

Proiect PCLP3

Duță Ștefan-Horia, 311CA

Descriere:

1. Tipul problemei:

Setul de date este destinat unei probleme de regresie, care consta in predictia veniturilor unei companii.

2. Structura setului de date:

Setul de date va avea 2000 de exemple si este impartit in doua subseturi:

- Subsetul de antrenare: are 1400 de exemple
- Subsetul de testare: are 600 de exemple

3. Caracteristici:

Fiecare instanta are 9 coloane relevante, dintre care “venit” este cea tinta.

Pentru acestea am folosit, 3 tipuri diferite de date: numere intregi(numar angajati, an infiintare, numar produse, numar clienti, numar parteneri), numere reale(capital social, venit) si valori categoriale(domeniu de activitate, tara)

4. Salvarea dataseturilor:

Subsetul de antrenare va fi salvat in “train.csv”, iar cel de testare in “test.csv”.

5. Explicarea modului de constructie a setului de date:

Problema noastra consta in predictia veniturilor unei companii. Setul de date a fost generat aleator, folosind distributia normala si uniforma, cu urmatoarele restrictii: media numarului de angajati este de 100, capitalul social este intre 50000 si 500000 de euro, anul de infiintare este intre 1970 si 2024, numarul de produse este intre 1 si 100, numarul de clienti este intre 50 si 1000, iar numarul de parteneri intre 1 si 20. Tara si domeniul de activitate au fost alese aleator dintr-un set fix de optiuni, acestea fiind Romania, Germania, Franta, Italia, Spania, respectiv, IT, Retail, Agricultura, Educatie, Financiar. Venitul a fost calculat folosind o formula la care contribuie toate caracteristicile unei instante si anume:

$$(\text{numar_angajati} * 3000 + \text{capital_social} * 0.5 + \text{vechime} * 1000 + \text{numar_clienti} * 500 + \text{numar_parteneri} * 1000) * \text{coef_domeniu} * \text{coef_tara}.$$

Apoi, alegem un numar intre 50 si 100 de valori din fiecare coloana pentru a le sterge.

6. Analiza exploratorie a datelor (EDA complex):

Analizare valori lipsa:

numar angajati	54
capital social	83
an infiintare	71
numar produse	70
domeniu de activitate	94
tara	86
clienti	58
parteneri	67
venit	0

In cazul coloanelor cu valori numerice, vom inlocui spatiile libere cu media aritmetica de pe acea coloana, iar in cazul celor cu valori categoriale, vom inlocui cu cea mai frecventa eticheta.

Statistici descriptive:

	numar angajati	capital social	an infiintare	numar produse	clienti	parteneri	venit
count	2000.000000	2000.000000	2000.000000	2000.000000	2000.000000	2000.000000	2.000000e+03
mean	100.405961	270572.452888	1996.191291	50.118653	533.848095	10.188308	1.020487e+06
std	19.442187	126352.652747	15.374407	28.087174	267.425996	5.337319	3.590860e+05
min	35.000000	50005.235640	1970.000000	1.000000	50.000000	1.000000	2.714512e+05
25%	87.000000	163007.283415	1984.000000	27.000000	305.000000	6.000000	7.477798e+05
50%	100.405961	270572.452888	1996.191291	50.118653	533.848095	10.188308	9.756660e+05
75%	113.000000	373262.840762	2009.000000	73.000000	758.000000	15.000000	1.234982e+06
max	177.000000	499800.966463	2023.000000	99.000000	999.000000	19.000000	2.219132e+06

Analiza distributiei variabilelor:

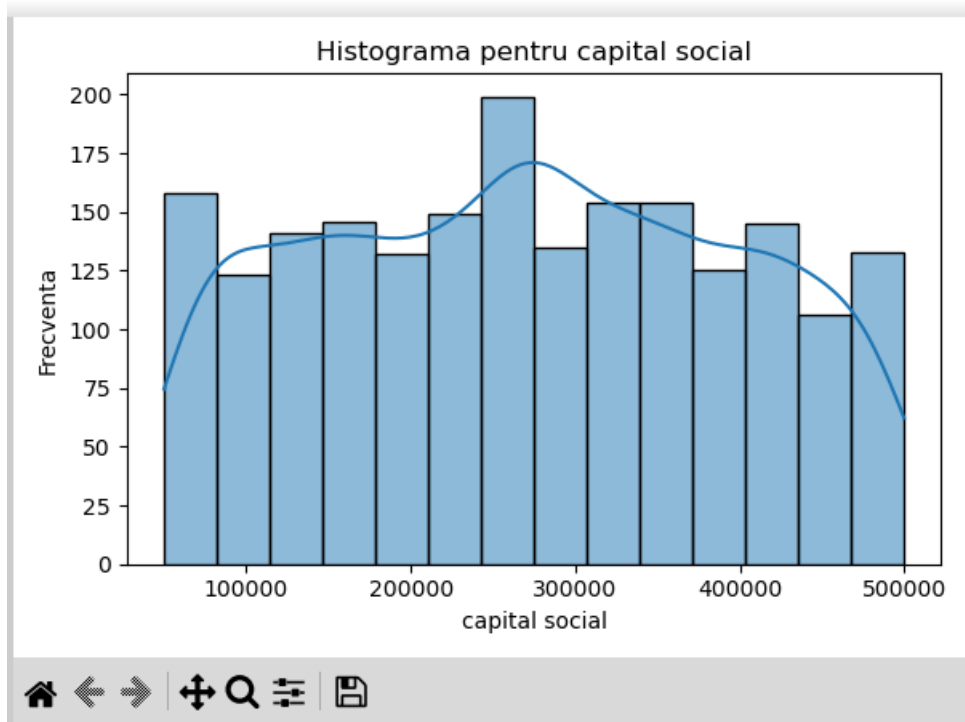
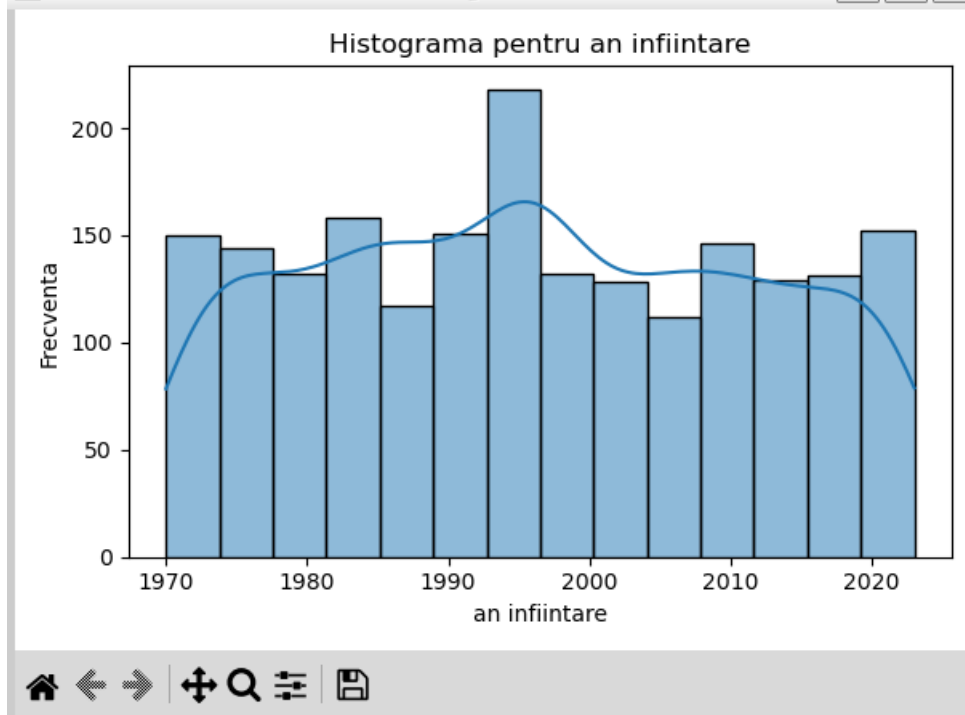
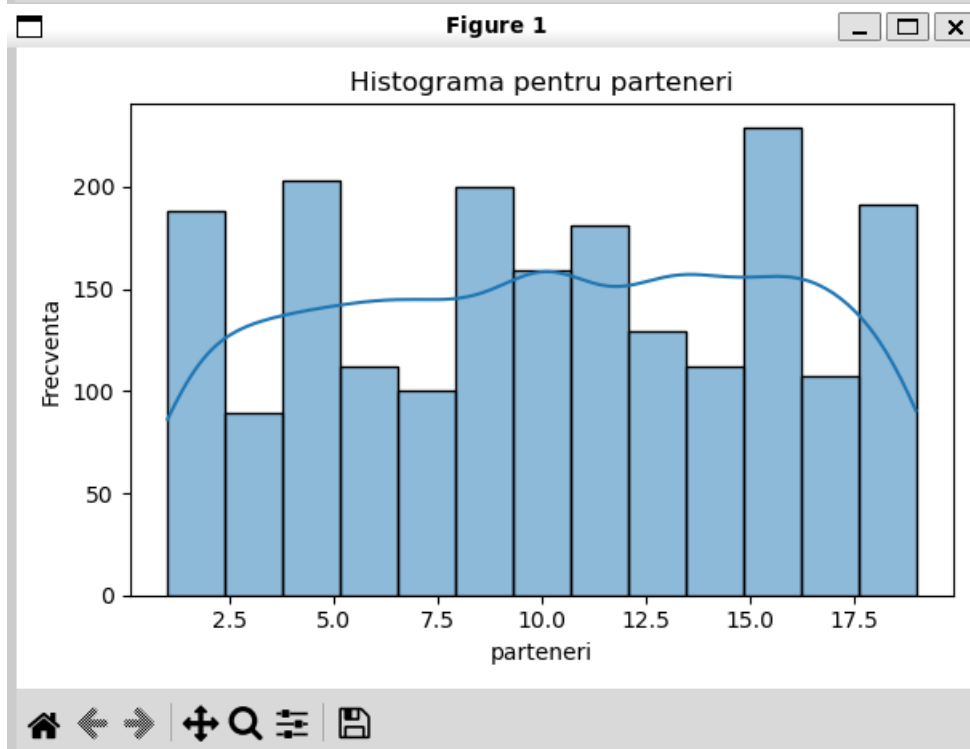
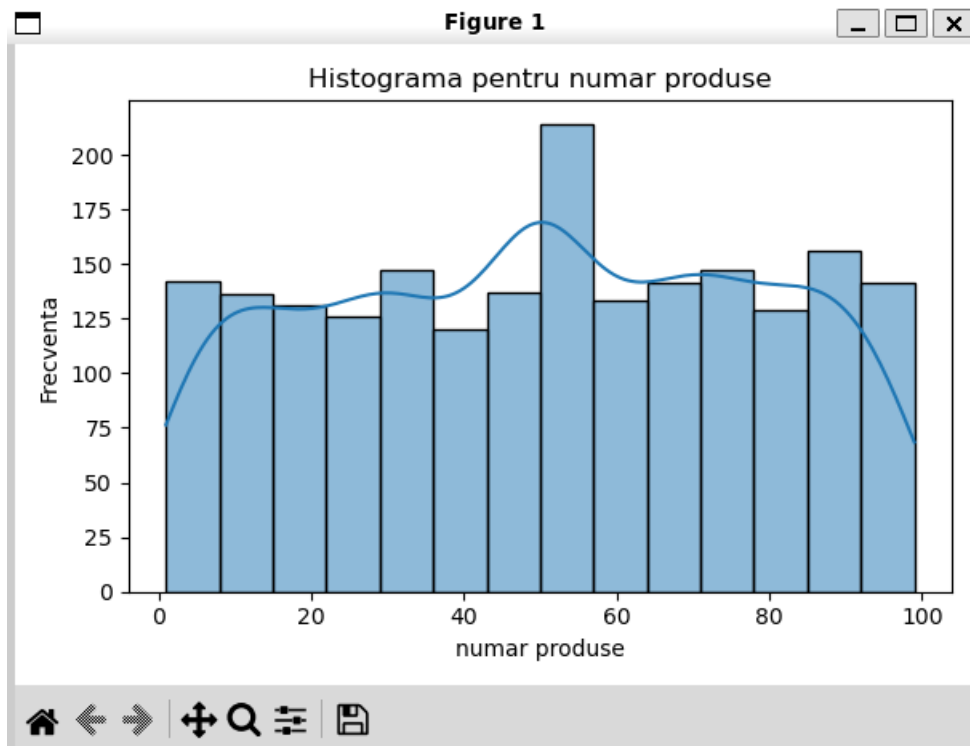
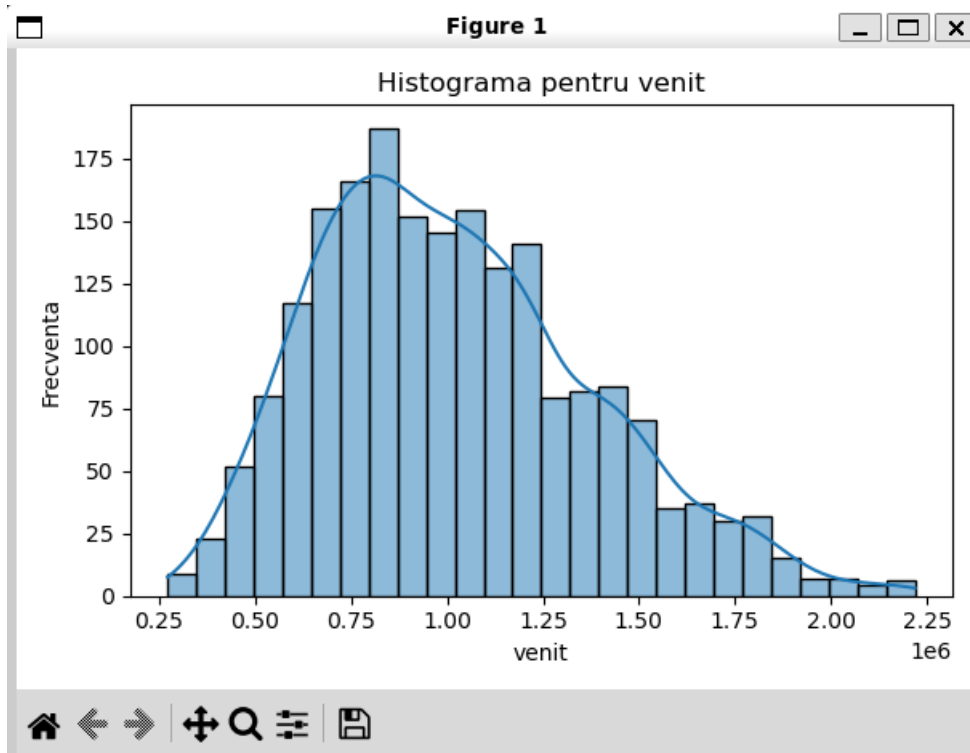
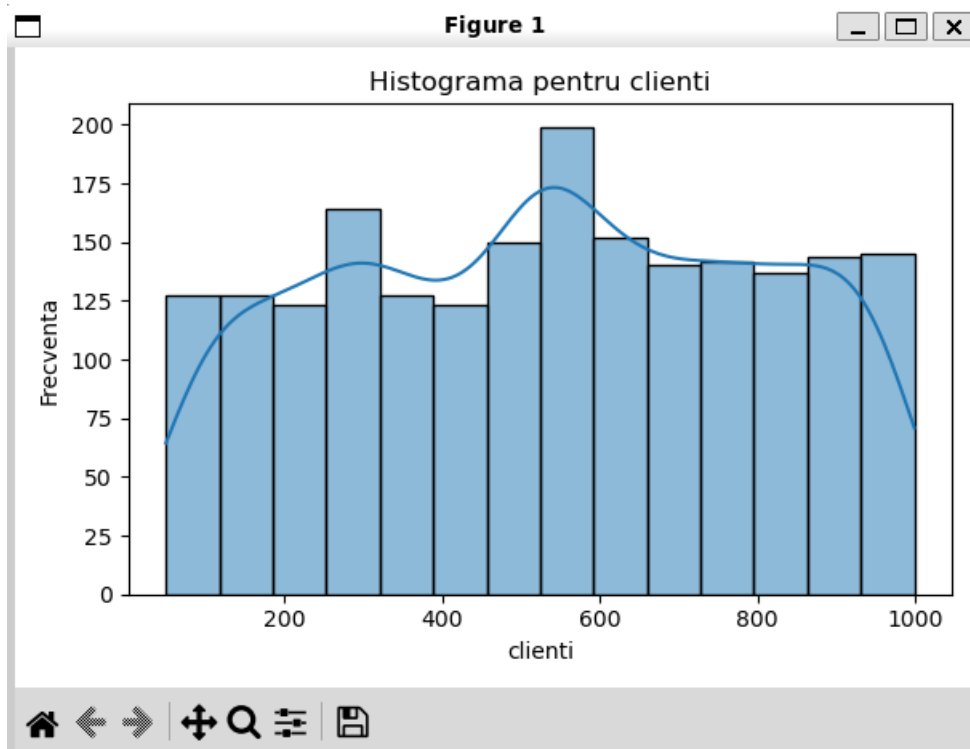
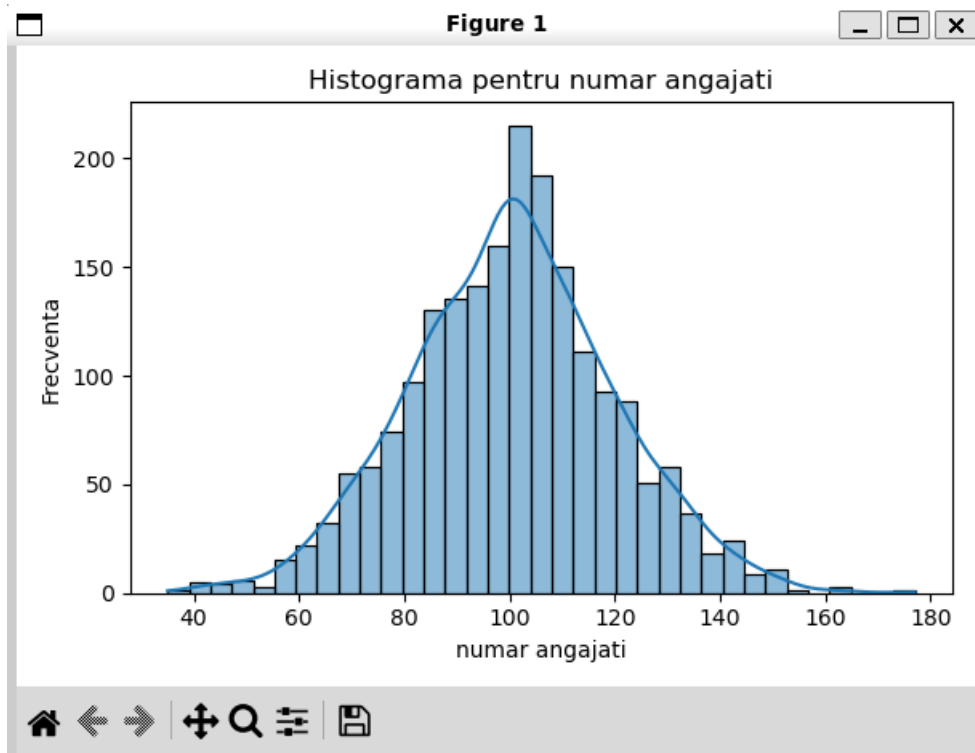


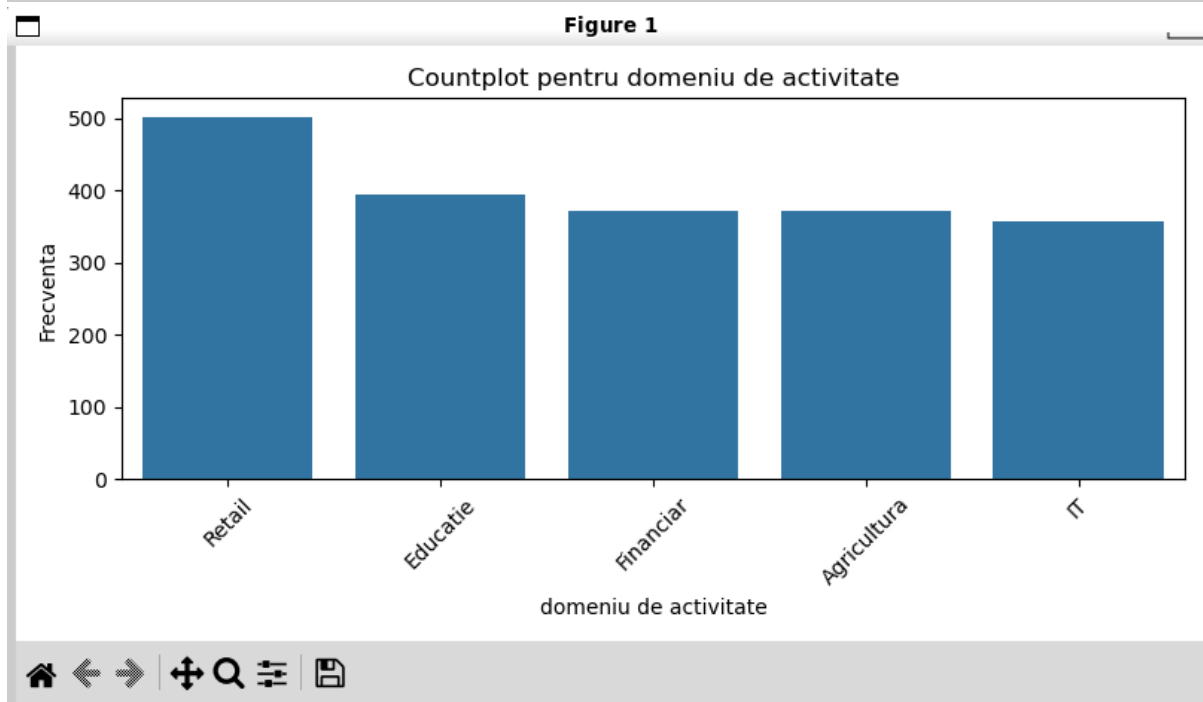
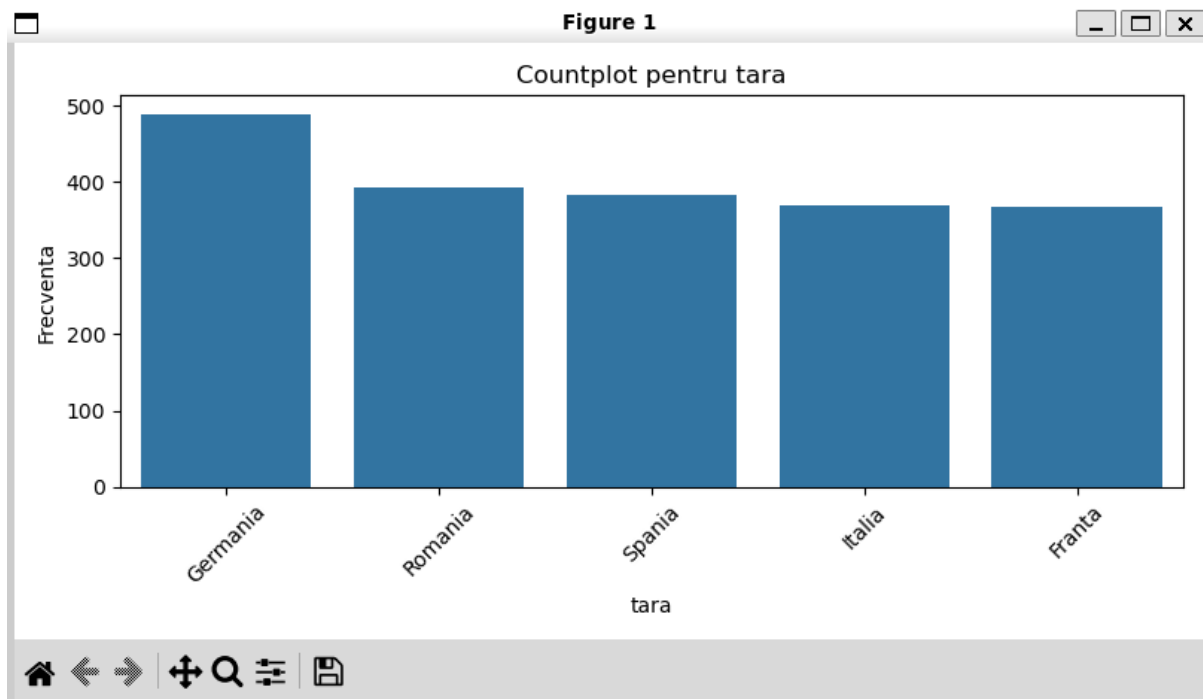
Figure 1





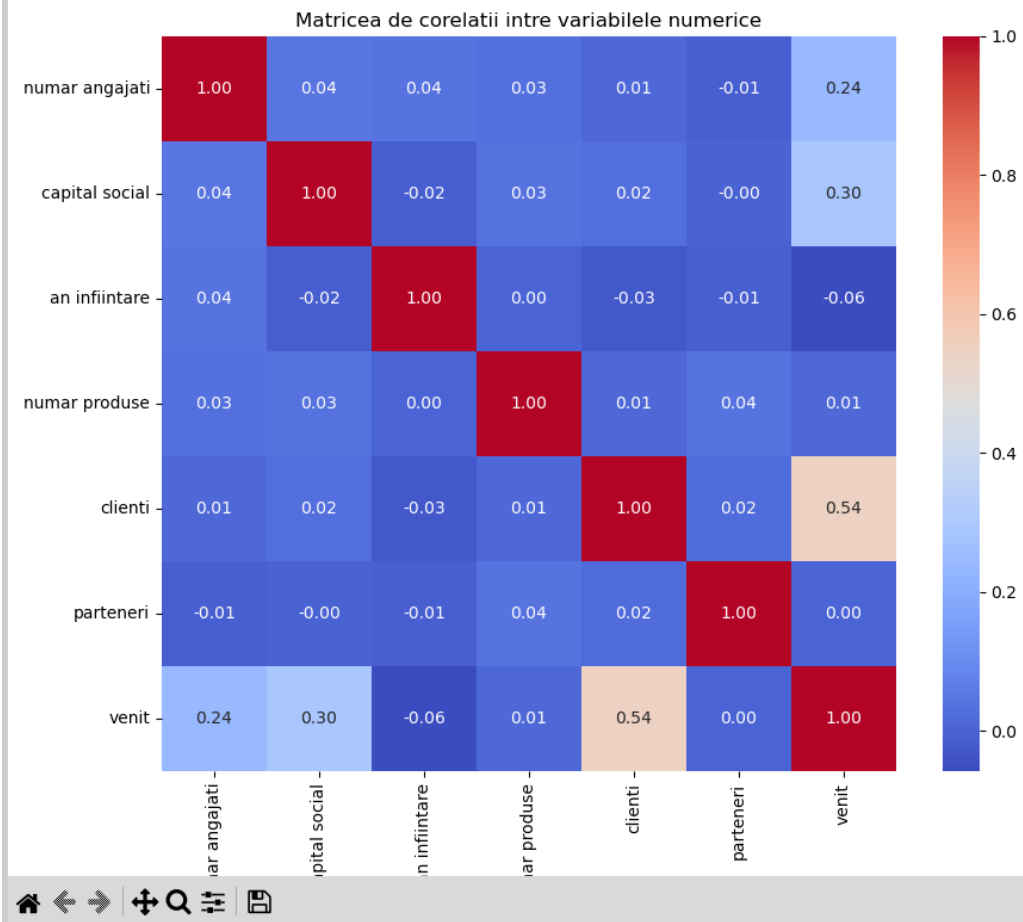




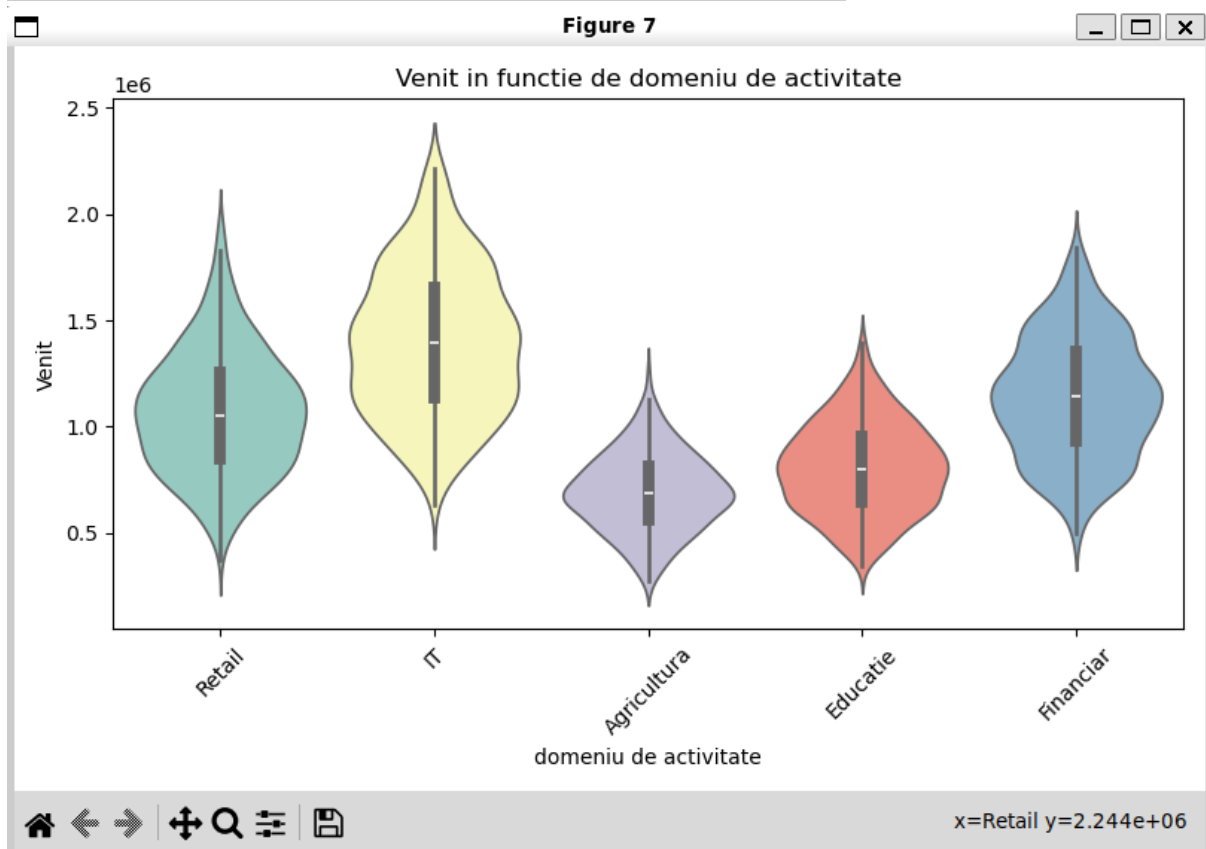
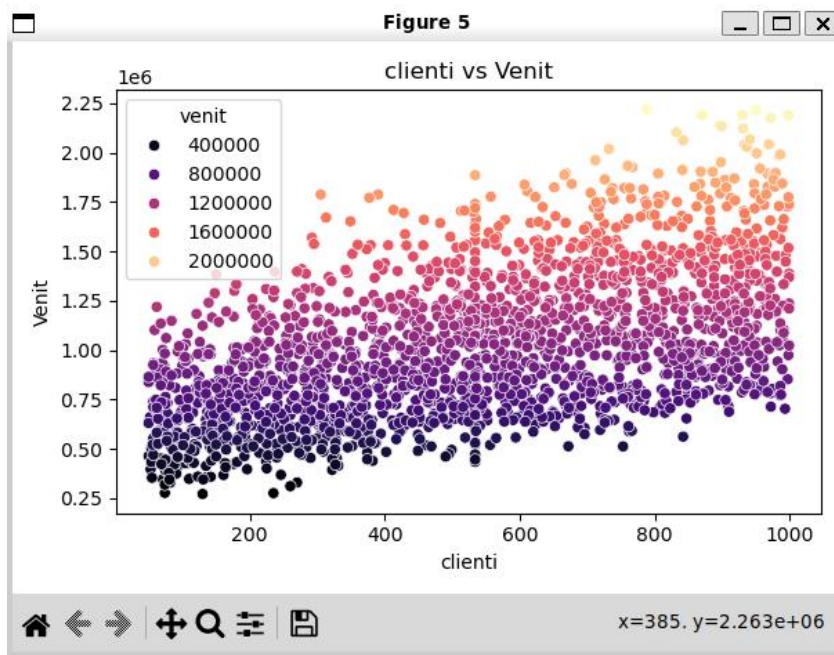


Analiza corelatiilor:

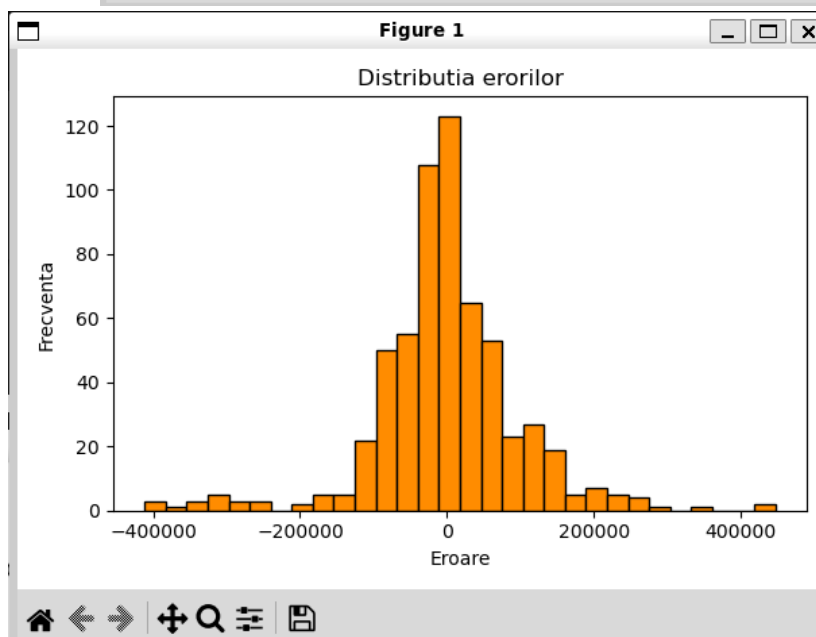
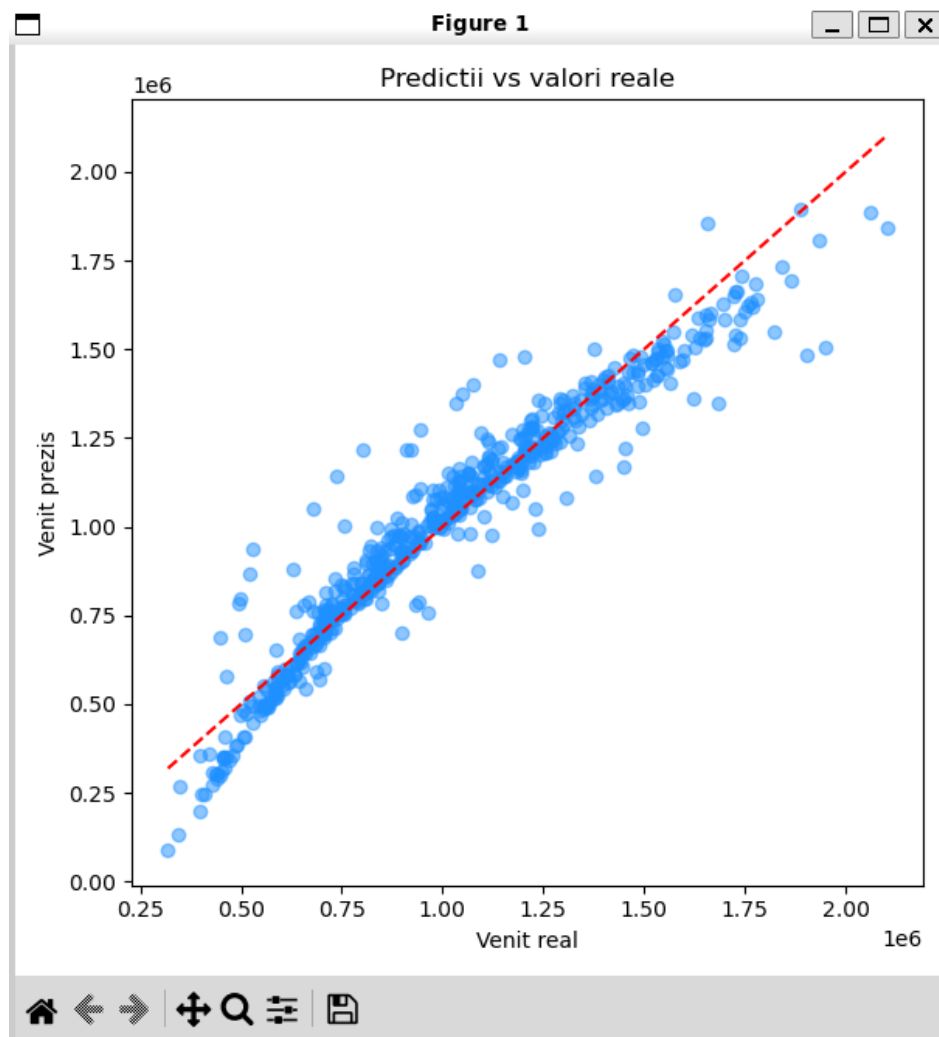
— □ ×



Analiza relatiilor cu variabila tinta:



7. Antrenarea si evaluarea unui model de baza:



MAE: 65922.13
RMSE: 98376.38
R²: 0.9205