

# Statistical methods in NLP

## Introduction



**UNIVERSITY OF  
GOTHENBURG**

Richard Johansson

January 19, 2016

today

- ▶ course matters
- ▶ analysing numerical data with Python
- ▶ basic notions of probability
- ▶ simulating random events in Python



















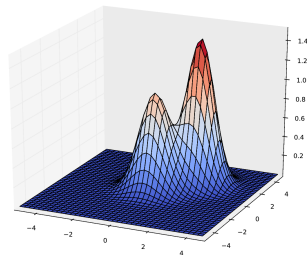






# some useful Python libraries we'll use in this course

- ▶ **SciPy**: a Python library for statistics and math in general
  - ▶ <http://www.scipy.org/>
- ▶ **NumPy**: efficient mathematical functions
  - ▶ <http://www.numpy.org/>
- ▶ **matplotlib**: using Python to draw diagrams
  - ▶ <http://matplotlib.org/>



## an example dataset

```
177 67 m
177 77 m
175 87 m
154 47 f
157 56 f
152 53 f
165 49 f
165 66 f
165 63 m
162 54 f
167 62 m
167 79 f
167 58 f
165 61 f

# read the data
data = []
with open('bodies.txt') as f:
    for l in f:
        h, w, s = l.split()
        data.append( (int(h), int(w), s) )

heights = [ h for h, _, _ in data ]
weights = [ w for _, w, _ in data ]

m_heights = [ h for h, _, s in data if s == 'm' ]
m_weights = [ w for _, w, s in data if s == 'm' ]
f_heights = [ h for h, _, s in data if s == 'f' ]
f_weights = [ w for _, w, s in data if s == 'f' ]
```

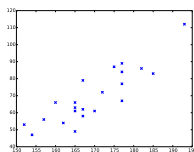
# plotting the data

- ▶ import the plotting library

```
from matplotlib import pyplot as plt
```

- ▶ plot a height/weight plot, each point as an 'x'

```
plt.plot(heights, weights, 'x')
```



- ▶ plot height/weight plot by gender

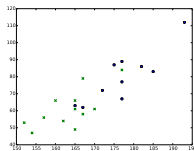
```
plt.plot(m_heights, m_weights, 'o', f_heights, f_weights, 'x')
```

- ▶ save the plot to a file

```
plt.savefig('myplot.pdf')
```

- ▶ alternatively, draw the plot on the screen

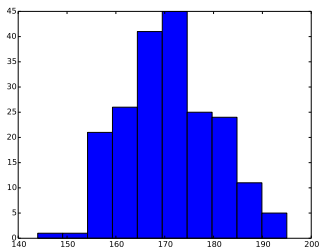
```
plt.show()
```



# plotting histograms

- ▶ a **histogram** is a diagram that shows how the data points are distributed
- ▶ the x axis shows “bins”, e.g. 165–170 cm, and the y axis shows the number of data points in that bin
- ▶ here's how we draw a histogram with matplotlib:

```
plt.hist(heights, bins=10)
```

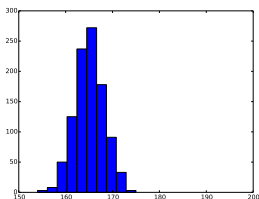




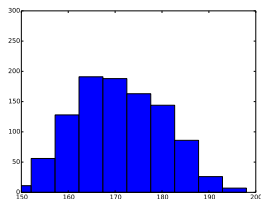


## example

- ▶ low variance: data concentrated near the mean
  - ▶ in the extreme case: all values are identical
- ▶ high variance: data spread out



$$\sigma = 2.9$$



$$\sigma = 9.1$$

## variance and standard deviation in SciPy

- ▶ variance:

```
var_height = scipy.var(heights)
```

- ▶ standard deviation:

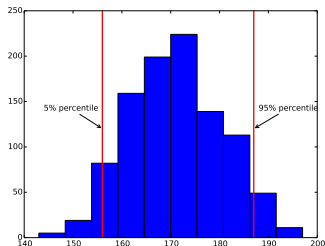
```
std_m_weight = scipy.std(m_weights)
```



# percentiles

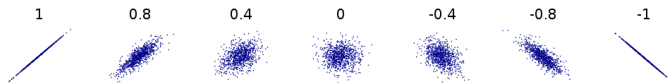
- ▶ how tall are the shortest 5% of the people in the dataset?
  - ▶ formally: what is the  $x$  such that 5% of the data is less than  $x$ ?
- ▶ this number is called the **5% percentile**
- ▶ in Python:

```
p5 = numpy.percentile(heights, 5)
```



## relations between two variables: correlation

- ▶ the **correlation coefficient** or the **Pearson  $r$**  measures how close the data is to a linear relationship
- ▶ it is a number that ranges between -1 and +1
  - ▶ example [Wikipedia]:

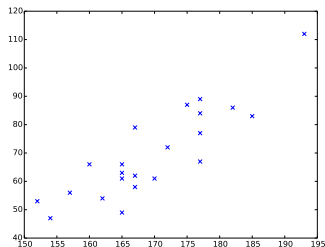


- ▶ it is defined

$$r(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sigma(x)\sigma(y)}$$

# correlation example

- ▶ for the height–weight data,  $r = 0.87$



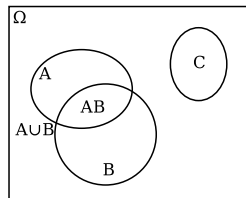
- ▶ with Python:  
`correlation = scipy.stats.pearsonr(heights, weights)[0]`





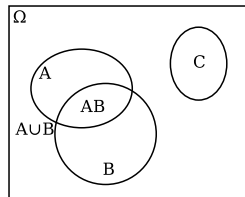
## some formal notation: events

- ▶ the theory of probability is built on the theory of sets
  - ▶ so we can draw Venn diagrams to make the notions more intuitive
- ▶  $\Omega$  is the **sample set**: the set of all possible situations
- ▶ an **event**  $A$  is a subset of  $\Omega$
- ▶ the **union** event  $A \cup B$  means that either  $A$  or  $B$  has happened
- ▶ the **joint** event  $AB$  (also written  $A \cap B$ ) means that  $A$  and  $B$  have both happened
- ▶ two events  $A$  and  $C$  that can't happen at the same time (that is, the intersection is empty) are called **disjoint**



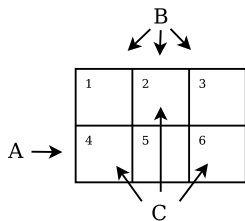
# the mathematical definition: the Kolmogorov axioms

- ▶ the probability  $P(A)$  is a number such that
  - ▶  $0 \leq P(A) \leq 1$  for every event  $A$
  - ▶  $P(\Omega) = 1$
  - ▶  $P(A \cup B) = P(A) + P(B)$  if  $A$  and  $B$  are disjoint
- ▶ in the illustrations,  $P(A)$  intuitively corresponds to the area covered by  $A$  in the Venn diagram



## example: Dice rolling

- ▶  $A = \text{"rolling a 4"}; P(A) = ?$
- ▶  $B = \text{"rolling 3 or lower"}; P(B) = ?$
- ▶  $C = \text{"rolling an even number"}; P(C) = ?$
- ▶  $P(A \cup B) = P(A) + P(B)?$
- ▶  $P(A \cup C) = P(A) + P(C)?$
- ▶  $P(\text{rolling 1, 2, 3, 4, 5, or 6}) = ?$



## some consequences

- ▶  $A' = \text{"everything but } A" = \Omega \setminus A$

$$P(A') = 1 - P(A)$$

$$P(\emptyset) = 0$$

1	2	3
	$A'$	
$^4$ $A$	5	6

- ▶  $A = \text{"rolling a 4"}; P(A) = ?$
- ▶  $A' = \text{"not rolling 4"}; P(A') = ?$
- ▶  $P(\text{rolling neither 1, 2, 3, 4, 5, 6}) = ?$

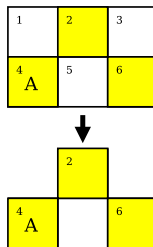
# joint and conditional probabilities

- ▶ the probability of both  $A$  and  $B$  happening is called the **joint probability**, written  $P(AB)$  or  $P(A, B)$
- ▶ definition: if  $P(B) \neq 0$ , then

$$P(A|B) = \frac{P(AB)}{P(B)}$$

is referred to as the **conditional probability of  $A$  given  $B$**

- ▶ intuitively in the Venn diagram: zoom in on  $B$ 
  - ▶ “what is the probability of a 4 if we know it’s an even number?”
- ▶ this is something we’ve already used in language models, taggers, etc





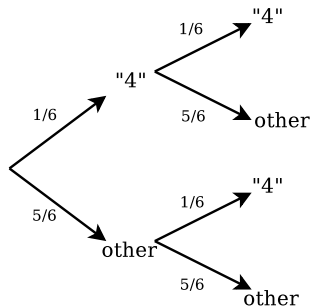




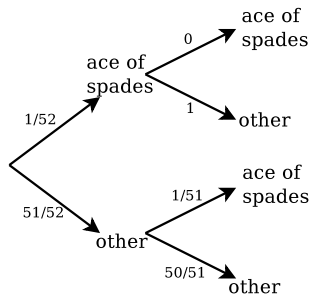




# drawing tree diagrams (1)



## drawing tree diagrams (2)





## a brief note on “random” numbers

- ▶ **pseudorandom numbers**: generating “random” numbers in a computer using a deterministic process
  - ▶ usually fast
  - ▶ we use a starting point called the **seed**
  - ▶ if we use the same seed, we'll get the same sequence
    - ▶ good for replicable experiments
  - ▶ might be a security risk in some situations
- ▶ hardware random numbers
  - ▶ sample noise from hardware devices
  - ▶ Linux: `/dev/random`





## simulating random events in Python

- ▶ `random.random` and `random.randint` can be used to simulate random events
- ▶ example of generating random words with different probabilities

```
import random

def random_word():
    r = random.random()
    if r < 0.4:
        return 'the'
    if r < 0.7:
        return 'and'
    if r < 0.9:
        return 'in'
    return 'is'

random_words = [ random_word() for _ in range(20) ]

print(random_words)
```

