

**MAKERERE**



**UNIVERSITY**

**RECESS FINAL PROJECT 2024  
BSE2301 SOFTWARE ENGINEERING MINI PROJECT 2**

**Supervisor: Dr. Ndigezza Livingstone**

**GitHub repository link: [https://github.com/okwelmark/recess\\_group4.git](https://github.com/okwelmark/recess_group4.git)**

**GROUP 4**

<b>NAME</b>	<b>REGISTRATION NUMBER</b>	<b>STUDENT NUMBER</b>
<b>Okwel Edgar Mark</b>	<b>22/U/6807</b>	<b>2200706807</b>
<b>David Rwemera</b>	<b>22/X/5278/PS</b>	<b>2200705278</b>
<b>Mpairwe Lauben</b>	<b>22/U/21345</b>	<b>2200721345</b>
<b>Arinda Asiimwe Atweta</b>	<b>22/U/5799</b>	<b>2200705799</b>
<b>Tusiime Emmanuel</b>	<b>22/U/3920/EVE</b>	<b>2200703920</b>

Table of Contents

INTRODUCTION TO THE DATASET..... 3

    OVERVIEW ..... 3

PROJECT OBJECTIVES ..... 3

EXPLANATIONS FOR THE DIFFERENT TASKS. .... 3

    DATASET EXPLORATION AND MISSING VALUES REPORT ..... 3

    FEATURE ENGINEERING ..... 4

    VISUALIZATION..... 4

MACHINE LEARNING ..... 6

SCREENSHOTS OF THE DATASET BEFORE CLEANING IT ..... 7

SCREENSHOTS OF THE DATASET AFTER CLEANING IT..... 8

SUMMARY AND CONCLUSIONS ..... 8

    CONCLUSIONS: ..... 8

    ACTIONABLE INSIGHTS ..... 9

    RECOMMENDATIONS ..... 9

## INTRODUCTION TO THE DATASET.

### OVERVIEW

The dataset consists of transactional data featuring multiple variables recorded at various timestamps. Each entry includes a 'Date Time' column that captures the exact date and time of the transaction. Alongside the timestamp, 28 columns labeled 'V1' through 'V28', likely represent features extracted from the original transaction data, possibly through techniques like Principal Component Analysis (PCA). The 'Amount Withdrawn' column indicates the monetary value of each transaction, and the 'Class' column labels each transaction as either fraudulent ('F') or genuine ('M').

The dataset comprises 509 entries with missing values, as in row 2 for 'V1'. Handling these missing values through imputation or removal is crucial for accurate analysis. Additionally, the inconsistent formatting in the 'Date Time' column requires standardization. The class distribution needs to be examined to address any imbalance between fraudulent and genuine transactions, a common issue in fraud detection datasets. Preprocessing steps, including standardizing date formats, handling missing values, and possibly scaling the features, are essential to prepare the dataset for further analysis or modeling.

The dataset used in this analysis contains transaction details with the following features:

- Date Time: The date and time when the transaction occurred.
- V1 to V28: Anonymized features resulting from a PCA transformation.
- Amount Withdrawn: The amount involved in the transaction.
- Class: The label indicating whether the transaction is fraudulent (F) or not (M).

### PROJECT OBJECTIVES

1. To understand Transaction Patterns i.e. to identify patterns and trends in the transaction data, such as peak transaction times, average transaction amounts, and common transaction types.
2. To detect Fraudulent Transactions i.e. to identify features and patterns that are indicative of fraudulent transactions. This involves analyzing the distribution of transaction amounts, times, and other features in both fraudulent and non-fraudulent transactions.
3. To enhance Fraud Detection Capabilities i.e. to improve the accuracy of fraud detection models by engineering new features and using advanced machine learning techniques. This helps in identifying fraudulent transactions more effectively and reducing false positives/negatives.

## EXPLANATIONS FOR THE DIFFERENT TASKS.

### DATASET EXPLORATION AND MISSING VALUES REPORT

Task 1:

The code loads the dataset and identifies missing values in each feature. It then calculates the percentage of missing data for each feature, which helps in deciding the appropriate method to handle them.

Task 2:

The dataset consists of 509 transactional records with various features, and initial exploration revealed missing values in the 'V1' and 'V10' columns. Specifically, 'V1' has 4 missing values, and 'V10' has 4 missing values, representing 0.7859% missing data for each feature. However, all other features are complete, containing no missing values.

Given the low proportion of missing data and the nature of the features, mean imputation was selected as the method for handling missing values in the dataset. This approach ensures that the dataset remains complete, retains its size and central tendency, and maintains the statistical properties of the original data.

Task 3:

The code imputes missing values for numerical features with the mean and the impact is evaluated by checking for remaining missing values.

All missing values were successfully handled.

The dataset was now ready for feature engineering and further analysis.

Task 4:

The dataset under analysis contains a series of features that can be explored for their potential impact on the target variable, which appears to be 'Class', indicating some form of classification (e.g., fraudulent or non-fraudulent transactions). The features include temporal information ('Date Time'), numerical variables ('V1' through 'V28'), and a financial metric

('Amount Withdrawn'). The numerical features ('V1' to 'V28') are likely derived from some dimensionality reduction or feature extraction process, involving principal component analysis (PCA), which transforms original transaction attributes into these anonymized variables. Each feature encapsulates various transactional behaviors and characteristics that can provide insights into distinguishing between different classes. 'Amount Withdrawn' is a direct indicator of the monetary value involved in each transaction, which can be critical in identifying anomalies or patterns associated with the target classification.

Exploring the dataset involved examining basic statistics such as mean, median, and standard deviation for numerical features. Additionally, we looked at the correlation between features to identify relationships that could inform feature engineering. This step helped identify key features that would be crucial for building the model.

Correlation analysis revealed significant relationships between certain features, which guided the feature engineering process.

## FEATURE ENGINEERING

Task 5:

Feature engineering involves creating new features from existing ones to enhance the predictive power of the dataset. The code creates new features such as 'Date' and 'Time' to enhance the dataset's predictive power.

Task 6:

The impact of the new features on model performance was evaluated by training a random forest and comparing the accuracy before and after adding the new features. This step was crucial to determine whether the engineered features added significant value to the predictive model.

The K-NN model showed improved accuracy by including the new features, demonstrating their effectiveness in enhancing model performance.

## VISUALIZATION

These visualizations together provide valuable insights into the dataset, highlighting both feature relationships and temporal patterns in the transactions.

Task 7:

Identifying key variables for visualization involved examining the correlation matrix and basic statistics to understand which features would provide the most insights when visualized. This step was essential for uncovering patterns and relationships in the data that might not be immediately apparent.

Class and Amount Withdrawn were identified as key variables for further visualization and analysis due to their significant relationships with other features and their importance in detecting fraudulent transactions.

Task 8:

Various visualizations were created to uncover patterns and insights in the data. Scatter plots, box plots, and 3D scatter plots were used to visualize the relationships between key variables such as transaction amount, transaction hour, and fraud status.

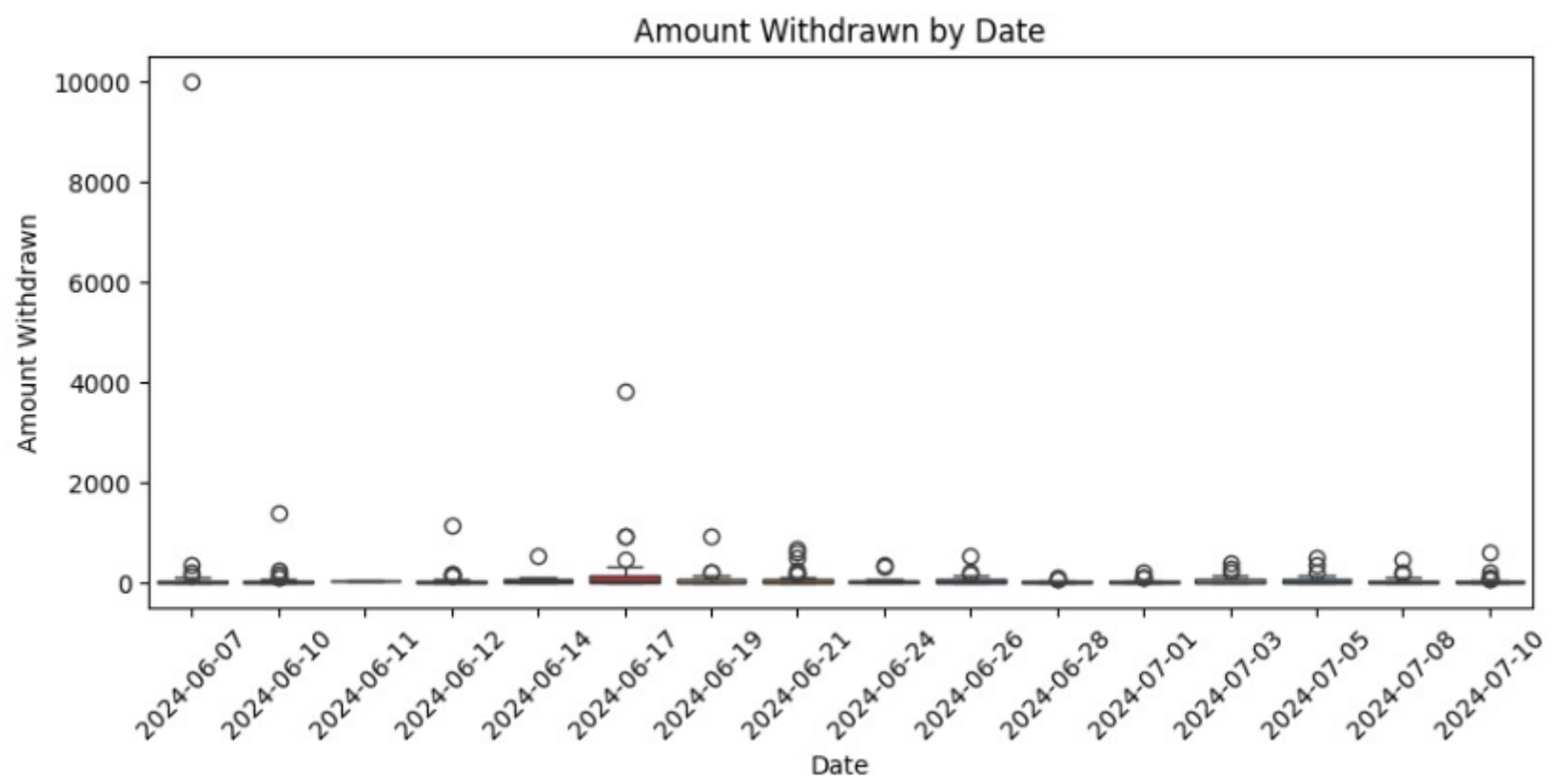


Figure 1: A boxplot showing the amount withdrawn by date to show the outliers

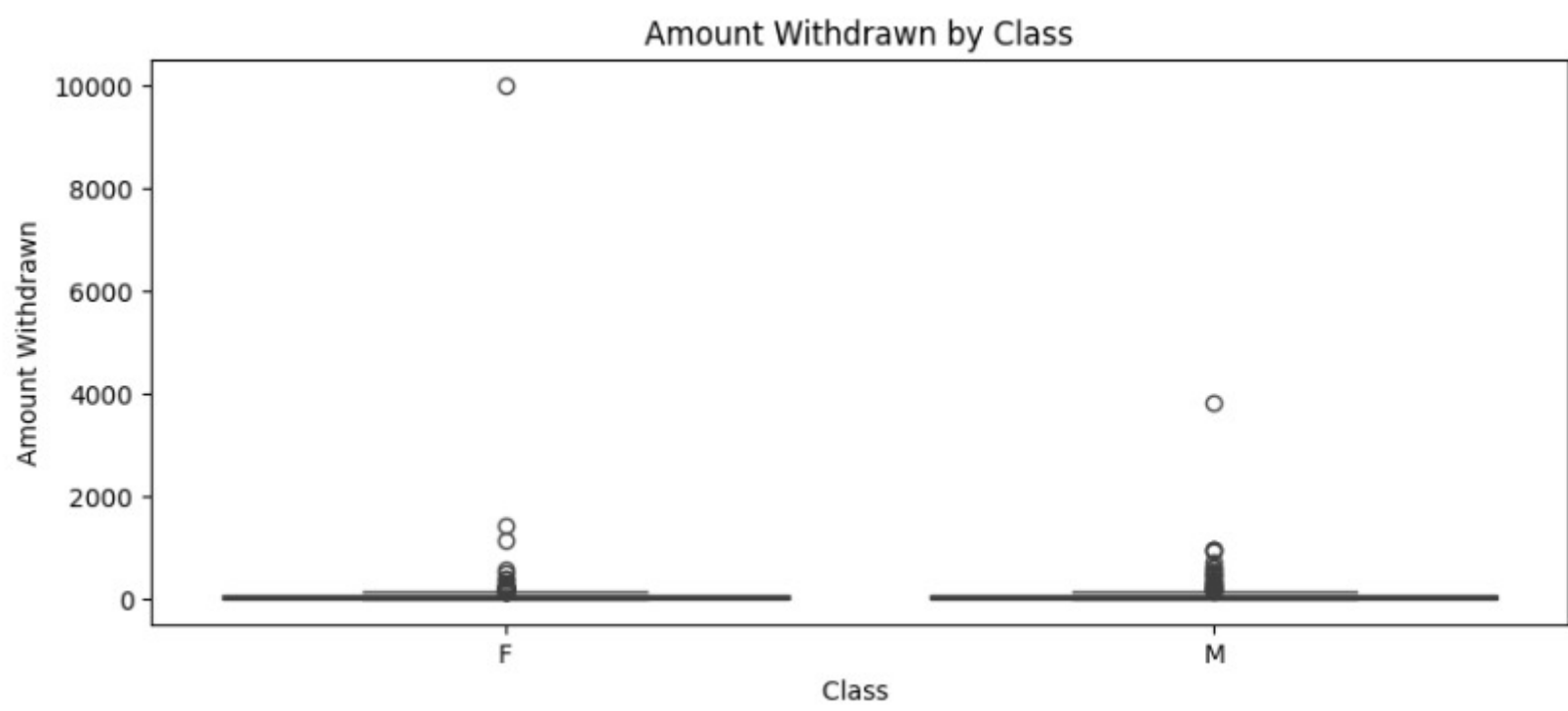


Figure 2: A boxplot showing the transactions by class

F - represents fraudulent transactions

M - represents legitimate transactions



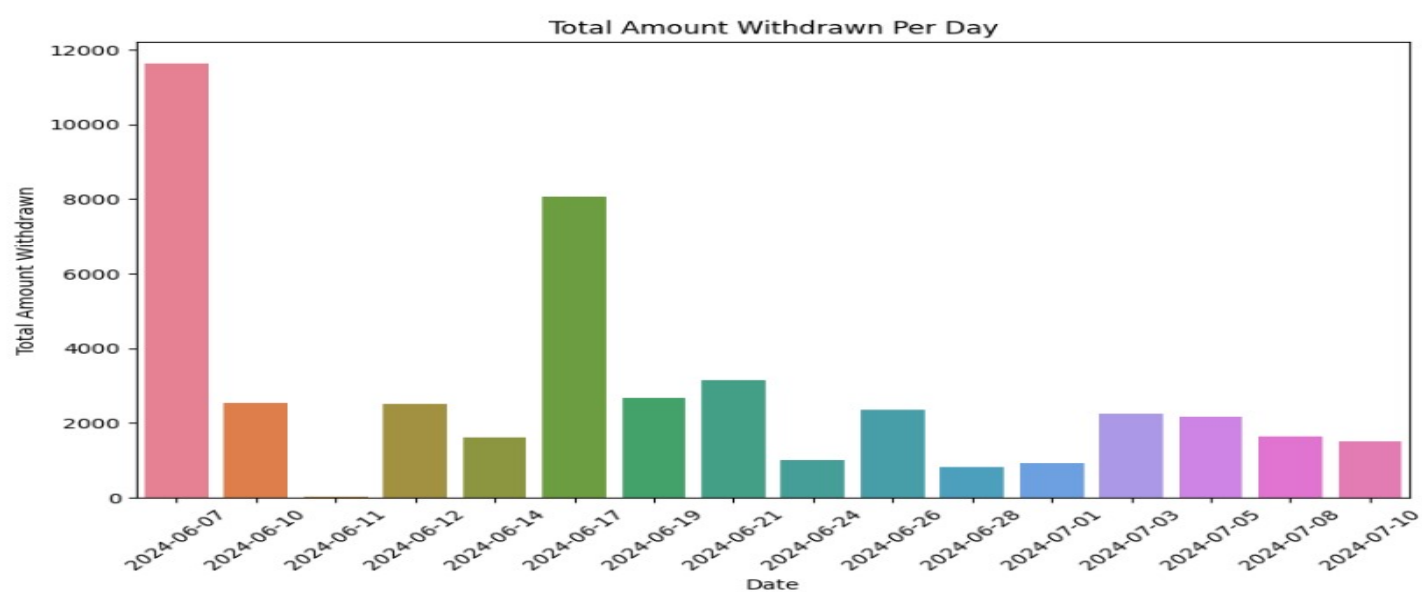


Figure 3: A line graph showing Total Amount Withdrawn Per Day

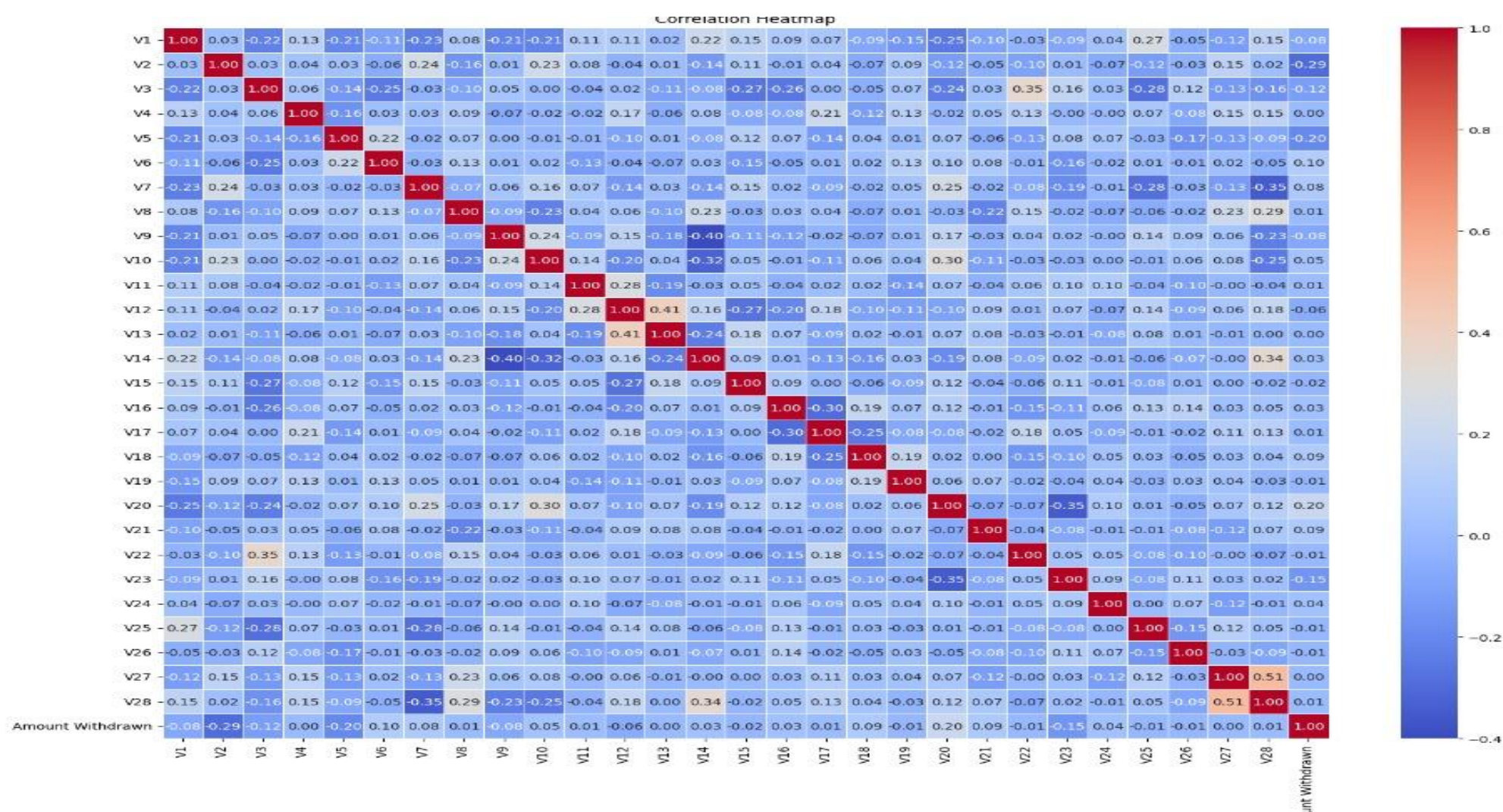


Figure 4: Correlation matrix

## Task 9:

The first Box plot revealed the amount withdrawn by date and the outliers.

The second Box plot revealed the amount withdrawn by each class.

A heatmap was created to visualize the correlation between various features in the dataset. This type of plot is particularly useful in identifying patterns and relationships between variables. In our dataset, the heatmap highlights the strength and direction of correlations among the numerical features ('V1' to 'V28') and 'Amount Withdrawn'. Darker colors indicate stronger positive or negative correlations, while lighter colors signify weaker correlations. This visualization aids in pinpointing features that have significant interrelationships, which can be pivotal in feature selection and engineering for predictive modeling. For instance, if certain features correlate highly with 'Amount Withdrawn', they might be crucial in determining the transaction's class.

## MACHINE LEARNING

### Task 10:

We split the dataset into training and testing sets to train our machine learning models. This step ensured that we could evaluate the models' performance on unseen data and avoid overfitting.



We also used Log Transformation on our Variables that were numerical because most of them were negatives, so we wanted to normalize the data before training our models on the data.

We trained different models including Random Forest, KNN, Support Vector Machine (SVM) and the Logistic regression model.

## Task 11:

We performed cross-validation to ensure our models' robustness and reliability. Cross-validation helped us evaluate the models' performance more accurately and identify the best-performing model.

The accuracy of our models was very low. The best model was the Random Forest that gave us the accuracy 0.522.

## SCREENSHOTS OF THE DATASET BEFORE CLEANING IT

Figure 5: This shows that the dataset had 510 rows before cleaning it

Figure 6: This shows how many missing values before cleaning it



## SCREENSHOTS OF THE DATASET AFTER CLEANING IT

	A	B	C	D	E	F	G	H	I	J	K	L	M
464	2024-07-08 12:20:11	1.27013608973005	-0.608139036121013	0.373101018688727	0.23375159974057	-0.721246758647321	0.396264693381749	-0.841541043223044	0.216319920355597	-0.734176366892903	1.03711082819275	-0.271628357355387	-0.622184300652378
465	2024-07-08 12:20:37	1.19549440194806	0.194929191243257	0.617510331478634	0.649717014766923	-0.474718219338476	-0.716084054893208	-0.0270776267494838	-0.0733847938675783	0.0572510527609135	-0.102522479593856	0.254488017068926	0.42798269864160
466	2024-07-08 12:20:40	-0.268621026179466	-0.233373932133886	1.02526288412106	-2.12932501874631	1.566375121693694	3.86306356459431	-0.732264663270428	0.807716476695413	-1.23005309133396	0.484128784110315	-0.0671846307387487	-1.2678523305118
467	2024-07-08 12:22:00	1.02166509245772	0.11057523002041	0.449479939310283	1.04142403576727	-0.213918613337462	-0.230678544777867	0.102705847458721	0.038497306399086	-0.319166006461889	0.0760029193491246	1.80514752412473	1.2558209365273
468	2024-07-08 12:38:17	0.779819724468341	-0.41580793896227	0.0881846163193077	1.34806553735487	-0.0664145468311553	0.41404476291984	0.246646687085592	0.0726392646089507	-0.0519651388631917	-0.0596743739274567	0.965016137989867	1.4105589465833
469	2024-07-08 12:40:17	1.0047604754428	0.423991944618589	0.579929072286713	2.47675356849986	-0.1243754773037013	-0.363935715693918	0.297698305770324	-0.104341836657113	-1.231239899685	0.812425561403615	1.49411135286794	1.1588988254529
470	2024-07-08 12:41:28	-3.4959837413741	-0.48841987249722	2.02484461194579	-0.74036273805294	-0.12813456738835	-1.23170169938276	-0.0865535583222919	0.157807359558066	1.67762074043078	-1.94685990082761	-0.741919716112826	0.46751341631713
471	2024-07-08 14:15:22	1.50757798120502	-1.09282023257786	0.360102102115976	-1.77093992289384	-0.913848623244812	0.678723287651672	-1.38022752331431	0.147689518439912	-1.97026082268998	1.51936867363724	0.577750072689378	0.15650027992757
472	2024-07-08 15:21:23	-0.928193196138154	0.67175801234988	1.92461595822292	0.273754662724616	1.04110848244839	0.1369301146374	0.36148391651264	0.143336967514427	-0.682134253112058	-0.650820411478348	-0.78952026239122	0.429949661383333
473	2024-07-10 10:16:08	1.0770785636955	0.284980216919471	0.0077313826582992	1.6570730005351	0.0520202406490099	0.446388674362298	-0.40703616617573	0.355703887753441	0.626039229744755	-0.92908713830021	1.09416510760764	0.579564773466065
474	2024-07-10 10:25:30	-1.53127111807147	1.39962086034665	-0.587061395411806	2.17500166489627	-2.13763992327444	-0.501576167574257	-1.215215050630479	0.956861641163217	-1.86656120022517	-2.31102388434087	2.77010079151076	-2.3610498498864
475	2024-07-10 11:14:57	-0.343986366495033	0.539788565141525	1.02466150982501	-0.371898822956345	1.03733561321496	1.26184919569894	0.473408900938114	0.227169255426177	0.248279161338581	-0.0862360584623006	-0.610694671521748	-0.041801180871005
476	2024-07-10 11:14:57	1.20544382671016	0.0084669857946875	0.953782320291429	1.14109301924441	-0.491215002839821	0.29730260419984	-0.503912857700374	0.0849482227756195	0.7964907027829493	-0.231653213875714	-1.36784403746149	0.74158089675758
477	2024-07-10 11:14:59	-0.703183200576619	1.21070441553077	0.713730881966961	1.14638078659484	0.196790068219018	0.468059751421935	0.20497737431824	0.62738883466996	-0.296261347312421	-0.309000901731636	-1.64599936511439	-0.35198285136093
478	2024-07-10 11:15:04	-1.04636191482801	0.720386219291532	1.60341258950027	0.608371108005076	0.460940564666008	-1.42340578494226	0.266078669898921	-0.141047826534776	-0.516792825520311	-0.20134654385537	0.21788409240458	0.2563895642171
479	2024-07-10 11:15:36	1.13584591864479	0.0841635452517641	0.269572728032023	1.21003437138404	0.0051279186390276	0.318050368299662	-0.082663815932095	0.19104874781851	0.157572985455954	0.0628041913787567	0.750879879966379	0.74339438494674
480	2024-07-10 11:15:54	0.150998913806338	-3.00211987419831	0.82430140640899	0.231721086027853	-2.62141526328587	0.128843083189467	-0.755233268800306	0.0959555802805319	0.638218705992502	0.0964622967819154	-1.22141951665696	-0.99340134111940
481	2024-07-10 11:15:58	-1.0894562543774	1.1971033223293	1.19731905139386	-1.17645349295576	-0.859954614304396	-1.47607543355181	0.30825431981024	0.19645417188019	0.542924753867272	-0.0324061967960648	-0.125285297156257	-0.14939513272302
482	2024-07-10 11:16:06	-0.687099302570821	0.790436270915636	2.24242432271366	2.406461630166905	0.359711754795171	0.279410298879838	0.10387751511258	0.39097184305469	-1.5187327616647	0.601268282242301	1.38977445144218	-0.1918927274938
483	2024-07-10 11:16:25	1.13562921316076	-0.17398536334862	0.730692302193073	0.711558021021266	-0.854208640751099	-0.849451355838813	-0.296619621489504	0.110503895931953	0.490269050061953	-0.0199644106303759	0.995819272110269	0.47583306199118
484	2024-07-10 11:16:30	0.019255815182034	0.47169811150898	0.869060279000251	-1.23912416338942	0.547897879724197	-0.50266107323217	0.808516316378471	-0.247667131076734	-0.130741129139825	-0.484993014601417	-0.70209446499769	-0.2134535927468
485	2024-07-10 11:16:31	1.08584018898175	0.126685010531901	0.498971594980159	1.3336995484128	-0.199128816165415	0.0326444603778757	0.0155736047636188	0.116090544559626	0.01599054002301	0.0386326141678344	1.3439204563189	1.2856096198821
486	2024-07-10 11:16:35	-1.24412629826194	1.10587092140905	2.70899427742413	3.0886201035228	-2.50579946457	0.992462333143792	-0.33746537947511	0.454982765029852	-0.205521576228304	1.21625494597924	0.57155167112302	0.57817854635585
487	2024-07-10 11:17:24	1.02139247472103	0.155450971619841	-0.312290950683655	0.961407335567952	1.046613111634237	0.121816269636434	0.525318014629569	-0.0538912610804885	-0.68686735338448	0.113058507536762	1.15862293590572	1.3064753349613
488	2024-07-10 11:17:28	1.5860925216242	-1.1690907377761	-1.35047670772023	-2.50579946457	0.508613111634237	0.313528169106679	-0.121769311985553	0.150875510370194	-0.0673070927524191	-0.208956269359136	-0.0218758187705562	-0.1061651330263
489	2024-07-10 11:17:34	-0.711928463741055	1.54742310558414	1.79975948782012	2.48905965234492	0.450496981199394	1.0737619805471	0.309374239228469	0.450853944132975	-1.56010692281051	1.23135612903401	1.7384492595567	0.56828285477376
490	2024-07-10 11:17:35	1.10625081303158	0.398625211642428	0.860420675089462	2.3888621035228	-0.365842960240857	-0.98934200478555	0.0654310855191308	-0.347445945013518	-0.109677745541789	-0.258586842026223	-0.32843275927937	0.3154853232327
491	2024-07-10 11:18:22	0.245656507097005	0.552493424117749	1.37423411846402	-0.447451226370809	0.2882475717152586	-1.0825359059151	0.969034687214964	-0.347445945013518	-0.109677745541789	-0.258586842026223	-0.32843275927937	0.3154853232327
492	2024-07-10 11:18:38	1.12619730203088	0.26345955260126	0.513253570696435	0.619053786672134	-0.472481831799073	0.98934200478555	0.0654310855191308	-0.347445945013518	-0.109677745541789	-0.258586842026223	-0.32843275927937	0.3154853232327
493	2024-07-10 11:18:40	1.29540611366961	0.296881362896136	0.068724824891864	0.572690232982924	-0.165698347876104	-0.781113471527394	0.0309534469772342	-0.133051027481164	0.165738866280832	-0.262785407465561	-0.655279963354797	-0.2806822320705
494	2024-07-10 11:19:30	-0.692698818107399	0.291548412085799	1.57522762626918	-1.21940041216185	-0.302254759513764	0.670790163275082	0.333059412564653	0.153874719975505	0.13854123068295	-1.47388205637055	-0.20219626758349	-0.3121875818770556
495	2024-07-10 11:19:32	-1.0286989631437	0.910514553316548	1.915180040516918	2.46938359686018	-0.0083750819742973	0.59758412370364	0.251531228603661	-0.331790263946744	-0.0956393317927111	1.34874966989555	-1.1261263894512	0.056429412345197
496	2024-07-10 11:19:38	-0.465895658702141	0.541707557441951	1.33129724259932	0.559447257178546	0.1840479048022	0.612792452974851	0.632758584309442	-0.15225961755782	-0.20219626758349	-0.510698499529772	-1.36392912418313	0.40718404363094
497	2024-07-10 11:20:00	-1.20943716011249	0.949446219081882	-0.429811066470978	1.63049352094142	2.60136860803058	0.35695653572267	0.4734941826269552	0.280728185601077	0.85304392864496	1.02389272369477	0.025354554381803	-0.38603422393421
498	2024-07-10 11:20:21	-1.1856978458655	-0.386597890901968	1.62307213323478	1.636405061013473	1.10250551000118	0.63316716784857	-0.303374045645056	0.28652538078265	-0.0560846500155522	-0.17392110656648	-0.98935953178869	-0.86621240874106
499	2024-07-10 11:20:23	1.0203989643983	-0.639478785995106	0.941567785978575	0.1063706684118969	-0.99631980555391	-0.33733376377325	-0.589656167545766	0.153583025302163	0.835978621256444	-0.347747200973645	0.749884414669471	1.5879608627510
500	2024-07-10 11:20:23	1.2562165510808	0.255027771814694	0.277716368985945	0.706677006371955	-0.309710062896315	-1.10763549562135	-0.0084930681798181	-0.043206426218129	0.24368732055133	-0.241481176592501	0.307360814471006	-0.50616317432414
501	2024-07-10 11:21:44	1.23917323388527	0.10336791843768	0.371710701187315	0.44003317388674	-0.37039355629073	-0.509383233340456	-0.129165339372591	0.030575679024732	-0.0105776510794315	0.3877803678566	0.9859230490693	0.3118778438921
502	2024-07-10 11:22:21	-0.40989592587	1.18308756526926	1.59896704783592	0.35308835423035	0.309710062896315	-0.31299991678796	0.707197130105536	-0.043206426218129	-0.08286898747464	0.648800184134323	0.23710117639282	0.1471815711851
503	2024-07-10 11:24:01	1.24060011229358	0.747735015779921	-0.214136192168044	1.10147536316146	0.124739708766457	-1.1721834169708	0.37105609831864	-0.306401437194445	-0.246672199166961	-0.597783921246271	0.227746612970307	0.3538686226900
504	2024-07-10 11:24:45	0.95391820343811	-0.760595477882238	1.09161059514784	1.30171485688749	1.43014758684862	-1.43014758684862	1.0704695818968	0.727964986419444	1.4327338699831525	-0.602645979204797	0.15224513416892	0.10998315715899
505	2024-07-10 11:31:21	1.25572940817485	0.297650061970427	0.28752640927889	0.6999020								



- Feature Engineering: New features significantly improved model performance.
- Model Performance: KNN outperformed other models in both accuracy and cross-validation scores.

#### ACTIONABLE INSIGHTS

- Real-Time Monitoring: Implement real-time monitoring systems that track transaction frequency and other derived features. Alerts can be triggered when unusual patterns, such as an unexpectedly high number of transactions in a short period, are detected.
- Model Improvement: Enhance existing fraud detection models by incorporating features derived from transaction patterns, such as frequency, time gaps between transactions, and transaction amounts. These features can significantly improve model accuracy and reduce false positives.
- Threshold Setting: Establish threshold values for critical features like transaction frequency, which can be adjusted dynamically based on the user's typical behavior profile. This adaptive approach helps in minimizing both false positives and false negatives in fraud detection..

#### RECOMMENDATIONS

- Regular Feature Update: Continuously update and refine the set of features used in fraud detection models. As fraudulent techniques evolve, new patterns and anomalies may emerge, necessitating ongoing feature engineering efforts.
- Integration with Other Data Sources: Consider integrating additional data sources, such as geolocation and device information, to enrich the feature set. This can provide a more comprehensive view of the transaction context and improve detection accuracy.