

## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

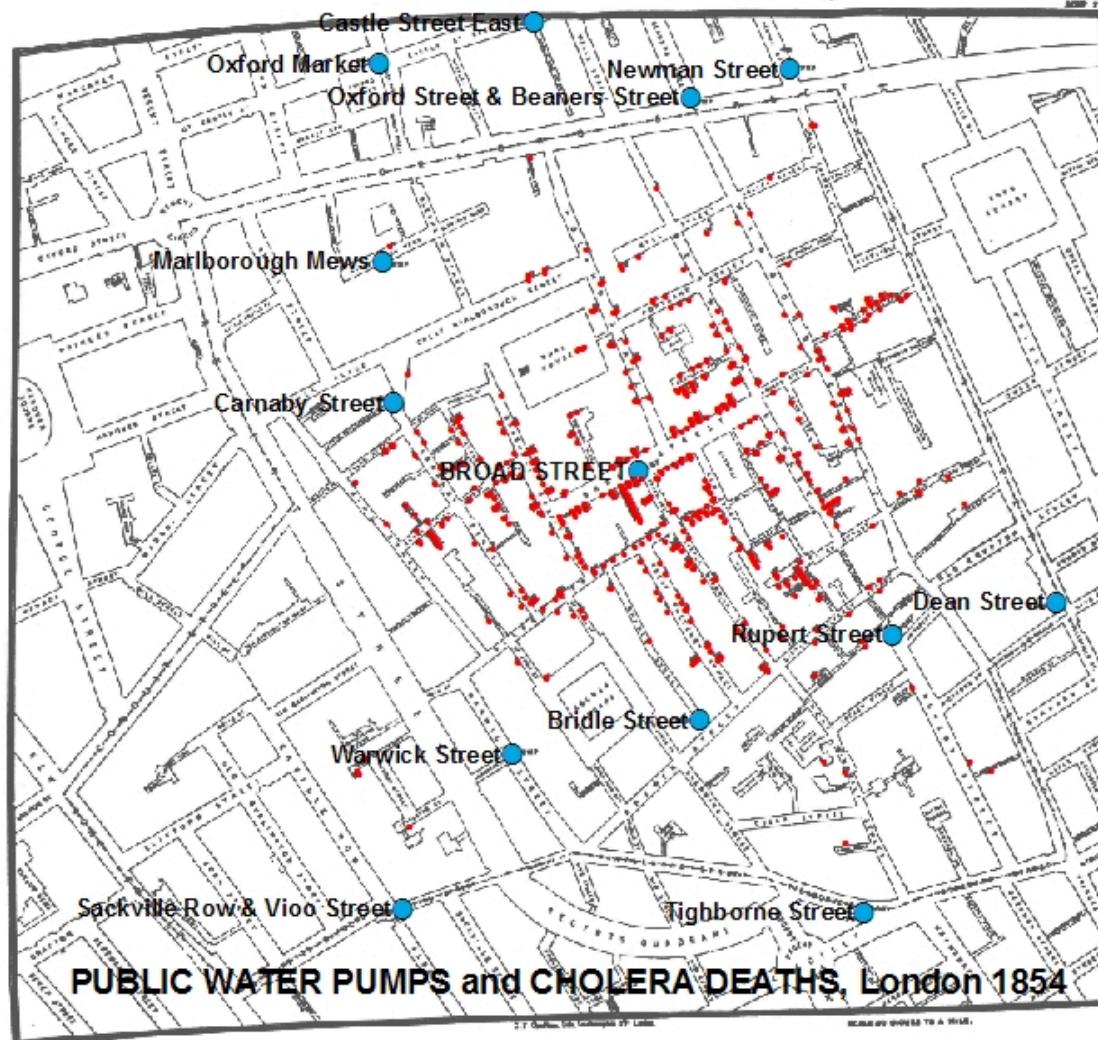
### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

#### 4.4.2 Regression

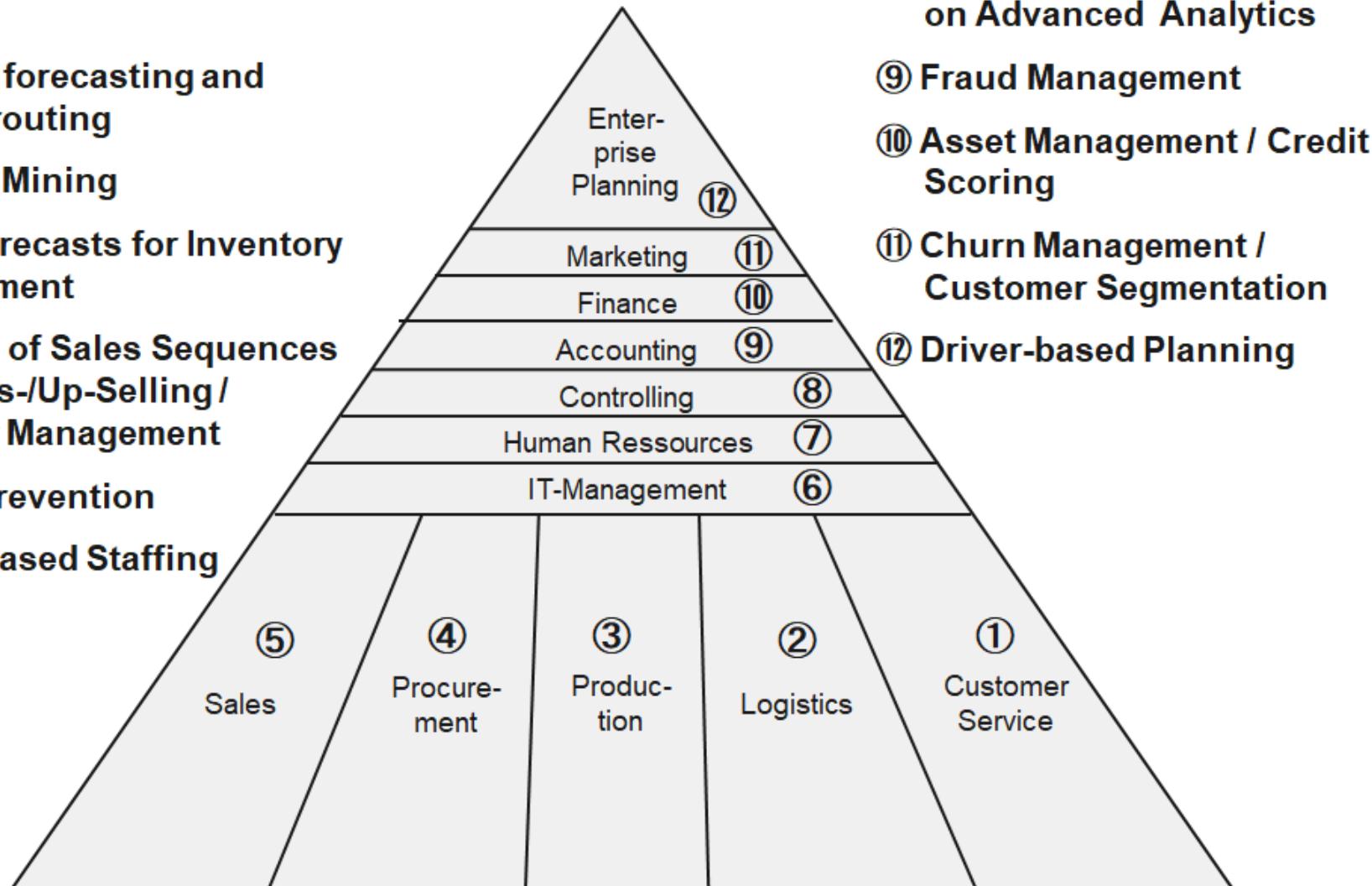
# First Application of Data Analytics

## Dr. John Snow's Map of the 1854 London Cholera Outbreak

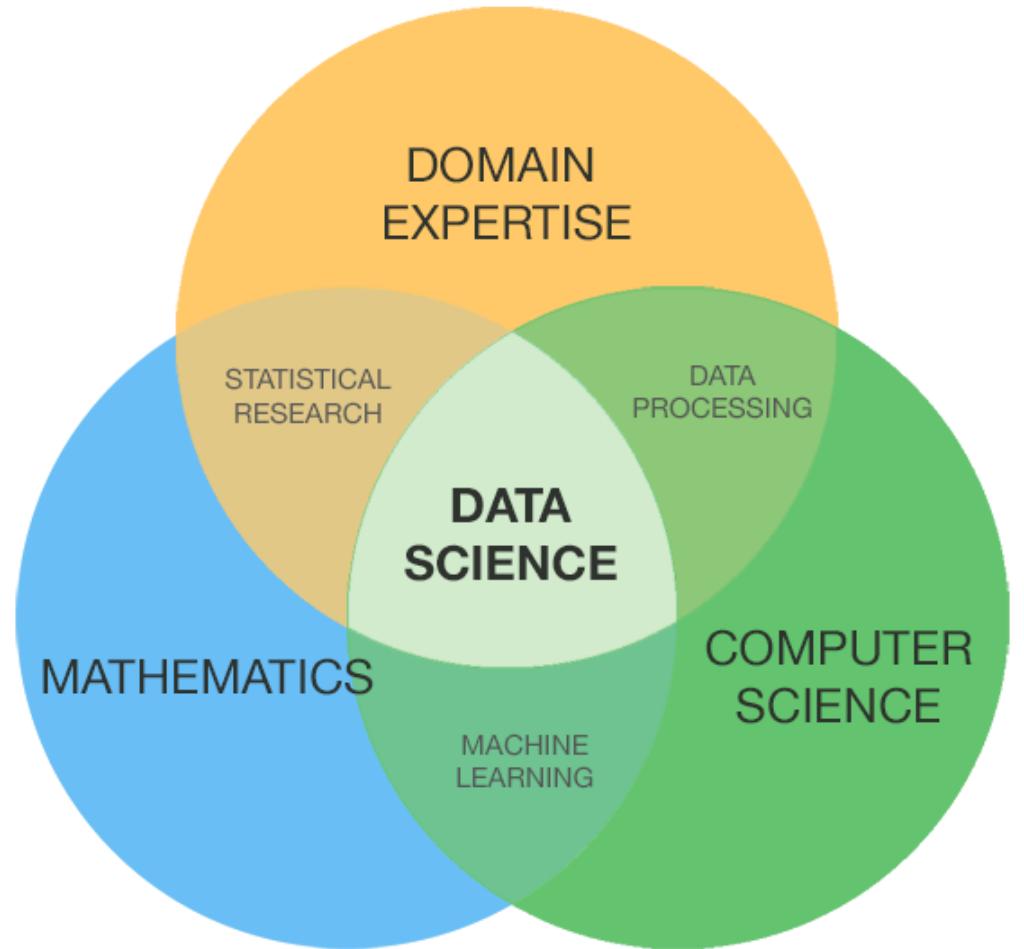


# Use Cases of Data Analytics - Examples

- ① Customer Support / Text Mining
- ② Demand forecasting and vehicle routing
- ③ Process Mining
- ④ Sales Forecasts for Inventory Management
- ⑤ Analysis of Sales Sequences for Cross-/Up-Selling / Voucher Management
- ⑥ Attack Prevention
- ⑦ Profile-based Staffing

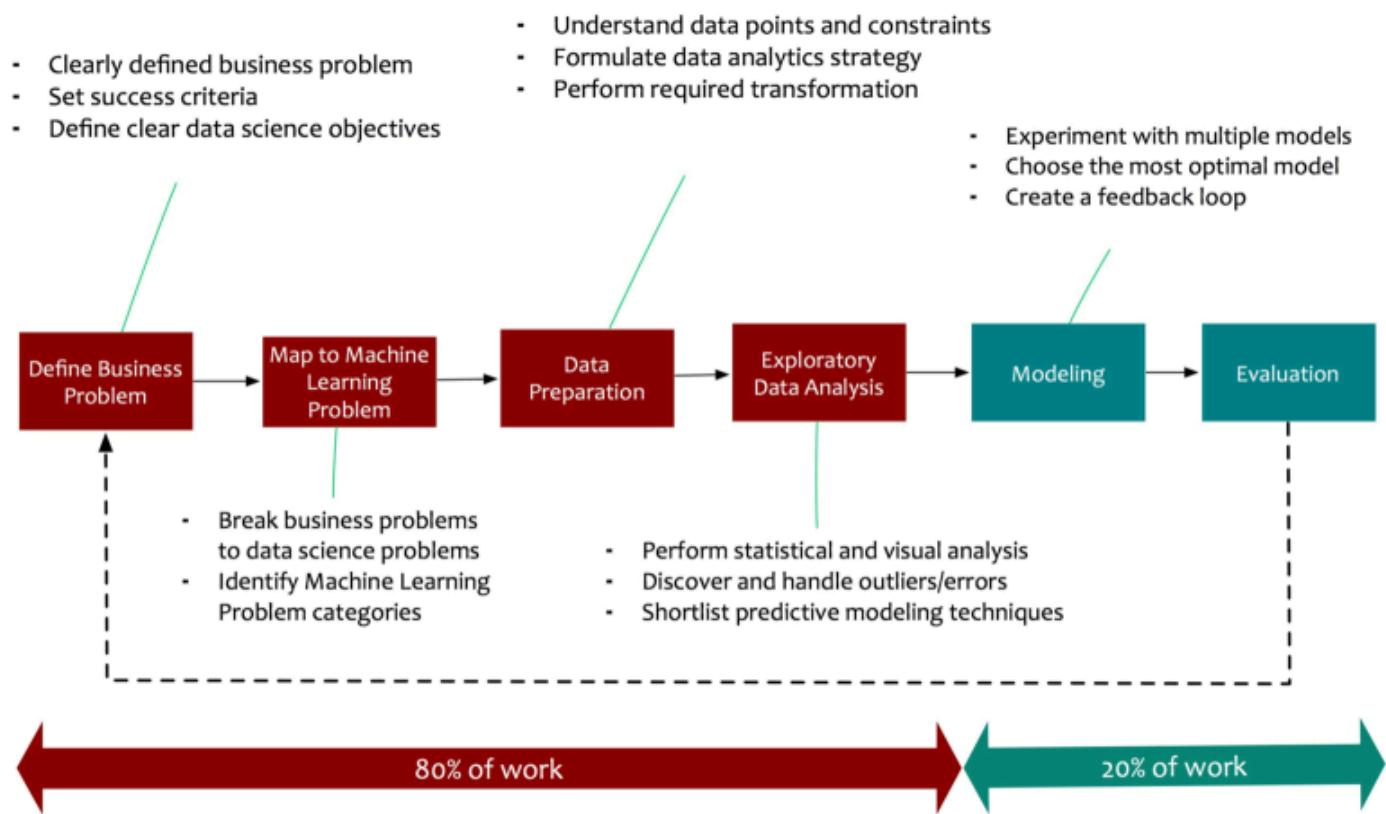


# Fundamental Skills of Data Analytics



Source: Palmer, Shelly. *Data Science for the C-Suite*.  
New York: Digital Living Press, 2015. Print.

# The Data Analytics Process



Source : <http://www.datasciencecentral.com/profiles/blogs/data-science-simplified-principles-and-process>

# Need for Knowledge about the Algorithms

The distance between using Excel and VBA for modeling in credit scoring, for example, and using machine learning algorithms and R or Python to enhance the results, is not that great, compared to the distance between someone running a packaged algorithm they don't really understand and someone who understands the mathematical and statistical operations within an algorithm and can optimize or adapt it as needed – and do so in the

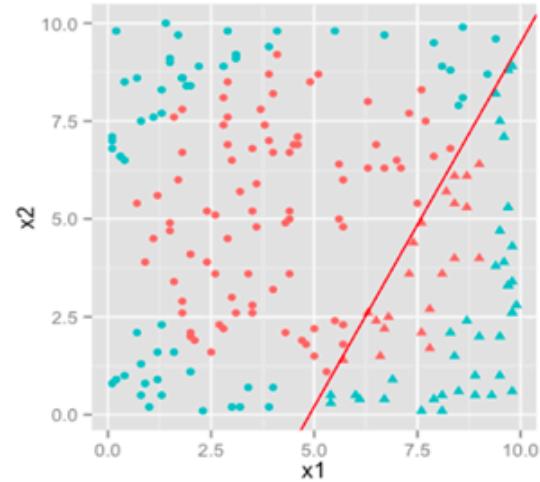
# Statistics vs. Machine Learning

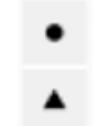
**Statistics** about finding valid conclusions about the underlying applied theory, and on the interpretation of parameters in their models. It insists on proper and rigorous methodology, and is comfortable with making and noting assumptions. It cares about how the data was collected and the resulting properties of the estimator or experiment (e.g. p-value). The focus is on hypothesis testing.

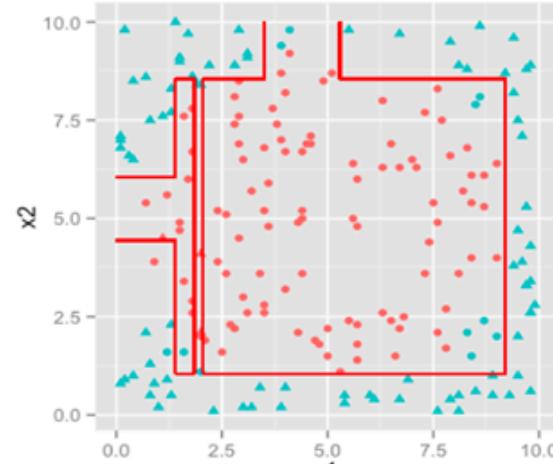
**Machine Learning (ML)** aims to derive practice-relevant

# Statistical Regression vs. Machine Learning Algorithms

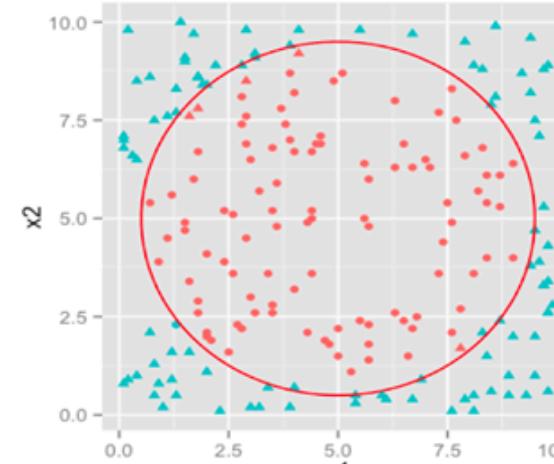
**Traditional Regression**



 **Real Values**  
 **Predicted Values**



**Decision Tree**



**Support Vector Machine**

# Explanation vs. Prediction (I)

Question 1: I have a headache. If I take an aspirin now, will it go away?

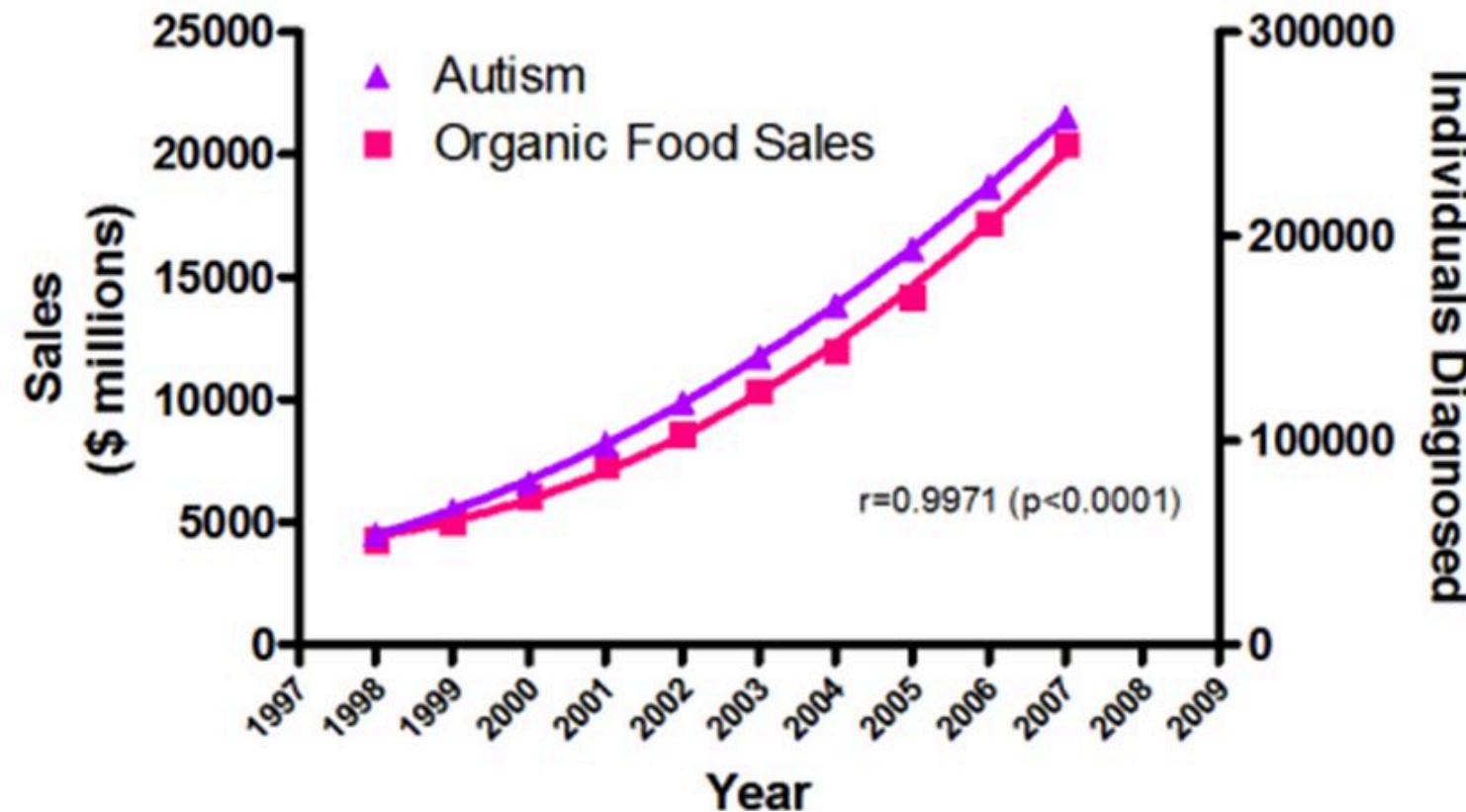
Question 2: I had a headache, but it passed. Was it because I took an aspirin two hours ago? Had I not taken such an aspirin, would I still have a headache?

## Explanation :

- Explanation is about understanding relationships and why certain things happen.
- It requires an understanding of cause and effect.
- Tests of causal hypotheses are fundamental.
- Measures of significance are central.
- A good explanatory model may also have predictive power.

## Prediction :

# Correlation



Source: Organic Trade Association, 2011 Organic Industry Survey, U.p.  
Department of Education, Office of Special Education Programs, Data  
Analysis System (DANS)

Organic food sales and the rate of autism seem to have a

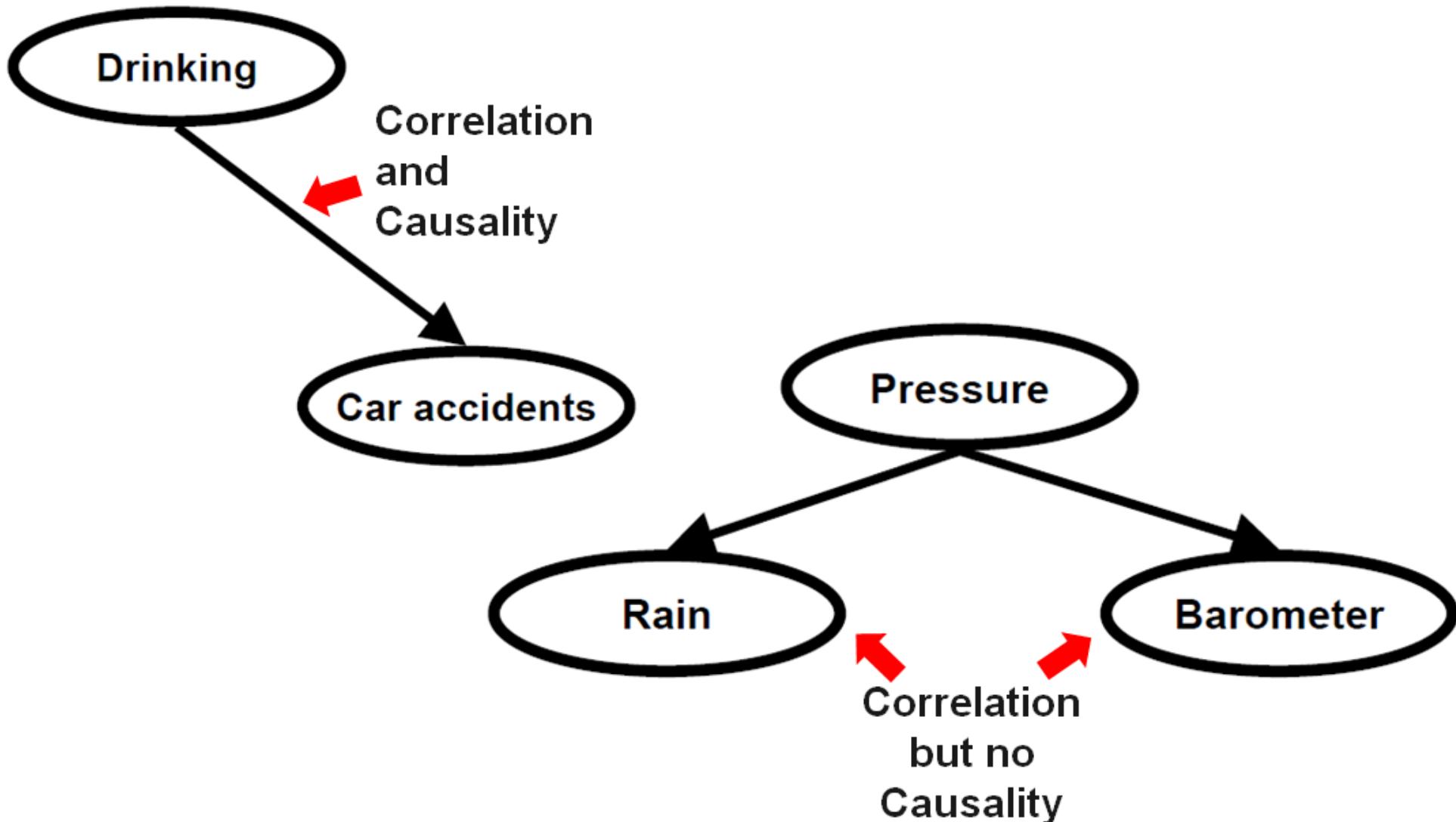
# Correlation vs. Causality (I)



**Correlation:** Two data series behave “similar”

**Causality:** Principle of Cause and Effect

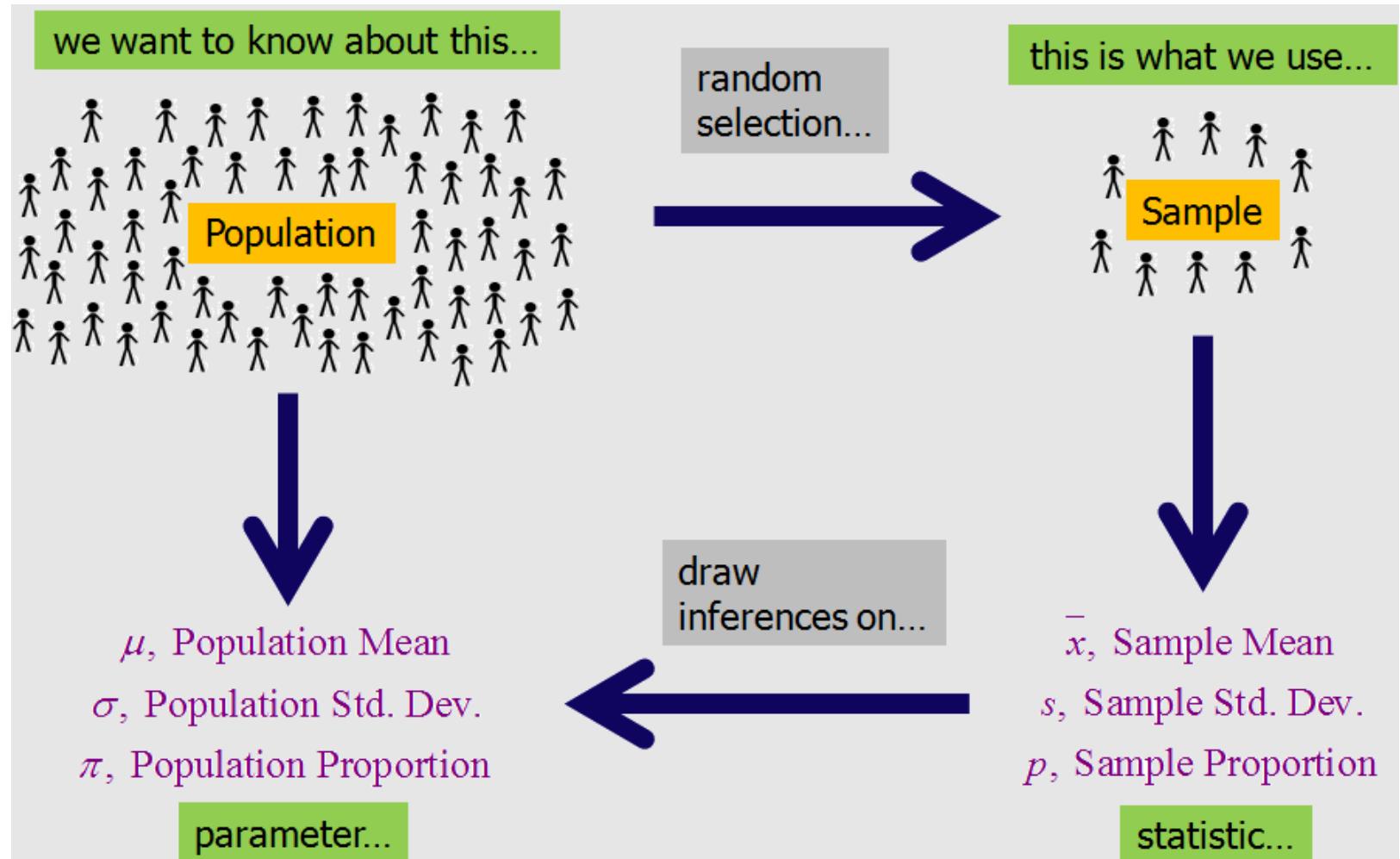
# Correlation vs. Causality (II)



But:

Sometimes it is better  
to know/predict  
something even if we

# Statistical Estimation



Source: <http://www.dxbydt.com/the-size-of-your-sample>

## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

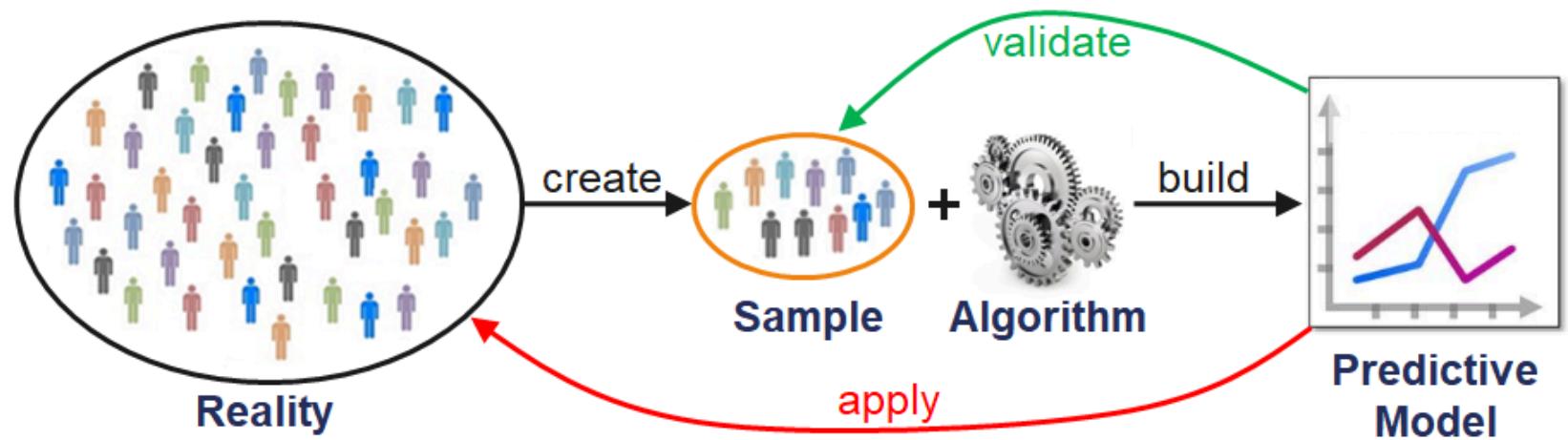
#### 4.4.2 Regression

A **method** is a composition of formalized principles that form the basis for a stringent calculation process.

An **algorithm** is a procedure or set of steps or rules to accomplish a task. It is usually the implementation of a method. Algorithms are used to build models.

In the given context, a **model** is the description of the relationship between variables. It is used to create output data from given input data, for example to make predictions.

# Traditional Analytics Process



# Example Regression - Fitting the model

Age	Rooms	Dist_School	Dist_Mall	Lakeside	Sales_Price
42	6	4.4	1.1	0	392000
35	6	4.7	12.4	0	308000
32	7	1.5	1.5	0	429000
39	5	2.8	10.2	1	338000
18	6	0.1	0.3	0	450000
43	6	2.5	13.5	0	304000
34	5	2.4	3.3	0	349000
49	6	1.4	10.2	0	326000
10	6	0.7	8.8	0	369000
8	6	2.2	1.1	0	413000
22	5	4.9	11.7	0	306000
3	7	4.7	11.5	0	378000
25	4	2.5	3.8	0	348000
38	4	1.4	0.7	0	367000
30	5	2.8	11.6	0	299000
26	6	2.2	8.1	0	350000
38	7	4.2	6.6	0	391000
2	6	4.6	10.4	0	343000
26	5	4.9	1	0	353000
31	5	2.5	5.6	0	365000
44	5	4.9	4.9	1	380000
25	7	1	14.1	0	353000
8	5	0.2	13.3	0	331000
29	4	1.4	2.7	0	351000
19	4	0.8	9.7	0	304000
26	7	2.2	3.2	1	466000
43	4	2.8	0.6	1	379000
24	4	2.6	12.5	1	301000
45	5	1.2	1.6	0	390000
28	6	4.4	3.2	0	383000
35	4	2	0.1	0	371000
48	6	2.4	11.3	0	301000
0	6	2.7	10.7	1	386000
18	4	2.4	0.7	0	387000
30	7	3.8	9.5	1	414000
42	6	0.9	1.4	0	397000
1	4	3	1.9	0	387000
26	6	0.4	1.1	1	391000

SUMMARY OUTPUT					
Regression Statistics					
Multiple R	0.965794107				
R Square	0.932758256				
Adjusted R Square	0.924753287				
Standard Error	11962.63096				
Observations	48				
ANOVA					
	df	SS	MS	F	Significance F
Regression	5	83374421845	16674884369	116.5223999	1.71591E-23
Residual	42	6010390655	143104539.4		
Total	47	89384812500			
	Coefficients	Standard Error	t Stat	P-value	
Intercept	309694.3587	10414.1276	29.73790707	7.64477E-30	
Age	-978.7070888	118.4386216	-8.263411678	2.40669E-10	
Rooms	26247.68844	1788.306679	14.67739776	3.88184E-18	
Dist_School	-6684.949298	1220.872194	-5.475552096	2.24688E-06	
Dist_Mall	-8002.138709	400.8778687	-19.96153775	4.72709E-23	
Lakeside	38025.024	4212.289608	9.027162788	2.17287E-11	

Proportion of variance explained by the model.

Probability that the null hypothesis that *all* of the regression coefficients are equal to zero.

Probability that the null hypothesis  $\beta = 0$  is true.  
The coefficients are highly likely to be nonzero and therefore significant.

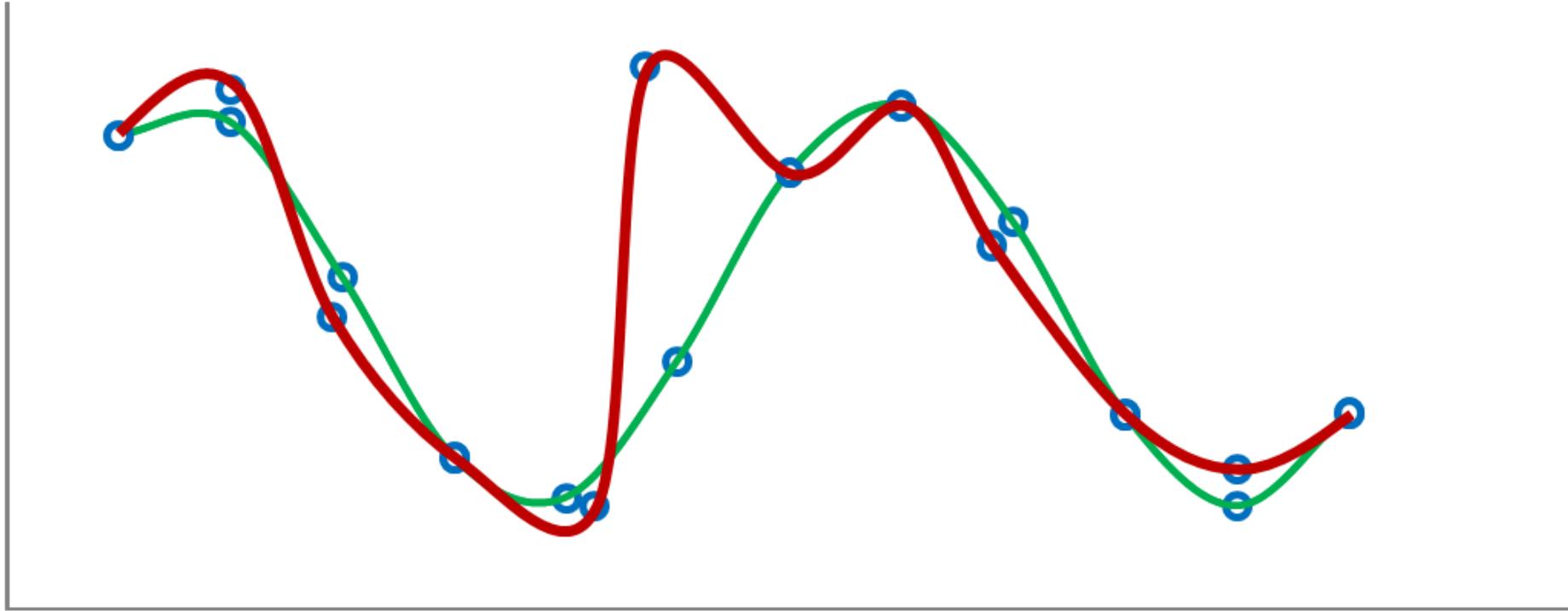
# Example Regression - Testing the model

Constant	Age	Rooms	Dist_School	Dist_Mall	Lakeside	Sales_Price	Estimation	Deviation	Dev_Mean
1	16	5	3.3	10.9	0	299000	315989.843	288654761.06	4364253906.25
1	31	4	0.6	12.8	0	270000	277906.848	62518240.02	9036878906.25
1	20	5	1	1.1	0	420000	405871.357	199618545.66	3018128906.25
1	13	5	2	3.8	0	377000	384431.583	55228426.91	142503906.25
1	31	5	3.7	7.1	0	323000	329043.384	36522489.25	1769253906.25
1	23	6	3	10	0	329000	344593.991	243172565.38	1300503906.25
1	36	7	3.5	6.1	0	340000	385984.354	2114560805.77	628128906.25
1	33	6	5	7.5	1	394000	379467.393	211196677.41	837378906.25
1	47	6	0	0.3	0	429000	418780.615	104435838.65	4088003906.25
1	5	5	2.3	1.7	0	390000	407060.246	291052003.12	621878906.25
1	8	4	2.5	12.4	0	299000	290916.563	65341961.39	4364253906.25
1	17	4	1.1	0.5	0	380000	386692.578	44790605.45	223128906.25
1	39	5	4.5	6.8	0	295000	318266.409	541325805.81	4908753906.25
1	17	7	3	14.6	0	260000	339904.084	6384662676.66	11038128906.25
1	6	7	4.8	6.7	1	447000	439878.873	50710445.81	6713753906.25
1	2	4	4.3	10	0	289000	303961.029	223832395.48	5785503906.25
1	13	4	3.3	12.9	0	289000	276673.998	151930318.16	5785503906.25
1	6	4	0.4	10.7	1	371000	358541.03	155225932.53	35253906.25
1	17	5	2.1	7.4	0	336000	351040.56	226218458.07	844628906.25
1	11	5	4.7	5	1	391000	396762.092	33201700.68	672753906.25
1	37	5	2.3	5	1	408000	387359.586	426026702.49	1843628906.25
1	19	7	1.3	9.9	1	417000	424946.16	63141455.48	2697503906.25
1	41	6	4.9	5.6	1	361000	387510.294	702795708.22	16503906.25
1	18	7	4.1	7.6	0	394000	387586.904	41127801.86	837378906.25
1	4	7	3.2	14.4	0	362000	352890.714	82979086.57	9378906.25
1	48	7	2.8	4.6	1	421000	428947.565	63163796.31	3129003906.25
1	47	7	4.3	8.2	0	349000	353066.125	16533374.30	258003906.25
1	25	4	2.8	13.4	0	295000	264270.919	944276448.70	4908753906.25
1	47	5	2.9	9	1	357000	341552.99	238610105.71	65003906.25
1	1	6	1.4	11	1	390000	406844.351	283732175.62	621878906.25
1	49	7	3.5	10.1	0	350000	341252.607	76516885.48	226878906.25
<b>Sum:</b> 14423104194.01 <b>80792496093.75</b>									
<b>Pseudo-R<sup>2</sup>:</b> 82.1480%									

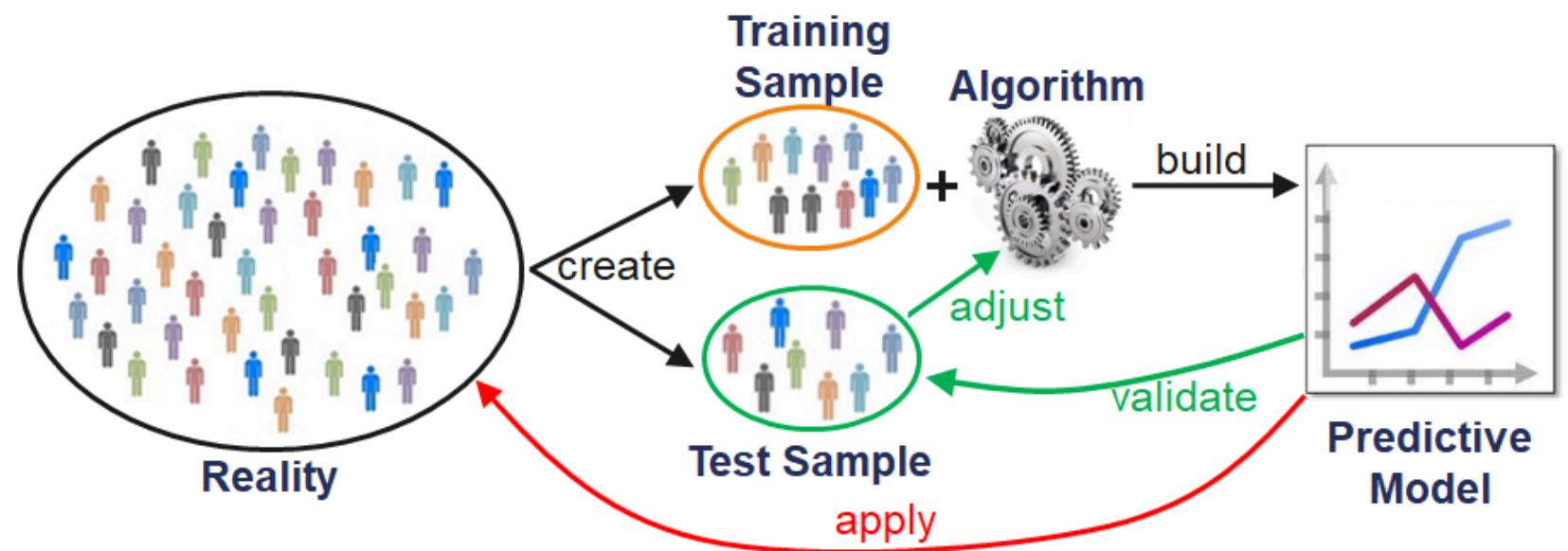
Conclusion:  
Model is not robust due to overfitting.



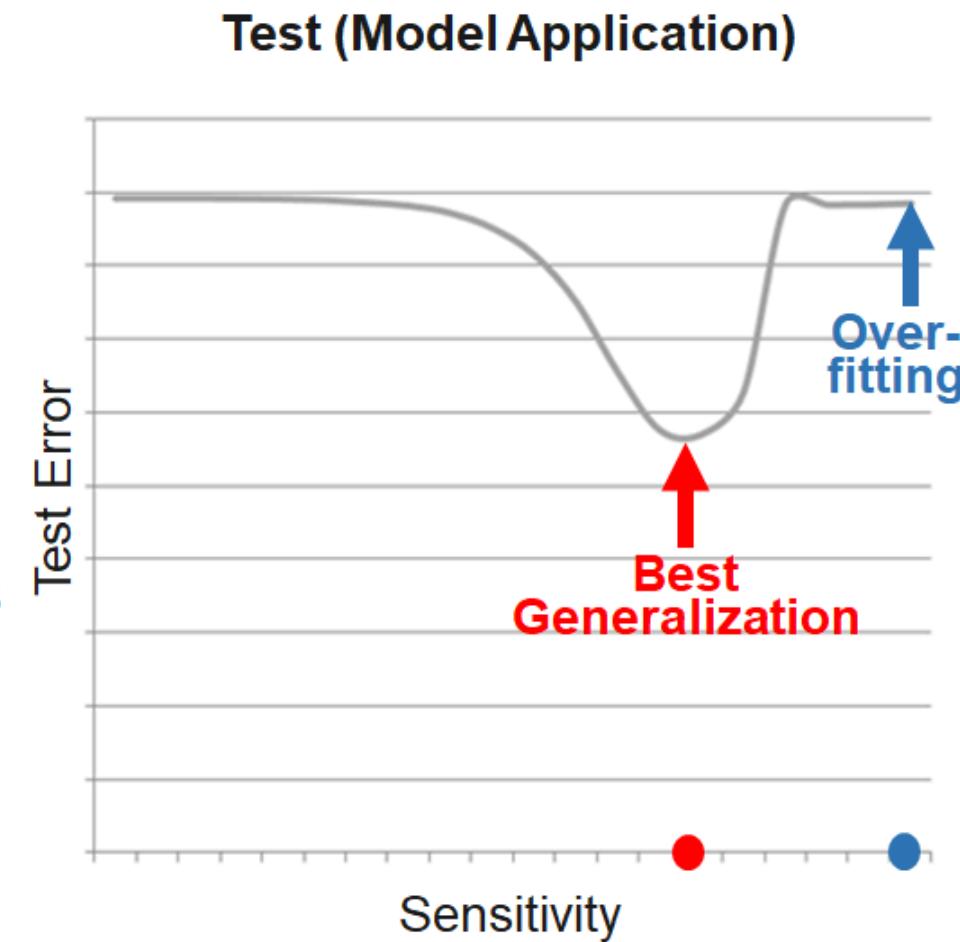
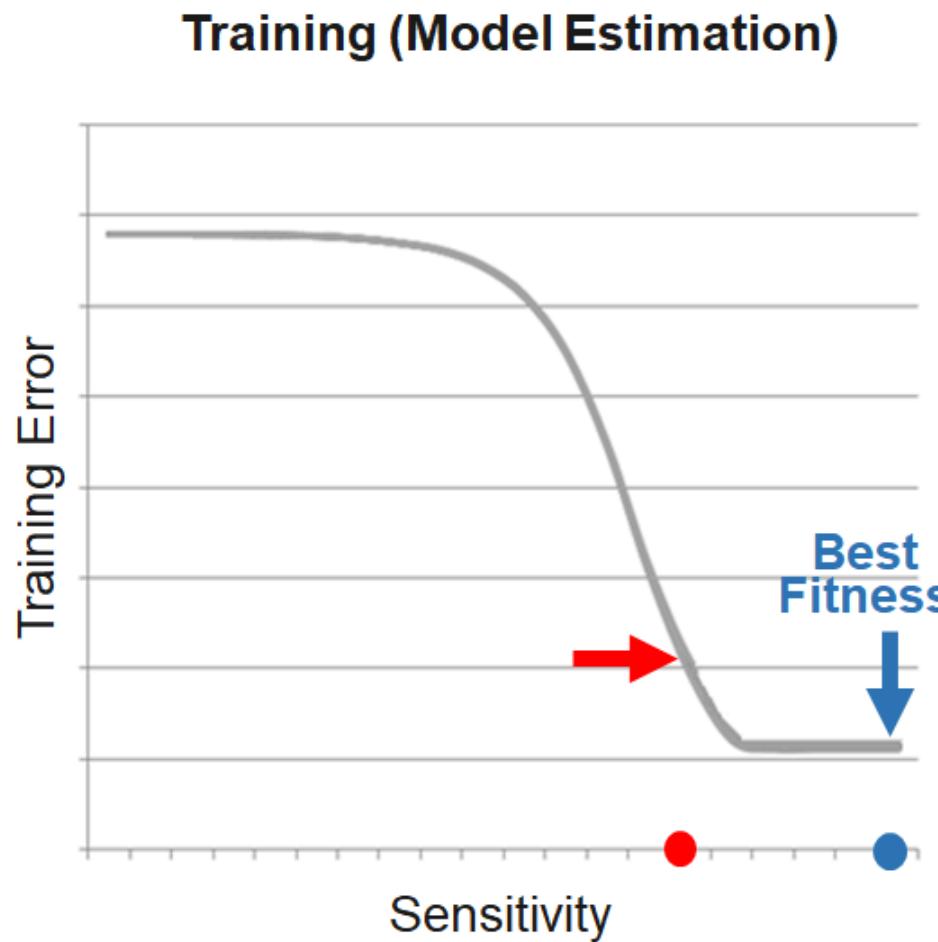
# Data Errors and their Consequences



# Modern Analytics Process



# Best Fit vs. Best Generalization



# Over- and Underfitting



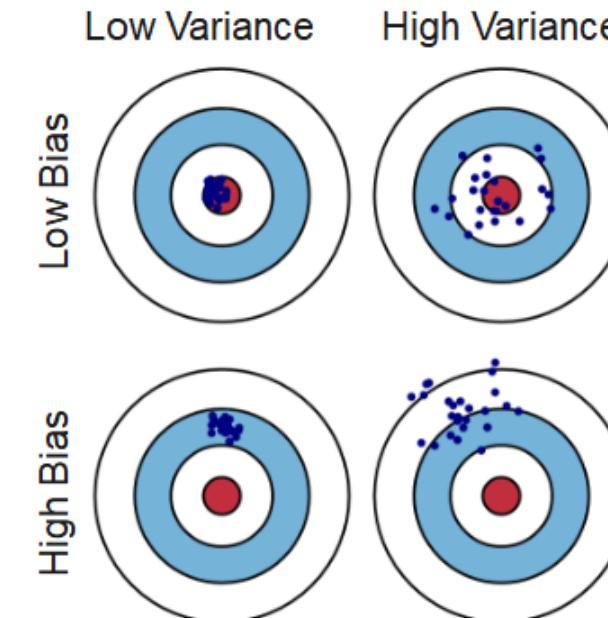
Due to the problem of overfitting, the main goal is to maximize the prediction quality and not to fit the data that is used for the model estimation as well as possible. This is

# The Bias-Variance Tradeoff

The prediction error is influenced by three components:

$$\text{Error} = \text{Bias} + \text{Variance} + \text{Noise}$$

Bias is the inability of the used method to learn the relevant relations between the inputs and the outputs. It reflects the method quality, e.g. if a method only produces linear models.



Variance is represents the deviation

# Summarizing: Statistics vs. Data Analytics

## Statistics

Goals: Explanation, Forecasting

Methods: Statistical Methods

Process:  → Method → Model → generalizable?  
(Statist. Tests)

## Data Analytics

Goals: Forecasting, Finding Structures

Methods: Statistical Methods, Machine Learning, Artificial Intelligence

Process:  → Method → Model →   

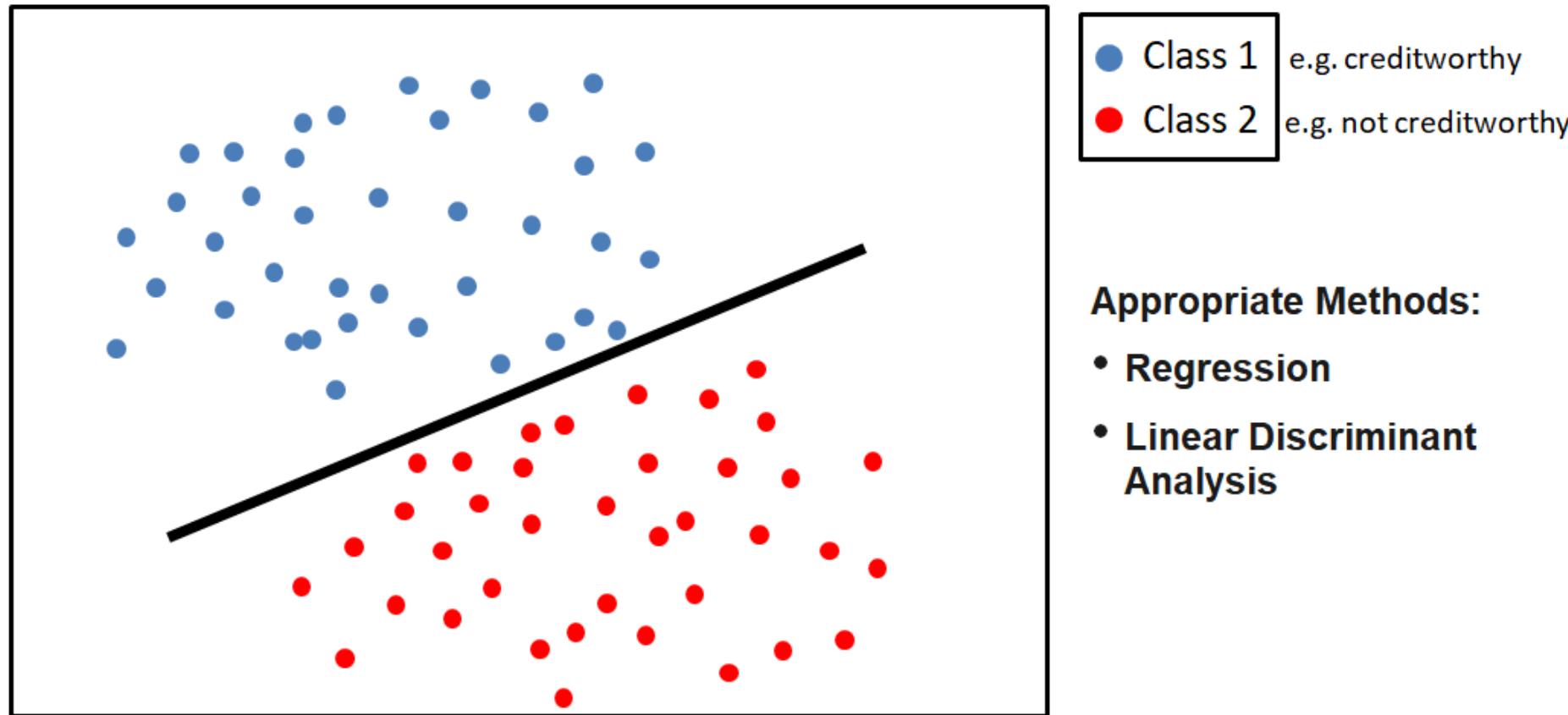

# Which Method should I choose?

The choice of the method of data analysis depends on the one hand on the scope of application, but on the other hand on the interrelationships of the data to be analyzed.

In the Big Data area, data spaces are often highly-dimensional, making it difficult to visualize the interrelationships.

For this reason, the choice of the method can often not be made ex ante. In these cases, different methods are

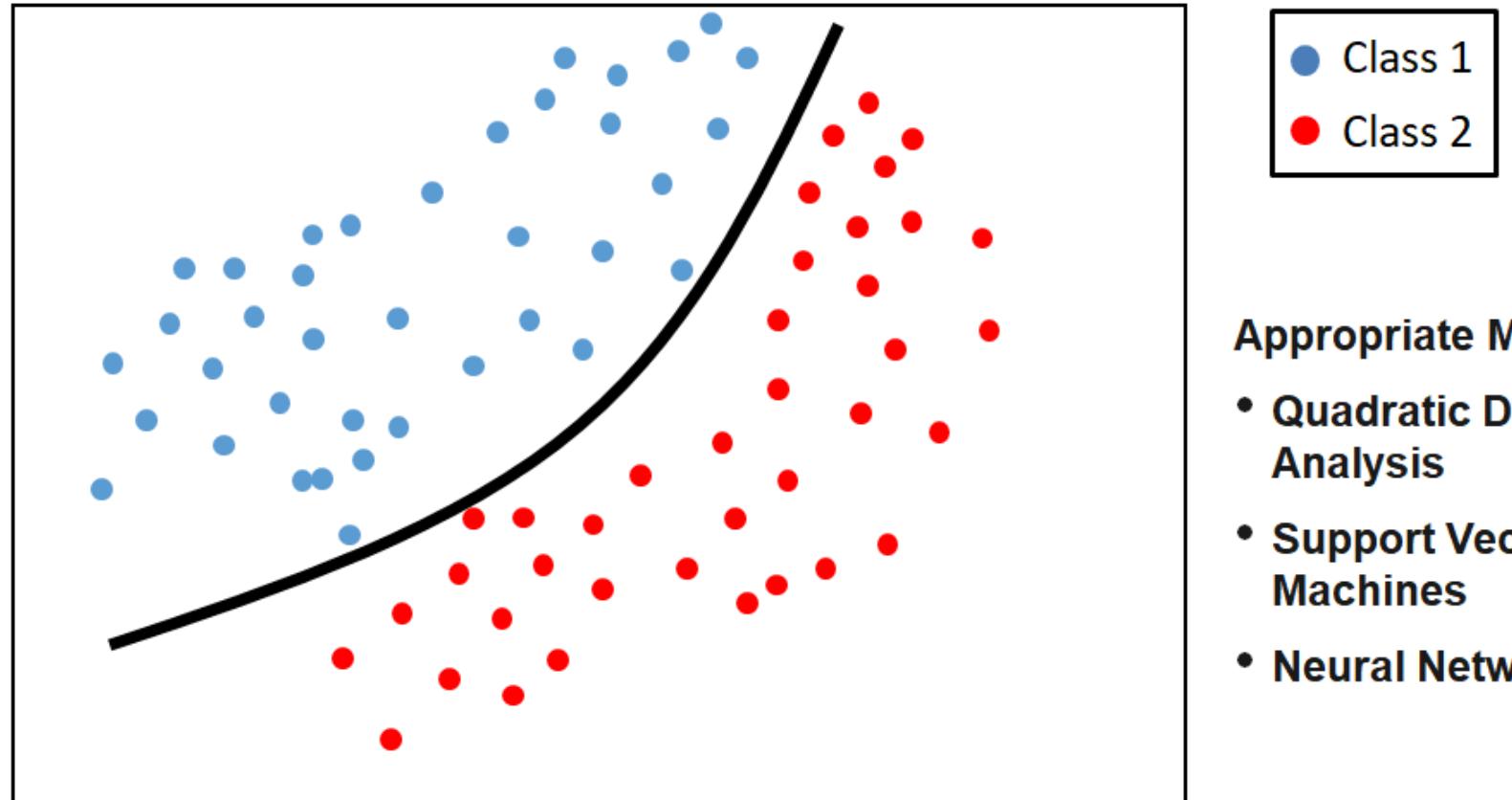
# Linear World



## Appropriate Methods:

- Regression
- Linear Discriminant Analysis

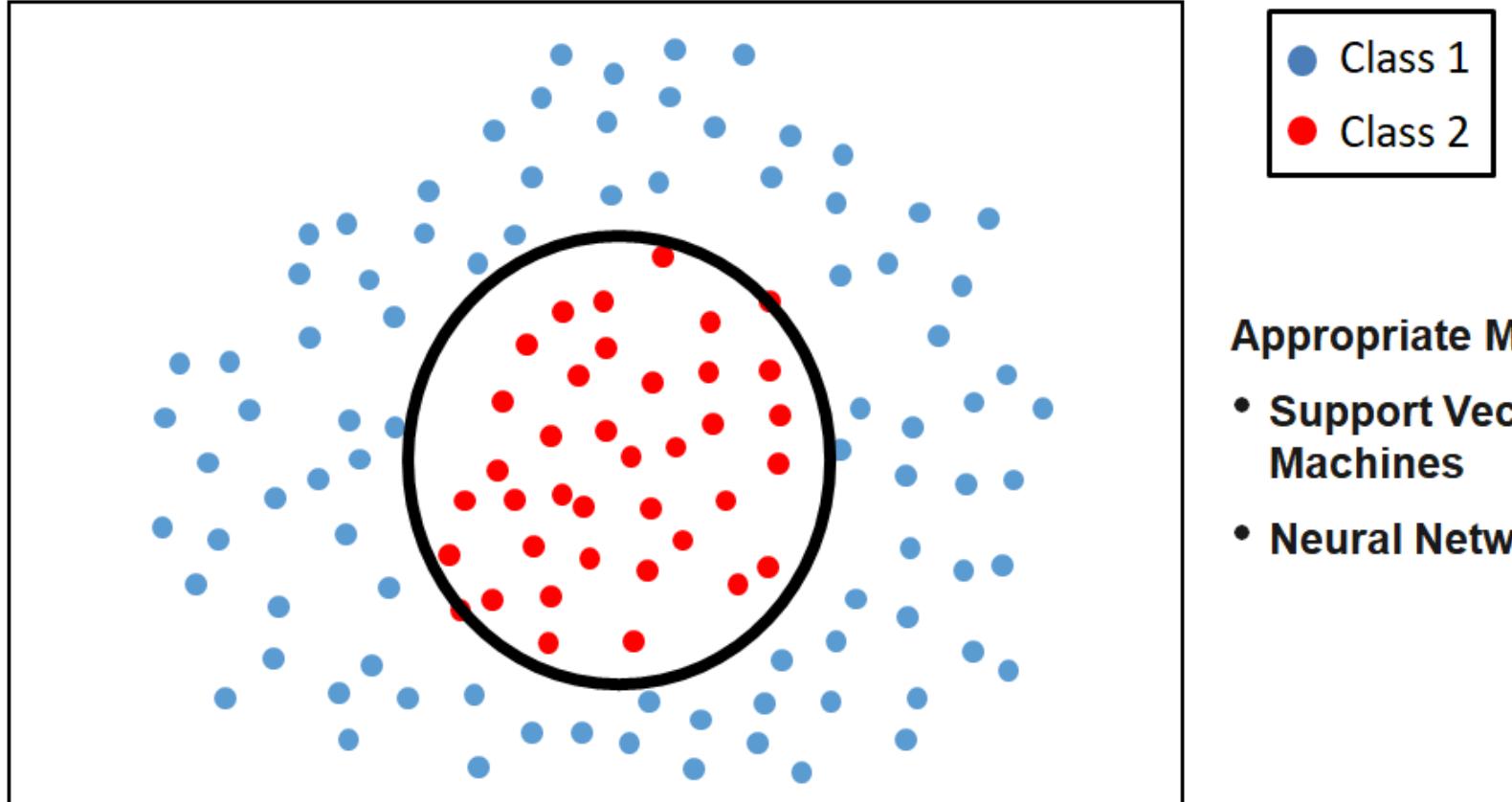
# Quadratic World



## Appropriate Methods:

- Quadratic Discriminant Analysis
- Support Vector Machines
- Neural Networks

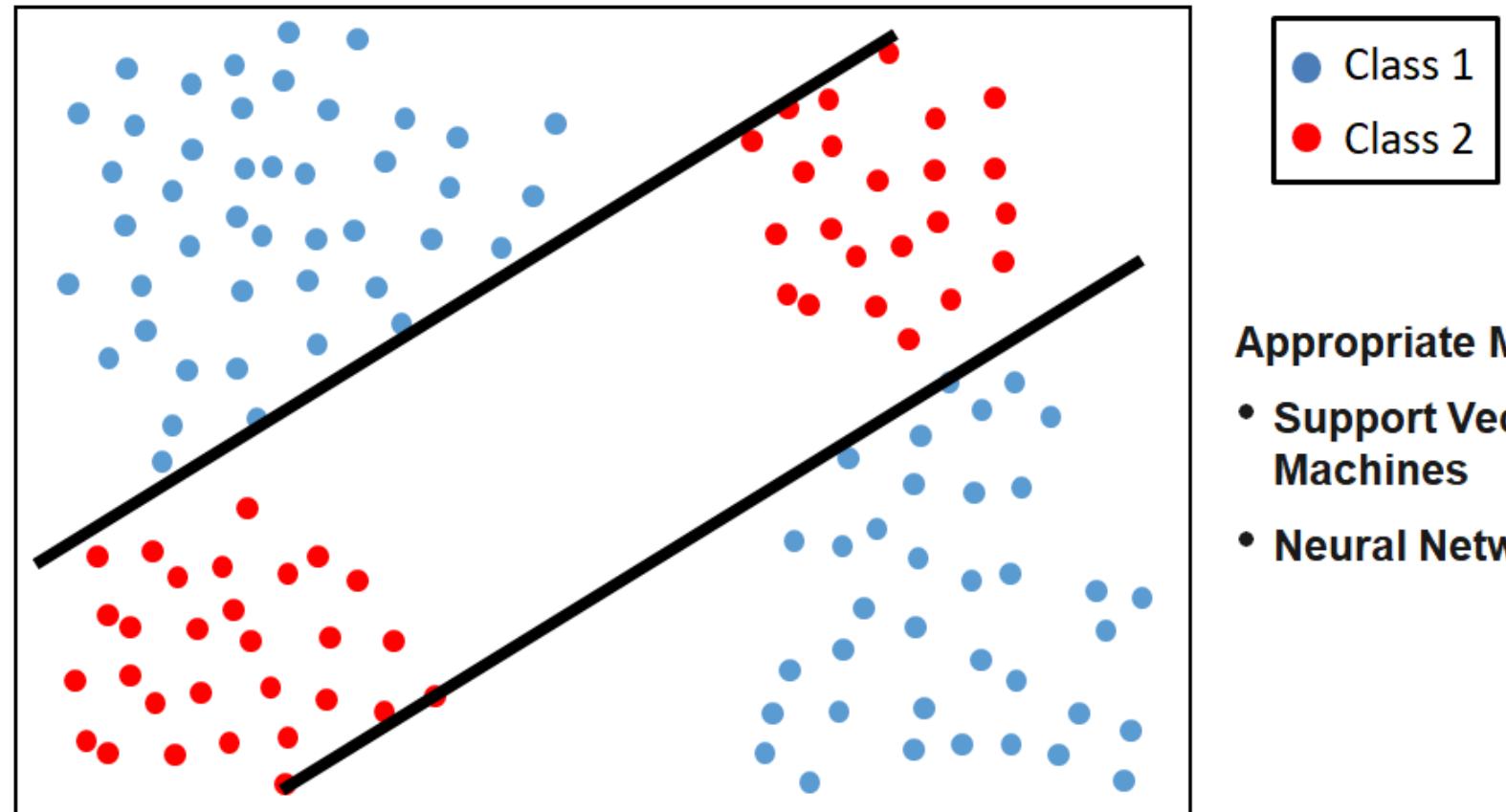
# Nonlinear World (Type 1)



## Appropriate Methods:

- Support Vector Machines
- Neural Networks

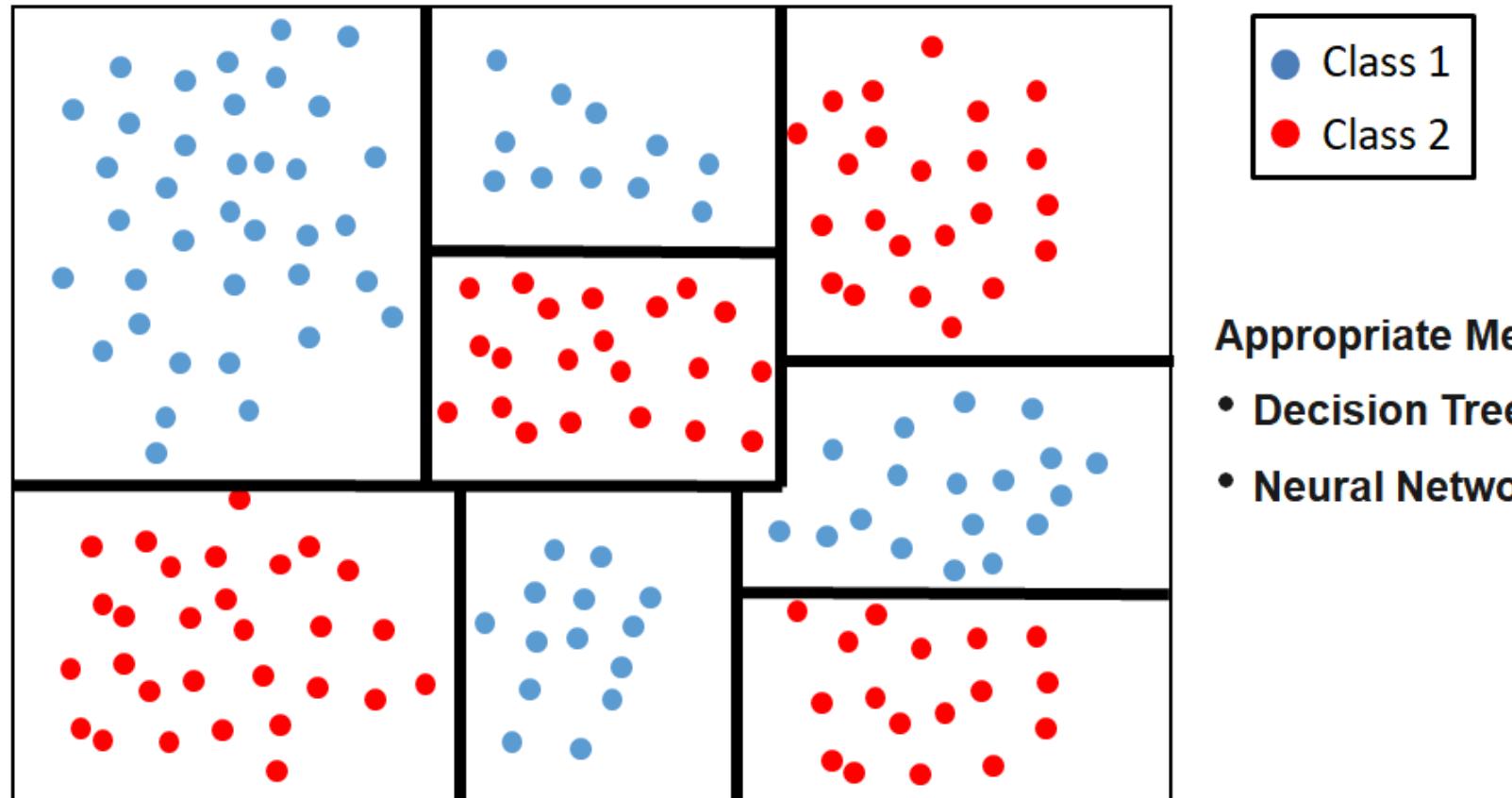
# Nonlinear World (Type 2)



Appropriate Methods:

- Support Vector Machines
- Neural Networks

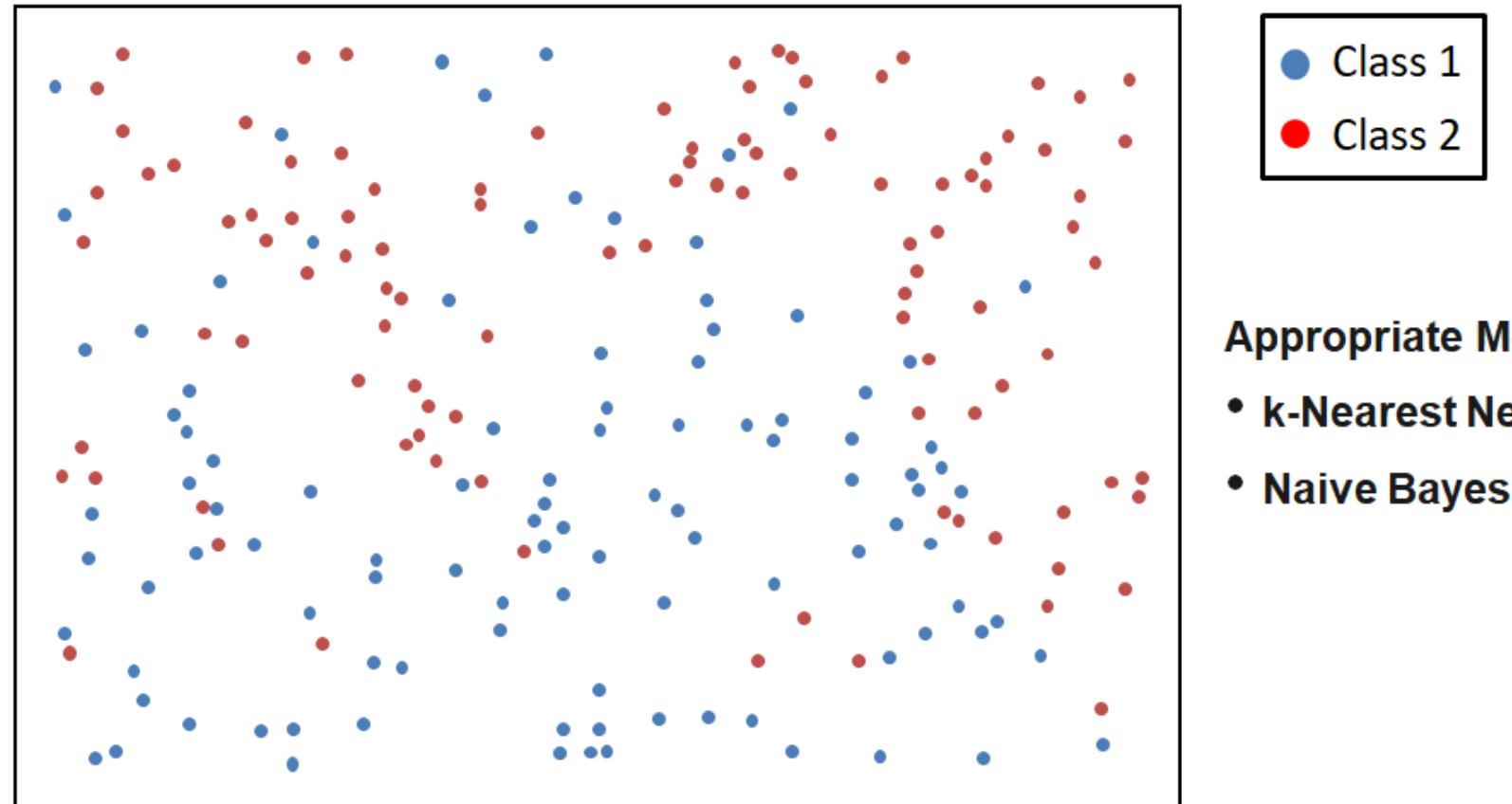
# Nonlinear World (Type 3)



## Appropriate Methods:

- Decision Trees
- Neural Networks

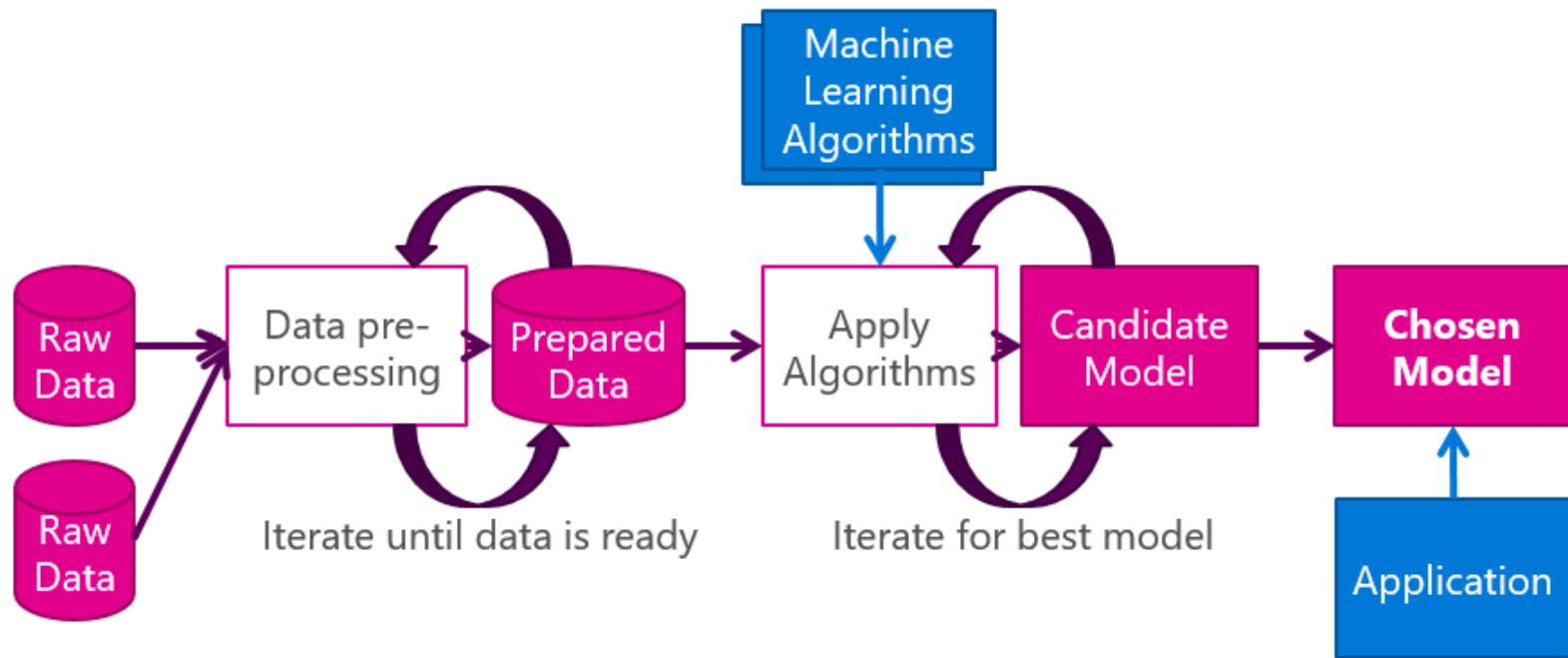
# Nonlinear World (Type 4)



**Appropriate Methods:**

- **k-Nearest Neighbors**
- **Naive Bayes**

# The Data Analytics Process - Technical View



Source: <http://blogs.msdn.microsoft.com/martinkearn/2016/03/01/machine-learning-is-for-muggles-too/>

## 4 Predictive Analytics

4.1 Subject of Predictive Analytics

4.2 The Analytics Process

4.3 Data Preparation

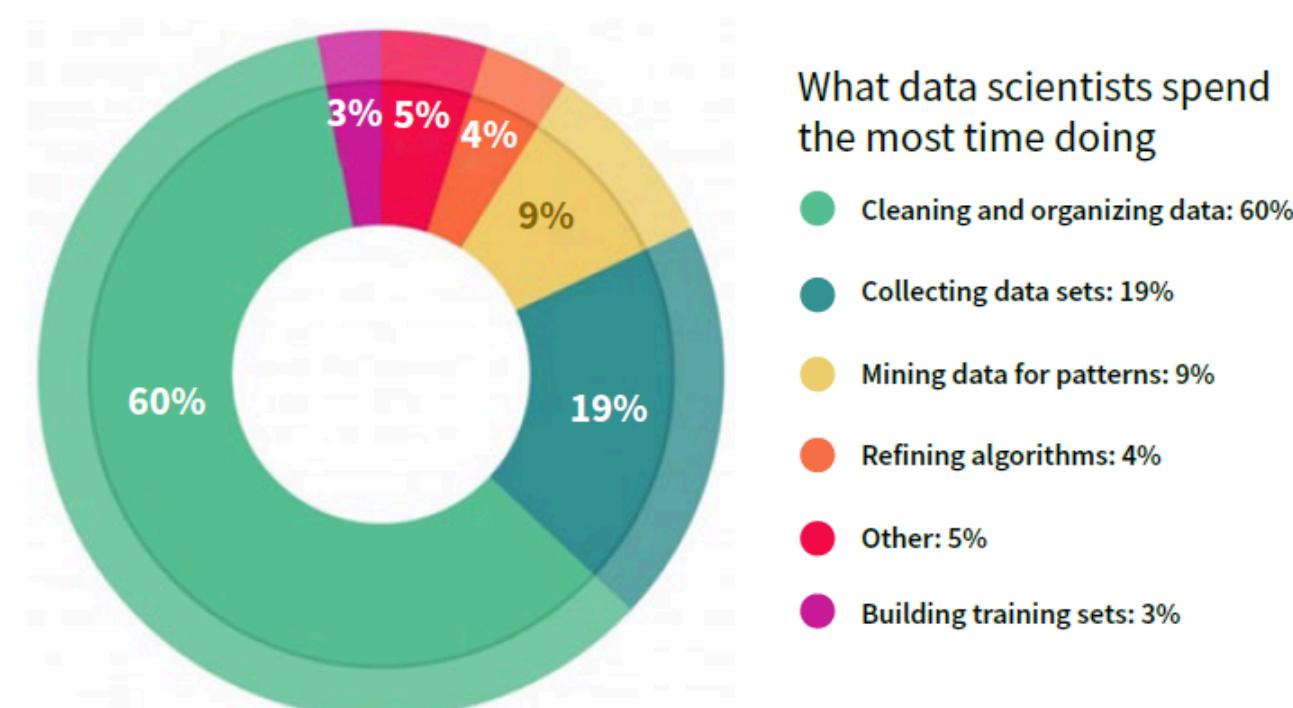
4.4 Methods, Algorithms and Applications

4.4.1 Classification

4.4.2 Regression

# Data Preparation and Enrichment

The data collection and preparation phase is the most labor-intensive one, consuming on average between 60-80% of a data scientist's time. It's critical therefore to select a tool that can automate or at least speed the workflows associated with data preparation.



## 1. Proof of correctness of the data

- examine for irregular outliers (e.g. Age=236)
- examine for typographical errors (e.g. Frankfrut)
- examine for different writing styles (e.g. behavior/behaviour)
- —

## 2. Handling missing values

eliminate → all objects having missing values (rowwise)  
variables having missing values first, then remaining objects

replace → with the true values  
with mean values

# Missing Values Strategies

## Initial Situation

Variable 1	Variable 2	Variable 3	Variable 4	Variable 5
8.95	0.99	NA	370.53	-34.97
NA	NA	168.02	617.46	-99.49
18.38	0.48	58.06	1177.66	-44.21
16.10	NA	35.39	NA	77.35
23.14	NA	230.44	408.97	-71.79
24.77	0.92	206.83	NA	46.13
20.72	NA	228.74	659.37	-67.77
15.52	0.04	22.46	NA	36.63



## Stupid Strategy

Variable 1	Variable 2	Variable 3	Variable 4	Variable 5
8.95	0.99	NA	370.53	-34.97
NA	NA	168.02	617.46	-99.49
18.38	0.48	58.06	1177.66	-44.21
16.10	NA	35.39	NA	77.35
23.14	NA	230.44	408.97	-71.79
24.77	0.92	206.83	NA	46.13
20.72	NA	228.74	659.37	-67.77
15.52	0.04	22.46	NA	36.63

## Smart Strategy

Variable 1	Variable 2	Variable 3	Variable 4	Variable 5
8.95	0.99	NA	370.53	-34.97
NA	NA	168.02	617.46	-99.49
18.38	0.48	58.06	1177.66	-44.21
16.10	NA	35.39	NA	77.35
23.14	NA	230.44	408.97	-71.79
24.77	0.92	206.83	NA	46.13
20.72	NA	228.74	659.37	-67.77
15.52	0.04	22.46	NA	36.63

### Step 1:

Variable 1	Variable 3	Variable 5
8.95	NA	-34.97
NA	168.02	-99.49
18.38	58.06	-44.21
16.10	35.39	77.35
23.14	230.44	-71.79
24.77	206.83	46.13
20.72	228.74	-67.77
15.52	22.46	36.63

### Step 2:

A population can be defined as including all people or items with the characteristic one wishes to understand.

Sampling is about to find a representative subset of that population.

Data represents the traces of the real-world processes, and exactly which traces we gather are decided by our sampling method.

There are two sources of randomness and uncertainty:

1. the randomness and uncertainty underlying the process

**Question:**

**Is there any need for sampling in times of Big Data?**

**Why not “N=ALL”?**

**Answer:**

**Data is not objective! Data does not speak for itself.**

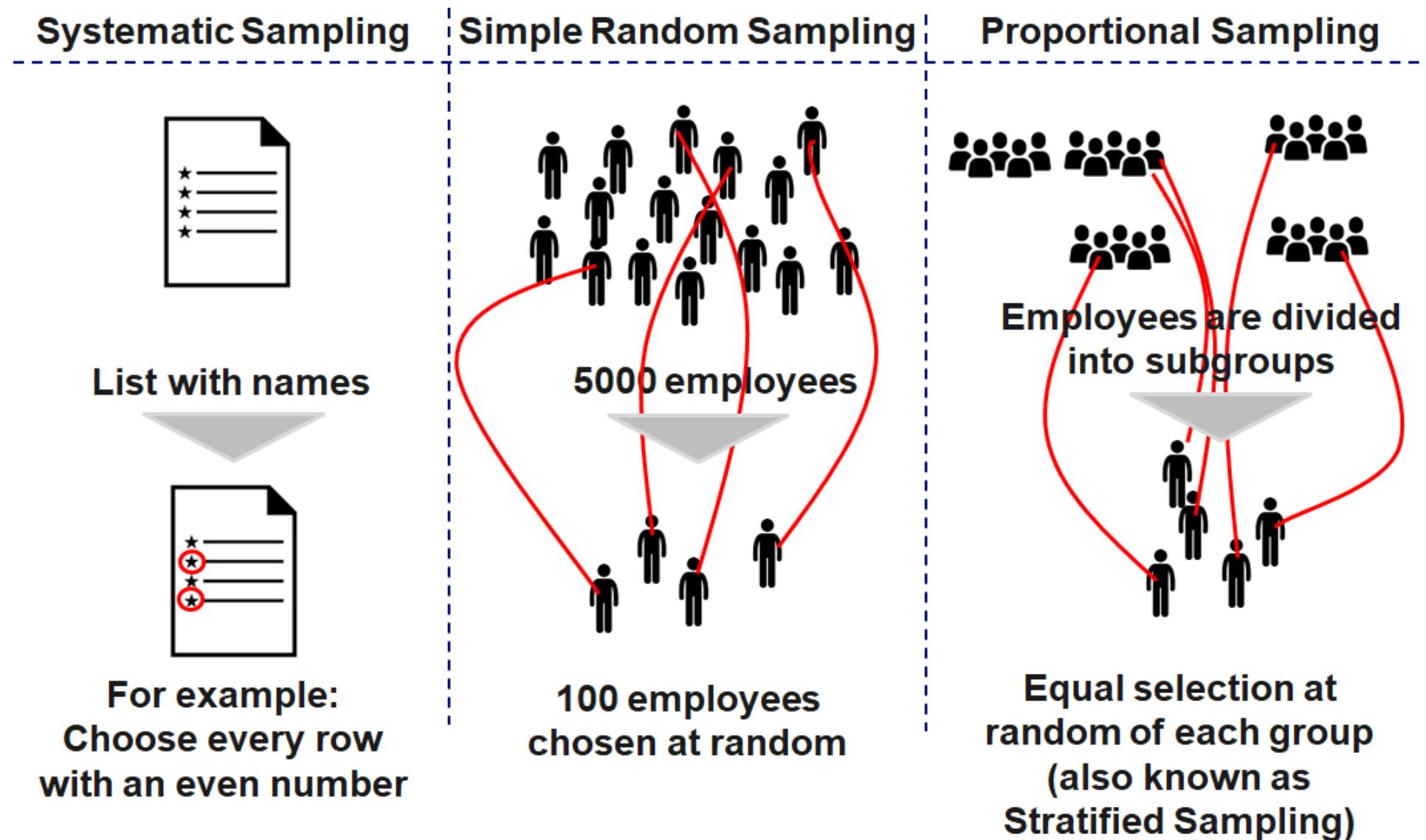
**Data is just a quantitative echo of the events of our society.**

**Examples:**

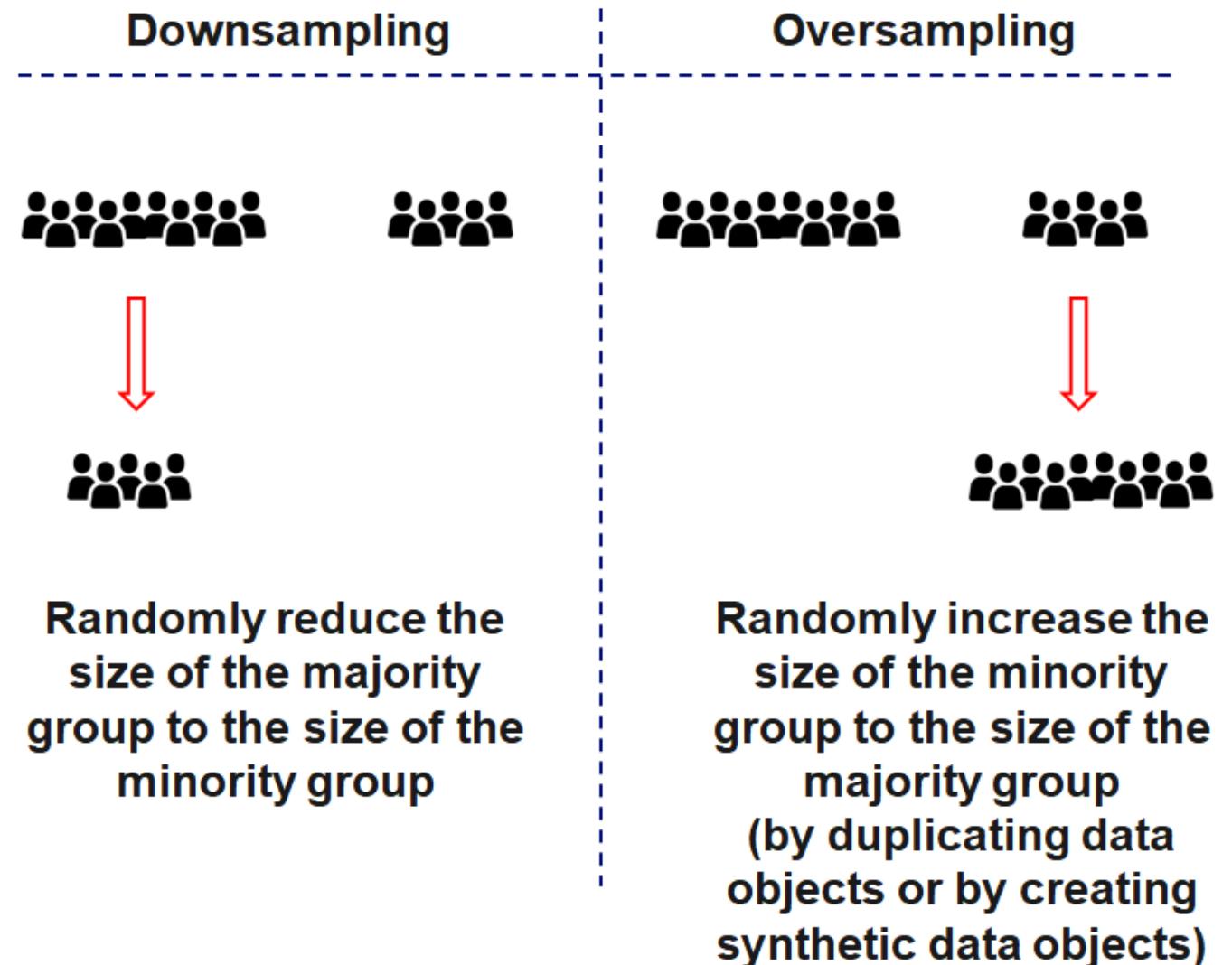
# Reasons for Sampling

- The volume of data is too large to capture and process
- Design the analytics process using a subset of the data for performance reasons. Later use the complete data set.
- The data set doesn't perfectly represent the target population.
- The data set is imbalanced.

# Popular Methods of Sampling (I)



# Popular Methods of Sampling (II)



# Systematic Sampling

Choose every  
row with an  
even number

day.mins	day.calls	day.charge	intl.mins	intl.calls	intl.charge	churn
157	79	26.69	7.1	6	1.92	No
184.5	97	31.37	8.7	4	2.35	No
258.6	84	43.96	11.2	5	3.02	No
129.1	137	21.95	12.7	6	3.43	Yes
187.7	127	31.91	9.1	5	2.46	No
128.8	96	21.9	11.2	2	3.02	No
156.6	88	26.62	12.3	5	3.32	No
120.7	70	20.52	13.1	6	3.54	No
332.9	67	56.59	5.4	9	1.46	Yes
196.4	139	33.39	13.8	4	3.73	No
190.7	114	32.42	8.1	3	2.19	No
189.7	66	32.25	10	5	2.7	No
224.4	90	38.15	13	2	3.51	No
155.1	117	26.37	10.6	4	2.86	No
62.4	89	10.61	5.7	6	1.54	Yes
84.8	95	14.42	14.2	6	3.83	No
226.1	105	38.44	10.3	5	2.78	No
212	121	36.04	12.6	10	3.4	No
249.6	118	42.43	11.8	3	3.19	Yes
176.8	94	30.06	8.3	4	2.24	No
140.4	94	23.87	11.1	9	3	No
126.3	102	21.47	9.4	2	2.54	No
173.1	85	29.43	14.6	15	3.94	Yes
124.8	82	21.22	10	4	2.7	No
85.8	77	14.59	9.2	4	2.48	No

Yes = 1

No = 11

# Random Sampling

Creates totally  
random samples

Example: 40%

day.mins	day.calls	day.charge	intl.mins	intl.calls	intl.charge	churn
157	79	26.69	7.1	6	1.92	No
184.5	97	31.37	8.7	4	2.35	No
258.6	84	43.96	11.2	5	3.02	No
129.1	137	21.95	12.7	6	3.43	Yes
187.7	127	31.91	9.1	5	2.46	No
128.8	96	21.9	11.2	2	3.02	No
156.6	88	26.62	12.3	5	3.32	No
120.7	70	20.52	13.1	6	3.54	No
332.9	67	56.59	5.4	9	1.46	Yes
196.4	139	33.39	13.8	4	3.73	No
190.7	114	32.42	8.1	3	2.19	No
189.7	66	32.25	10	5	2.7	No
224.4	90	38.15	13	2	3.51	No
155.1	117	26.37	10.6	4	2.86	No
62.4	89	10.61	5.7	6	1.54	Yes
84.8	95	14.42	14.2	6	3.83	No
226.1	105	38.44	10.3	5	2.78	No
212	121	36.04	12.6	10	3.4	No
249.6	118	42.43	11.8	3	3.19	Yes
176.8	94	30.06	8.3	4	2.24	No
140.4	94	23.87	11.1	9	3	No
126.3	102	21.47	9.4	2	2.54	No
173.1	85	29.43	14.6	15	3.94	Yes
124.8	82	21.22	10	4	2.7	No
85.8	77	14.59	9.2	4	2.48	No

Yes = 1  
No = 9

# Proportional Sampling

Creates  
proportional  
samples

Example: 40%

day.mins	day.calls	day.charge	intl.mins	intl.calls	intl.charge	churn
157	79	26.69	7.1	6	1.92	No
184.5	97	31.37	8.7	4	2.35	No
258.6	84	43.96	11.2	5	3.02	No
129.1	137	21.95	12.7	6	3.43	Yes
187.7	127	31.91	9.1	5	2.46	No
128.8	96	21.9	11.2	2	3.02	No
156.6	88	26.62	12.3	5	3.32	No
120.7	70	20.52	13.1	6	3.54	No
332.9	67	56.59	5.4	9	1.46	Yes
196.4	139	33.39	13.8	4	3.73	No
190.7	114	32.42	8.1	3	2.19	No
189.7	66	32.25	10	5	2.7	No
224.4	90	38.15	13	2	3.51	No
155.1	117	26.37	10.6	4	2.86	No
62.4	89	10.61	5.7	6	1.54	Yes
84.8	95	14.42	14.2	6	3.83	No
226.1	105	38.44	10.3	5	2.78	No
212	121	36.04	12.6	10	3.4	No
249.6	118	42.43	11.8	3	3.19	Yes
176.8	94	30.06	8.3	4	2.24	No
140.4	94	23.87	11.1	9	3	No
126.3	102	21.47	9.4	2	2.54	No
173.1	85	29.43	14.6	15	3.94	Yes
124.8	82	21.22	10	4	2.7	No
85.8	77	14.59	9.2	4	2.48	No

Yes = 2

No = 8

# Downsampling

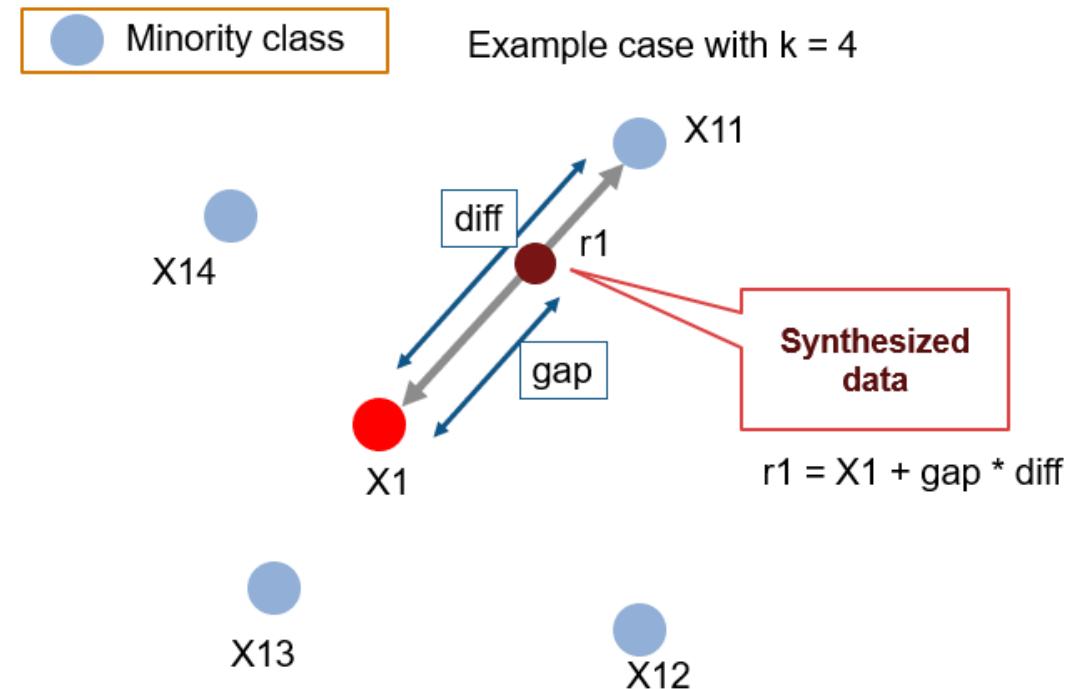
Creates  
balanced  
samples

day.mins	day.calls	day.charge	intl.mins	intl.calls	intl.charge	churn
157	79	26.69	7.1	6	1.92	No
184.5	97	31.37	8.7	4	2.35	No
258.6	84	43.96	11.2	5	3.02	No
129.1	137	21.95	12.7	6	3.43	Yes
187.7	127	31.91	9.1	5	2.46	No
128.8	96	21.9	11.2	2	3.02	No
156.6	88	26.62	12.3	5	3.32	No
120.7	70	20.52	13.1	6	3.54	No
332.9	67	56.59	5.4	9	1.46	Yes
196.4	139	33.39	13.8	4	3.73	No
190.7	114	32.42	8.1	3	2.19	No
189.7	66	32.25	10	5	2.7	No
224.4	90	38.15	13	2	3.51	No
155.1	117	26.37	10.6	4	2.86	No
62.4	89	10.61	5.7	6	1.54	Yes
84.8	95	14.42	14.2	6	3.83	No
226.1	105	38.44	10.3	5	2.78	No
212	121	36.04	12.6	10	3.4	No
249.6	118	42.43	11.8	3	3.19	Yes
176.8	94	30.06	8.3	4	2.24	No
140.4	94	23.87	11.1	9	3	No
126.3	102	21.47	9.4	2	2.54	No
173.1	85	29.43	14.6	15	3.94	Yes
124.8	82	21.22	10	4	2.7	No
85.8	77	14.59	9.2	4	2.48	No

Yes = 5  
No = 5

**SMOTE (Synthetic Minority Oversampling Technique) is an oversampling technique where the synthetic samples are generated for the minority class.**

At first the total number of oversampling observations N is set up. Usually, it is selected such that the resulting class distribution is 1:1. Now the iteration starts by



# Feature Engineering

Feature engineering is the process of using domain knowledge of the data to create features that make machine learning algorithms work. If feature engineering is done correctly, it increases the predictive power of machine learning algorithms by creating features from raw data that help facilitate the machine learning process.

A feature (variable, attribute) is depicted by a column in a dataset. Considering a generic two-dimensional dataset, each observation is depicted by a row and each feature by a column, which will have a specific value for an observation:

## 1. Transformation

- convert features (e.g., birth date → age)
- build lag structures (e.g., time-lags)
- normalization / standardization / scaling

## 2. Type Conversion

- if numerical type is needed, transform categorical into numerical data using dummy features
- if categorical type is needed or more informative, discretize numerical features (e.g., income → poor / rich classes)

## 3. Feature Combination

# Scaling

Most datasets contain features highly varying in magnitudes, units and range.

Most machine learning algorithms have problems with this because they use distance measures or calculate gradients. The features with high magnitudes will weigh in a lot more in the distance calculations than features with low magnitudes and gradients may end up taking a long time or are not accurately calculable.

To overcome this effect, we scale the features to bring them to the same level of magnitudes. The two most discussed scaling methods are Normalization and Standardization.

# Type Conversion (Encoding)

Many machine learning algorithms cannot work with categorical data directly. To convert categorical data to numbers, there exist two variants:

**Label encoding** refers to transforming the word labels into numerical form so that the algorithms can understand how to operate on them. Every categorical value is assigned to one numerical value, e.g. young -> 1, middle\_age -> 2, old -> 3. This only works in specific situations where you have somewhat continuous-like data, e.g. if the categorical feature is ordinal.

**One hot encoding** is a representation of a categorical variable as binary vectors. Every categorical value is assigned to an artificial

# Example of Feature Engineering (I)

Data sets often contain date/time features. These features are rarely useful in their original form because they only contain ongoing values. However, they can be useful for extracting cyclical factors, such as weekly or seasonal effects. Suppose, we are given a data “flight date time vs status”. Then, given the date-time data, we have to predict the status of the flight.

	Date_Time_Combined	Status		Hour_Of_Day	Status
0	2018-02-14 20:40	Delayed		0	20 Delayed
1	2018-02-15 10:30	On Time	→	1	10 On Time
2	2018-02-14 07:40	On Time		2	7 On Time
3	2018-02-15 18:10	Delayed		3	18 Delayed
4	2018-02-14 10:20	On Time		4	10 On Time

# Example of Feature Engineering (II)

Suppose we are given the latitude, longitude and other data with the objective to predict the target feature “ Price\_Of\_House “. Latitude and longitude are not of use in this context if they are alone. So, we will combine the latitude and the longitude to make one feature.

In other cases, it might be appropriate to transform latitude and longitude into categories which reflect regions, for example



# Example of Feature Engineering (III)

Suppose we are given a feature “`Marital_Status`” and other data with the objective to classify customers into “`Creditworthy`” and “`Not_Creditworthy`”. In the data set the marital status has many different values, for example

- single living alone
- single living with his parents
- married living together
- married living separately
- divorced
- divorced but living together
- registered partnerships
- living in marriage-like community

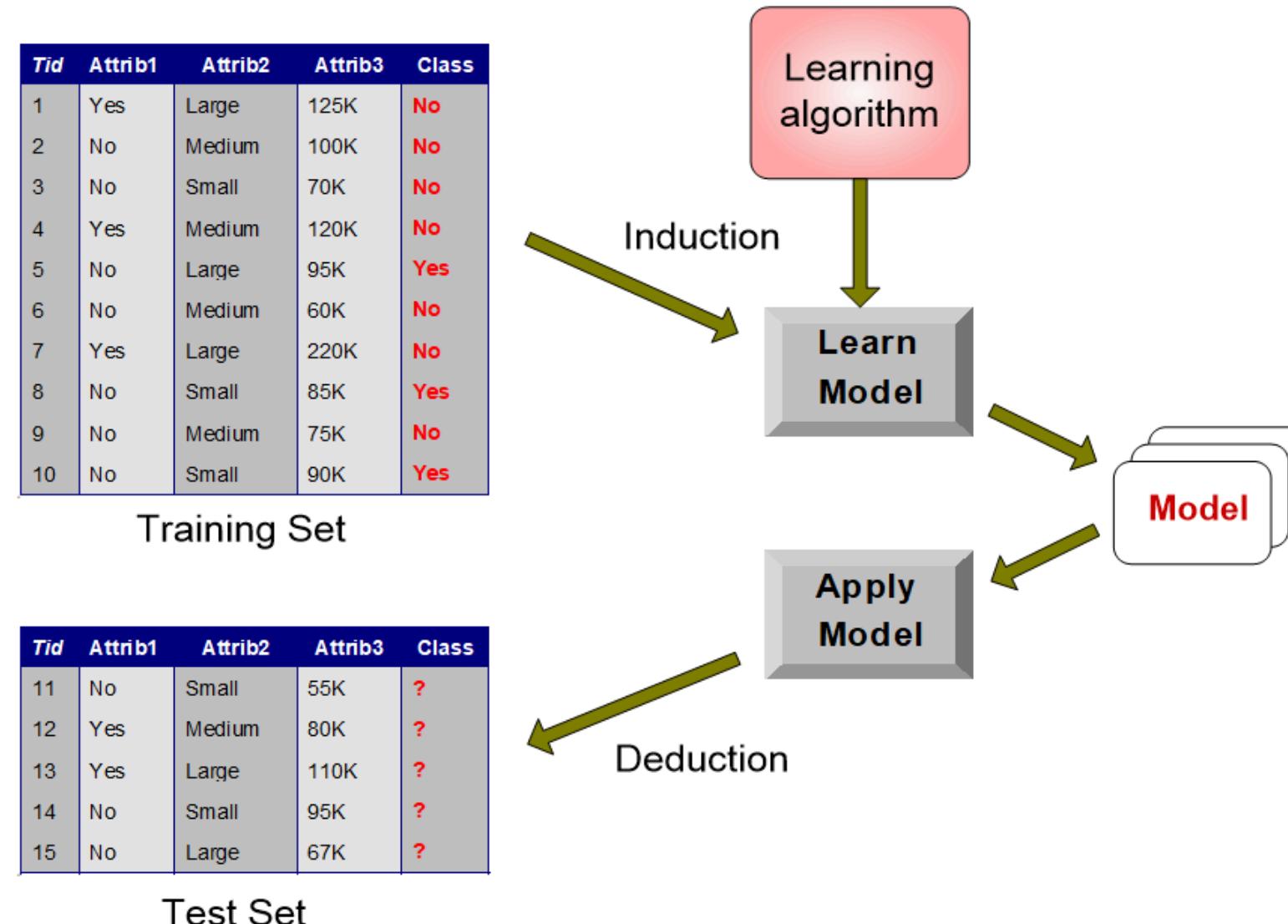
# Partitioning the Data

The partitioning of the data in **Training and Test Data** has the aim to proof if the analytical results can be generalized. The analysis (e.g. the development of a classifier) is carried out on the basis of training data. Subsequently, the results are applied to the test data. If the results are significantly worse than the training data, the model is not generalizable, which is called overfitting.



The partitioning of the data in training and test data can be

# Applying Training and Test Data



Source: <http://www.cs.kent.edu/~jin/BigData/Lecture10-ML-Classification.pptx>

# Partitioning

**Partitions into  
training and test  
data**

**Example:**  
**Training = 60%**  
**Test = 40%**

day.mins	day.calls	day.charge	intl.mins	intl.calls	intl.charge	churn
258.6	84	43.96	11.2	5	3.02	No
129.1	137	21.95	12.7	6	3.43	Yes
156.6	88	26.62	12.3	5	3.32	No
332.9	67	56.59	5.4	9	1.46	Yes
190.7	114	32.42	8.1	3	2.19	No
62.4	89	10.61	5.7	6	1.54	Yes
226.1	105	38.44	10.3	5	2.78	No
249.6	118	42.43	11.8	3	3.19	Yes
126.3	102	21.47	9.4	2	2.54	No
173.1	85	29.43	14.6	15	3.94	Yes



# Exploratory Data Analysis

In Exploratory Data Analysis (EDA), there is no hypothesis and there is no model.

People are not very good at looking at a column of numbers or a whole data table and then determining important characteristics of the data. EDA techniques have been devised as an aid in this situation.

Reasons for EDA:

- gain intuition about the data
- make comparisons between distributions
- sanity checking (making sure the data is on the scale you

**Non-graphical exploratory data analysis is the first step when beginning to analyze the data. This preliminary data analysis step focuses on four points:**

- measures of central tendency, i.e. mean and median. The median, known as 50th percentile, is more resistant to outliers.
- measures of spread, i.e. variance, standard deviation, and interquartile range
- the shape of the distribution
- the existence of outliers

Outlier are data objects, which are clearly different from the others.

Usually, the detection of outliers is an unsupervised process, because they are not known before analyses.

In the case of **numerical attributes** the Interquartil Range can be used. Here, an outlier is defined if the attribute lies outside the interval

$$[Q_{0.25} - k \cdot (Q_{0.75} - Q_{0.25}), Q_{0.75} + k \cdot (Q_{0.75} - Q_{0.25})]$$

Usually, k has a value between 1.5 and 3. The bigger k, the more different the values must be to be classified as outliers

- Outlier have to be **eliminated** if they
  - 1. would bias the analysis, e.g. if 9 persons have an age between 20 and 30 and the 10th person is 80 years old.
  - 2. are erroneous data, e.g. as a result of input errors or a defect sensor.
- It is not always acceptable to drop an observation just because it is an outlier. They can be legitimate observations and are sometimes interesting ones. It's important to investigate the nature of the outlier before deciding.
- In those cases where you shouldn't drop the outlier one option

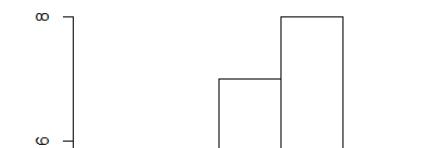
# Univariate Graphical EDA

Non-graphical and graphical EDA methods complement each other, they have the same focus. While the non-graphical methods are quantitative and objective, they do not give a full picture of the data. The distribution of a variable tells us what values the variable takes and how often each value occurs.

Types of displays:

for numerical variables: Histograms, Boxplots, Quantile-normal plots, ...

for categorical variables: Pie charts, Bar graphs, ...



# Multivariate Non-Graphical EDA

Multivariate non-graphical EDA techniques generally show the relationship between two or more variables in the form of either cross-tabulation for categorical variables or correlation statistics for numerical variables.

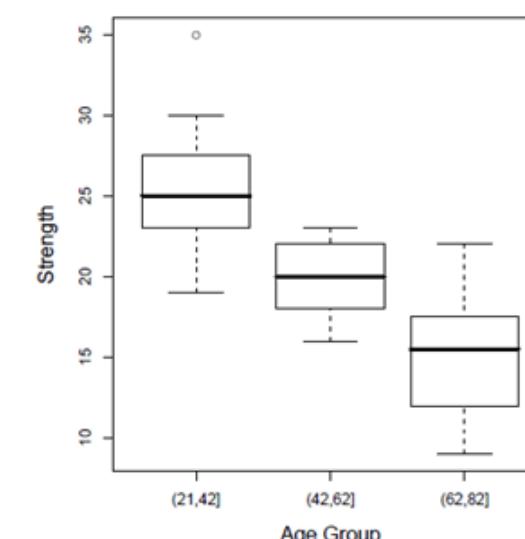
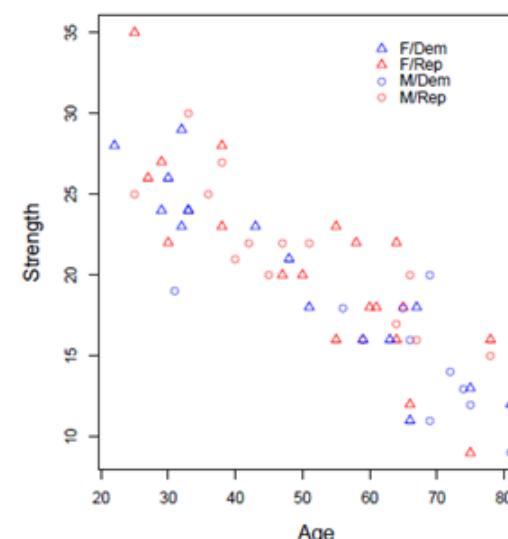
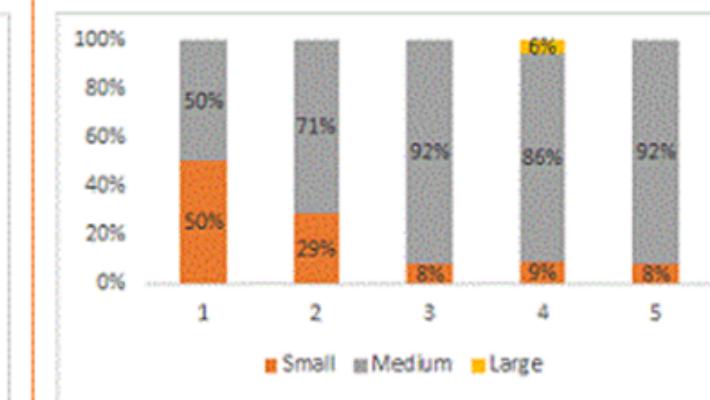
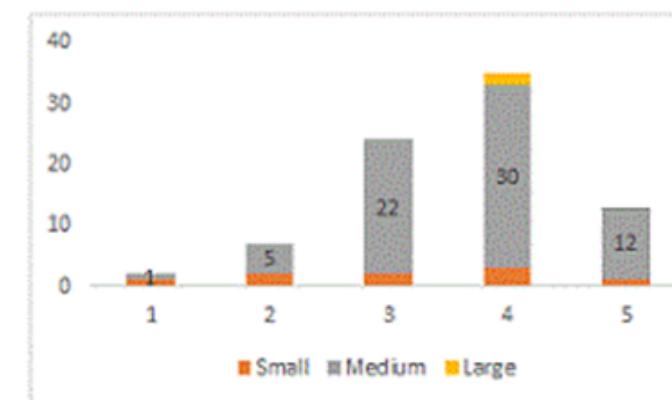
		Product Category				
Frequency	Row Pct	1	2	3	4	5
Small	11.11	22.22	22.22	33.33	11.11	11.11
Medium	1.43	7.14	31.43	42.86	12.14	12.14
Large	0.00	0.00	0.00	100.00	0.00	0.00
Total		2	7	24	35	13

Frequency Missing = 77

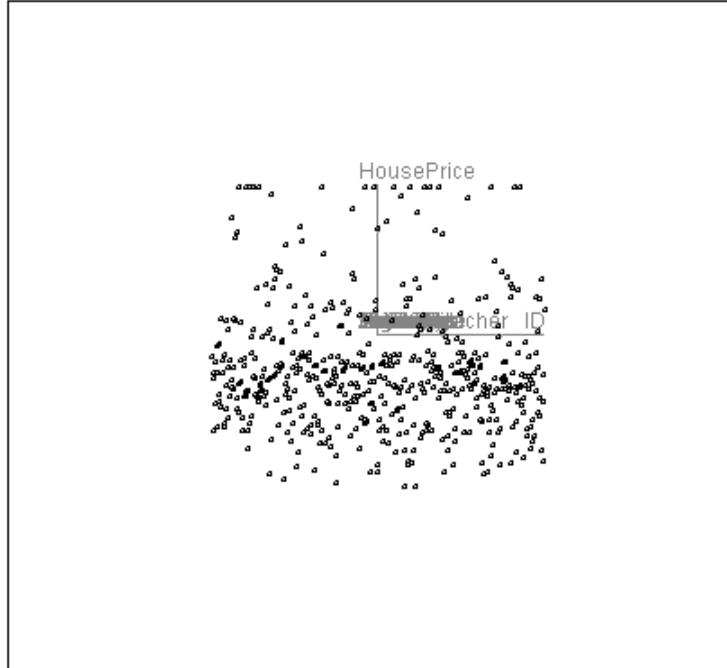
	V2	V3	V4	V5	V6
V2	-0.388				
V3	0.549	0.189			
V4	0.824	0.758	-0.619		
V5	0.802	0.833	-0.469	-0.946	
V6	-0.186	-0.659	0.349	0.555	0.649

# Multivariate Graphical EDA

Multivariate graphical EDA techniques are scatterplots for numerical variables, Barcharts for categorical variables, or Boxplots for mixed types.



# Touring Diagram



## 4 Predictive Analytics

4.1 Subject of Predictive Analytics

4.2 The Analytics Process

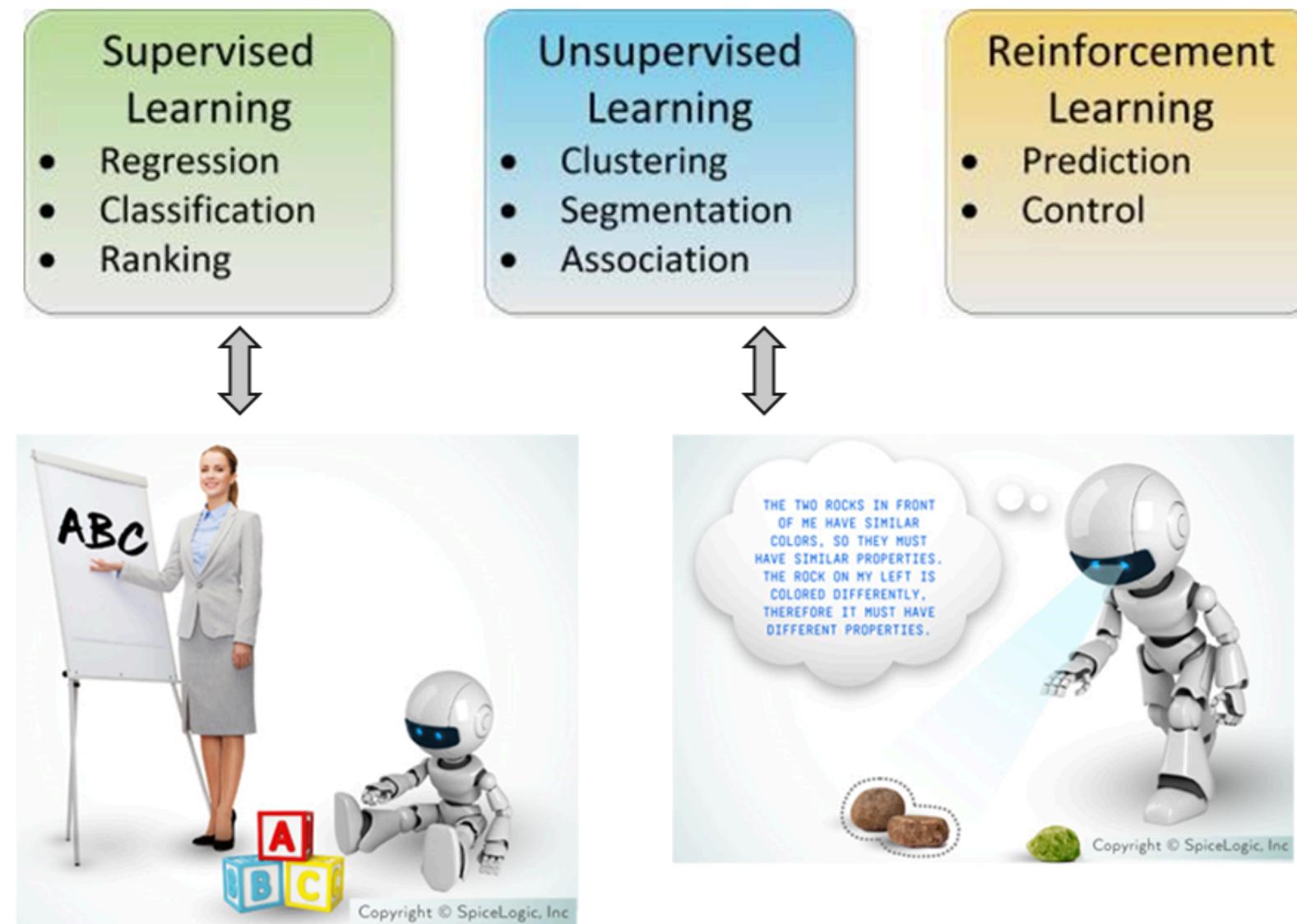
4.3 Data Preparation

4.4 Methods, Algorithms and Applications

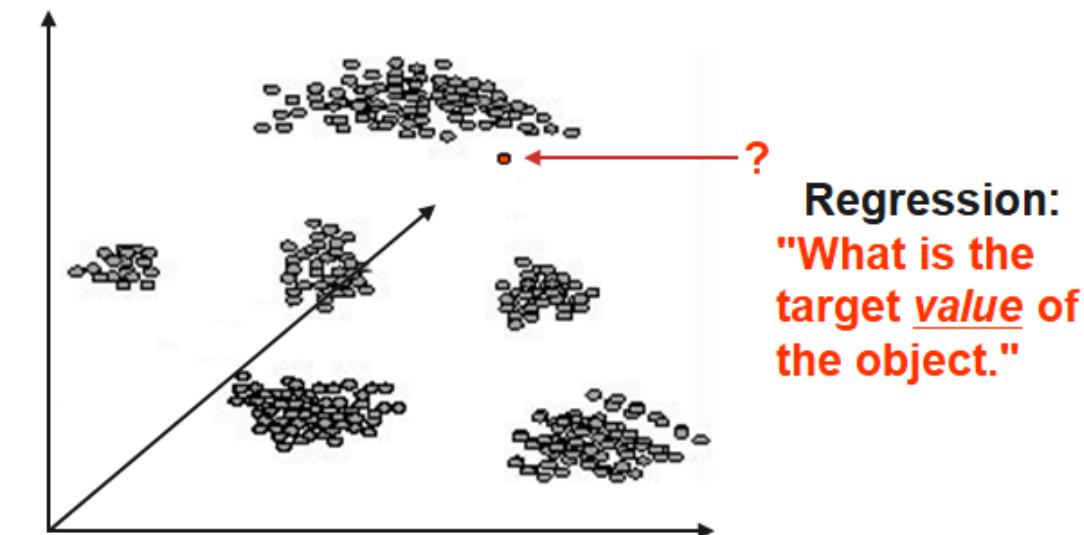
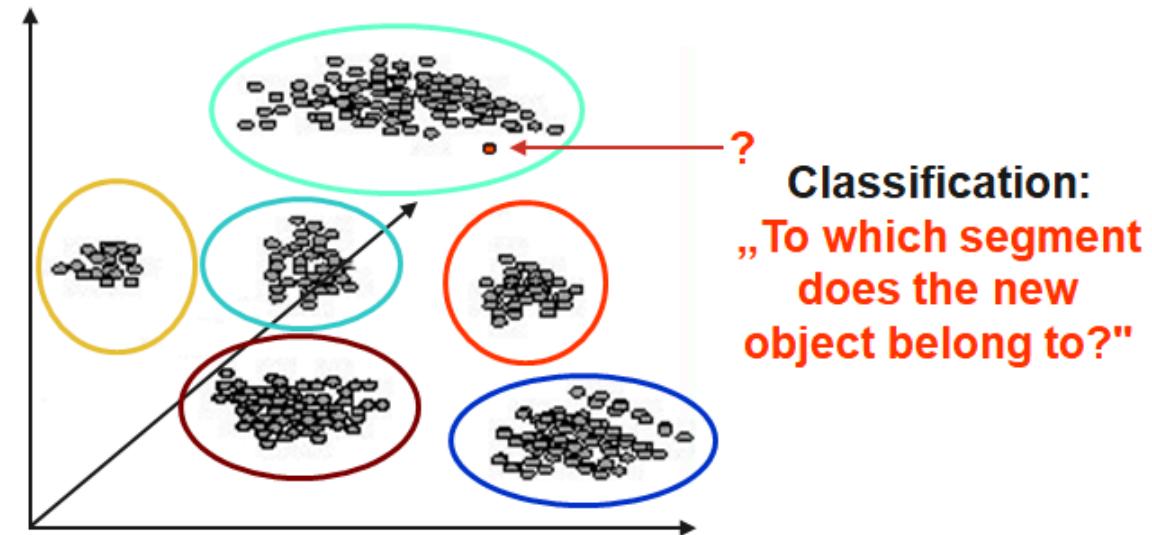
4.4.1 Classification

4.4.2 Regression

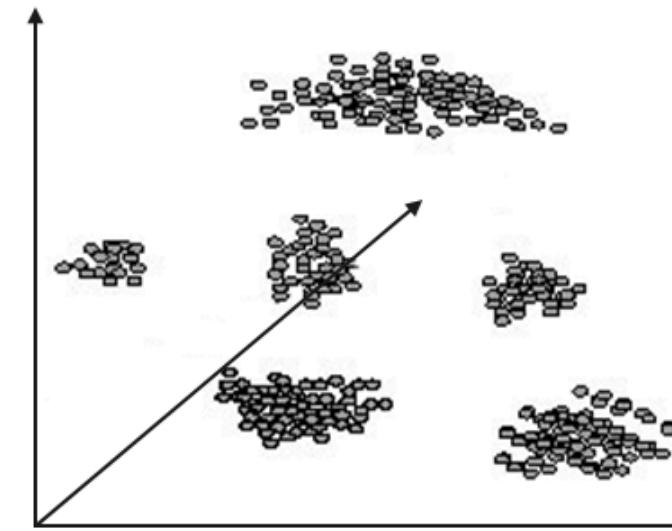
# Categories in Machine Learning



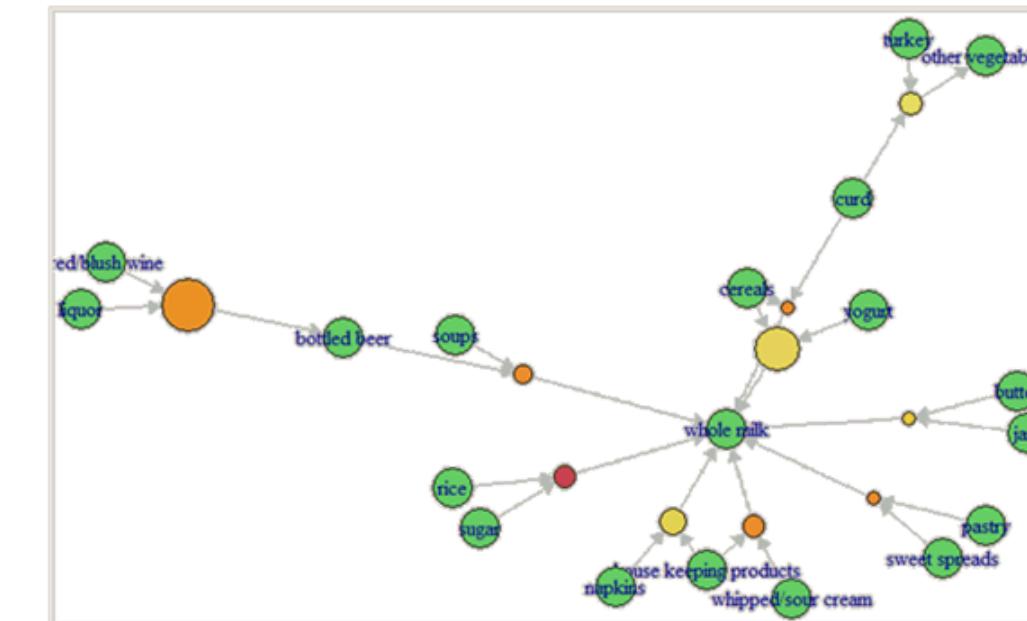
# Supervised Learning



# Unsupervised Learning



**Segmentation (Clustering):**  
**"In how many and which segments can the data be partitioned?"**



**Association:**  
**"Detect the strength and the direction of relationships in my data."**

# Supervised and Unsupervised Learning

## Supervised Learning:

- Goal: predict a **target or outcome** variable
- Training data where target value is known
- Score to test data where value is not known

DEP_TIME	DEST	DISTANCE	FL_DATE	ORIGIN	Weather	DAY_WEEK	DAY_OF_MONTH	FlightStatus
1455	JFK	184	01.01.2004	BWI	0	4	1	ontime
1640	JFK	213	01.01.2004	DCA	0	4	1	ontime
1245	LGA	229	01.01.2004	IAD	0	4	1	ontime
1709	LGA	229	01.01.2004	IAD	0	4	1	ontime
1035	LGA	229	01.01.2004	IAD	0	4	1	ontime
839	JFK	228	01.01.2004	IAD	0	4	1	ontime
1243	JFK	228	01.01.2004	IAD	0	4	1	ontime
1644	JFK	228	01.01.2004	IAD	0	4	1	ontime
1710	JFK	228	01.01.2004	IAD	0	4	1	ontime
2129	JFK	228	01.01.2004	IAD	0	4	1	ontime
2114	LGA	229	01.01.2004	IAD	0	4	1	ontime
1458	JFK	213	01.01.2004	DCA	0	4	1	ontime
932	LGA	214	01.01.2004	DCA	0	4	1	ontime
1228	LGA	214	01.01.2004	DCA	0	4	1	ontime

HousePrice	HouseSize	HouseAge	LotSize	RiverSide	CrimeRate	Industrial
244	6.064	59.1	0	1	.13587	10.59
212	6.176	72.5	0	0	.13158	10.01
222	6.358	52.9	30	0	.1029	4.93
134	6.103	85.1	0	0	705.042	18.1
238	5.877	79.2	0	0	180.028	19.58
174	5.594	36.8	12.5	0	.13554	6.07
222	6.316	38.1	0	0	.05083	5.19
262	6.718	17.5	22	0	.19073	5.86
140	6.174	93.6	0	0	.2909	21.89
190	5.985	45.4	0	0	.05497	5.19
178	7.393	99.3	0	0	824.809	18.1

## Unsupervised Learning:

- Goal: detect patterns
- Grouping of cases based on similarities in input values
- There is no target (outcome) variable

# Use Cases Quiz

## Online Shop for Electronic products:

- Management of Customer Termination (Churn Management)
- Returns Management
- Identify similar Regions for Advertising
- Sales Forecasts for Supply and Inventory Management
- Automatically assign new Products into existing Product Categories
- Grouping of similar Customers
- Voucher Management
- Analysis of frequently sold Product Combinations
- Fraud Detection
- Structuring the Competitors
- Analysis of Sales Sequences
- ...

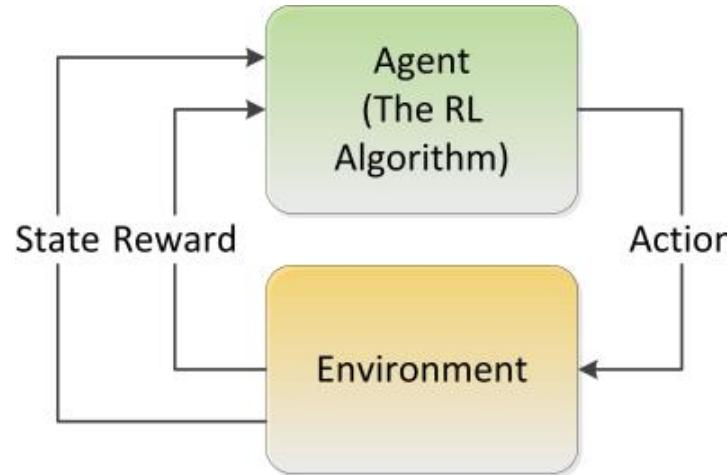
# Reinforcement Learning

The solution to many of the problems in our lives cannot be automated. This is not because current computers are too slow, but simply because it is too difficult for humans to determine what the program should do.

Supervised learning is a general method for training an approximator. However, supervised learning requires sample input-output pairs from the domain to be learned.

For example, we might not know the best way to program a computer to recognize an infrared picture of a tank, but we do have a large collection of infrared pictures, and we do know

# Reinforcement Learning



**The agent learns how to achieve a given goal by trial-and-error interactions with its environment by maximizing a reward.**

Go is one of the hardest games in the world for AI because of the huge number of different game scenarios and moves. The number of potential legal board positions is greater than the number of atoms in the universe.

The core of AlphaGo is a deep neural network. It was initially trained to learn playing by using a database of around 30 million recorded historical moves.



An artificial intelligence called Libratus has beaten four of the world's best poker players in a grueling 20-day tournament in January 2017.

Poker is more difficult because it's a game with imperfect information. With chess and Go, each player can see the entire board, but with poker, players don't get to see each other's hands. Furthermore, the AI is required to bluff and correctly interpret misleading information in order to win.

“We didn’t tell Libratus how to play poker. We gave it the rules of poker and said ‘learn on your own.’” The AI started playing randomly but over the course of playing trillions of hands was able

# Types of Artificial Intelligence

**Discriminative AI** is designed to differentiate and classify input, but not to create new content. Examples include image or speech recognition, credit scoring or stock price prediction.

**Generative AI** is able to generate new content based on existing information and user specifications. This includes texts, images, videos, program code, etc. The generated content can often hardly be distinguished from human-generated content. As things stand at

**ChatGPT is a generative AI that produces human-like text and communicates with humans.**

The “GPT” in ChatGPT comes from the language model of the same name, which was extended for ChatGPT with various components for communication and quality assurance.

GPT is based on a huge neural network that essentially represents the language model. While the first GPT-3 has 175 billion parameters, the

# ChatGPT - Approach

**ChatGPT generates its response word by word via a sequence of probabilities, with each new word depending on the previous ones.**

*The best thing about AI is its ability to*

learn	4.5%
predict	3.5%
make	3.2%
understand	3.1%
do	2.9%



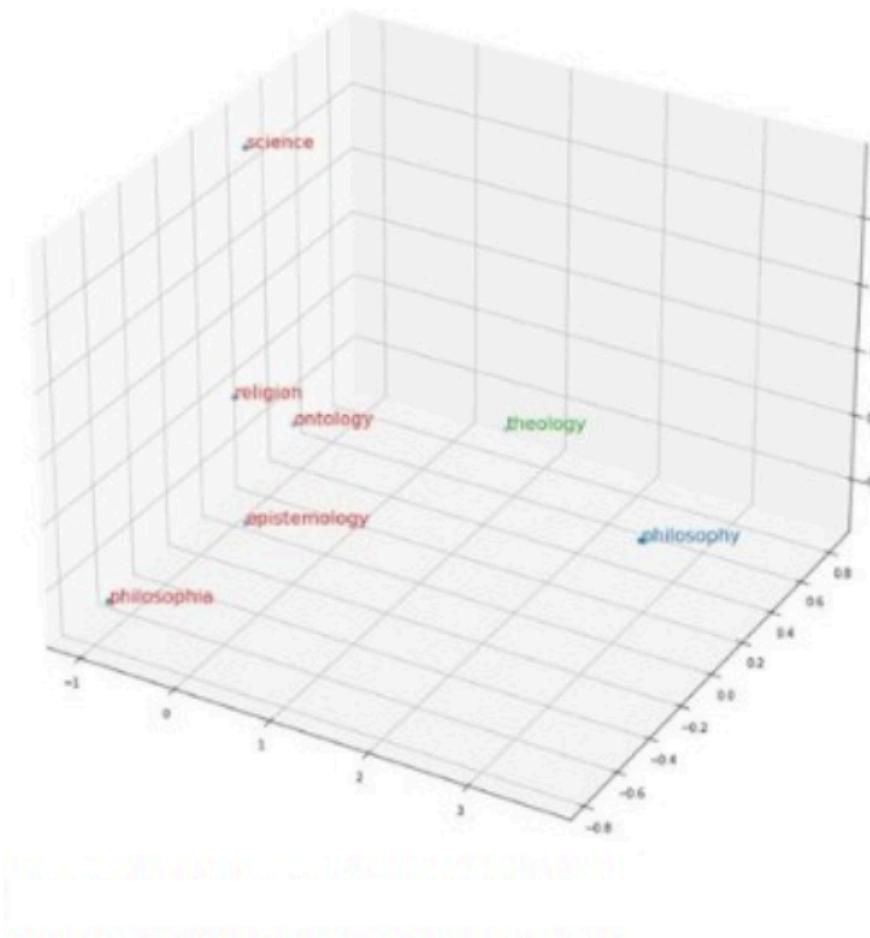
The best thing about AI is its ability to,  
The best thing about AI is its ability to learn,  
The best thing about AI is its ability to learn from,  
The best thing about AI is its ability to learn from experience,  
The best thing about AI is its ability to learn from experience.,  
The best thing about AI is its ability to learn from experience. It,  
The best thing about AI is its ability to learn from experience. It's,  
The best thing about AI is its ability to learn from experience. It's not

**The most probable word is not always selected; instead, randomization takes place. This means that different variants can be created for the same task.**

# ChatGPT - Semantic Spaces (I)



# ChatGPT - Semantic Spaces (II)

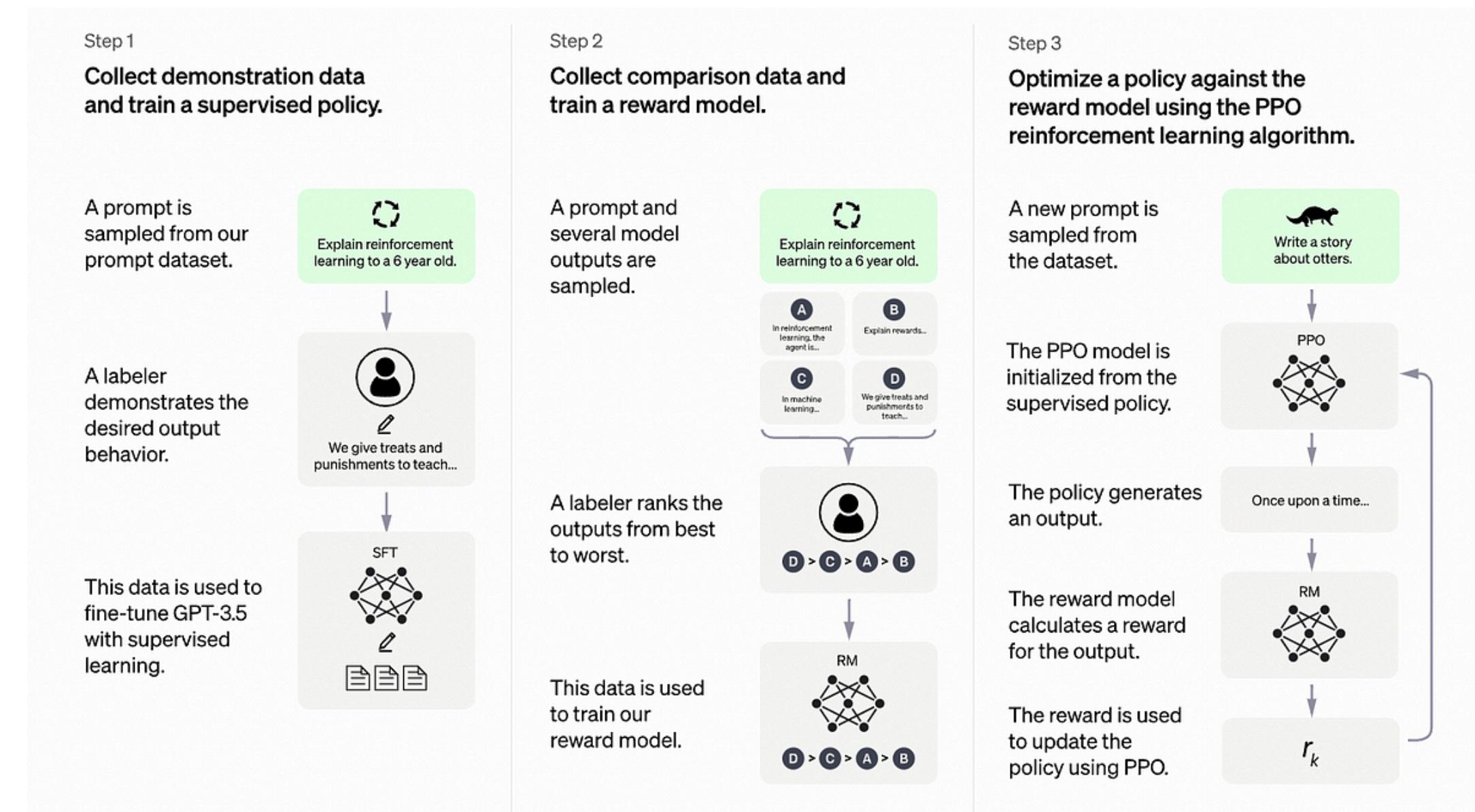


```
1 print_neighbors('ontology', model5)
```

## Approximate Neighbors

```
('Eros', 0.8966988027095795)
('Astronomy', 0.8926706984639168)
('Philosophy', 0.886092908680439)
('mathematics', 0.8845213651657104)
('Psychology', 0.8803791999816895)
('Night', 0.8794359490275383)
('subject-object', 0.8780221417546272)
('king', 0.8750540241599083)
('epistemology', 0.8732230812311172)
('Dionysus', 0.8667336702346802)
```

# ChatGPT - Evaluation Component



## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

##### 4.4.1.1 K-Nearest Neighbors

##### 4.4.1.2 Evaluating the Quality of Classification

##### 4.4.1.3 Decision Tree Approaches

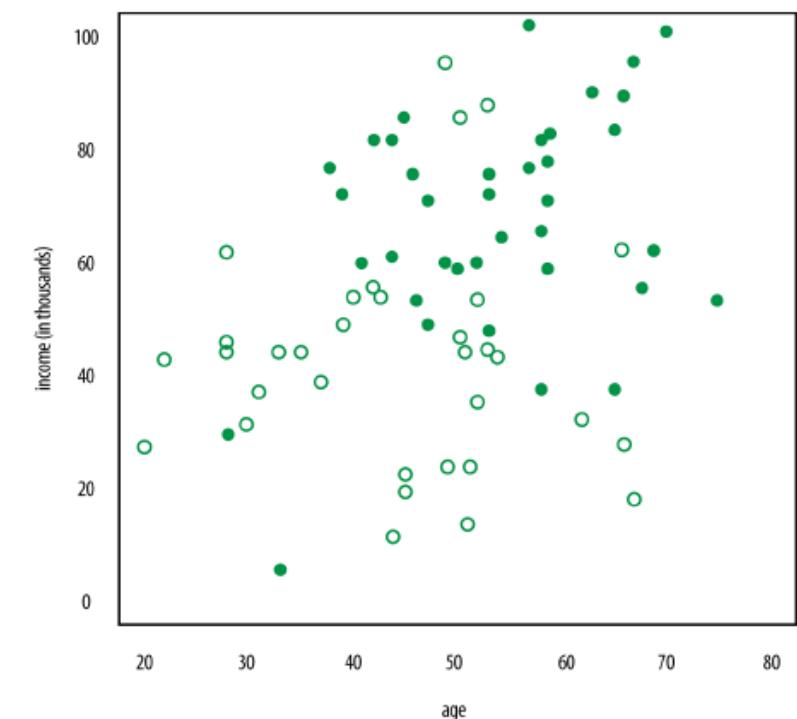
##### 4.4.1.4 Logistic Regression

# Introductory Example

Credit-Scoring is a typical example for a classification problem. A bank wants to determine the creditworthiness of a customer.

Assume you have the age, income, and a creditworthiness category of “yes” or “no” for a bunch of people and you want to use the age and income to predict the creditworthiness for a new person.

You can plot people as points on the plane and label people with an empty circle if they have low credit ratings.



# 4 Predictive Analytics

## 4.1 Subject of Predictive Analytics

## 4.2 The Analytics Process

## 4.3 Data Preparation

## 4.4 Methods, Algorithms and Applications

### 4.4.1 Classification

#### 4.4.1.1 K-Nearest Neighbors

#### 4.4.1.2 Evaluating the Quality of Classification

#### 4.4.1.3 Decision Tree Approaches

#### 4.4.1.4 Logistic Regression

#### 4.4.1.5 Neural Networks

#### 4.4.1.6 Resampling

# k-Nearest Neighbors

- k-Nearest Neighbors (k-NN) is an algorithm that can be used when you have a bunch of objects that have been classified or labeled in some way, and other similar objects that have not gotten classified or labeled yet, and you want a way to automatically label them.
- The intuition behind k-NN is to consider the most similar other items defined in terms of their attributes, look at their labels, and give the unassigned item the majority vote. If there's a tie, you randomly select

# Measuring Similarity

**The purpose of a measure of similarity is to compare two vectors of data, and compute a single number which evaluates their similarity.**

**Usually the similarity  $s(x, y)$  of two objects  $x$  and  $y$  is measured using the inverse of some distance measure  $d(x, y)$ :**

$$s(x, y) = 1 - d(x, y)$$

**Euclidean distance is most often used for numerical data vectors:**

$$d(x, y) = \|x - y\| = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

**with  $n$  as size of the vectors.**

**Note:** The variables must be scaled to avoid unwanted weighting resulting from different level of measurement.

# Unnormalized vs. Normalized

Unnormalized:

Objects \ Features	Age	Income	Rent	Years in Job
Objects				
Peter	47	4000	1400	22
Thomas	33	2700	1050	15

$$d = \sqrt{(47 - 33)^2 + (4000 - 2700)^2 + (1400 - 1050)^2 + (22 - 15)^2} = 1346.38$$

Normalized:

Objects \ Features	Age	Income	Rent	Years in Job
Objects				
Peter	0.32	0.21	0.85	0.53
Thomas	0.26	0.13	0.54	0.34

$$d = \sqrt{(0.32 - 0.26)^2 + (0.21 - 0.13)^2 + (0.85 - 0.54)^2 + (0.53 - 0.34)^2} = 0.3771$$

# Example (I)

**Data:**

Customer	Age	Monthly Income	Monthly Costs	Creditworthy
A	24	1800	1400	yes
B	26	1700	1300	no
C	30	2800	2200	yes
D	33	4200	3800	yes
E	35	2300	1800	no
F	39	3000	2800	no
G	41	2400	1600	yes
H	47	5000	2500	yes
I	53	4700	3200	no
J	59	2100	1700	yes

**New Customer:**

X	32	2900	1800	?
---	----	------	------	---

# Example (II)

1. **k = 5** is chosen

2. Calculate the distances (Normalization and then Euclidian Distance)

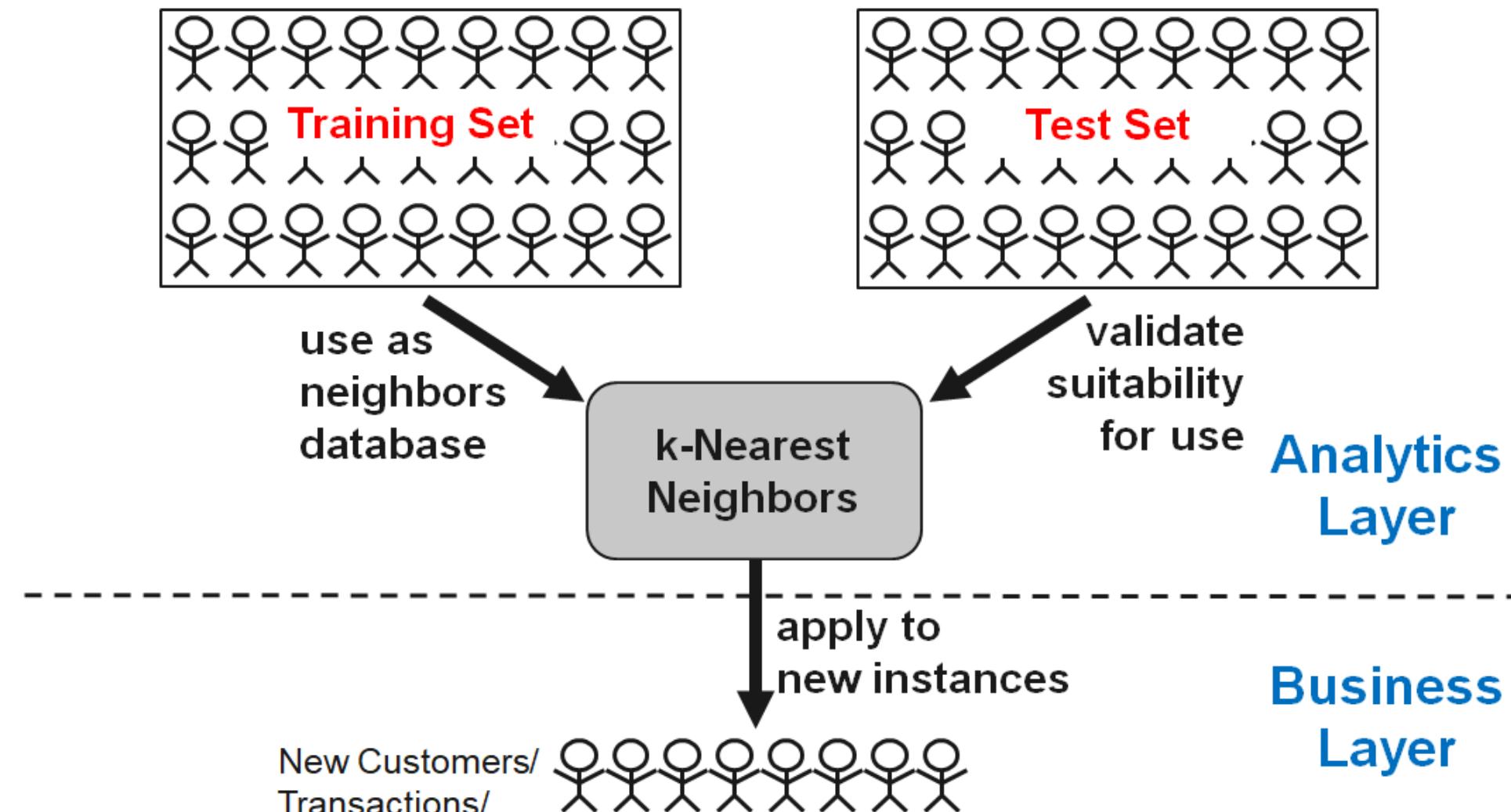
Customer	Age	Monthly Income	Monthly Costs	Creditworthy	Distance
A	0.0000	0.0303	0.0400	yes	0.4347
B	0.0571	0.0000	0.0000	no	0.4490
C	0.1714	0.3333	0.3600	yes	0.1726
D	0.2571	0.7576	1.0000	yes	0.8922
E	0.3143	0.1818	0.2000	no	0.2010
F	0.4286	0.3939	0.6000	no	0.4482
G	0.4857	0.2121	0.1200	yes	0.3090
H	0.6571	1.0000	0.4800	yes	0.8167
I	0.8286	0.9091	0.7600	no	0.9855
J	1.0000	0.1212	0.1600	yes	0.8096
X	0.2286	0.3636	0.2000	?	

## Example (III)

### 3. Choose the k nearest neighbors

Customer	Age	Monthly Income	Monthly Costs	Creditworthy	Distance
A	0.0000	0.0303	0.0400	yes	0.4347
C	0.1714	0.3333	0.3600	yes	0.1726
E	0.3143	0.1818	0.2000	no	0.2010
F	0.4286	0.3939	0.6000	no	0.4482
G	0.4857	0.2121	0.1200	yes	0.3090
X	0.2286	0.3636	0.2000	?	

# Creation and Use of Models



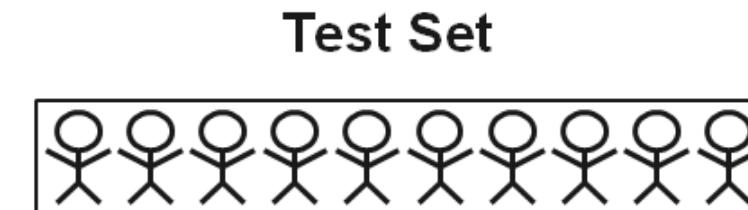
# Calculating Accuracies



+ + - + - - + - + -  
+ - - + - + + - + -

The training set is used to calculate the classifier model.

The accuracy of the training set indicates the goodness of fit of the model.



- + - - + + + - + - real

- + + - + - + - + + predicted

The test set is used to evaluate the classifier model.

The accuracy of the test set indicates the robustness/predictability of the model.

$$\text{Accuracy} = \frac{\#\text{Match}}{\#\text{All}}$$

# Determining Parameter k

1. Split the original labeled dataset into training and test data.
2. Pick an evaluation metric. Misclassification rate or accuracy are good ones.
3. Run k-NN a few times, changing k and checking the evaluation measure.
4. Optimize k by picking the one with the best evaluation measure.

## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

##### 4.4.1.1 K-Nearest Neighbors

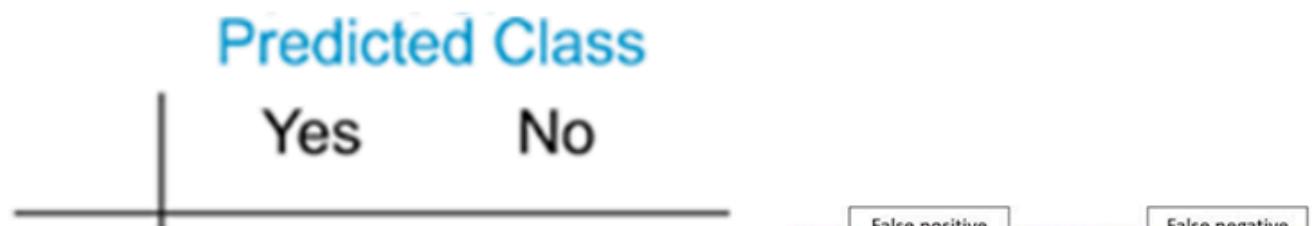
##### 4.4.1.2 Evaluating the Quality of Classification

##### 4.4.1.3 Decision Tree Approaches

##### 4.4.1.4 Logistic Regression

# Evaluating the quality of Classification (I)

True positives (TP), true negatives (TN), false positives (FP), and false negatives (FN), are the four different possible outcomes of a single prediction for a two-class case. A false positive is when the outcome is incorrectly classified as “yes”, when it is in fact “no”. A false negative is when the outcome is incorrectly classified as negative, when it is in fact positive. True positives and true negatives are obviously correct classifications.



# Evaluating the quality of Classification (II)

Test metrics are used to assess how accurately the model predicts the known values:

$$\text{accuracy} = \frac{TN + TP}{TP + FP + TN + FN}$$

$$\text{misclassification rate} = \frac{FN + FP}{TP + FP + TN + FN} = 1 - \text{accuracy}$$

Most classification algorithms pursue to minimize the misclassification rate. They implicitly assume that all misclassification errors cost equally. In many real-world applications, this assumption is not true. Cost-sensitive learning takes costs, such as the misclassification cost,

# Evaluating the quality of Classification (III)

Misclassification rate and accuracy can be misleading, for example in the case of imbalanced samples. Extreme case:

|     | yes  | no |
|-----|------|----|
| yes | 1000 | 3  |
| no  | 45   | 2  |

→  $\frac{1000 + 2}{1000 + 45 + 3 + 2} = 0.954$

For problems like, this additional measures are required to evaluate a classifier.

**Sensitivity** (true positive rate, recall) measures the proportion of positives that are correctly identified as such. **Specificity** (true negative rate) measures the proportion of negatives that are correctly identified as such.

# Problem of Imbalancing and Accuracy

Assume the following case: A credit card company wants to create a fraud detection system to include it into their transactional systems. The outcomes should be “Accept” (Y) and “Reject” (N). Because fraud rarely occurs, the data set consists of 320 observations for Y and 139 for N. They are partitioned into training and test set. Finally, the model is trained and tested.

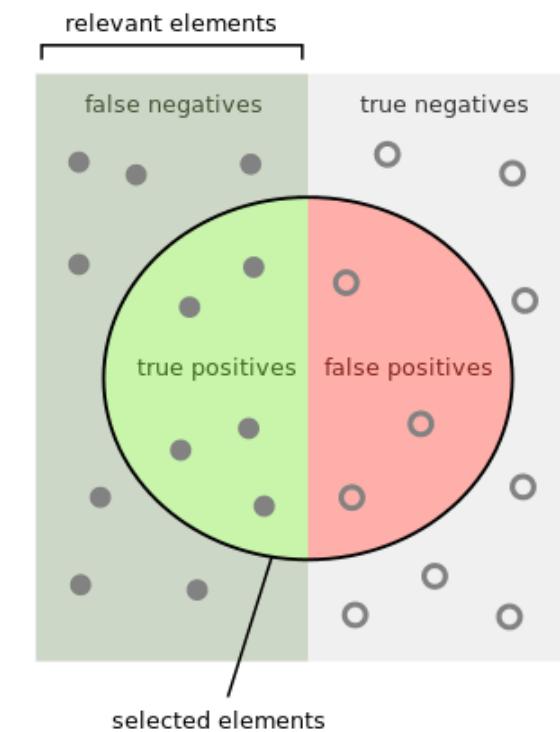
Because of the majority of the Y class, the training process concentrates on these cases because their correct classification promises the highest accuracy.

# Evaluating the quality of Classification (IV)

**Precision** measures the proportion of predicted positives who are true positives. A precision of 0.5 means that whenever the model classifies a positive, there is a 50% chance of it really being a positive. The higher the precision the smaller the number of false positives.

$$precision = \frac{TP}{TP + FP}$$

$$recall = \frac{TP}{TP + FN}$$



How many selected items are relevant?

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

**Recall** measures the percentage of positives the model is able to catch.

# Evaluating the quality of Classification (V)

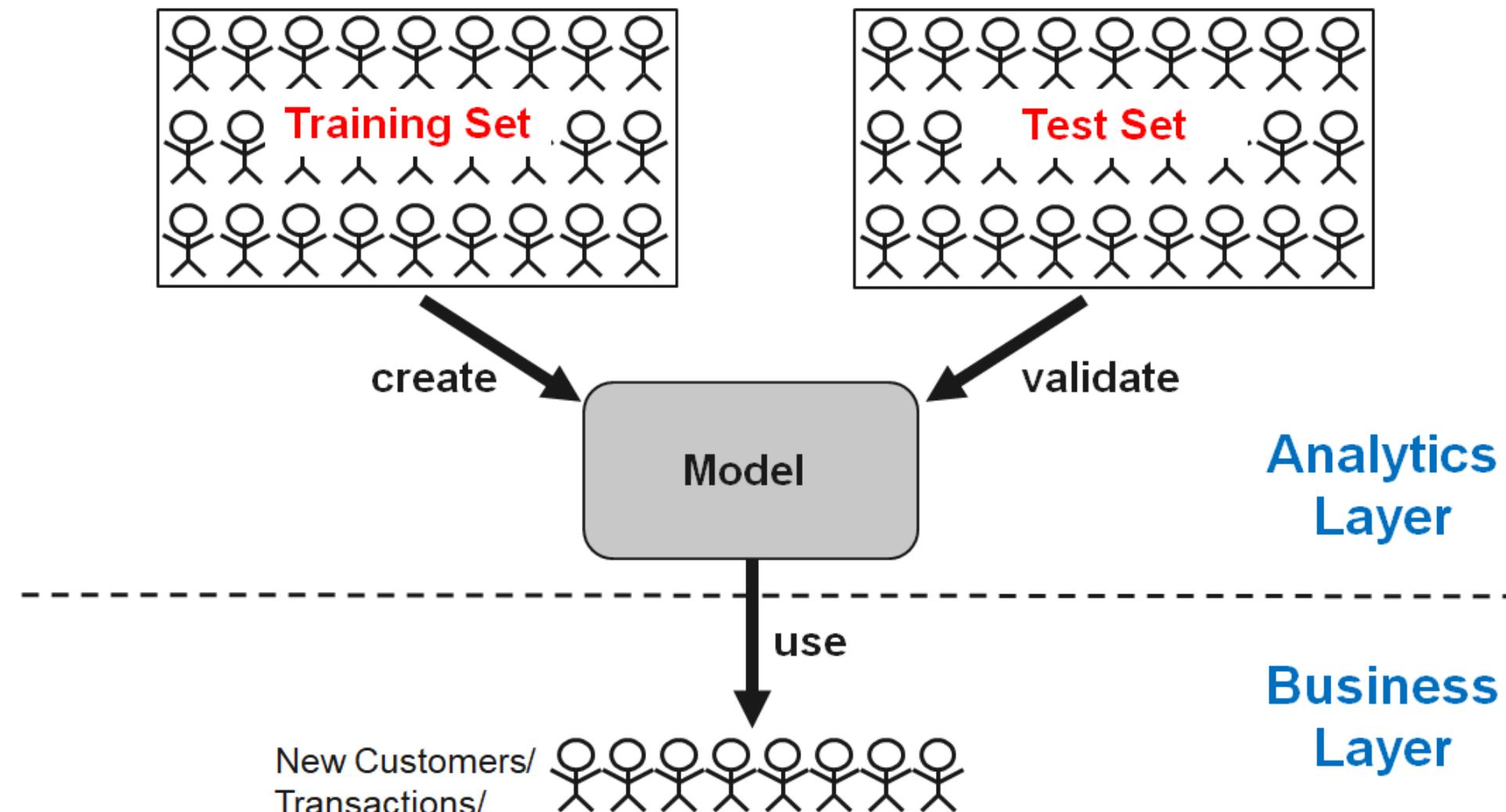
The F1 Score can be interpreted as the weighted average of both precision and recall. The main idea of the F1 Score is to strike a balance between both precision and recall and measure it in a single metric.

9      
$$F1 = 2 * \frac{precision * recall}{precision + recall}$$

A F1 score reaches its best value at 1 (perfect precision and recall) and worst at 0.

It is commonly used in cases of high class imbalance.

# Creation and Use of Models



## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

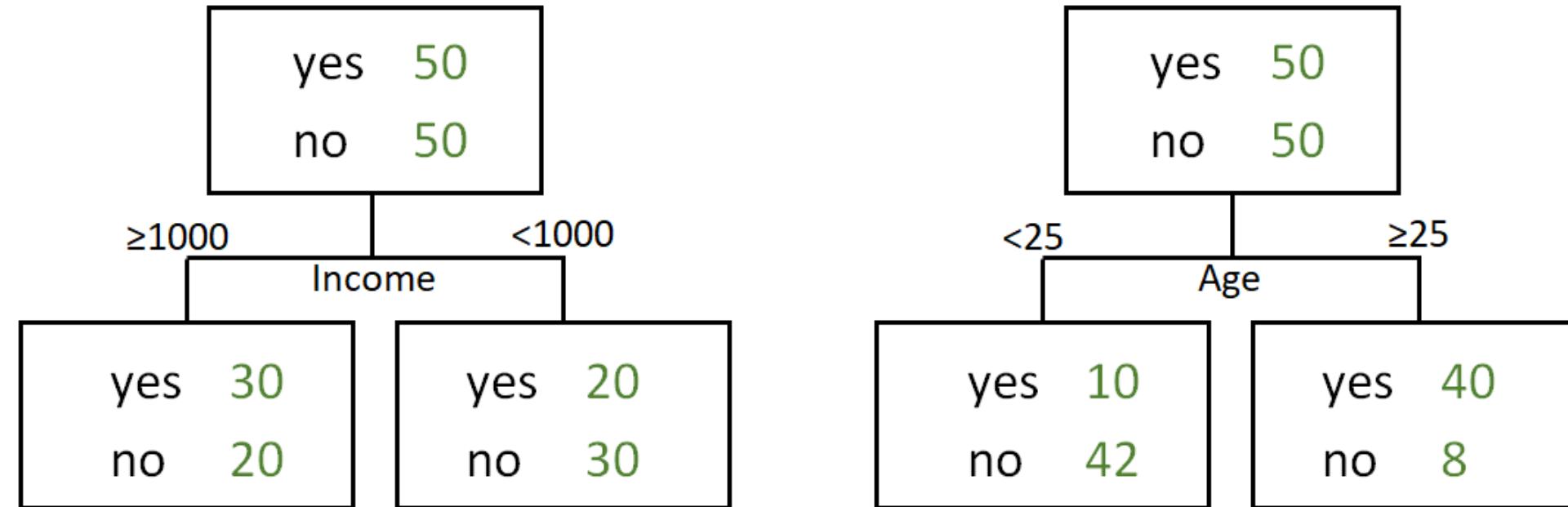
##### 4.4.1.1 K-Nearest Neighbors

##### 4.4.1.2 Evaluating the Quality of Classification

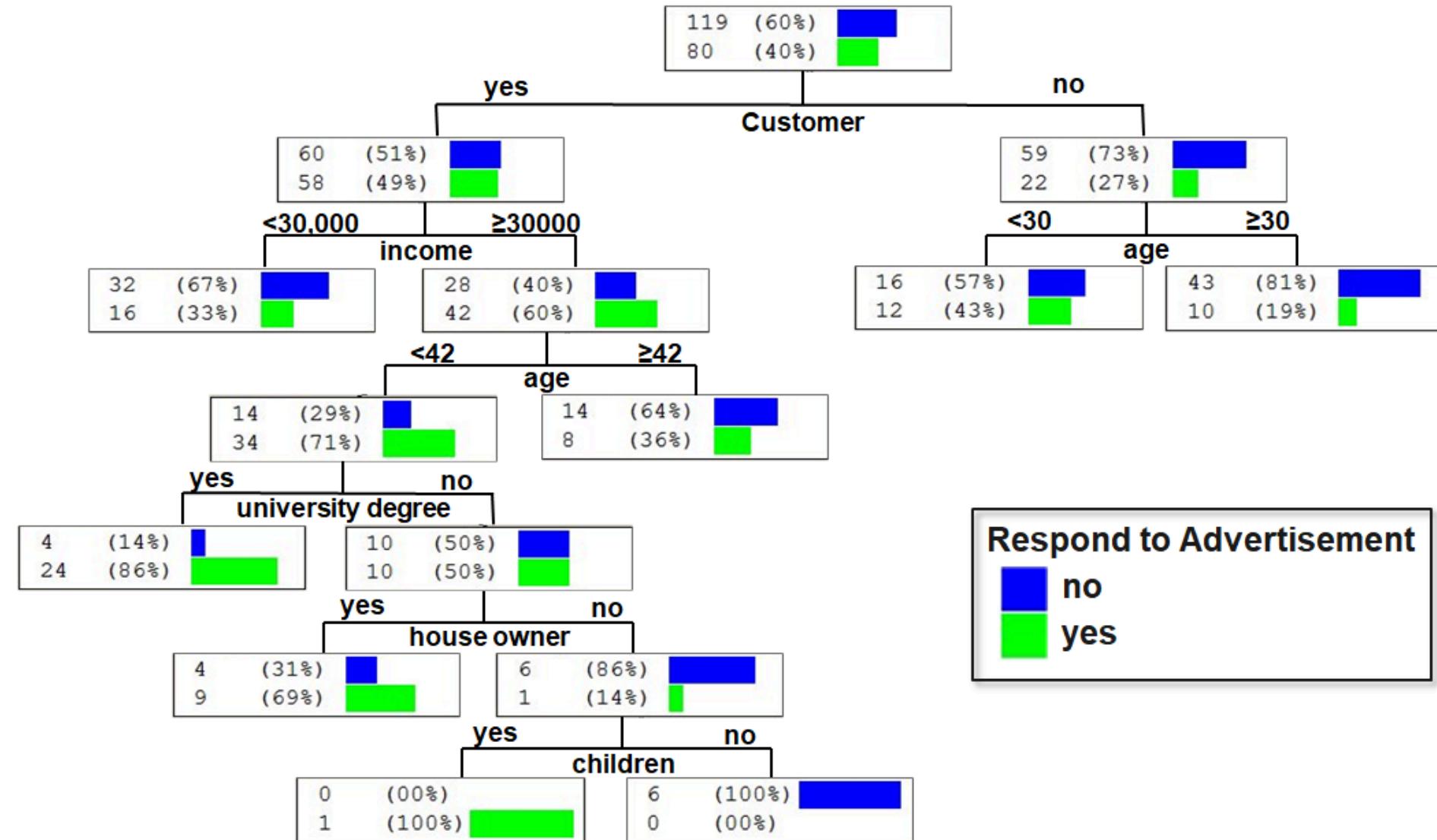
##### 4.4.1.3 Decision Tree Approaches

##### 4.4.1.4 Logistic Regression

# Which one is better?



# Introductory Example



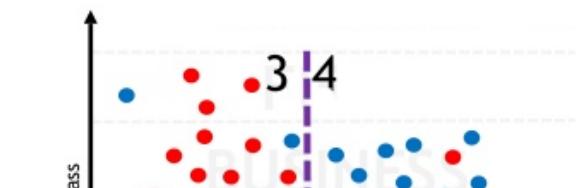
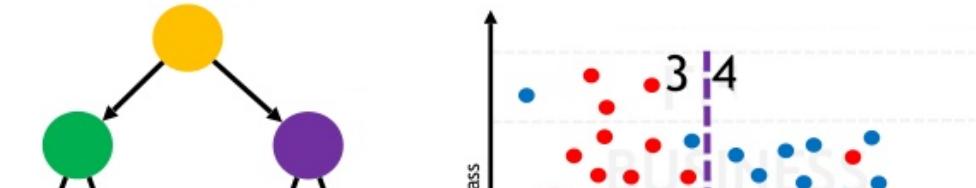
Decision trees belong to the hierarchical methods of classification. They analyze step-by-step (recursive partitioning).

A decision tree consists of nodes and borders. The topmost node (without any parent node) is called “root”. A node without a child node is called “leaf”. Nodes that have parent and child nodes are called “interior nodes”. The interior nodes represent the splitting of the included object sets. An interior node has at least two child nodes (sons). If every interior node has exactly two child nodes

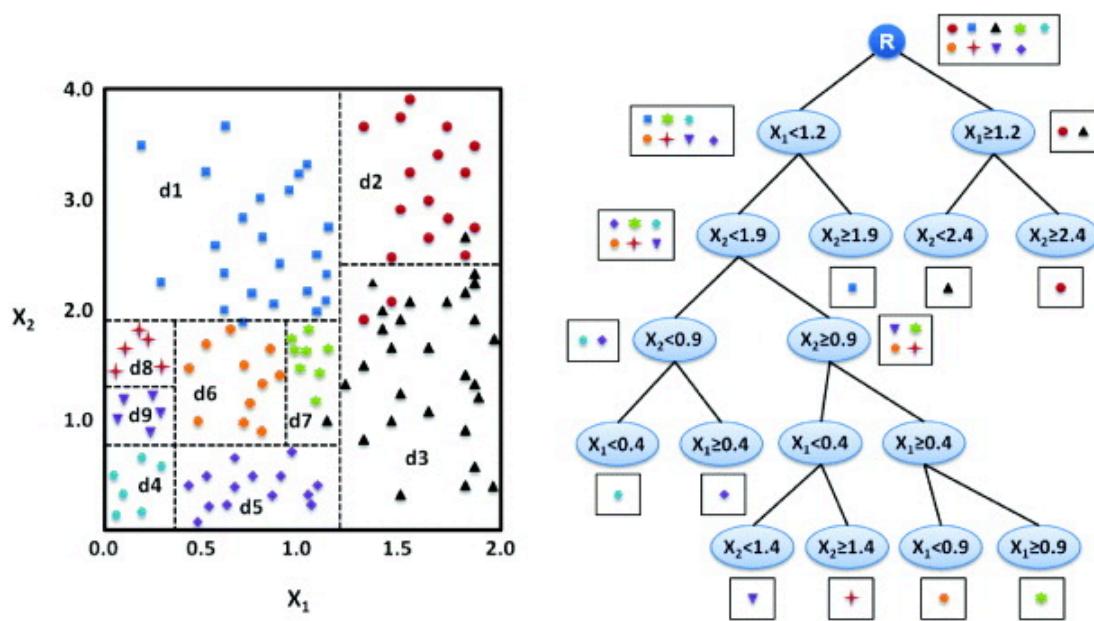
# Decision Trees (II)

Graphically, decision tree models divide the dataspace in a large number of subspaces and search for the variables which are able to split the dataspace with the greatest homogeneity. We can think of the decision tree as a map of different path. For a distinct combination of predictor variables and their observed values, we would enter a specific path, which gives the classification in the leaf of the decision tree.

The decision tree approach does not require any



# Decision Trees (III)



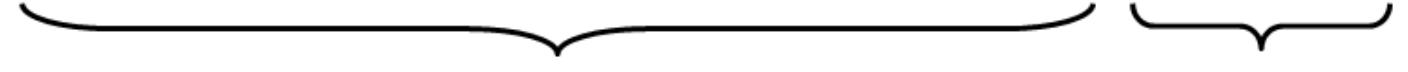
Source: <http://iopscience.iop.org/article/10.1088/1749-4699/5/1/015004>

# Overview of important Decision Tree Methods

| Name              | CART  | ID3   | C5.0  | CHAID  | Random Forests   |
|-------------------|---|---|---|--|--|
| Idea              | Choose the attribute with the highest information content | One of the first methods from Quinlan; uses the concept of information gain | Like ID3 based on the concept of information gain | Choose the attribute that is most dependent on the target variable | Construct many trees with different sets of features and samples (randomly). Result by voting. |
| Measure used      | Gini-Index  | Information gain (entropy)  | Ratio of information gain                         | Chi-square split   | Optional, mostly Gini-Index  |
| Type of Splitting | Binary  | Complete, pruning   | Complete, pruning                                 | Complete, pruning  | Complete   |

# Introductory Example

| Day | outlook  | temperature | humidity | wind   | play tennis |
|-----|----------|-------------|----------|--------|-------------|
| 1   | sunny    | hot         | high     | light  | no          |
| 2   | sunny    | hot         | high     | strong | no          |
| 3   | overcast | hot         | high     | light  | yes         |
| 4   | rainy    | mild        | high     | light  | yes         |
| 5   | rainy    | cold        | normal   | light  | yes         |
| 6   | rainy    | cold        | normal   | strong | no          |
| 7   | overcast | cold        | normal   | strong | yes         |
| 8   | sunny    | mild        | high     | light  | no          |
| 9   | sunny    | cold        | normal   | light  | yes         |
| 10  | rainy    | mild        | normal   | light  | yes         |
| 11  | sunny    | mild        | normal   | strong | yes         |
| 12  | overcast | mild        | high     | strong | yes         |
| 13  | overcast | hot         | normal   | light  | yes         |
| 14  | rainy    | mild        | high     | strong | no          |

  
**features**      **class variable**

# Splitting with Entropy in ID3

| Day | wind   | play tennis |
|-----|--------|-------------|
| 1   | light  | no          |
| 2   | strong | no          |
| 3   | light  | yes         |
| 4   | light  | yes         |
| 5   | light  | yes         |
| 6   | strong | no          |
| 7   | strong | yes         |
| 8   | light  | no          |
| 9   | light  | yes         |
| 10  | light  | yes         |
| 11  | strong | yes         |
| 12  | strong | yes         |
| 13  | light  | yes         |
| 14  | strong | no          |

ID3 uses the entropy as a measure for the impurity of a node t:

$$\text{entropy}(t) = - \sum_{i=1}^k p_i \cdot \log_2 p_i$$

where  $p_i$  = relative frequency of cases in a node, that belong to a class i

Attribute Wind:

**light: 6 x yes, 2 x no**

$$\Rightarrow \text{entropy}(\text{light}) = - \left( \frac{6}{8} \cdot \log_2 \frac{6}{8} + \frac{2}{8} \cdot \log_2 \frac{2}{8} \right) = 0.811$$

**strong: 3 x yes, 3 x no**

$$\Rightarrow \text{entropy}(\text{strong}) = - \left( \frac{3}{6} \cdot \log_2 \frac{3}{6} + \frac{3}{6} \cdot \log_2 \frac{3}{6} \right) = 1.0$$

# Calculating the Information Gain

The information gain is a measure, that shows (by combination of the entropies) the appropriateness of an attribute for splitting:

$$\text{Information gain} = \text{entropy}(t) - \sum_{i=1}^m \frac{t_i}{t} \cdot \text{entropy}(t_i)$$

where m = number of values (here two: light, strong),  $t_i$  = number of data sets with strong or light wind (8 resp. 6), t = total number of data sets (14) and  $\text{entropy}(t)$  = entropy before splitting.



# Decision using ID3

Information gain (outlook) = 0.246

Information gain (humidity) = 0.151

Information gain (wind) = 0.048

Information gain (temperature) = 0.029

We choose the attribute with the largest information gain (here: outlook) for the first splitting.

As solution we obtain the following tree:

# Decision using C5.0

ID3 tends to favor attributes that have a large number of values, resulting in larger trees. For example, if we have an attribute that has a distinct value for each record, then the entropy is 0, thus the information gain is maximal.

To compensate for this, C5.0 is a further development that uses the information gain ratio as a splitting criterion:

$$\text{Information gain ratio} = \frac{\text{entropy}(t) - \sum_{i=1}^m \frac{t_i}{t} \cdot \text{entropy}(t_i)}{- \sum_{i=1}^m \frac{t_i}{t} \cdot \log_2 \frac{t_i}{t}}$$

In the case of our example the GainRatio of Windy is

$$\text{Information gain ratio(Windy)} = \frac{0.048}{-\frac{6}{14} \cdot \log_2 \frac{6}{14} - \frac{8}{14} \cdot \log_2 \frac{8}{14}} = 0.049$$

# Handling Numerical Attributes

Numerical attributes are usually splitted binary. In contrast to categorical attributes many possible splitting points exist .

The splitting point with the highest information gain is looked for. For this, the potential attribute is sorted according to its values first and then all possible splitting point and the corresponding information gains are calculated. In extreme cases there exists  $n-1$  possibilities.

| day | temperature | play |
|-----|-------------|------|
| 5   | 64          | yes  |
| 6   | 65          | no   |
| 7   | 66          | yes  |

Test the splitting point temperature = 71.5:

temperature < 71.5: 4 x yes, 2 x no

# The CART Algorithm

The CART algorithm (Classification And Regression Trees) constructs trees that have only binary splits. Like C5.0, it is able to handle categorical and numerical attributes.

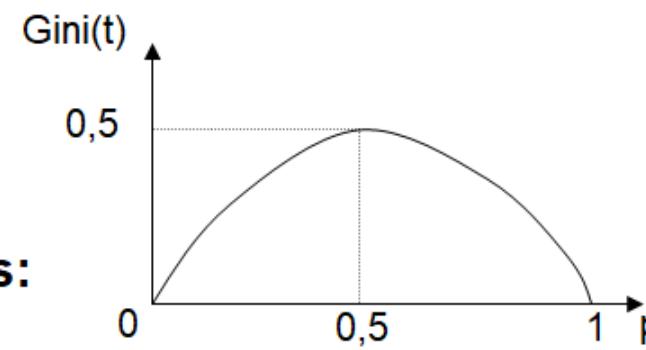
As a measure for the impurity of a node  $t$ , CART uses the Gini Index. In the case of two classes the Gini Index is defined as:

$$\text{Gini}(t) = 2 \cdot p_1 \cdot p_2 = 2 \cdot p_1 \cdot (1 - p_1)$$

where  $p_i$  = relative frequency of cases in a node, that belong to a class  $i$ ,  $i = 1, 2$ .

A node  $t$  is splitted by the attribute that minimizes:

$$\frac{t_1}{t} \cdot \text{Gini}(t_1) + \frac{t_2}{t} \cdot \text{Gini}(t_2) \rightarrow \text{Min}$$



# Splitting in CART

| Day | wind   | play tennis |
|-----|--------|-------------|
| 1   | light  | no          |
| 2   | strong | no          |
| 3   | light  | yes         |
| 4   | light  | yes         |
| 5   | light  | yes         |
| 6   | strong | no          |
| 7   | strong | yes         |
| 8   | light  | no          |
| 9   | light  | yes         |
| 10  | light  | yes         |
| 11  | strong | yes         |
| 12  | strong | yes         |
| 13  | light  | yes         |
| 14  | strong | no          |

**Attribute Wind:**

**light: 6 x yes, 2 x no**

$$\Rightarrow \text{Gini}(\text{light}) = 2 \cdot \frac{6}{8} \cdot \left(1 - \frac{6}{8}\right) = 0.375$$

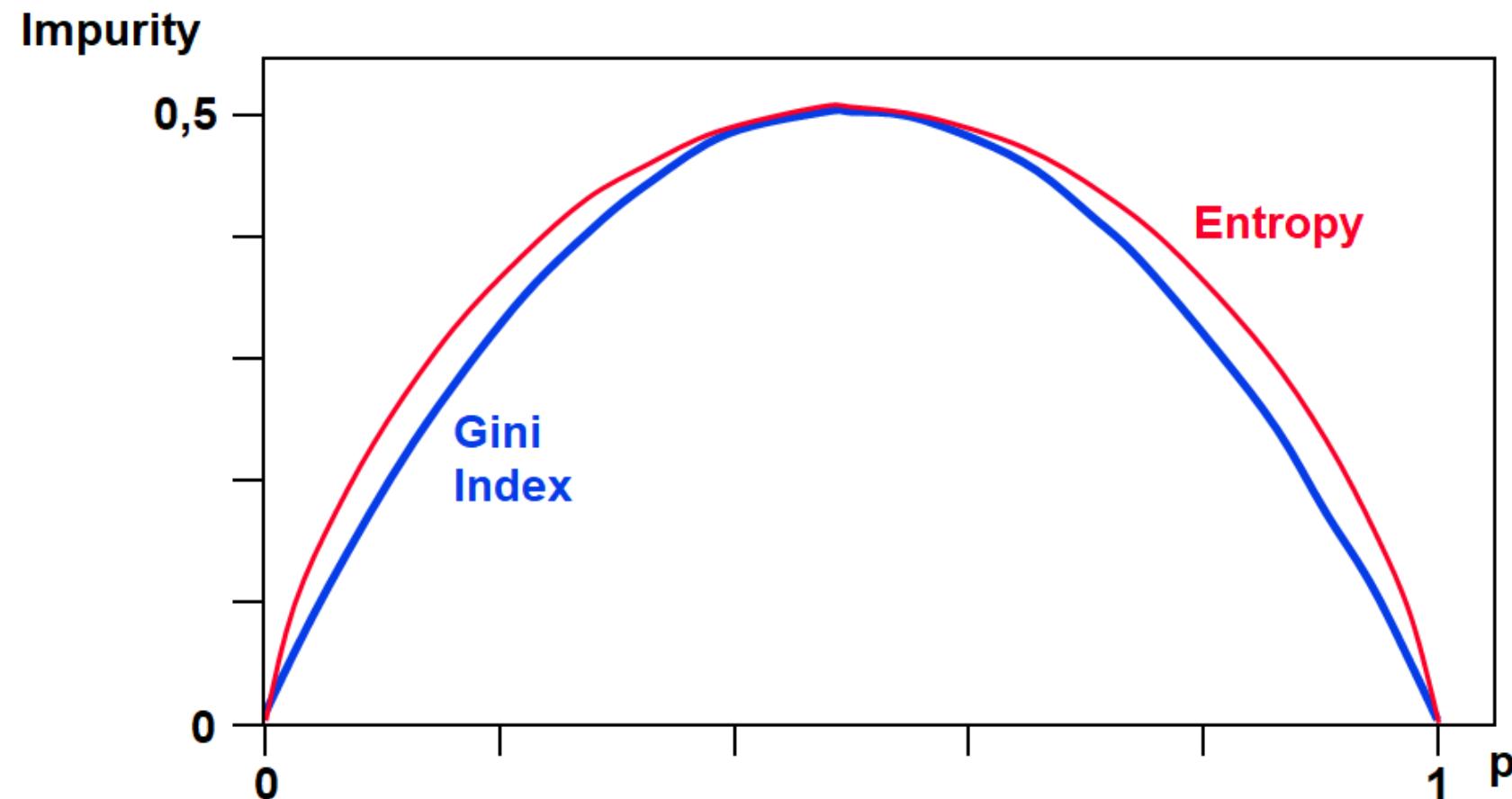
**strong: 3 x yes, 3 x no**

$$\Rightarrow \text{Gini}(\text{strong}) = 2 \cdot \frac{3}{6} \cdot \left(1 - \frac{3}{6}\right) = 0.5$$

**resulting in:**

$$\begin{aligned} \text{Information gain} &= 0.46 - \left( \frac{8}{14} \cdot 0.375 + \frac{6}{14} \cdot 0.5 \right) \\ &= 0.0306 \end{aligned}$$

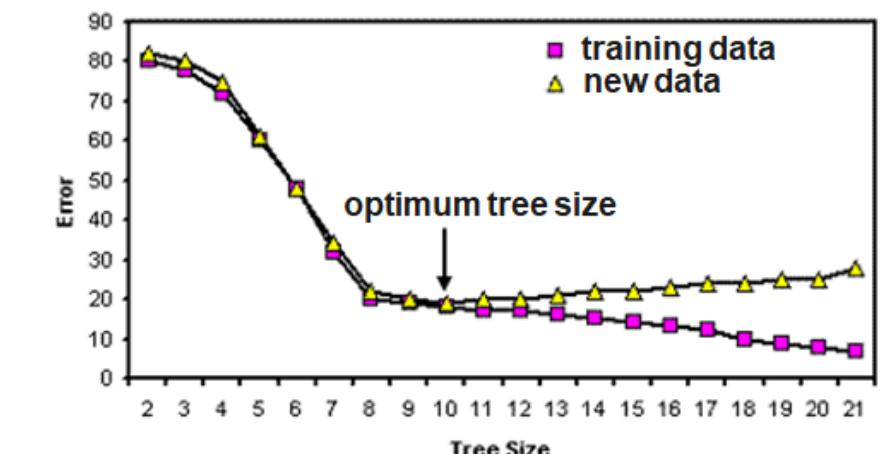
# Coherence between Entropy and Gini Index



Remark: Entropy has been scaled from (0, 1) to (0, 0.5)!

Most decision tree algorithms partition training data until every node contains objects of a single class, or until further partitioning is impossible because two objects have the same value for each attribute but belong to different classes. If there are no such conflicting objects, the decision tree will correctly classify all training objects.

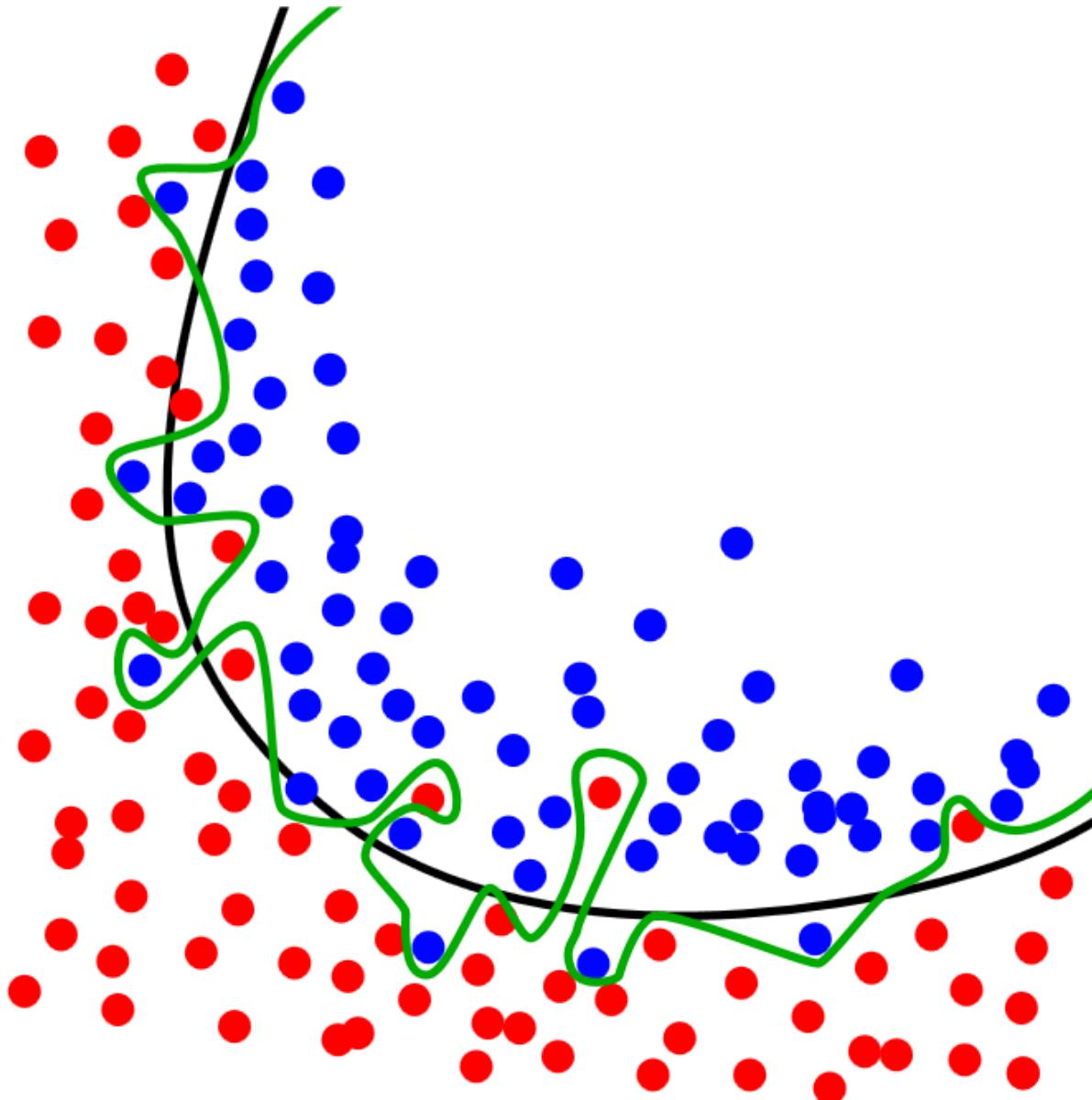
If tree performance is measured from the number of correctly classified cases it is common to find that the training data gives an over-optimistic



## Overfitting (II)

The learner overfits to correctly classify, the noisy data objects

Noisy or dirty data objects

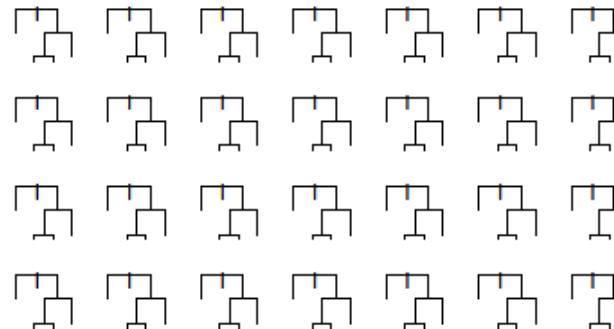


## Random Forest (I)

Random forest is an ensemble classifier that consists of many decision trees.

For every tree a subset of the data objects and a subset of features is randomly chosen. Then the tree is constructed usually using the Gini Index.

In the end, a simple majority vote is taken for prediction.



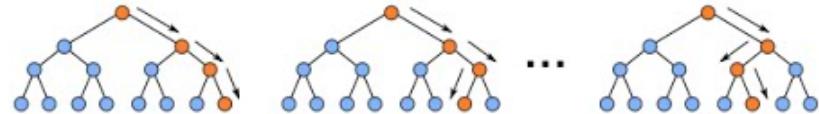
Algorithm :

1. Create  $n$  samples from the original data. Frequent sample size is  $2/3$ .
2. For each of the samples, grow a tree, with the following modification: at each node, rather than choosing the best split among all predictors, randomly sample  $m^*$  of the  $m$  predictors and choose the best split from among those variables.
3. Predict by aggregating the predictions of the  $n$  trees (majority votes).

## Random Forest (II)

Voting-Principle of Random Forest:

To avoid overfitting effects , the size and the depth of the trees can be restricted .



## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

##### 4.4.1.1 K-Nearest Neighbors

##### 4.4.1.2 Evaluating the Quality of Classification

##### 4.4.1.3 Decision Tree Approaches

##### 4.4.1.4 Logistic Regression

##### 4.4.1.5 Neural Networks

##### 4.4.1.6 Resampling

##### 4.4.1.7 Ensemble Learning

#### 4.4.2 Regression

# Introductory Example

|    | A    | B    | C    | D    | E    | F    | G    |
|----|------|------|------|------|------|------|------|
| 1  | URL1 | URL2 | URL3 | URL4 | URL5 | URL6 | URL7 |
| 2  | 1    | 0    | 0    | 0    | 0    | 0    | 0    |
| 3  | 0    | 0    | 0    | 0    | 0    | 0    | 0    |
| 4  | 1    | 0    | 0    | 0    | 0    | 0    | 0    |
| 5  | 0    | 1    | 0    | 1    | 1    | 0    | 0    |
| 6  | 1    | 0    | 0    | 0    | 0    | 0    | 0    |
| 7  | 0    | 0    | 0    | 0    | 0    | 0    | 0    |
| 8  | 1    | 0    | 0    | 0    | 0    | 0    | 0    |
| 9  | 0    | 0    | 0    | 0    | 0    | 1    | 0    |
| 10 | 0    | 0    | 0    | 0    | 0    | 0    | 0    |
| 11 | 1    | 0    | 0    | 0    | 0    | 0    | 0    |
| 12 | 0    | 0    | 0    | 0    | 1    | 0    | 1    |
| 13 | 1    | 0    | 1    | 0    | 0    | 0    | 0    |
| 14 | 0    | 0    | 0    | 0    | 0    | 0    | 0    |
| 15 | 0    | 0    | 0    | 0    | 0    | 1    | 0    |
| 16 | 0    | 0    | 0    | 0    | 0    | 0    | 0    |
| 17 | 1    | 0    | 0    | 0    | 0    | 0    | 0    |
| 18 | 0    | 0    | 0    | 0    | 0    | 0    | 0    |
| 19 | 1    | 0    | 0    | 0    | 0    | 0    | 0    |
| 20 | 1    | 0    | 0    | 0    | 0    | 0    | 0    |

| BGT  | BGU  | BGV   |
|------|------|-------|
| URL1 | URL2 | click |
| 0    | 0    | 0     |
| 0    | 0    | 0     |
| 0    | 0    | 1     |
| 0    | 1    | 1     |
| 0    | 0    | 0     |
| 1    | 0    | 0     |
| 0    | 0    | 1     |
| 0    | 0    | 0     |
| 0    | 0    | 1     |
| 0    | 0    | 1     |
| 0    | 0    | 0     |
| 0    | 0    | 0     |
| 0    | 0    | 1     |
| 1    | 0    | 1     |
| 0    | 1    | 1     |
| 0    | 0    | 0     |
| 0    | 0    | 0     |
| 0    | 0    | 0     |

Why not classical linear regression?

It is possible to implement a linear regression on such a dataset where  $Y=\{0,1\}$ .

Problems:

The predicted values of the linear model can be greater than 1 or less than 0

$\epsilon$  is not normally distributed because  $Y$  takes on only two values

The error terms are heteroscedastic (the error variance is not constant for all values of  $X$ )

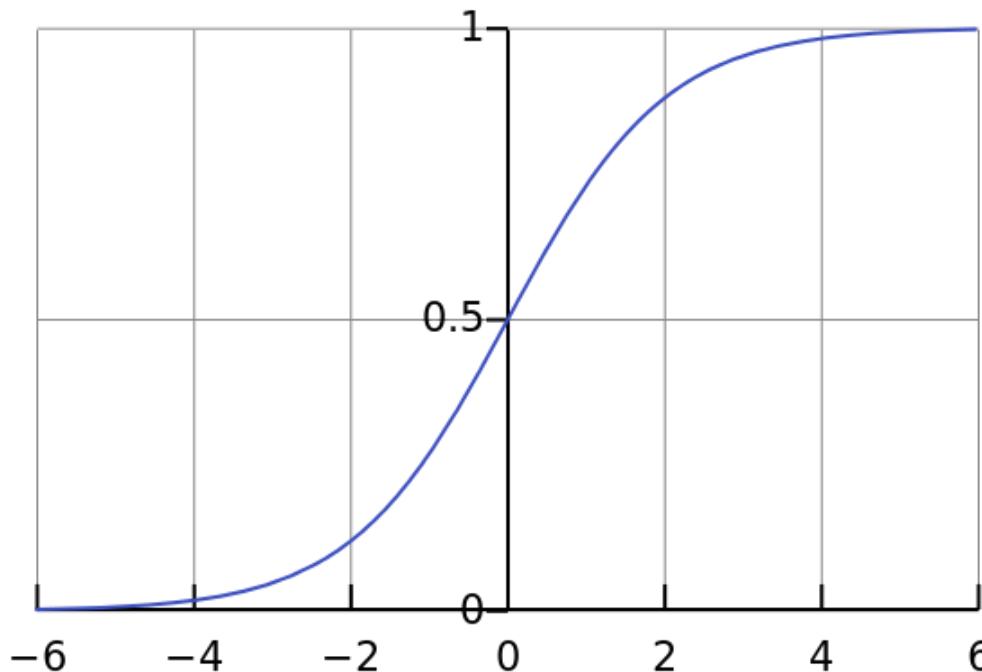
Source: Bichler (2015): Course Business Analytics, TU München

## Logistic regression (I)

Logistic regression is a regression model where the dependent variable is categorical. The classical logistic regression is a binary classifier, where the dependent variable has two states. The output of a logistic regression model ranges between 0 and 1.

Logistic regression uses the logistic function (or Sigmoid function) because it can take an input with any value from negative to positive infinity, whereas the output always takes values between zero and one and hence is interpretable as a probability.

It is defined as:



## Logistic regression (II)

If we set

the logistic function can now be written as:

We interpret  $F(x)$  as the conditional probability that the class attribute has the value 1 with the given input vector  $x$ .

The coefficients  $\beta_0$  and  $\beta$  can be estimated via Maximum Likelihood Estimation.

The parameter  $\beta_0$  represents the unconditional probability of “ $Y=1$ ” knowing nothing about the feature vector  $x$ .

The parameter vector  $\beta$  defines the slope of the logit function. It determines the extent to which certain features contribute for increased or decreased likelihood to “ $Y=1$ ”.

The output of a logistic model is a probability. To use this for classification purposes:

If the predicted probability is  $> 0.5$  the label is 1  
and otherwise 0.

## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

##### 4.4.1.1 K-Nearest Neighbors

##### 4.4.1.2 Evaluating the Quality of Classification

##### 4.4.1.3 Decision Tree Approaches

##### 4.4.1.4 Logistic Regression

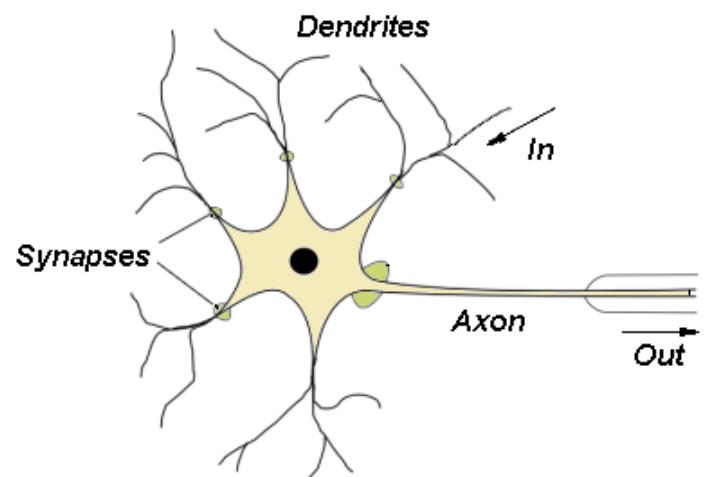
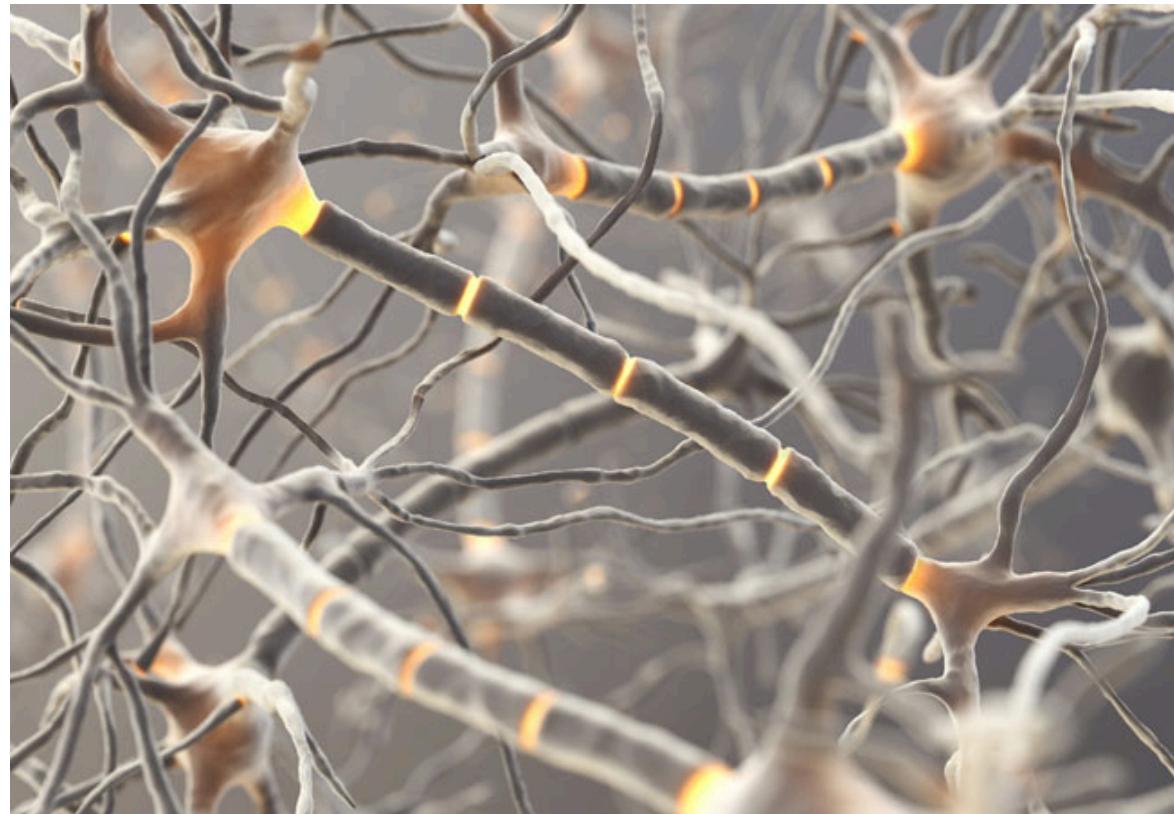
##### 4.4.1.5 Neural Networks

##### 4.4.1.6 Resampling

##### 4.4.1.7 Ensemble Learning

#### 4.4.2 Regression

# Functionality of Human Neurons



An Easy Example (I)

**$f(x) = \text{Activation function}$**

e.g.

where  $t$  = Stimulus threshold

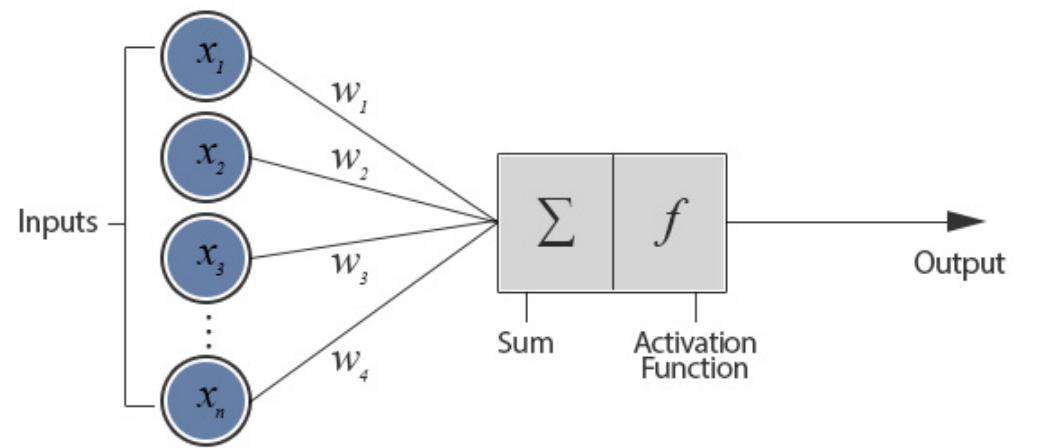
An Easy Example (II)

**$f(x) = \text{Activation function}$**

e.g.

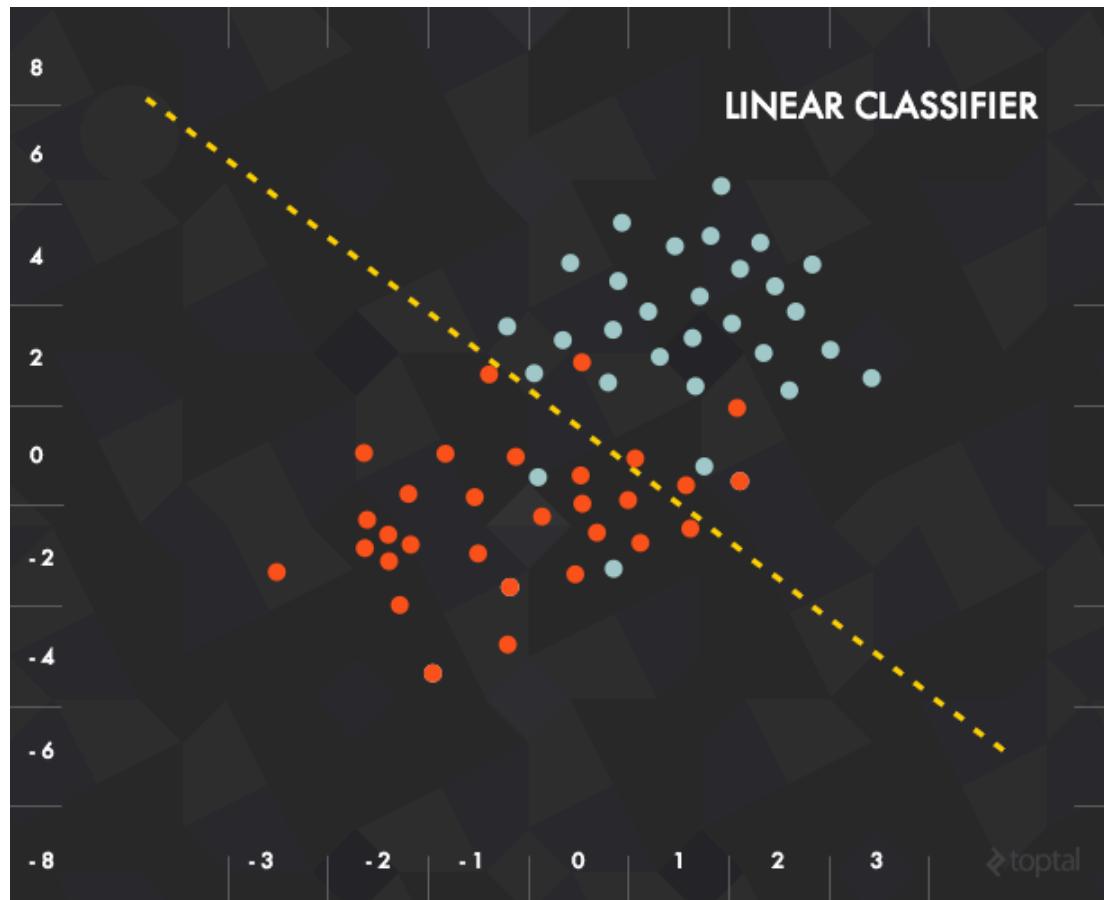
where  $t = \text{Stimulus threshold}$

# Functionality of a Neuron



For the case of  $n$  inputs, we can rewrite the neuron's function to

with  $b = -t$ .  $b$  is known as the perceptron's bias. The result of this function would then be fed into an activation function to produce a labeling



This results in a linear classifier. Finally, we have to pick a line that best separates the labeled data. The training of the perceptron consists of feeding it multiple training samples and calculating the output for each of them. After each sample, the weights  $w$  are adjusted in such a way so as to minimize the output error, defined for example as accuracy or MSE.

## The Multilayer Perceptron

The single perceptron approach has a major drawback: it can only learn linear functions. To address this problem, we'll need to use a multilayer perceptron, also known as feedforward neural network. Here, we add layers between the input and the output layer, so-called **hidden layers**. The hidden layer is where the network stores its internal abstract representation of the training data.

Input Neurons : receive signals from the outer world .

Hidden Neurons : have an internal representation of the outer world .

Output Neurons : pass signals to the outer world .

## Types of Activation Functions

A linear composition of linear functions is still just a linear function, so most neural networks use non-linear activation functions:

**tangens hyperbolicus**

logistic function (sigmoid)

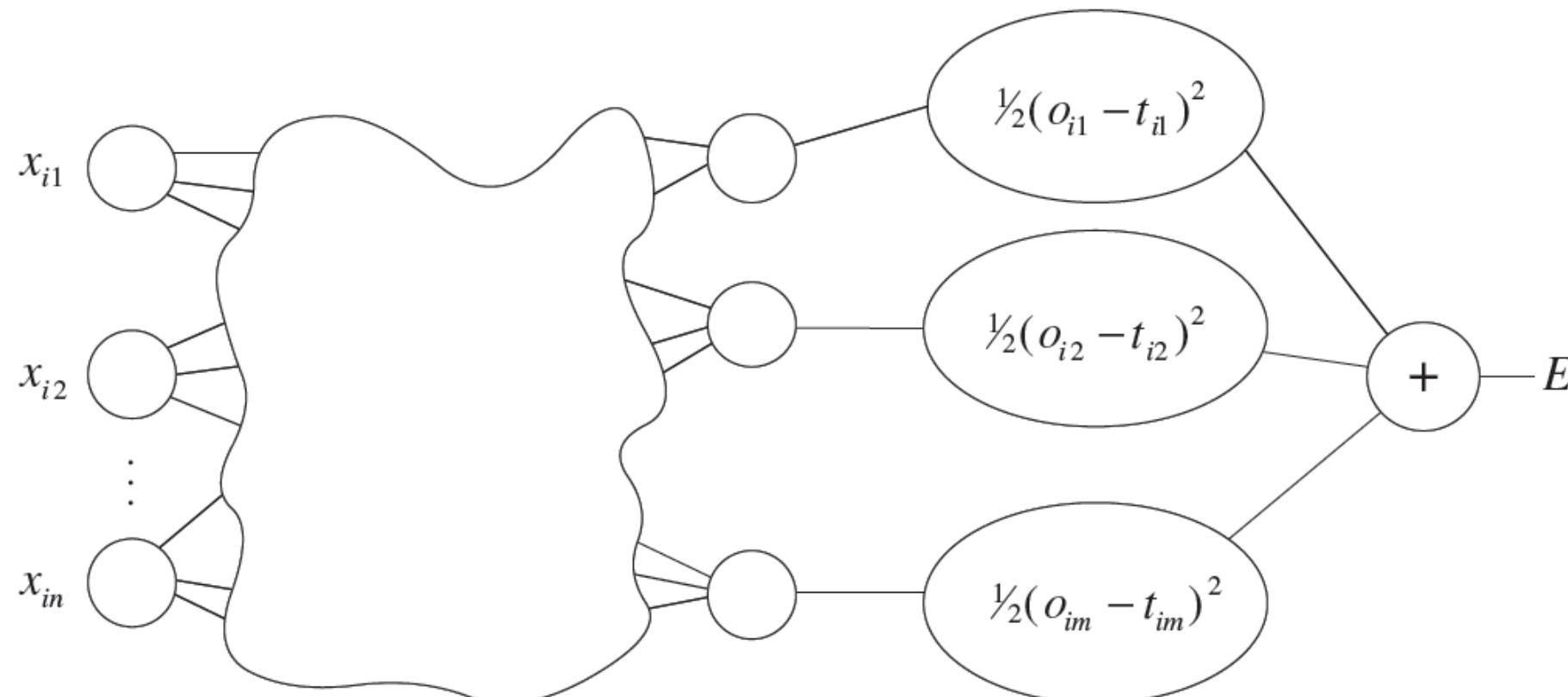
# Design of a Multilayer Perceptron

- The Backpropagation algorithm is used for calculating the weights. In a training set, weights are iteratively calculated using training data sets in such a way that the difference between the calculated and the expected (true) results is minimized. Because the simultaneous calculation of all weights is not possible, they must be found via a learning process. The backpropagation algorithm looks for the minimum of the error function in weight space using the method of gradient descent.
- The procedure in principle:
  - (1) Define the initial weights
  - (2) Put the training set into the input layer
  - (3) Calculate the result (value of the output layer) via successive processing one layer after the other
  - (4) Compare the output values and target values and calculate the difference
  - (5) Iterate steps (2) to (4) for every training set
  - (6) Calculate the total error. Adjust the weights beginning with the output layer towards the input layer (backpropagation)
  - (7) Iterate steps (2) to (6) until the total error reaches the defined error level or the number

## Adjusting the Weights (I)

The error of a training set  $i$  is calculated using the quadratic deviation between the neurons of the output layer and their corresponding true values  $t_{ij}$ .

The sum of the errors of all  $h$  training objects is the total error value  $E$ :

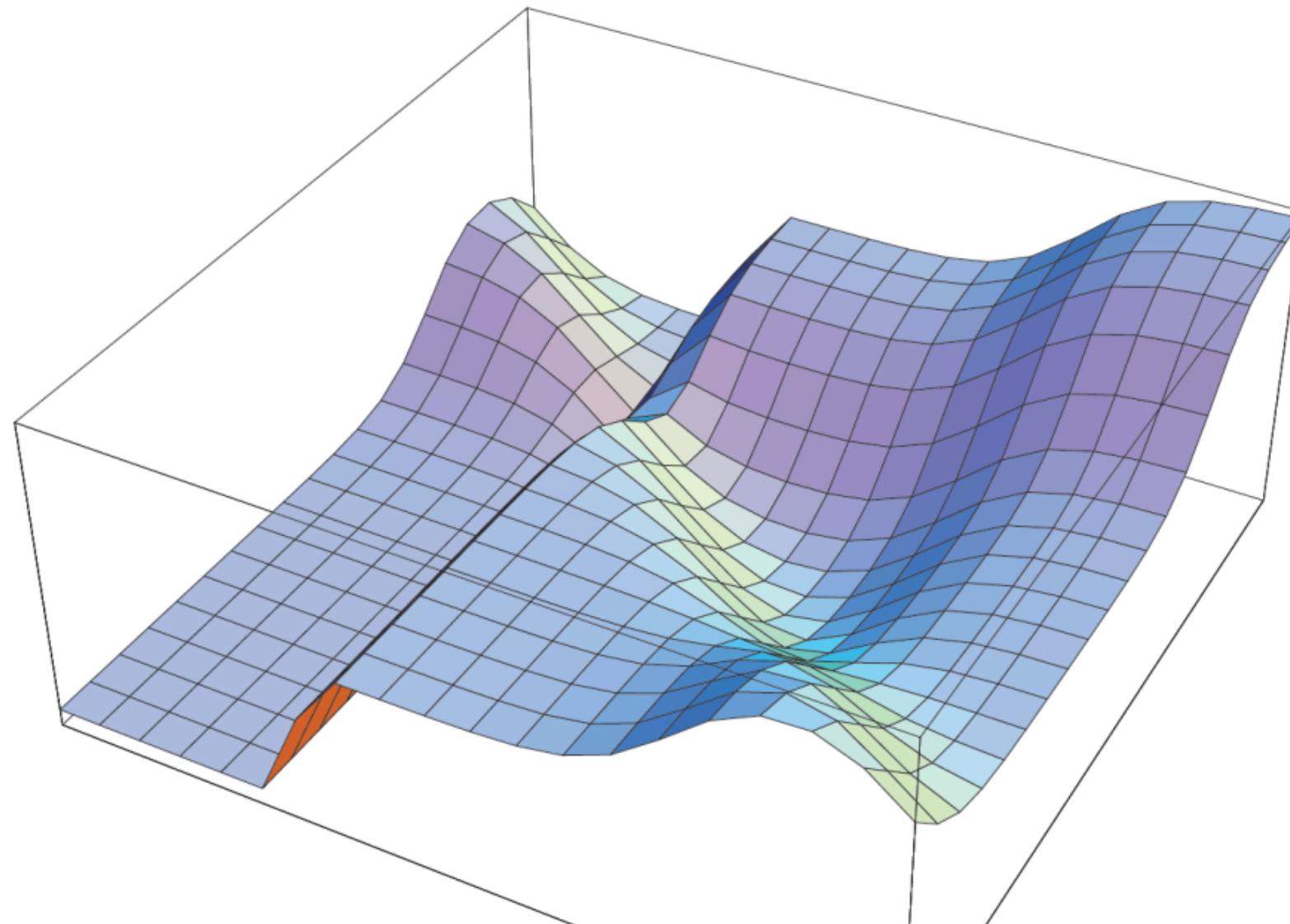


## Adjusting the Weights (II)

The function  $E$  has to be minimized. Because it depends on the output neurons  $o_j, i_l$  automatically depends on their weights to the precedent layer(s) :

Thus, the weights have to be found where  $E$  is minimal.

Examples of Error functions with two weights:



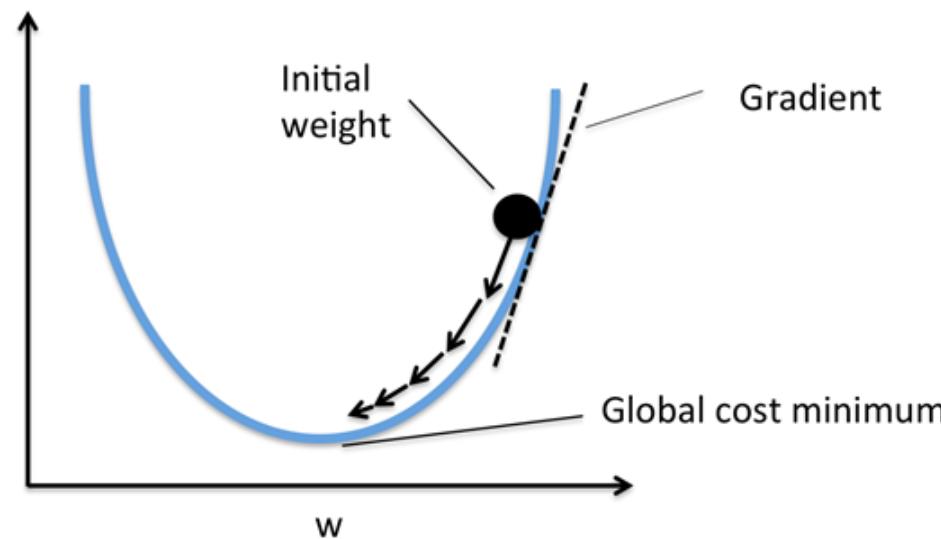
## Adjusting the Weights (III)

To minimize the error (cost) function  $E$  the backpropagation algorithm uses the method of gradient descent . This method searches those weights, where the vector containing the partial first derivatives of the error function (gradient) is equal to the zero vector (minimum):

To adjust the weight  $w_{ij}$  , which connects neurons i to j, the formula is:

where  $\alpha$  represents a predefined learning rate , which defines the step length of each iteration in the negative gradient direction and  $x_i$  denote the output value of neuron i .

The adjusted weight is then computed via



## Principle of Gradient Descent (I)

Gradient descent is used to find the minimum of the error function . It works iterative. In an 1-dimensional world, we define the error by

The error function is at minimum if the error is equal to zero.

The prediction is the result of a combination of input and weight

The weight as the dynamic component is now adjusted until the error is at minimum. Starting with an initial weight, gradient descent jumps step by step into the minimum by adjusting the weight. The adjustment is done by calculating the direction and the amount for a step via

Now, the weight is adjusted via

After repeating this several times, the minimum is reached.

## Principle of Gradient Descent (II)

### The formula

represents the derivative of the error to the weight.

A derivative is a term that is calculated as the slope (or gradient) of a graph at a particular point. The slope is described by drawing a tangent line to the graph at the point. So, if we are able to compute this tangent line, we might be able to compute the desired direction to reach the minima.

Since the weight only indirectly affects the error, the chain rule must be applied

## Principle of Gradient Descent (III)

Gradient Descent isn't perfect. When the gradients are too big we might overshoot so much that we're even farther away than we started

This problem is destructive because overshooting this far means we land at an even steeper slope in the opposite direction. This causes us to overshoot again even farther.

If the gradients are too big, we can make them smaller. We do this by multiplying them by a single number between 0 and 1 (such as 0.01). This fraction is typically named alpha.

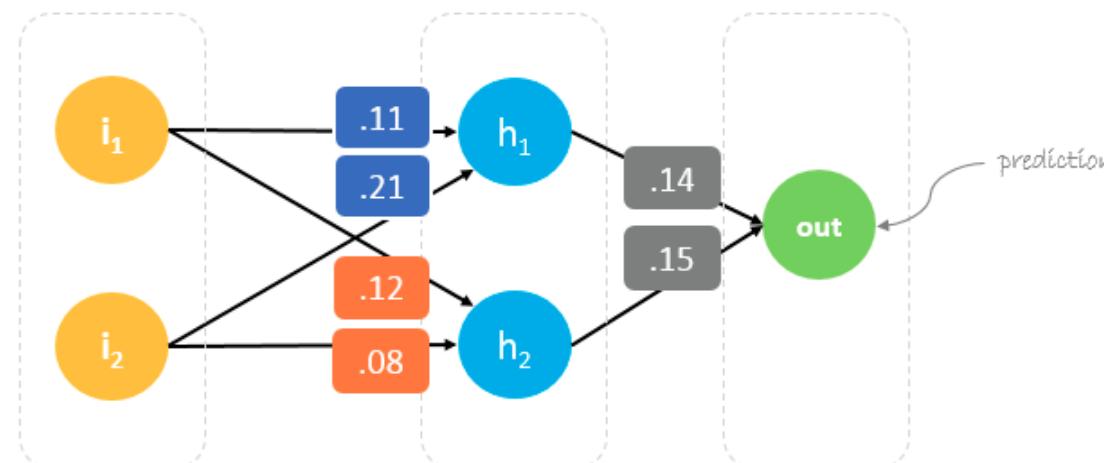
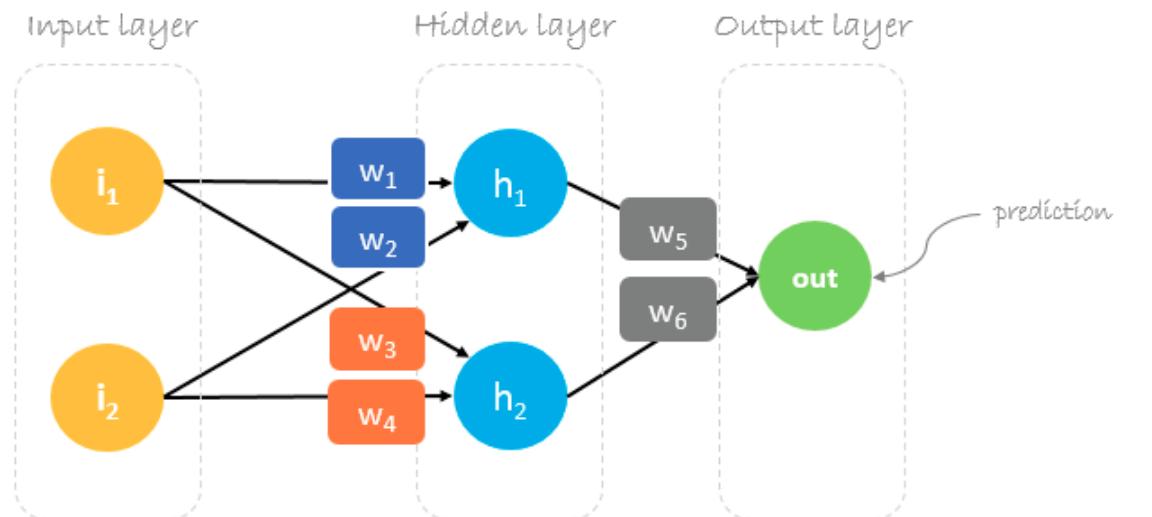
Thus, the adjustment of the weights is done by

Source: <https://iamtrask.github.io/2015/07/27/python-network-part2/>

## Backpropagation Step by Step (I)

In the following, the backpropagation process will be demonstrated using a simple Neural Network consisting of three layers: Input layer with two inputs neurons, one hidden layer with two neurons, and output layer with a single neuron:

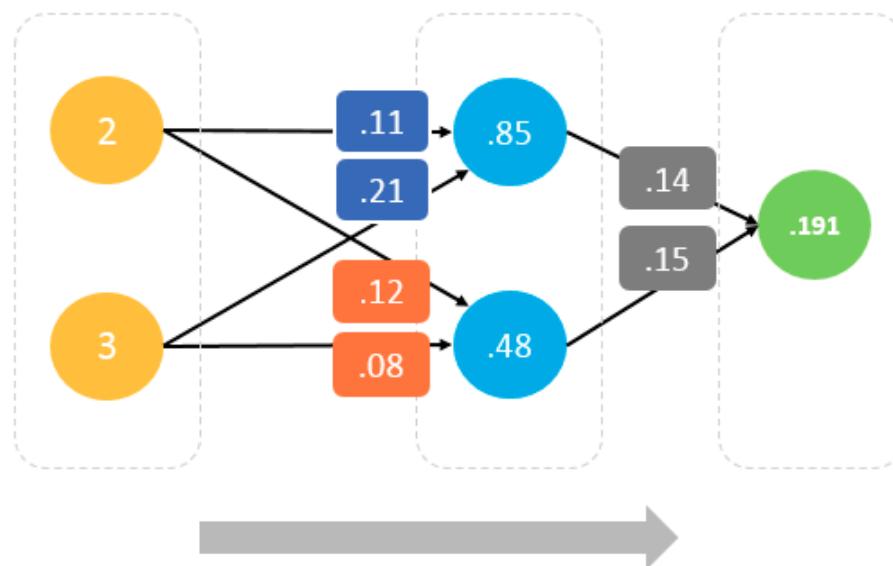
Our initial weights will be:  $w_1 = 0.11$ ,  $w_2 = 0.21$ ,  $w_3 = 0.12$ ,  $w_4 = 0.08$ ,  $w_5 = 0.14$  and  $w_6 = 0.15$ .



## Backpropagation Step by Step (II)

Our dataset has one sample with two inputs and one output with the values  $\text{inputs}=[2, 3]$  and  $\text{output}=[1]$ . We will use given weights and inputs to predict the output. Inputs are multiplied by weights; the results are then passed forward to next layer:

**For reasons of simplification, no activation function is used in the neurons.**



$$[2 \ 3] \cdot \begin{bmatrix} 0.11 & 0.12 \\ 0.21 & 0.08 \end{bmatrix} \cdot \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix} = [0.85 \ 0.48] \cdot \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix} = [0.191]$$

$$2 \times 0.11 + 3 \times 0.21 = 0.85$$

$$2 \times 0.12 + 3 \times 0.08 = 0.48$$

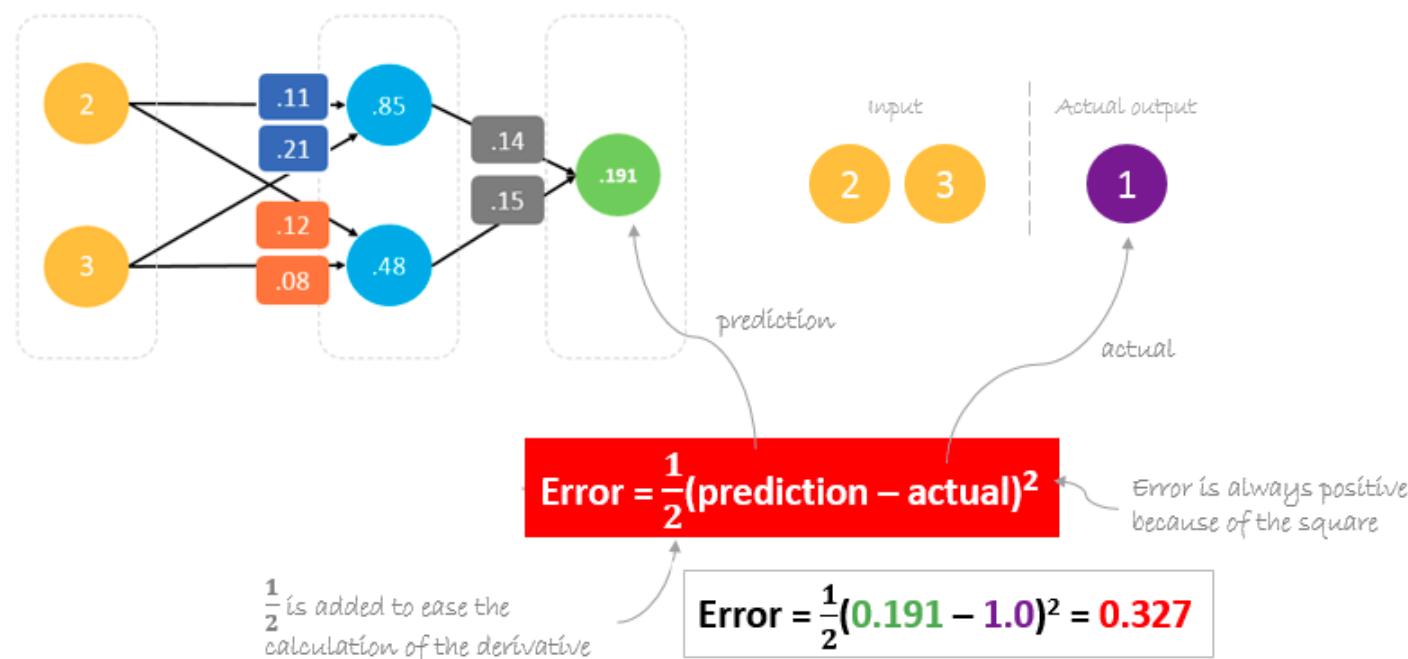
$$0.85 \times 0.14 + 0.48 \times 0.15 = 0.191$$

Source: <http://hmkcode.github.io/ai/backpropagation-step-by-step>

## Backpropagation Step by Step (III)

The network output, or prediction, is not even close to actual output. We can calculate the difference or the error as following:

Our main goal of the training is to reduce the error or the difference between prediction and actual output. Since actual output is constant, “not changing”, the only way to reduce the error is to change prediction value. The question now is, how to change prediction value?



Source: <http://hmkcode.github.io/ai/backpropagation-step-by-step>

## Backpropagation Step by Step (IV)

By decomposing prediction into its basic elements we can find that weights are the variable elements affecting prediction value. To change prediction value, we need to adjust the weights:

We do this using Backpropagation. To find a local minimum of a function using gradient descent, one takes steps proportional to the negative of the gradient of the function at the current point:

For example, we update w<sub>6</sub>:

$$\begin{aligned}\text{prediction} &= (\textcolor{blue}{h_1}) w_5 + (\textcolor{blue}{h_2}) w_6 \\ &= (\textcolor{orange}{i_1} \textcolor{blue}{w_1} + \textcolor{orange}{i_2} \textcolor{blue}{w_2}) w_5 + (\textcolor{orange}{i_1} \textcolor{red}{w_3} + \textcolor{orange}{i_2} \textcolor{red}{w_4}) w_6\end{aligned}$$

$$*W_x = W_x - \textcolor{orange}{a} \left( \frac{\partial \text{Error}}{\partial W_x} \right)$$

↓ Old weight      ↓ Derivative of Error  
 New weight      with respect to weight  
 ↑                  ↑ Learning rate

We can picture gradient descent optimization as a hiker (the weight coefficient) who wants to climb down a mountain (cost function) into a valley (cost minimum), and each step is determined by the negative of the slope (gradient) of the hill (loss function).

## Backpropagation Step by Step (V)

The derivation of the error function is evaluated by applying the chain rule:

To update  $w_6$  we can apply the following formula:

Similarly, we can derive the update formula for  $w_5$  and any other weights existing between the output and the hidden layer:

$$*W_6 = W_6 - \alpha \Delta h_2$$

$$*W_5 = W_5 - \alpha \Delta h_1$$

Source: <http://hmkcode.github.io/ai/backpropagation-step-by-step>

## Backpropagation Step by Step (VI)

When moving backward to update  $w_1$ ,  $w_2$ ,  $w_3$  and  $w_4$  existing between input and hidden layer, the partial derivative for the error function with respect to  $w_1$ , for example, will be as following:

We can find the update formula for the remaining weights  $w_2$ ,  $w_3$  and  $w_4$  in the same way.

Source: <http://hmkcode.github.io/ai/backpropagation-step-by-step>

## Backpropagation Step by Step (VII)

In summary, the update formulas for all weights will be:

We can rewrite the update formulas in matrices:

$$w_6 = w_6 - a (\Delta h_2)$$

$$w_5 = w_5 - a (\Delta h_1)$$

$$w_4 = w_4 - a (i_2 \cdot \Delta w_6)$$

$$w_3 = w_3 - a (i_1 \cdot \Delta w_6)$$

$$w_2 = w_2 - a (i_2 \cdot \Delta w_5)$$

$$w_1 = w_1 - a (i_1 \cdot \Delta w_5)$$

$$\begin{bmatrix} w_5 \\ w_6 \end{bmatrix} = \begin{bmatrix} w_5 \\ w_6 \end{bmatrix} - a \Delta \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = \begin{bmatrix} w_5 \\ w_6 \end{bmatrix} - \begin{bmatrix} a \Delta h_1 \\ a \Delta h_2 \end{bmatrix}$$

$$\begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} = \begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} - a \Delta \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} \cdot \begin{bmatrix} w_5 & w_6 \end{bmatrix} = \begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} - \begin{bmatrix} a i_1 \Delta w_5 & a i_1 \Delta w_6 \\ a i_2 \Delta w_5 & a i_2 \Delta w_6 \end{bmatrix}$$

Source: <http://hmkcode.github.io/ai/backpropagation-step-by-step>

## Backpropagation Step by Step (VIII)

With the derived formulas we can now adjust the weights:

$$\Delta = 0.191 - 1 = -0.809 \quad \text{Delta} = \text{prediction} - \text{actual}$$

$$a = 0.05$$

$$\begin{bmatrix} w_5 \\ w_6 \end{bmatrix} = \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix} - 0.05(-0.809) \begin{bmatrix} 0.85 \\ 0.48 \end{bmatrix} = \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix} - \begin{bmatrix} -0.034 \\ -0.019 \end{bmatrix} = \begin{bmatrix} 0.17 \\ 0.17 \end{bmatrix}$$

$$\begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} = \begin{bmatrix} .11 & .12 \\ .21 & .08 \end{bmatrix} - 0.05(-0.809) \begin{bmatrix} 2 \\ 3 \end{bmatrix} \cdot \begin{bmatrix} 0.14 & 0.15 \end{bmatrix}$$

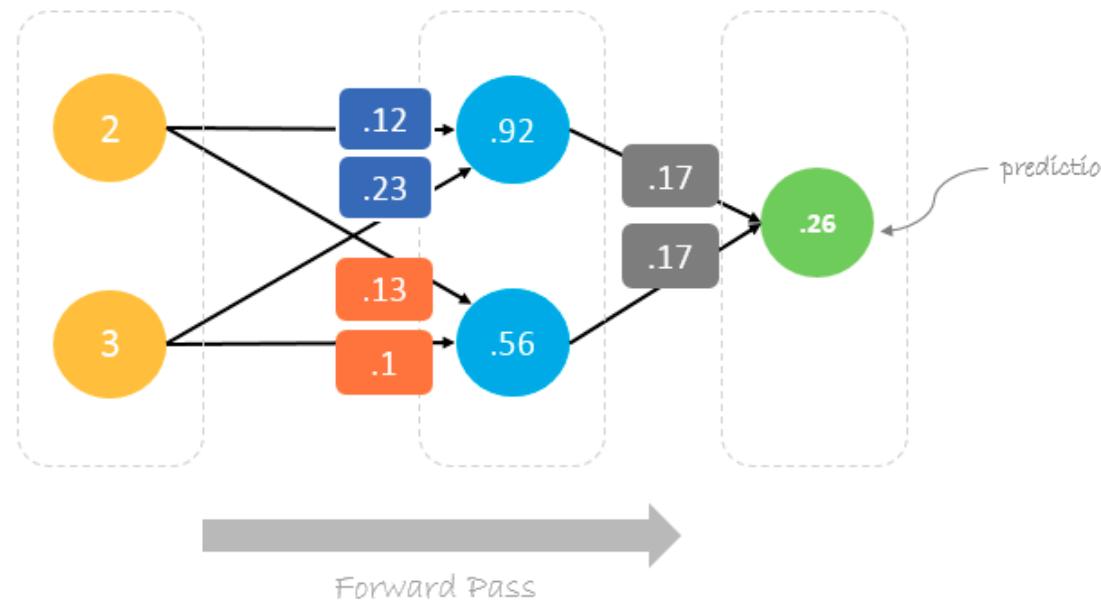
$$= \begin{bmatrix} .11 & .12 \\ .21 & .08 \end{bmatrix} - \begin{bmatrix} -0.011 & -0.012 \\ -0.017 & -0.018 \end{bmatrix} = \begin{bmatrix} .12 & .13 \\ .23 & .10 \end{bmatrix}$$

Source: <http://hmkcode.github.io/ai/backpropagation-step-by-step>

## Backpropagation Step by Step (IX)

... and use the new weights to recalculate the example:

The new prediction 0.26 is bit closer to the output than the previously predicted one 0.191. We repeat now the same process until error is close or equal to zero.



$$[2 \ 3] \cdot \begin{bmatrix} 0.12 & 0.13 \\ 0.23 & 0.10 \end{bmatrix} = [0.92 \ 0.56] \cdot \begin{bmatrix} 0.17 \\ 0.17 \end{bmatrix} = [0.26]$$

$$2 \times .12 + 3 \times .23 = .85 \quad .92 \times .17 + .56 \times .17 = .26$$

$$2 \times .13 + 3 \times .10 = .48$$

Source: <http://hmkcode.github.io/ai/backpropagation-step-by-step>

## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

##### 4.4.1.1 K-Nearest Neighbors

##### 4.4.1.2 Evaluating the Quality of Classification

##### 4.4.1.3 Decision Tree Approaches

##### 4.4.1.4 Logistic Regression

##### 4.4.1.5 Neural Networks

##### 4.4.1.6 Resampling

##### 4.4.1.7 Ensemble Learning

#### 4.4.2 Regression

## Problems with fixed Training and Test Samples

Method 1 optimize

Test data is used for two things:

Optimize the model training

Select the best model via testing the model quality

Method 2 optimize

Method 3 optimize

This contradicts the idea of independent testing and results in:

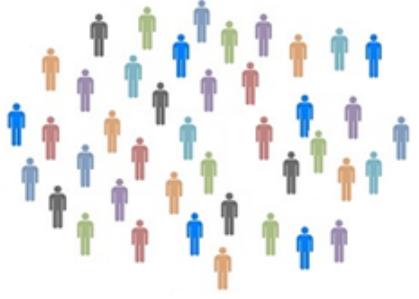
Endogenization of the test data

Selection Bias

... optimize

Rule : NEVER use any information from the test data for model training !

# Addressing the Endogeneity Problem



*Predictive \_ Model\_*

*Validation Sample*



- Training and test error can be highly variable, depending on precisely which observations included in the training set and which observations are included in the validation set (**Selection Bias**).
  - Example of different OLS models as a result of different samples:
  - To avoid such problems, one can use so-called resampling methods.

- In Data Science cross validation can be used for model selection and adjustment. In these cases, cross validation is applied to the training data set. For every iteration,  $k-1$  folds are used for model fitting and the remaining fold for testing the model (Validation). Every time, the quality measure (e.g. accuracy) for the validation fold is captured. At the end of this step, the average and the standard deviation of the measures are calculated. The best model is the one with the best ratio in high average and low standard deviation.
- Once the model type and its optimal parameters have been selected, a final model is trained using these hyper-parameters on the *full* training set, and the generalization quality is measured on the test set.



- Partition the Dataset into a training and test set
- For every hyperparameter value combination apply cross validation
- For the combination with the highest (mean) quality calculate the final model with the complete training set
- Test the final model with the test set
- Compare the accuracies of training and test with regard to overfitting

Calculate the mean quality of the validation folds, e.g. mean accuracy or mean F1

# Cross Validation and Grid Search in Python

- Performs cross validation with the given hyperparameter combinations and manages the evaluation process

Using the original libraries and functions

`KNeighborsClassifier()`

`cross_val_score()`

- Performs cross validation

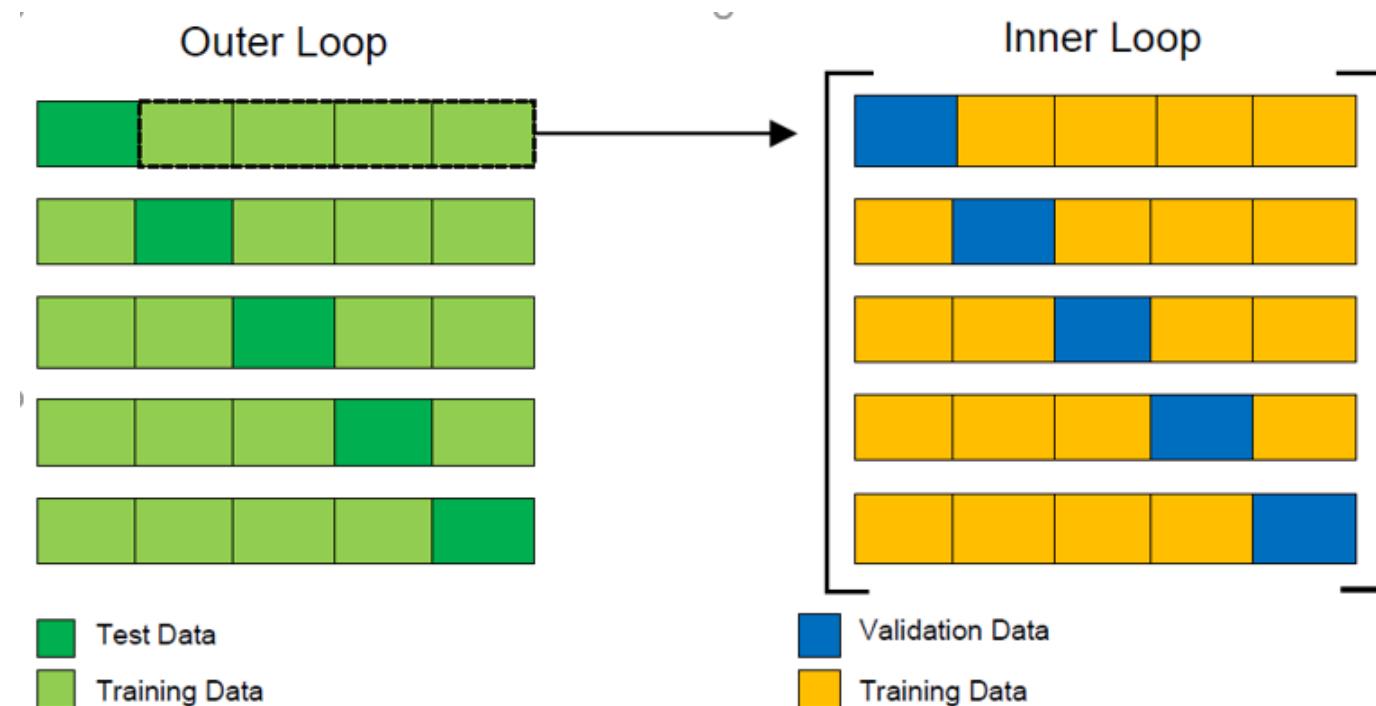
`DecisionTreeClassifier()`

`RandomForestCl...()`

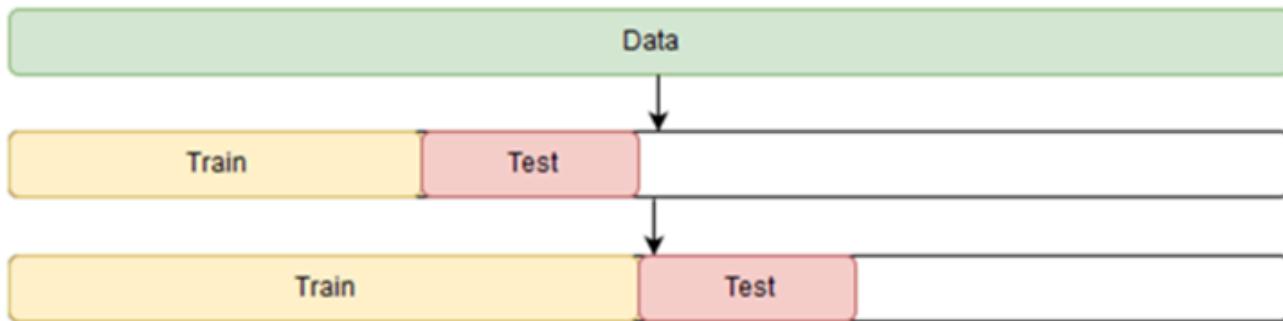
- 1. Grid Search
- Grid search sequentially goes through a preselected list of permutations for each hyperparameter and evaluates the entire search space.
- 2. Random Search
- Random search selects values for hyperparameters at random within a predefined distribution.
- While a grid search is able to find the best model given the provided options, limited compute resources means that in practice, the search space selected will have to be limited. A random search on the other hand does not iterate over the entire search space.

# Other Variants of Cross Validation

- 1. Repeated Cross Validation
- 2. Nested Cross Validation
- Different composition of the folds by random selection.



- In the case of time series, classical cross validation may cause problems. If we choose random samples and assign them to either the test set or the training set we are quickly in the situation of using values from the future to forecast values in the past. But we want to avoid future-looking when we train our model. If there is a temporal dependency between observations, we must preserve that relation during training and testing.
- A procedure that can be used for cross validating a time series model is cross validation on a rolling basis. Start with a small subset of data for training purpose, forecast for the later data points and then check the accuracy for the forecasted data points. The time frame for the forecast is then included as part of the next training dataset and subsequent data points are forecasted and so on.
- Scikit-learn provides a class *TimeSeriesSplit* to do this.



## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

##### 4.4.1.1 K-Nearest Neighbors

##### 4.4.1.2 Evaluating the Quality of Classification

##### 4.4.1.3 Decision Tree Approaches

##### 4.4.1.4 Logistic Regression

##### 4.4.1.5 Neural Networks

##### 4.4.1.6 Resampling

##### 4.4.1.7 Ensemble Learning

#### 4.4.2 Regression

## Ensemble Methods

Ensemble methods use different models (created via different data sets, feature sets or methods) that are simultaneously applied to the same problem. The results are sent to an aggregating operation that produces the final result.

The most widely used classes of ensemble methods are:

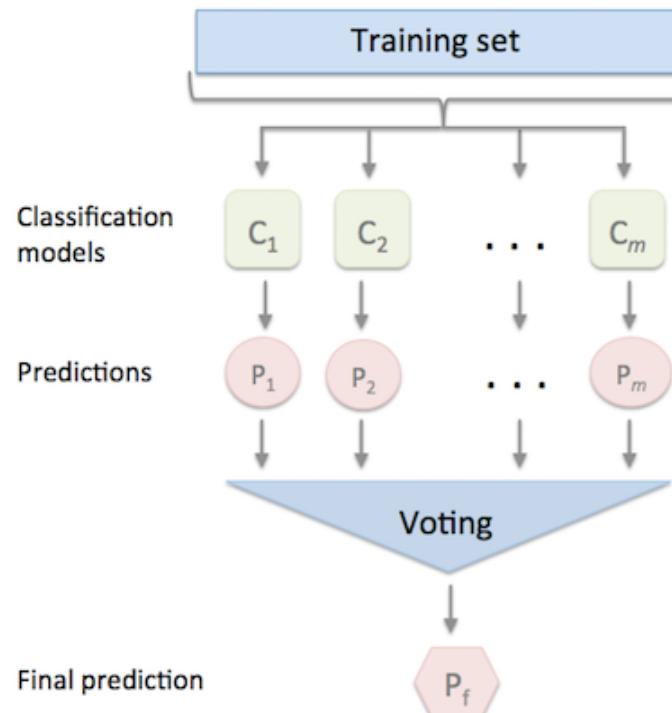
Bagging

Boosting

Stacking

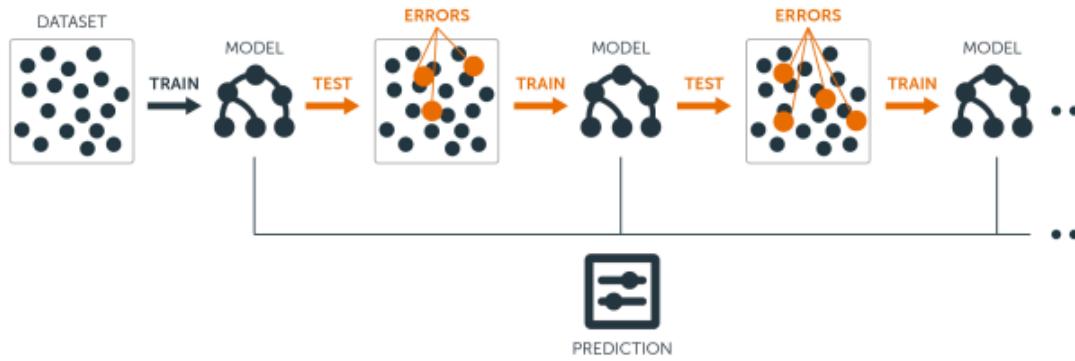
**Bagging** means to build multiple models from different subsamples of the training set and/or with different methods. The results are sent to an (weighted) voting operation that produces the final result.

Source: [http://rasbt.github.io/mlxtend/user\\_guide/classifier/EnsembleVoteClassifier/](http://rasbt.github.io/mlxtend/user_guide/classifier/EnsembleVoteClassifier/)



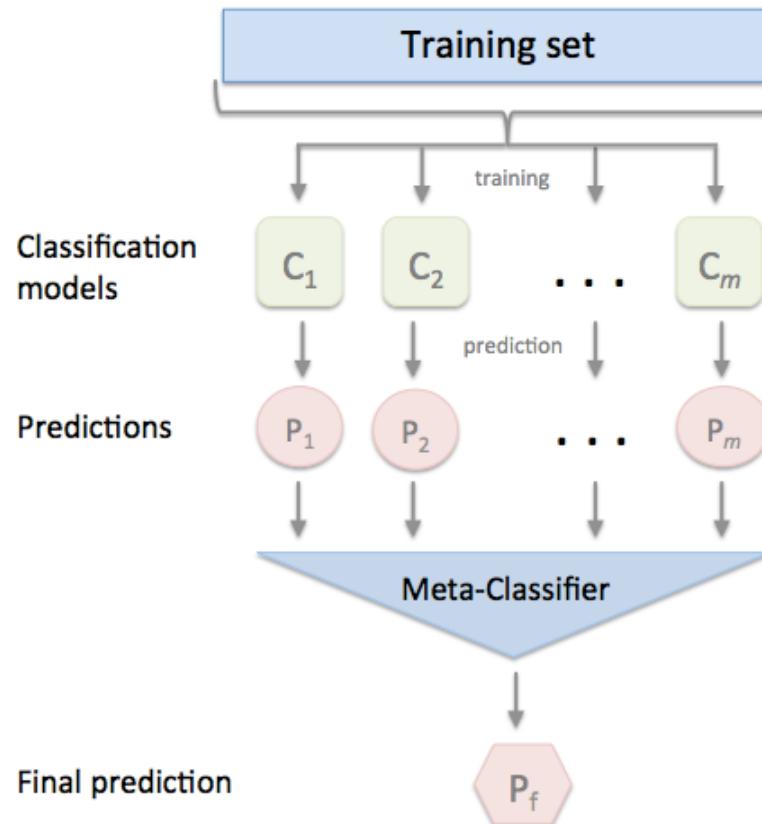
**Boosting** involves sequentially building an ensemble by training each new model in ~~influence~~ variants exist, emphasize the training instances that previous models mispredict . Different variants exist, mostly based on tree methods. In general, any method can be used. This involves the usage of different methods at the different iterations when building the sequence of models.

Source: <https://blog.bigml.com/2017/03/14/introduction-to-boosted-trees/>



**Stacking** means to build multiple models (typically of differing types) and a supervisor model that learns how to best combine the predictions of the primary models. The inputs of the supervisor model (meta-classifier) are the outputs of the other models:

Source: [http://rasbt.github.io/mlxtend/user\\_guide/classifier/StackingClassifier/](http://rasbt.github.io/mlxtend/user_guide/classifier/StackingClassifier/)



# Types of Ensembles

## Type 1:

consists of only a few models

each is a strong model

like few professional experts

risk of diverging opinions

risk of experts being biased to their experiences

## Type 2:

consists of many models

each is a weak model as a principle

based on the idea of the wisdom of the masses

Random Forest and Gradient Boosted Trees are examples



## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

#### 4.4.2 Regression

##### 4.4.2.1 OLS

##### 4.4.2.2 Ridge Regression

##### 4.4.2.3 Support Vector Regression

##### 4.4.2.4 Neural Networks

##### 4.4.2.5 Decision Trees

##### 4.4.2.6 K-Nearest Neighbors

# Predicting using Regression Methods

Groundtruth

↓

| House No. | Square Footage | Bedrooms | Age | School Rating | Final Price |
|-----------|----------------|----------|-----|---------------|-------------|
| H1        | 1000           | 4        | 3   | 2             | \$100,000   |
| H2        | 800            | 3        | 1   | 4             | \$90,000    |
| H3        | 1200           | 5        | 3   | 5             | \$125,000   |
| H4        | 600            | 2        | 5   | 1             | \$60,000    |
| H5        | 1500           | 6        | 3   | 3             | \$150,000   |

Example: Predicting House Prices

Function: Price = f(SquareFootage, Bedrooms, Age, SchoolRating)

Source:

<http://www.sclgsummit.org/uploads/presentation/8934b2d0be055a2261f5d0320f5b59bb.pdf>

## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

#### 4.4.2 Regression

##### 4.4.2.1 OLS

##### 4.4.2.2 Ridge Regression

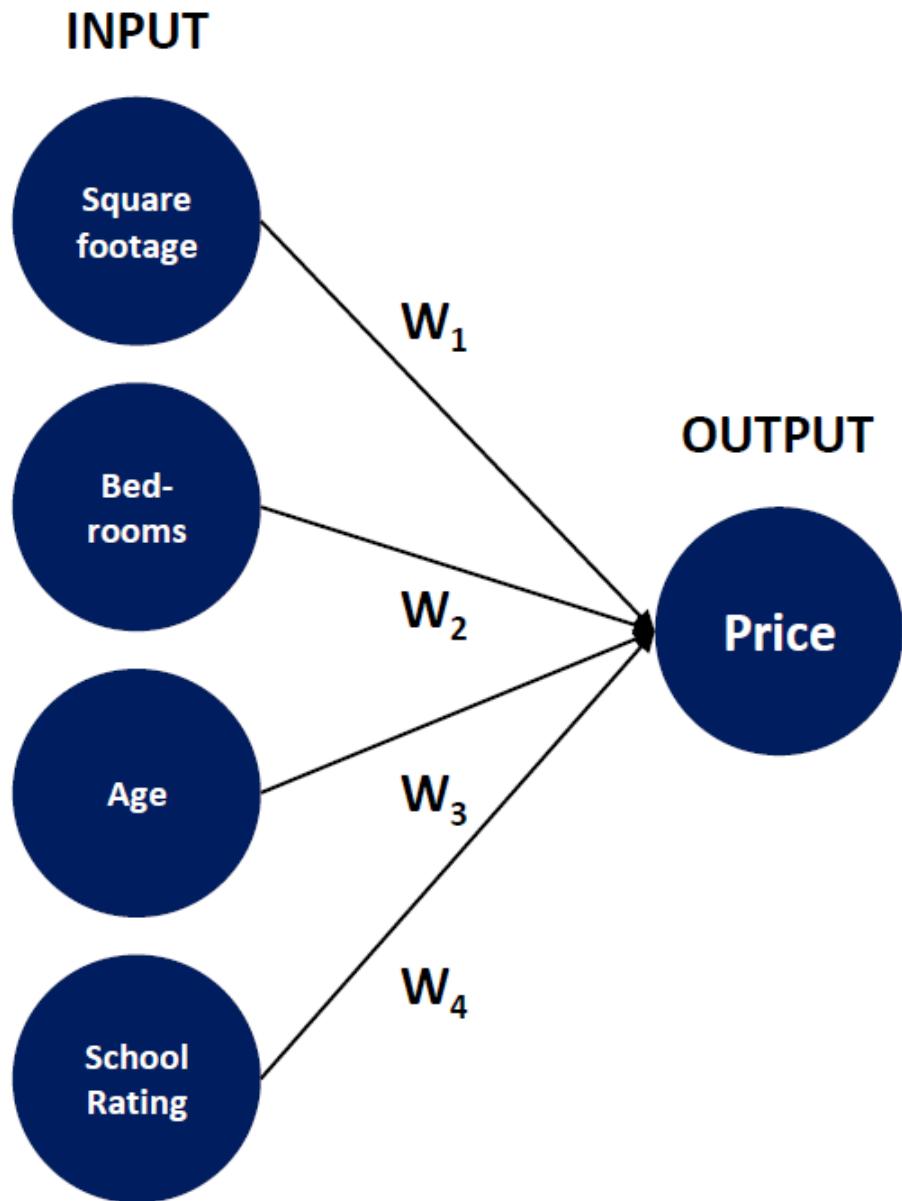
##### 4.4.2.3 Support Vector Regression

##### 4.4.2.4 Neural Networks

##### 4.4.2.5 Decision Trees

##### 4.4.2.6 K-Nearest Neighbors

# Traditional OLS Regression Approach



Function:

$$\text{Price} = \beta_0 + \beta_1 * \text{SquareFootage} + \beta_2 * \text{Bedrooms} + \beta_3 * \text{Age} + \beta_4 * \text{SchoolRating}$$

# Ordinary Least Squares Regression

## Measuring the Quality of Fit (I)

Measuring the quality of fit means to measure how well the predictions of a model match the observed data.

A commonly-used measure is the Mean Absolute Error (MAE) which can be calculated for the training and the test set

A variant is the Mean Absolute Percentage Error (MAPE) which expresses the error in percent

While MAE and MAPE are easily interpretable, using the absolute value of the error often is not as desirable as squaring this difference. Depending on how you want your model to treat outliers, or extreme values, in your data, you may want to bring more attention to these outliers or downplay them.

Consequently, the most used measure in regression is the Mean Squared Error (MSE) or its variant the Root Mean Squared Error (RMSE), which is the square root of the MSE.

# Measuring the Quality of Fit (II)

## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

#### 4.4.2 Regression

##### 4.4.2.1 OLS

##### 4.4.2.2 Ridge Regression

##### 4.4.2.3 Support Vector Regression

##### 4.4.2.4 Neural Networks

##### 4.4.2.5 Decision Trees

##### 4.4.2.6 K-Nearest Neighbors

# Ridge Regression

Complexity can be measured as the size of the set of possible outputs for a given set of inputs.

In this example the interval 0 to  $x^*$  represents the set of possible inputs. Function  $h_0$  has the lowest complexity because there is just one output independent of the inputs.  $h_2$  has the highest complexity because here the set of possible outputs is the biggest one.

# Complexity und Generalisation

## Mean Squared Error

# Different Complexities

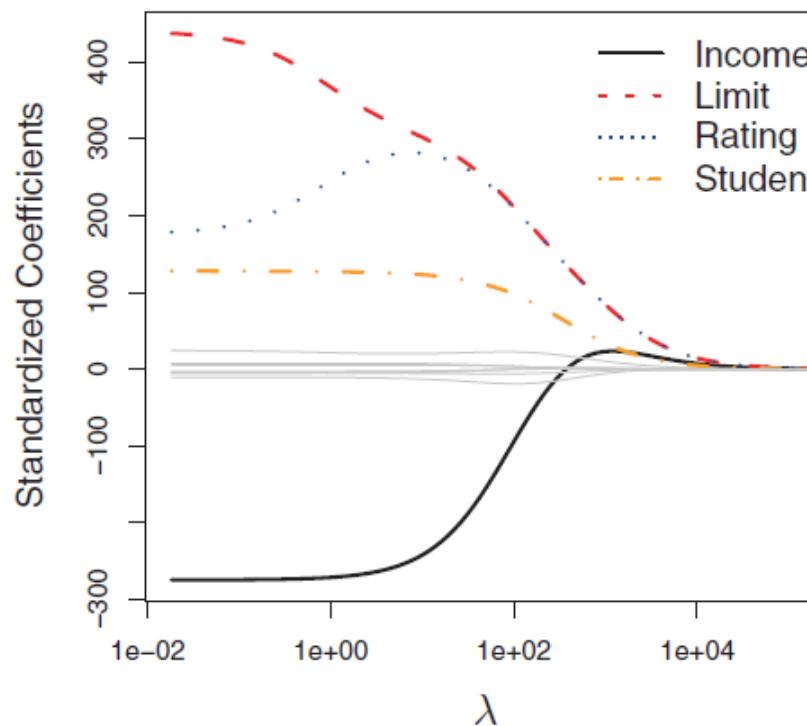
$\lambda \rightarrow \infty$  : Lowest Complexity

the ridge regression coefficients are equal to zero. For every input, the result is  $\beta_0$ .

$\lambda = 0$  : Relative High Complexity (linear Model)

the penalty term has no effect, and ridge regression will produce the least squares estimates.

Example:



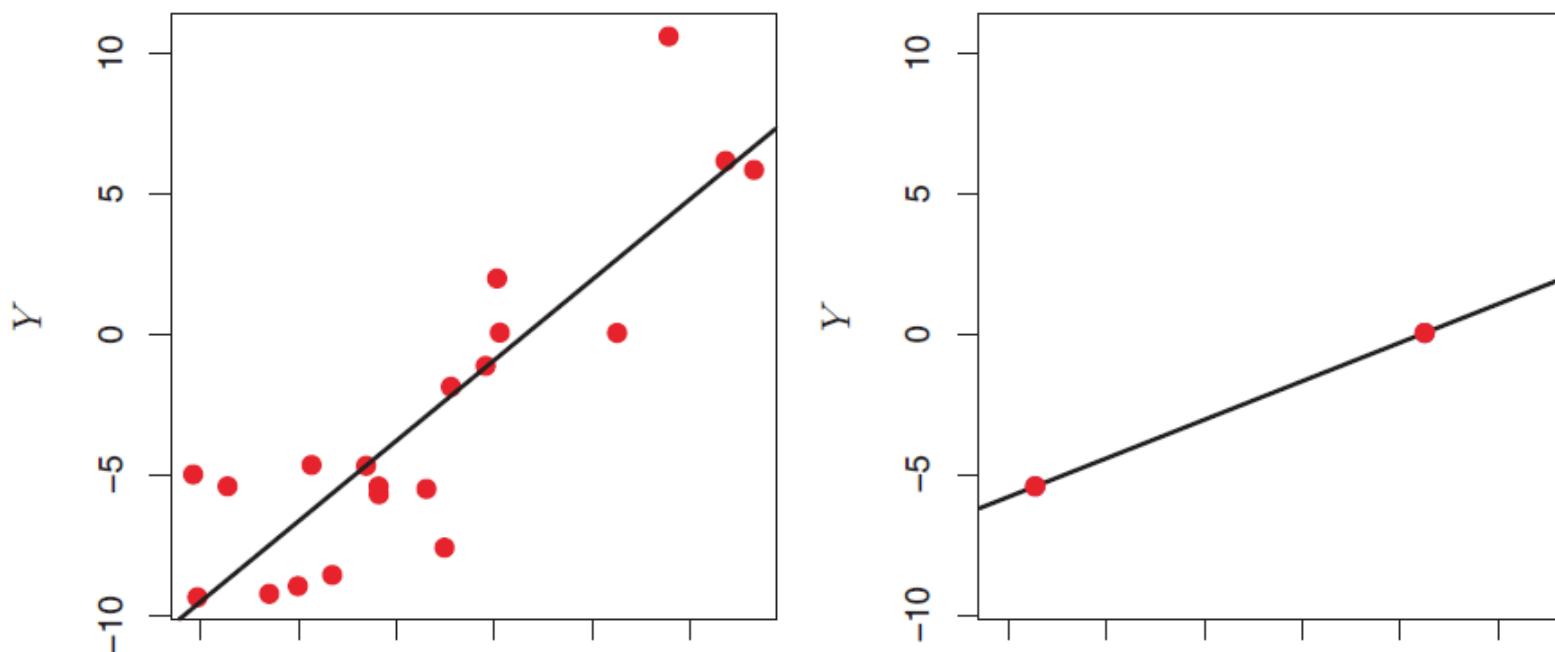
Source:

James et al. (2013): An Introduction to Statistical Learning with R Applications, p. 215f.

## Handling High-Dimensionality (I)

OLS is not suitable for high-dimensional data. Especially when the number of features  $p$  is as large as, or larger than, the number of observations, OLS cannot be applied. Regardless of whether or not there truly is a relationship between the features and the response, OLS will yield a set of coefficient estimates that result in a perfect fit to the data, such that the residuals are zero.

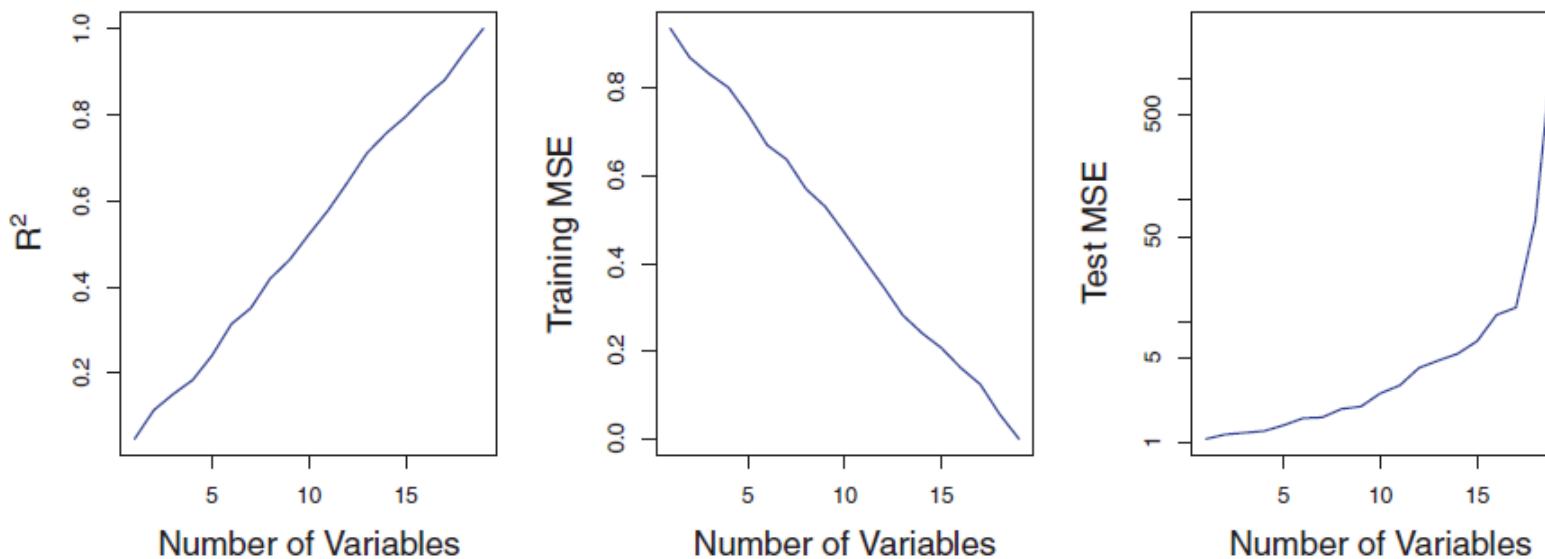
The figure shows two cases. When there are 20 observations,  $n > p$  and the OLS line does not perfectly fit the data. When there are only two observations, then regardless of the values of those observations, the regression line will fit the data exactly. This is problematic because this perfect fit will almost certainly lead to overfitting of the data.



## Handling High-Dimensionality (II)

The figure illustrates the risk of applying OLS when the number of features  $p$  is large. The model R<sup>2</sup> increases to 1 as the number of features increases, and the training set MSE decreases to 0. At the same time, the MSE on a test set becomes extremely large as the number of features increases.

In contrast, methods like ridge regression are particularly useful for performing regression in the high-dimensional setting. Essentially, these approaches avoid overfitting by using a less flexible fitting approach than least squares.



Source: James et al. (2013): An Introduction to Statistical Learning with R Applications, p. 240f.

## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

#### 4.4.2 Regression

##### 4.4.2.1 OLS

##### 4.4.2.2 Ridge Regression

##### 4.4.2.3 Support Vector Regression

##### 4.4.2.4 Neural Networks

##### 4.4.2.5 Decision Trees

##### 4.4.2.6 K-Nearest Neighbors

# Support Vector Regression

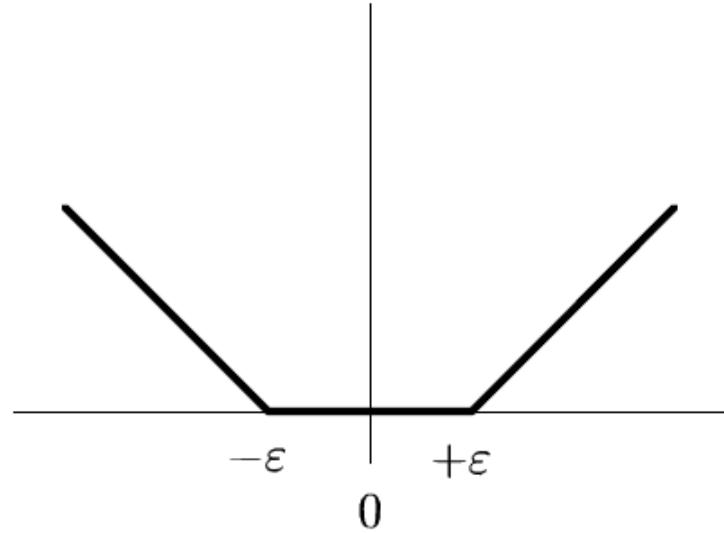
The Goal is to find a robust model with a high generalization ability.

SVR regards two sources of Robustness:

1. Eliminating Noise
2. Handling Complexity

## In insensitive Loss Function (I)

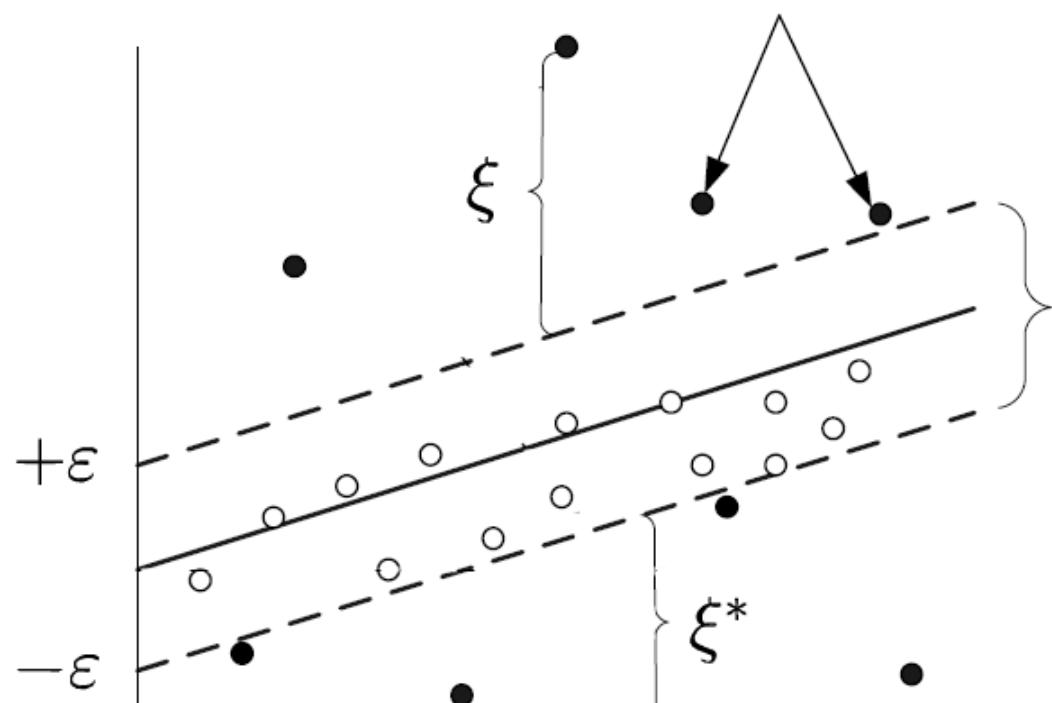
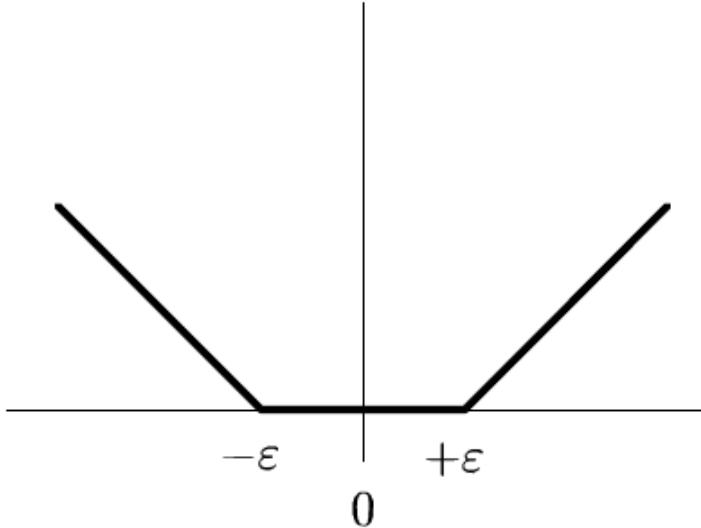
### -insensitive Loss



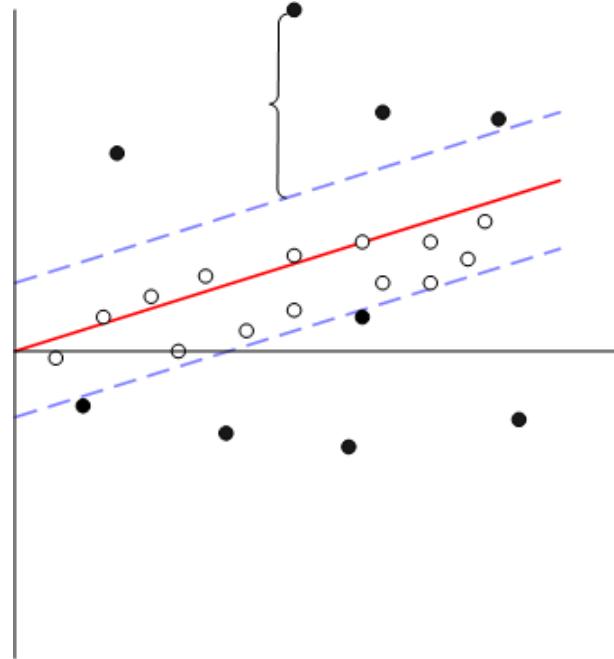
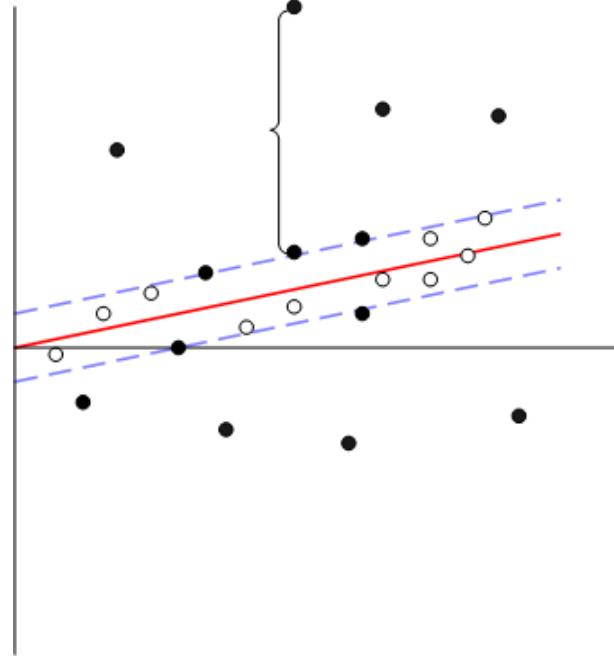
does not penalize acceptable deviations (defined by )

## In insensitive Loss Function (II)

Using the  $\epsilon$ -insensitive loss function, only those data objects are considered in the estimation, which have a distance greater than  $\epsilon$  from the regression function:



# Support Vector Regression (I)



# Estimating the SVR (Linear Case)

Nonlinearity (I)

The linear case :

The nonlinear \_\_\_ case :

# Nonlinearity (II)

## Kernel Functions (I)

Kernel Functions are used to project n-dimensional input to m-dimensional input, where m is higher than n:

Any point  $x$  in the original space is mapped into the higher dimensional space. For reason of efficiency, the mapping is not performed in real but instead embedded in the model building process via the kernel function:

Instead of  $\beta_0 + \beta \cdot x = y$  the following is used  $\beta_0 + \beta \cdot F(x) = y$

The main idea to use a kernel is: A linear regression curve in higher dimensions becomes a non-linear regression curve in lower dimensions.

# Estimating the SVR (Nonlinear Case)

## Kernel Functions (II)

A frequently used kernel function is the Polynomial Kernel Function:

where  $x$  and  $z$  are vector points in any fixed dimensional space and  $n$  is the order of the kernel.

In the case of order equal to 2, we get:

Source: <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1000173>

## Kernel Functions (III)

A nother frequently \_\_\_\_ used \_\_\_\_ kernel \_\_\_\_ function \_\_\_\_ is \_\_\_\_ the \_\_\_\_ Function \_\_\_\_ ( \_\_\_\_ RBF):

It maps the data according a Gaussian function where Sigma ( s ) is a streching factor.

Different Sigmas

= Euclidean distance between x and z

Source: <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1000173>

## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

#### 4.4.2 Regression

##### 4.4.2.1 OLS

##### 4.4.2.2 Ridge Regression

##### 4.4.2.3 Support Vector Regression

##### 4.4.2.4 Neural Networks

##### 4.4.2.5 Decision Trees

##### 4.4.2.6 K-Nearest Neighbors

## Using Neural Network for Regression

Artificial neural networks are often used for classification because of the relationship to logistic regression. Neural networks typically use a logistic activation function and output values from 0 to 1 like logistic regression.

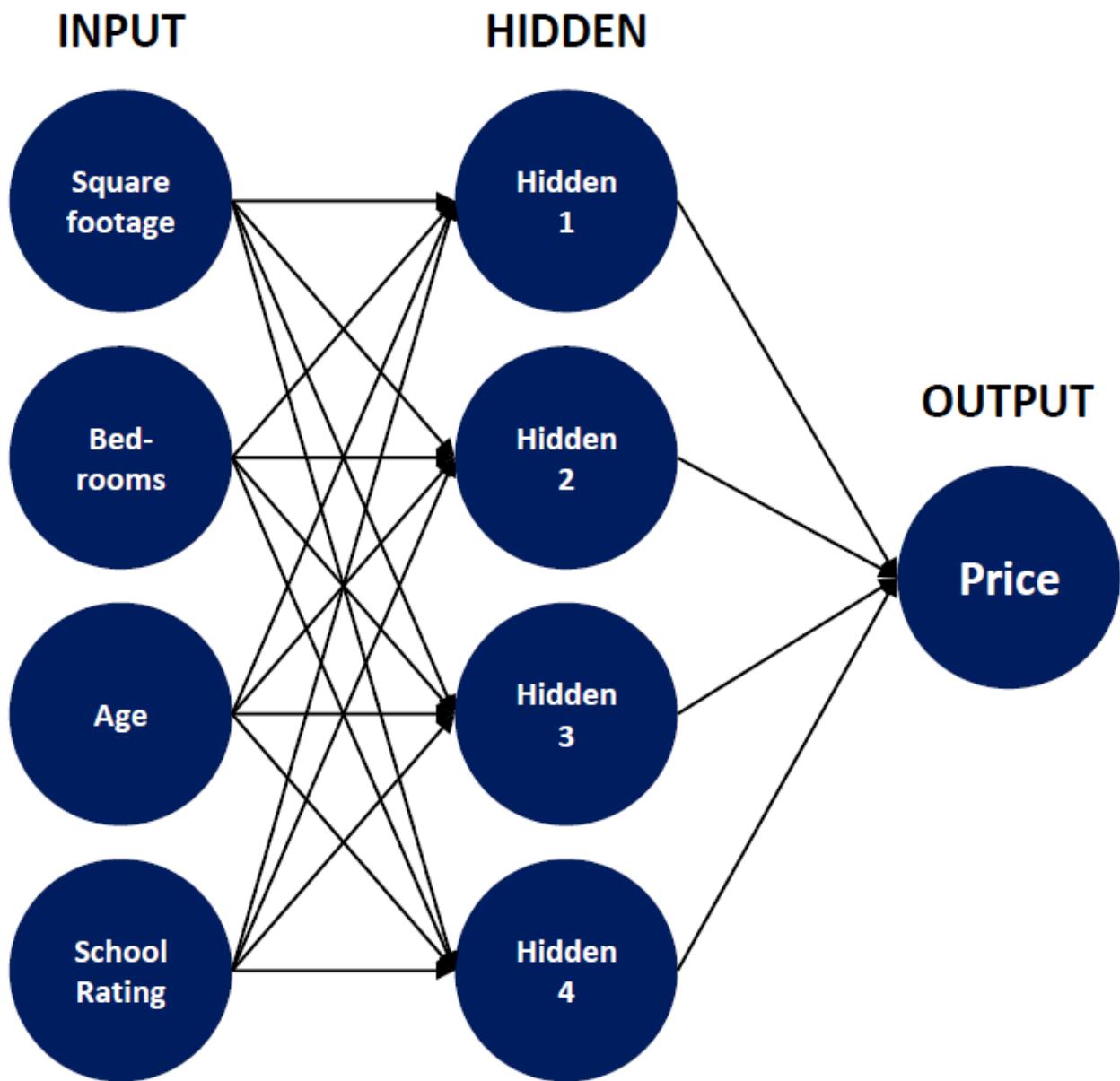
But the continuous output of a net must not be interpreted as a probability, so neural networks can be used too for regression, to model complex and non-linear relationships.

The Singlelayer Perceptron corresponds to a linear regression while a Multilayer Perceptron is able to approximate nearly any function regard-less of the complexity and nonlinearity.

Because of the high complexity of the MLP, the models are usually very sensitive and have a tendency to overfitting.

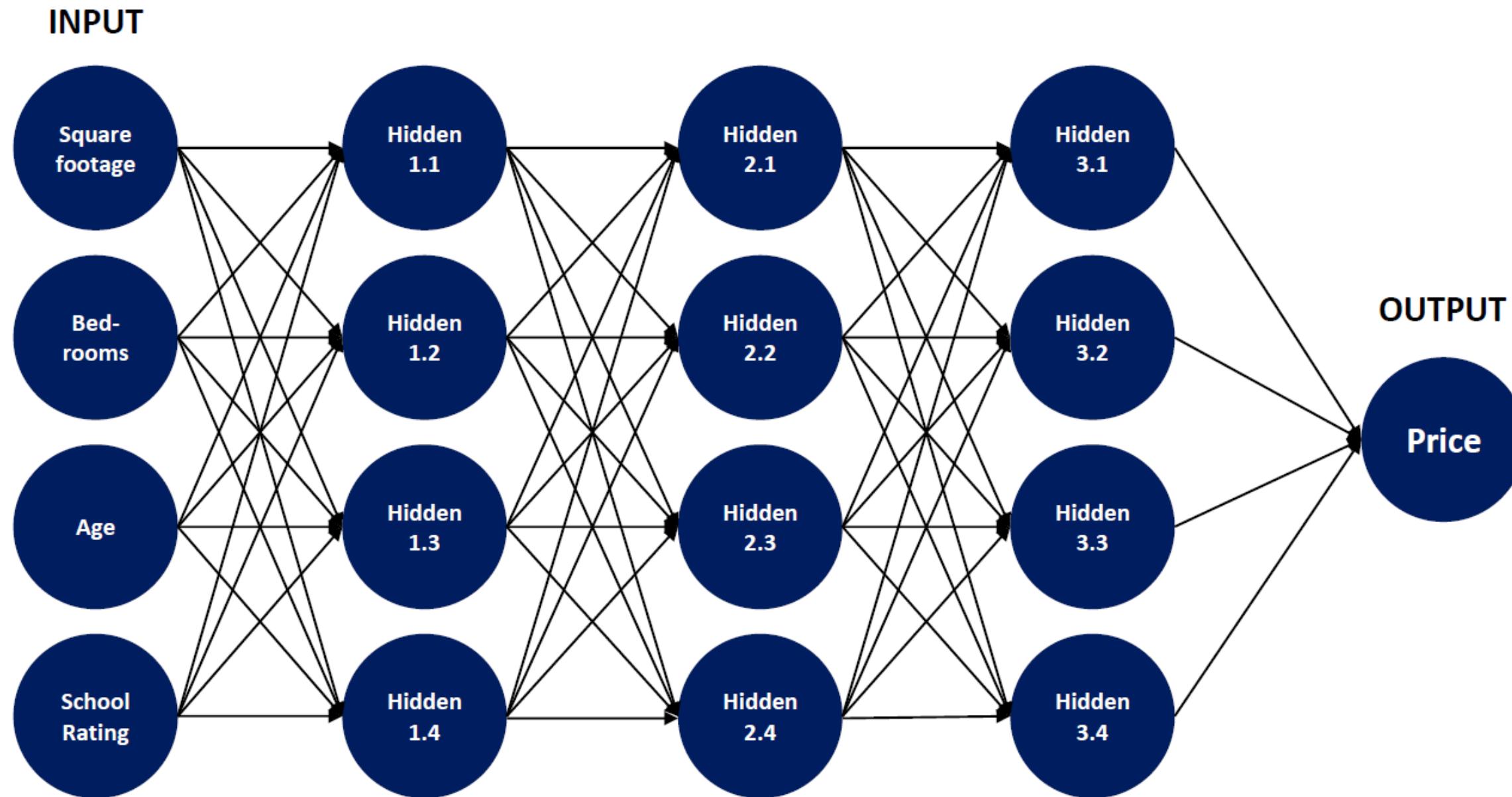
There exist regularization methods, which make the networks better at generalizing beyond the training data.(see <http://neuralnetworksanddeeplearning.com/chap3.html>)

# Neural Network (Multilayer Perceptron)



Source:

<http://www.sclgsummit.org/uploads/presentation/8934b2d0be055a2261f5d0320f5b59bb.pdf>



Source:

<http://www.sclgsummit.org/uploads/presentation/8934b2d0be055a2261f5d0320f5b59bb.pdf>

## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

#### 4.4.2 Regression

##### 4.4.2.1 OLS

##### 4.4.2.2 Ridge Regression

##### 4.4.2.3 Support Vector Regression

##### 4.4.2.4 Neural Networks

##### 4.4.2.5 Decision Trees

##### 4.4.2.6 K-Nearest Neighbors

## Introductory Example

### Decision Tree for Predicting Fuel Consumption of Cars(in Miles-per-Gallon )

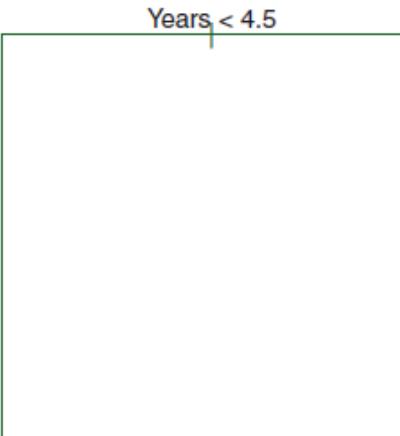
## Regression Trees

Some of the tree approaches can be used for regression too. They can be used for nonlinear multiple regression. The output must be numerical.

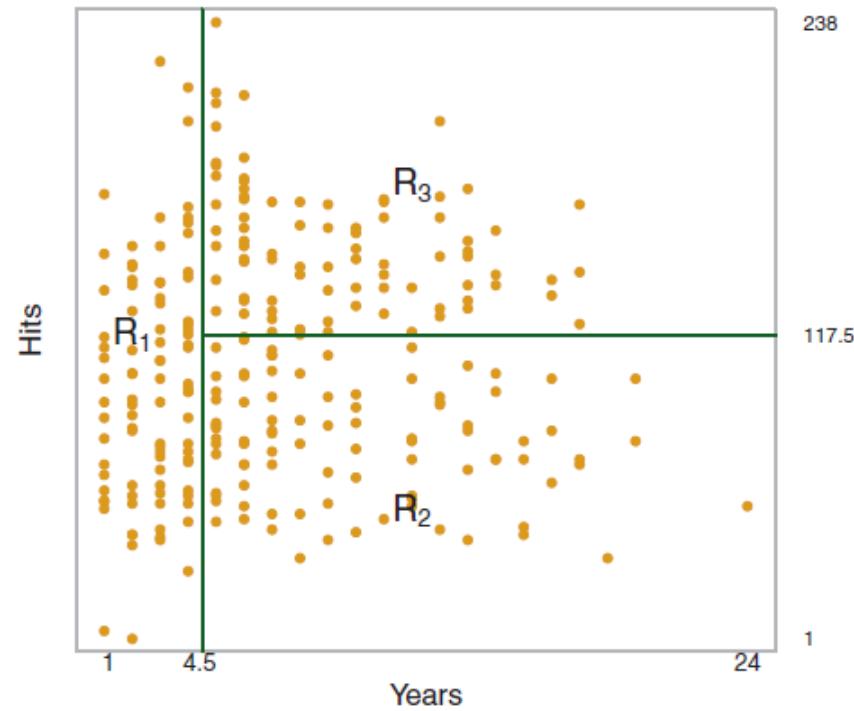
The figure shows a regression tree for predicting the salary of a baseball player, based on the number of years that he has played in the major leagues and the number of hits that he made in the previous year.

The predicted salary is given by the mean value of the salaries in the corresponding leaf, e.g. for the players in the data set with  $\text{Years} < 4.5$ , the mean (log-scaled) salary is 5.11, and so we make a prediction of  $e^{5.11}$  thousands of dollars, i.e. \$165,670, for these players.

Players with  $\text{Years} \geq 4.5$  are assigned to the right branch, and then that group is further subdivided by Hits. The predicted salaries for the resulting two groups are  $1,000 \cdot e^{6.00} = \$403,428$  and  $1,000 \cdot e^{6.74} = \$845,346$ .



# Constructing a Regression Tree (I)



Source: James et al. (2013): An Introduction to Statistical Learning with R Applications, p. 305f.

# Constructing a Regression Tree (II)

## Random Forests for Regression

Due to the usage of means as predictors a regression tree usually simplifies the true relationship between the inputs and the output. The advantage over traditional statistical methods is, that it can give valuable insights about which variables are important and where. But the prediction ability is poor compared to other regression approaches.

A much better prediction quality can be achieved with the creation of an ensemble of trees, use them for prediction and averaging their results. This is done, when applying the Random Forests approach to a regression task.

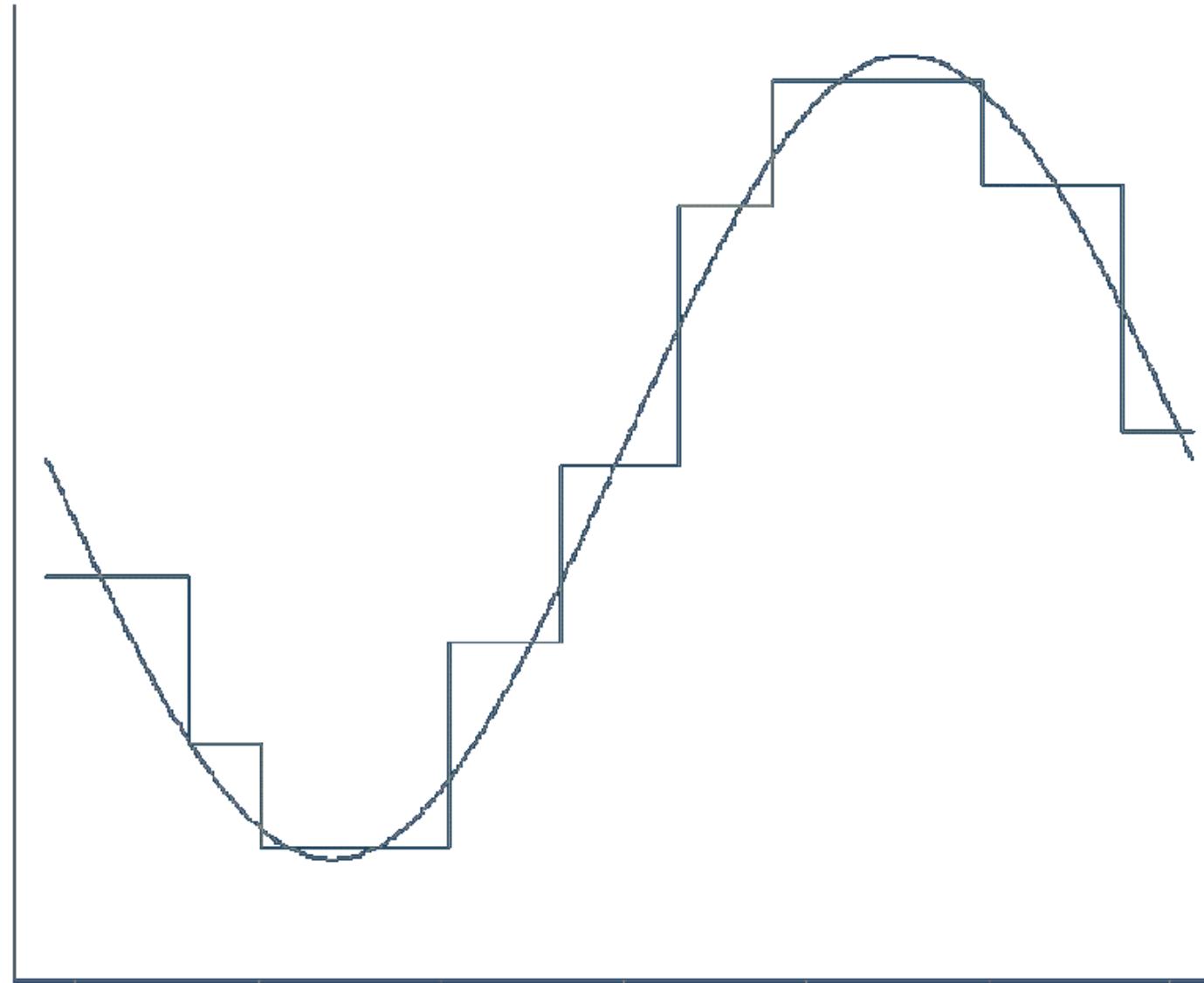
Regression Forests are an ensemble of different regression trees and are used for nonlinear multiple regression. The principle is the same as in classification, except that the output is not the result of a voting but instead of an averaging process.

The disadvantage of Random Forests is that the analysis, which aggregates over the results of many bootstrap trees, does not produce a single, easily interpretable tree diagram.

# Comparing the Fitting Ability of one vs. many Regression Trees

Single Regression Tree

Average of 100 Regression Trees



## Limitations of Tree Methods in Regression

When applied to regression problems, tree methods have the limitation that they cannot exceed the range of values of the target variable used in training. The reason for this lies in their design principle, how the leaves of the trees are created.

Thus, Random Forests may perform poorly when the target data is out of the range of the original training data, e.g. in the case of data with persistent trends. A solution may be a frequent re-training in this case.

An important strength of Random Forests is that they are able to perform still well in the case of missing data. According to their construction principle, not every tree is using the same features.

If there is any missing value for a feature during the application there usually are enough trees remaining that do not use this feature to produce accurate predictions.

## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

#### 4.4.2 Regression

##### 4.4.2.1 OLS

##### 4.4.2.2 Ridge Regression

##### 4.4.2.3 Support Vector Regression

##### 4.4.2.4 Neural Networks

##### 4.4.2.5 Decision Trees

##### 4.4.2.6 K-Nearest Neighbors

## k-Nearest Neighbors for Regression

k-Nearest Neighbors cannot only be used for classification but also for regression. The only difference in regression is that the prediction is not the result of a majority vote but of an averaging process.

A simple implementation of KNN regression is to calculate the average of the numerical target of the k-nearest neighbors. Another approach uses an inverse distance weighted average of the K-nearest neighbors. KNN regression uses the same distance functions as KNN classification.

Example:

## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

#### 4.4.2 Regression

#### 4.4.3 Segmentation

##### 4.4.3.1 K-Means

##### 4.4.3.2 Hierarchical Cluster Analysis

## Introductory Example

Assume you are a wholesale distributor and each row of your dataset corresponds to a customer showing the following attributes:

1) FRESH: annual spending on fresh products (Continuous); 2) MILK: annual spending on milk products (Continuous); 3) GROCERY: annual spending on grocery products (Continuous); 4) FROZEN: annual spending on frozen products (Continuous) 5) DETERGENTS\_PAPER: annual spending on detergents and paper products (Continuous) 6) DELICATESSEN: annual spending on delicatessen products (Continuous); 7) CHANNEL: customers buying channel (Nominal) 8) REGION: customers region (Nominal)

Your goal is to segment the users. That means finding similar types of users and bunching them together.

Why would you want to do this?

You might want to give different users different experiences. Marketing often does this; for example, to offer toner to people who are known to own printers.

You might have a model that works better for specific groups. Or you might have different models for different groups.

Cluster analysis is a type of multivariate statistical analysis. It is used to group data into separate clusters. The main objective of clustering is to find similarities between data objects, and then group similar objects together to assist in understanding relationships that might exist among them. Cluster analysis is based on a mathematical formulation of a measure of similarity.

There are different types of cluster analysis methods:

## Clustering Methods

## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

#### 4.4.2 Regression

#### 4.4.3 Segmentation

##### 4.4.3.1 K-Means

##### 4.4.2.2 Hierarchical Cluster Analysis

## Partitioning Cluster Methods

The partitioning cluster methods divide the data into a predetermined number of clusters.

The most popular technique is the K-Means algorithm.

Given a set of observations ( $x_1, x_2, \dots, x_n$ ), where each observation is a  $m$ -dimensional real vector,  $k$ -means clustering aims to partition the  $n$  observations into ( $k \leq n$ ) segments  $S = \{S_1, S_2, \dots, S_k\}$  so as to minimize the within-cluster sum of squares (WCSS).

The objective is to find

where  $\bar{x}_i$  is the mean of points in  $S_i$ .

Procedure of K-Means:

**Step 1: Randomly partition the data objects into k clusters.**

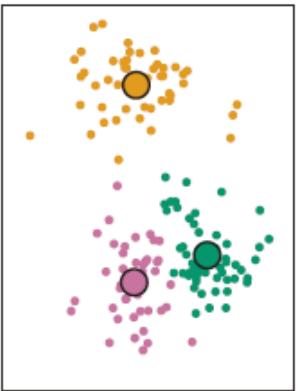
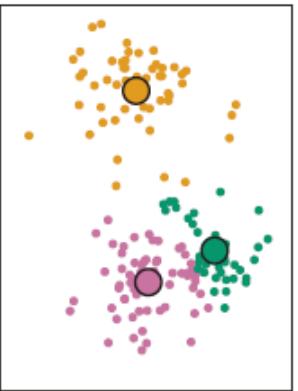
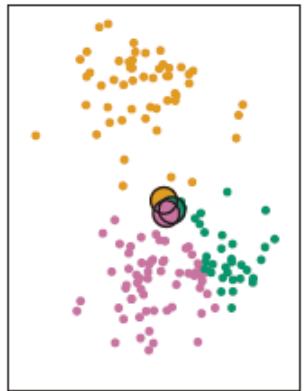
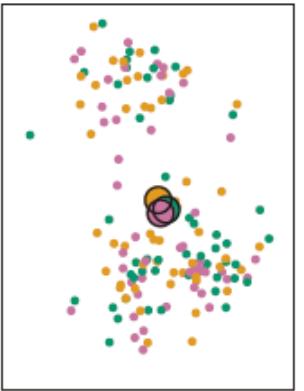
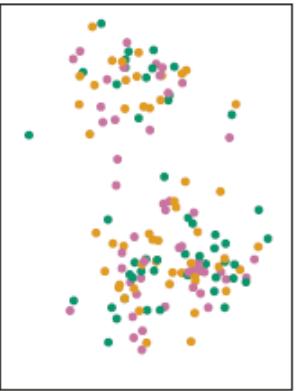
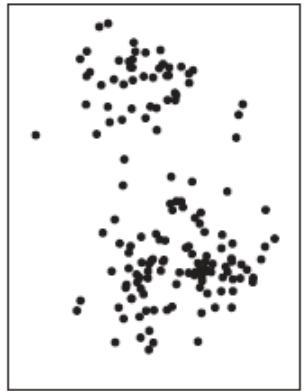
**Step 2: Calculate the cluster centroids.**

**Step 3: Calculate the distance from every data point to all centroids**

**Step 4: If a data point is closest to its own centroid, leave it where it \_\_\_\_ is. If the data point is not closest to its own centroid, assign \_\_\_\_ it to the cluster with the closest centroid.**

**Step 5: Repeat the step 2 to 4 until a complete pass through of all \_\_\_\_ the data points results in no data point changing from one \_\_\_\_ cluster to another.\_\_\_\_**

# Example of a K-Means Cluster Analysis



**Between cluster variance:**

**Within cluster variance:**

Finding the Optimal Number of Clusters (I)

The aim of the cluster analysis is the segmentation of objects into clusters, which are preferably homogeneous in it selves and heterogeneous to each other. The less variance exists within the clusters and the more variance exists between the clusters, the better is the number of clusters.

**Total variance:**

**Accumulated variance within the k clusters:**

**This results in the variance between the clusters:**

**with  $n$  = number of objects**

\_\_  $m$  = number of attributes \_\_

\_\_\_\_  $n_k$  \_\_ = number of objects in cluster  $k$  \_\_

\_\_\_\_  $c_k$  \_\_ = cluster  $k$  \_\_

## Finding the Optimal Number of Clusters (II)

If you put  $V$  in \_\_ on the ordinate and the number of cluster  $k$  on the abscissa, it often results in a curve with one or several kinks. At the point where exists the (first) significant kink, you can find the optimal number of clusters:\_\_

**Total variance  $V_{tot}$**

**Between \_\_ \_\_ cluster variance  $V_{betw}$**

**Within cluster variance  $V_{in}$**

**Number of clusters**

## Finding the Optimal Number of Clusters (III)

Instead of visually identifying the optimal cluster number, we can calculate the *distances from the points on the elbow curve to a straight line linking the first and the last point on the curve.* The cluster number with the largest distance is then chosen as the one with the strongest kink.

### Number of clusters

## 4 Predictive Analytics

### 4.1 Subject of Predictive Analytics

### 4.2 The Analytics Process

### 4.3 Data Preparation

### 4.4 Methods, Algorithms and Applications

#### 4.4.1 Classification

#### 4.4.2 Regression

#### 4.4.3 Segmentation

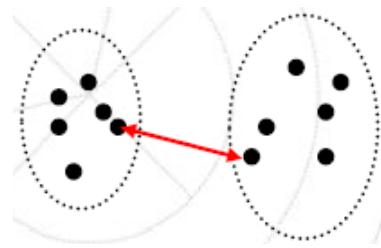
##### 4.4.3.1 K-Means

##### 4.4.2.2 Hierarchical Cluster Analysis

- There are two types of hierarchical cluster methods:
  - Agglomerative hierarchical clustering is a bottom-up clustering method. It starts with every single data object in a single cluster. Then, in each iteration, it agglomerates (merges) the closest pair of clusters by satisfying some similarity criteria, until all of the data is in one cluster.
  - Divisive hierarchical clustering is a top-down clustering method. It works in a similar way to agglomerative clustering but in the opposite direction. This method starts with a single cluster containing all data objects, and then successively splits resulting clusters until only clusters of individual data objects remain.

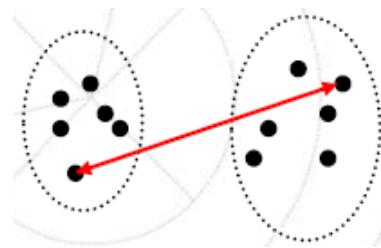
# Process of the Hierarchical Cluster Analysis

## Measuring Similarity between Clusters (I)

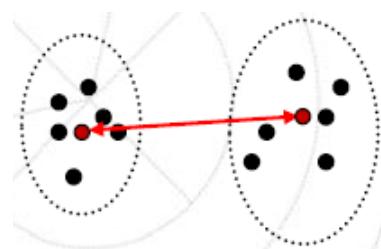


Distance between two clusters is the distance between the closest points:

Complete Linkage:



Distance between two clusters is the distance between the farthest pair of points:

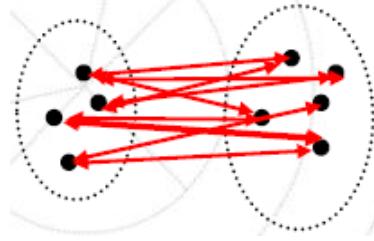


Distance between two clusters  $i$  and  $j$  is the distance between their centroids :

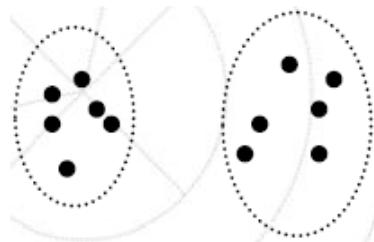
## Measuring Similarity between Clusters (II)

Average Linkage:

Distance between clusters is the average distance between the cluster points:



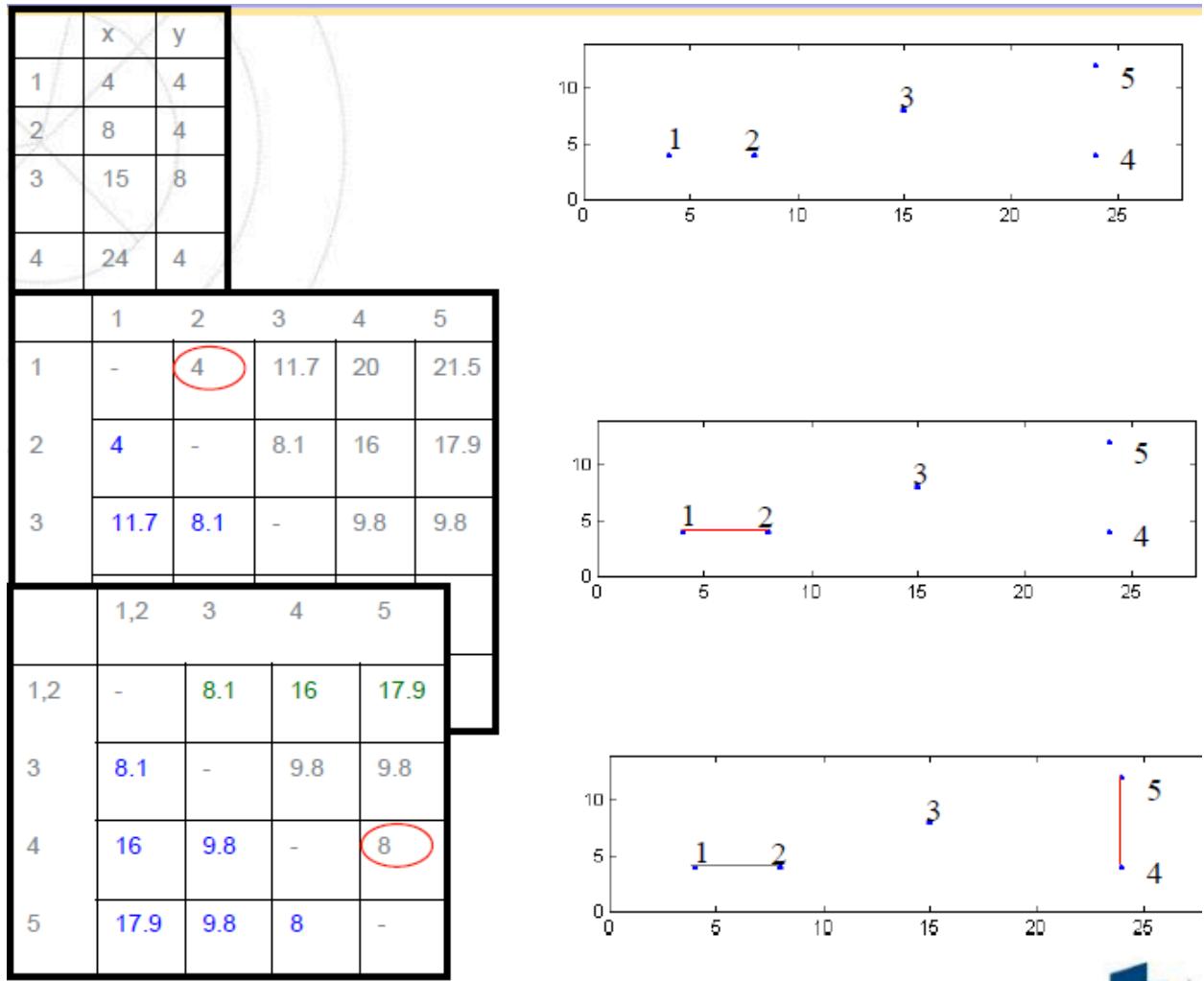
Ward's Method / Minimum Variance Method (only Agglomerative):



Ward's minimum variance criterion minimizes the total within-cluster variance. At each step the pair of clusters is merged that leads to minimum increase in total within-cluster variance after merging. This can be calculated as the square of the distance between cluster means divided by the sum of the reciprocals of the number of observations in each cluster:

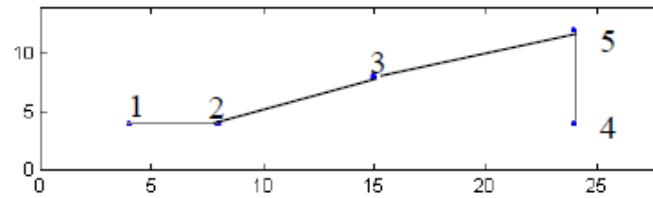
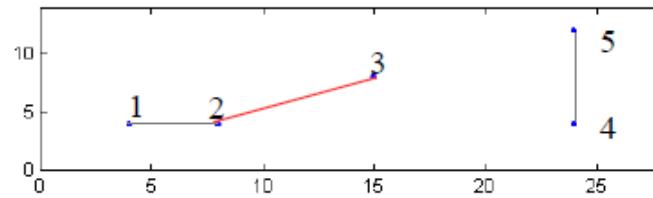
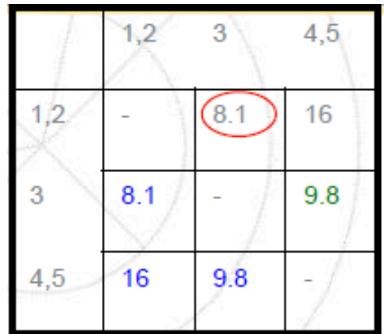
For a comparison of the methods see: Ferreira, L.; Hitchcock, D. B. (2009): A Comparison of Hierarchical Methods for Clustering Functional Data,

# Single Linkage Example (I)

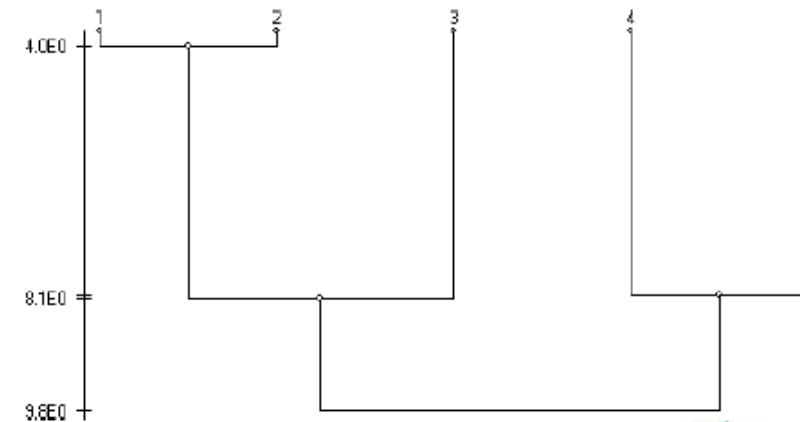


Source: Fred, Ana: Unsupervised Learning, Universidade Técnica de Lisboa

## Single Linkage Example (II)



Dendrogram



Source: Fred, Ana: Unsupervised Learning, Universidade Técnica de Lisboa

A dendrogram is a tree diagram frequently used to illustrate the arrangement of the clusters produced by hierarchical clustering. The y-axis represents the value of this distance metric (e.g. euclidean distance) between the clusters.

In a dendrogram the widths of the horizontal lines give an impression about the dissimilarity of the merging object. Thus, a good cluster number might be at a point from where the width of the following horizontal lines is significantly smaller in length. The red line in the graph below shows such a point:

Counting the points that cut this line might be a good answer for the number of clusters the data can have. It is the number 6 in this case.

