

Prezime Ime: \_\_\_\_\_ Broj indeksa: \_\_\_\_\_

## R uputstvo:

- U prilogu uz ovaj zadatak dat je materijal koji obuhvata dataset i cheatsheet. R projekat kreirati u pomenutom direktorijumu na desktopu i nazvati ga RProjekat
- Zadatak se priznaje samo ako je kreiran PROJEKAT, inače se ne pregleda. Nije dovoljno kreirati samo R script fajl. R script fajl ne nazivati isto kao ceo projekat!
- Folder u kojem se nalazi kreirani projekat mora da sadrži SVE elemente koji su neophodni da se na nekom DRUGOM računaru (na kojem postoje instalirani R i RStudio) projekat otvori, da se u okviru otvorenog projekta otvori odgovarajući R script, kao i da se taj R script izvršava bez dodatnih intervencija

## Zadatak (KMeans): Grupisanje pesama

U fajlu **spotify.csv** dati su podaci o najpopularnijim pesmama na Spotify platformi 2023. godine. Varijable dataset-a su:

- |  |   |
|--|---|
| • <b>track_name:</b> naziv pesme   | • <b>bpm:</b> mera brzine pesme (beats per minute)                |
| • <b>released_year:</b> godina izdanja   | • <b>mode:</b> karakter pesme (minor / major)                     |
| • <b>in_spotify_playlists:</b> na koliko playlista na spotify platformi se nalzi pesma | • <b>energy:</b> percipirani energetski nivo pesme                |
| • <b>in_spotify_charts:</b> pozicije pesme na spotify rang listi                       | • <b>instrumentalness:</b> nivo instrumentalnosti                 |
| • <b>streams:</b> broj puštanja pesme  | • <b>liveness:</b> nivo prisustva elemenata uživo izvođenja       |
| • <b>in_apple_playlists:</b> na koliko playlista na appleplatformi se nalzi pesma      | • <b>speechiness:</b> Skalirana količina izgovorenih reči u pesmi |
| • <b>in_apple_charts:</b> pozicije pesme na applerang listi                            |   |

## Potrebno je uraditi sledeće:

1. Napraviti podskup podataka koji ne sadrži pesme čija je skalirana količina izgovorenih reči (**speechiness**) preko 60. Za sve naredne zahteve zadatka koristiti ovako dobijen dataset. **(1 poen)**
2. Odabrati atribut koji će biti uključeni u **KMeans model** i **navesti razlog** za njihov odabir ili neodabir. **(6 poena)**
3. U dataset-u dobijenom na osnovu prethodna dva zahteva proveriti da li postoje nedostajuće vrednosti (NA, "-", " ", ili "") i ako je moguće, zameniti ih adekvatnim vrednostima. Prokomentarisati postupak zamene vrednosti, odnosno zašto je određen metod zamene vrednosti odabran. Tako dobijeni dataset eventualno **dodatno obraditi** da bi bio pogodan za primenu metode KMeans. **(8 poena)**
4. Primenom **Elbow metode** utvrditi najbolju vrednost za broj klastera (**k**) **(8 poena)**
5. Izvršiti klasterizaciju za izabranu (tj. utvrđenu najbolju) vrednost za **k**. **(2 poena)**
6. **Interpretirati**<sup>1</sup> dobijene klastere (grupisane pesme) na osnovu: broja pesama po klasteru, centara klastera, disperzije od centra. **(10 poena)**

**Napomena:** Varijable koje **očigledno** nema smisla da ostavi u modelu **mogu izbaciti bez analiziranja** (npr. Izbacivanje ID-a) - potrebno je samo napisati komentar zasto se ta kolona izbacuje.

<sup>1</sup> Značenje tog izraza je: **ne** navoditi u komentarima samo brojke koje se dobiju, već protumačiti šta one **znače** u kontekstu problema koji se rešava u zadatku i dataset-a koji se koristi.