

# Documentație proiect ML Fundamentals

## Task 2 - Stackoverflow Data Analysis

*Tița Andrei-Petruț*

*Februarie 2021*

## Partea 1

Pentru prima parte, unde trebuia să selectez top 3 topicuri dintr-un anumit post, am selectat din baza de date, doar coloanele „title”, „body”, „tags” pentru că toate 3 pot conține un topic. Cred totuși că cele mai multe topicuri se află în coloana „tags” (cel puțin numele limbajului de programare).

În acest exemplu datele nu vin însoțite de label-uri. Deci voi face o verificare luând la întâmplare câteva postări.

În prima parte o să filtrez datele cu aceleași procedee de la task-ul 1 (Twitter Analysis). În plus am șters tagurile de html si pipe-urile de la tags.

Când am vizualizat datele mi s-a părut că elementele din tag-uri au mai mare relevanță și pare ca ar putea fi topicuri așa că la construirea stringului cu title+body+tags m-am gândit să adaug de două ori tags pentru a crește aparițiile cuvintelor din tags și a le acorda o șansă mai mare de a deveni topicuri.

Inițial am încercat să folosesc tf-idf pentru găsirea topicurilor însă apoi mi-am dat seama, din scoruri, că multe topicuri au foarte multe cuvinte în comun pentru că datele provin de pe un site cu întrebări de programare. Deci nu cred că este o variantă bună să decidem topicul unei întrebări în funcție de frecvența inversă a cuvintelor din alte întrebări. Așa că mă voi baza doar pe frecvența cuvintelor din întrebarea curentă.

Pentru asta am folosit metoda `get_top_n_words` care calculează frecvența cuvintelor dintr-un string (fiecare linie din setul meu de date) și returnează un vector cu primele `n` elemente (și cu frecvența lor).

Am pus și un `ngram_range` la `CountVectorizer` pentru ca observasem că la un text puseseră cele mai frecvente cuvinte „visual” și „studio” care ar trebui să fie un singur topic.

Am adăugat toate aceste rezultate într-un vector pe care l-am pus ca o nouă coloană în dataframe.

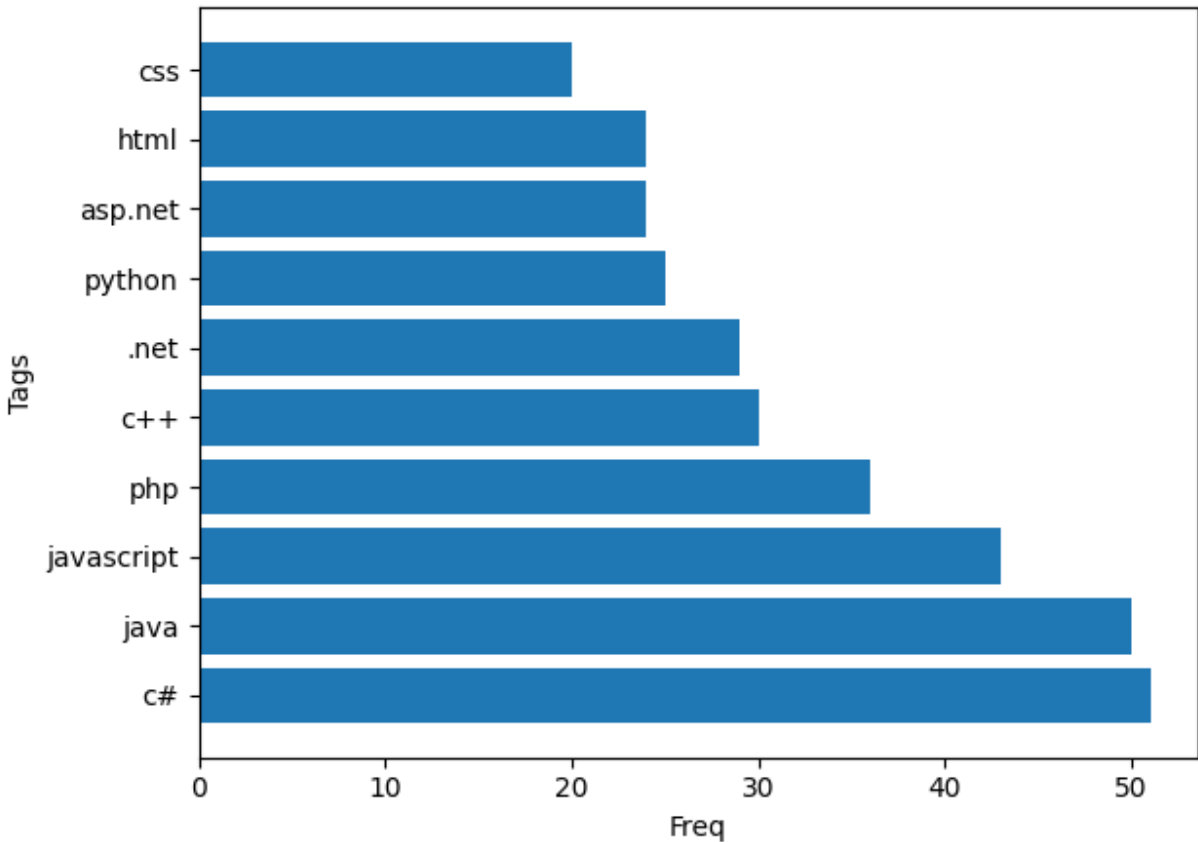
## Partea 2

Ma gândisem să încerc să fac un model de ML care să prezică în funcție de mai multe criterii (precum scor, vizualizări și altele) ce badge ar trebui să primească un utilizator dar după ce am verificat pe site ce condiții sunt la badges mi-am dat seama că ar fi prea complicat așa că am ales să fac niște taskuri mici de data analysis. Așa că am vrut să aflu care sunt cele mai folosite 10 taguri. Observ că cel mai folosit tag este „C#” urmat de Java și JavaScript.

Am folosit același fișier .csv pentru că am avut nevoie doar de aceleași câmpuri de date.

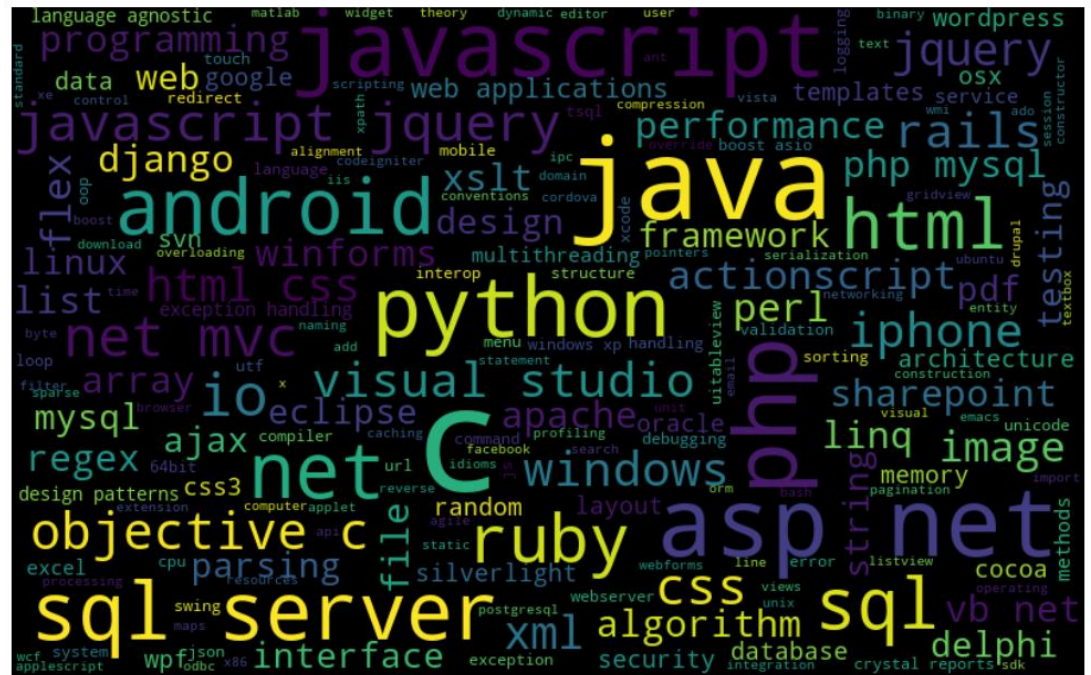
Nu am explicat fiecare parte din cod în documentație pentru că am comentat codul și mi se pare destul de self explanatory.

### Top 10 most freq tags

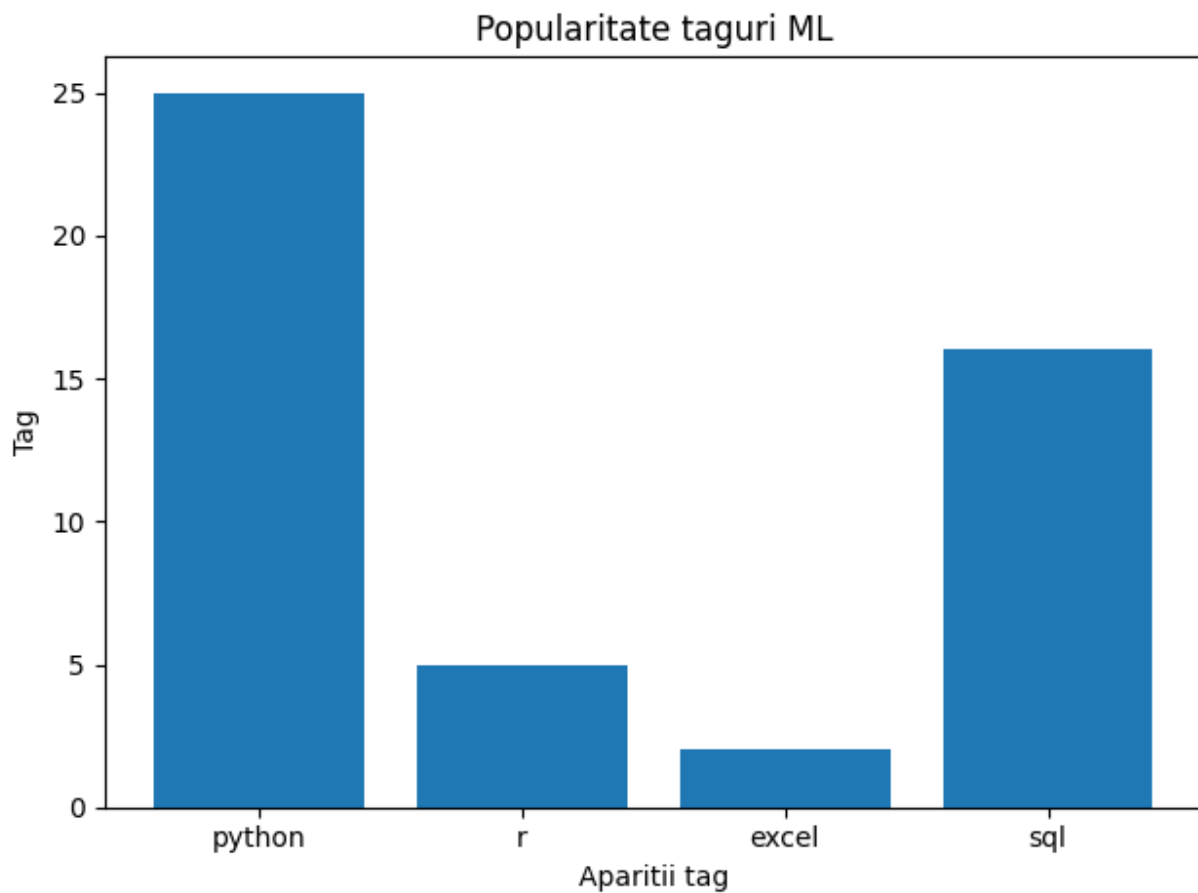


Pentru o vizualizare a mai multor taguri importante am combinat vectorul cu toate tagurile intr-un singur string si am folosit wordcount astfel am obtinut:

Observ ca C (fara #, nu stiu de ce, probabil a fost eliminat), java si javascript apar destul de mari in imagine deci se confirma graficul de mai sus.



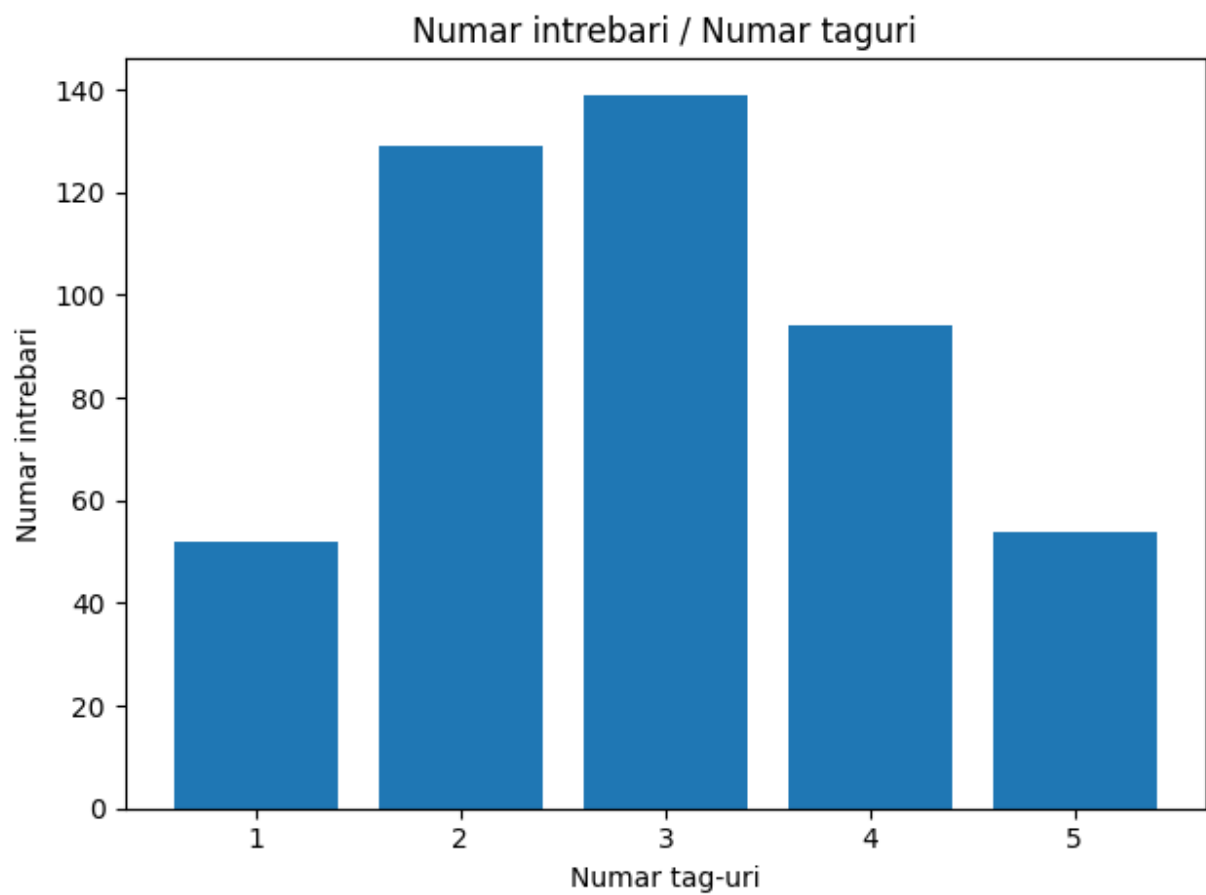
În continuare m-am gândit să vad cât de populare sunt tagurile de ML în baza de date. Astfel am definit un vector cu taguri care au legătură cu machine learningul: `ml_tags = ['python', 'r', 'excel', 'sql']`. După definirea valorilor de pe axa x și axa y și plotare am obținut următorul grafic:



Observăm că python și sql sunt destul de întâlnite, python chiar apare pe locul 1 la frecvența tuturor cuvintelor.

Mai departe am încercat să aflu câte taguri ar trebui să folosesc ca să am cea mai mare șansă să mi se răspundă la o întrebare, deci am făcut un grafic cu numărul de întrebări în funcție de numărul de taguri.

După afișarea graficului observ că cele mai multe întrebări au 3 taguri deci probabil șansele de a primi un răspuns sunt cele mai mari când folosim 3 taguri.



Am aflat astfel, prin tehnicile de mai sus câteva inside-uri despre baza de date folosind numai 3 coloane.