

Tema PCLP3

Partea II

Nume: Solomon Stefan

Coleg: Carp Andrei Costin

Grupa: 312 CA

Github: [https://github.com/krpandrei05/Tema1\\_PCLP3](https://github.com/krpandrei05/Tema1_PCLP3)

## 1. Prelucrarea datelor

1. a) Encodarea valorilor categorice: Datele contin coloane cu valori de tip text (aplicatie, daca clientul utilizeaza cupoane, ce dispozitiv foloseste, metoda de plata). Am folosit LabelEncoder pentru coloane binare (aplicatia, cupoane\_de\_reducere) si OneHotEncoding pentru coloane cu mai multe valori posibile (dispozitiv, metoda\_de\_plata).

1. b) Normalizarea si standardizarea valorilor numerice: Normalizare cu MinMaxScaler (aduce valorile intre 0 si 1) si Standardizare cu StandardScaler (transforma datele intr-o distributie cu media 0 si deviatie standard 1). Coloane prelucrate: valoare\_comanda\_medie, timp, frecventa\_lunara, varsta, recenzii.

Datele prelucrate au fost salvate in doua fisiere separate: train\_encodat.csv (pentru antrenare) si test\_encodat.csv (pentru testare).

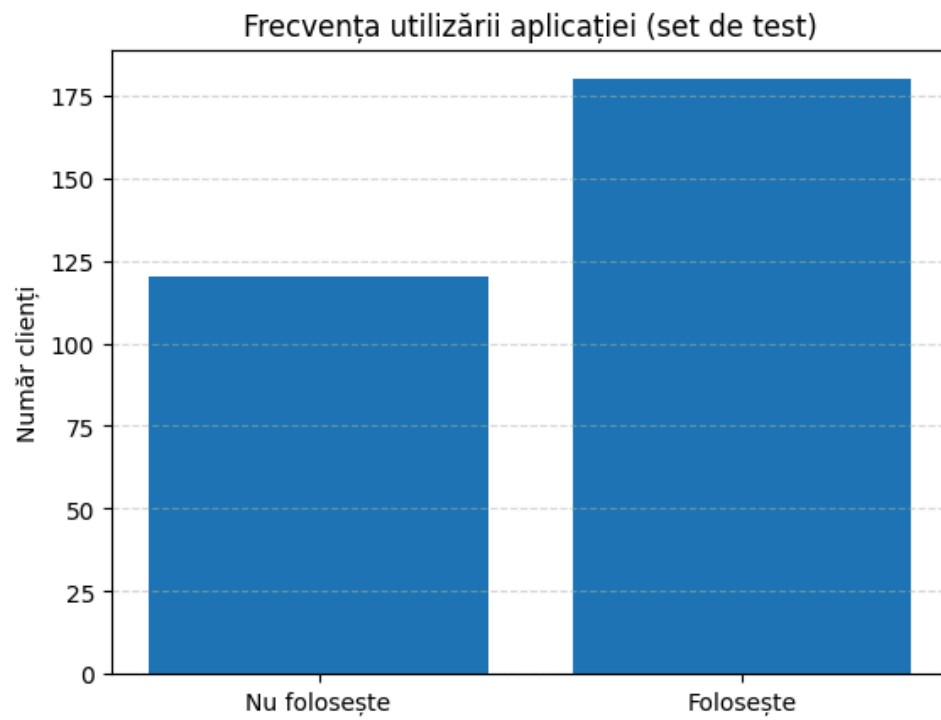
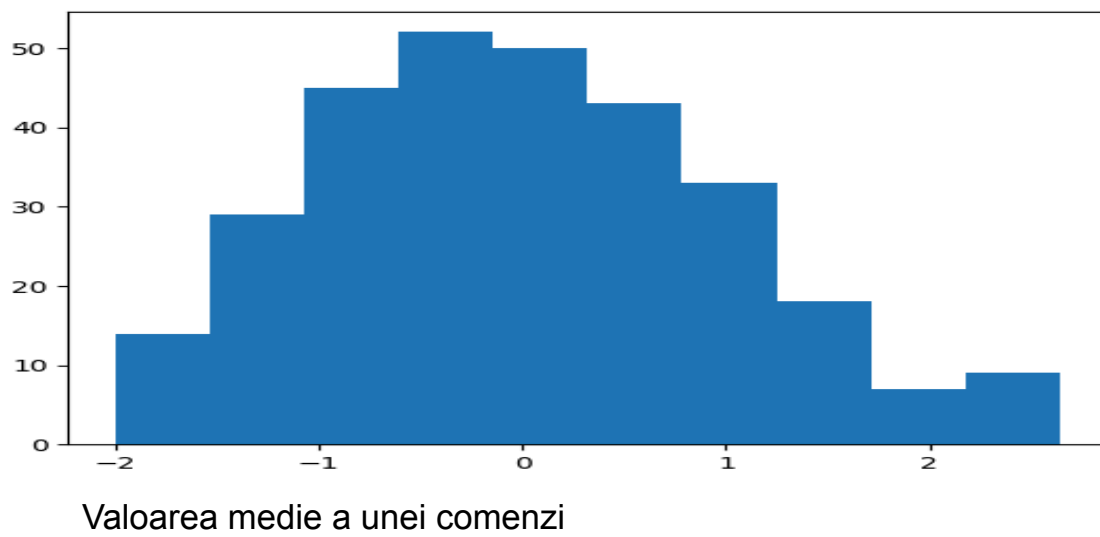
## 2. Analiza exploratie a datelor (EDA complex) dupa aplicarea prelucrarilor

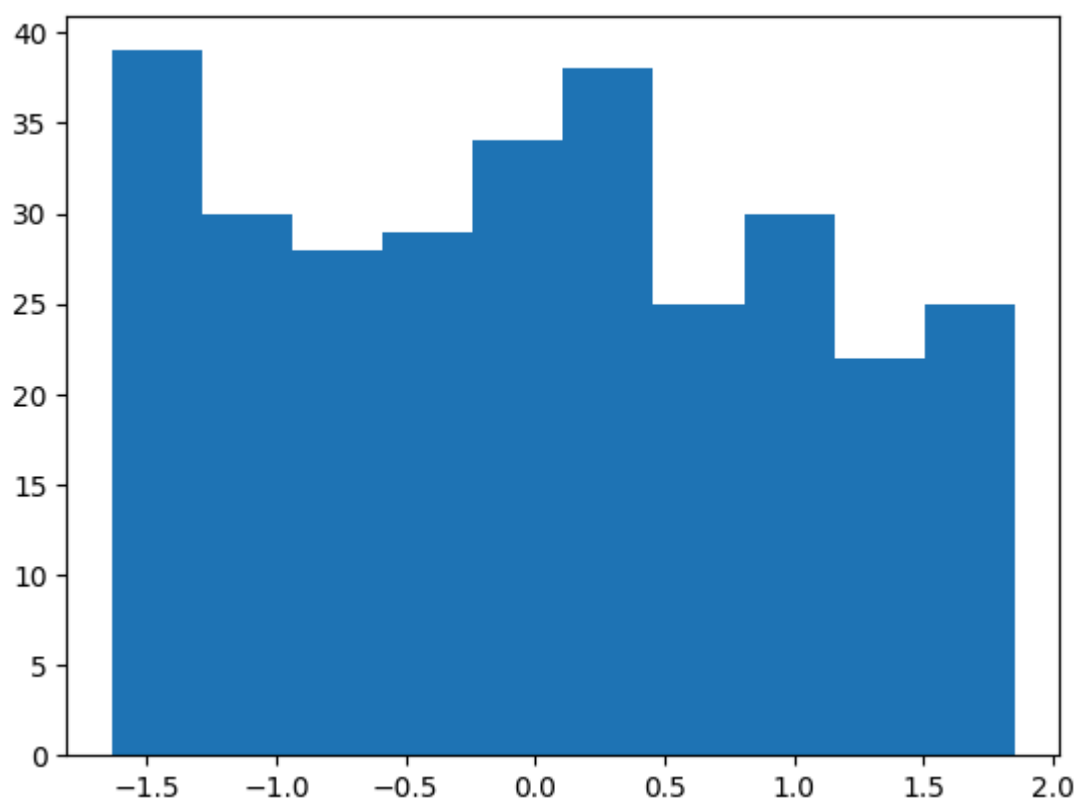
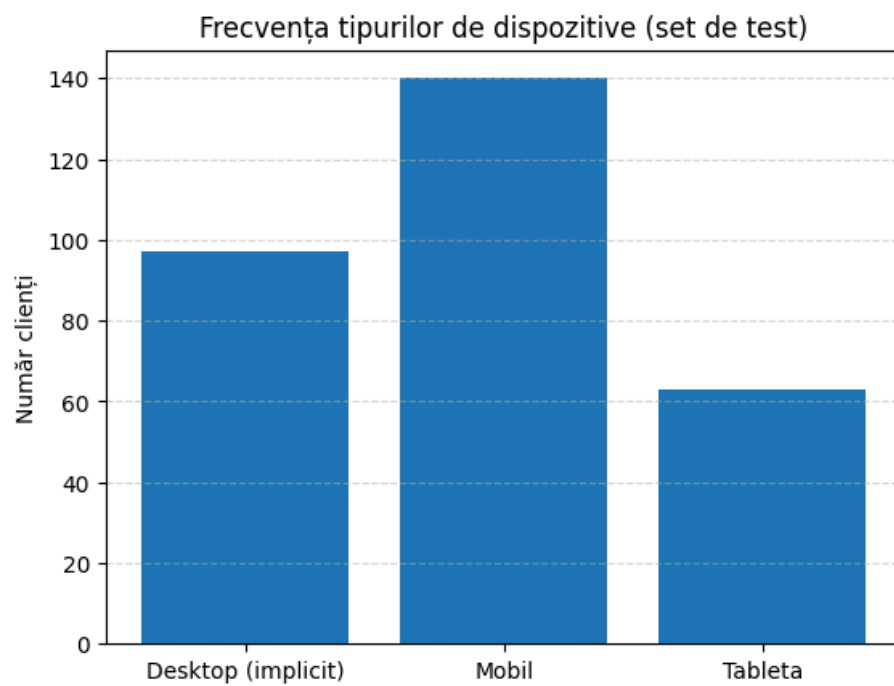
2. a) Analiza valorilor lipsa: Am verificat fiecare coloana pentru a vedea daca exista date lipsa. Am calculat valoarea absoluta si procentul valorilor lipsa. Rezultat: Nu s-au identificat valori lipsa in seturile de date prelucrate.

2. b) Statistici descriptive: Am generat statistici pentru fiecare coloana: media, deviatia standard, valorile minime si maxime etc.

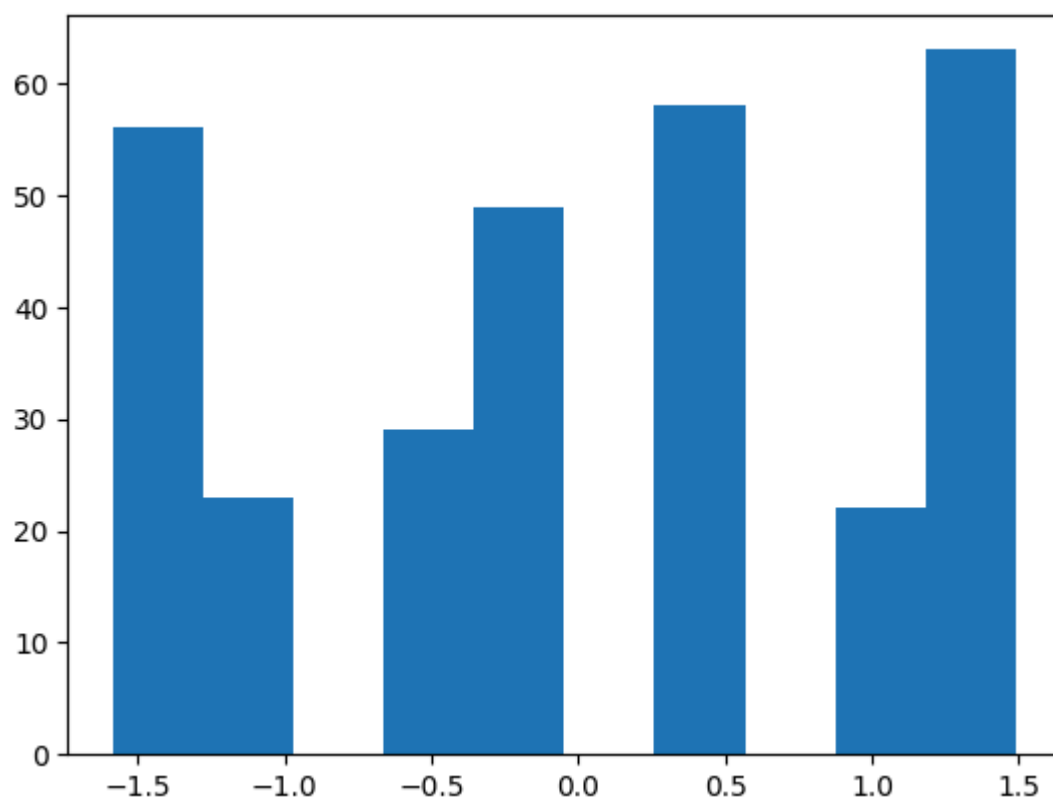
2. c) Analiza distributiei variabilelor: Pentru fiecare coloana am folosit histograme pentru valoare\_comanda\_medie, timp, frecventa\_lunara, etc. si am folosit grafice bara pentru variabile care sunt categorice la baza (aplicatia, dispozitiv, cupoane\_de\_reducere). Am reconstruit si valorile eliminate prin

One-Hot Encoding (ex: dispozitiv\_desktop).

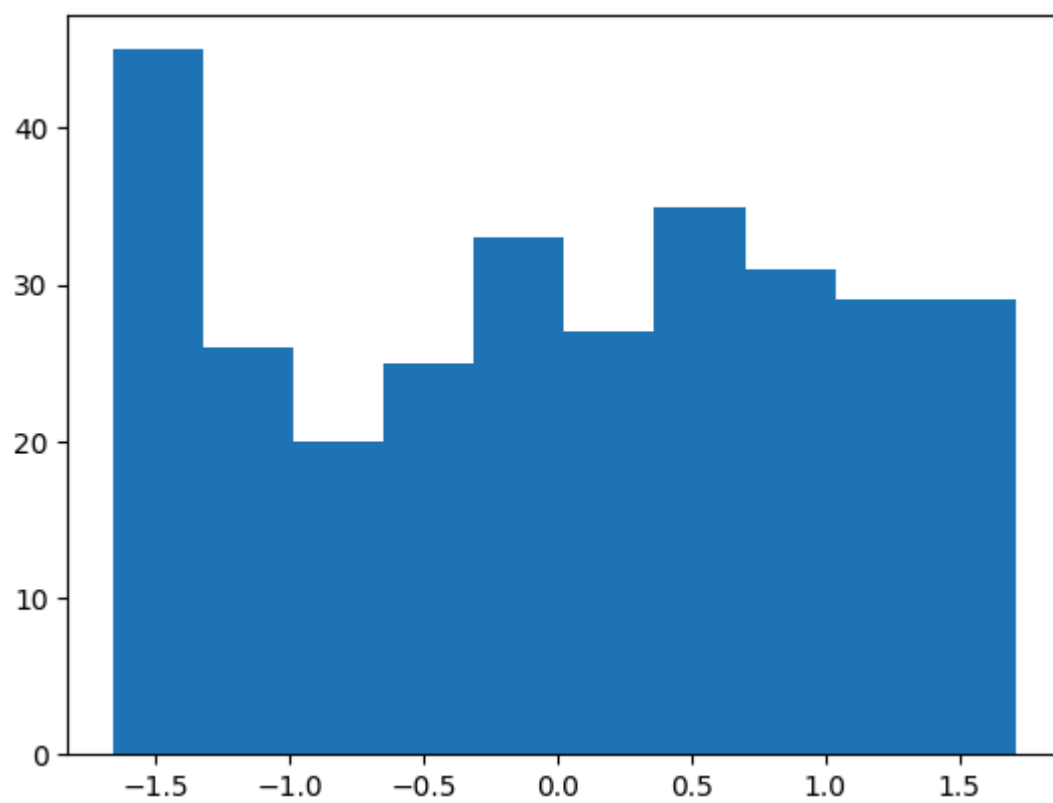




Cat timp petrece pe magazin

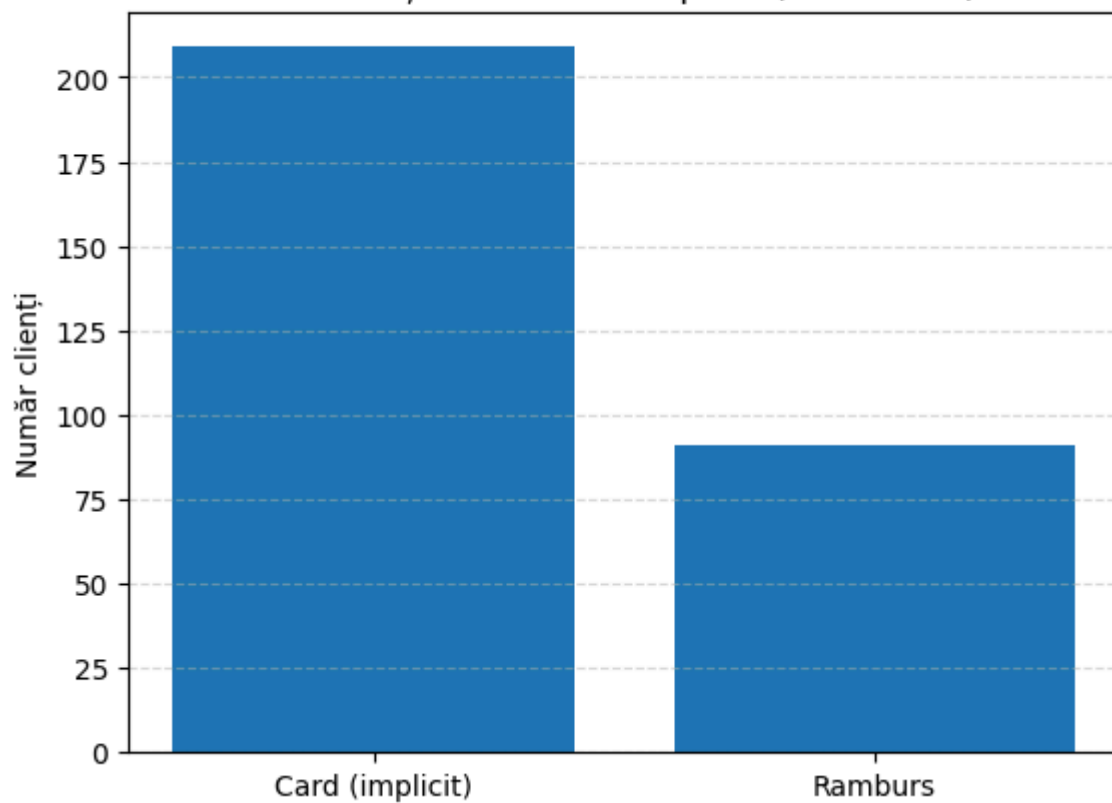


Frecventa lunara

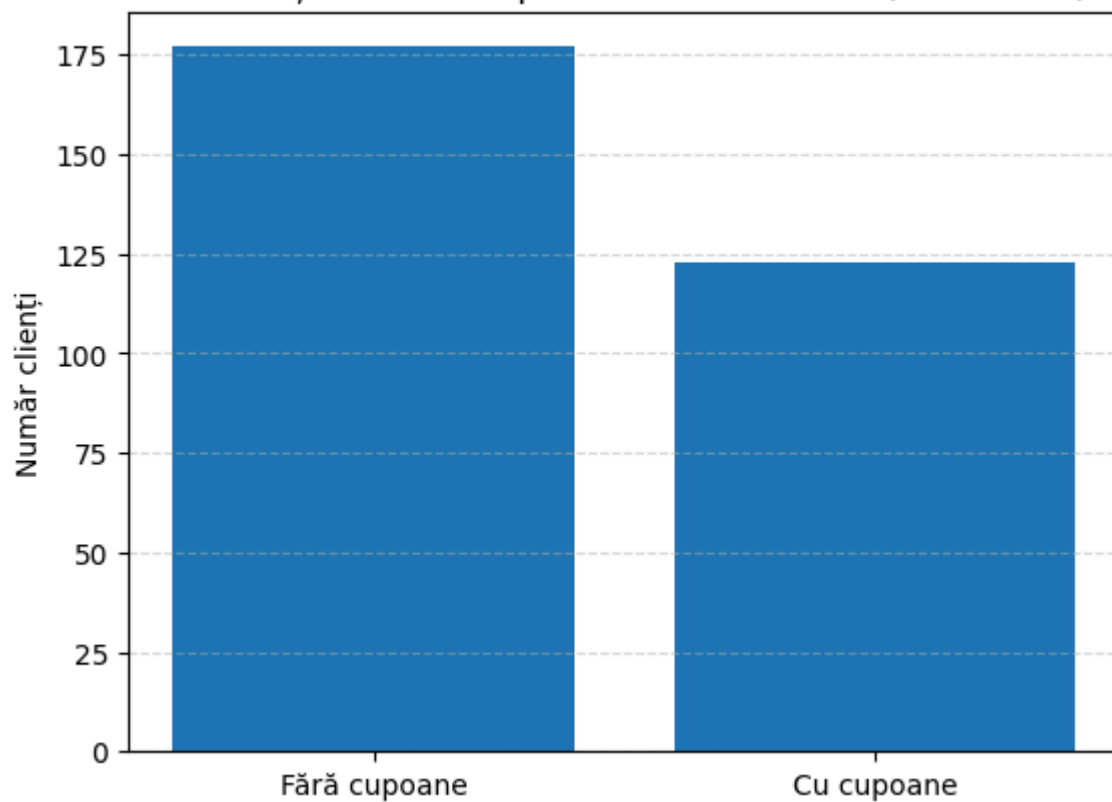


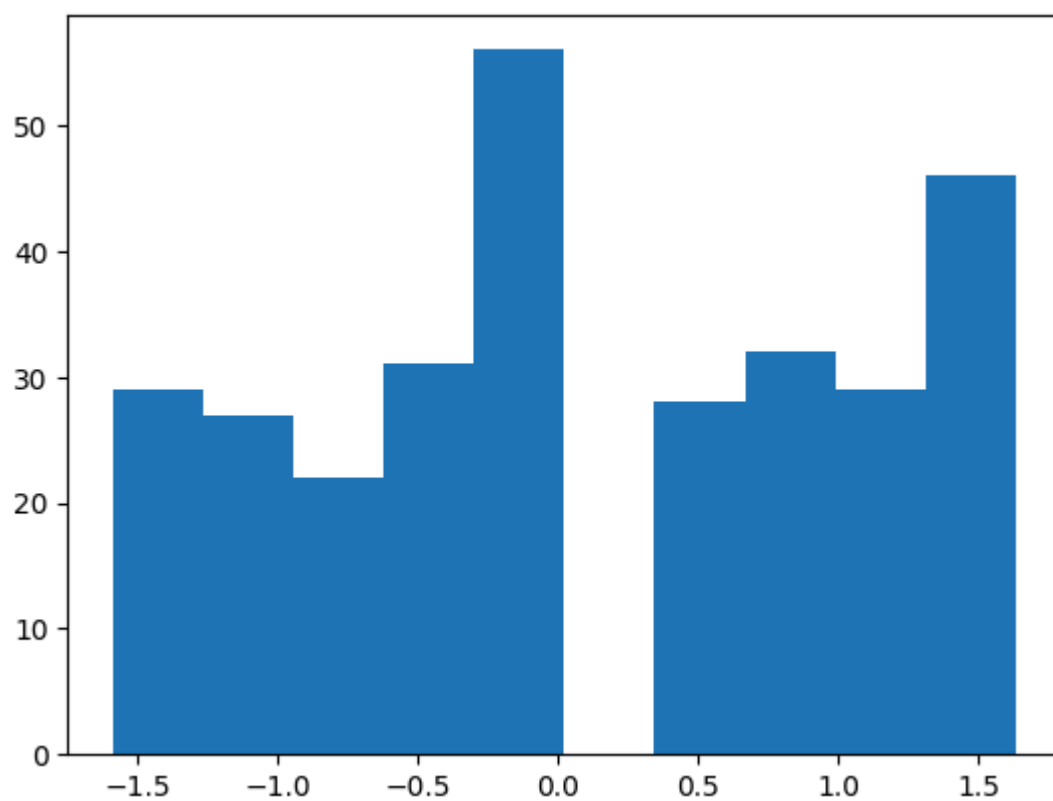
Varsta

Frecvența metodelor de plată (set de test)

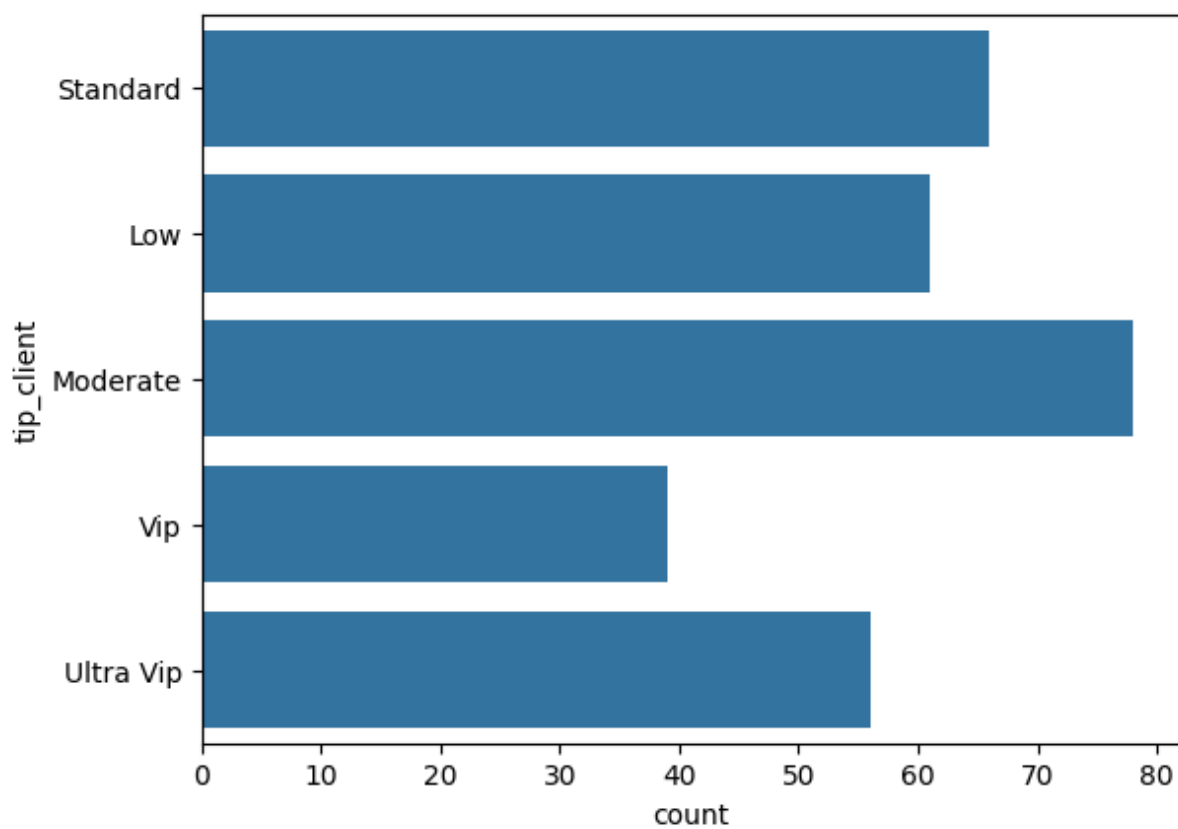


Frecvența utilizării cupoanelor de reducere (set de test)



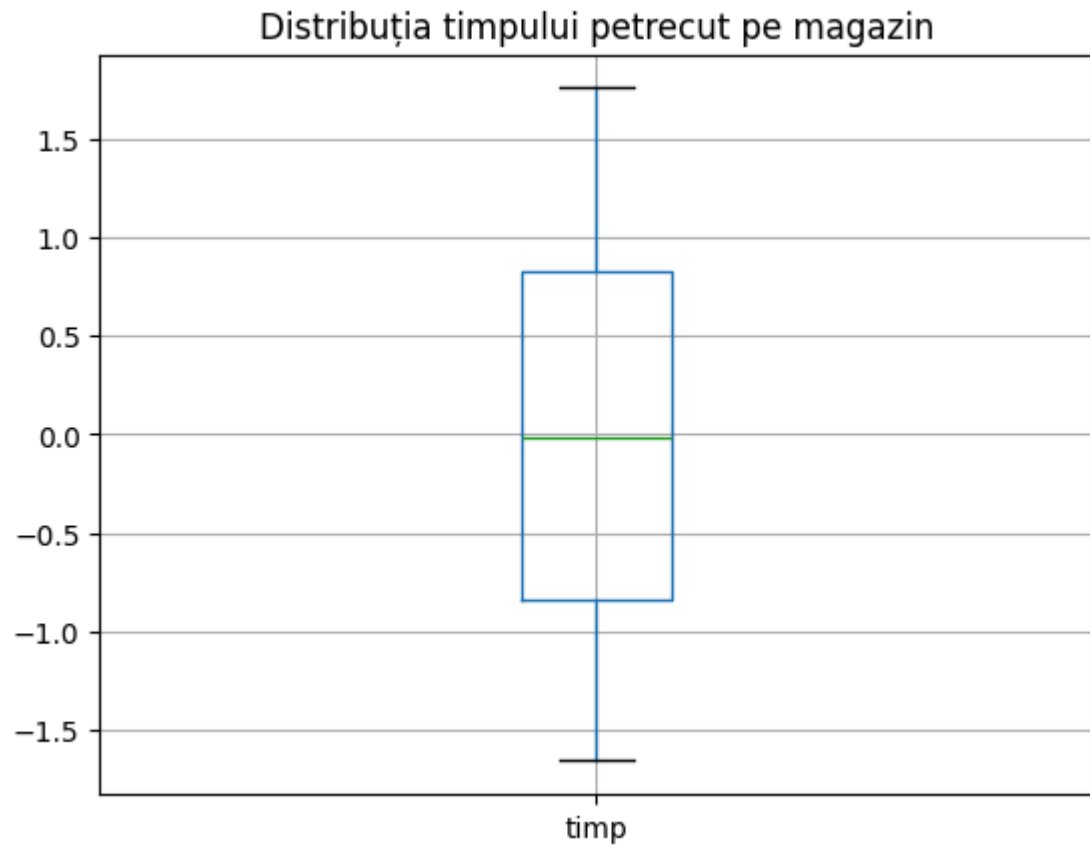


Cate recenzii acorda

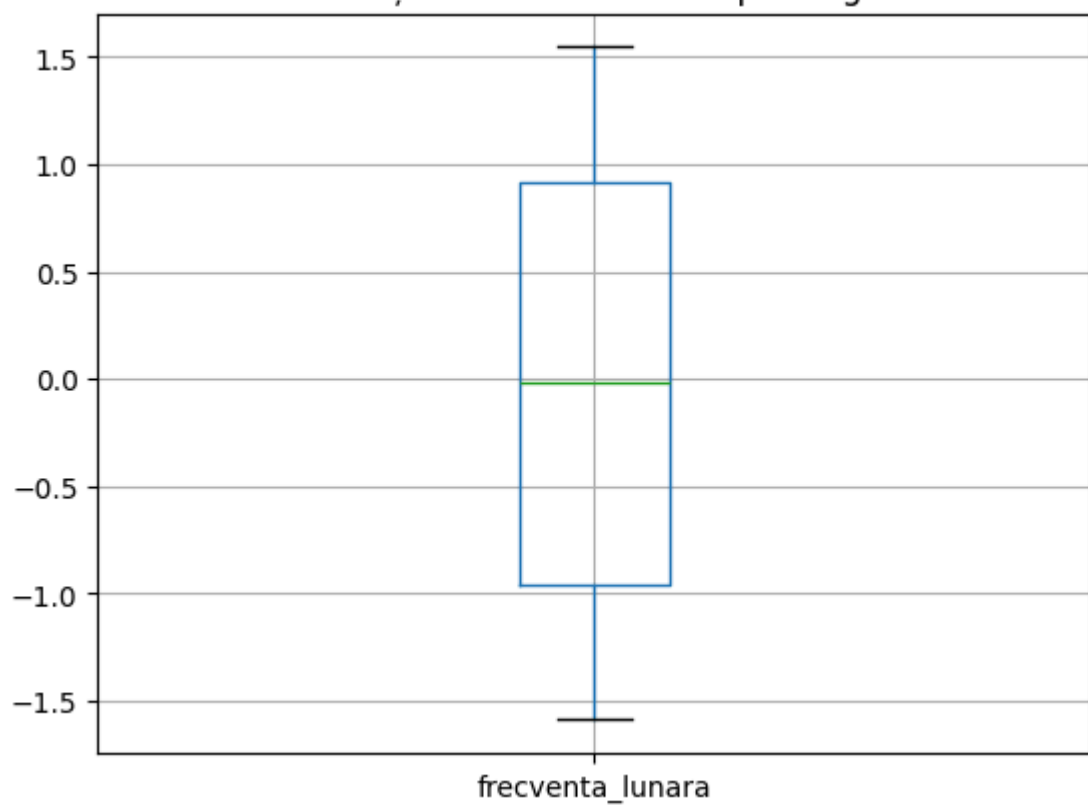


Tipuri de clienti

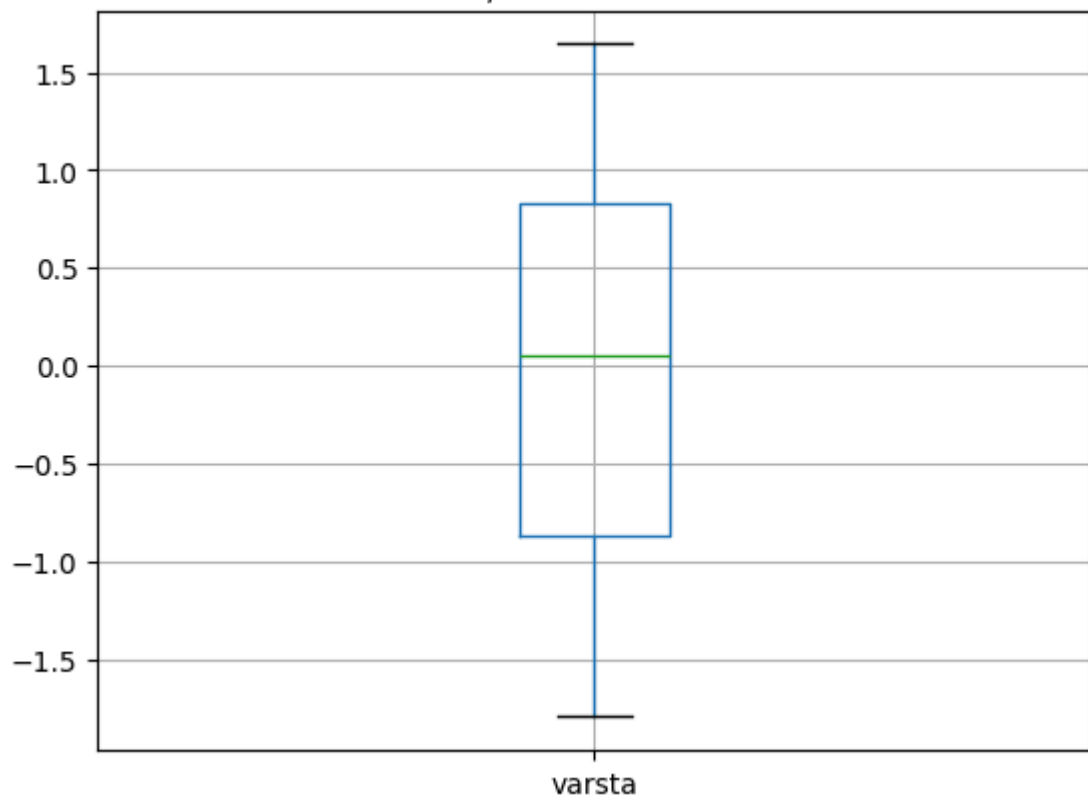
2. d) Detectarea outlierilor: Am utilizat boxplot-uri pentru a vizualiza distribuția fiecărei variabile numerice. Am analizat `valoare_comanda_medie`,  `timp`, `frecventa_lunara`, `varsta`, `recenzii`.



Distribuția frecvenței lunare pe magazin

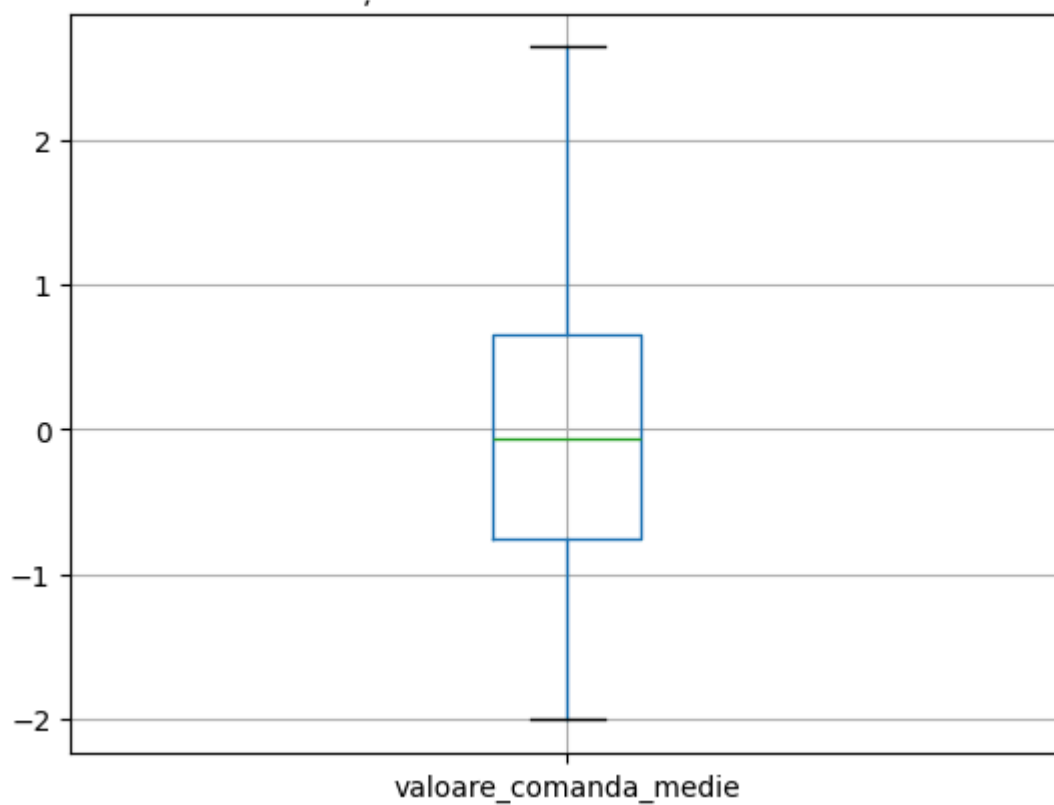


Distribuția varstelor clientilor

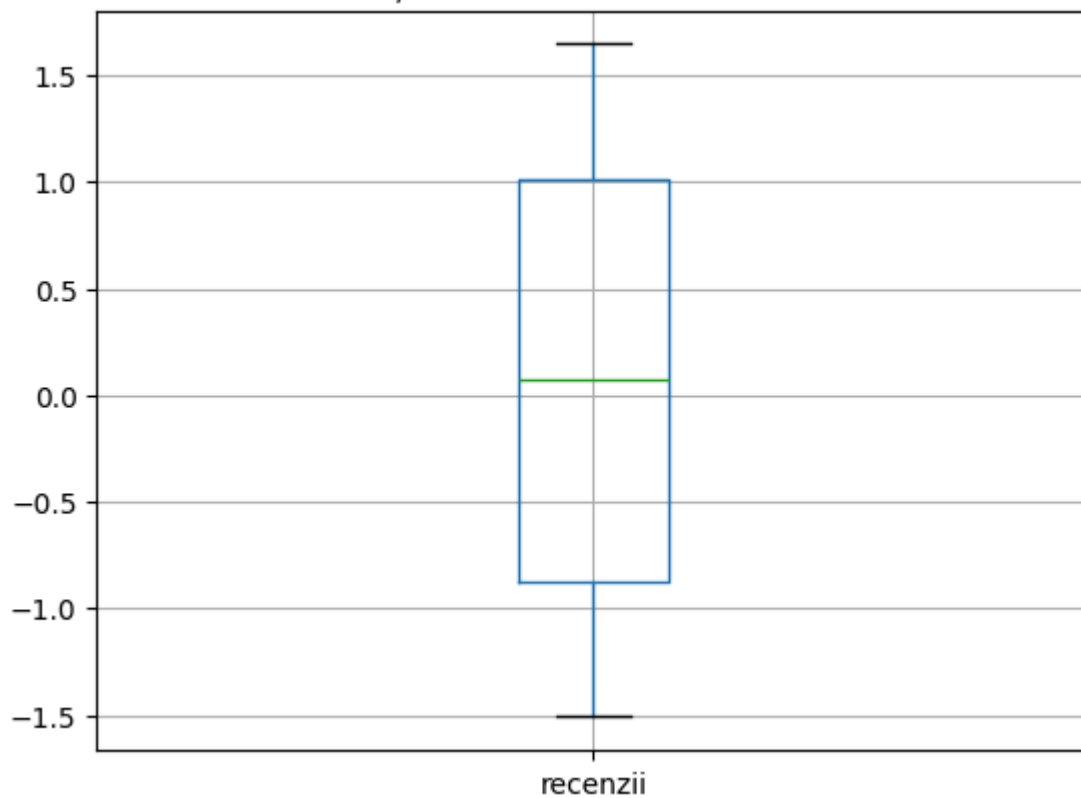




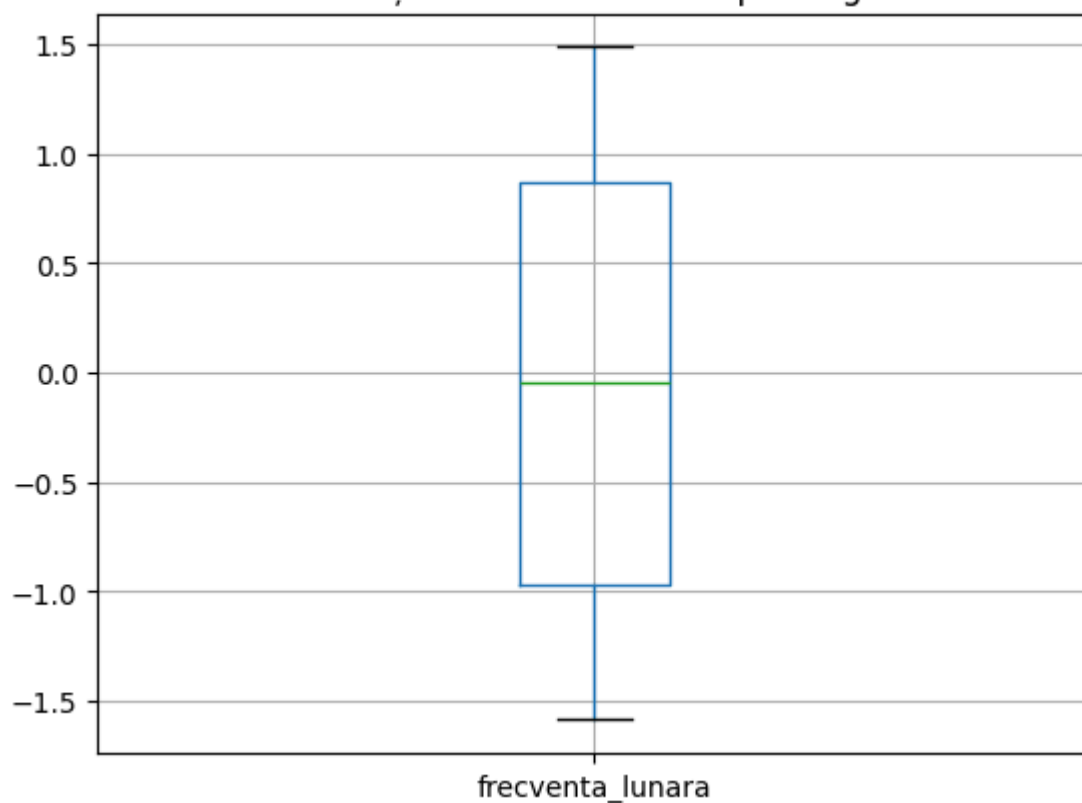
Distribuția valorilor medii ale comenzilor



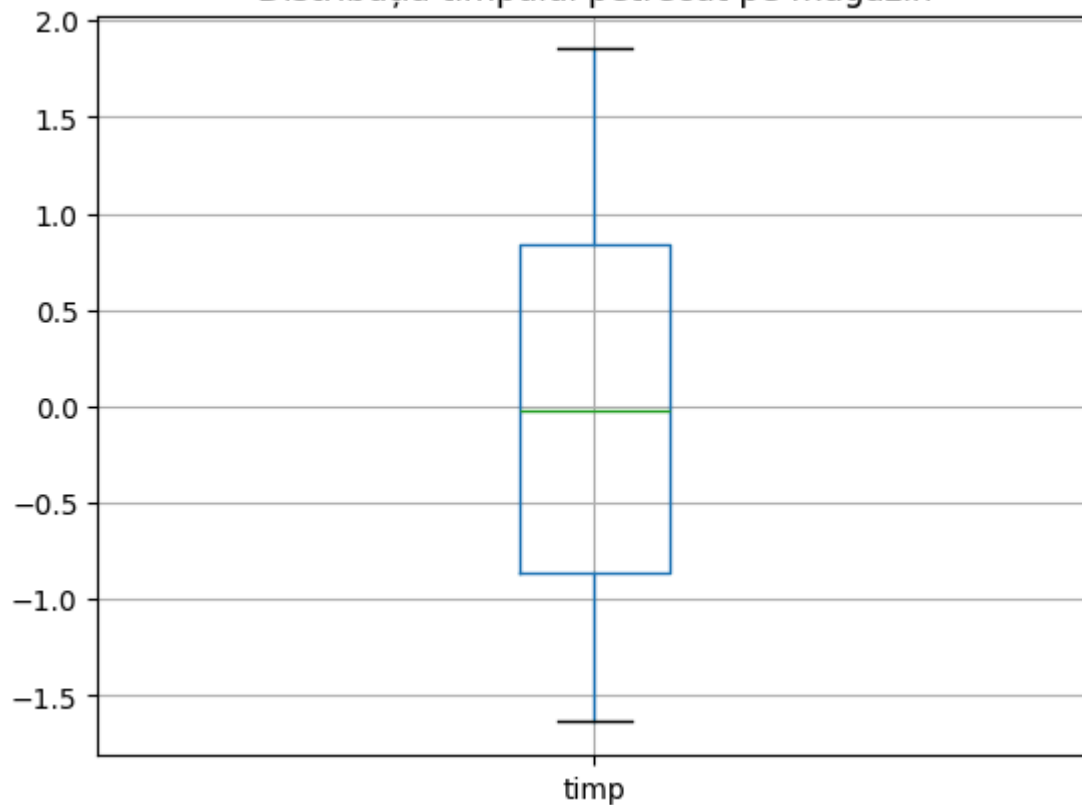
Distribuția recenziilor acordate de client

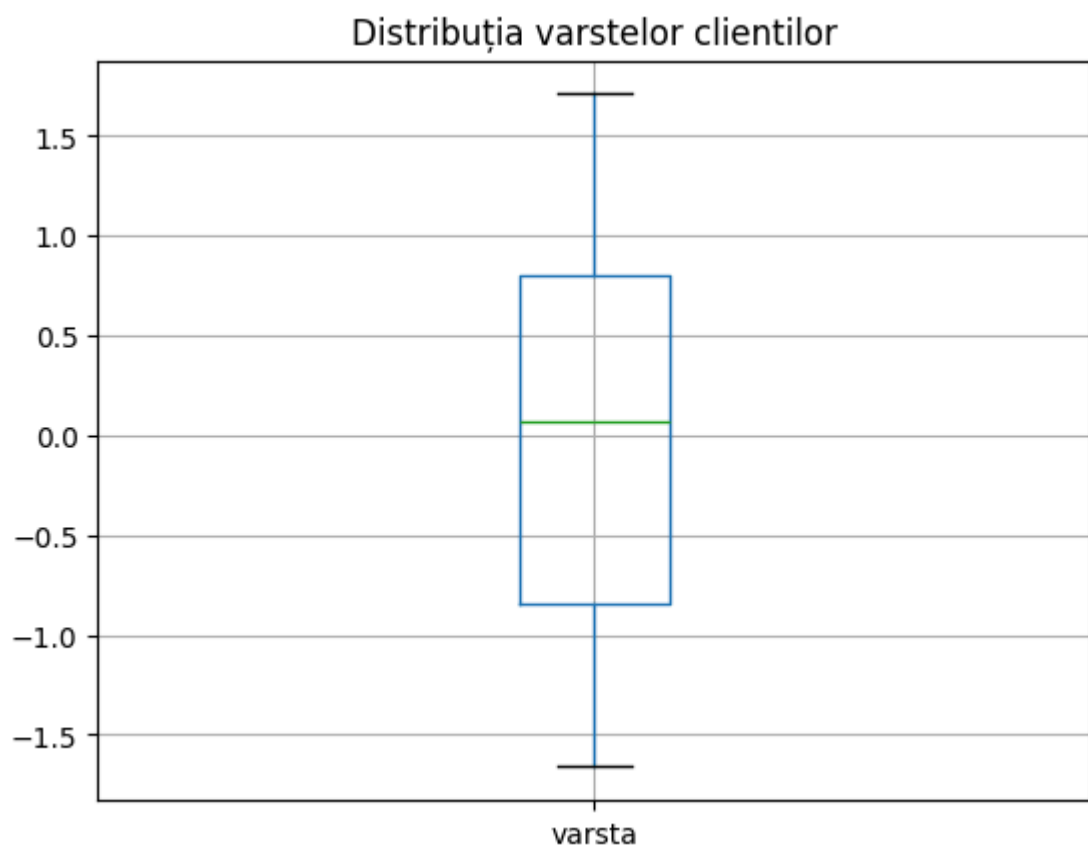
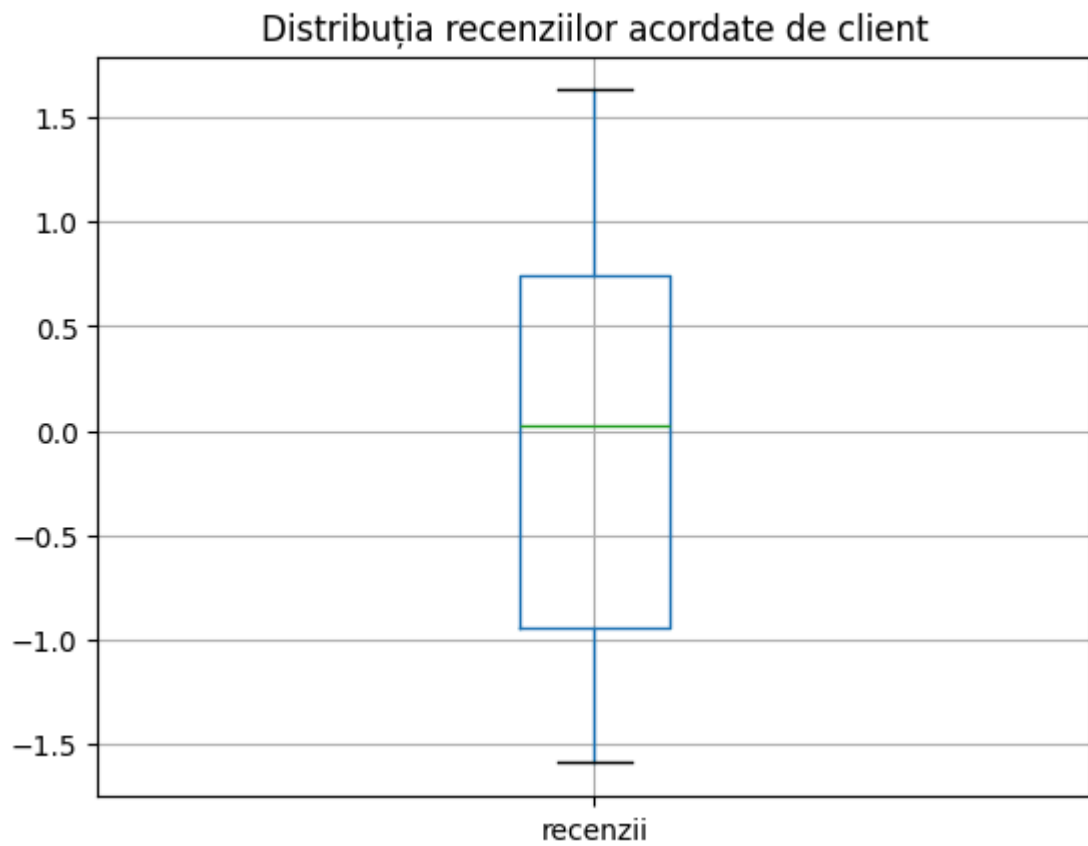


Distribuția frecvenței lunare pe magazin

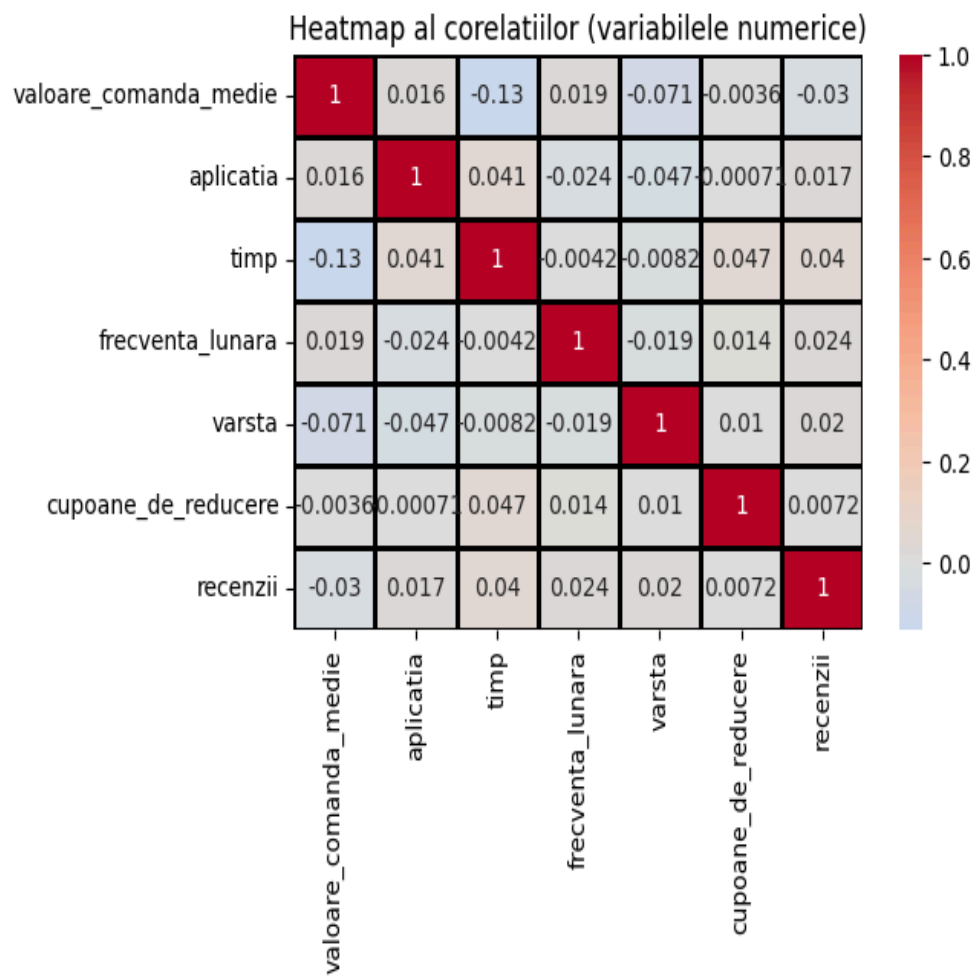


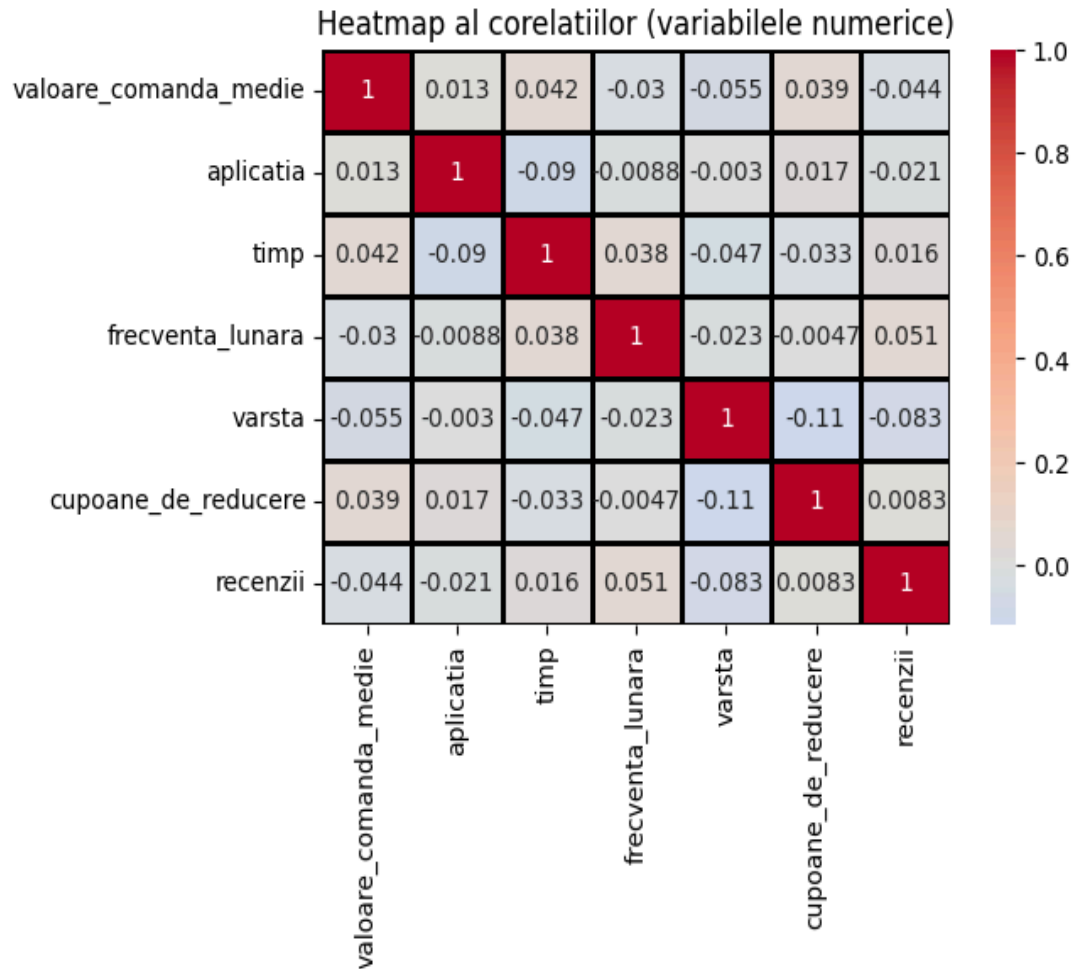
Distribuția timpului petrecut pe magazin





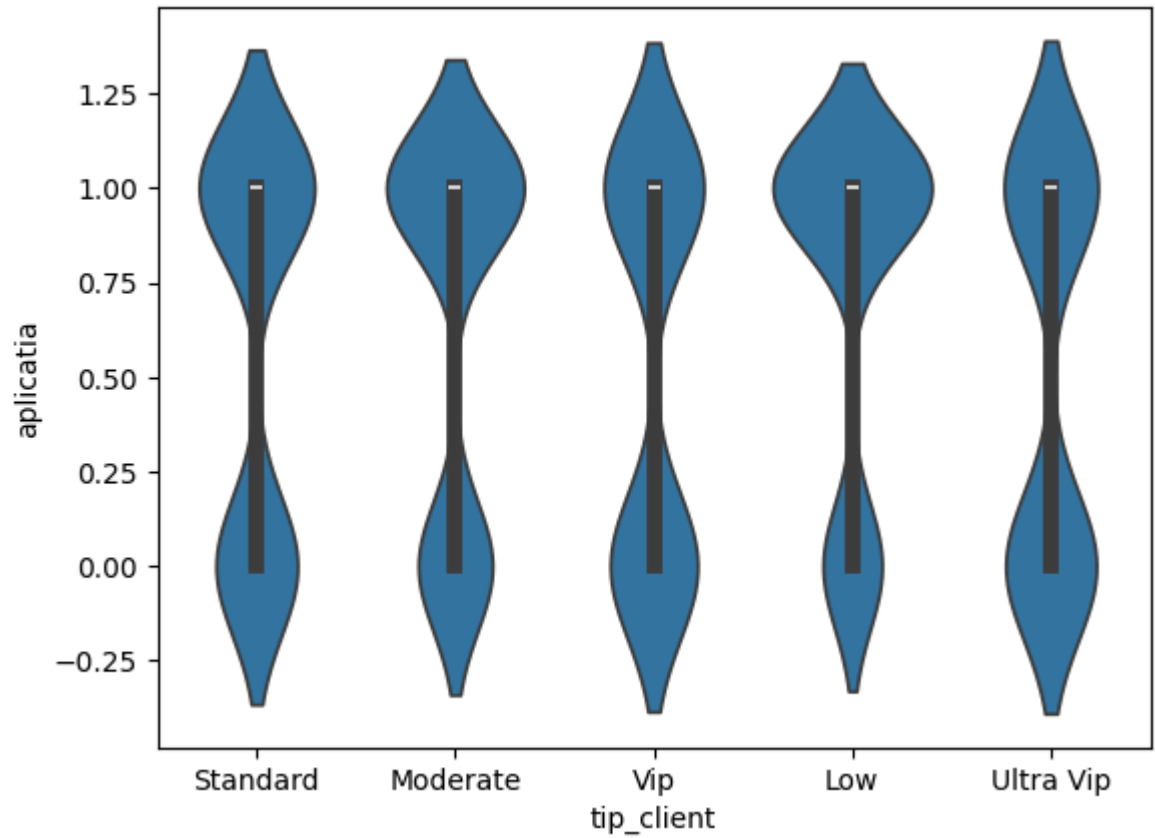
2 .e) Analiza corelatiilor: Am generat un heatmap cu coeficientii de corelatie.



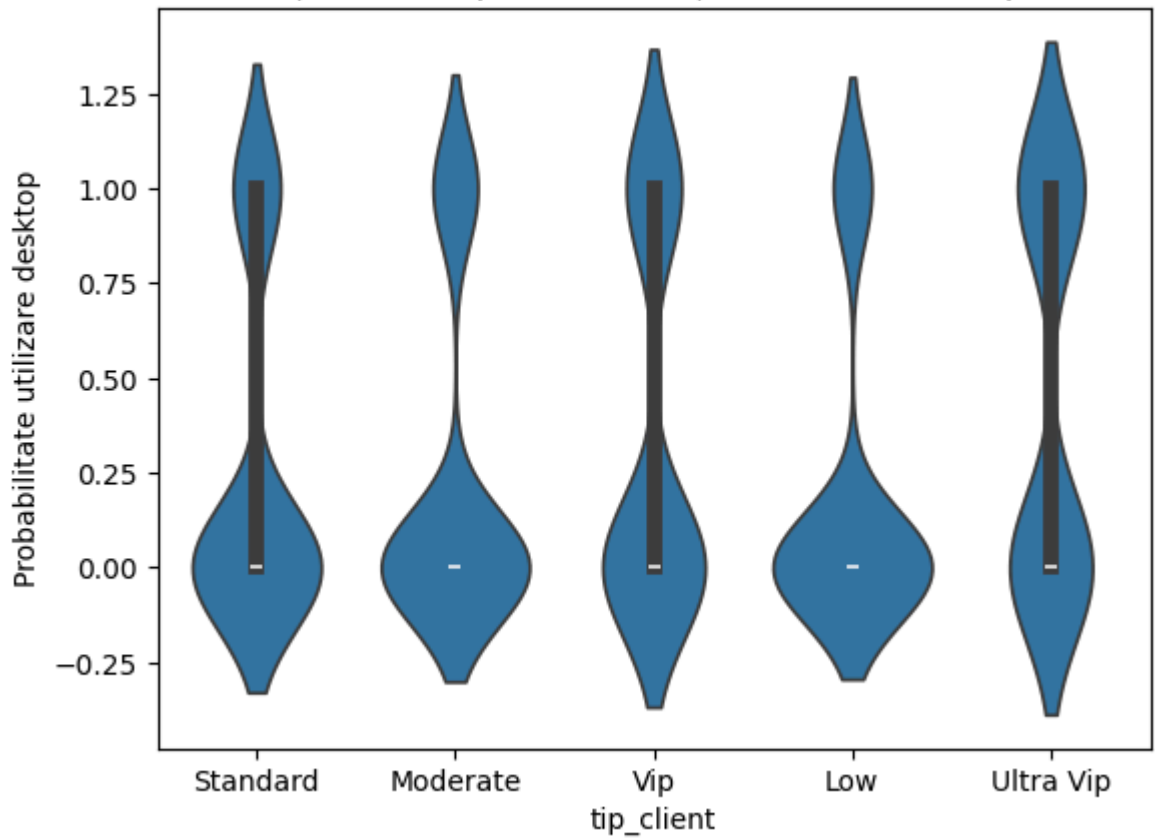


2. f) Relatii între trasaturi și coloana tinta (tip\_client): Pentru a analiza mai profund legătura dintre fiecare trasatură și tip\_client, am folosit Violin plots (pentru a vedea distribuțiile în funcție de clasă) și am inclus atât datele de train, cât și cele de test pentru comparație (relatia dintre tip\_client și utilizarea aplicației, relatia dintre tip\_client și varsta/frecvența comenzilor, impactul metodei de plată sau a cupoanelor de reducere).

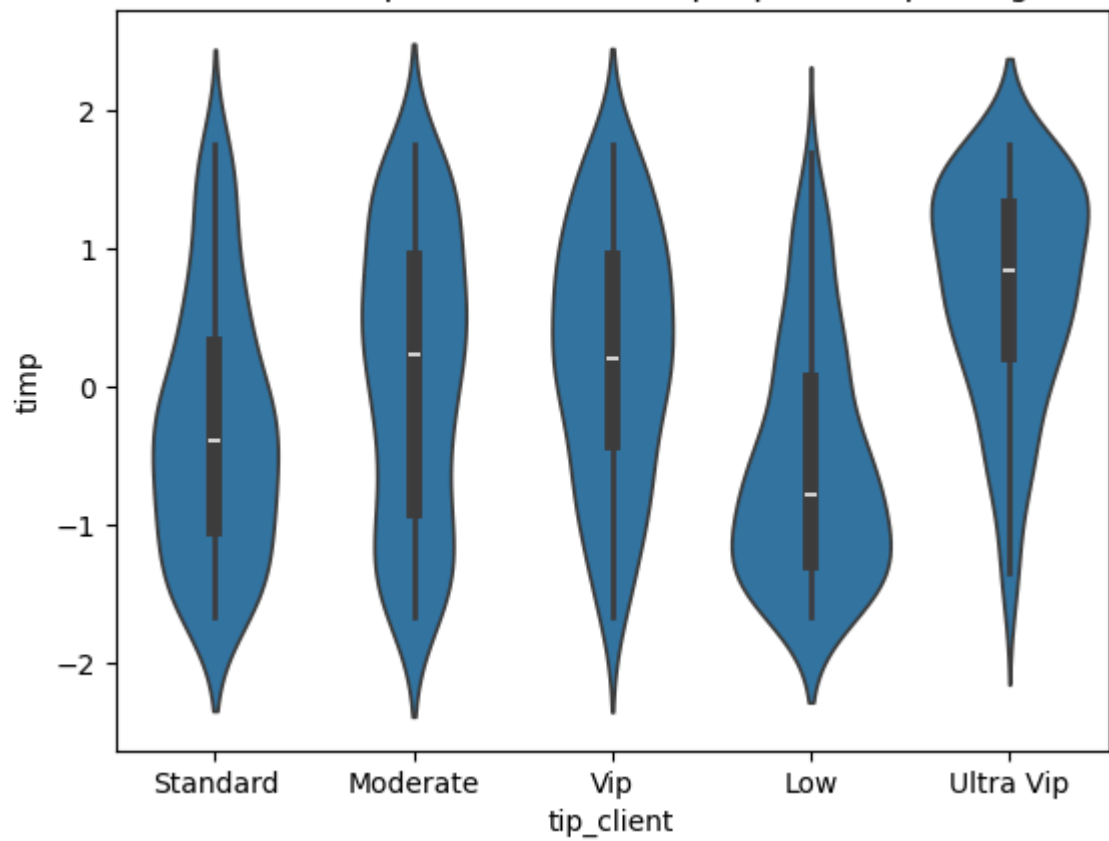
Relatia dintre tipul de client si aplicatia magazinului



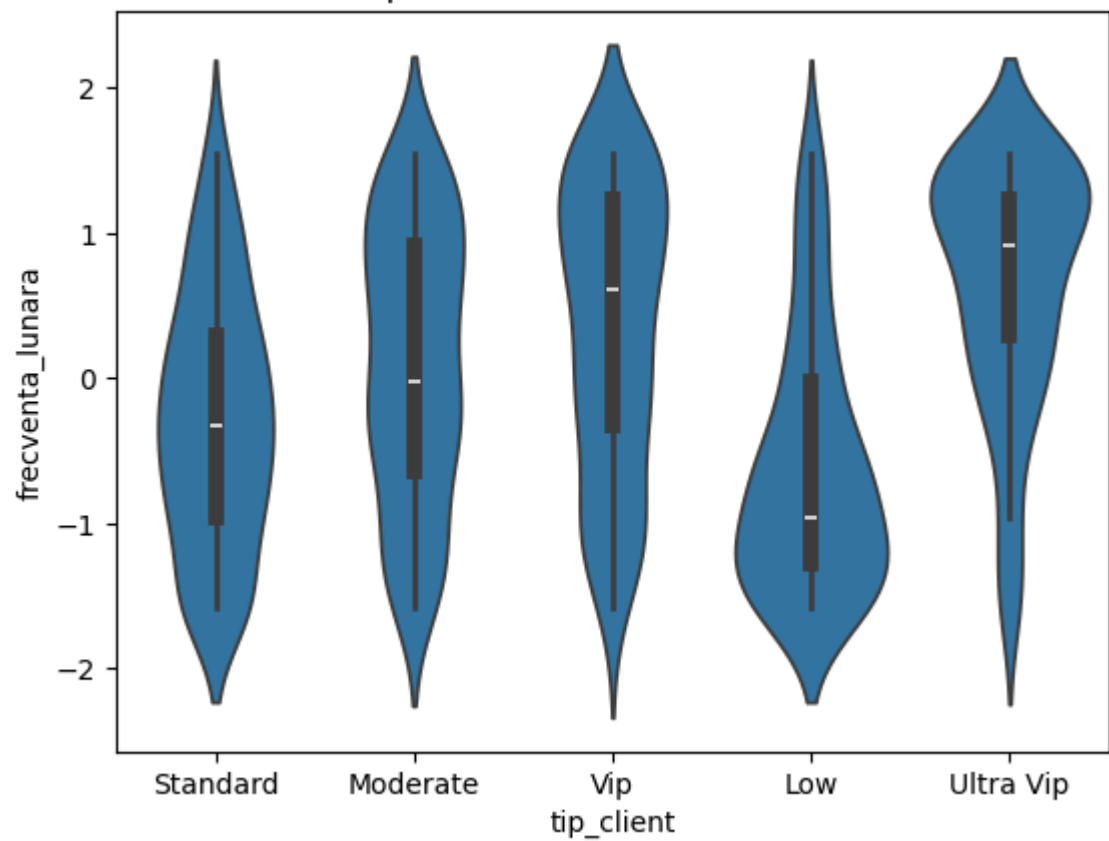
Relația dintre tipul de client și utilizarea desktopului



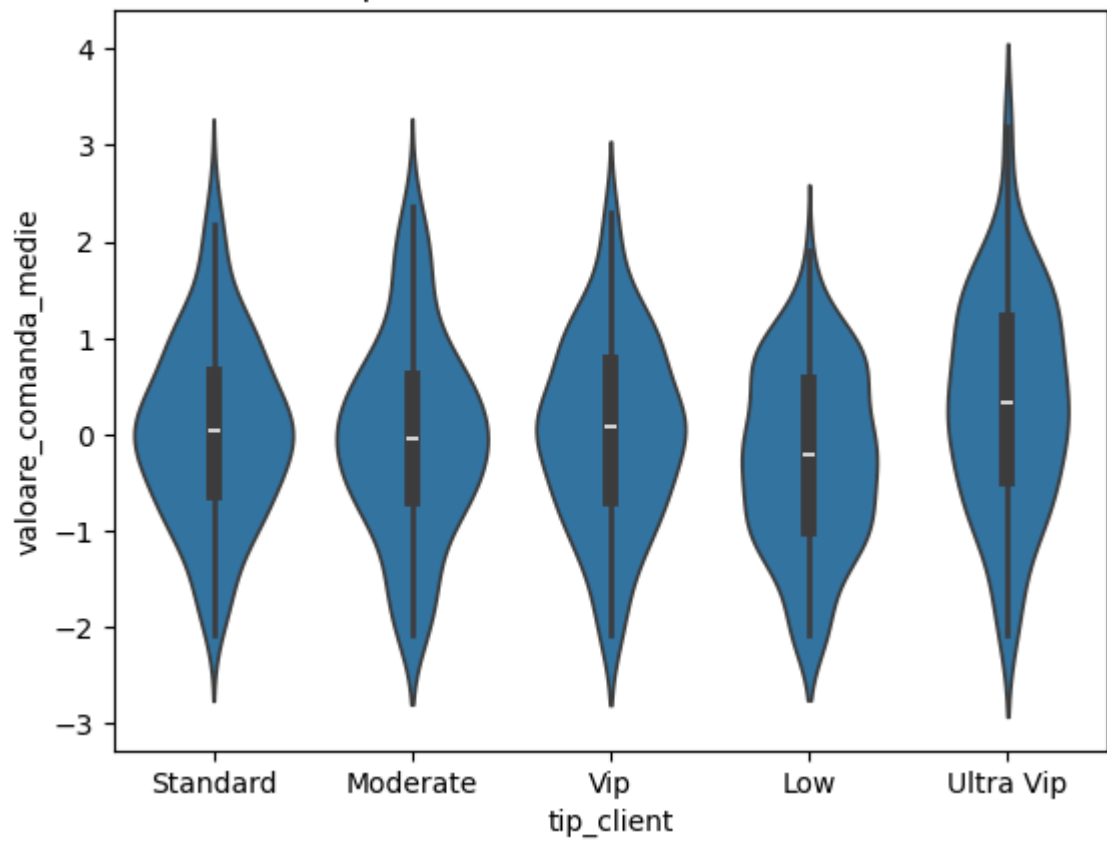
Relatia dintre tipul de client si timpul petrecut pe magazin



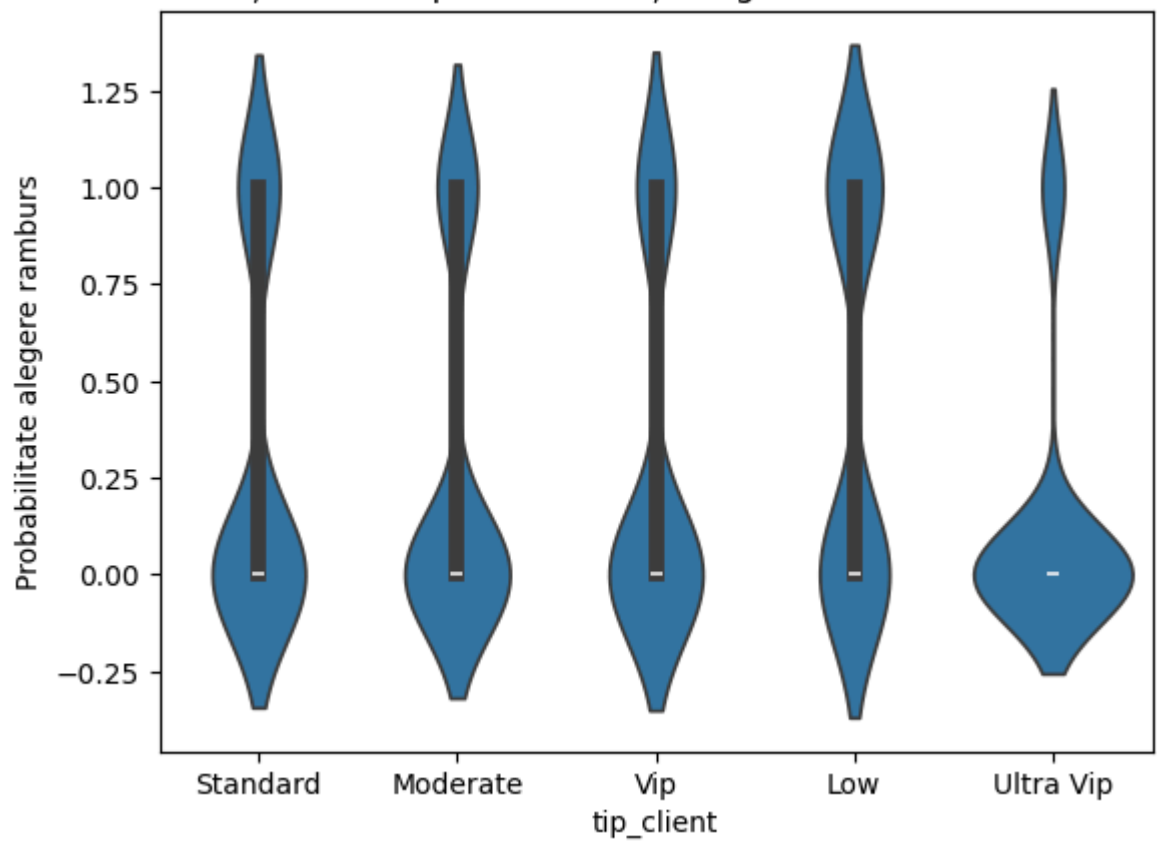
Relatia dintre tipul de client si frecventa lunara a acestuia



Relatia dintre tipul de client si valoarea medie a unei comenzi

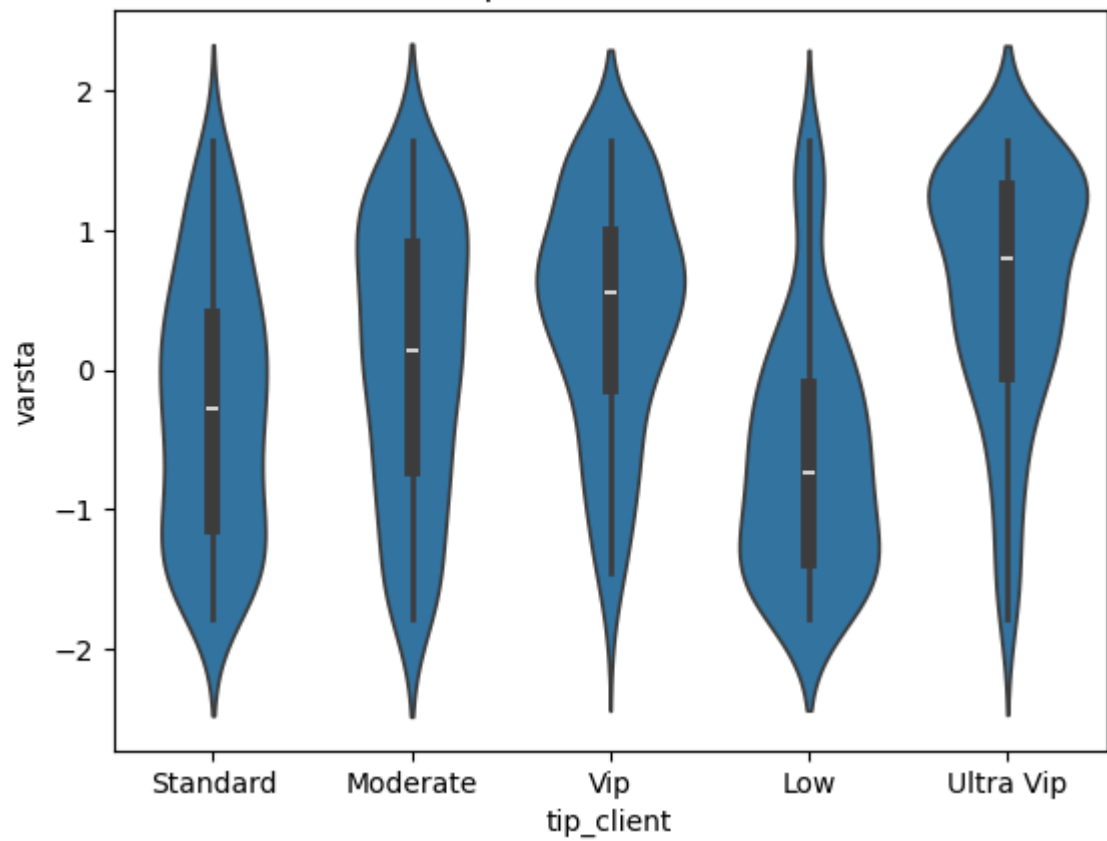


Relația dintre tipul de client și alegerea metodei Ramburs

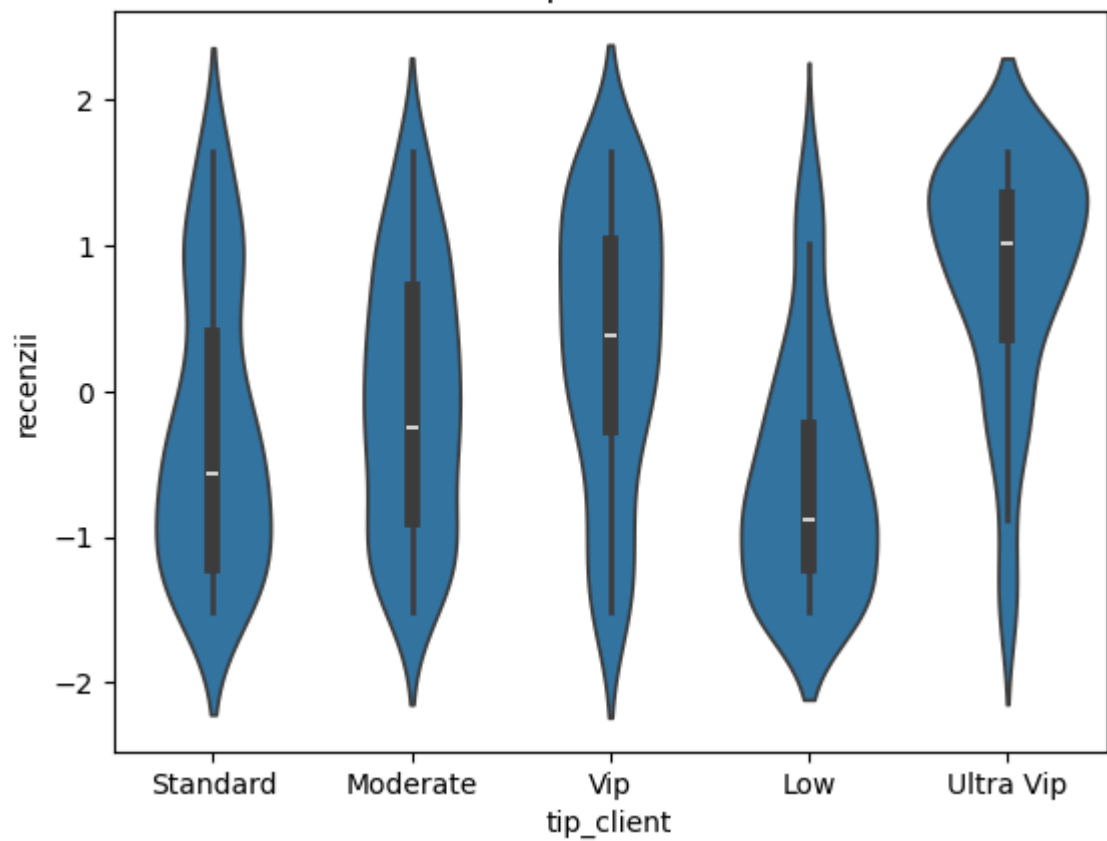




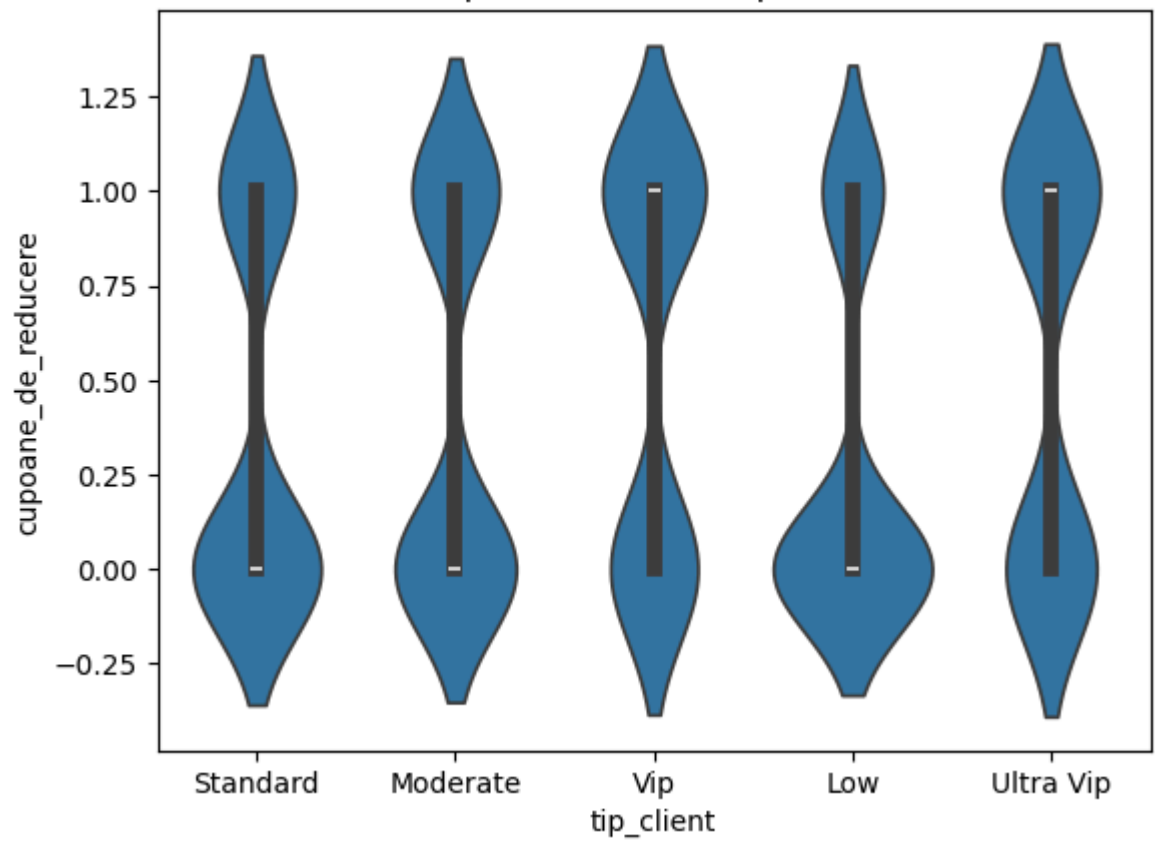
Relatia dintre tipul de client si varsta acestuia



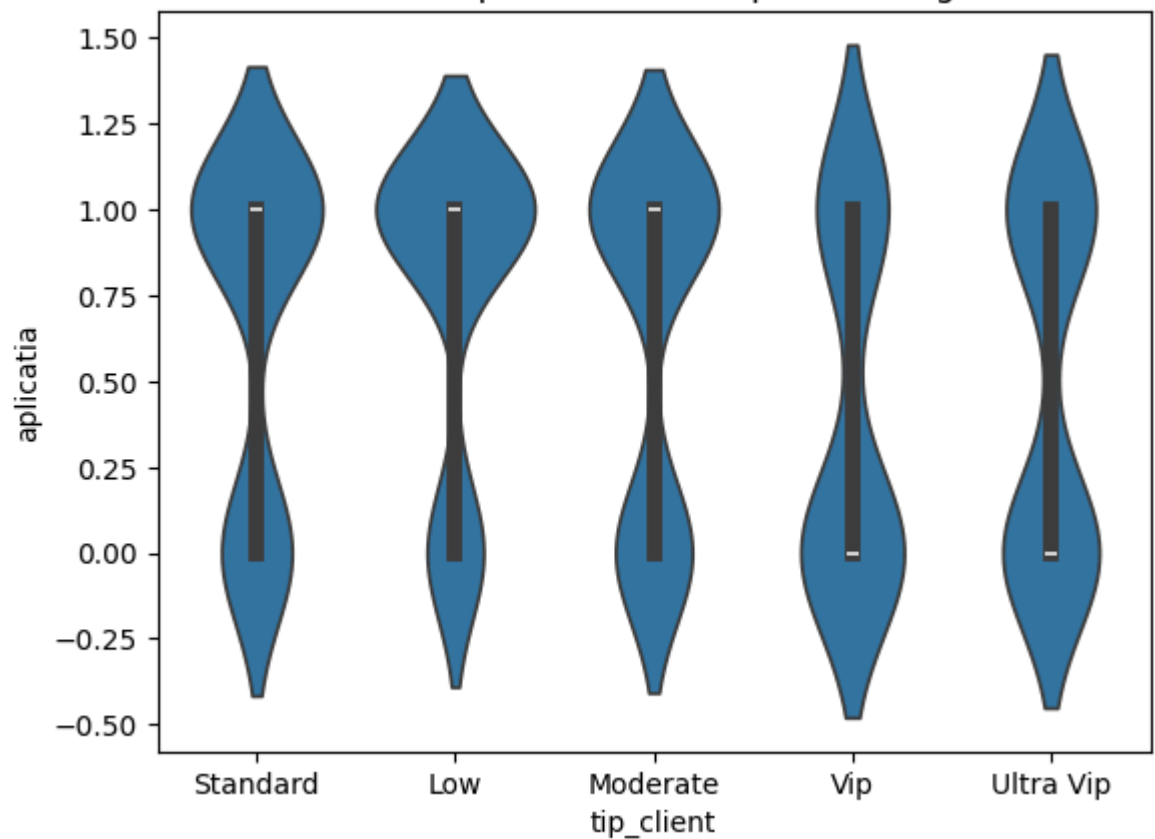
Relatia dintre tipul de client si recenzii



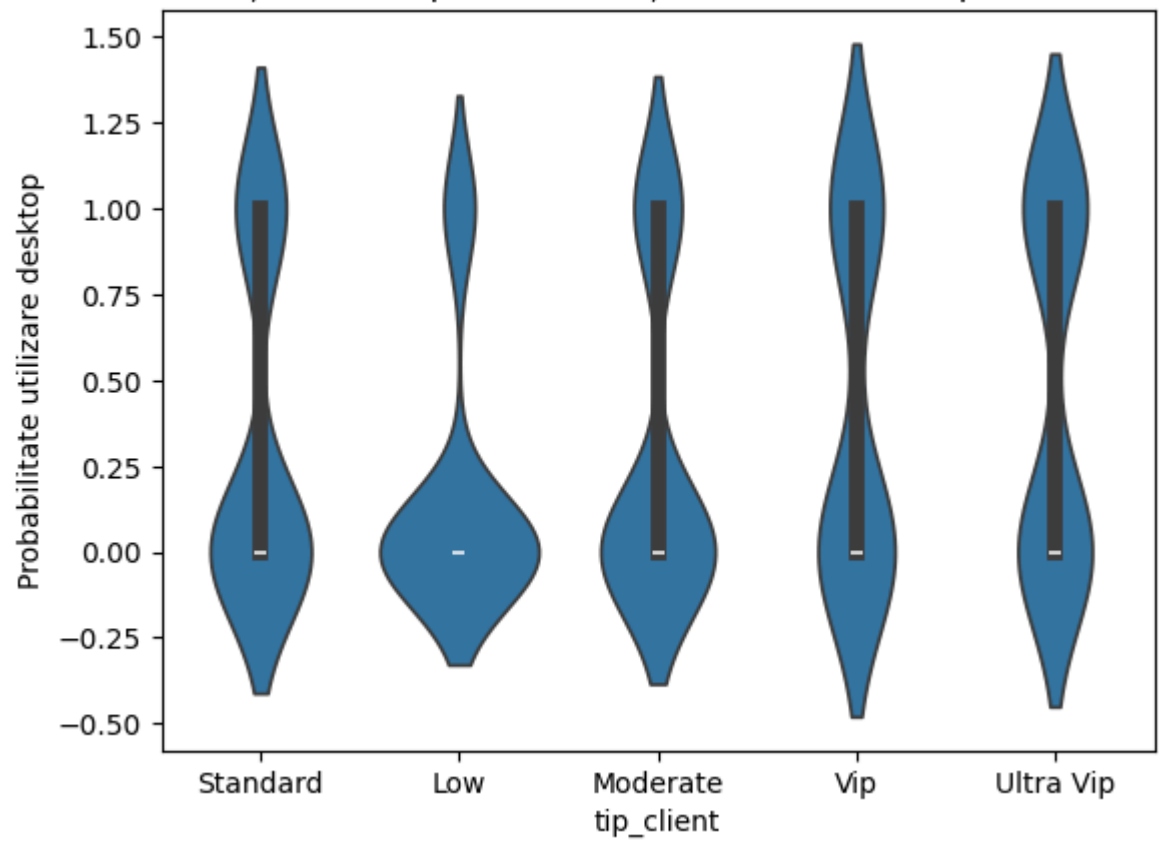
Relatia dintre tipul de client si cupoanele de reducere



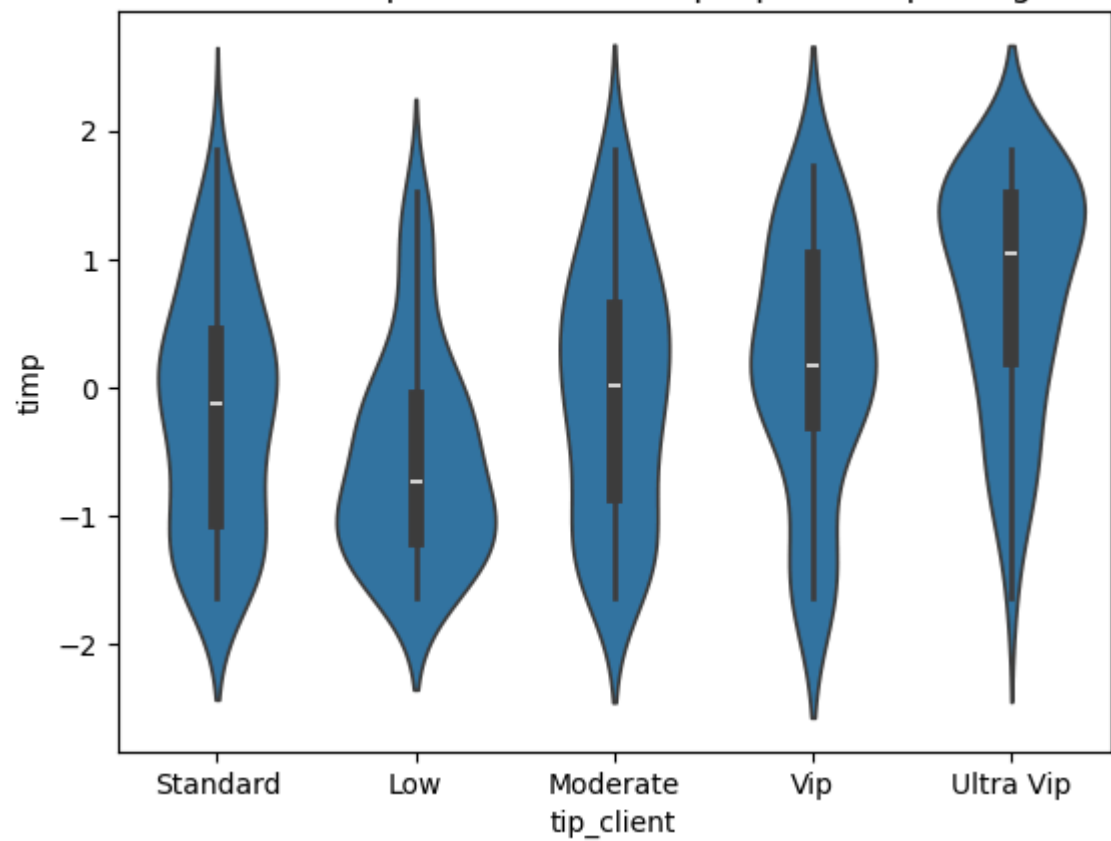
Relatia dintre tipul de client si aplicatia magazinului



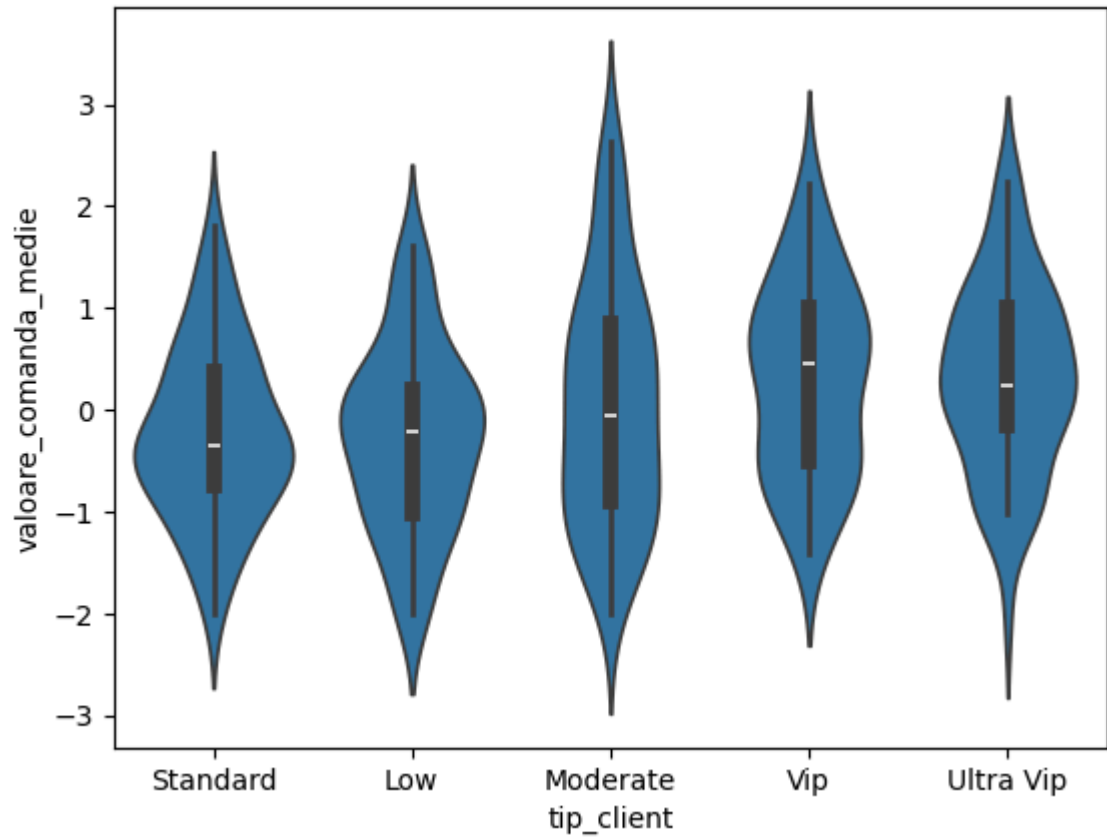
Relația dintre tipul de client și utilizarea desktopului (test)



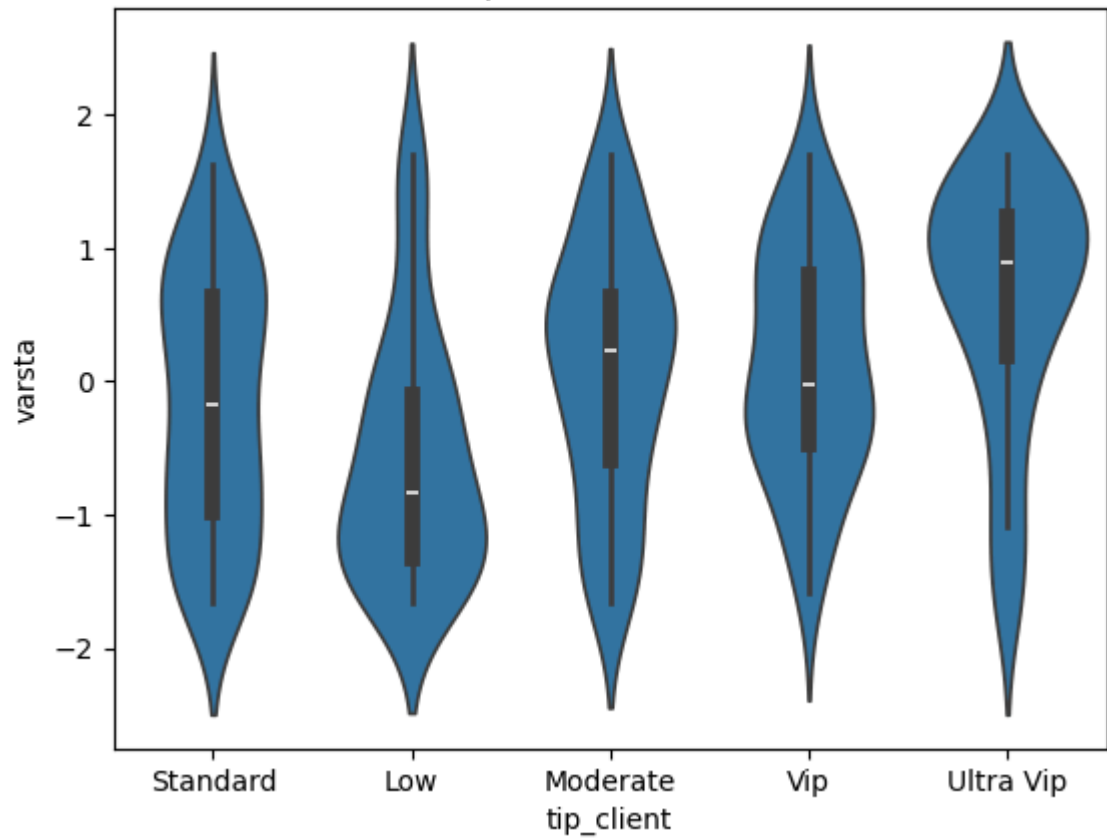
Relatia dintre tipul de client si timpul petrecut pe magazin



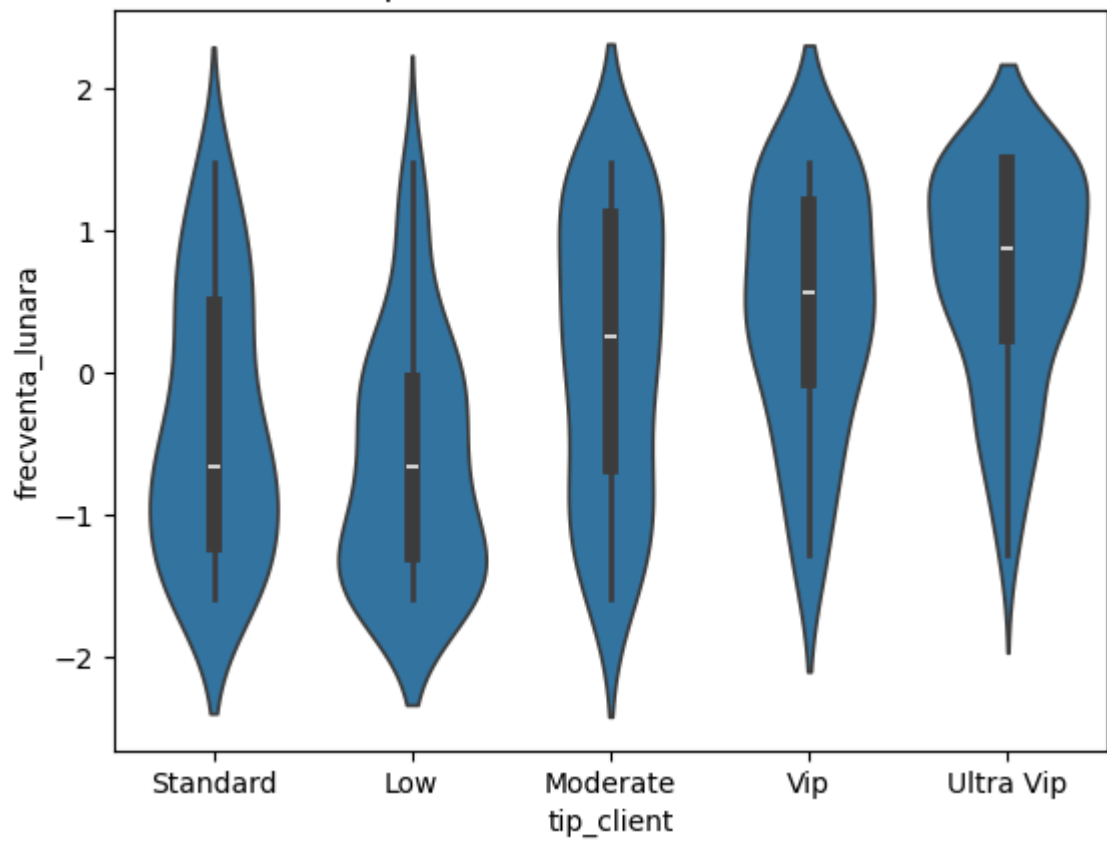
Relatia dintre tipul de client si valoarea medie a unei comezi



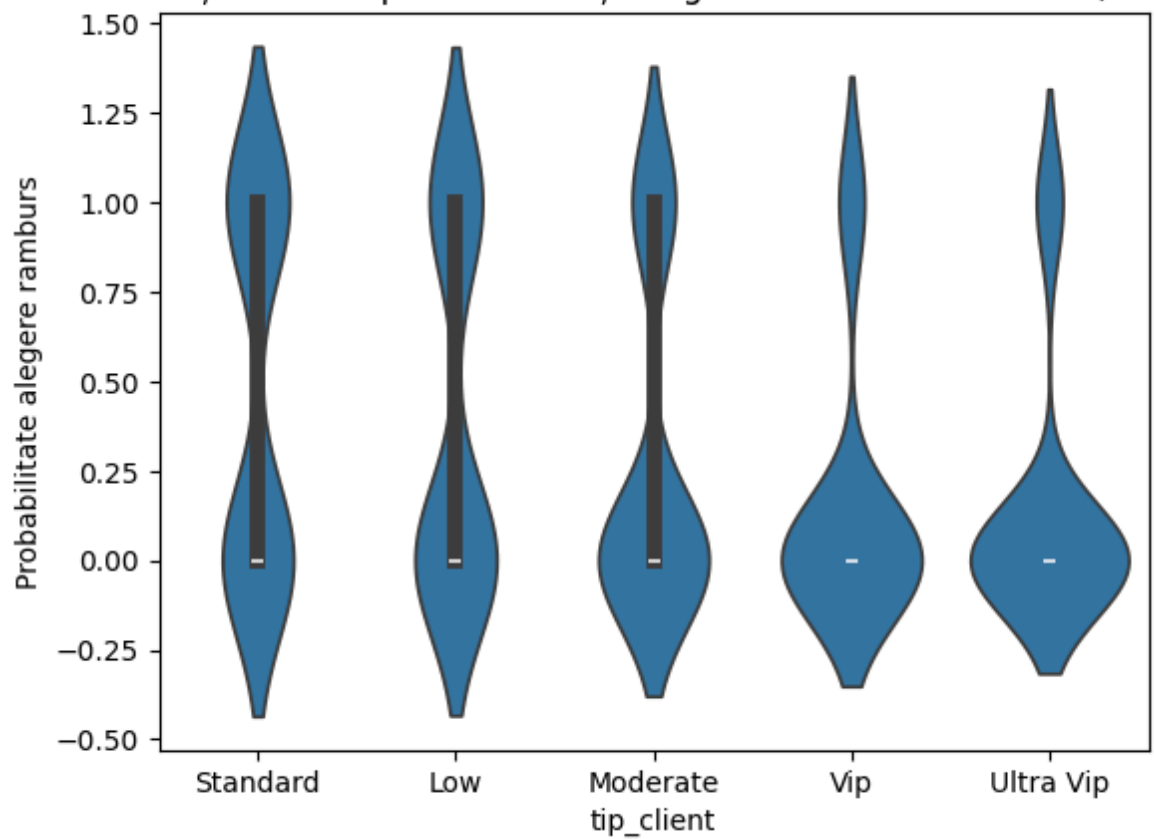
Relatia dintre tipul de client si varsta acestuia



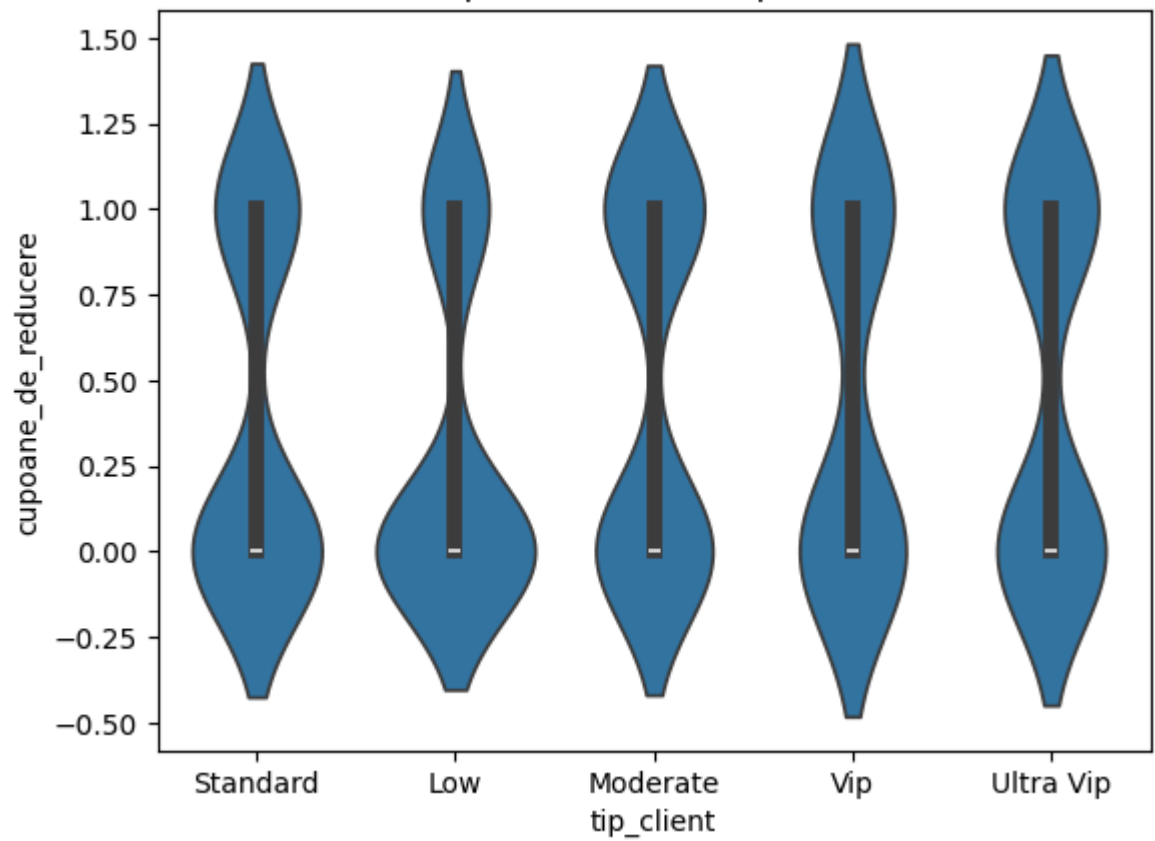
Relatia dintre tipul de client si frecventa lunara a acestuia



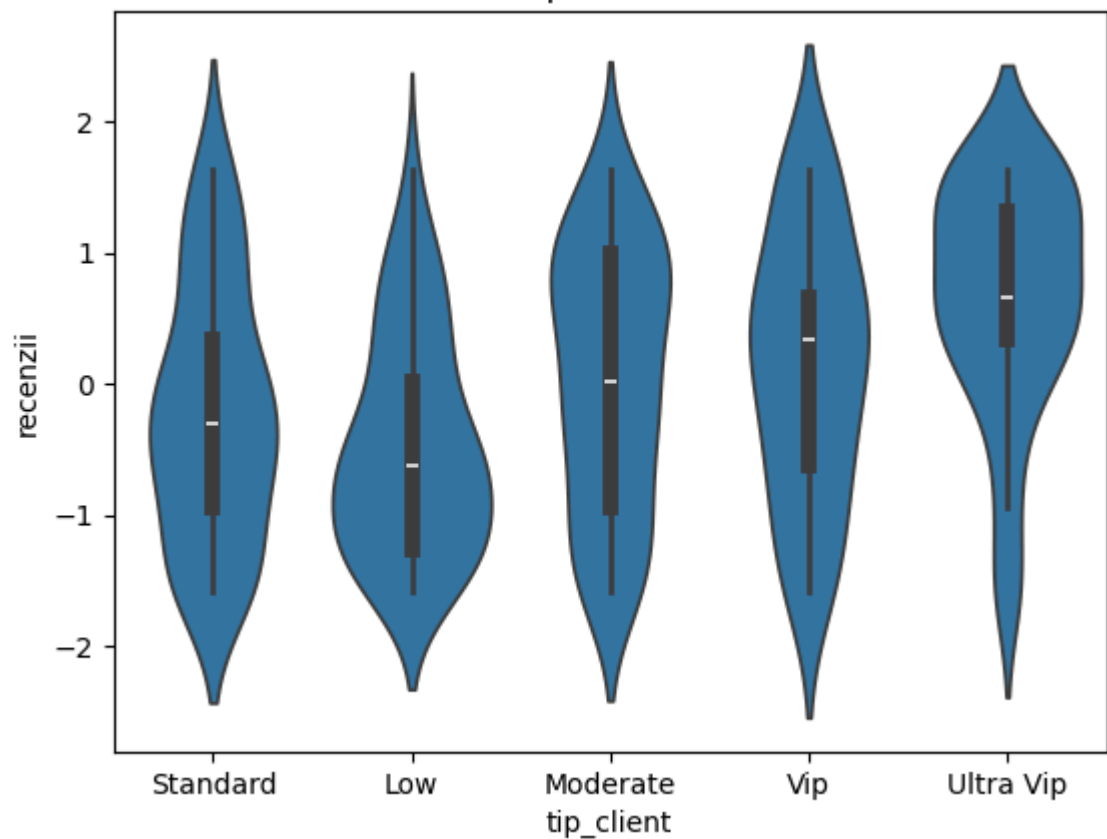
Relația dintre tipul de client și alegerea metodei Ramburs (test)



Relatia dintre tipul de client si cupoanele de reducere



Relatia dintre tipul de client si recenzii



### 3. Antrenarea si compararea a 3 algoritmi diferiti si evaluarea performantei

3. a) Impartirea datelor in antrenament si test: Setul de date train\_encodat.csv a fost impartit in X\_train (variabilele independente fara coloana tip\_client) si y\_train (variabila dependenta tip\_client). Setul de date test\_encodat.csv a fost pregatit in acelasi mod pentru testare (X\_test, y\_test).

3. b) Alegerea si antrenarea modelelor: Am ales trei modele clasice: Regresie Logistica, Random Forest, KNN. Toate cele trei modele au fost antrenate folosind datele din X\_train si y\_train.

3. c) Precizarea rezultatelor: Fiecare model a fost folosit pentru a face predictii asupra setului de test X\_test. Primele 10 predictii au fost comparate vizual cu valorile reale din y\_test.

3. d) Evaluarea performantei: Au fost calculate acuratetea (cat de multe clasificari au fost corecte) si scorul F1 (media armonica intre precizie si recall)

3. e) Matricea de confuzie: Pentru fiecare model am generat o matrice de confuzie care arata cate instante din fiecare clasa au fost corect sau fressit clasificate si distributia erorilor.

3. f) Tabel comparativ de scoruri: Un tabel final comparativ a fost creat cu: numele fiecarui model, acuratetea si scorul F1 corespunzator.