

NASLOVNICA

INFORMACIJE

6 faz (podrobno v e-učilnici)

Programska orodja

- WEKA
- MS Excel / LO Calc
- Python

Ocena:

- Izpit: 50% (lahko s kolokviji, vsak vsaj 40%, avg. 50%)
- Kratke DN: 15% (9-10 nalog a.k.a. kvizov)
- Projekt: 30% → 1. md na GH repo
- Ustni izpit: 5%
- Izračun: $0,15 \times DN + 0,5 \times PI + 0,3 \times PP + 0,05 \times UP$

Kolokviji:

1. 21.11.2024
2. 16.01.2025

- Posebne prijave (mimo ŠISa)

Prijava na izpit v ŠIS

10.10.29

STROJNO UČENJE, AI, PODATKOVNO RUDARJENJE

Trendi, ki vodijo v podatkovno poplavo

- Generiranje velikega št. podatkov:
 - Poslovne ustanove (npr. Banke)
 - Znanstveni podatki (npr. biologija)
 - Splet (npr. socialna omrežja)

Primeri ogromnega št. podatkov

- Evropska VLBI mreža teleskopov generira 1 gigabit podatkov na sekundo za vsak teleskop v 25 dneh.
- AT&T posreduje več kot milijardo klicev na dan

Največje zbirke podatkov leta 2003

- Komercialne DB:
 - France Telecom: 30 TB
 - AT&T: 26 TB
- Splet:
 - Alexa internetni arhiv: 7 let = 500 TB

Kaj pa danes?

- Konec junija 2017 CERN shranjeval 200 TB podatkov

Področja uporabe

- Znanost (npr. bioinformatiki)
- Poslovni svet - profiliranje strank
- Splet
- Državna uprava

Prikladni problemi za podatkovno rudarjenje

- Zahtevajo odločitve na podlagi znanih
- Imajo opreminjajoče okolje
- Imajo ne-optimalne trenutne rešitve
- Imajo veliko poplačilo za prave odločitve
- Paziti moramo na zasebnost!

Def. "Odkrivanja zakonitosti"

- ... v podatkih je:
 - Netrivialni proces identifikacije veljavnih, novih potencialno uporabnih in nazadnje razumljivih vzorcev v podatkih

Statistika, str. učenje in PR

- Statistika:
 - Bolj usmerjena v teorijo
 - Večji fokus na testiranje hipotez
- Strojno učenje
 - Uporaba heuristik
 - Usmerjeno v izboljševanje performans učnih agentov
 - Se ukvarja z učenjem v realnem času in robotiko
- Podatkovno rudarjenje
 - Statistika + Strojno učenje
- Razmejitev so "mekke"

Zgodovina

- 1960 - Data fishing, dredging → uporabljeno s strani statistikov
- 1990 - Data mining - uporabljajo DB strokovnjaki
- 1999 - Knowledge Discovery in DB - Uporabljajo AI, ljudje s področja strojnega učenja

Tehnike str. učenja

- Klasifikacija
- Razvrščanje v skupine
- Asociacije
- Vizualizacija
- Povzemanje
- Odkrivanje odstopanj
- Regresija / ocenjevanje
- Analiza povezav

Klasifikacija

- S pomočjo statistike, odločitvenega drevesa, nevronske mreže,...

Razvrščanje v skupine

- Razvrstimo vs skupine po razdalji

PR PO CRISP-DM STANDARDU

CRISP-DM

- CROSS-Industry Standard Process for Data Mining

Zakaj standardizirati?

- PR je proces, ki mora biti ponovljiv in zanesljiv
- Nudenje okvira za shranjevanje preteklih izkušenj
 - Projekte lahko repliciramo
- Lažje planiranje in vodenje projektov
- Faktor udobja za nove uporabnike
 - Dokazuje zrelost PR
 - Zmanjšuje odvisnost od t.i. super-strokovnjakov

Kako je prišlo do standardizacije

- Inicijativa s strani izkušenejših podjetij
- Razvoj in dopolnitve na workshop-ih
- K procesnemu modelu prispevalo 300+ organizacij
- Leta 1999 izide CRISP-DM 1.0

CRISP-DM je:

- Ne-lastniški
- Neodvisen od aplikacije / problema
- Neodvisen od orodja / programa
- Osredotočen na poslovne probleme
- Ogradnje, ki služi kot vodilo
- Baza iz preteklih izkušenj

CRISP-DM pregled:

- Podatkovno rudarjenje kot metodologija
- Za vsakega
- S "komplet" navodili
- V ciklu iz 6 faz

CRISP-DM faze

1. Razumevanje problema
2. Razumevanje podatkov
3. Priprava podatkov
4. Modeliranje
5. Evalvacija (ocena modelov)
6. Predaja končnemu uporabniku

Razumevanje problema

- Določiti cilj projekta
- Oceniti trenutno situacijo
- Določiti cilje pod. rudarjenja
- Narediti projektni plan

Razumevanje podatkov

- Zbrati začetne podatke
- Opisati podatke
- Raziskati podatke
- Preveriti kakovost podatkov

Priprava podatkov

- Izbrati podatke
- Precistiti podatke
- Sestaviti podatke