

Obiectivele lucrării

Obiectivul principal al acestei lucrări este de a utiliza regresia liniară pentru a prezice sumele facturate pe baza datelor istorice de vânzări. În acest sens, vom urmări următoarele obiective specifice:

- Pregătirea și explorarea setului de date de vânzări.
- Implementarea unui model de regresie liniară pentru a prezice sumele facturate.
- Evaluarea performanței modelului și interpretarea rezultatelor obținute.

Prin atingerea acestor obiective, ne propunem să evidențiem utilitatea și eficacitatea regresiei liniare în anticiparea cererii și în luarea deciziilor într-un mediu comercial.

Metodologie

În această secțiune, vom oferi o descriere detaliată a metodelor și tehnicilor utilizate în lucrare pentru pregătirea datelor, implementarea modelului de regresie liniară și evaluarea performanței acestuia.

- Pregătirea datelor

Pregătirea datelor este un pas esențial în procesul de dezvoltare a modelului. În acest scop, am efectuat următoarele etape:

1. Importul datelor: Am încărcat setul de date de vânzări într-un cadru de date utilizând biblioteca pandas în Python.
2. Explorarea datelor: Am analizat și am explorat setul de date pentru a înțelege structura acestuia, tipurile de date ale fiecărei variabile și pentru a identifica eventualele probleme sau lipsuri.
3. Pregătirea datelor: Am efectuat operațiuni de curățare a datelor, inclusiv tratarea valorilor lipsă, eliminarea sau înlocuirea datelor anormale și transformarea datelor non-numerice într-o formă numerică utilizabilă pentru model.

- Modelul de regresie liniară

Modelul de regresie liniară este utilizat pentru a estima relația liniară între variabila dependentă (suma facturată) și una sau mai multe variabile independente (baza impozitare, valoare TVA). În implementarea noastră, am urmat următoarele etape:

1. Definirea variabilelor dependente și independente: Am selectat variabilele relevante din setul de date pentru a fi utilizate în model.
2. Împărțirea datelor: Am împărțit datele într-un set de date de antrenament și un set de date de testare pentru a evalua performanța modelului pe date nevăzute.
3. Antrenarea modelului: Am utilizat setul de date de antrenament pentru a antrena modelul de regresie liniară.

4. Testarea modelului: Am testat performanța modelului folosind setul de date de testare și am evaluat metrici cum ar fi Mean Squared Error (MSE) și R^2 Score.

- Evaluarea modelului

Pentru a evalua performanța modelului, am folosit următoarele metrici:

- Mean Squared Error (MSE): Este o măsură a dispersiei dintre valorile reale și cele prezise ale variabilei dependente. Cu cât este mai mic MSE, cu atât modelul este considerat mai precis.

- R^2 Score: Reprezintă coeficientul de determinare și oferă o măsură a calității ajustării modelului la datele de antrenament. Valoarea R^2 apropiată de 1 indică o bună potrivire a modelului la date.

- Predicția

Pentru a realiza predicții cu privire la sumele facturate în viitor, am utilizat modelul de regresie liniară antrenat pe datele existente. Procesul de predicție a urmat următoarele etape:

****Selectarea caracteristicilor****: Am selectat caracteristicile relevante care să fie utilizate pentru a face predicțiile. În cazul nostru, aceste caracteristici au inclus baza impozitare și valoarea TVA.

****Pregătirea datelor de predicție****: Am pregătit datele de predicție, asigurându-ne că acestea au aceeași structură și formă ca și datele folosite pentru antrenamentul modelului.

****Realizarea predicțiilor****: Am aplicat modelul antrenat pe datele de predicție pentru a obține estimări cu privire la sumele facturate.

****Evaluarea predicțiilor****: Am comparat predicțiile obținute cu datele reale existente pentru a evalua acuratețea și performanța modelului în a prezice sumele facturate.

Rezultate

Setul de date are o structură tabulară și este organizat într-un format tip tabel cu rânduri și coloane. Fiecare rând reprezintă o înregistrare individuală sau o observație, în timp ce fiecare coloană reprezintă o variabilă sau un atribut specific.

Structura setului de date include următoarele caracteristici:

- Dimensiuni: 6461 de înregistrări și 6 coloane.

- Coloane: Fiecare coloană reprezintă o variabilă specifică sau un atribut al datelor. În cazul acestui set de date, coloanele includ:

1. Nr. factura: Numărul facturii
2. Data factura: Data la care a fost emisă factura.
3. Client: Numele sau identificatorul clientului asociat facturii.

4. Suma factura: Valoarea totală a facturii.

5. Baza impozitare: Valoarea de bază a facturii, fără TVA.

6. Valoare TVA: Valoarea taxei pe valoarea adăugată (TVA) asociată facturii.

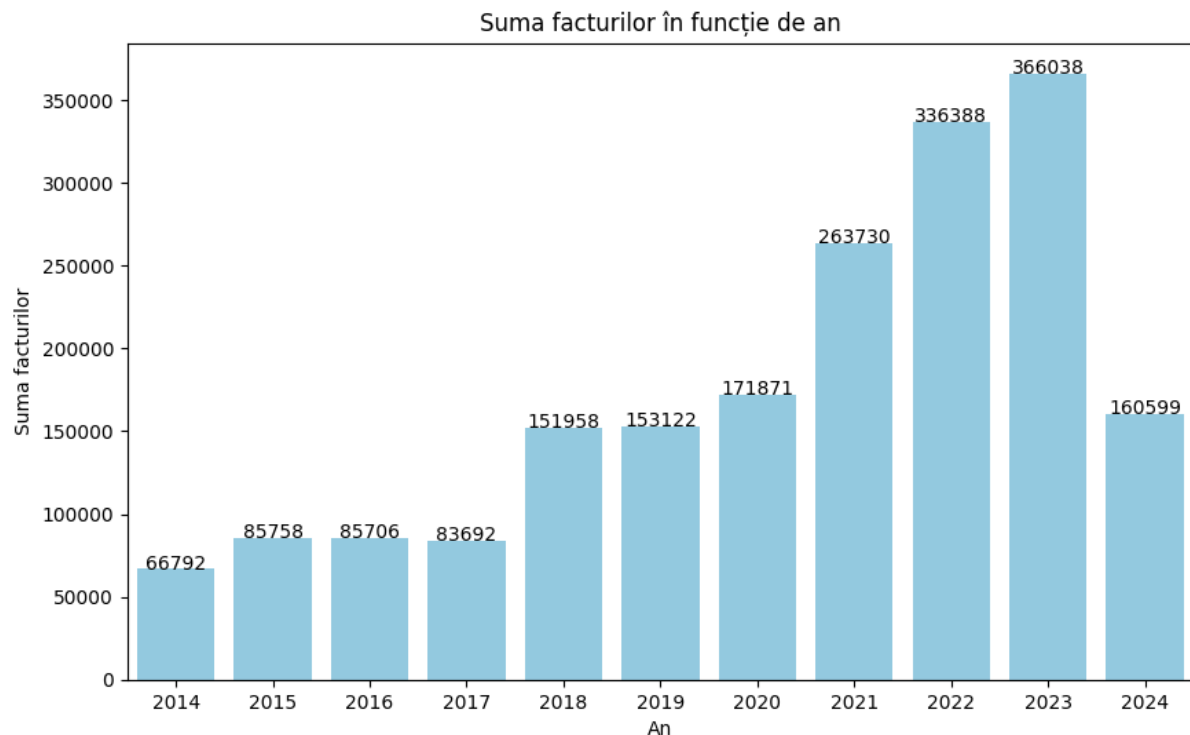
- Tipuri de date: Fiecare coloană poate avea un tip de date specific, cum ar fi întregi, șiruri de caractere sau valori reale. De exemplu, Nr. factura poate fi de tip întreg, Data factura și Client pot fi șiruri de caractere, iar Suma factura, Baza impozitare și Valoare TVA pot fi valori reale.

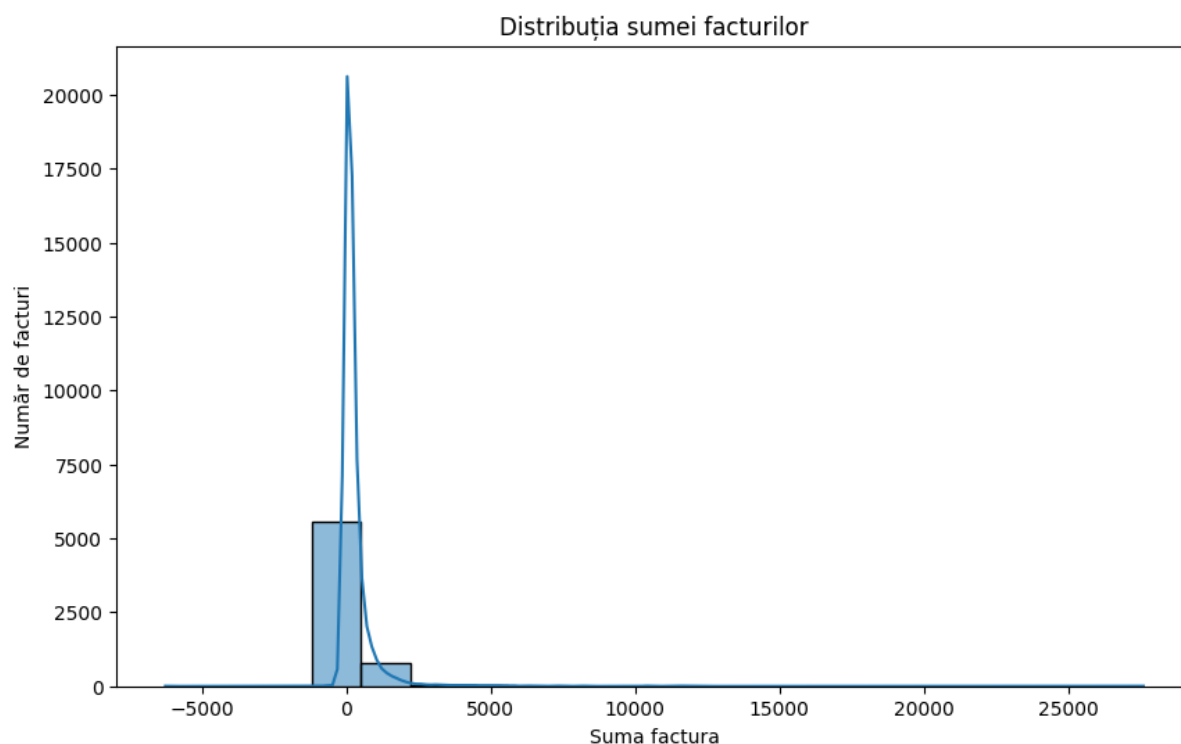
Vizualizare primele 5 rânduri din set:

	Nr. factura	Data factura	Client \
0	53962	02.05.2014 00:00:00	MEGA AUTO S.R.L.
1	53963	02.05.2014 00:00:00	GIP EST SRL
2	53964	05.05.2014 00:00:00	COLEGIUL NATIONAL DE ARTA OCTAV BANCILA
3	53971	05.05.2014 00:00:00	EPTISA ROMANIA SRL
4	53972	08.05.2014 00:00:00	PALATUL COPIILOR - IASI

	Suma factura	Baza impozitare	Valoare TVA
0	30.24	24.39	5.85
1	407.51	328.64	78.87
2	56.00	45.16	10.84
3	270.44	218.10	52.34
4	180.00	145.16	34.84

Vizualizare grafice:



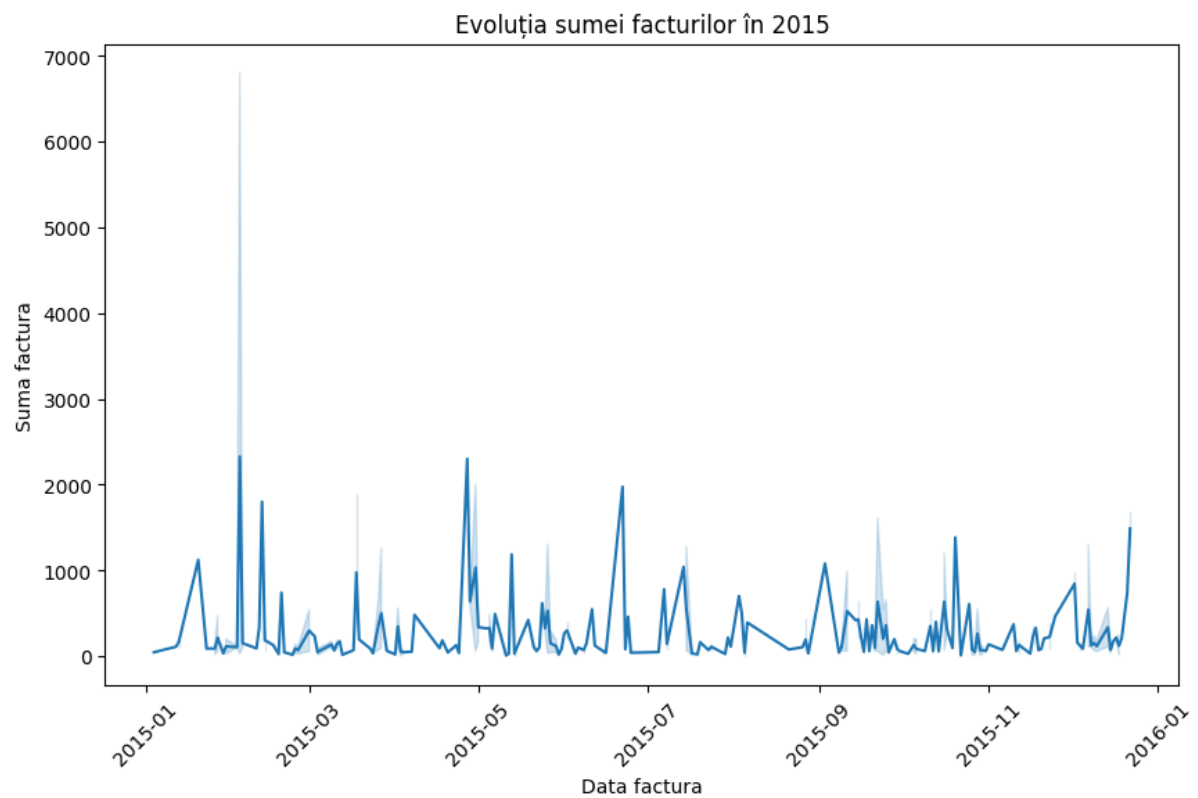
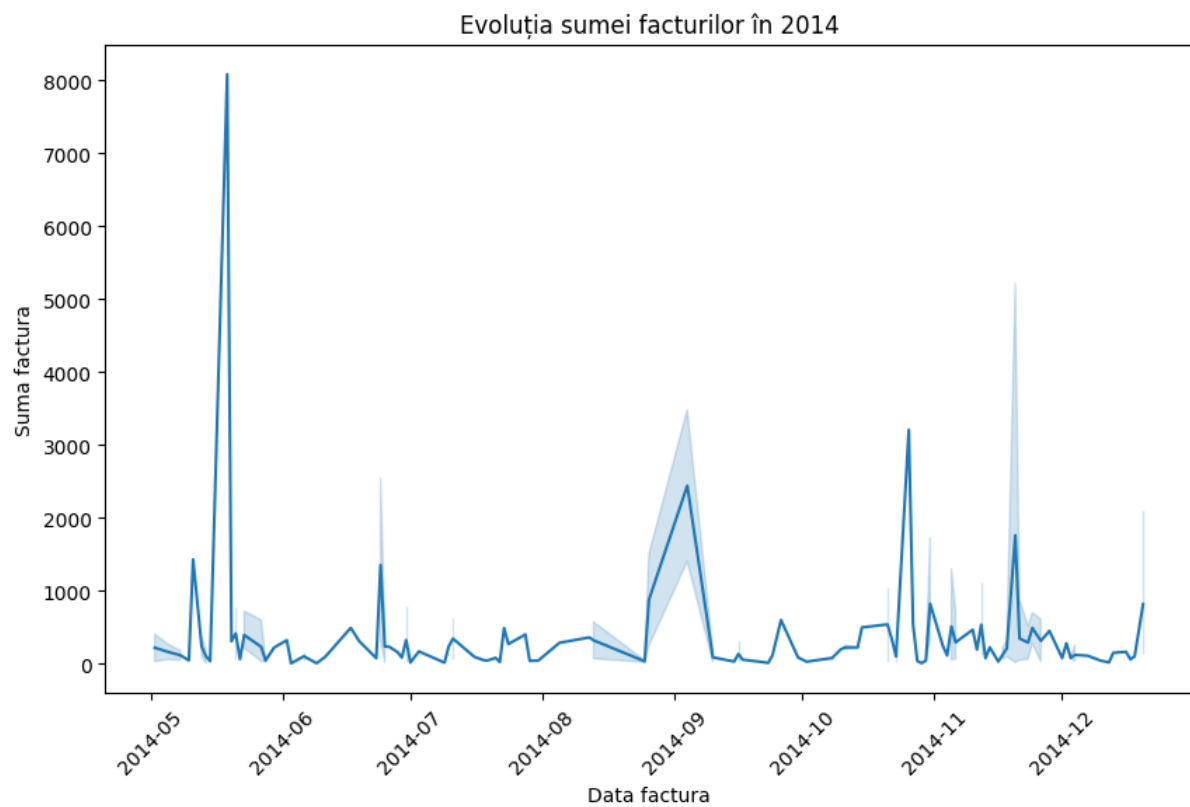


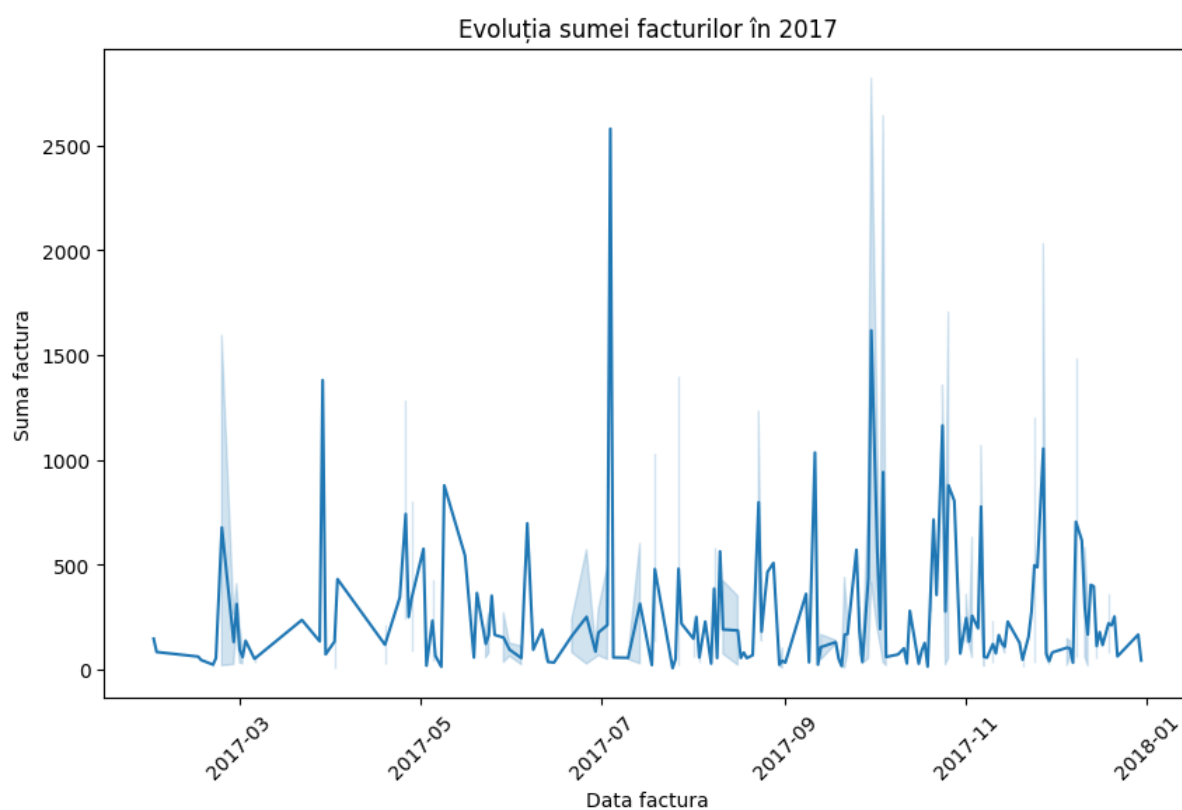
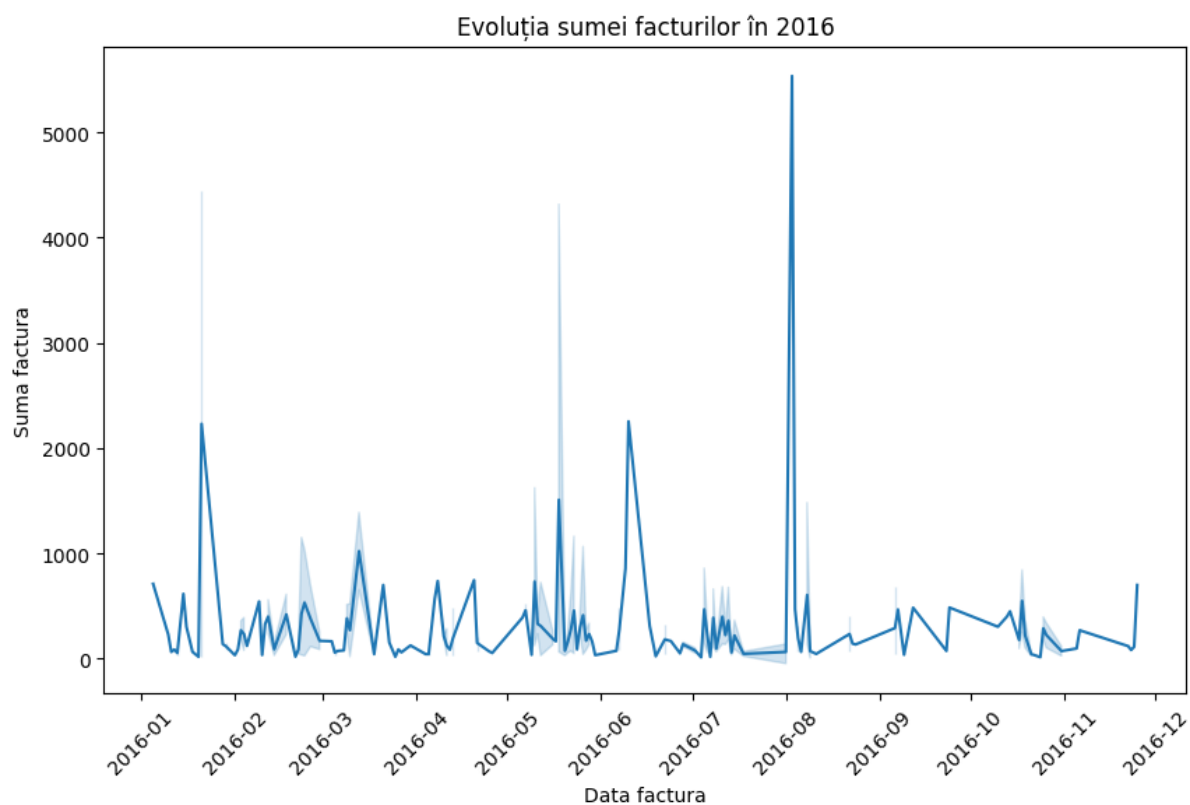
Graficul prezentat arată distribuția sumei facturilor, având pe axa orizontală sumele facturilor și pe axa verticală numărul de facturi.

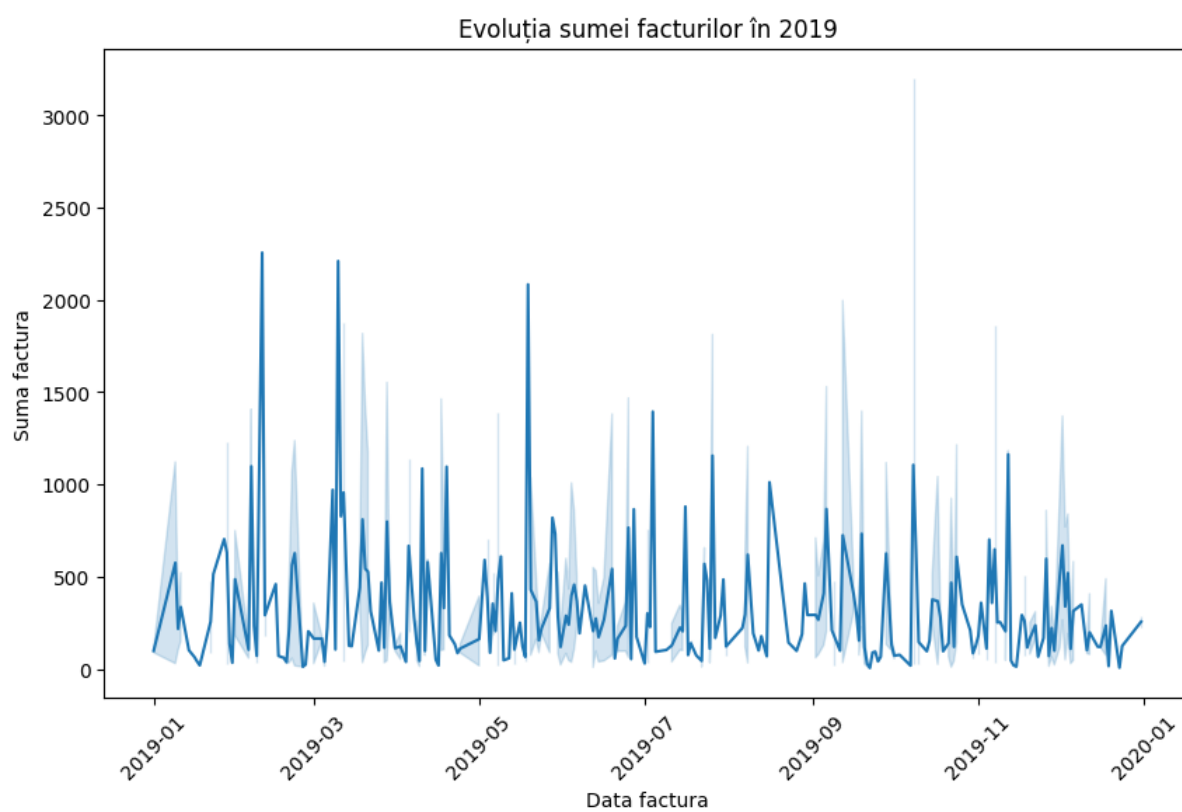
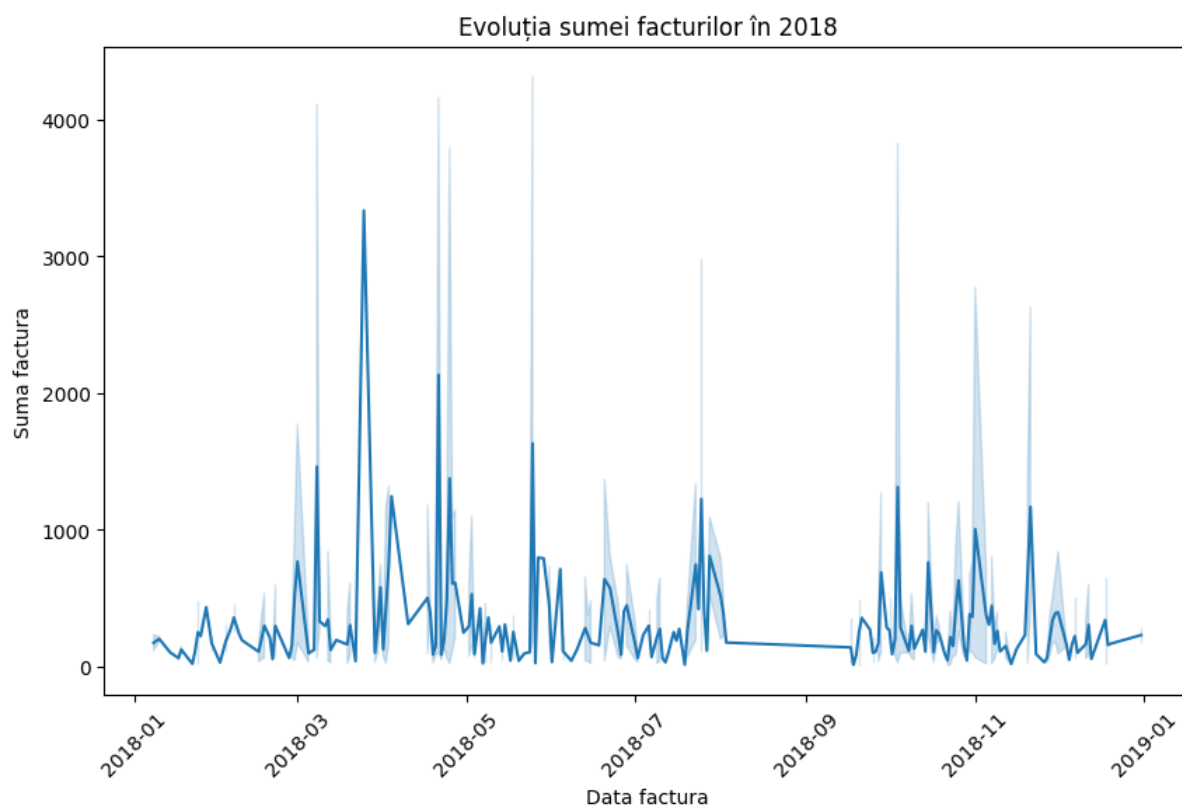
Observăm că majoritatea facturilor se concentrează în jurul valorii zero, unde vedem un vârf foarte pronunțat, indicând o frecvență mare a facturilor cu sume foarte mici sau zero. Există de asemenea o mică frecvență pentru facturi cu valori negative, posibil reprezentând rambursări sau corecții.

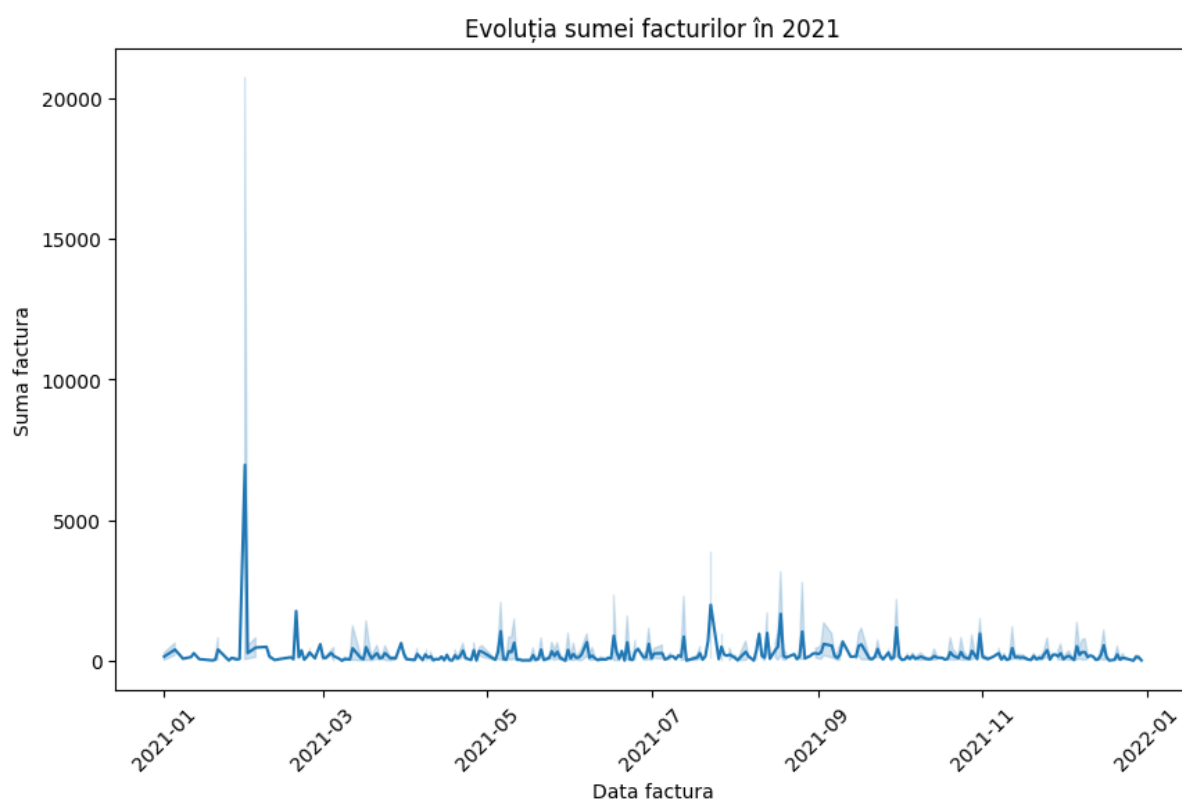
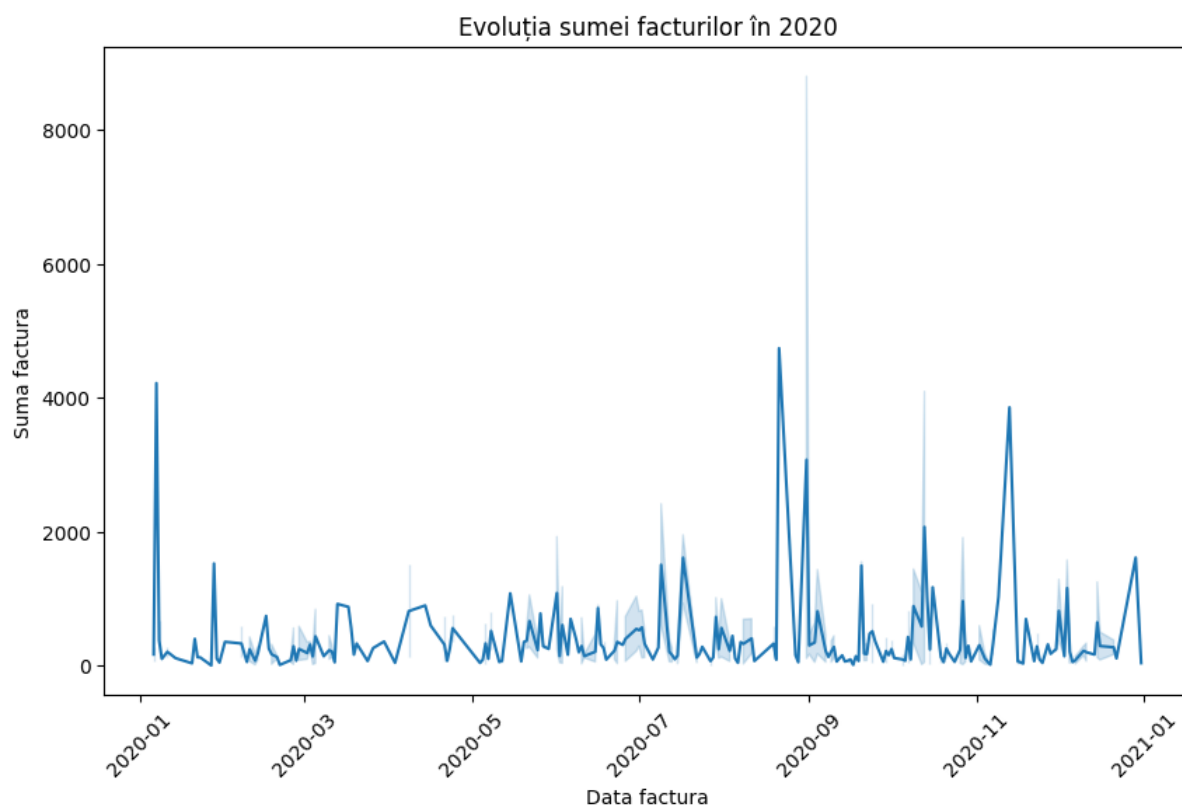
După acest vârf pronunțat, numărul de facturi scade dramatic odată cu creșterea sumei facturii, sugerând că facturile cu sume mari sunt mult mai rare.

Evoluția sumei facturilor pe ani:

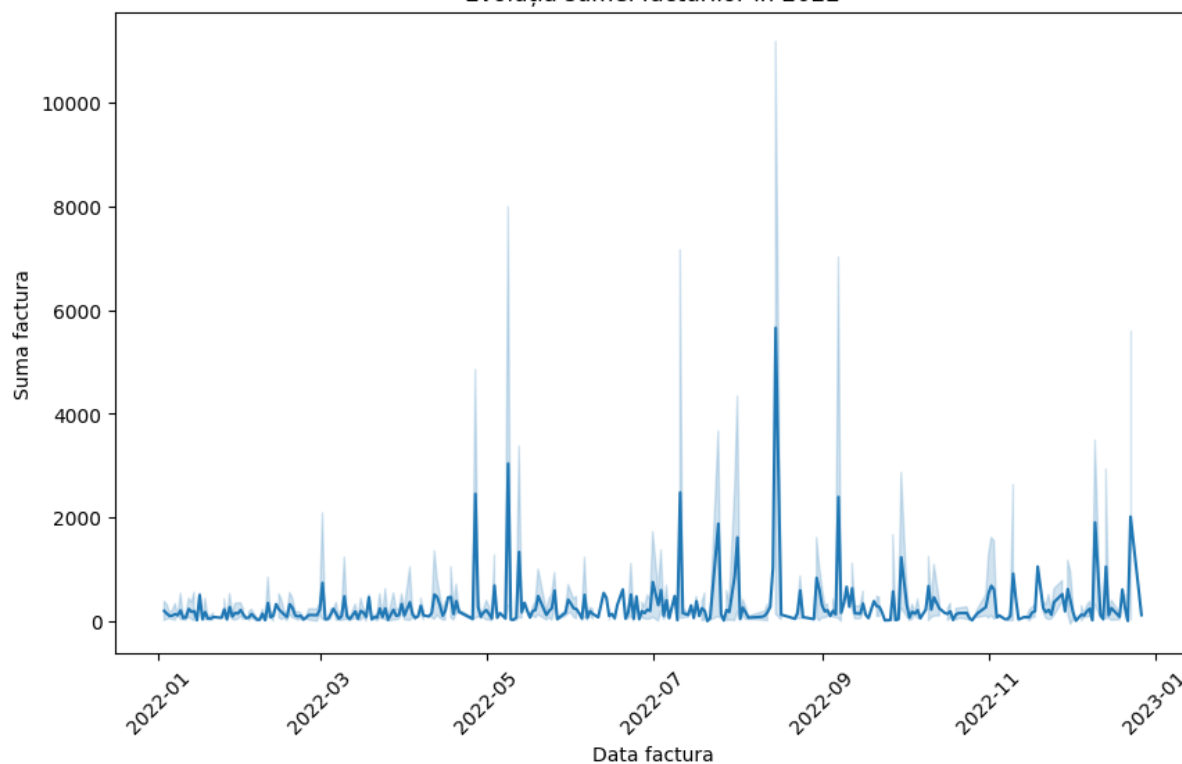




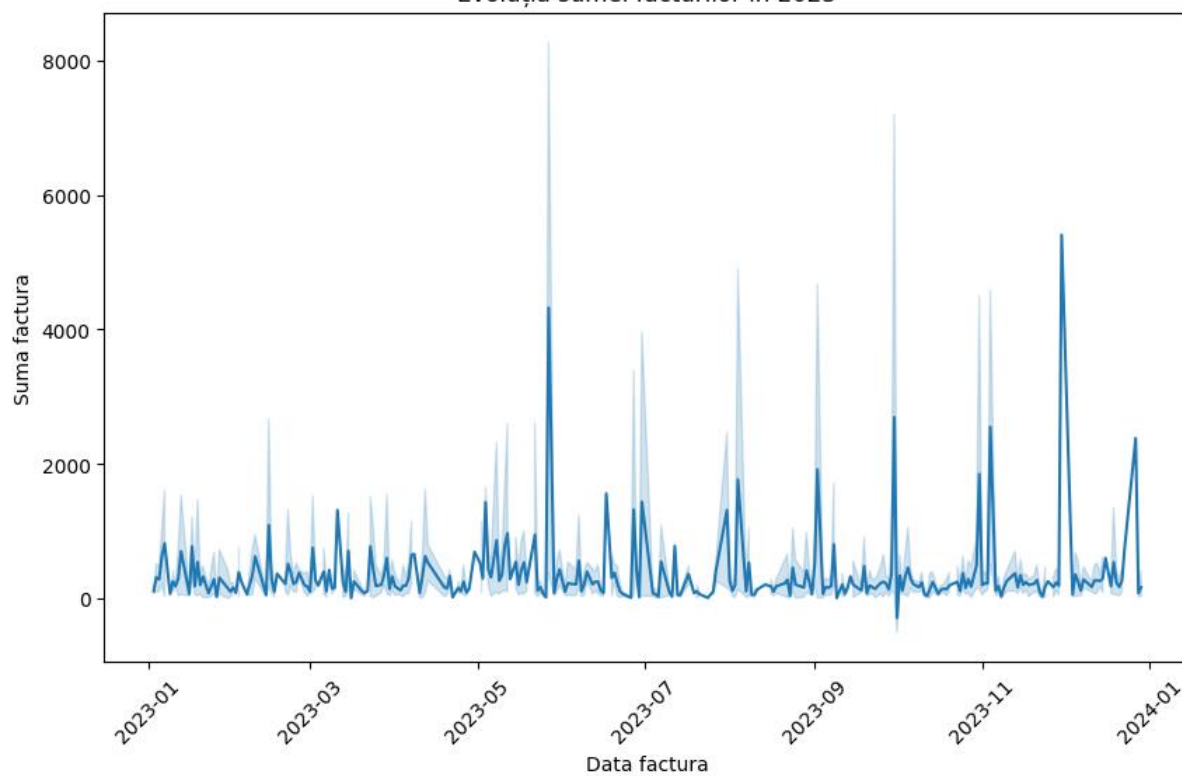


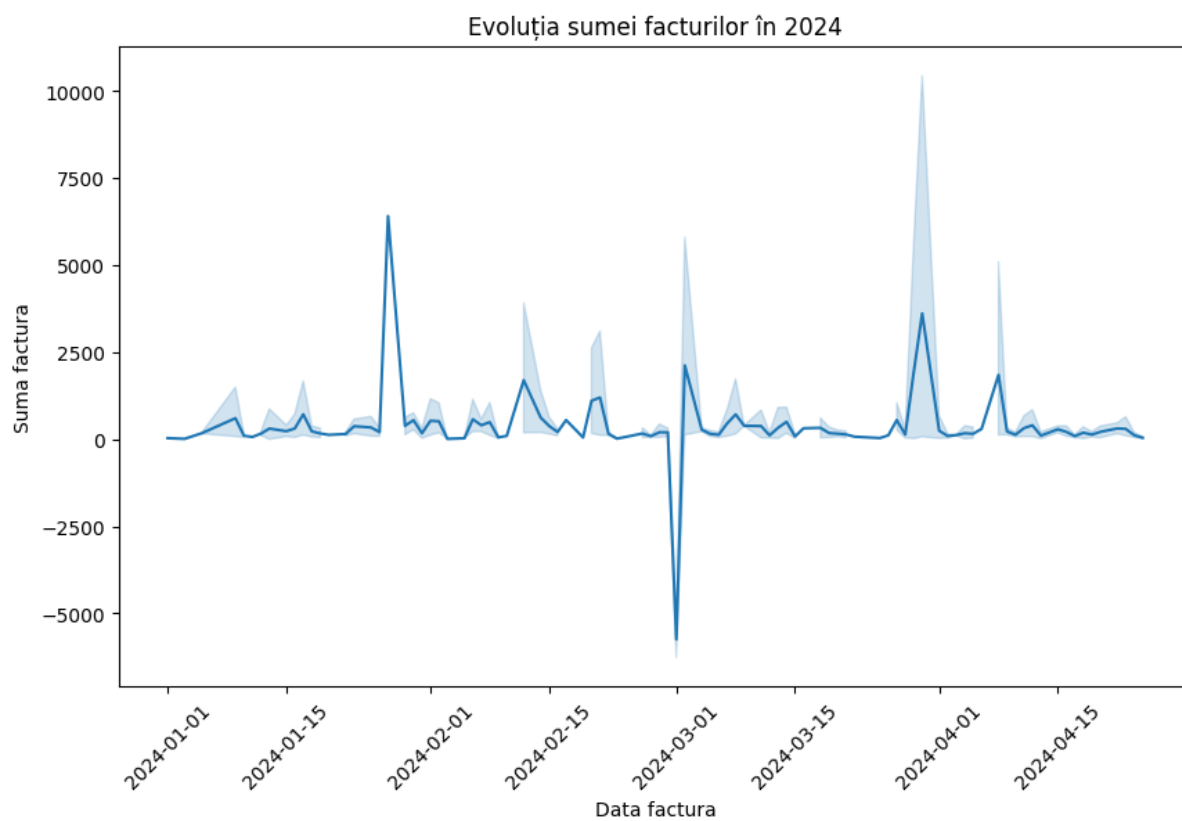


Evoluția sumei facturilor în 2022

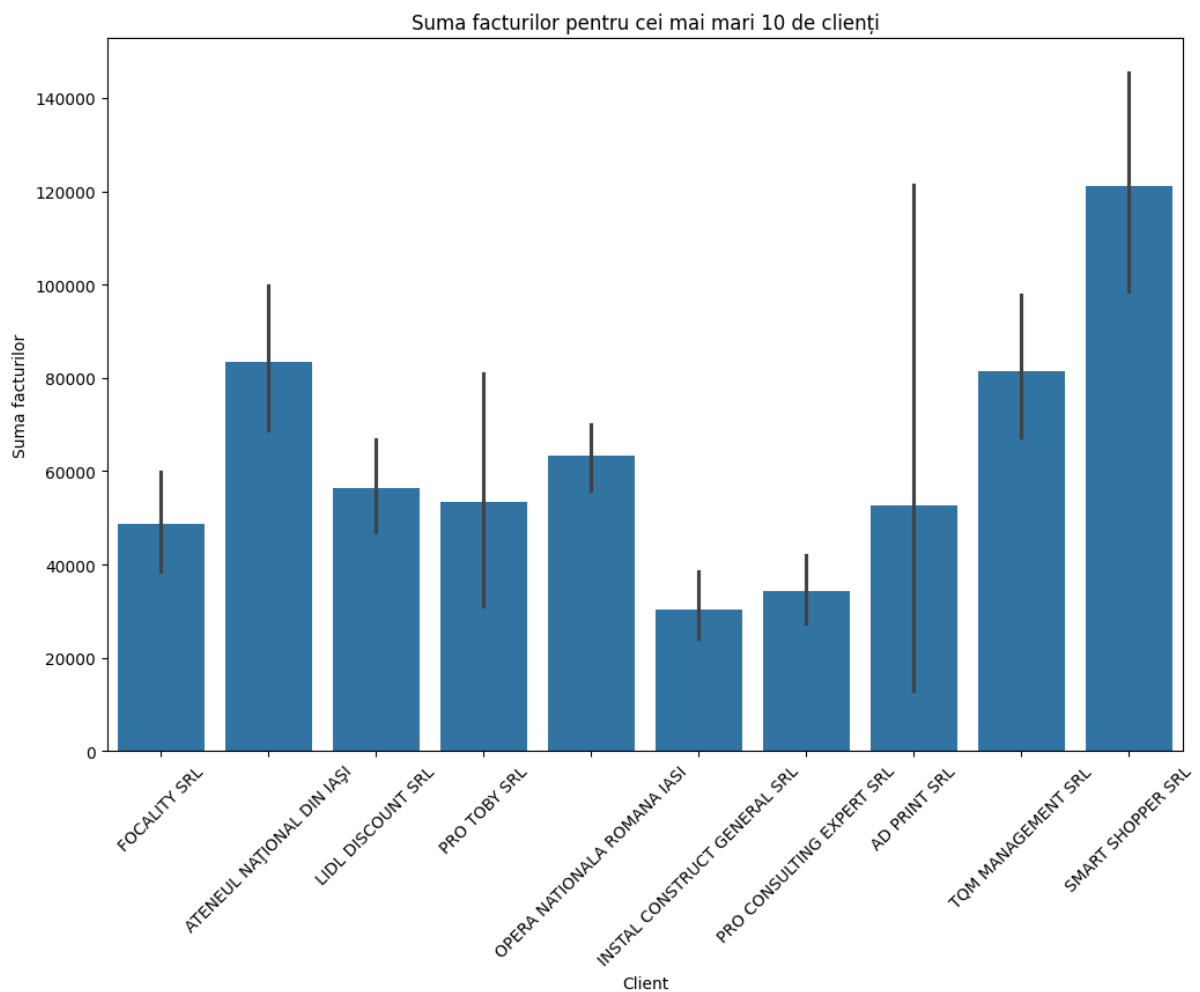


Evoluția sumei facturilor în 2023





Suma facturilor pentru top 10 clienti.



Partea de Machine Learning – Antrenare Model:

```
# Exemplu de pregătire a datelor pentru regresie liniară
X = df[['Baza impozitare', 'Valoare TVA']]
y = df['Suma factura']
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)

model = LinearRegression()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)

# Evaluate model
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
print("Mean Squared Error:", mse)
print("R^2 Score:", r2)
```

Rezultatul obținut:

Mean Squared Error: 3.652068269258797e-05
R² Score: 0.9999999999353562

Interpretare:

1. Mean Squared Error (MSE): Acesta este o măsură a diferenței medii pătratice între valorile prezise de model și valorile reale. Cu cât MSE este mai mic, cu atât modelul prezice mai bine datele. În cazul meu, valoarea MSE este foarte mică (3.652068269258797e-05), ceea ce indică faptul că modelul prezice foarte bine valorile.

2. R-squared Score (R²): Acesta este o măsură a cât de bine variabilitatea în datele de intrare este explicată de model. R-squared variază între 0 și 1, unde 1 înseamnă că modelul explică perfect variația datelor și 0 înseamnă că modelul nu explică deloc variația datelor. În cazul meu, valoarea R² este aproape de 1 (0.9999999999353562), ceea ce indică faptul că modelul explică foarte bine variația datelor.

Predicția:

```
# Folosim modelul pentru a face predicții pentru cererea viitoare
cerere_viitoare_pred = model.predict(df[['Baza impozitare', 'Valoare
TVA']])

# Afișează predicțiile pentru cererea viitoare
print(cerere_viitoare_pred)
```

Rezultatul obținut:

```
# Selecția primelor 50 de valori prezise și cele reale
primele_50_valori = df[['Predictie', 'Suma factura']].head(50)

# Afișarea într-un tabel
print(primele_50_valori.to_string(index=False))
```

Predictie	Suma facturi reală
30.246328	30.24
407.514988	407.51
56.006236	56.00
270.445474	270.44
180.005795	180.00
51.106253	51.10
45.236274	45.23
1429.501357	1429.50
371.345116	371.34
94.006101	94.00
12.306391	12.30
195.305741	195.30
32.856318	32.85
8082.757724	8082.78
306.305347	306.30

57.006233	57.00
771.003696	771.00
62.006215	62.00
214.005675	214.00
245.205564	245.20
720.003878	720.00
55.806237	55.80
34.026314	34.02
598.254310	598.25
26.606340	26.60
58.686226	58.68
216.005667	216.00
225.005636	225.00
321.265294	321.26
4.806418	4.80
23.046353	23.04
15.126381	15.12
193.505748	193.50
151.155898	151.15
124.505992	124.50
5.766415	5.76
92.906105	92.90
491.004691	491.00
308.165340	308.16
74.906169	74.90
2542.947402	2542.95
311.345329	311.34
1200.522170	1200.52
537.004527	537.00
19.326366	19.32
147.005913	147.00
232.835608	232.83
157.485875	157.48
84.606135	84.60
797.993600	797.99

Există diferențe între valorile prezise și cele reale:

