

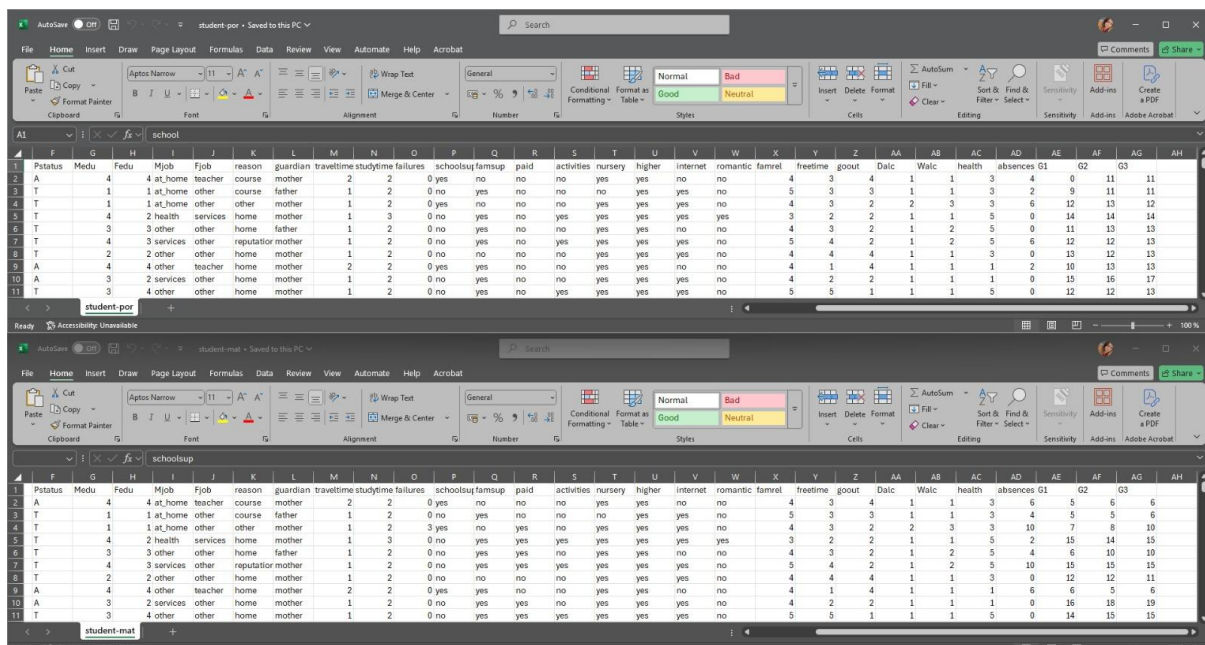
Raport Proiect „Student Performance”

Acest proiect are ca scop analiza setului de date Student Performance, folosind tehnici de preprocesare a datelor, clasificare, clustering si vizualizare. Prin aceste etape, se urmareste intelegerea factorilor care influenteaza performanta scolara si evaluarea comportamentului elevilor din perspective diferite.

1. Preprocesarea Datelor

Pentru aceasta tema am ales setul de date Student_Performance, bazat pe performanta elevilor, asupra caruia am aplicat mai multe operatii pentru a-l aduce intr-o forma potrivita pentru analiza.

Dupa descarcarea locala, am importat fisierul student-mat.csv, care contine informatii despre elevi, precum scoala, varsta, absentele si notele. Dupa incarcarea datelor, am verificat structura tabelului si tipurile de valori pentru a ma asigura ca informatiile sunt citite corect.



Am verificat daca exista valori lipsa, insa setul de date este complet. In cazul in care ar fi existat, valorile numerice ar fi putut fi completate cu media, iar valorile text cu moda, insa aceasta decizie ar fi fost luata doar daca distributia datelor ar fi permis acest lucru, deoarece aceste metode nu sunt potrivite in toate situatiile.

```
...
Verificarea valorilor lipsă:
school      0
sex         0
age         0
address     0
famsize     0
Pstatus     0
Medu        0
Fedu        0
Mjob        0
Fjob        0
reason      0
guardian    0
traveltime  0
studytime   0
```

Setul contine coloane cu valori text, cum ar fi “school”, “address” sau raspunsurile “yes” si “no”. Pentru a le folosi in analiza, le-am transformat in valori numerice folosind **LabelEncoder**, astfel incat fiecare categorie sa aiba un cod propriu.

```
28 # in acest set avem variabile de tip text ("school", "address", "yes"/"no"), trebuie
29 # aplicam LabelEncoder pe toate coloanele de tip "object"
30 label_encoder = LabelEncoder()
31 cat_cols = df.select_dtypes(include="object").columns.tolist()
32
33 for col in cat_cols:
34     df[col] = label_encoder.fit_transform(df[col])
35
36 print("După codare: ")
37 print(df[cat_cols + ['age', 'G1', 'G2', 'G3']].head(), "\n")
38
```

PROBLEMS 84 OUTPUT DEBUG CONSOLE TERMINAL PORTS SPELL CHECKER 84

(.venv) PS C:\Users\mihae\Desktop\Master anul 2\WADDM\Proiect_Student_Performance> python

dtype: int64

După codare:

	school	sex	address	famsize	Pstatus	Mjob	...	internet	romantic	age	G1	G2	G3
0	0	0	1	0	0	0	...	0	0	18	5	6	6
1	0	0	1	0	1	0	...	1	0	17	5	5	6
2	0	0	1	1	1	0	...	1	0	15	7	8	10
3	0	0	1	0	1	1	...	1	1	15	15	14	15
4	0	0	1	0	1	2	...	0	0	16	6	10	10

[5 rows x 21 columns]

Pentru coloanele text cu mai mult de doua valori (de exemplu Mjob, reason, guardian), am folosit **One-Hot Encoding**, pentru a evita introducerea unei ordini false intre valori.

Apoi am aplicat scalarea datelor cu metoda **Min-Max**, pentru a aduce toate valorile numerice pe acelasi interval, intre 0 si 1. Am aplicat transformarea doar pe variabilele de intrare, fara a modifica notele G1, G2 si G3.

```
39 # scalare Min-Max (NUMAI pe features, NU pe G1,G2,G3)
40 target_cols = ['G1', 'G2', 'G3']
41 feat_cols = [c for c in df.columns if c not in target_cols]
42
43 scaler = MinMaxScaler()
44 df_scaled = df.copy()
45 df_scaled[feat_cols] = scaler.fit_transform(df_scaled[feat_cols])
46 df_scaled[feat_cols] = df_scaled[feat_cols].round(3) # lizibil in Excel
47
48 print("După scalare: ")
49 print(df_scaled.head(), "\n")
50
```

PROBLEMS 84 OUTPUT DEBUG CONSOLE TERMINAL PORTS SPELL CHECKER 84

(.venv) PS C:\Users\mihae\Desktop\Master anul 2\WADDM\Proiect_Student_Performance> python

...

[5 rows x 21 columns]

După scalare:

	school	sex	age	address	famsize	Pstatus	...	Walc	health	absences	G1	G2	G3
0	0.0	0.0	0.429	1.0	0.0	0.0	...	0.00	0.5	0.080	5	6	6
1	0.0	0.0	0.286	1.0	0.0	1.0	...	0.00	0.5	0.053	5	5	6
2	0.0	0.0	0.000	1.0	1.0	1.0	...	0.50	0.5	0.133	7	8	10
3	0.0	0.0	0.000	1.0	0.0	1.0	...	0.00	1.0	0.027	15	14	15
4	0.0	0.0	0.143	1.0	0.0	1.0	...	0.25	1.0	0.053	6	10	10

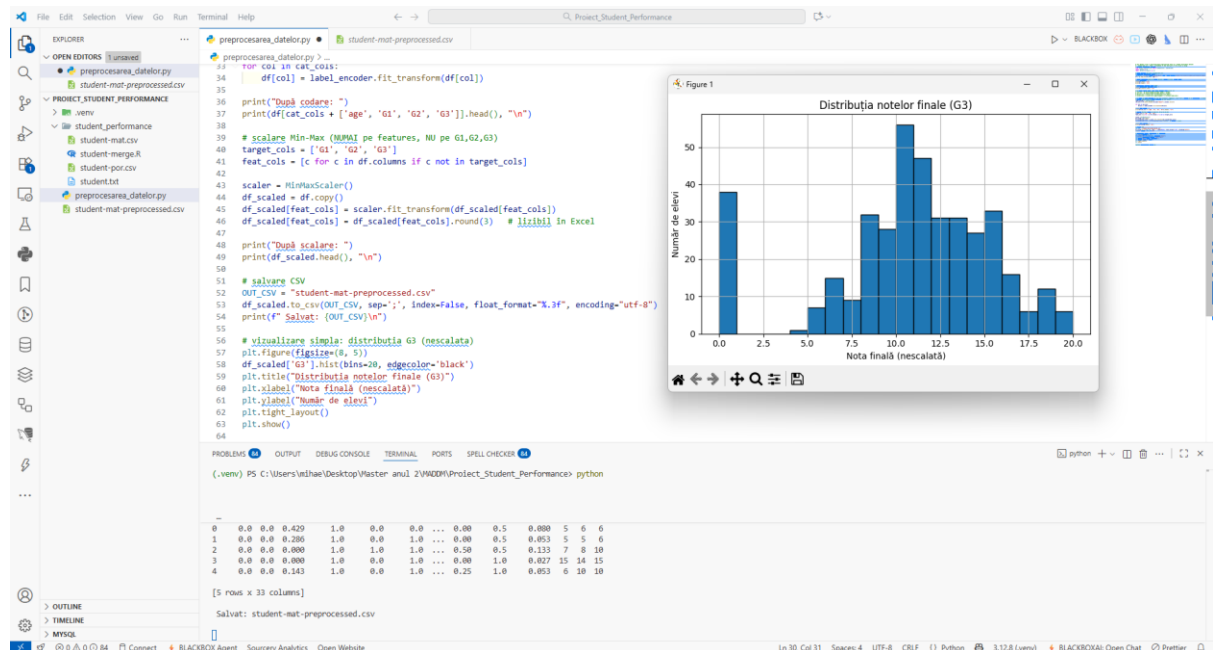
[5 rows x 33 columns]

Salvat: student-mat-preprocessed.csv

Pentru comparatie, am aplicat si scalarea **Z-score**, care aduce datele la media 0 si deviatia standard 1.

		school																																		
A1	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37
	age	address	famsize	status	Medu	Fedu	traveltime	studytime	failures	schoolout	famout	paid	activities	nursery	higher	internet	romantic	formal	freetime	gout	Dalc	Walc	health	absences	G1	G2	G3	Mjob								
0	1.023.046	1	0	0	1.143.856	1.360.371	0.792251	-0.042286	-0.449944	0	1	0	0	0	1	1	0	0.062194	-0.236010	0.801479	-0.540699	-1.003.789	-0.399289	0.036424	5	6	6									
0	0.238380	1	0	0	1.600.009	1.399.970	0.643249	-0.042286	-0.449944	0	1	0	0	0	1	1	0	0.1178.860	-0.236010	-0.097295	-0.540699	-1.003.789	-0.399289	-0.213796	5	5	6									
0	1.330.954	1	1	1	1.600.009	1.399.970	0.643249	-0.042286	3.589.323	1	0	1	0	1	1	1	1	0.062194	-0.236010	-0.097295	0.583385	0.551100	-0.399289	0.536865	7	8	10									
0	1.330.954	1	0	0	1.143.856	-0.479857	0.643249	1.150.779	-0.449944	0	1	1	1	1	1	1	1	1	-1.054.472	-1.238.419	-0.097295	-0.540699	-1.003.789	1.041.070	-0.464016	15	14	15								
0	-0.546287	1	0	0	0.229234	0.440257	0.643249	-0.042286	-0.449944	0	1	1	0	1	1	1	0	0.062194	-0.236010	-0.097295	-0.540699	-0.226345	1.041.070	-0.213796	6	10	10									
1	-0.546287	1	1	1	1.143.856	0.440257	0.643249	-0.042286	-0.449944	0	1	1	1	1	1	1	1	0.1178.860	0.766399	-0.097295	-0.540699	-0.226345	1.041.070	0.536865	15	15	15									
1	-0.546287	1	1	1	-0.685387	-0.479857	0.643249	-0.042286	-0.449944	0	0	0	0	1	1	1	1	0.062194	0.766399	0.801479	-0.540699	-1.003.789	-0.399289	-0.714236	12	12	11									
0	0.238380	1	0	0	1.143.856	1.360.371	0.792251	-0.042286	-0.449944	1	1	0	0	1	1	1	0	0.062194	-2.240.828	0.801479	-0.540699	-1.003.789	-1.839.649	0.036424	6	5	6									
0	-0.546287	1	0	0	0.229234	-0.479857	0.643249	-0.042286	-0.449944	0	1	1	1	0	1	1	1	0.062194	-1.238.419	-0.097295	-0.540699	-1.003.789	-1.839.649	-0.714236	16	18	19									
1	1.330.954	1	0	0	0.229234	1.360.371	0.643249	-0.042286	-0.449944	0	1	1	1	1	1	1	1	0.1178.860	1.768.808	-1.896.683	-0.540699	-1.003.789	1.041.070	-0.714236	14	15	15									
0	1.330.954	1	0	1	1.143.856	1.360.371	0.643249	-0.042286	-0.449944	0	1	1	0	1	1	1	1	0.1054.472	0.236010	-0.097295	-0.540699	-0.226345	-1.119.469	-0.714236	10	8	9									
0	1.330.954	1	0	0	-0.685387	-1.399.970	2.227.751	1.150.779	-0.449944	0	1	0	1	1	1	1	1	0.1178.860	-1.238.419	-0.097295	-0.540699	-1.003.789	0.320890	-0.213796	10	12	12									
1	1.330.954	1	1	1	1.143.856	1.360.371	0.643249	-1.235.351	-0.449944	0	1	1	1	1	1	1	1	0.062194	-0.236010	-0.097295	-0.540699	0.551100	1.041.070	-0.464016	14	14	14									
1	1.330.954	1	0	0	1.143.856	0.440257	0.792251	-0.042286	-0.449944	0	1	1	0	1	1	1	1	0.1178.860	0.766399	-0.097295	-0.540699	-0.226345	-0.399289	-0.464016	10	10	11									
16	1.330.954	1	0	0	-0.685387	-0.479857	0.643249	1.150.779	-0.449944	0	1	0	0	1	1	1	1	1.062194	1.768.808	-0.097295	-0.540699	-1.003.789	-0.399289	-0.714236	14	16	16									
0	-0.546287	1	0	0	1.143.856	1.360.371	0.643249	-1.235.351	-0.449944	0	1	0	0	0	1	1	1	0.062194	0.766399	0.801479	-0.540699	-0.226345	-1.119.469	-0.213796	14	14	14									
18	0	-0.546287	1	0	1.143.856	1.360.371	0.643249	1.150.779	-0.449944	0	1	1	1	1	1	1	1	0.1054.472	-1.238.419	-0.097295	-0.540699	-0.226345	-1.119.469	0.036424	13	14	14									
0	-0.546287	1	0	0	0.229234	0.440257	2.227.751	-0.042286	-0.449944	1	1	0	1	1	1	1	0	0.1178.860	-0.236010	-0.097295	-0.540699	-1.003.789	0.320890	-0.213796	8	10	10									
19	1.229390	1	0	0	0.229234	-0.479857	0.643249	-1.235.351	3.589.323	0	1	0	1	1	1	1	1	0.1178.860	1.768.808	-1.700.867	0.583385	1.328.545	1.041.070	1.287.526	6	5	5									
21	-0.546287	1	1	1	1.143.856	0.440257	0.643249	-1.235.351	-0.449944	0	0	1	1	1	1	1	1	0.1054.472	-2.240.828	-0.097295	-0.540699	0.551100	1.041.070	-0.213796	8	10	10									
22	1	1.330.954	1	0	1.143.856	0.440257	0.643249	-0.042286	-0.449944	0	0	0	0	0	1	1	1	0.062194	0.766399	-1.896.683	-0.540699	-1.003.789	-1.839.649	-0.714236	13	14	15									
1	1.330.954	1	0	0	1.143.856	1.360.371	0.643249	-1.235.351	-0.449944	0	1	1	0	1	1	1	1	0.1178.860	0.766399	-0.097295	-0.540699	-1.003.789	1.041.070	-0.464016	12	15	15									
1	-0.546287	1	1	1	1.143.856	-0.479857	0.643249	-0.042286	-0.449944	0	1	1	0	1	1	1	1	0.062194	1.768.808	-1.896.683	-0.540699	0.551100	1.041.070	-0.464016	15	15	16									
25	-0.546287	1	1	1	-0.685387	-0.479857	0.792251	-0.042286	-0.449944	0	1	0	1	1	1	1	1	0.1178.860	0.766399	0.801479	0.583385	1.328.545	1.041.070	-0.714236	13	13	12									
0	1.330.954	0	0	0	-0.685387	1.360.371	0.643249	1.150.779	-0.449944	1	1	1	1	1	1	1	1	0.062194	-0.236010	-0.097295	-0.540699	-1.003.789	1.041.070	-0.464016	10	9	8									
27	0	-0.546287	1	0	-0.685387	-0.479857	0.643249	-1.235.351	2.242.901	0	1	1	0	0	1	1	1	0.3287.804	-1.238.419	-0.097295	-0.540699	0.551100	1.041.070	1.037.305	6	9	8									
28	1.330.954	1	0	0	-0.685387	-0.479857	0.643249	-1.235.351	-0.449944	0	1	1	0	1	1	1	1	0.062194	-1.238.419	-0.097295	-0.540699	-0.226345	1.041.070	-0.464016	12	12	11									
1	1.330.954	1	0	0	1.143.856	-0.479857	0.643249	-1.235.351	-0.449944	0	1	0	1	1	1	1	1	0.2171.138	-1.238.419	0.801479	0.583385	1.328.545	-1.839.649	-0.213796	15	16	15									
30	-0.546287	1	1	1	0.229234	1.360.371	0.643249	-0.042286	-0.449944	1	1	1	1	1	1	1	1	0.1178.860	-0.236010	-0.097295	-0.540699	-1.003.789	1.041.070	-0.213796	11	11	11									
1	-0.546287	1	0	1	1.143.856	1.360.371	0.643249	-0.042286	-0.449944	0	1	1	1	1	1	1	1	1.062194	0.766399	1.700.867	3.955.638	2.105.989	1.041.070	1.287.526	10	12	11									
32	1	1.330.954	1	0	1.143.856	1.360.371	0.643249	-0.042286	-0.449944	0	1	1	0	0	1	1	1	0.1178.860	0.766399	-0.097295	1.707.469	1.328.545	1.041.070	-0.714236	9	11	12									
33	1	1.330.954	1	0	1.143.856	1.360.371	0.792251	-0.042286	-0.449944	0	1	0	1	1	1	1	1	0.062194	-0.236010	-1.896.683	-0.540699	-1.003.789	1.041.070	-0.714236	17	16	17									
1	1.330.954	0	0	0	1.143.856	0.440257	0.643249	-0.042286	-0.449944	0	1	0	1	1	1	1	1	1.062194	1.768.808	-0.097295	-0.540699	-1.003.789	1.041.070	-0.714236	17	16	16									
1	1.330.954	1	1	1	0.229234	0.440257	0.643249	-0.042286	-0.449944	0	0	0	1	0	1	1	1	0.1178.860	-0.236010	-0.097295	-0.540699	-1.003.789	-1.119.469	-0.714236	8	10	12									
35	-0.546287	1	0	0	0.229234	-0.479857	0.643249	-1.235.351	-0.449944	0	1	1	0	0	1	1	1	0.1178.860	0.766399	-0.097295	-0.540699	-1.003.789	1.041.070	-0.714236	12	14	15									
37	0	1.330.954	1	0	-0.685387	0.440257	0.792251	-1.235.351	-0.449944	0	1	0	1	1	1	1	0	0.1054.472	1.768.808	-1.896.683	-0.540699	-1.003.789	1.041.070	-0.714236	8	7	6									

Pentru o vizualizare rapida, am realizat o histograma a coloanei G3, care arata ca majoritatea elevilor au note intre 8 si 15.



In final, am salvat datele preprocesate in fisierul student-mat-preprocessed-minmax.csv si student-mat-preprocessed-zscore.csv.

student-mat-preprocessed-minmax																																				
	school	sex	age	address	famsize	Pstatus	Medu	Fedu	travelttime	studytime	failures	schoolsup	famsup	paid	activities	nursery	higher	internet	romantic	famrel	freetime	goutot	Dalc	Walc	health	absences	G1	G2	G3	M						
1	0	0	0.429	1	0	1	1.000	1.000	0.333	0.333	0.000	1	0	0	0	1	1	0	0.750	0.500	0.750	0.000	0.000	0.500	0.080	5	6	6								
2	0	0	0.296	1	0	1	0.250	0.250	0.000	0.333	0.000	1	0	0	0	0	0	1	1	0	1.000	0.500	0.000	0.500	0.053	5	5	6								
3	0	0	0.000	1	1	1	0.250	0.250	0.000	0.333	1.000	1	0	1	0	1	1	1	1	0.750	0.500	0.250	0.250	0.500	0.133	7	8	10								
4	0	0	0.000	1	0	1	1.000	0.500	0.000	0.667	0.000	0	1	1	1	1	1	1	1	1.500	0.250	0.250	0.000	0.000	1.000	0.027	15	14	15							
5	0	0	0.143	1	0	1	0.750	0.750	0.000	0.333	0.000	0	1	1	0	1	1	0	0.750	0.500	0.250	0.000	0.250	1.000	0.053	6	10	10								
6	0	1	0.143	1	1	1	1.000	0.750	0.000	0.333	0.000	0	1	1	1	1	1	1	1	0	1.000	0.750	0.250	0.250	1.000	0.133	15	15	15							
7	0	1	0.143	1	1	1	0.500	0.500	0.000	0.333	0.000	0	0	0	0	1	1	1	1	0.750	0.750	0.750	0.000	0.000	0.500	0.000	12	12	11							
8	0	0	0.296	1	0	0	1.000	1.000	0.333	0.333	0.000	1	1	0	0	1	1	1	0	0.750	0.000	0.750	0.000	0.000	0.080	6	5	6								
9	0	1	0.000	1	1	1	0.750	0.500	0.000	0.333	0.000	0	1	1	1	1	1	1	1	0.750	0.250	0.250	0.000	0.000	0.000	16	18	19								
10	0	1	0.000	1	0	1	0.750	1.000	0.000	0.333	0.000	0	1	1	1	1	1	1	0	1.000	1.000	0.000	0.000	1.000	0.000	14	15	15								
11	0	0	0.000	1	0	1	1.000	1.000	0.000	0.333	0.000	0	1	1	0	1	1	1	1	0.500	0.500	0.500	0.000	0.250	0.000	10	8	9								
12	0	0	0.000	1	0	1	0.500	0.250	0.667	0.667	0.000	0	1	0	1	1	1	1	1	0	1.000	0.250	0.250	0.000	0.250	0.053	10	12	12							
13	0	1	0.000	1	1	1	1.000	1.000	0.000	0.000	0.000	0	1	1	1	1	1	1	1	0.750	0.500	0.500	0.000	0.000	0.000	10	14	14								
14	0	1	0.000	1	0	1	1.000	0.750	0.333	0.333	0.000	0	1	1	0	1	1	1	1	0	1.000	0.750	0.500	0.000	0.250	0.500	0.027	10	11							
15	0	1	0.000	1	0	0	0.500	0.500	0.000	0.667	0.000	0	1	0	0	1	1	1	1	0.750	1.000	0.250	0.000	0.500	0.000	13	14	16								
16	0	0	0.143	1	0	1	1.000	1.000	0.000	0.000	0.000	0	1	0	0	1	1	1	1	0.750	0.750	0.750	0.000	0.250	0.053	14	14	14								
17	0	0	0.143	1	0	1	1.000	1.000	0.000	0.667	0.000	0	1	1	1	1	1	1	1	0.500	0.250	0.500	0.000	0.250	0.080	13	14	14								
18	0	0	0.143	1	0	1	0.750	0.750	0.667	0.333	0.000	0	0	1	0	1	1	1	0	1.000	0.500	0.250	0.000	0.750	0.053	8	10	10								
19	0	0	0.296	1	0	1	0.750	0.500	0.000	0.000	1.000	0	1	0	1	1	1	1	1	0	1.000	1.000	0.250	0.750	1.000	0.213	6	5	5							
20	0	0	0.143	1	1	1	1.000	0.750	0.000	0.000	0.000	0	0	1	1	1	1	1	1	0.500	0.000	0.500	0.000	0.500	1.000	0.053	8	10	10							
21	0	1	0.000	1	0	1	1.000	0.750	0.000	0.333	0.000	0	0	0	0	1	1	1	1	0.750	0.750	0.000	0.000	0.000	0.000	13	14	15								
22	0	1	0.000	1	0	1	1.000	1.000	0.000	0.000	0.000	0	1	1	1	1	1	1	1	0	1.000	0.750	0.250	0.000	0.000	1.000	0.000	12	15	15						
23	0	0	0.143	1	1	1	1.000	0.500	0.000	0.333	0.000	0	0	0	1	1	1	1	1	0.750	1.000	0.000	0.000	0.500	1.000	0.027	15	15	16							
24	0	0	0.143	1	1	1	0.500	0.500	0.333	0.333	0.000	0	1	0	1	1	1	1	1	0	1.000	0.750	0.250	0.250	0.750	1.000	0.000	13	13	12						
25	0	0	0.000	0	0	1	0.500	1.000	0.000	0.667	0.000	1	1	1	1	1	1	1	1	0.750	0.500	0.250	0.000	0.000	1.000	0.027	10	9	8							
26	0	0	0.143	1	0	1	0.500	0.500	0.000	0.000	0.667	0	1	1	0	0	1	1	1	0.000	0.250	0.250	0.000	0.500	1.000	0.187	6	9	8							
27	0	1	0.000	1	0	1	0.500	0.500	0.000	0.000	0.000	0	1	1	0	1	1	1	1	0.750	0.250	0.250	0.000	0.250	1.000	0.027	12	12	11							
28	0	1	0.000	1	0	1	1.000	0.500	0.000	0.000	0.000	0	0	1	0	1	1	1	1	0.250	0.250	0.750	0.250	0.750	0.000	0.053	15	16	15							
29	0	0	0.143	1	1	1	0.750	1.000	0.000	0.333	0.000	1	1	0	1	1	1	1	1	0	1.000	0.500	0.500	0.000	0.000	1.000	0.053	11	11	11						
30	0	0	0.143	1	0	1	1.000	1.000	0.000	0.333	0.000	0	1	1	1	1	1	1	1	1.750	0.750	1.000	1.000	1.000	0.213	10	12	11								
31	0	1	0.000	1	0	1	1.000	1.000	0.000	0.333	0.000	0	1	1	0	0	1	1	1	0	1.000	0.750	0.250	0.500	0.750	1.000	0.000	9	11	12						
32	0	1	0.000	1	0	1	1.000	1.000	0.333	0.333	0.000	0	1	1	0	1	1	1	1	0.750	0.500	0.000	0.000	0.000	1.000	0.000	17	16	17							
33	0	1	0.000	0	0	1	1.000	0.750	0.000	0.333	0.000	0	1	0	1	1	1	1	1	1.750	1.000	0.250	0.000	0.000	1.000	0.000	17	16	16							
34	0	1	0.000	1	1	1	0.750	0.750	0.000	0.333	0.000	0	0	0	1	0	1	1	1	0	1.000	0.500	0.250	0.000	0.250	0.000	8	10	12							
35	0	0	0.143	1	0	1	0.750	0.500	0.000	0.000	0.000	0	1	1	0	0	1	1	1	0	1.000	0.750	0.500	0.000	0.000	1.000	0.000	12	14	15						
36	0	0	0.000	1	0	1	0.500	0.750	0.333	0.000	0.000	0	1	0	1	1	1	1	0	0.500	1.000	0.000	0.000	1.000	0.000	8	7	6								

student-mat-preprocessed-zscore																												
File Home Insert Draw Page Layout Formulas Data Review View Automate Help Acrobat																												
A B C D E F G H I J K L M N O P Q R S T U V W X Y Z AA AB AC																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100																												
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 4																												

2. Clasificarea

Pentru aceasta parte am continuat sa lucrez pe setul de date preprocesat mai sus, Student_Performance, bazat pe performanta elevilor, asupra caruia am aplicat un algoritm de clasificare.

Am ales algoritmul **Random Forest Classifier**, deoarece ofera o precizie buna si permite observarea influentei parametrilor diferiti asupra predictiei.

Am creat o noua coloana numita target, care indica daca un elev este promovat ($1 = \text{nota finala (G3)} \geq 10$) sau nepromovat ($0 = \text{G3} < 10$).

```
# -----
cv # rulam 10-fold cross-validation pentru fiecare configuratie
me for params in configs:
    model = RandomForestClassifier(random_state=42, **params)
    print(f"\nConfiguratie: {params}")
    res = {}
    for m in metrics:
        scores = cross_val_score(model, X, y, cv=cv, scoring=m)
        mean_score = scores.mean()
        print(f"{m:>9}: {mean_score:.3f}")
        res[m] = mean_score
    results.append({'config': params, **res})
```

Pentru evaluare am folosit metoda **10-fold cross-validation**, prin care setul de date este impartit in zece parti pentru antrenare si testare succesiva. Am analizat performanta modelului folosind mai multe metrice: **accuracy**, **precision**, **recall** si **F1-score**.

Am modificat parametrii `n_estimators` (numărul de arbori) și `max_depth` (adâncimea maximă a arborilor) pentru a observa efectul asupra performanței modelului. În general, creșterea numărului de arbori a dus la rezultate mai stabile, însă după un anumit prag îmbunătățirile au fost minore.

Pentru datele scalate cu **Min-Max**, cea mai bună configurație a fost `n_estimators = 100` și `max_depth = 10`, cu un scor **F1** mediu de aproximativ **0.82**.

Pentru datele scalate cu **Z-score**, performanțele au fost ușor mai scăzute, scorul **F1** maxim fiind de aproximativ **0.81**.

```
PS C:\Users\mihae\Desktop\master_an2\WADDM\Project_Student_Performance> .venv\Scripts\activate
(.venv) PS C:\Users\mihae\Desktop\master_an2\WADDM\Project_Student_Performance> python clasificare.py
Evaluare RandomForest pentru fisierul: student-mat-preprocessed-minmax.csv

Configuratie: {'n_estimators': 50, 'max_depth': 5}
accuracy: 0.694
precision: 0.706
recall: 0.932
f1: 0.803

Configuratie: {'n_estimators': 100, 'max_depth': 10}
accuracy: 0.724
precision: 0.736
recall: 0.921
f1: 0.817

Configuratie: {'n_estimators': 200, 'max_depth': None}
accuracy: 0.717
precision: 0.733
recall: 0.910
f1: 0.811

Rezumat configuratii pentru: student-mat-preprocessed-minmax.csv
{'config': {'n_estimators': 50, 'max_depth': 5}, 'accuracy': np.float64(0.6935256410256411), 'precision': np.float64(0.7063813416110649), 'recall': np.float64(0.9319088319088319), 'f1': np.float64(0.8030195395041811)}
{'config': {'n_estimators': 100, 'max_depth': 10}, 'accuracy': np.float64(0.7241666666666667), 'precision': np.float64(0.736251119097609), 'recall': np.float64(0.9207977207977208), 'f1': np.float64(0.817426814407975)}
{'config': {'n_estimators': 200, 'max_depth': None}, 'accuracy': np.float64(0.7166666666666666), 'precision': np.float64(0.7334388202904028), 'recall': np.float64(0.9095441595441596), 'f1': np.float64(0.8114628648553617)}

[OK] Rezultate salvate in 'rezultate_clasificare_rf_cv10_student-mat-preprocessed-minmax.csv.csv'
```

Evaluare RandomForest pentru fisierul: student-mat-preprocessed-zscore.csv

Configuratie: {'n_estimators': 50, 'max_depth': 5}
f1: 0.810

Configuratie: {'n_estimators': 100, 'max_depth': 10}
accuracy: 0.706
precision: 0.723
recall: 0.913
f1: 0.807

Configuratie: {'n_estimators': 200, 'max_depth': None}
accuracy: 0.699
precision: 0.715
recall: 0.917
f1: 0.803

Rezumat configuratii pentru: student-mat-preprocessed-zscore.csv

```
{'config': {'n_estimators': 50, 'max_depth': 5}, 'accuracy': np.float64(0.7061538461538461), 'precision': np.float64(0.7149027312011058), 'recall': np.float64(0.9354700854700855), 'f1': np.float64(0.8100239378128171)}  
{'config': {'n_estimators': 100, 'max_depth': 10}, 'accuracy': np.float64(0.7062820512820512), 'precision': np.float64(0.7232473830196791), 'recall': np.float64(0.9131054131054132), 'f1': np.float64(0.8068443498658853)}  
{'config': {'n_estimators': 200, 'max_depth': None}, 'accuracy': np.float64(0.6987820512820513), 'precision': np.float64(0.7151239283592224), 'recall': np.float64(0.917094017094017), 'f1': np.float64(0.8033247206488662)}
```

[OK] Rezultate salvate in 'rezultate_clasificare_rf_cv10_student-mat-preprocessed-zscore.csv.csv'
(.venv) PS C:\Users\mihae\Desktop\master_an2\VAADDH\Proiect_Student_Performance>

În concluzie, algoritmul Random Forest Classifier s-a dovedit potrivit pentru problema de clasificare a promovării elevilor, oferind rezultate stabile în cadrul validării 10-fold.

Compararea celor două metode de scalare a arătat că Min-Max a oferit rezultate ușor mai bune decât Z-score în acest caz, însă diferențele nu sunt foarte mari (0.01).

Eliminarea notelor G1 și G2 a permis o evaluare corectă a modelului, fără informații care influențează direct nota finală.

3. Clustering

3.1 Rezumat Articol „K-means clustering algorithm and Python implementation”

Articolul "K-means clustering algorithm and Python implementation" prezinta algoritmul K-means, un algoritm clasic de invatare nesupravegheata folosit pentru gruparea automata a datelor in clustere. Principiul de baza este ca elementele din acelasi cluster trebuie sa fie similare, iar distanta, de obicei cea Euclidiană, este utilizata pentru a masura aceasta apropiere.

Autorul descrie propria implementarea a algoritmului, realizata printr-un API si cod in Python, folosind Java pentru generarea datelor experimentale. Procedura urmeaza pasii standard: alegerea aleatorie a K centroizi, atribuirea fiecarui obiect la cel mai apropiat centru si recalcularea centrozilor pe baza mediei, repetand procesul pana la convergenta.

Experimentele au fost efectuate pe un set de date cu un milion de inregistrari, iar rezultatele obtinute au fost similare cu cele generate de implementarea K-means din Apache Spark. Articolul subliniaza si limitarile algoritmului, precum necesitatea stabilirii valorii K si sensibilitatea la valori anormale.

3.2 Aplicarea algoritmului K-means pe setul Student Performance (Python, sklearn)

Am aplicat algoritmul **K-means** pe setul Student Performance preprocesat la inceput, folosind variantele scalate cu **Min-Max** si **Z-score**.

Datele au fost incarcate din fisiere CSV, iar coloanele G1, G2 si G3 au fost eliminate, deoarece K-means este un algoritm nesupravegheat.

Modelul a fost antrenat cu $K = 2, 3$ si 4 , folosind implementarea KMeans din scikit-learn, cu initializare aleatorie a centroizilor. Pentru fiecare valoare K au fost calculate metricele **WCSS** si **Silhouette Score**.

Cele mai bune rezultate au fost obtinute pentru datele scalate cu **Min-Max**, in cazul $K = 4$, unde **Silhouette Score** a fost aproximativ **0.30**, ceea ce arata ca studentii pot fi impartiti mai clar in patru grupuri.

Pentru varianta scalata cu **Z-score**, cea mai buna configuratie a fost $K = 2$, dar cu o separare mai slaba intre grupuri, unde **Silhouette Score** a fost aproximativ **0.15**.

Pentru evaluarea modelului am ales **WCSS** deoarece este metrica folosita direct de algoritmul K-means pentru a forma grupuri cat mai compacte.

Am utilizat si **Silhouette Score** pentru a verifica cat de bine sunt separate grupurile intre ele.

Rezultatele arata ca studentii pot fi impartiti in mai multe categorii, in functie de nivelul de implicare scolara, cum ar fi prezenta la cursuri, timpul de studiu si suportul familial. Chiar daca diferentele nu sunt foarte mari, modelul ajuta la intelegerea modului in care studentii se grupeaza in functie de comportamentul lor scolar.

```
PS C:\Users\mihae\Desktop\master_an2\MADDM\Proiect_Student_Performance> .venv\Scripts\activate
(.venv) PS C:\Users\mihae\Desktop\master_an2\MADDM\Proiect_Student_Performance> python kmeans_clustering.py

Clustering K-means pentru fisierul: student-mat-preprocessed-minmax.csv
Dimensiune set de date: (395, 39)

Rulare K-means cu K = 2
Silhouette score: 0.2610633291079305
WCSS: 49486287.903445095

Rulare K-means cu K = 3
Silhouette score: 0.26005901931818287
WCSS: 40159874.09994477

Rulare K-means cu K = 4
Silhouette score: 0.29596813929123594
WCSS: 33104090.517095327
Rezultate metrice salvate in: rezultate_clustering_kmeans_student-mat-preprocessed-minmax.csv.csv
Clustere finale salvate in: clustere_kmeans_K4_student-mat-preprocessed-minmax.csv

Clustering K-means pentru fisierul: student-mat-preprocessed-zscore.csv
Dimensiune set de date: (395, 39)

Rulare K-means cu K = 2
Silhouette score: 0.15302502011299593
WCSS: 4510302412718668.0

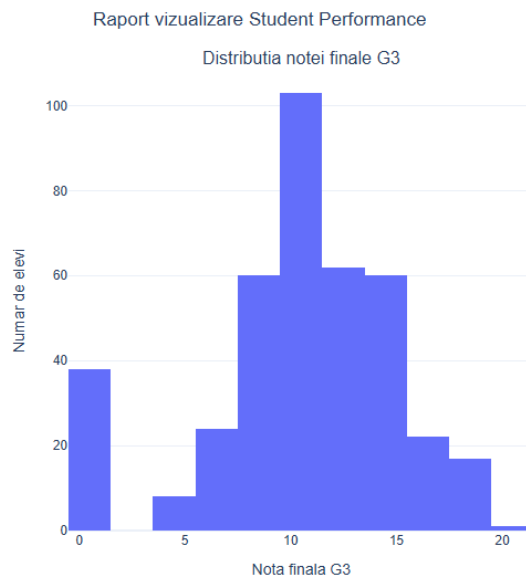
Rulare K-means cu K = 3
Silhouette score: 0.1045560419200294
WCSS: 4126518451452694.0

Rulare K-means cu K = 4
Silhouette score: 0.09622254738234111
WCSS: 3885829183780169.0
Rezultate metrice salvate in: rezultate_clustering_kmeans_student-mat-preprocessed-zscore.csv.csv
Clustere finale salvate in: clustere_kmeans_K2_student-mat-preprocessed-zscore.csv
```

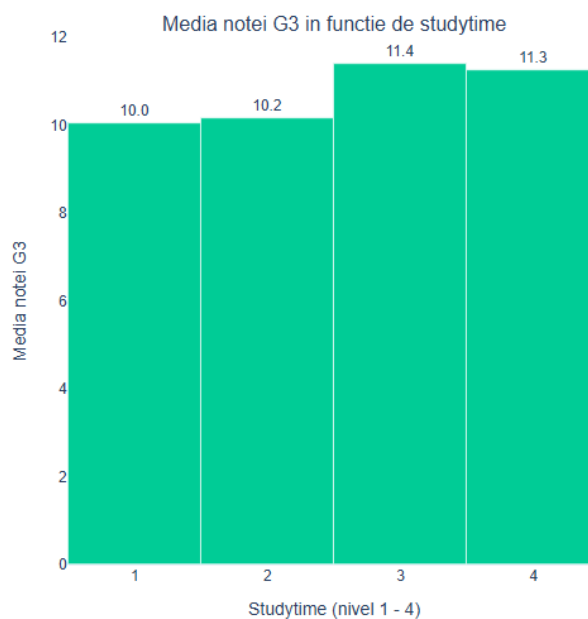
4. Vizualizarea Datelor

Setul de date Student Performance contine informatii despre elevi si nota finala G3.

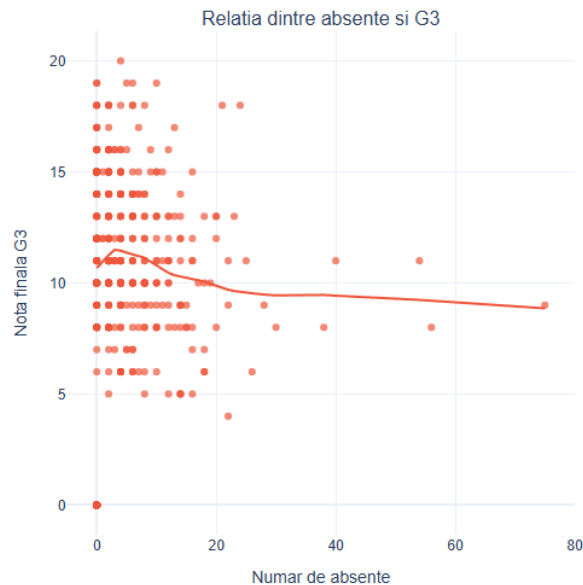
In primul grafic, histograma notei G3 arata o distributie concentrata in zona valorilor medii, cu cei mai multi elevi avand note intre aproximativ 8 si 14. Se observa si un numar vizibil de elevi cu nota finala G3 egala cu 0, care corespund cazurilor de nepromovare sau neprezentare la evaluare si nu reprezinta erori in setul de date.



Al doilea grafic, prezinta media notei G3 pe fiecare nivel de studytime, evidentiaza o crestere a mediei notei finale pe masura ce creste timpul de studiu. Elevii din grupele cu studytime 3 si 4 obtin, in medie, rezultate mai bune decat cei din grupele 1 si 2, ceea ce sugereaza ca timpul de studiu are un impact pozitiv asupra performantei scolare.



Al treilea grafic surprinde relatia dintre numarul de absente si nota finala G3, impreuna cu o curba de trend. Curba indica o relatie slaba si neliniara între absente si performanta, fara a sugera o legatura cauzala directa. In ansamblu, graficele arata ca timpul de studiu influenteaza mai clar rezultatele finale decat numarul de absente.



5. Concluzii

In concluzie, analiza setului de date Student Performance a aratat ca timpul de studiu este un factor mai relevant pentru performanta scolara decat absentele.

Algoritmul Random Forest a oferit rezultate stabile pentru clasificarea promovarii elevilor, iar clustering-ul K-means a permis identificarea unor grupuri de elevi cu comportamente similare.

Vizualizarea datelor a completat analiza, oferind o interpretare clara si intuitiva a rezultatelor obtinute.