

Statistics (II)

SLIDES BY:

JIANNAN WANG

<https://www.cs.sfu.ca/~jnwang/>

Outline

Correlation Analysis

- Big Picture
- How to do correlation analysis

Hypothesis Testing

- Big Picture
- A/B Testing

Outline

Correlation Analysis

- Big Picture
- How to do correlation analysis

Hypothesis Testing

- Big Picture
- A/B Testing

Correlation Analysis

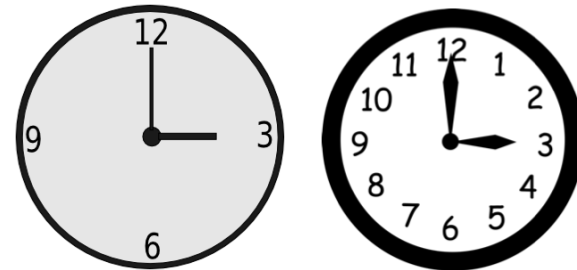
Correlation

- It is a measure of relationship between two variables

Why is correlation analysis useful?

- For understanding data better
- For making predictions better

Correlation \neq Causation



Case Study:

How to do correlation analysis

Height and weight are correlated

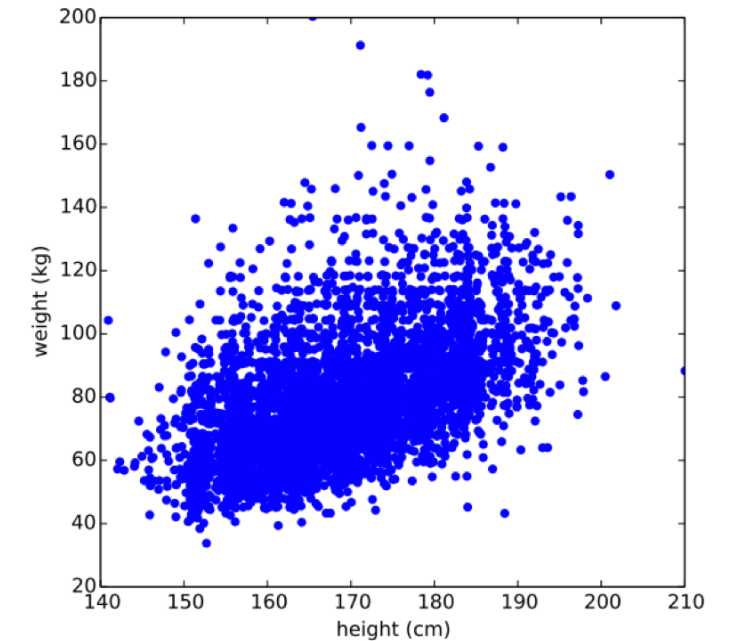
1	height	weight	age	male
2	151.765	47.8256065	63	1
3	139.7	36.4858065	63	0
4	136.525	31.864838	65	0
5	156.845	53.0419145	41	1
6	145.415	41.276872	51	0
7	163.83	62.992589	35	1
8	149.225	38.2434755	32	0

Source: *Think Stats -- Exploratory Data Analysis in Python*

Idea 1. Visualization

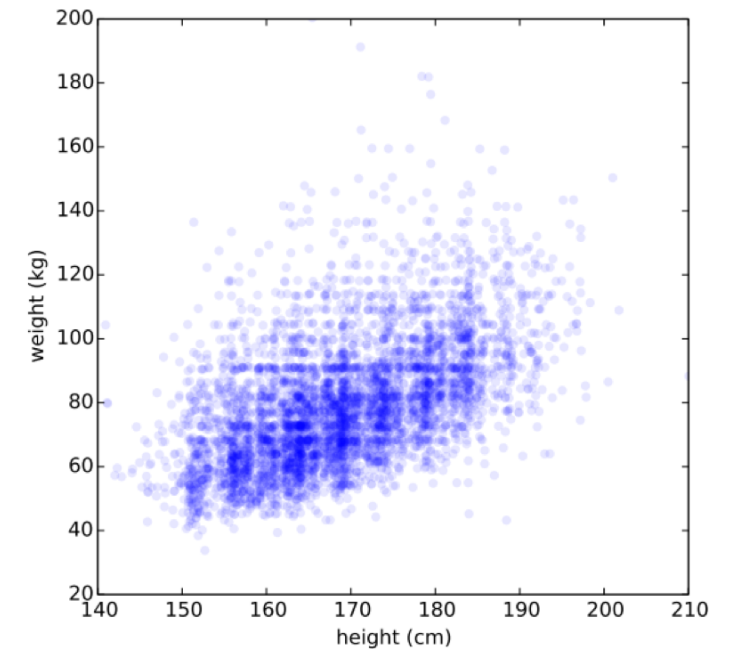
Scatter Plot

1	height	weight	age	male
2	151.765	47.8256065	63	1
3	139.7	36.4858065	63	0
4	136.525	31.864838	65	0
5	156.845	53.0419145	41	1
6	145.415	41.276872	51	0
7	163.83	62.992589	35	1
8	149.225	38.2434755	32	0



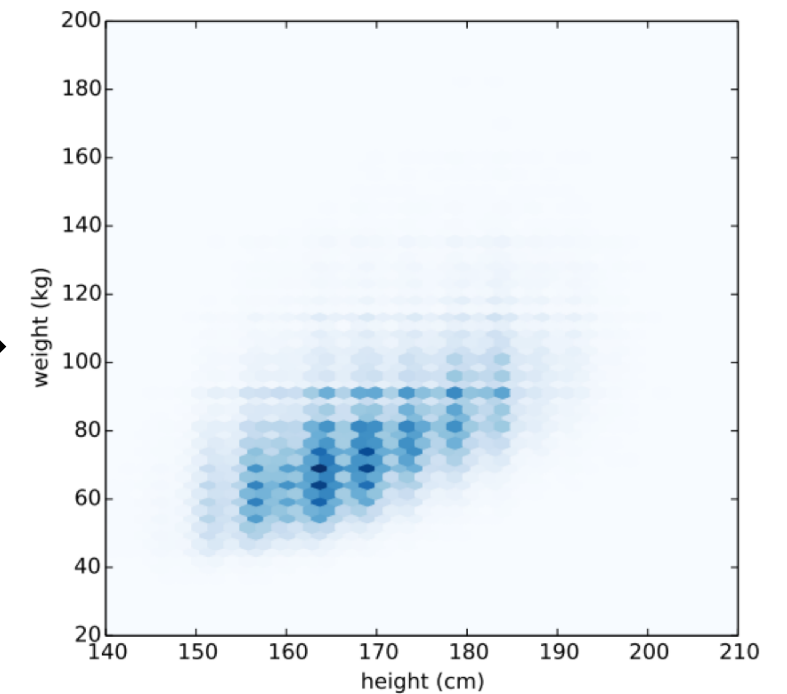
Scatter Plot (with transparency)

1	height	weight	age	male
2	151.765	47.8256065	63	1
3	139.7	36.4858065	63	0
4	136.525	31.864838	65	0
5	156.845	53.0419145	41	1
6	145.415	41.276872	51	0
7	163.83	62.992589	35	1
8	149.225	38.2434755	32	0



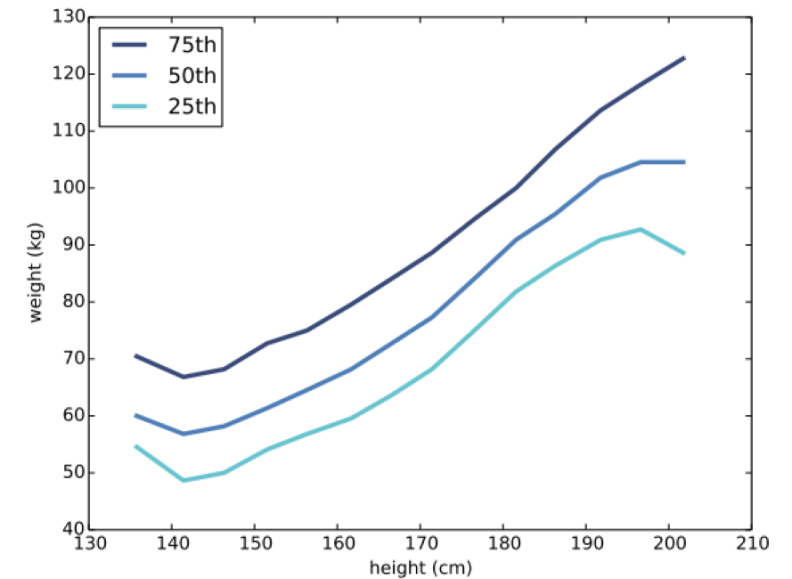
Hexbin Plot

1	height	weight	age	male
2	151.765	47.8256065	63	1
3	139.7	36.4858065	63	0
4	136.525	31.864838	65	0
5	156.845	53.0419145	41	1
6	145.415	41.276872	51	0
7	163.83	62.992589	35	1
8	149.225	38.2434755	32	0



Characterizing relationships

1	height	weight	age	male
2	151.765	47.8256065	63	1
3	139.7	36.4858065	63	0
4	136.525	31.864838	65	0
5	156.845	53.0419145	41	1
6	145.415	41.276872	51	0
7	163.83	62.992589	35	1
8	149.225	38.2434755	32	0



Idea 2. Correlation Coefficient

Covariance

Covariance is a measure of the **tendency** of two variables to vary together.

$$\text{cov}(X, Y) = \text{E} [(X - \text{E}[X])(Y - \text{E}[Y])]$$

$$\text{cov}(X, Y) = \text{E}[XY] - \text{E}[X] \text{E}[Y]$$

Hard to interpret
113 kilogram-centimeters

Pearson's correlation

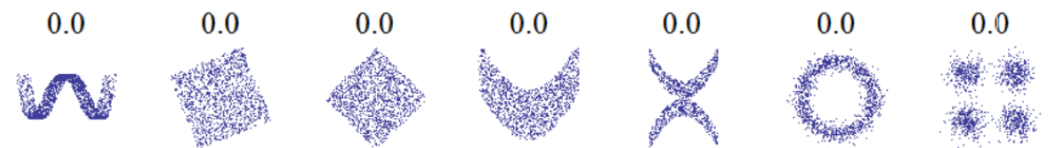
Pearson's correlation is a measure of the **linear relationship** between two variables

$$\rho_{X,Y} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$$

Easy to Interpret

- $[-1, 0) \rightarrow$ Negative Correlated
- $(0, +1] \rightarrow$ Positive Correlated
- -1 or $+1 \rightarrow$ Perfectly Correlated

What about non-linear relationship?



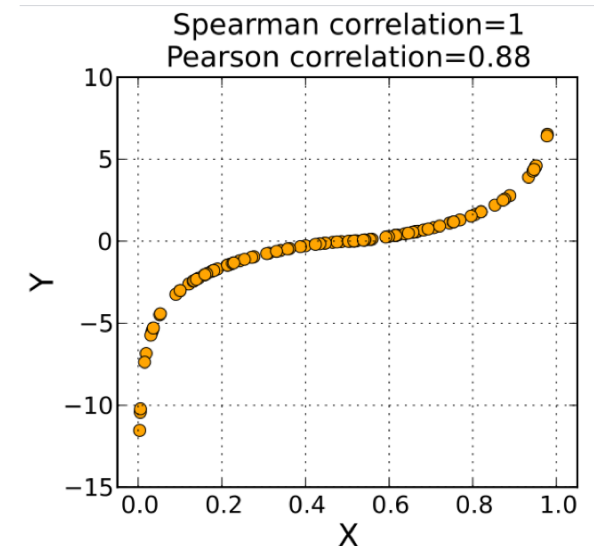
Spearman's rank correlation

Spearman's rank correlation is a measure of **monotonic relationship** between two variables

$$r_s = \rho_{r_X, r_Y} = \frac{\text{cov}(r_X, r_Y)}{\sigma_{r_X} \sigma_{r_Y}}$$

Advantages

- Mitigate the effect of outliers
- Mitigate the effect of skewed distributions



Outline

Correlation Analysis

- Big Picture
- How to do correlation analysis

Hypothesis Testing

- Big Picture
- A/B Testing

Why Hypothesis Testing?

We want to make a claim from our data
But, data is just a sample
How to prove our claim in this situation?

Using Hypothesis Testing

Example

- Claim: A data scientist earns more money than a software engineer
- Data: A sample of 50 data scientists and 50 software engineers
- Result: 100K vs. 70k

Can we use this result to prove that our claim is correct?

Hypothesis Testing

Equivalent Terms

- Hypothesis == Claim
- Hypothesis Testing == Claim Proving

Key Idea

- Prove by contradiction

Analogy

- How to prove: There is no smallest rational number greater than zero.
- Hint: a rational number is any number that can be expressed as the fraction a/b of two integers

Alternative and Null Hypotheses

Alternative Hypothesis (H_a)

- This is the claim that you want to prove it's correct

Null Hypothesis (H_0)

- The opposite side of H_a

Possible Outcomes

- Reject H_0 (a contradiction is found) \rightarrow Accept H_a
- Fail to reject H_0 (no contradiction is found)

Example

Alternative Hypothesis (H_a)

- A data scientist earns **more** money than a software engineer

NULL Hypothesis (H_0)

- A data scientist earns **less (or equal)** money than a software engineer

If H_0 is true, what's the probability of seeing:

- Data Scientist (100 K) vs. Software Engineer (70 K)



This is called P-value

Make a decision based on p-value

We hope that

- p-value is as low as possible so that we can reject H_0 (i.e., accept H_a)

Level of Significance (e.g., $\alpha = 0.01$)

- How low do we want p-value to be?

Level of Confidence (e.g., $c = 1 - \alpha = 99\%$)

- How confident are we in our decision?

P-Hacking (Cheating on a P-Value)

Common Mistakes

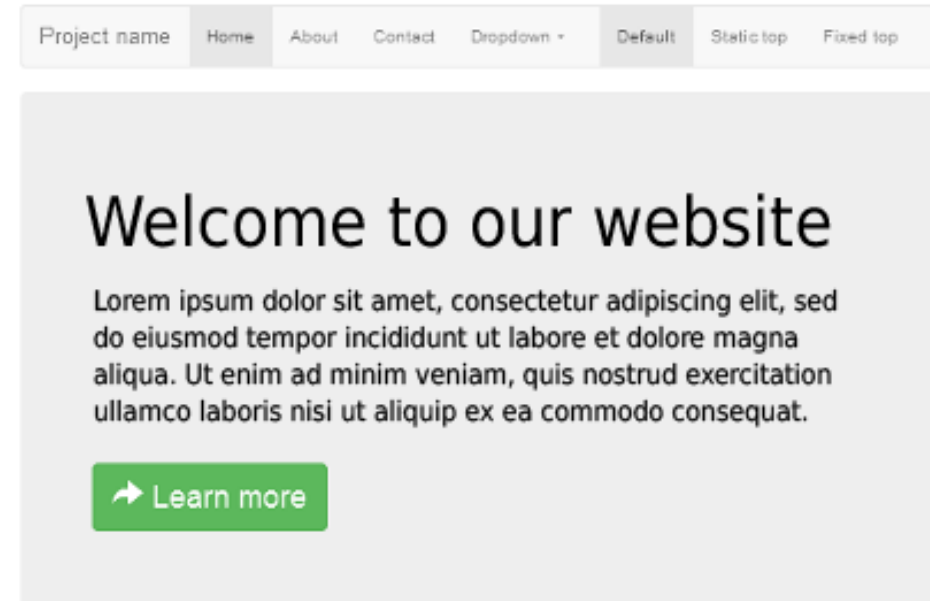
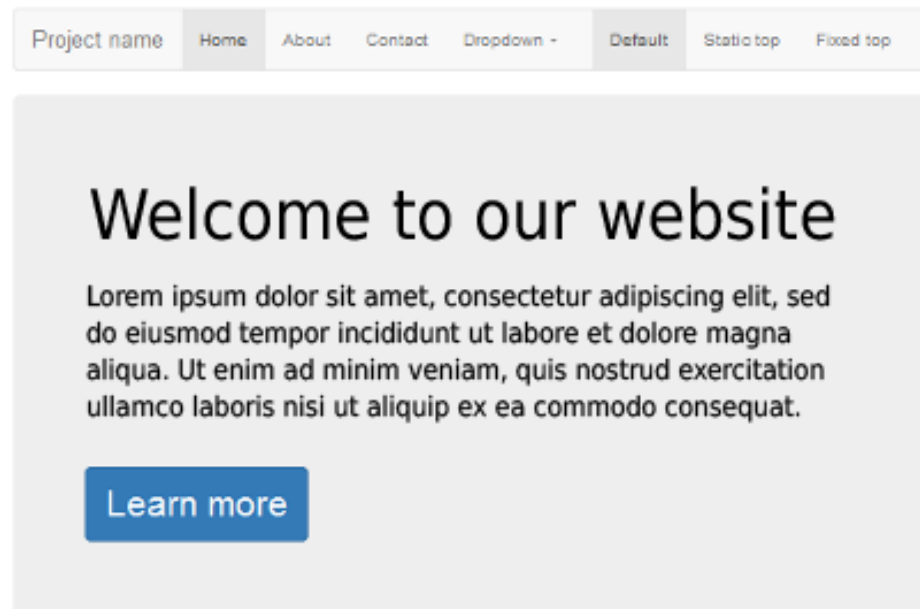
1. Collect data until the hypothesis testing is passed
2. Keep doing analysis on the same data until you find something significant

Solution

- You should know what you're looking for (H_0 and H_a) before you start
- Decrease the level of significance (e.g., $\alpha/2$ for two hypothesis tests on the same data)

A/B Testing

What UI is better?



Surprising A/B Tests

A. Get \$10 off the first purchase. Book online now!

B. Get an additional \$10 off. Book online now.

Control Button



Experiment Button




<https://www.wordstream.com/blog/ws/2012/09/25/a-b-testing>

Permutation Test

<https://youtu.be/lq9DzN6mvYA?t=8m9s>

**Sneetches:
Stars and
Intelligence**



Test Scores

★		×	
84	72	81	69
57	46	74	61
63	76	56	87
99	91	69	65
		66	44
		62	69

★ mean: 73.5
× mean: 66.9
difference: 6.6

PYCON 2016
ROSE CITY
PORTLAND, OREGON
MAY 28TH - JUNE 5TH

3/12/18 8:51 / 40:44

Conclusion

Correlation Analysis

- Using visualizations (scatter plot, hexbin plot)
- Using correlation coefficients (Pearson, Spearman's rank)

Hypothesis Testing

- Null Hypothesis (H_0) and Alternative Hypothesis (H_a)
- P-value and P-hacking
- A/B Testing