

# README

Maxim Anca Stefania, 314CC

Git: [https://github.com/theonlytruealex/PCLP3/tree/main\\_I](https://github.com/theonlytruealex/PCLP3/tree/main_I)

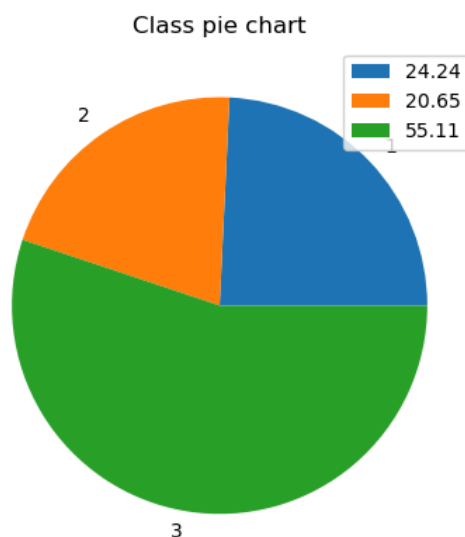
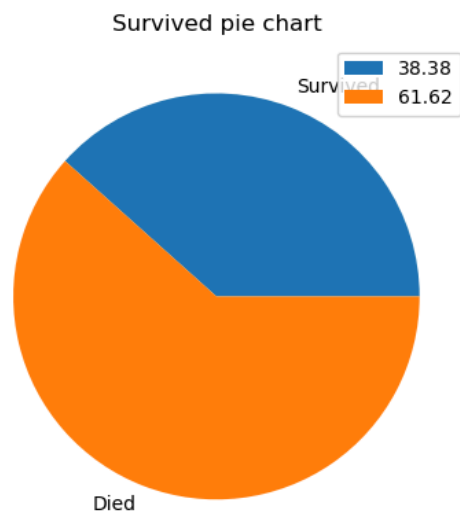
Partea I se afla pe branch-ul main\_I.

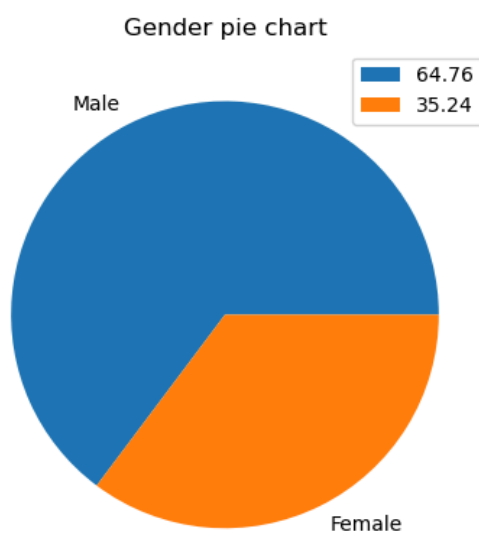
Cerința 1:

Am determinat, conform cerinței, numărul de coloane, tipul lor, numărul de spații din fiecare coloana, numărul de linii și numărul de linii duplicate.

Cerința 2:

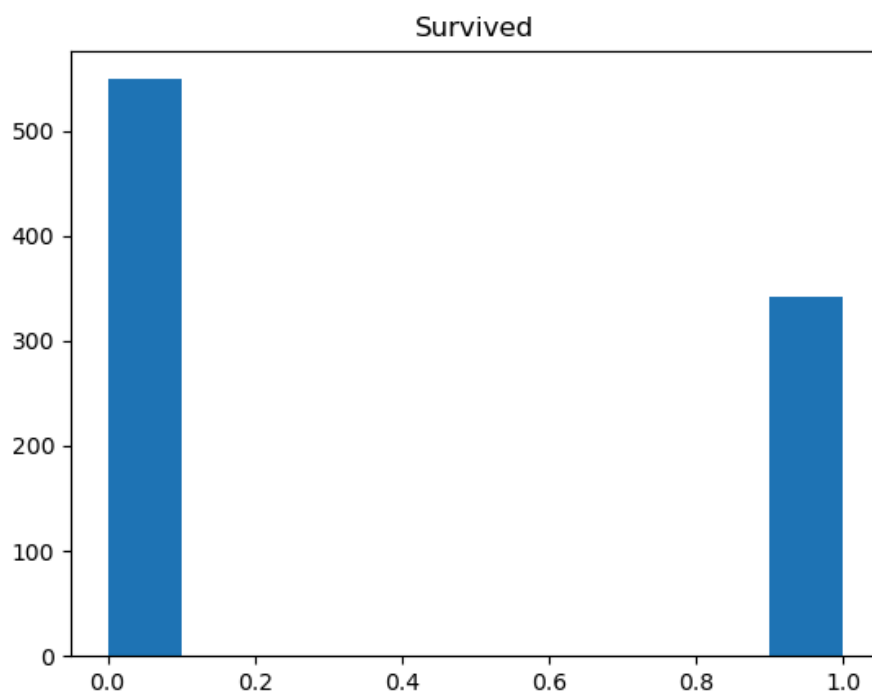
Am determinat separat numărul de campuri care aveau valoarea 1 pe coloana Survived și calculat procentul aferent. Pentru fiecare tip de clasa (1, 2, 3), am aflat numărul și procentul pasagerilor corespunzatori clasei respective (aplicand aceeași metoda ca pt Survived). Similar, am calculat numărul și procentul de bărbați și femei aflați la bord. Graficele rezultate sunt:

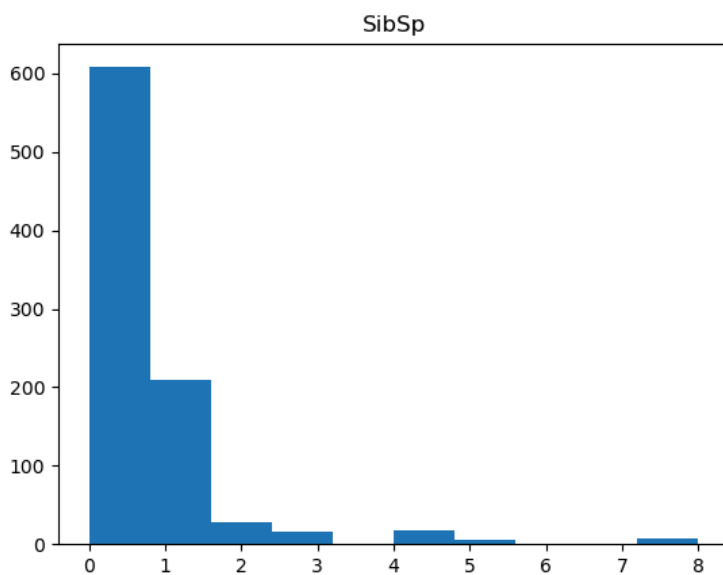
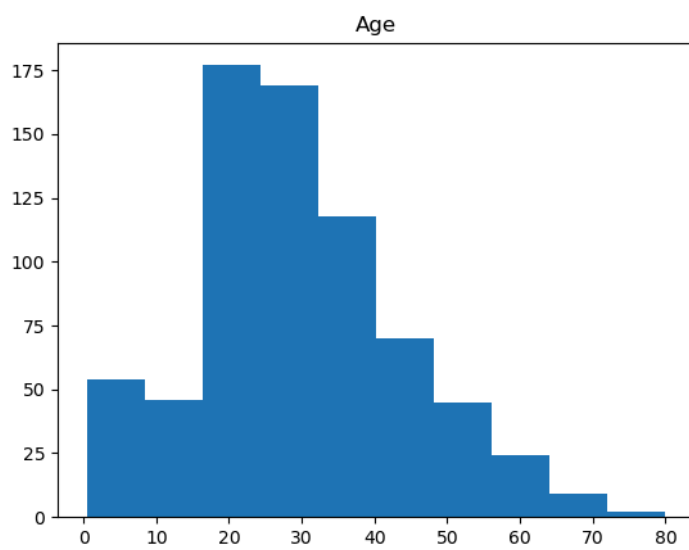
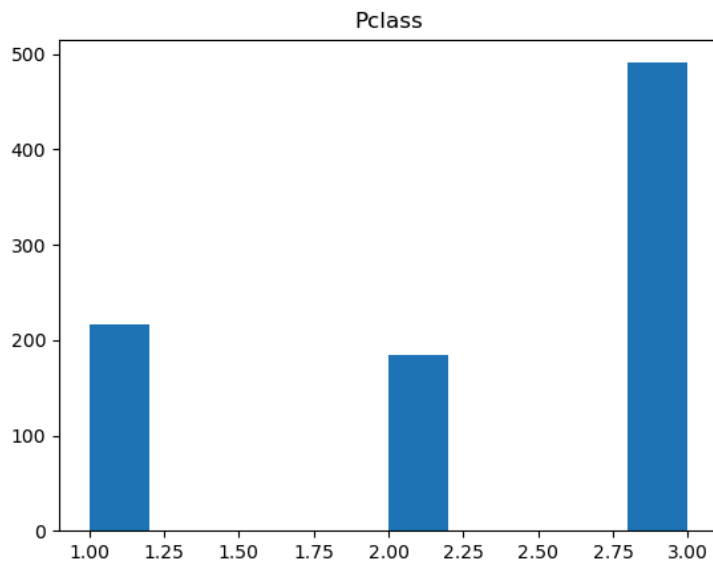


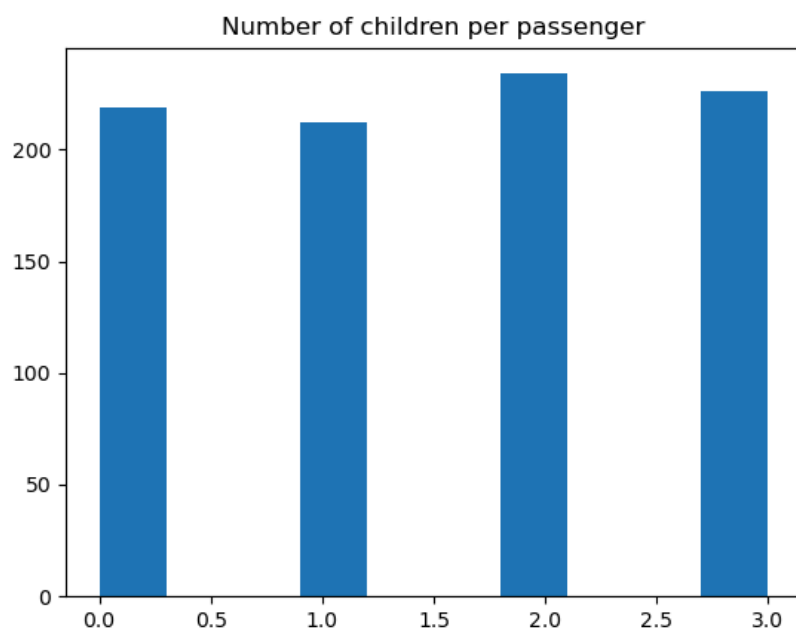
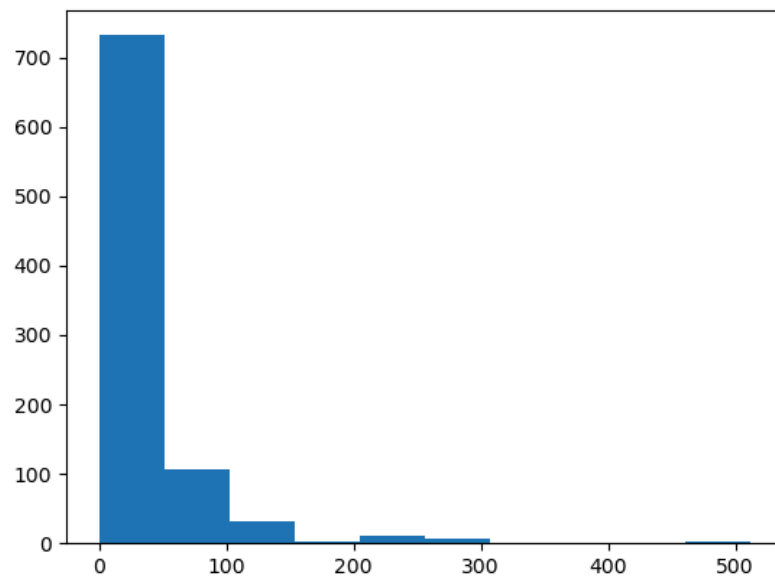
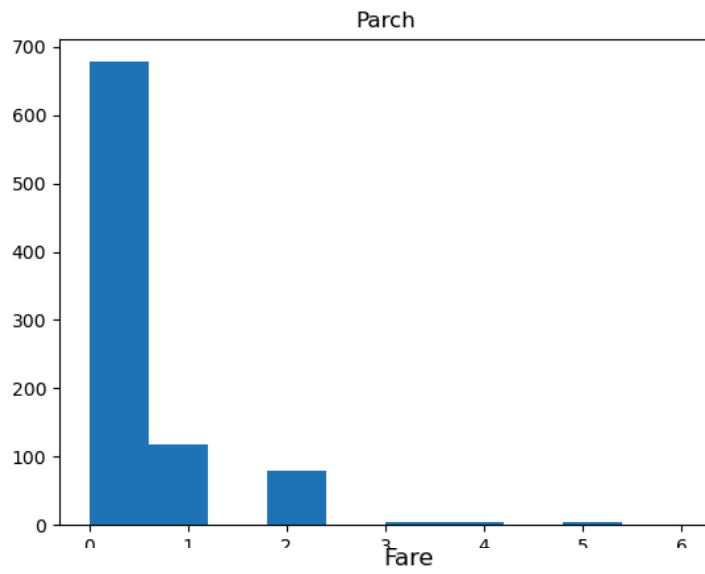


Cerința 3:

Am preluat doar coloanele cu date numerice din fișierul train.csv, asemanator ,  
și am realizat cate o histograma pentru fiecare





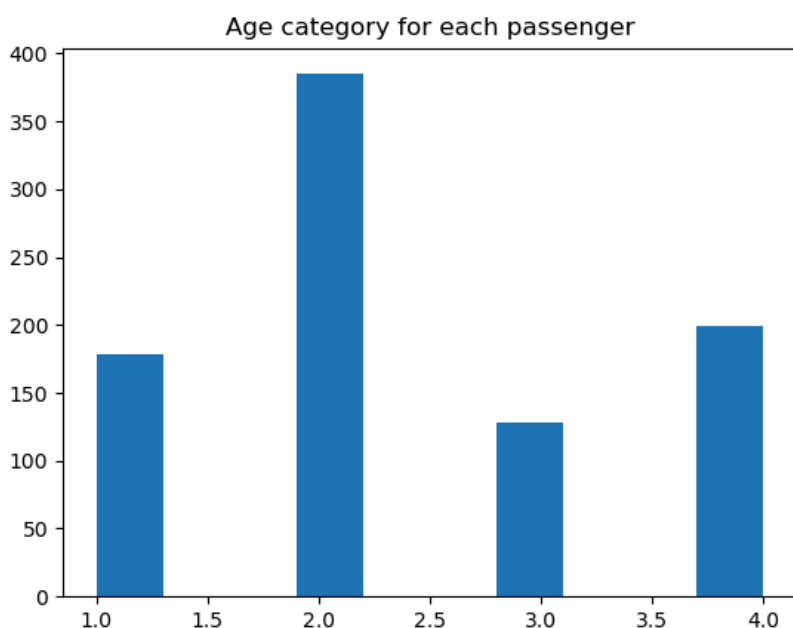


#### Cerința 4:

Am selectat coloanele cu date lipsa folosind funcția gaps (creata de mine), unde am construit o lista cu coloanele care contin date lipsa, in urma folosirii funcției isnull() și am calculat numărul de date lipsa si procentajul acestora, pentru fiecare coloana obtinuta. Am apelat în mod repetat funcția gaps pentru a obține datele aferente pentru pasagerii supraviețuitori și nesupraviețuitori.

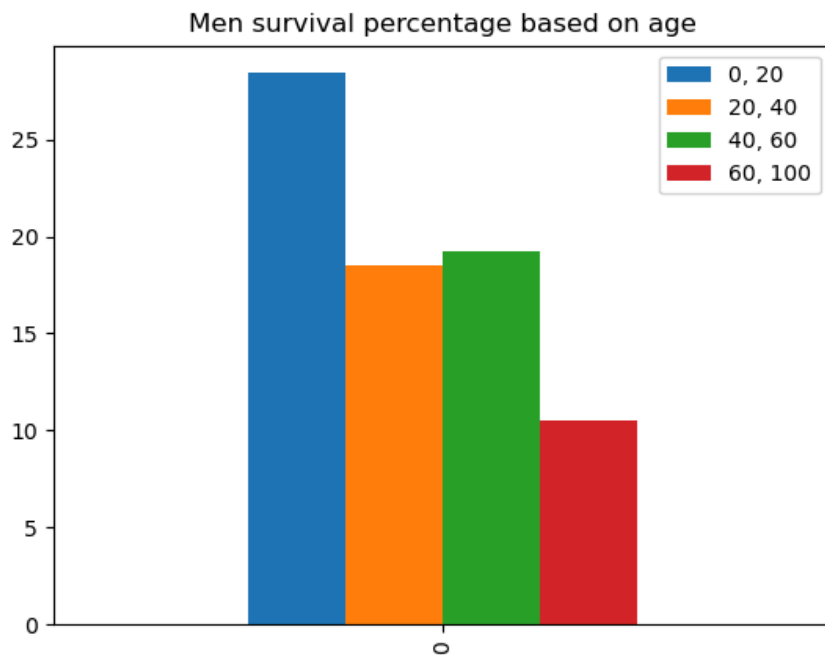
#### Cerința 5:

Am calculat numărul de persoane aparținând fiecărei categorii de varsta, după care am construit coloana suplimentară, unde, iterand prin fiecare rand al DataFrame ului, am determinat indexul corespunzător categoriei de varsta. La final, am afisat histograma corespunzătoare indecsilor obtinuti anterior.



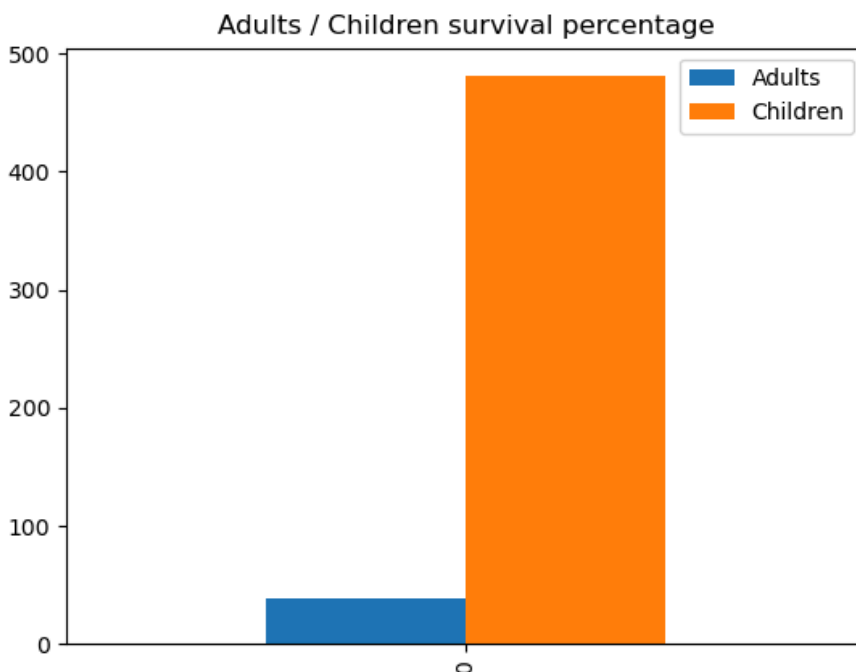
#### Cerința 6:

Folosind rezultatele obținute anterior, am calculat numărul bărbaților supraviețuitori din fiecare categorie. Pentru fiecare dintre cele 4 valori obținute, am creat un nou dataframe în care am asociat numărul obtinut cu intervalul de varsta din care face parte și am reprezentat rezultatele sub forma unui grafic de tip 'bar'.



Cerința 7:

Am calculat numărul și procentul copiilor aflați la bord și al copiilor supraviețuitori, după care, numărul și procentul adulților și al adulților supraviețuitori. Similar cu Cerința 6, am construit un nou dataframe pentru a reprezenta datele obținute, mapând procente cu categoria din care fac parte (Adults / Children) și le-am reprezentat sub formă de grafic de tip 'bar'.



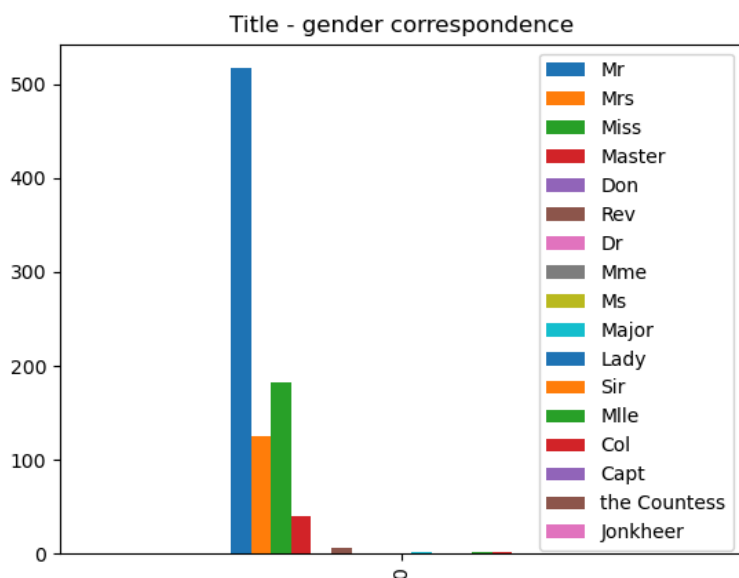
### Cerința 8:

Pentru a obține golurile din fiecare coloana, iterez prin fiecare coloana, preiau golurile cu metodele `.isnull()` și `.any()` (`.isnull()` îmi afișează o serie de valori bool, aleg doar seria care are măcar o valoare True cu `.any()`). Pentru fiecare pasager determinat anterior, calculez media valorilor pasagerilor care aparțin aceleiași clase și înlocuiesc rezultatul în poziția valorii lipsa (o localizez folosind metoda `.isna()`). Asemănător coloanelor cu valori numerice, fac și pentru cele cu valori categoricale, folosind metoda `.mode()` (care îmi returnează cel mai frecvent rezultat).

Salvez noul `dataFrame` într-un fișier `csv`.

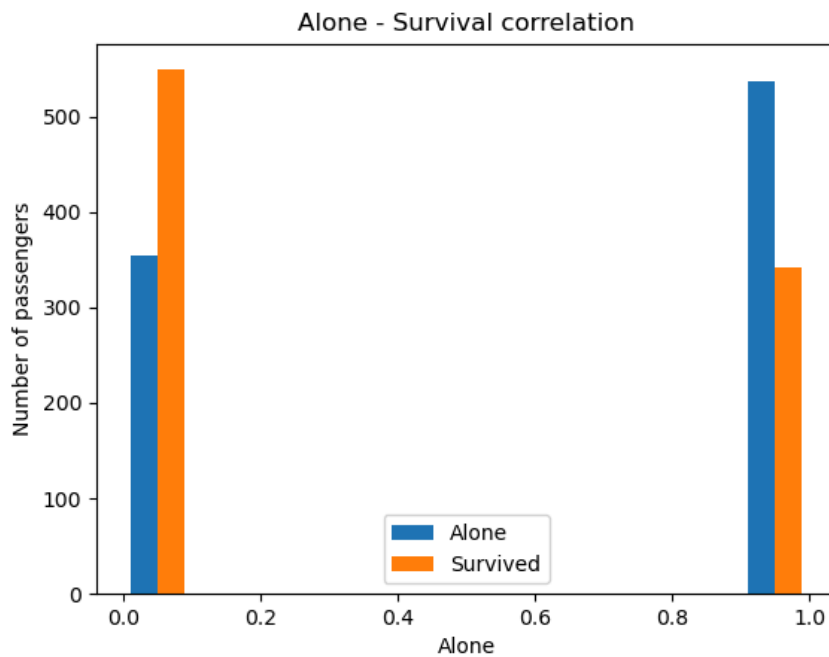
### Cerința 9:

Pentru a extrage toate titlurile posibile, am folosit metoda `split` în mod repetat asupra coloanei `Name` din `dataFrame`. După ce am determinat titlurile, le-am mapat cu sexul corespunzător (manual). Pentru fiecare coloana din setul de date, am verificat corespondența și am păstrat numărul acestora într-un `dataFrame`. Graficul rezultat este:



Cerința 10:

Preiau doar pasagerii care nu au rude la bord (coloanele SibSp si Parch au valorile 0) si creez un nou dataframe in care pun si valorile de pe coloana Survived. Am reprezentat corelatia celor doua valori prin urmatoarea histograma:



Pentru primele 100 de valori din fisierul train.csv, evidentiez corelatia valorilor de pe coloanele Fare, PClass si Survival printr-un grafic de tip catplot, astfel:

