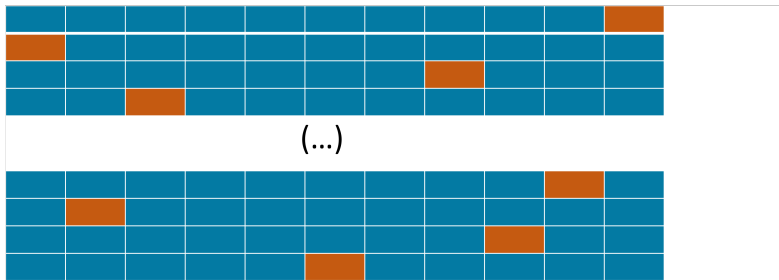
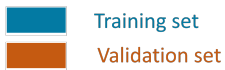
- Cross validation is used to select  $m$ -(training-validation) sets from the original data set (here again randomly)



## Many DIFFERENT SETS!

Using resampling methods to estimate the error on the prediction

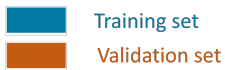
- Cross validation is used to select  $m$ -(training-validation) sets from the original data set (here again randomly)



## Many DIFFERENT SETS!

Using resampling methods to estimate the error on the prediction

- Cross validation is used to select  $m$ -(training-validation) sets from the original data set (here again randomly)



## *Many* DIFFERENT SETS!

*m-fold* Cross-Validation estimates the average prediction error on  $m$  different(training-validation) sets

## Many DIFFERENT SETS!

*m-fold* Cross-Validation estimates the average prediction error on  $m$  different(training-validation) sets

- For each training-validation) set  $i$ , one can compute the  $MSE_i$  since the true  $y_i$  is known on the validation set!

## Many DIFFERENT SETS!

*m-fold* Cross-Validation estimates the average prediction error on  $m$  different(training-validation) sets

- ▶ For each training-validation) set  $i$ , one can compute the  $MSE_i$  since the true  $y_i$  is known on the validation set!
- ▶ Cross Validation error is then:

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

## Many DIFFERENT SETS!

*m*-fold Cross-Validation estimates the average prediction error on *m* different(training-validation) sets

- ▶ For each training-validation) set *i*, one can compute the  $MSE_i$  since the true  $y_i$  is known on the validation set!
- ▶ Cross Validation error is then:

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

- ▶  $CV_{(m)}$  is a good estimate of the prediction error of the model

## Many DIFFERENT SETS!

*m*-fold Cross-Validation estimates the average prediction error on *m* different(training-validation) sets

- ▶ For each training-validation) set *i*, one can compute the  $MSE_i$  since the true  $y_i$  is known on the validation set!
- ▶ Cross Validation error is then:

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

- ▶  $CV_{(m)}$  is a good estimate of the prediction error of the model
- ▶  $CV_{(m)}$  can serve to select and compare models (example: select *k* in k-NN regression)



## Many DIFFERENT SETS!

*m*-fold Cross-Validation estimates the average prediction error on *m* different(training-validation) sets

- ▶ For each training-validation) set *i*, one can compute the  $MSE_i$  since the true  $y_i$  is known on the validation set!
- ▶ Cross Validation error is then:

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

- ▶  $CV_{(m)}$  is a good estimate of the prediction error of the model
- ▶  $CV_{(m)}$  can serve to select and compare models (example: select *k* in k-NN regression)
- ▶ In practice  $m \in 5, \dots, 10$  shows good performances

# TASKS FOR MACHINE LEARNING

Machine Learning involves several tasks, some are time consuming

# TASKS FOR MACHINE LEARNING

Machine Learning involves several tasks, some are time consuming

- ▶ Data collection (not treated here)

# TASKS FOR MACHINE LEARNING

Machine Learning involves several tasks, some are time consuming

- ▶ Data collection (not treated here)
- ▶ Data organization (not treated here)

# TASKS FOR MACHINE LEARNING

Machine Learning involves several tasks, some are time consuming

- ▶ Data collection (not treated here)
- ▶ Data organization (not treated here)
- ▶ Data cleaning (not treated here)

# TASKS FOR MACHINE LEARNING

Machine Learning involves several tasks, some are time consuming

- ▶ Data collection (not treated here)
- ▶ Data organization (not treated here)
- ▶ Data cleaning (not treated here)
- ▶ Data visualization

# TASKS FOR MACHINE LEARNING

Machine Learning involves several tasks, some are time consuming

- ▶ Data collection (not treated here)
- ▶ Data organization (not treated here)
- ▶ Data cleaning (not treated here)
- ▶ Data visualization
- ▶ Data analysis ← this is the core of this course

# WRAP-UP



## WRAP-UP

- To understand the data, we can go beyond linearity using polynomial or nonparametric model ( $k$ -NN)

## WRAP-UP

- ▶ To understand the data, we can go beyond linearity using polynomial or nonparametric model ( $k$ -NN)
- ▶ More complex models allow for more accuracy, but introduce variance in the estimation

## WRAP-UP

- ▶ To understand the data, we can go beyond linearity using polynomial or nonparametric model ( $k$ -NN)
- ▶ More complex models allow for more accuracy, but introduce variance in the estimation
- ▶ There is an unavoidable **bias-variance** trade-off

## WRAP-UP

- ▶ To understand the data, we can go beyond linearity using polynomial or nonparametric model ( $k$ -NN)
- ▶ More complex models allow for more accuracy, but introduce variance in the estimation
- ▶ There is an unavoidable **bias-variance** trade-off
- ▶ Theory helps understanding but not choosing the right model

## WRAP-UP

- ▶ To understand the data, we can go beyond linearity using polynomial or nonparametric model ( $k$ -NN)
- ▶ More complex models allow for more accuracy, but introduce variance in the estimation
- ▶ There is an unavoidable **bias-variance** trade-off
- ▶ Theory helps understanding but not choosing the right model
- ▶ The train + validation sets approach is central in machine learning