

Junior Data Scientist Akadémia 3. alkalom

- Adatgyűjtés

- Bash

Adatgyűjtés: "Garbage in, garbage out ..."

-> 6 lépés -> 1. Adatgyűjtés -> 1st lehet el.

2. Adattisztítás

3. Adatintegrálás

4. Adatkezelés

5. Visualizálás -> 1st lehet el.

6. Döntés -> 1st lehet el.

Google Tag Manager -> követő adat tud beszerezni

Felbőr -> adatminőségért

Mit gyűjtünk? -> Vannak kockázatok: -> nem az adathoz való hozzájárulás

-> "nem juttat el"

-> titkosság, jelszó kockázatok

-> programozói hiba miatt kockázatok

-> adat betöltési sebesség (villanás)

- Hogyan döntsd el, hogy mit akarsz gyűjteni? -> mindig mindig a tényleges

-> gondolkodás visszakér (1-on-1-el)

=> Workshop! Workshop! Workshop!

Magyar málkó?

Website → Tooling scripts → Data warehouse

↓
Frontend → Production server → termé

↳ Data server → Data server → ahát több server is lehet

Adatgyűjtési szint → SQL → struktúrált adatok

→ noSQL → nem struktúrált adatok

→ JSON → ábrázolási formátum

→ graph adatbázis rendszar → PL: neo4j, neo4j

Bash :

→ nem adatbázis nyelvé

→ ez egy megosztó a nyelvé kint

→ direkt kérésre nyelvé, vagy kérésre / adatok nyelvé

→ Fájlszerkezet

→ Data Science at the Command line (könyv)

→ SSH ⇒ távoli munka

clear → törlés a képernyőn

cd . → jelenlegi kó

ls → lista

pwd → Hol vagy?

cd → visszatérés az előzőre

Command prompt.

- grep (grep lehet, csv => használható is)

- head -100 lehet.csv | wc -c

cut -d'|' -f3,4,11 lehet.csv > osdopok.csv

- cut -d'|' -f4 lehet.csv | cut -d'|' -f1 > uszame.csv

- grep "Volkswagen" lehet.csv | head

- cut -d'|' -f11 lehet.csv | sort -n

- sort -n -t'|' -k11 lehet.csv | tail

- sort -n -> duplikátumokat tünteti el

- uniq -> összegezés

- cut -d'|' -f2 lehet.csv | sort | uniq -c

- cut -d'|' -f2 lehet.csv | sort -n | grep -v 'city' | grep -v

'UNKNOWN' | wc -l

- Script -ing is automaticus

- accident.sh

- contact -l

- psql -U user14 -d postgres -f proba.sql -A'|'

Loreskop:

- Adatgyűjtés

1. Milyen adatokat gyűjtünk le?

2. Milyen paramétereket?

3. Milyen struktúrában?

Típusátok

Erőforrások

Vektori gyorsaság változó adatok

Adatok minőségi állapot

Star is Snowflake név