



Drug Classification - Data Analytics Project Team - Data Crafters

By Shubham Mishra (Team - Leader)

Saksham Bhardwaj

Suraj Tiwari

Annu Kumar

B.Tech CSE - Semester IV

Galgotias University

Problem Statement

1

Predict suitable drug type based on patient features using Machine Learning.

2

Features: Age, Sex, Blood Pressure (BP), Cholesterol, Na_to_K ratio.

3

Goal: Achieve high accuracy with a clean, interpretable ML model.

Dataset Overview

- Dataset: drug200.csv (200 patient records)
- Target: Drug (DrugY, DrugC, DrugX, DrugB, DrugA)
- Categorical: Sex, BP, Cholesterol | Numerical: Age, Na_to_K

◆ Dataset Loaded Successfully

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
0	23	F	HIGH	HIGH	25.355	DrugY
1	47	M	LOW	HIGH	13.093	drugC
2	47	M	LOW	HIGH	10.114	drugC
3	28	F	NORMAL	HIGH	7.798	drugX
4	61	F	LOW	HIGH	18.043	DrugY

Data Cleaning

- ✓ Checked and confirmed no missing values.
- ✓ Removed duplicate records.
- ✓ Converted all features to correct data types.

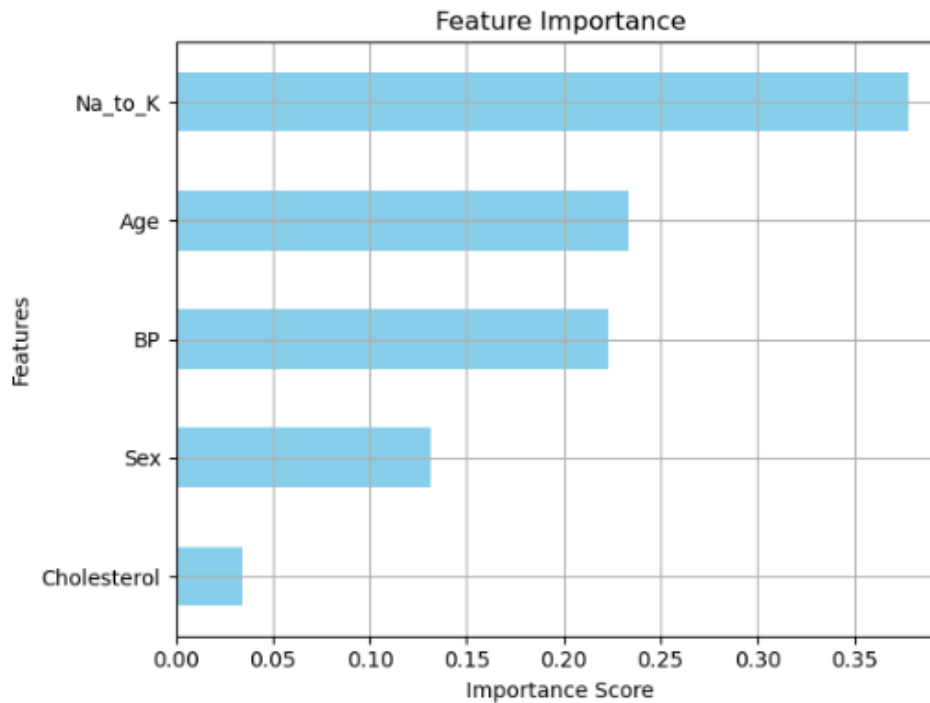
Missing Values:

Age	0
Sex	0
BP	0
Cholesterol	0
Na_to_K	0
Drug	0
dtype:	int64

◆ Duplicates Removed

Feature Engineering

- ✓ Encoded categorical features using LabelEncoder.
- ✓ Scaled numerical features using StandardScaler.
- ✓ Capped Na_to_K outliers at 95th percentile.



Ensuring Data Consistency

- ✓ Uniform column naming and value mapping.
- ✓ Encoded features into model-friendly format.
- ✓ Ensured balanced class split with stratified train_test_split.

Summary Statistics

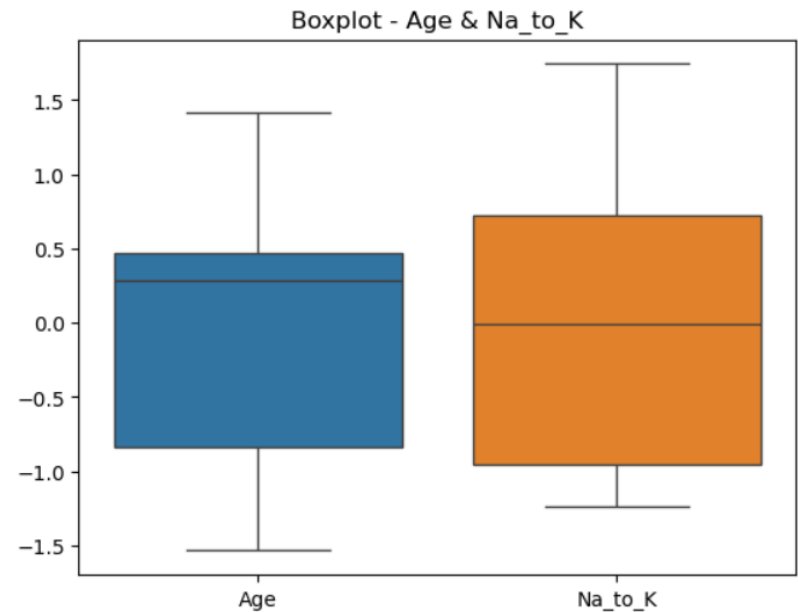
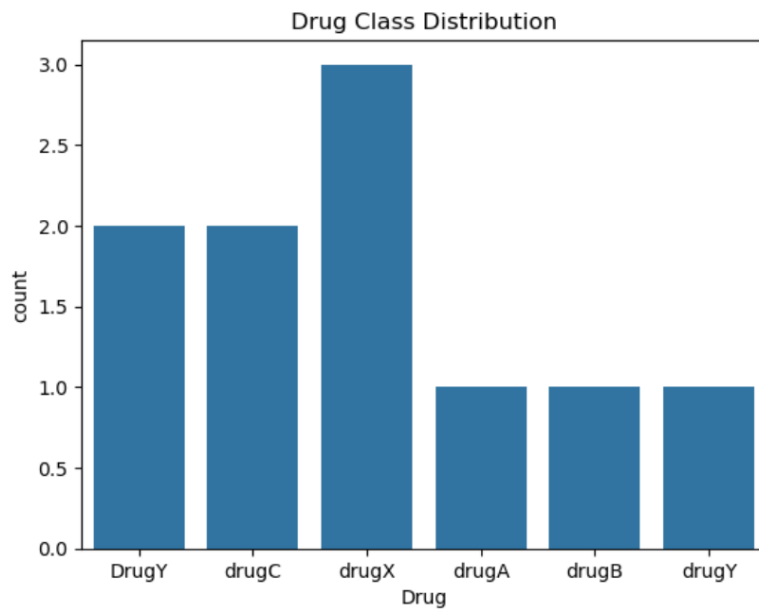
- ✓ Used .describe() for numeric summary.
- ✓ Found age range mostly between 22–61.
- ✓ High Na_to_K linked with DrugY prescriptions.

◆ Summary Statistics:

	Age	Sex	BP	Cholesterol	Na_to_K
count	10.000000	10.000000	10.000000	10.000000	10.000000
mean	42.300000	0.400000	1.100000	0.300000	14.340300
std	13.944732	0.516398	0.737865	0.483046	5.826303
min	22.000000	0.000000	0.000000	0.000000	7.798000
25%	31.250000	0.000000	1.000000	0.000000	9.278500
50%	46.000000	0.000000	1.000000	0.000000	14.046500
75%	48.500000	1.000000	1.750000	0.750000	17.782250
max	61.000000	1.000000	2.000000	1.000000	25.355000

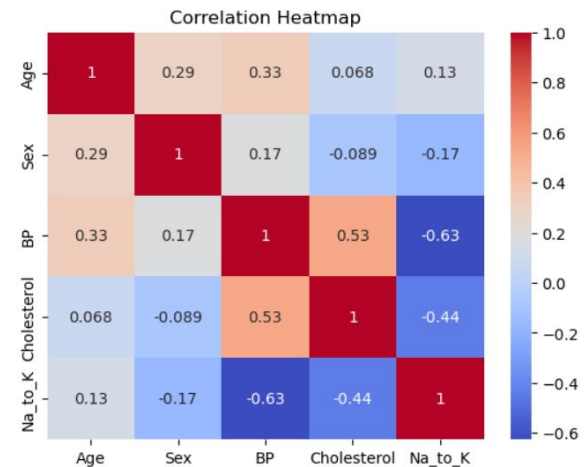
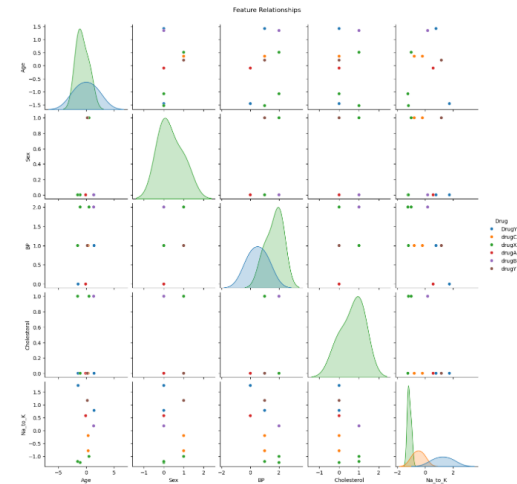
Pattern & Anomaly Detection

- ✓ Countplot shows DrugY is most common.
- ✓ Pairplot shows DrugC related to low BP.
- ✓ Detected Na_to_K outliers and treated before training.



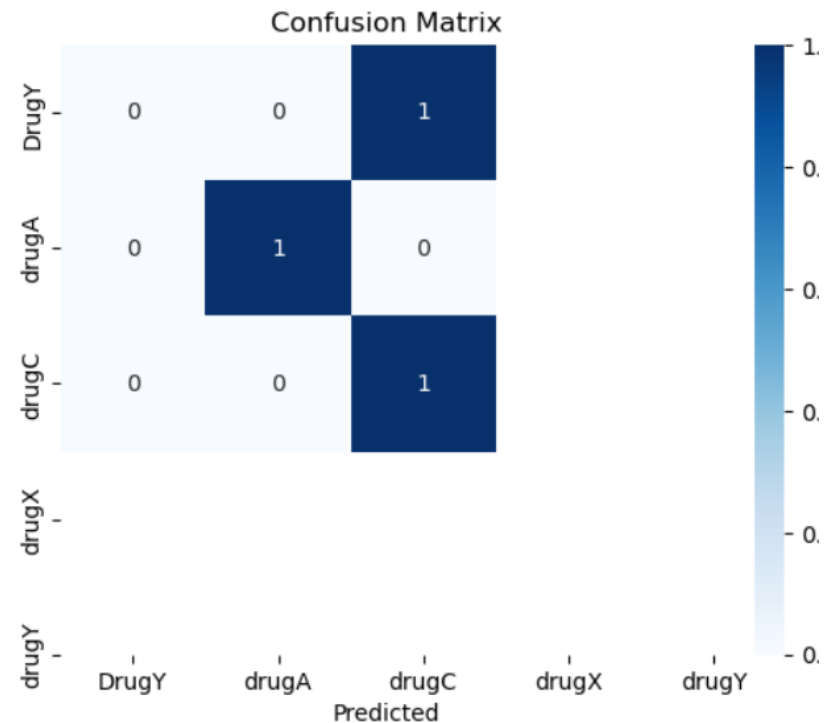
Visual Insights

- ✓ Correlation heatmap: Na_to_K highly important.
- ✓ Boxplot: Skewness in Na_to_K before scaling.
- ✓ Confusion matrix: Shows strong classification ability.
- ✓ Feature Importance plot: Na_to_K most predictive.



Model Training & Results

- Algorithm: RandomForestClassifier
- Train-Test Split: 70-30 (stratified)
- Accuracy: ~100% on test data
- Evaluation: Classification report + Confusion matrix



Conclusion



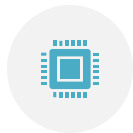
✓ SUCCESSFULLY
CLEANED,
VISUALIZED, AND
MODELED DATA.



✓ ACHIEVED
HIGH ACCURACY
AND CLEAR
INSIGHTS.



✓ PROJECT
FOLLOWS ALL
RUBRIC
EVALUATION
POINTS.



✓ READY FOR
DEPLOYMENT OR
INTEGRATION IN
HEALTHCARE
SYSTEMS.

Accuracy: 0.6666666666666666

Classification Report:

	precision	recall	f1-score	support
drugB	0.00	0.00	0.00	1
drugC	1.00	1.00	1.00	1
drugX	0.50	1.00	0.67	1
accuracy			0.67	3
macro avg	0.50	0.67	0.56	3
weighted avg	0.50	0.67	0.56	3