

- skúsť otázky: najbližšie dnu
- ANKETA

Ak projekt vyriešite skôr, napíšte mi, prosím mail

Minule: Testy

perm. test: H_0 : rozdelenie X = rozdelenie Y (vzhľadom na stat. \bar{S})
 H_1 : *

napr. $\bar{S}_* = \bar{x} - \bar{y}$ X : muži, Y : ženy

• $x_1, \dots, x_n, y_1, \dots, y_m$

preusporiadame znaky $H/\bar{z} \Rightarrow \bar{S} = \dots \rightarrow$ * preusporiadanie

• H_0 tam, ak \bar{S}_* je v 5% najextrémnejších

Pr.: 3 muži, 2 ženy

H_1 : muži sú väčší ako ženy

valky $\hookrightarrow 87, 95, 72 \hookrightarrow 63, 79$

$\bar{S}_* = \bar{x} - \bar{y} = \underline{13,7}$

$\frac{87, 95, 72, 63, 79}{m_1 \quad m_2 \quad m_3 \quad z_1 \quad z_2} \rightarrow \bar{S} = \dots$

$m_1, m_2, z_1, z_2, m_3 \rightarrow \bar{S} = \dots$

\downarrow
 $\bar{S}!$ \bar{S} -ieč



p -val = % extrémnejších \bar{S}

- časťnosť: iba K náhodných permutácií: aproximácia

Big Data

- také, ak. naviac spracovať na 1 počítači
 \rightarrow Twitter, FB, Google...

- rýchlosť - statické: fixné

- "historické": pomaly pribúdajú - denné výnosy

- streamovacie: real-time pribúdajú

: nemáme vždy kapacitu uložiť

\rightarrow snímanie kamier, Twitter...

} Big data

- kvalita - rozbiere na účelom (chceme zistiť účinnosť lieku)

- výsledok plošného zberu: všetko, čo čomu sa dostane } Big data

→ nemusia byť reprezentatívne...

Hlavný problém: nemôžeme použiť klasické postupy (parametrické)

Riešenia: I. vyhnúť sa práci s veľkými dátami

→ 1.) Vzorovanie (sampling, subsampling) - iba časť dát

→ 2.) Sumarizácia (sufficient statistics) - pracovať iba s tým, čo potrebujeme

II. Špec. postupy pre Big data

1.) Vzorovanie:

- modely typicky majú rozumný # parametrov → môže stačiť menej dát

a) Náhodný výber: reproducibilita?

: môže byť ťažké získať viac náh. vzoriek, kt. sa nepodávajú
(train, val, test)

Nenáhodne!

b) Prvých k záznamov → môže byť veľmi nerepresentatívne

c) "Každé 1-té": 1, 1001, 2001, 3001, ...

: take $i \equiv 1 \pmod{1000}$

alebo $\equiv 34 \dots$

: takže nepredstavujúce sa výbery: $i \equiv 1 \pmod{1000} \rightarrow 1,$
 $i \equiv 2 \pmod{1000} \rightarrow 2,$

: môže byť ok, ak dáta majú špec. štruktúru:

napr. zdravý - chorý - zdravý - chorý ...

všed.: pozor na reprezentatívnosť

- nedáme štatist. relevantných outlierov

→ "malé firmy" (napr. nezbankrotované firmy)

→ stratifikované vzorkovanie:

- najprv rozdeľme na mužov a ženy → potom z každej skupiny: výber

→ váhované vzorkovanie: môže pr. rôznym záznamom

(napr. vysoká pp. pre objekty zo vzácnnej triedy)

2.) Sumarizácia

- uložiť si len to, čo nutne potrebujeme pre model
 - sufficient statistics (postacujúca štat.)
- napr. kontingenčná tabuľka namiesto veľa záznamov

II. Keď to inak nejde: špecializované postupy

- relačné databázy (SQL)
 - "shamelessly" paralelizovateľné výpočty
 - ak sa dá rozložiť na menšie výpočty, kt. sú spúšťateľné samostatne
 - každý výpočet na 1 počítači / jadre...
 - paralelné výpočty
 - každá jednotka svoje, ale musí sa to priebežne spájať / aktualizovať
 - treba minimalizovať množstvo dát, kt. si počítače medzi sebou preštvávajú
- 3 "jazýky": frameworky - MapReduce, Hadoop
- treba používať overené nástroje, neimplementovať vlastné riešenia
 - pri veľkých dátach sa nič nepočíta, frameworky to majú oštiepené