

- Information Visualization -

Caracteristicile principale ale deceselor datorate problemelor cardiace.

- Project Design –

Cuprins

1. Obiectivele proiectului.....	3
2. Prezentarea setului de date.....	4
3. Alegerea graficele potrivite.....	6
4. Tehnologiile utilizate.....	7
5. Etapele principale	8
6. Concluzii	8

1. Obiectivele proiectului

Scopul principal al acestui proiect este de a analiza principalele motive pentru care oameni se îmbolnăvesc de probleme cardiace.

Vom observa distributia fiecărei variabile și corelatia acestora pentru stabili cât de important este să menținem anumite măsurători într-un anumit interval ce ne va fi oferit în urma vizualizării diferitelor aspecte ale setului de date. În urma acestui studiu vom lua în considerare și identificarea valorilor extreme (outliers) pentru fiecare variabila si identificarea pacientilor care prezinta aceste valori extreme.

În urma aflării mediei, deviatiei standard si varianței pentru fiecare variabila ne vom face și o părere despre distributia variabilelor care ne poate ajuta la identificarea unor tendinte sau anomalii.

Prin identificarea corelatiilor intre variabilele din setul tau de date vom înțelege mai bine cum variabilele afecteaza una pe alta.

În urma unei analize preliminare a setului de date am identificat anumite caracteristici precum grupul sau starea finala a pacientului. Acestea ne pot îndemna să divizăm setul de date în anumite subdiviziuni pentru a analiza mai bine situația anumitor tipuri de pacienți.

După ce am realizat un top al principalelor deficiențe care ar trebui să ne îngrijore, vom realiza o analiza a unor modele de regresie pentru a identifica variabilele care sunt asociate cu riscul de a dezvolta boli de inima, diabet sau alte deficiențe.

2. Prezentarea setului de date

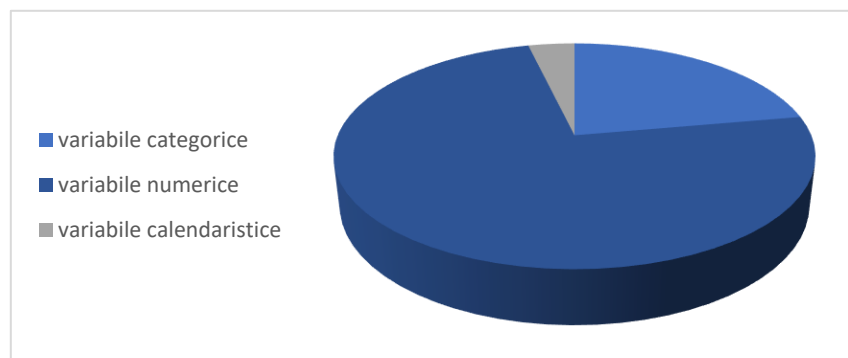
Setul de date ales de noi a fost extras dintr-o bază de date postgres numită **MIMIC III Clinical Test Database** ce conține o serie de înregistrări medicale variate pornind de la diverse tulburări cardiace până la o listă de tratamente oferite pacienților. Datele noastre au fost extrase prin intermediul unei operații de join asupra tabelelor ce dețin informații cu interes pentru noi. În urma procesării tabelor, s-a putut extrage un set de date cu 8239 înregistrări care se poate regăsi și studia la următoarea adresă.

În urma analizei acestui dataset, am extras următoarele date:

caracteristici demografice (varsta la momentul internării în spital, sex, etnie, greutate și înălțime), semne vitale (frecvență cardiacă[HR], tensiunea arterială sistolică[SBP], tensiunea arterială diastolică[DBP], tensiunea arterială medie[MBP], frecvența respiratorie[RR], temperatura corpului[BT], saturatia pulsului de oxygen[SPO2], urină), comorbidități (hipertensiune arterială, fibrilație atrială, boala cardiacă ischemică, diabet zaharat[DM], depresie, anemie hipoferică[HA], hiperlipidemie, boală renală cronică[CKD], boală pulmonară obstructivă cronică[BPOC]) și variabile de laborator (hematocrit, globule roșii, hemoglobina corpusculară medie [MCH], concentrația medie a hemoglobinei corpusculare [MCHC], volumul corpuscular mediu [MCV], lățimea distribuției globulelor roșii [RDW], numărul de trombocite, globule albe, neutrofile, bazofile, limfocite, timp de protrombină [PT], raport internațional normalizat [INR], NT-proBNP, creatin kinază, creatinină, azot ureic din sânge [BUN], glucoză, potasiu, sodiu, calciu, clorură, magneziu, gap anionic, bicarbonat, lactate, concentrația ionilor de hidrogen [pH], presiunea parțială a CO2 în sângele arterial și LVEF).

Caracteristicile demografice și semnele vitale extrase au fost înregistrate în primele 24 de ore de la fiecare internare și au fost măsurate variabilele de laborator pe toată durata șederii la unitățile de la terapie intensivă. Comorbiditățile au fost identificate folosind codurile ICD-9. Pentru variabile de laborator au fost realizate măsurători multiple.

50 + features		
10	40	2
categorical columns	numerical columns	date columns
8239 rows		



	<i>data type</i>
<i>group</i>	<i>categorical data</i>
<i>ID</i>	<i>numerical data</i>
<i>dod</i>	<i>time series</i>
<i>outcome</i>	<i>categorical data</i>
<i>age</i>	<i>numerical data</i>
<i>gendera</i>	<i>categorical data</i>
<i>BMI</i>	<i>numerical data</i>
<i>dofa</i>	<i>time series</i>
<i>hypertensive</i>	<i>categorical data</i>
<i>atrialfibrillation</i>	<i>categorical data</i>
<i>CHD with no MI</i>	<i>categorical data</i>
<i>diabetes</i>	<i>categorical data</i>
<i>deficiencyanemias</i>	<i>categorical data</i>
<i>depression</i>	<i>categorical data</i>
<i>Hyperlipemia</i>	<i>categorical data</i>
<i>Renal failure</i>	<i>categorical data</i>
<i>COPD</i>	<i>categorical data</i>
<i>heart rate</i>	<i>numerical data</i>
<i>Systolic blood pressure</i>	<i>numerical data</i>
<i>Diastolic blood pressure</i>	<i>numerical data</i>
<i>Respiratory rate</i>	<i>numerical data</i>
<i>temperature</i>	<i>numerical data</i>
<i>SP O2</i>	<i>numerical data</i>
<i>Urine output</i>	<i>numerical data</i>
<i>hematocrit</i>	<i>numerical data</i>
<i>RBC</i>	<i>numerical data</i>
<i>MCH</i>	<i>numerical data</i>
<i>MCHC</i>	<i>numerical data</i>

MCV	<i>numerical data</i>
RDW	<i>numerical data</i>
Leucocyte	<i>numerical data</i>
Platelets	<i>numerical data</i>
Neutrophils	<i>numerical data</i>
Basophils	<i>numerical data</i>
Lymphocyte	<i>numerical data</i>
PT	<i>numerical data</i>
INR	<i>numerical data</i>
NT-proBNP	<i>numerical data</i>
Creatine kinase	<i>numerical data</i>
Creatinine	<i>numerical data</i>
Urea nitrogen	<i>numerical data</i>
glucose	<i>numerical data</i>
Blood potassium	<i>numerical data</i>
Blood sodium	<i>numerical data</i>
Blood calcium	<i>numerical data</i>
Chloride	<i>numerical data</i>
Anion gap	<i>numerical data</i>
Magnesium ion	<i>numerical data</i>
PH	<i>numerical data</i>
Bicarbonate	<i>numerical data</i>
Lactic acid	<i>numerical data</i>
PCO2	<i>numerical data</i>
EF	<i>numerical data</i>

3. Alegerea graficelor potrivite

În urma stabilirii tipurilor de date întâlnite în setul de date ales trebuie să stabilim principalele grafice pe care le vom realiza pe parcursul studiului unei anumite variabile.

Atunci când vine vorba de reprezentarea unei variabile într-un mod vizual, alegerea tipului de grafic potrivit este esențială pentru a transmite informațiile în mod eficient și ușor de înțeles. Alegerea unui grafic nepotrivit poate duce la o interpretare incorectă a datelor și poate afecta deciziile care sunt luate în urma analizei.

Alegerea celui mai bun tip de grafic se realizează prin stabilirea tipurilor de date ale variabilelor și analizarea mai multor tipuri de grafice.

În setul de date ales de noi există 3 tipuri principale de variabile: categorice, numerice și interval de timp. În funcție de tipul de date stabilit avem următoarele posibilități de grafice ce se pot regăsi în tabelul alăturat.

	<i>variabile categorice</i>	<i>variabile numerice</i>	<i>variabile de tip date calendaristice</i>
<i>line chart</i>		<i>X</i>	<i>X</i>
<i>pie chart</i>	<i>X</i>		
<i>bar chart</i>	<i>X</i>		
<i>frequency chart</i>	<i>X</i>		
<i>strip chart</i>	<i>X</i>	<i>X</i>	<i>X</i>
<i>box and whisker chart</i>		<i>X</i>	<i>X</i>
<i>quartile chart</i>	<i>X</i>	<i>X</i>	
<i>histogram</i>		<i>X</i>	<i>X</i>
<i>scatter plot</i>		<i>X</i>	
<i>percentage chart</i>		<i>X</i>	
<i>radar chart</i>		<i>X</i>	
<i>Gantt chart</i>			<i>X</i>
<i>heatmap calendar chart</i>			<i>X</i>

4. Tehnologiile utilizate

Tehnologiile folosite de echipa noastră pentru proiect vor fi **Pyspark**(python) si **PowerBI**.

Pyspark este o librărie Python pentru procesarea datelor în cadrul platformei Apache Spark. Apache Spark este un framework de procesare a datelor în memorie, conceput pentru a procesa cantități mari de date în mod eficient și rapid, utilizând un model distribuit de calcul.



Pyspark este util pentru proiectul nostru, deoarece oferă o mulțime de funcționalități puternice pentru procesarea și analiza datelor. Printre aceste funcționalități, vom aminti de:

1. **Preprocesarea datelor:** Pyspark poate fi utilizat pentru a încărca, curăța și transforma datele într-un format adecvat pentru analiza ulterioară. Datorită capacității sale de a procesa cantități mari de date, pyspark este potrivit pentru prelucrarea datelor în timp real sau a datelor istorice.
2. **Analiza datelor:** Pyspark oferă o serie de funcții pentru analiza datelor, cum ar fi agregarea, filtrarea și sortarea datelor. De asemenea, permite utilizarea de algoritmi de machine learning pentru modelarea datelor și dezvoltarea de modele predictibile.
3. **Vizualizarea datelor:** Pyspark poate fi integrat cu alte librării de vizualizare a datelor, cum ar fi Matplotlib sau Seaborn, pentru a crea grafice și diagrame care să ajute la înțelegerea datelor.



Power BI este un instrument puternic pentru proiectele de data science, care permite utilizatorilor să vizualizeze, să analizeze și să partajeze datele într-un mod interactiv și ușor de înțeles. Power BI poate fi utilizat în cadrul proiectului nostru pentru a integra și analiza date pentru a crea vizualizări personalizate.

5. Etapele principale

step 0	extragerea datelor din tabelele bazei de date			
PySpark & python			Power BI	
step 1	_analizarea datelor lipsă	PROCESAREA DATELOR		
step 2	stabilirea modalităților de completare a setului de date			
step 3	procesarea datelor extrase			
step 4	stabilirea strategiei de analiza a datelor și tipurile de grafice folosite			
step 5	calcularea statisticilor descriptive (mediana,media, deviatia standard)	ANALIZA EXPLORATORIE A DATELOR	step 6	calcularea statisticilor descriptive (mediana,media, deviatia standard)
step 7	distributia fiecărei variabile (histograms, density plot, bar plot, box plot)		step 8	distributia fiecărei variabile (histograms, density plot, bar plot, box plot)
step 9	identificarea anomaniilor si valorilor extreme		step 10	identificarea anomaniilor si valorilor extreme
step 11	relatiile dintre variabile setului de date (line plots scatter plots, correlation matrix)		step 12	relatiile dintre variabile setului de date (line plots scatter plots, correlation matrix)
step 13	crearea unui model de pentru a identifica variabilele care sunt asociate cu riscul de a dezvolta boli de inima, diabet sau alte deficiențe.	REZULTATE FINALE	step 14	crearea unui dashboard interactiv în care să prezentăm informațiile descoperite în urma analizei exploratorii a datelor.

6. Concluzii

Rezultatul final al acestui proiect este de a scoate în evidență caracteristicile importante care sunt semnificative pentru decesul indivizilor care au participat la testele realizate, de a realiza o serie de modele de regresie care să ne ofere mai multe detalii cu privire la caracteristicile corelate în prezicerea bolilor cardiovasculare și nu numai.