

PavelloanCristian 311CAb

#### #Task\_1

am citit informatiile din fisier (my\_file = pandas.read\_csv('train.csv'))  
functia shape returneaza nr de linii si coloane sub forma de tuple  
nr\_col = my\_file.shape[0] #coloane  
nr\_lin = my\_file.shape[1] #linii  
metoda .dtypes returneaza un series cu numele coloanelor si data type ul din fiecare coloana  
metoda .isnull() creeaza un dataframe cu true / false pentru valori lipsa  
metoda .sum() transforma df ul intr-un series si aduna acele valori de true  
metoda .duplicated().sum() un series cu numele col si True / False, unde valorile de True vor fi insumate

#### #Task\_2

am accesat coloana cu supravietuitori cu my\_file['Survived'] si folosit functia .mean() am facut media  
am determinat apoi procentul inmultind cu 100; nonsupravietuitorii vor fii 100 - supravietuitori  
pentru Pclass si Sex am calculat procentajul de clasificare;  
metoda .value\_counts() calculeaza de cate ori apare fiecare valoare din coloana respectiva  
normalize='True' transforma numarul de valori intr-o proportie  
am creat apoi o figura 'fig' cu 3 subploturi, axs  
pentru primul grafic:  
['Supravietuitori', 'nonsupravietuitori'] - labelurile  
[supravietuitori\_percentage, nonsupravietuitori\_percentage] - valorile de reprezentat prin bars  
color=['blue', 'orange'] - culorile pentru fiecare bar  
salvam graficul ca imagine fig.savefig('grafic\_task\_2.png');

#### #Task\_3

metoda .select\_dtypes selecteaza din data frame coloanele cu tipurile de date precizate in include  
pentru fiecare coloana selectata, se creeaza un grafic care reprezinta valoarea in functie de frecventa  
ei in aparitia in coloana

#### #Task\_4

my\_file.isnull().any() seteaza cu True coloanele care au elemente lipsa  
pasand valoarea ca argument la my\_file.columns, selectam doar acele coloane din lista de coloane  
metoda .tolist() transforma lista pandas intr-o lista python  
pentru fiecare coloana din cele cu valori lipsa se calculeaza procentul supravietuitorilor si al nonsupravietuitorilor, la fel ca la taskul 2; de asemenea, se calculeaza si procentul valorilor lipsa,

#### #Task\_5

pentru a reprezenta histograma trebuie sa definim labelurile cu range urile de varsta, cat si bin urile

label urile vor fi stringurile prezentate in cerinta; pentru binuri, avem nevoie de un capat superior drept,  
reprezentat de limita maxima de varsta atinsa in data frame; o calculam prin metoda .max()  
apoi folosim metoda .cut() pentru a imparti coloana Age in bin uri, cu labeluri din labels  
numaram apoi valorile unice (persoanele) prin value\_counts()  
desenam graficul folosind bars, unde elemente de pe axa X sunt etichetele din dataframe  
(indexurile),  
iar pe Y avem valorile efective

#### #Task\_6

```
persoane_pe_categorie = my_file.groupby(['categorie-varsta', 'Sex'], observed=False)
['Survived'].mean() * 100
```

se face gruparea persoanelor in functie de categoria de varsta si separarea dupa gen  
observed=false nu permite adaugarea unor categorii care nu sunt prezente in setul de date  
['Survived'].mean() \* 100 acceseaza coloana survived din data frame si calculeaza procentul  
se deseneaza apoi graficul, asemanator ca la task 5; plot este o metoda pandas folosita  
pentru grafice;

#### #Task\_7

```
my_file[my_file['Age'] < 18]
```

selecteaza doar randurile cu persoanele cu varsta < 18  
procentajul copiilor se calculeaza impartind numarul de copii la numarul total de persoane si  
inmultind cu 100  
procentul de supravietuire se calculeaza ca la subpunctul anterior  
graficul se deseneaza ca la taskul 2;

#### #Task\_8

```
metoda .groupby(['Pclass', 'Survived'])
```

grupeaza persoanele care au aceeasi clasa si stare de supravietuire  
metoda .transform aplica o functie asupra fiecarei valori din coloanele selectate  
pentru metoda transform avem nevoie de o functie, de aceea cream o functie lambda  
metoda .fillna(x.mean()) inlocuieste valorile lipsa (care sunt setate ca NaN -> Not a Number)  
cu media valorilor (x.mean())

#### #Task\_9

```
nume_noblime = my_file['Name'].str.split(',') #desparte textul in functie de virgule
nume_noblime = nume_noblime.str[1] #extrage partea care contine numele si prenumele
nume_noblime = nume_noblime.str.split('.') #desparte textul in functie de .
nume_noblime = nume_noblime.str[0].str.strip()
#str[0] este titlu de noblime, iar str.strip() sterge eventualele spatii albe de la inceput si sfarsit
se creeaza apoi coloana titlu_de_noblime
folosind metoda .groupby se grupeaza data frame ul dupa coloanele title_de_noblime si Sex
graficul se creeaza folosind metoda .plot() din pandas, ca la task 6
```

#### #Task\_10

folosind catplot din seaborn se creeaza histograma;  
data=my\_file.head(100) specifica ca lucram cu primele 100 de randuri

pclass e coloana plasata pe axa x a graficului  
 fare este coloana plasat pe axa y a graficului  
 survived este coloana folosita pentru a colora punctele  
 strip specifica ca vrem sa cream un grafic de puncte  
 jitter adauga o mica distantare intre puncte pentru a le face vizibile in caz ca se suprapun















