

# Statistics for data science

Tuesday 22<sup>nd</sup> November, 2016

1 Introduction - Statistics, what is it good for?

2 Data

3 Relations between data elements

4 Conclusion

Israellëch challenge

A data science perspective

**Data** → **Knowledge** 

Israellëch challenge

A data science perspective

Data → Knowledge

Israellëch challenge

A data science perspective

### $Data \longrightarrow Knowledge \\$

- there is randomness in the data
  - a telephonic survey



#### A data science perspective

### $Data \longrightarrow Knowledge \\$

- ▶ there is randomness in the data
  - a telephonic survey
- not enough data to achieve certainty
  - choosing the best restaurant when there are very few reviews



#### A data science perspective

### **Data** → **Knowledge**

- ▶ there is randomness in the data
  - a telephonic survey
- not enough data to achieve certainty
  - choosing the best restaurant when there are very few reviews
- the knowledge required is of a statistical nature
  - What are my odds when betting on a soccer match



Statistics is an important component of a data scientist's toolbox



Statistics is an important component of a data scientist's toolbox

"raw" statistics - "getting a feel of the data" data exploration



Statistics is an important component of a data scientist's toolbox

- "raw" statistics "getting a feel of the data" data exploration
- many data science tools are of a statistical nature



Statistics is an important component of a data scientist's toolbox

- "raw" statistics "getting a feel of the data" data exploration
- many data science tools are of a statistical nature
- a knowledge of statistics is essential
  - choose the right tool
  - correctly interpret results

1 Introduction - Statistics, what is it good for?

2 Data

3 Relations between data elements

4 Conclusion



#### **Structured Data**

data is organized in well defined fields ...



#### Structured Data

data easily fits in a relational database - tables

Examples?



#### **Structured Data**

Name	Height (H)	$\textbf{Weight} \ (W)$
Prince Humperdinck	6'	166lb
Inigo Montoya	5'3"	142lb
Princess Buttercup	5'4"	136lb



#### **Unstructured Data**

data is not organized in well defined fields...



#### **Unstructured Data**

data does not easily fit in a relational database - tables what's easily? examples?



#### **Unstructured Data**











Google	(Section 1 Decrees)
To Design the same	- 555
	2000
Good	e News
- acogs	
- HEEV-	
THE PERSON NAMED IN COLUMN TWO IS NOT THE PERSON NAMED IN COLUMN TWO IS NAM	
Total of New Arrange 10" and Street	COLUMN TWO IS NOT THE OWNER.
27	

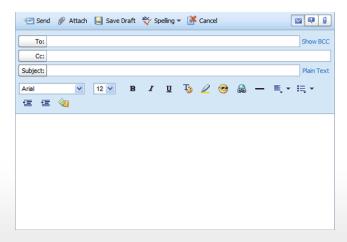


#### Semi-structured Data

Data is partially structured and partially unstructured examples?



#### Semi-structured Data





► Numerical - discrete/continuous



- Numerical discrete/continuous
  - Height
  - Temprature
  - Number of children



- Numerical discrete/continuous
  - Height
  - Temprature
  - Number of children
- Categorical



- Numerical discrete/continuous
  - Height
  - Temprature
  - Number of children
- Categorical
  - Gender
  - Nationality



- Numerical discrete/continuous
  - Height
  - Temprature
  - Number of children
- Categorical
  - Gender
  - Nationality
- Ordinal



- Numerical discrete/continuous
  - Height
  - Temprature
  - Number of children
- Categorical
  - Gender
  - Nationality
- Ordinal
  - 1st, 2nd, 3rd...
  - low, medium, high



- Arbitrary (string)
  - natural language quotes, reviews, phrases
  - hashes



- Arbitrary (string)
  - natural language quotes, reviews, phrases
  - hashes
- Time series
  - Navigation data (Waze)
  - Stock market data

1 Introduction - Statistics, what is it good for?

2 Data

3 Relations between data elements

4 Conclusion

## A "practical" example



#### Physical measurements of N test subjects

Name	Height (H)	$\textbf{Weight} \ (W)$
Prince Humperdinck	6'	166lb
Inigo Montoya	5'3"	142lb
Princess Buttercup	5'4"	136lb

:

## A "practical" example



### Physical measurements of N test subjects

- What can you say about height?
- What can you say about weight?



#### Physical measurements of N test subjects

mean

$$\mu_X$$

$$E[X] = \sum_{x} x \cdot Pr(x)$$

$$E[X] = \int_{x} x \cdot f(x) dx$$

variance

$$\sigma_X^2$$

$$V(X) = E\left[\left(X - E[X]\right)^{2}\right]$$

## A "practical" example



#### Physical measurements of N test subjects

- What can you say about height?
- What can you say about weight?
- ▶ How would you assess the connection between the two?



#### Physical measurements of N test subjects

Covariance

$$cov(X, Y) = E\left[\left(X - E[X]\right)\left(Y - E[Y]\right)\right]$$



#### Physical measurements of N test subjects

Covariance

$$cov(X, Y) = E\left[\left(X - E[X]\right)\left(Y - E[Y]\right)\right]$$

Covariance of height and weight

$$cov(H, W) = 6.27$$

## A "practical" example



### Physical measurements of N test subjects

News from the citadel...

they've invented the metric system

## A "practical" example



### Physical measurements of N test subjects

News from the citadel...

Name	Height $(H)$	$\textbf{Weight}\ (W)$
Prince Humperdinck	183cm	75.3kg
Inigo Montoya	160cm	64.41kg
Princess Buttercup	162.5cm	61.69kg

:

What happens to the connection between height and weight?

### A "practical" example



### Physical measurements of N test subjects

Covariance of height and weight

$$cov(H, W) = 87.1125$$

### Discussion



▶ What's the problem with using covariance to assess the connection between height and weight?



- ▶ What's the problem with using covariance to assess the connection between height and weight?
- We need a measure that is invariant to a change of units (scale)



- ▶ What's the problem with using covariance to assess the connection between height and weight?
- We need a measure that is invariant to a change of units (scale)
- ▶ How about normalizing the covariance?

## Israellëch challenge

 $\rho_{X,Y}$ 

informally: normalized covariance

$$corr(X,Y) = \frac{cov(X,Y)}{\sigma_X \cdot \sigma_Y}$$

$$= \frac{E\left[\left(X - E[X]\right)\left(Y - E[Y]\right)\right]}{\sqrt{E\left[\left(X - E[X]\right)^2\right]}\sqrt{E\left[\left(Y - E[Y]\right)^2\right]}}$$

## Israellëch challenge

 $\rho_{X,Y}$ 

- ▶  $-1 \le \rho_{X,Y} \le 1$
- ▶ high direct correlation  $\rho_{X,Y} \approx 1$
- ▶ high inverse correlation  $\rho_{X,Y} \approx -1$

#### Pearson Correlation



 $\rho_{X,Y}$ 

#### **Examples**

- direct correlation
  - lacksquare X number of days it |rains| in the morning
  - lacksquare Y number of days sidewalks are  $\lfloor wet \rfloor$  at noon
- inverse correlation
  - lacksquare X number of days it |rains| in the morning
  - lacksquare Y number of days sidewalks are |dry| at noon



statistical dependence...



### Does a lack of correlation imply independence

$$\rho_{X,Y} = 0 \stackrel{?}{\Rightarrow} X \perp\!\!\!\perp Y$$



### Does a lack of correlation imply independence

$$\rho_{X,Y} = 0 \stackrel{?}{\Rightarrow} X \perp \!\!\!\perp Y$$

 $\mathbf{No}$ , Correlation = Linear dependence

$$X, X^2$$



### Does a lack of correlation imply independence

$$\rho_{X,Y} = 0 \stackrel{?}{\Rightarrow} X \perp \!\!\!\perp Y$$

 $\mathbf{No}$ , Correlation = Linear dependence

$$X, X^2$$



#### Does correlation imply dependence

$$\rho_{X,Y} \in \{-1,1\} \stackrel{?}{\Rightarrow} X \not\perp Y$$



#### Does correlation imply dependence

$$\rho_{X,Y} \in \{-1,1\} \stackrel{?}{\Rightarrow} X \not\perp Y$$

Yes



#### Does correlation imply dependence

$$\rho_{X,Y} \in \{-1,1\} \stackrel{?}{\Rightarrow} X \not\perp \!\!\! \perp Y$$

**Yes**, Linear dependence = a private case of dependence

## A word about samples



statistical operators - mean, variance,...



statistical operators - mean, variance,...

$$E\left[X\right] = \sum_{x} x \cdot Pr\left(x\right)$$

## A word about samples



statistical operators - mean, variance,...

$$E\left[X\right] = \sum_{x} x \cdot Pr\left(x\right)$$

how are these computed from a sample?



statistical operators - mean, variance,...

$$E[X] = \sum_{x} x \cdot Pr(x)$$

- how are these computed from a sample?
- sample operators sample mean, sample variance, ...
- all based on averaging

## A word about samples



$$X = \{x_1, x_2, \dots, x_n\}$$

#### Sample Mean

$$\overline{X} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

## A word about samples



$$X = \{x_1, x_2, \dots, x_n\}$$

#### Sample Variance

$$\sigma_X^2 = \frac{1}{n} \sum_{i=1}^n \left( x_i - \overline{X} \right)^2$$



$$X = \{x_1, x_2, \dots, x_n\}$$

#### Sample covariance

$$cov(X,Y) = \frac{1}{n} \sum_{i=1}^{n} (x_i - \overline{X}) (y_i - \overline{Y})$$
$$= \frac{1}{n} (X - \overline{X})^T (Y - \overline{Y})$$

1 Introduction - Statistics, what is it good for?

2 Data

3 Relations between data elements

4 Conclusion



- statistics is important
- knowledge is statistical
- data comes in diverse forms
- simple relations between data elements covariance, correlation