Mechatronic Engineering program

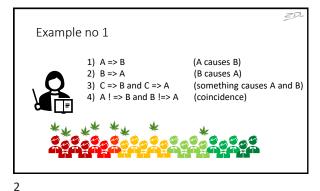
Python for machine learning and data science:

2: Scientific method, exploratory data analysis

Ziemowit Dworakowski

AGH University of Krakow

1



Example no 1

(A causes B)

1) A => B 2) B => A

(B causes A)

3) C => B and C => A

(something causes A and B)

ZD.

4) A! => B and B!=> A (coincidence)

Maybe gather more data

and observe whether the relations is maintained?

Example no 1

1) A => B (A causes B) 2) B => A (B causes A)

3) C => B and C => A (something causes A and B)

ZD.

ZD.

Maybe calculate more features and check for other correlated features?

Maybe try to artificially cause A or B and see if the other follows?

4

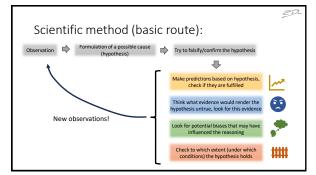
Example no 1

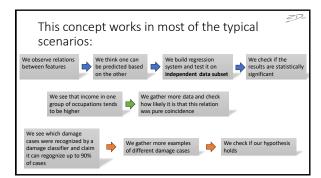
A => B (A causes B)
 B => A (B causes A)

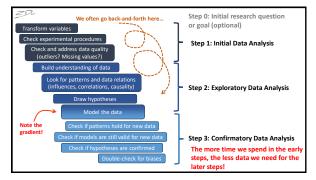
Maybe look for time-dependent relations? If something was observed first, it might have been the root cause

Maybe look for physical model to infer causality based on expert knowledge?

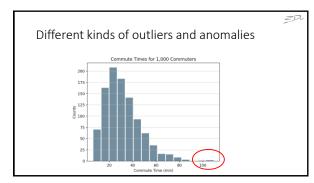
5

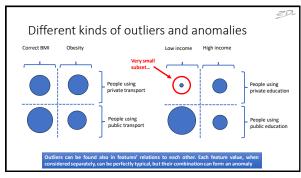


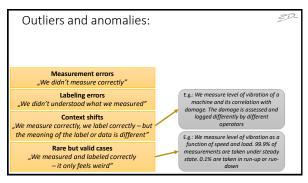


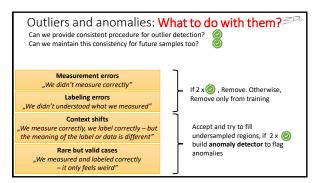


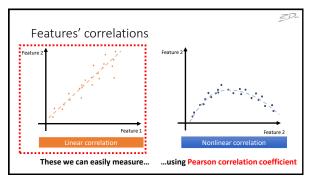
| Employment area | Earnings / mo | Age | Education | Opened term deposit? |
|--------------------|------------------|-----|------------|----------------------|
| Agriculture | 3900\$ | 34 | MSc | Yes |
| Agriculture | 2400\$ | 45 | College | Yes |
| Law | 5400\$ | 41 | Illiterate | Yes |
| Transportation | 3000\$ | 30 | College | No |
| Science & Ed | 4200\$ | 36 | MSc | No |
| Food industry | 5100\$ | 32 | BSc | Yes |

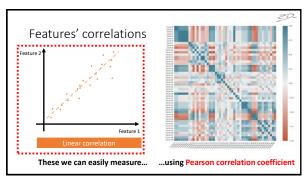


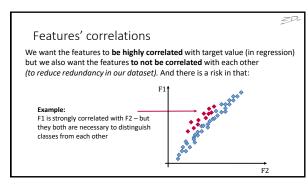


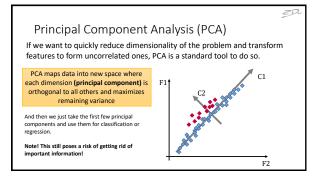


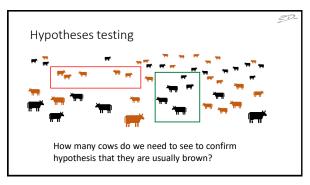












Hypotheses testing

P-value: Probability of obtaining observed results by coincidence (while the assumed hypothesis is not true)

In a (gross) simplification:

- 1) Start from a formulated hypothesis (e.g. more than 90% of cows are brown)
- 2) Assume a threshold for "hypothesis acceptance" -> α = 0.05 This means that we accept the fact, that 5% of our studies will end in a false positive conclusion (in strict tests α probably should be higher)
- 3) Consider experimental data (e.g.: we observed 5 brown cows)
- Calculate how likely it is, that a test result returned 5 brown cows if the underlying probability of a brown cow is **at most** 90% (That we obtained such results despite hypothesis being not true)

 $p = 0.9^5 = 0.59$

ZD.

5) If this probability (**p-value**) is lower than **0.05**, we say that the test was statistically significant (It is unlikely to get this data given false hypothesis)

19

| | | | | | | | <i>ZD.</i> |
|--------------------|---------------|---------------|---------------------|---|--|----------|---|
| Hypotheses testing | | | | result | oility of obtaining observed is by coincidence med hypothesis is not true) | | |
| | Brown cows | Black cows | Assumed hypothesis | Null hypothesis (alternative to assume | d?) | p-value | _ |
| | 5 | 0 | More than 90% brown | 90% or less are brown | | 0,590 | $p = 0.9^5 = 0,590$ |
| | 50 | 0 | More than 90% brown | 90% or less are brow | vn | 0,005 | $p = 0.9^{50} = 0,005$ |
| | 5 | 0 | More than 50% brown | 50% or less are brow | m | 0,031 | $p = 0.5^5 = 0.031$ |
| | 4 | 1 | More than 90% brown | 90% or less are brow | m | 0.328 | $p = 5 \cdot (0.9^4 \cdot 0.1) = 0.328$ |
| | | | | better: 80% are brow | vn | 0.409 | $p = 5 \cdot (0.8^4 \cdot 0.2) = 0.409$ |
| | 0 | 5 | More than 90% brown | 90% or less are brow | m | 0.00001? | $p < 0.1^5 = 0.00001$ |
| | | | | better: 0% are brown | า | 1 | $p < 1^5 = 1$ |
| | | | | | | | |

20

Hypotheses testing

Low enough p-value tells us only that data support hypothesis in a statistically significant way. It is **not** a final confirmation of the hypothesis.

It works under the assumption that the tests are independent to each other (often not true! We can just happen to observe one pasture with cows of the same color and it will tell us nothing regarding the whole cow population)

High p-value does not tell us that the hypothesis is false.

p-value **should not** be used for early stopping of the experiment or to select a subset of data to confirm a hypothesis <- this is p-hacking

| Things to | remember: |
|-----------|-----------|
|-----------|-----------|

ZDL

- 1. Explain an overview on scientific method
- 2. Explain steps of data analysis (IDA, EDA, CDA)

- 2. Explain seeps of data drialysis (IDA, EDA, CDA)

 3. Explain sources for anomalies, explain how they affect model preparation

 4. What does it mean that data are correlated? How do we measure correlation?

 5. How does PCA work?

 6. What are the risks associated with looking blindly into correlation information or using PCA?
 7. How do we test hypotheses?
- 8. Explain what a p-value is and how is it used. Note what p-value does **not** allow.