

Machine Learning for Official Statistics and SDGs

First Live Lecture (Webinar):

Starts in **15** minutes

Christophe Bontemps, UN SIAP



Machine Learning for Official Statistics and SDGs

First Live Lecture (Webinar):

Starts in **10** minutes

Christophe Bontemps, UN SIAP



Machine Learning for Official Statistics and SDGs

First Live Lecture (Webinar):

Starts in 5 minutes

Christophe Bontemps, UN SIAP



Machine Learning for Official Statistics and SDGs

Statistical learning: *vs* Machine Learning



[- REMINDER -]

- ▶ Mute yourself **always!**

[- REMINDER -]

- ▶ Mute yourself **always!**
- ▶ The lecture is recorded

[- REMINDER -]

- ▶ Mute yourself **always**!
- ▶ The lecture is recorded
- ▶ Ask questions in the chat

[- AGENDA -]

► Introduction

[- AGENDA -]

- ▶ Introduction
- ▶ Statistical learning *vs* Machine Learning

[- AGENDA -]

- ▶ Introduction
- ▶ Statistical learning *vs* Machine Learning
- ▶ Q&A

[- AGENDA -]

- ▶ Introduction
- ▶ Statistical learning *vs* Machine Learning
- ▶ Q&A
- ▶ Next week

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)

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?

“ 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?

“ 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** an *outcome* y and *explanatory* variables x s

WHAT IS STATISTICAL LEARNING?

Two main learning problems:

- ▶ We observe **both** an *outcome* y and *explanatory* variables x s
- ↪ **Supervised** learning

WHAT IS STATISTICAL LEARNING?

Two main learning problems:

- ▶ We observe **both** an *outcome* y and *explanatory* variables x s

↪ **Supervised** learning

Most of the examples and applications are supervised learning

WHAT IS STATISTICAL LEARNING?

Two main learning problems:

- ▶ We observe **both** an *outcome* y and *explanatory* variables x s
- ↪ **Supervised** learning
 - Most of the examples and applications are supervised learning*
- ▶ We **do not** observe an outcome y but **only** several x s

WHAT IS STATISTICAL LEARNING?

Two main learning problems:

- ▶ We observe **both** an *outcome* y and *explanatory* variables x s
- ↪ **Supervised** learning
 - Most of the examples and applications are supervised learning*
- ▶ We **do not** observe an outcome y but **only** several x s
- ↪ **Unsupervised** learning (or *cluster analysis*)

WHAT IS STATISTICAL LEARNING?

Two main learning problems:

- ▶ We observe **both** an *outcome* y and *explanatory* variables x s

↪ **Supervised** learning

Most of the examples and applications are supervised learning

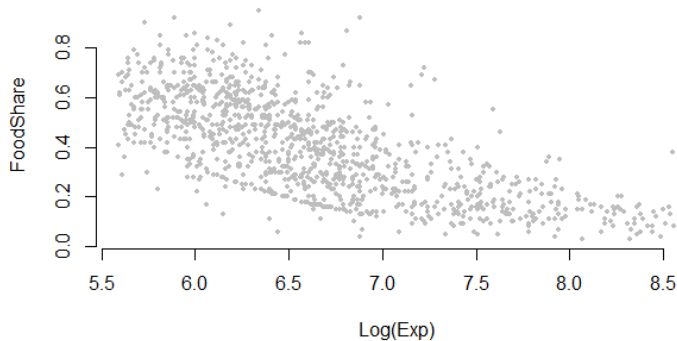
- ▶ We **do not** observe an outcome y but **only** several x s

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

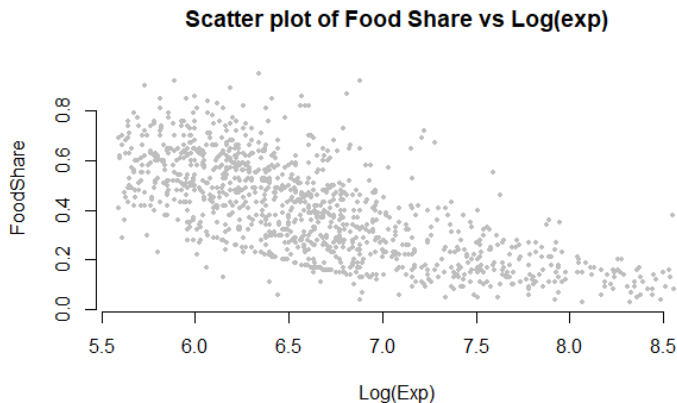
More complex models we'll see at the end of the course

STATISTICAL LEARNING ON AN EXAMPLE

Scatter plot of Food Share vs Log(exp)

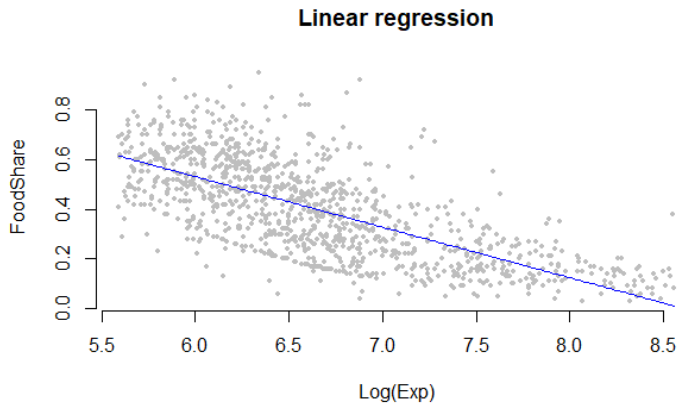


STATISTICAL LEARNING ON AN EXAMPLE

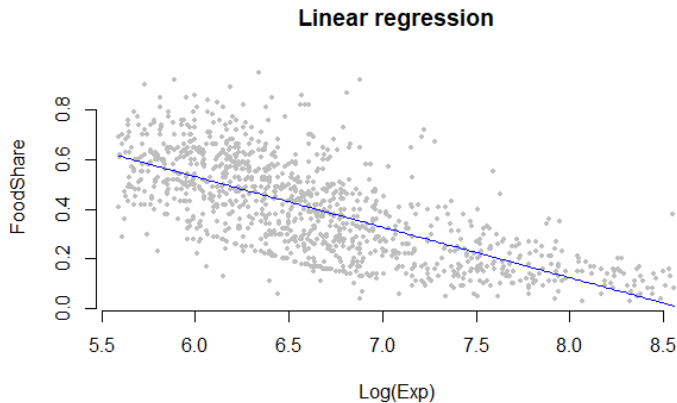


We may be interested in the **relationship** between the two variables

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 β_0 and β_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):

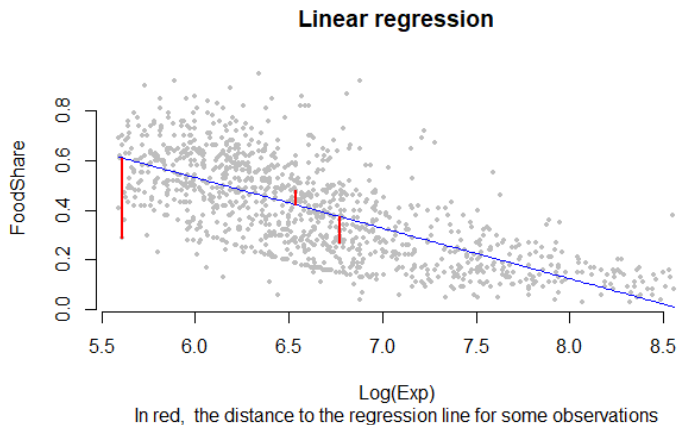
$$MSE(\beta_0, \beta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

- The regression line, defined by β_0 and β_1 , is simply the solution of:

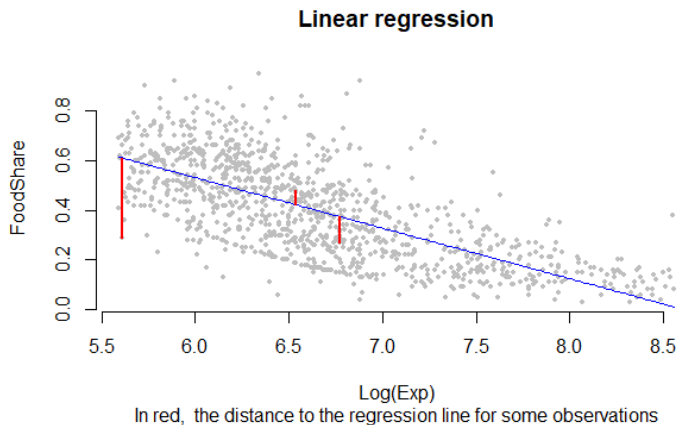
$$\text{Min}_{(\beta_0, \beta_1)} MSE(\beta_0, \beta_1)$$

The MSE it is the *cost function* minimized to determine $(\hat{\beta}_0, \hat{\beta}_1)$

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

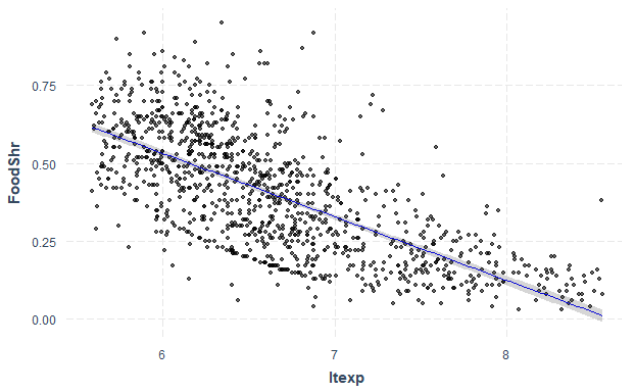
The quality of the adjustment may be measured by the $R^2 = 0.478$

BEYOND LINEARITY

BEYOND LINEARITY

- ▶ A linear model may be unadapted or too simple

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

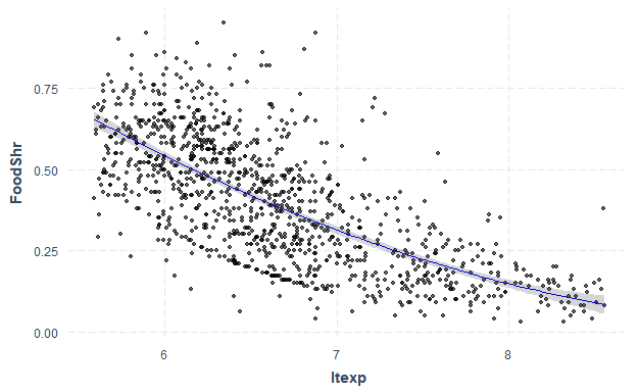


The fit (measured by R^2) is: $R^2 = 0.478$

BEYOND LINEARITY

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

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

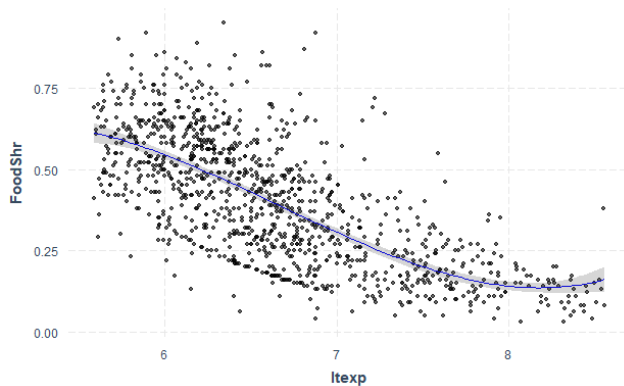


Do we have a better fit? $R^2 = 0.484$

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

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

IN PRACTICE

- ▶ Linear and polynomial models are determined by parameters $(\beta_0, \beta_1, \dots, \beta_p)$

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

IN PRACTICE

- ▶ Linear and polynomial models are determined by parameters $(\beta_0, \beta_1, \dots, \beta_p)$

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

- ▶ How to choose the degree p ?

IN PRACTICE

- ▶ Linear and polynomial models are determined by parameters $(\beta_0, \beta_1, \dots, \beta_p)$

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

- ▶ How to choose the degree p ?
- ▶ Collinearity of x^p and x^q for $p \neq q$?

IN PRACTICE

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

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

- ▶ How to choose the degree p ?
- ▶ Collinearity of x^p and x^q for $p \neq q$?
- ...

IN PRACTICE

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

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

- ▶ How to choose the degree p ?
- ▶ Collinearity of x^p and x^q for $p \neq q$?
- ...

↪ *How does that relates to the learning exercise?*

NEAREST NEIGHBORS (K-NN)

- ▶ Other methods more flexible

NEAREST NEIGHBORS (K-NN)

- ▶ Other methods more flexible
- ▶ Nearest neighbors (or k-NN)

NEAREST NEIGHBORS (K-NN)

- ▶ Other methods more flexible
- ▶ Nearest neighbors (or k-NN)
 - ↪ The goal is to estimate $f(\cdot)$ not β_s !

NEAREST NEIGHBORS (K-NN)

- ▶ Other methods more flexible
- ▶ Nearest neighbors (or k-NN)
 - ↪ The goal is to estimate $f(\cdot)$ not β_s !
Similar in spirit to "moving average" estimator

NEAREST NEIGHBORS (K-NN)

► Other methods more flexible

► Nearest neighbors (or k-NN)

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

Similar in spirit to "moving average" estimator

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

NEAREST NEIGHBORS (K-NN)

► Other methods more flexible

► Nearest neighbors (or k-NN)

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

Similar in spirit to "moving average" estimator

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

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

NEAREST NEIGHBORS (K-NN)

► Other methods more flexible

► Nearest neighbors (or k-NN)

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

Similar in spirit to "moving average" estimator

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

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

NEAREST NEIGHBORS (K-NN)

- ▶ Other methods more flexible
- ▶ Nearest neighbors (or k-NN)

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

Similar in spirit to "moving average" estimator

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

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

- ⊙ The method follows a very general idea:

NEAREST NEIGHBORS (K-NN)

► Other methods more flexible

► Nearest neighbors (or k-NN)

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

Similar in spirit to "moving average" estimator

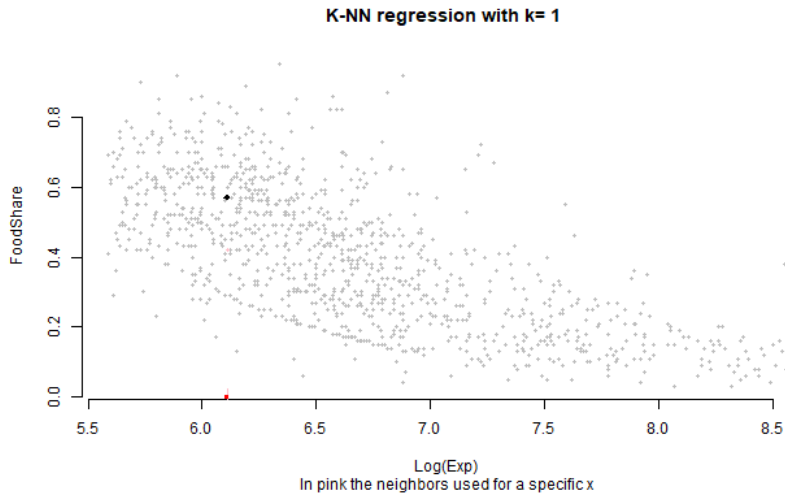
$$\hat{f}(x_i) = \frac{1}{k} \sum_{j \in \{k\text{-nearest neighbours of } x_i\}} y_j$$

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

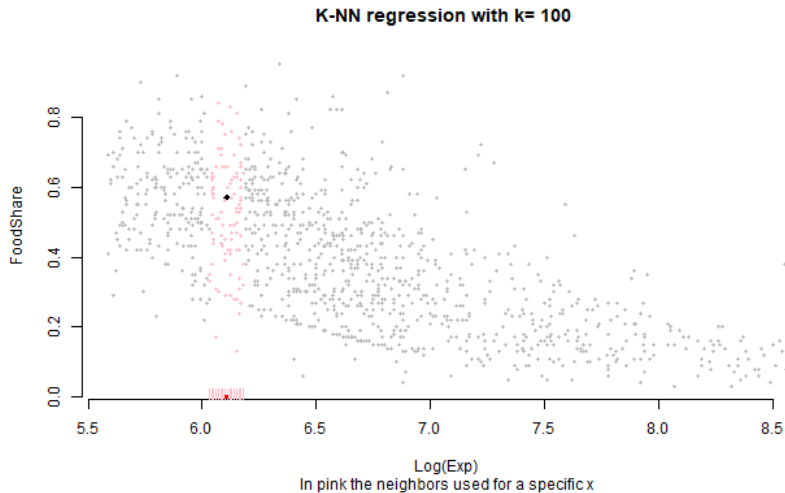
⊙ The method follows a very general idea:

"Observations close in the x dimension should be close in the y dimension"

K-NN IN PRACTICE



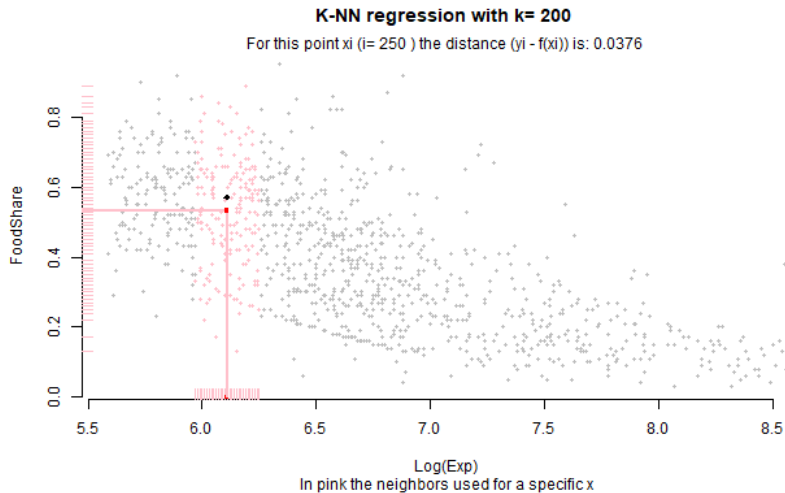
K-NN IN PRACTICE



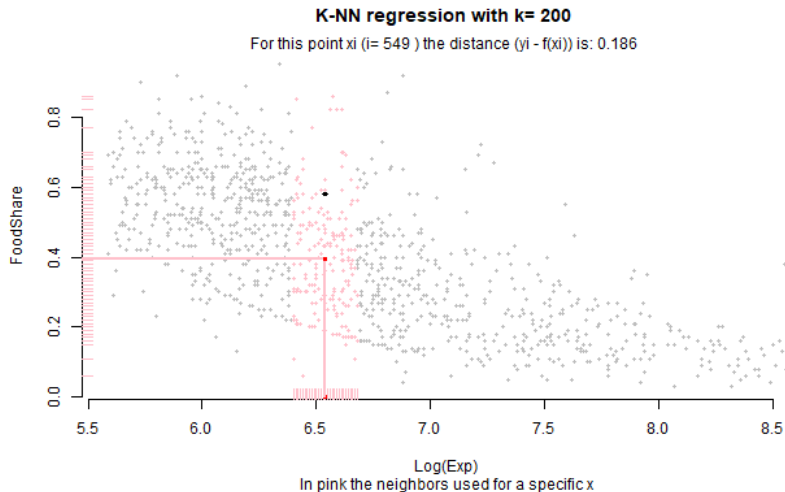
K-NN IN PRACTICE



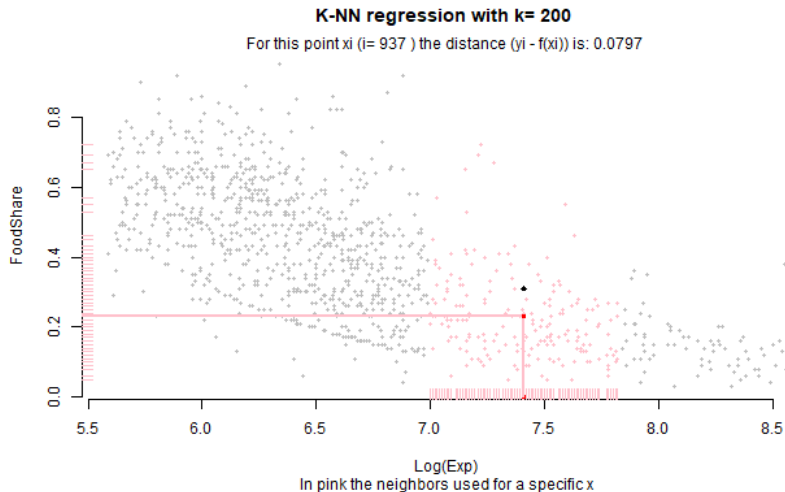
K-NN IN PRACTICE



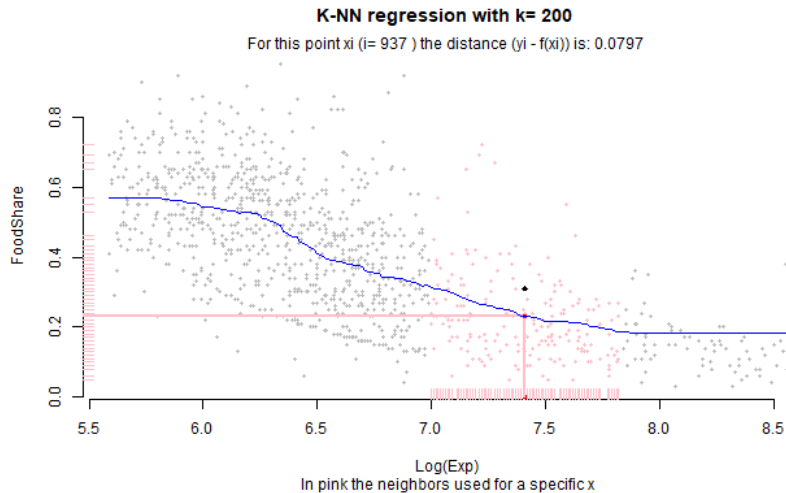
K-NN IN PRACTICE



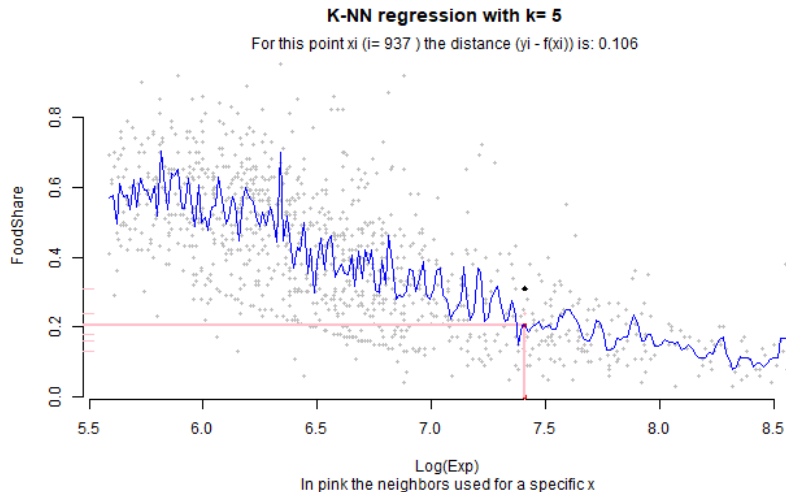
K-NN IN PRACTICE



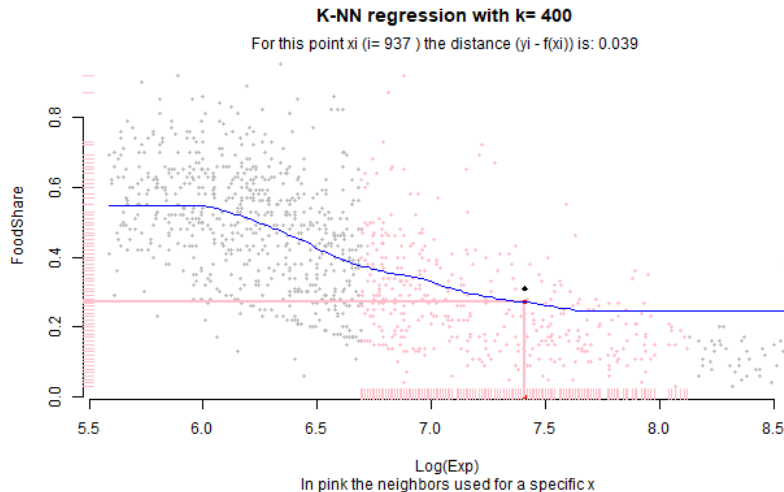
K-NN IN PRACTICE



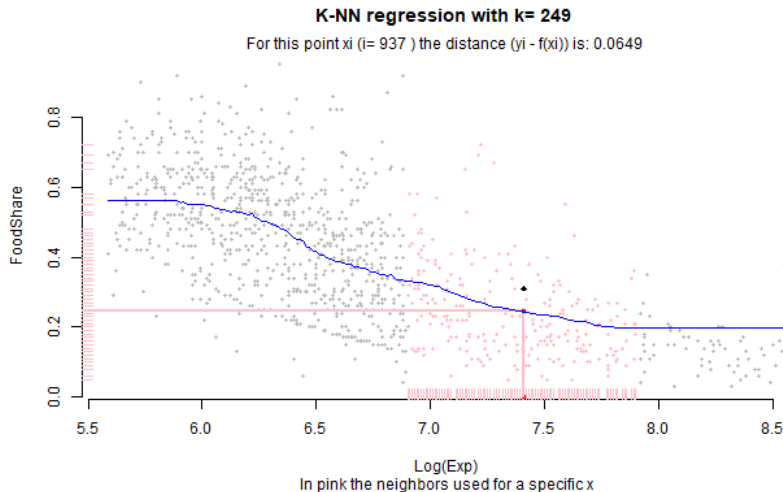
K-NN IN PRACTICE: CHOOSING K



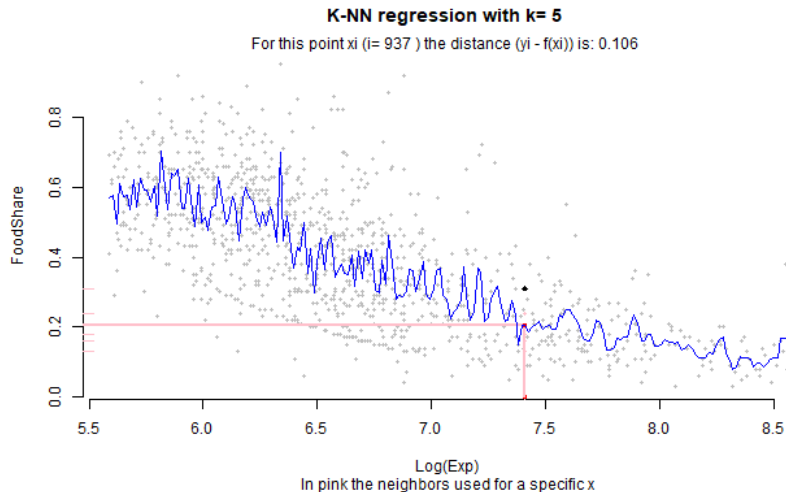
K-NN IN PRACTICE: CHOOSING K



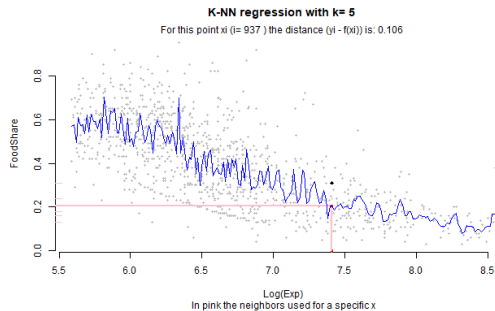
K-NN IN PRACTICE: CHOOSING K



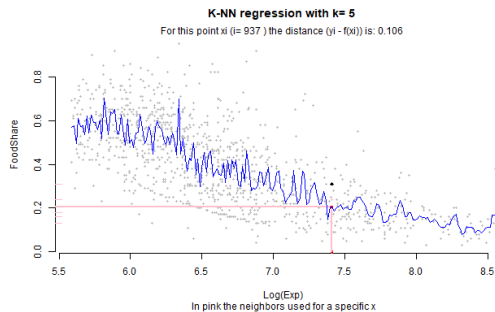
K-NN IN PRACTICE: CHOOSING K



K-NN IN PRACTICE: OVERFITTING

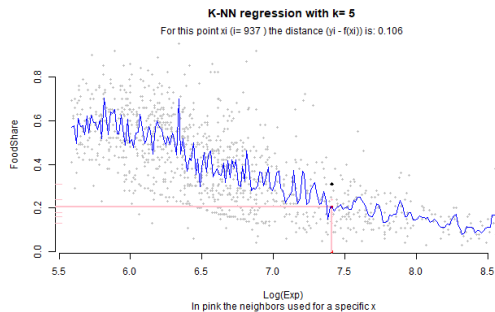


K-NN IN PRACTICE: OVERFITTING



Overfitting has many consequences

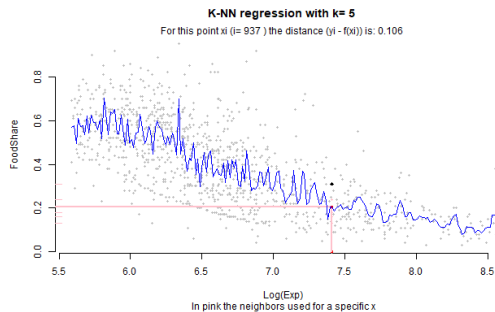
K-NN IN PRACTICE: OVERFITTING



Overfitting has many consequences

- The estimated curve follows the **data set** too closely

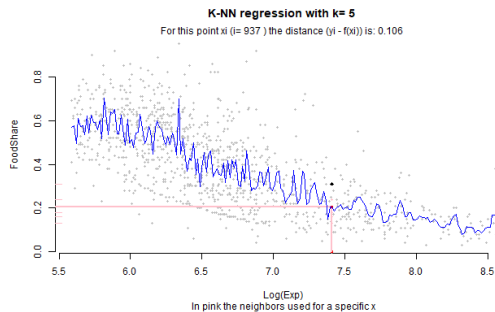
K-NN IN PRACTICE: OVERFITTING



Overfitting has many consequences

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

STATISTICAL *vs* MACHINE LEARNING

What is the goal?

STATISTICAL *vs* MACHINE LEARNING

What is the goal?

- ▶ If the goal is to formalize a model, one may focus on testing statistical properties, significance, relationships, ...

STATISTICAL *vs* MACHINE LEARNING

What is the goal?

- ▶ If the goal is to formalize a model, one may focus on testing statistical properties, significance, relationships, ...
 ↪ this is the purpose of **Statistical Learning**

STATISTICAL *vs* MACHINE LEARNING

What is the goal?

- ▶ If the goal is to formalize a model, one may focus on testing statistical properties, significance, relationships, ...
↪ this is the purpose of **Statistical Learning**
- ▶ If the goal is to predict, one may focus on prediction accuracy

STATISTICAL *vs* MACHINE LEARNING

What is the goal?

- ▶ If the goal is to formalize a model, one may focus on testing statistical properties, significance, relationships, ...
 ↪ this is the purpose of **Statistical Learning**
- ▶ If the goal is to predict, one may focus on prediction accuracy
 ↪ this is the purpose of **Machine Learning**

STATISTICAL *vs* MACHINE LEARNING

What is the goal?

- ▶ If the goal is to formalize a model, one may focus on testing statistical properties, significance, relationships, ...
↪ this is the purpose of **Statistical Learning**
- ▶ If the goal is to predict, one may focus on prediction accuracy
↪ this is the purpose of **Machine Learning**
- ▶ Many statistical learning methods are relevant and useful to estimate $f(\cdot)$

STATISTICAL *vs* MACHINE LEARNING

What is the goal?

- ▶ If the goal is to formalize a model, one may focus on testing statistical properties, significance, relationships, ...
↪ this is the purpose of **Statistical Learning**
- ▶ If the goal is to predict, one may focus on prediction accuracy
↪ this is the purpose of **Machine Learning**
- ▶ Many statistical learning methods are relevant and useful to estimate $f(\cdot)$
- ▶ 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

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
- ▶ The data serve **both** for estimating $f(\cdot)$ **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



- ▶ Compare y_i with the prediction based on the validation set x_s

WHY DIFFERENT SETS?

Estimating parameters using predictions accuracy

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

WHY DIFFERENT SETS?

Estimating parameters using predictions accuracy

- ▶ 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?

Estimating parameters using predictions accuracy

- ▶ 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?

Estimating parameters using predictions accuracy

- ▶ 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.



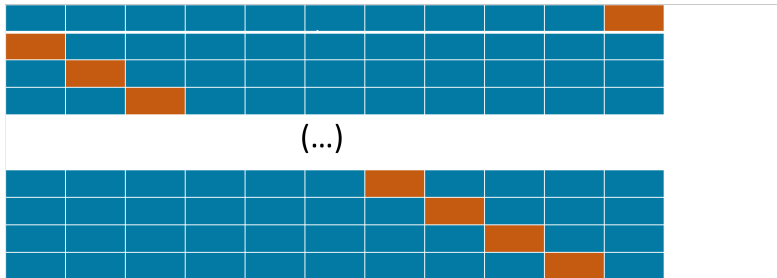
Many DIFFERENT SETS!

Using resampling methods to estimate the error on the prediction

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

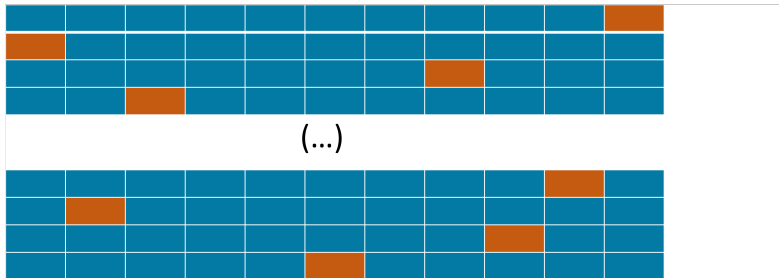
- ▶ Cross validation is used to select m -(training-validation) sets from the original data set (here again randomly)



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 *j*, one can compute the MSE_j since the true *y*s are 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 *j*, one can compute the MSE_j since the true *ys* are known on the validation set!
- ▶ Cross Validation error is then:

$$CV_{(m)} = \frac{1}{m} \sum_{j=1}^m MSE_j$$

Many DIFFERENT SETS!

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

- ▶ For each (*training – validation*) set *j*, one can compute the MSE_j since the true *ys* are known on the validation set!
- ▶ Cross Validation error is then:

$$CV_{(m)} = \frac{1}{m} \sum_{j=1}^m MSE_j$$

- ▶ $CV_{(m)}$ 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 *j*, one can compute the MSE_j since the true *y*s are known on the validation set!
- ▶ Cross Validation error is then:

$$CV_{(m)} = \frac{1}{m} \sum_{j=1}^m MSE_j$$

- ▶ $CV_{(m)}$ is a good estimate of the prediction error of the model
- ▶ $CV_{(m)}$ can serve to select and compare models

Many DIFFERENT SETS!

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

- ▶ For each (*training* – *validation*) set *j*, one can compute the MSE_j since the true *ys* are known on the validation set!
- ▶ Cross Validation error is then:

$$CV_{(m)} = \frac{1}{m} \sum_{j=1}^m MSE_j$$

- ▶ $CV_{(m)}$ 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

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 use linear, polynomial, nonparametric models (k -NN) or other complex methods (including those for classification)

WRAP-UP

- ▶ To *understand* the data, we use linear, polynomial, nonparametric models (k -NN) or other complex methods (including those for classification)
- ▶ More complex models ↗ accuracy, but → variability (variance) in the estimation

WRAP-UP

- ▶ To *understand* the data, we use linear, polynomial, nonparametric models (k -NN) or other complex methods (including those for classification)
- ▶ More complex models ↗ accuracy, but → variability (variance) in the estimation
- ▶ There is an unavoidable **bias-variance** trade-off

WRAP-UP

- ▶ To *understand* the data, we use linear, polynomial, nonparametric models (k -NN) or other complex methods (including those for classification)
- ▶ More complex models ↗ accuracy, but → variability (variance) in the estimation
- ▶ There is an unavoidable **bias-variance** trade-off
- ▶ Theory helps understanding but not in choosing the right model

WRAP-UP

- ▶ To *understand* the data, we use linear, polynomial, nonparametric models (k -NN) or other complex methods (including those for classification)
- ▶ More complex models \nearrow accuracy, but \rightarrow variability (variance) in the estimation
- ▶ There is an unavoidable **bias-variance** trade-off
- ▶ Theory helps understanding but not in choosing the right model
- ▶ The (*train* + *validation*) sets approach is central in machine learning

WRAP-UP

- ▶ To *understand* the data, we use linear, polynomial, nonparametric models (k -NN) or other complex methods (including those for classification)
 - ▶ More complex models ↗ accuracy, but → variability (variance) in the estimation
 - ▶ There is an unavoidable **bias-variance** trade-off
 - ▶ Theory helps understanding but not in choosing the right model
 - ▶ The (*train* + *validation*) sets approach is central in machine learning
- ↪ In a machine learning framework, the efficiency of the prediction will guide the choices, not the statistical properties!

[Q&A]

Write your questions in the chat

[NEXT WEEK]

[NEXT WEEK]

- ▶ Module 2: "Classification" (Examples of classifiers, Measures of fit, *Logit* as a classifier)

[NEXT WEEK]

- ▶ Module 2: "Classification" (Examples of classifiers, Measures of fit, *Logit* as a classifier)
- ▶ Webinar on "Classification" Thursday, same time

[NEXT WEEK]

- ▶ Module 2: "Classification" (Examples of classifiers, Measures of fit, *Logit* as a classifier)
- ▶ Webinar on "Classification" Thursday, same time
- ▶ Complete the activities before the webinar!

[NEXT WEEK]

- ▶ Module 2: "Classification" (Examples of classifiers, Measures of fit, *Logit* as a classifier)
- ▶ Webinar on "Classification" Thursday, same time
- ▶ Complete the activities before the webinar!

Have a nice week!