

Fundamente de Big Data

Proiect

Termen limita de predare: cel puțin trei zile înaintea datei de examen anunțată în sesiunea de examen

Privire de ansamblu

Obiectivul acestui proiect este sa permită studenților sa experimenteze întreg ciclul de cercetare asociat studiului unei probleme de afaceri, pornind de la identificarea unei întrebări de cercetare interesante, obținerea setului de date, pregătirea acestuia pentru analiza, executarea sarcinii de analiza a datelor utilizând mai multe metode alternative, extragerea concluziilor si pregătirea unui raport de comunicare a rezultatelor.

Pregătirea sarcinii de lucru

Pentru acest proiect sugerăm studenților sa lucreze in RStudio, folosind limbajul de programare R. Studenții pot sa utilizeze si alte medii de programare dedicate R, sau să folosească Python ca si limbaj alternativ.

Pentru editarea raportului final, sugestia este sa folosiți Microsoft Word.

Studenții vor lucra in echipe de maxim câte 3, cu mențiunea ca toti studenți făcând parte dintr-o echipă trebuie sa cunoască în detaliu tot proiectul, incluzând aici toate experimentele executate in mediul de programare precum si concluziile desprinse de pe urma execuției acestor experimente. Proiectele vor fi susținute oral, si membrii echipei vor răspunde individual întrebărilor puse de examinatori. Fiecare student membru al unei echipe trebuie să fie capabil să răspundă la întrebări care vizează tot proiectul trimis spre evaluare. În funcție de modul in care studenții reușesc sa răspundă la întrebările puse pe parcursul susținerii, ei pot să contribuie pozitiv sau negativ la notarea finală a proiectului. Notarea finală este individuală pentru fiecare student, si nu pe proiect. Studenții care nu participă la susținerea proiectului, chiar dacă proiectul a fost predat pe Moodle, se consideră a fi absenți la examen.

Studenții înrolați la învățământ la distanță pot să lucreze și in echipe de 2 studenți sau individual.

Proiectul trebuie predat pe Moodle la termenul specificat, cu cel puțin **trei zile** înaintea datei alese pentru prezentarea la examen.

Instrucțiuni detaliate

In acest proiect aveți de îndeplinit mai multe sarcini.

1. Prima sarcina este sa identificați una sau mai multe întrebări de cercetare de interes pentru o audienta de afaceri.

2. Să identificați și obțineți un set de date care să vă permită să răspundeți la întrebarea de cercetare pusă la pasul 1. Setul de date trebuie pregătit de așa manieră încât să puteți aplica metode științifice de analiză a datelor pentru a obține modele alternative
3. Pentru întrebările de cercetare alese și setul de date pregătit în pasul 2 veți selecta mai multe metode de analiză a datelor și veți conduce experimente cu aceste metode pentru a crea modele alternative. Pe baza metodologiei de validare potrivită, veți determina metoda și modelul cel mai potrivit
4. Pe baza experimentelor realizate la pasul 3 veți extrage concluzii de business, încercând să răspundeți cât mai clar și mai convingător la întrebările stabilite la pasul inițial.

Pentru punctele menționate anterior sunt de interes atât setul de date, codul sursă care conține experimentele executate cât și raportul care prezintă întrebările și concluziile extrase în urma experimentelor.

Sfat: Puteți consulta secțiunea Resurse a acestui document pentru a desprinde idei posibile de întrebări de cercetare, precum și seturi de date relevante pentru acestea.

Structura raportului pe care trebuie să îl realizați este următoarea:

- **Introducere.** Reprezintă prima secțiune a raportului. Aceasta trebuie să furnizeze informații contextuale despre aria de cercetare aleasă, să identifice în mod clar întrebările de cercetare alese, să explice de ce aceste întrebări sunt relevante și importante și care este audiența care va beneficia de pe urma studiului, și dacă întrebările au fost abordate în trecut, care sunt rezultatele altor studii răspunzând la aceleași întrebări sau lucrând pe același set de date. Pentru introducere, trebuie să fiți concisi, și în general, să finalizați introducerea în maxim o pagină. Dacă este necesar mai mult spațiu, aveți libertate de expresie în acest sens.
- **Setul de date.** Folosiți câteva paragrafe ca să explicați setul de date utilizat. Precizați de ce acest set de date este relevant pentru întrebările de cercetare alese, descrieți pașii realizați pentru curățarea datelor sau preprocesarea acestora, precum și caracteristicile de bază ale datelor, înainte ca acestea să intre în procesul de analiză. Puteți folosi tehnici de vizualizare a datelor și să printați grafice relevante, cu mențiunea că acestea trebuie explicate pentru a face raportul inteligibil.
- **Rezultate și discuții.** Aceasta reprezintă partea cea mai importantă a raportului. Trebuie să prezentați în detaliu analiza realizată. Prezentați care sunt metodele de analiză alese, ce setări ați testat pentru aceste metode, care a fost strategia de validare selectată, ce rezultate ați obținut pentru metodele selectate și cu parametri testați. Comparați rezultatele obținute, și pe baza acestora precizați care este metoda și modelul final considerat și cum răspunde acesta întrebărilor de cercetare descrise în prima parte a raportului. În această parte a raportului, puteți motiva alegerile realizate prin fraze de tipul: “am folosit metoda deoarece ...”, etc. Aceasta secțiune trebuie să fie o combinație de text, tabele și grafice. Este absolut necesar să descrieți și interpretați rezultatele, nu doar să le afișați în tabele sau grafice. De asemenea, trebuie să discutați limitările acestor rezultate, dacă credeți că puteați

obține rezultate mai bune și ce anume v-a împiedicat în acest sens. Puteți să creați subsecțiuni care să structureze mai bine această parte a raportului.

- **Concluzia.** Cel mult 2 paragrafe în care să se sintetizeze întrebările de cercetare formulate precum și rezultatele obținute. Această secțiune trebuie să fie scurtă și concisă, însă nu trebuie să supraliciteze (adică să extragă concluzii mai puternice decât cele obținute din analiză și justificate în secțiunea precedentă)

Întregul raport va avea între 10 și 15 pagini, va fi redactat cu fontul Times New Roman de 12 caractere și spațiere maximă de 1.2 între rânduri.

Pe lângă raportul realizat, va trebui să furnizați și setul de date, precum și codul / codurile sursă folosite pentru procesarea și analiza datelor.

Pe Moodle veți încărca o arhivă care va conține:

- Raportul cerut mai sus
- Setul de date
- Fișierele cu codul sursă folosit pentru procesarea și analiza datelor.

Exemple de întrebări de cercetare

- În mediul bancar, cât de bine putem să identificăm clienții care nu vor putea să își ramburseze creditul luat?
- Cât de bine putem prezice vânzările unui magazin, pe un anumit domeniu comercial?
- Cât de bine putem prezice succesul box-office al unui film?
- Putem identifica un grup de persoane care să fie mai receptivi la o anumită formă de publicitate pentru un produs?
- Care sunt predictorii cei mai importanți pentru a caracteriza mișcarea de persoane pe piața muncii într-un domeniu particular?

Trimiterea proiectului

Termenul pentru trimiterea proiectului este: cel puțin trei zile înainte de data aleasă în sesiune pentru prezentare la examen

Veți trimite un fișier Zip care conține:

- Raportul în format Doc / PDF (neprotejate)
- Un director care să conțină codul / codurile sursă
- Setul de date. Dacă este un set de date foarte mare, puteți indica un link web de unde setul de date poate fi descărcat.

Atenție!!! Rapoartele vor fi verificate pentru similitudine folosind Turnitin. Proiectele cu grad mare de similaritate vor fi descalificate (nota 1 final).

Grila de notare – total 100 pct.

Următoarele criterii vor fi folosite pentru notarea proiectului (se acorda 10 puncte din oficiu):

- **Introducerea: 10 puncte**

- Daca se furnizează cititorului suficienta informație pentru a înțelege restul raportului
- Daca întrebările de cercetare sunt clar stabilite
- Relevanta respectiv importanta întrebărilor de cercetare. Daca contribuția propusa prin raport este clar prezentata

- **Setul de date (15 pct)**

- Daca datele culese sunt potrivite pentru a răspunde la întrebările de cercetare
- Daca datele sunt descrise corespunzător?

- **Rezultate si discuții (30 pct)**

- dacă analiza realizata este potrivita pentru a răspunde întrebărilor de cercetare alese
- dacă s-au ales metode potrivite de analiză, daca sarcinile de analiza au fost rulate corespunzător
- daca rezultatele obținute sunt interpretate corespunzător
- daca rezultatele prezentate sunt clare si aceasta prezentarea are o ordine logica, potrivita
- daca tabelele si graficele realizate au puterea de a informa asupra concluziilor si interpretărilor textuale
- daca sunt prezentate limitări ale studiului si munca adiționala care se poate face pentru a obține rezultate si mai pertinente?

- **Concluzia (5 pct)**

- daca se furnizează un sumar scurt si concis potrivit pentru raport?
- daca concluziile sunt potrivite întrebărilor de cercetare alese si sunt susținute de analiza realizata

- **Codul sursa furnizat (30 pct)**

- daca codul sursa este complet, susține analiza, si poate fi rulat cu ușurință, pentru a reproduce rezultatele prezentate in raport
- cât de eficient este codul sursa (daca sunt taskuri duplicat care pot fi evitate)?

În final, calitatea scrisului contează. Deci încercați să fiți siguri că exprimările sunt clare și concise și se înțelege ceea ce doriți să transmiteți, să eliminați posibilele confuzii de interpretare.

Resurse

Mai jos aveți câteva seturi de date care pot fi folosite ca și sursă de date și întrebări de cercetare:

- Google's dataset search (<https://toolbox.google.com/datasetsearch>)
- Kaggle (<https://www.kaggle.com/datasets>)
- OpenML (<https://www.openml.org>)
- UCI ML (<https://archive.ics.uci.edu/ml/index.php>)
- KDNuggets (<https://www.kdnuggets.com/datasets/index.html>)