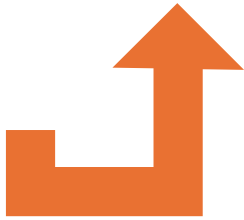


# 1.1 Dataset –Definition of ML Task



Predicting whether a newly IPO'd company will be delisted within the next five quarters.



Initial Public Offering (IPO): Process in which a privately owned company lists its shares on a stock exchange to raise capitals from investors







Binary Classification

# 1.2 Dataset – Use of Different File

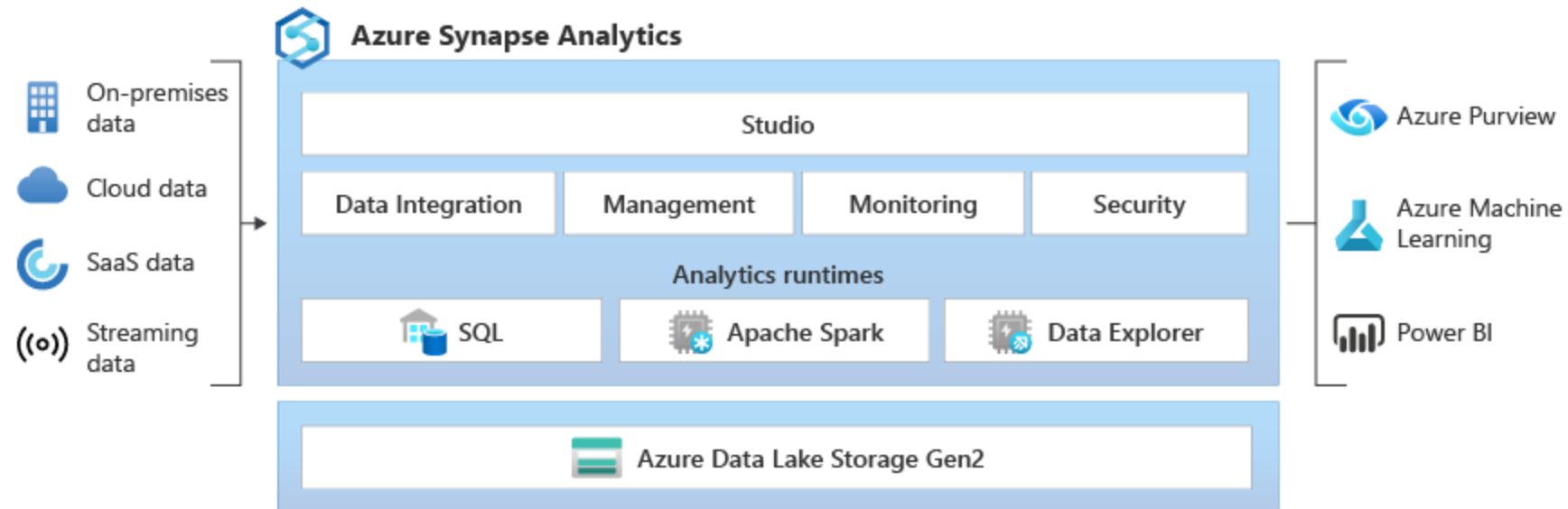
File Name	Source	Number of Rows	Number of Columns	File Size (MB)	Main Entities	Usage
DS1-IPODataFull.csv	Kaggle	3673	1669	28.8MB	Name, LastSale, Market Cap, Sector, Industry, Summary Quote	Not Used
DS2-CompaniesRanked.csv	Companies MarketCap	8091	6	474KB	Rank, Name, Symbol, Employees Count, Price (USD), Country	Not Used
DS3-company_ipo.csv	Kaggle	1766	7	109KB	ID, IPO Date, Symbol, Company Name, IPO Price, Current Return	Not Used
DS4-final_merged_ipos.csv	StockAnalysis	2307	19	185KB	Company Name, Country, Current Deal Size, Employees, Exchange, Founded, IPO Date, IPO Price, Industry, Is SPAC, Market Cap, Open Price, Return, Return(Open), Sector, Shares Offered, Symbol, Volume	Used
DS5-delisted_companies.csv	AlphaVantage	7854	7	637KB	Symbol Name, Exchange, Asset Type, IPO Date, Delisting Date, Status	Used

## 1.2 Dataset – Data sources

	Data source for listed companies	Data source for delisted companies
<i>Company names &amp; stock symbols</i>		
<i>Financial Data</i>		

# 1.3 Dataset – Use of Storage Systems

A mix of everything



# 1.3 Dataset – Use of Storage Systems

[+](#) Create [⚙️](#) Manage view [✖](#) [🗑️](#) Delete resource group [🔄](#) Refresh [⬇️](#) Export to CSV [🔗](#) Open query | [🏷️](#) Assign tags [➡️](#) Move [✖](#) [🗑️](#) Delete [⬇️](#) Export template [📱](#) Open in mobile

## ^ Essentials

[JSON View](#)

Subscription ([move](#)) : [Azure subscription 1](#)

Deployments : [3 Succeeded](#)

Subscription ID : c01f036f-8552-4c98-9ba8-6c766f97492d

Location : UK South

Tags ([edit](#)) : [Add tags](#)

## Resources

[Recommendations](#)

Filter for any field...

Type equals **all** ✖

Location equals **all** ✖

[+🔍](#) Add filter

Showing 1 to 3 of 3 records. ☐ Show hidden types ⓘ


No grouping ▼

☰ List view ▼

☐ Name ↑↓

Type ↑↓

Location ↑↓

☐  comp333-workspace

Synapse workspace

UK South

...

☐  comp333Spark (comp333-workspace/comp333Spark)

Apache Spark pool

UK South

...

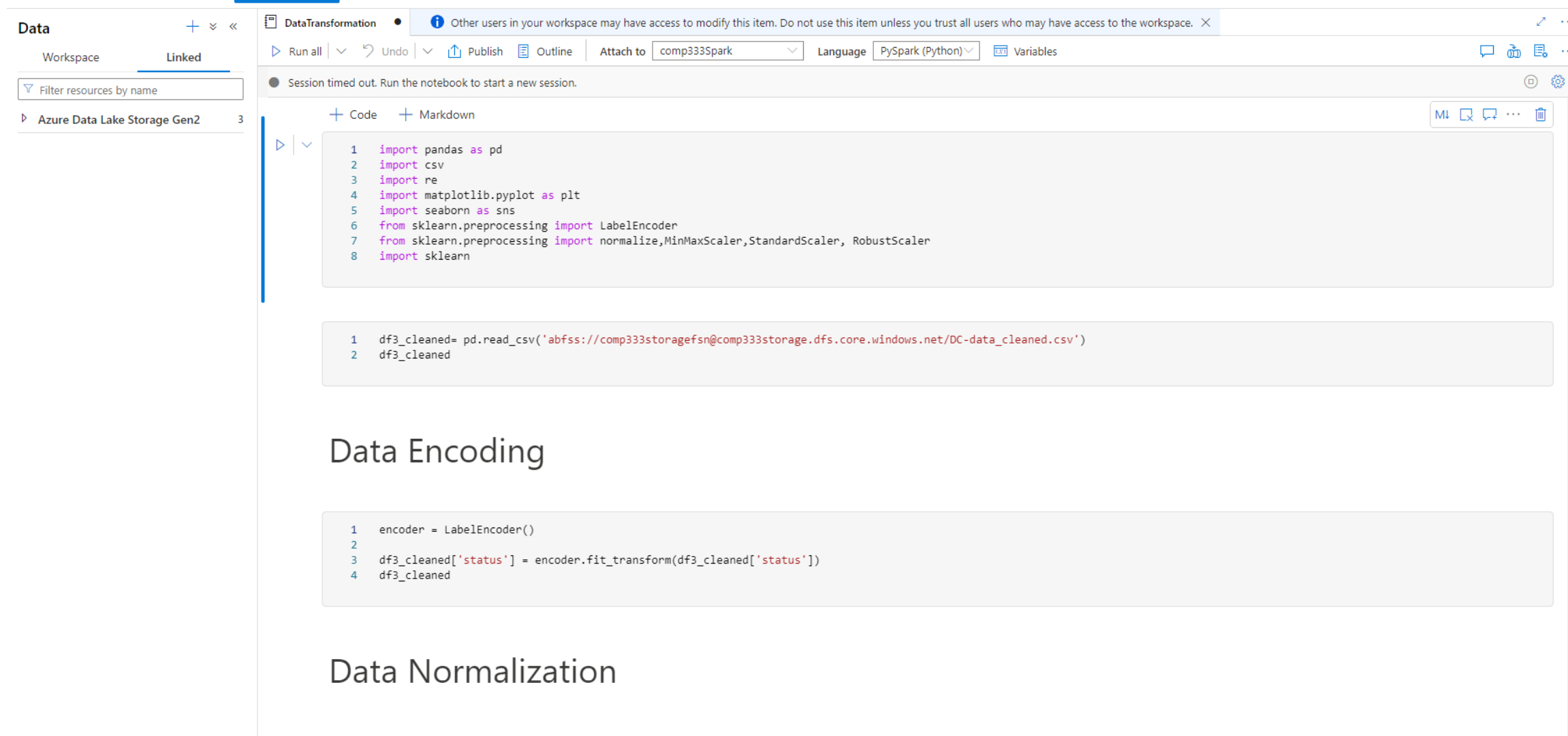
☐  comp333storage

Storage account

UK South

...

# 1.3 Dataset – Example with Data Transformation



The screenshot displays a Jupyter Notebook environment. On the left, a sidebar titled 'Data' shows a 'Workspace' tab and a 'Linked' tab. Under 'Linked', there is a search bar 'Filter resources by name' and a resource 'Azure Data Lake Storage Gen2' with a count of 3. The main area is titled 'DataTransformation' and includes a warning: 'Other users in your workspace may have access to modify this item. Do not use this item unless you trust all users who may have access to the workspace.' Below this, there are controls for 'Run all', 'Undo', 'Publish', 'Outline', 'Attach to' (set to 'comp333Spark'), 'Language' (set to 'PySpark (Python)'), and 'Variables'. A message states 'Session timed out. Run the notebook to start a new session.' The notebook content is divided into two sections: 'Code' and 'Markdown'. The 'Code' section contains two code blocks. The first code block imports libraries: pandas, csv, re, matplotlib.pyplot, seaborn, sklearn.preprocessing (LabelEncoder, normalize, MinMaxScaler, StandardScaler, RobustScaler), and sklearn. The second code block reads a CSV file from 'abfss://comp333storagefsn@comp333storage.dfs.core.windows.net/DC-data\_cleaned.csv' and assigns it to 'df3\_cleaned'. The 'Markdown' section contains two sections: 'Data Encoding' and 'Data Normalization'. The 'Data Encoding' section contains a code block that creates a 'LabelEncoder' object and applies it to the 'status' column of 'df3\_cleaned'.

```
1 import pandas as pd
2 import csv
3 import re
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from sklearn.preprocessing import LabelEncoder
7 from sklearn.preprocessing import normalize, MinMaxScaler, StandardScaler, RobustScaler
8 import sklearn
```

```
1 df3_cleaned= pd.read_csv('abfss://comp333storagefsn@comp333storage.dfs.core.windows.net/DC-data_cleaned.csv')
2 df3_cleaned
```

## Data Encoding

```
1 encoder = LabelEncoder()
2
3 df3_cleaned['status'] = encoder.fit_transform(df3_cleaned['status'])
4 df3_cleaned
```

## Data Normalization

# 1.3 Dataset – Data Lake Storage (Before)

Data

Workspace

Linked

Filter resources by name

Azure Data Lake Storage Gen2

3

comp333-workspace (Primary - co...

comp333storagefsn (Primary)

(Attached Containers)

AzureDataLakeStorage1 (comp333...

comp333storagefsn

Other users in your workspace may have access to modify this item. Do not use this item unless you trust all users who may have access to the workspace.

New SQL script

New data flow

New integration dataset

Upload

Download

New folder

Select all

Copy link

Rename

Manage access

More

comp333storagefsn

Name	Last Modified	Content Type
synapse	4/24/2024, 3:50:17 PM	Folder
DC-data_cleaned.csv	4/24/2024, 6:36:49 PM	

# 1.3 Dataset – Data Lake Storage (After)

The screenshot shows the Azure Data Lake Storage Gen2 interface. On the left, the 'Data' pane shows the 'Workspace' and 'Linked' tabs. Under 'Linked', there is a search bar 'Filter resources by name' and a list of resources: 'Azure Data Lake Storage Gen2' (3 items), 'comp333-workspace (Primary - co...', 'comp333storagefsn (Primary)', '(Attached Containers)', and 'AzureDataLakeStorage1 (comp333...'. The main pane shows the 'comp333storagefsn' container. At the top, there is a warning: 'Other users in your workspace may have access to modify this item. Do not use this item unless you trust all users who may have access to'. Below the warning, there are buttons for 'New SQL script', 'New notebook', 'New data flow', 'New integration dataset', 'Upload', 'Download', 'New folder', and 'Select all'. The main area displays a table of files and folders:

Name	Last Modified
synapse	4/24/2024, 3:50:17 PM
DC-data_cleaned.csv	4/24/2024, 6:36:49 PM
DT-minmax_scaling.csv	4/24/2024, 6:40:57 PM
DT-normalized.csv	4/24/2024, 6:40:49 PM
DT-robust_scaling.csv	4/24/2024, 6:40:53 PM
DT-standard_scaling.csv	4/24/2024, 6:40:55 PM

Outputted files from our DataTransformation code file can be seen directly in the data lake



DT-normalized.csv

Path

https://comp333storage.dfs.core.windows.net/comp333storagefsn/DT-normalized.csv

Modified

4/24/2024, 6:40:49 PM

With column header

☒ On

STATUS	SALESQ5	SALESQ4	SALESQ3
0	0.004744900661...	0.004260704379...	0.005315800...
0	0.454580712883...	0.408192746759...	0.50927520...
0	0.454581864047...	0.408193780452...	0.50927649...
0	0.004744900661...	0.004260704379...	0.005315800...
0	0.004744900661...	0.004260704379...	0.005315800...
0	0.004744900661...	0.004260704379...	0.005315800...
0	0.004744900661...	0.004260704379...	0.005315800...
0	0.004744900661...	0.004260704379...	0.005315800...
0	0.004744900661...	0.004260704379...	0.005315800...
0	0.004744900661...	0.004260704379...	0.005315800...
0	0.004744900661...	0.004260704379...	0.005315800...
0	0.058761507889...	0.072576909431...	0.08715983...
0	0.004744900661...	0.004260704379...	0.005315800...
0	0.004744900661...	0.004260704379...	0.005315800...

OK

DT-minmax\_scaling.csv

Path

https://comp333storage.dfs.core.windows.net/comp333storagefsn/DT-minmax\_scaling.csv

Modified

4/24/2024, 6:40:57 PM

With column header

☒

On

STATUS	SALESQ5	SALESQ4	SALESQ3
0	0.293188593162...	0.053360077968...	0.00532465...
0	0.612653510540...	0.440620108310...	0.50929682...
0	0.612654328075...	0.440621099338...	0.50929811...
0	0.293188593162...	0.053360077968...	0.00532465...
0	0.293188593162...	0.053360077968...	0.00532465...
0	0.293188593162...	0.053360077968...	0.00532465...
0	0.293188593162...	0.053360077968...	0.00532465...
0	0.293188593162...	0.053360077968...	0.00532465...
0	0.293188593162...	0.053360077968...	0.00532465...
0	0.293188593162...	0.053360077968...	0.00532465...
0	0.293188593162...	0.053360077968...	0.00532465...
0	0.331550169756...	0.118856579522...	0.08717076...
0	0.293188593162...	0.053360077968...	0.00532465...
0	0.293188593162...	0.053360077968...	0.00532465...

OK

DT-robust\_scaling.csv

Path

https://comp333storage.dfs.core.windows.net/comp333storagefsn/DT-robust\_scaling.csv

Modified

4/24/2024, 6:40:53 PM

With column header

☒

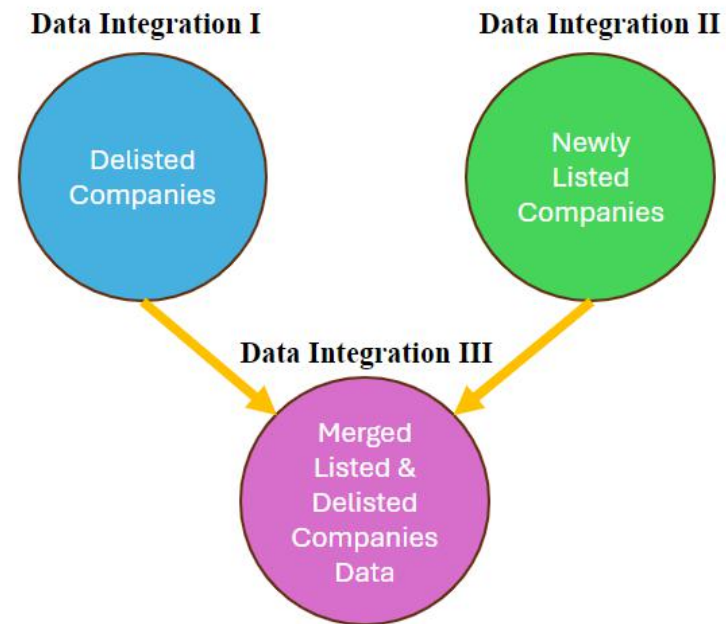
On

STATUS	SALESQ5	SALESQ4	SALESQ3
0	-0.06941170248...	-0.07726521722...	-0.05673455...
0	1.421105690206...	1.426997887601...	1.52292285...
0	1.421109504554...	1.427001737125...	1.52292689...
0	-0.06941170248...	-0.07726521722...	-0.05673455...
0	-0.06941170248...	-0.07726521722...	-0.05673455...
0	-0.06941170248...	-0.07726521722...	-0.05673455...
0	-0.06941170248...	-0.07726521722...	-0.05673455...
0	-0.06941170248...	-0.07726521722...	-0.05673455...
0	-0.06941170248...	-0.07726521722...	-0.05673455...
0	-0.06941170248...	-0.07726521722...	-0.05673455...
0	0.109570695939...	0.177147743223...	0.19980502...
0	-0.06941170248...	-0.07726521722...	-0.05673455...
0	-0.06941170248...	-0.07726521722...	-0.05673455...

OK

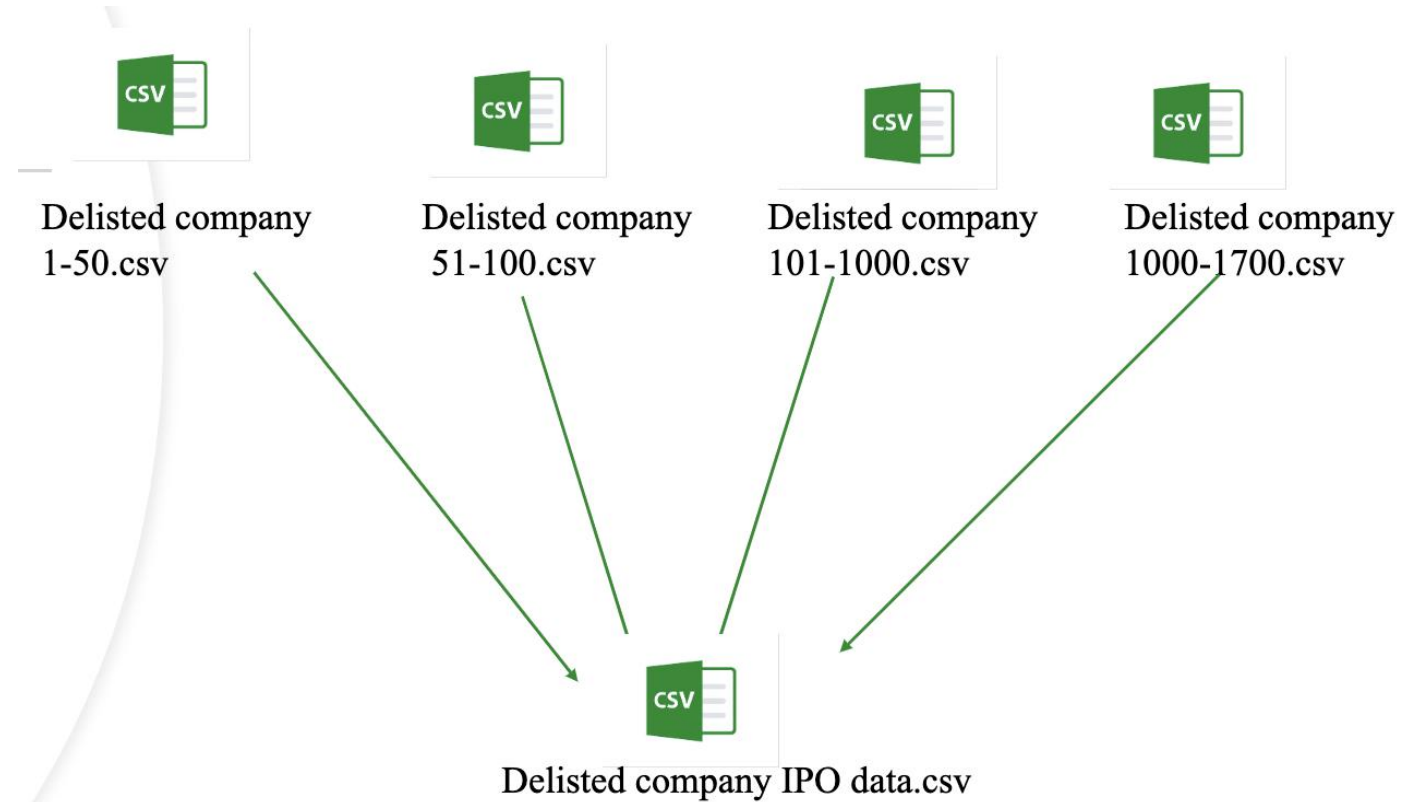
# 2. 1 Schema Integration- Overview

- Three main groups of Integration:



# Schema Integration 1 and Mapping: Delisted IPO Data

- Key: 'symbol'
- **Initial Aggregation:** Add monthly financial data from the first two files—'salesJan', 'salesFeb', 'salesMar' to calculate 'salesQ1'.
- **Subsequent Mapping:** Correspond 'salesQ1' from initial files with 'salesQ1' in subsequent datasets, where quarterly data already exists.
- **Uniform Naming Convention:** Consolidate figures under a standardized label 'XXQN' across all files for clarity and consistency in financial analysis.



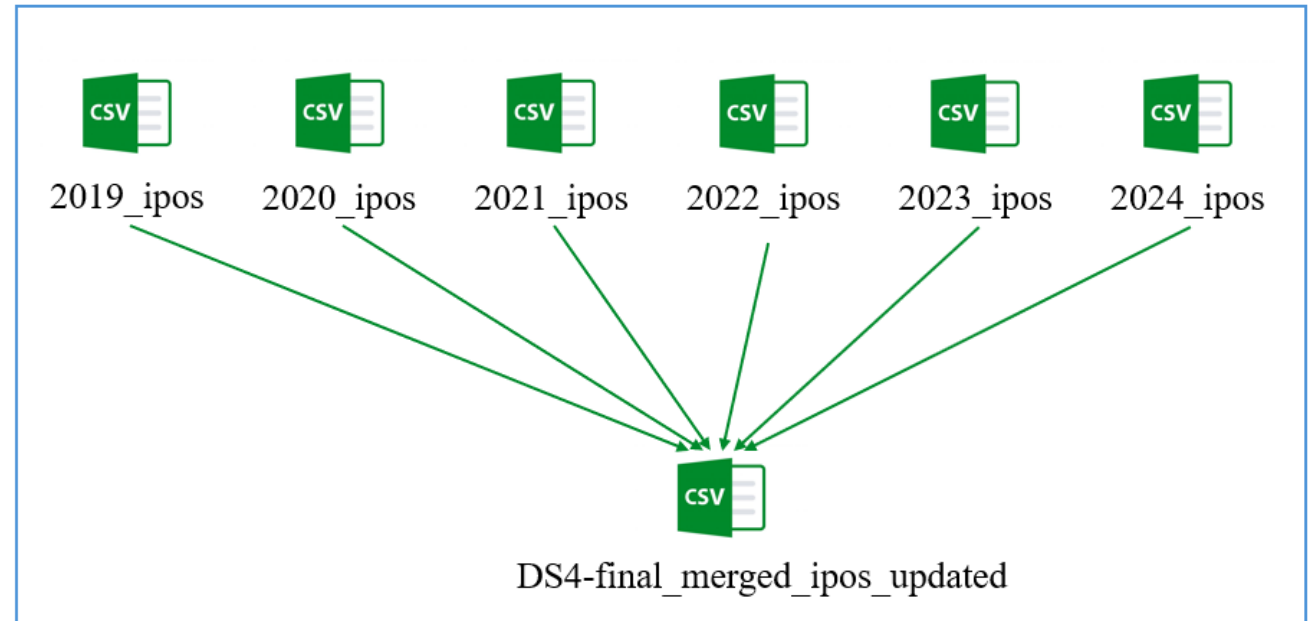
## Schema Integration and Mapping 2: Listed IPO Data

-Key: '**Symbol**'

-**Date Harmonization:** Align '**IPO Date**' from 2019 and 2020 files with '**ipo date**' in subsequent datasets.

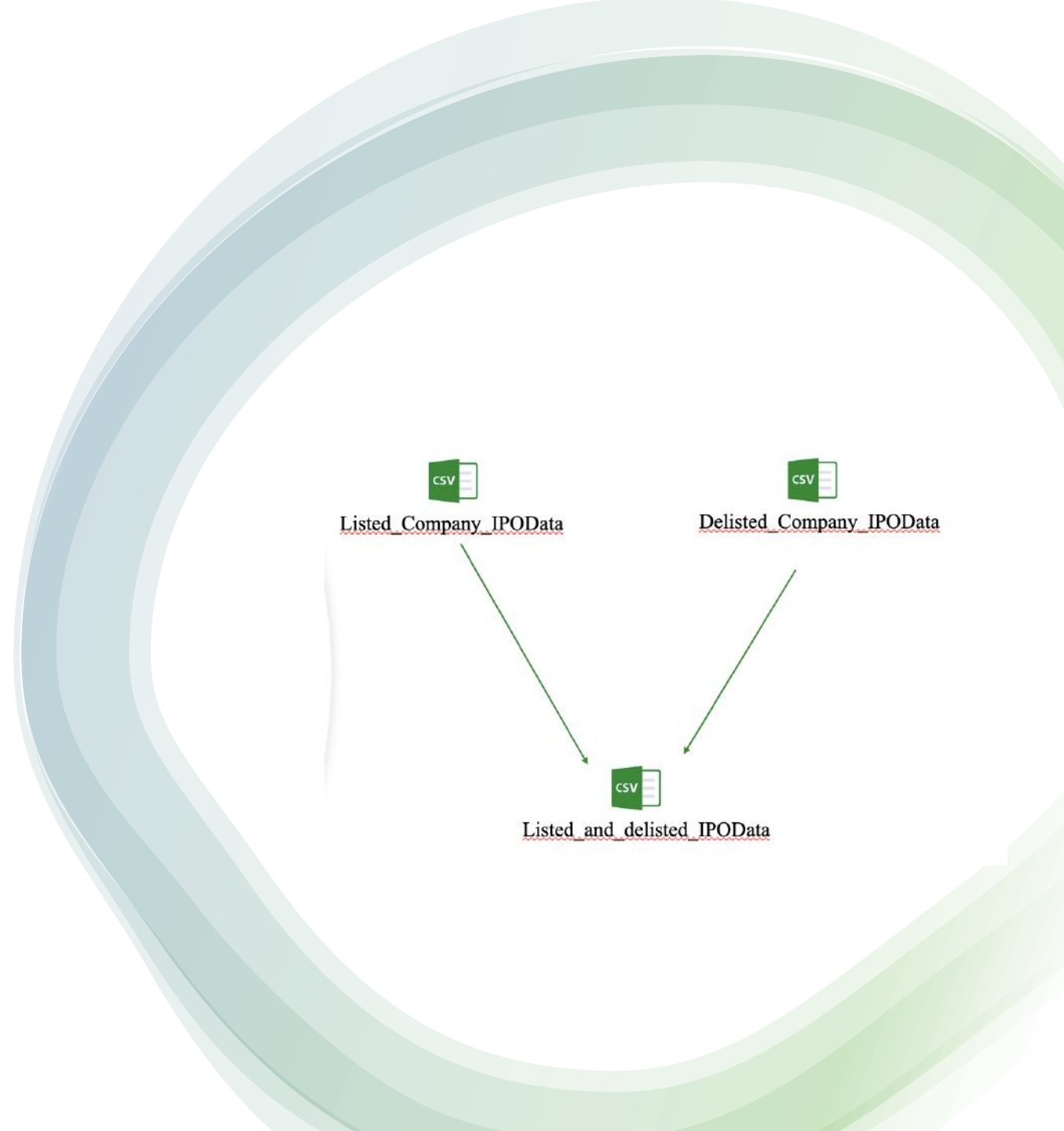
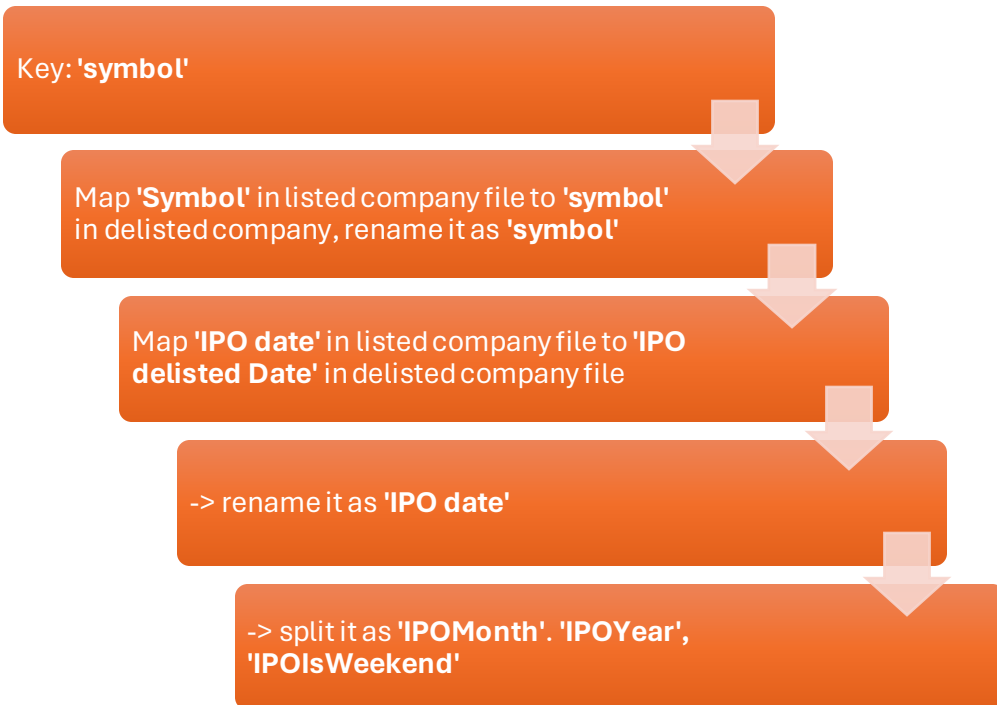
-**Renaming for Consistency:** Standardize field names by renaming all instances to '**IPO Date**'.

-**Format Unification:** Ensure date structure adheres to a consistent 'Year-Month-Date' format across all records.

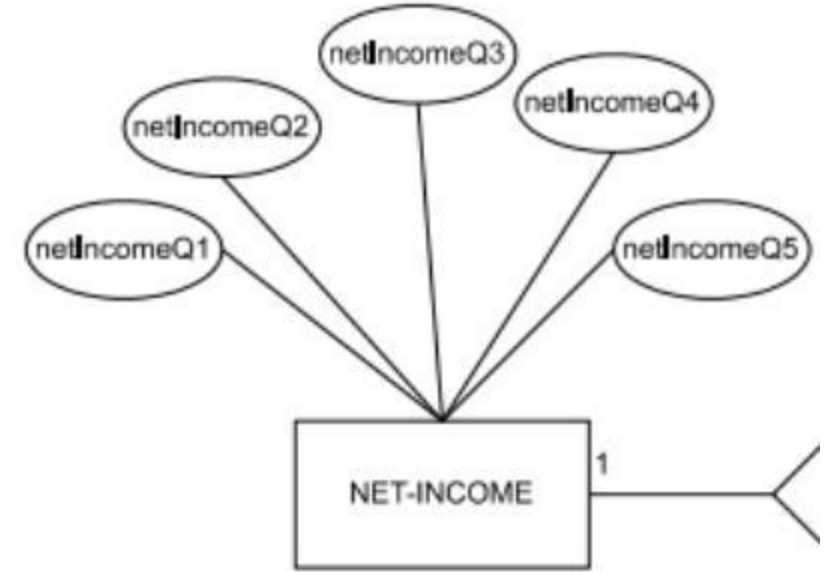
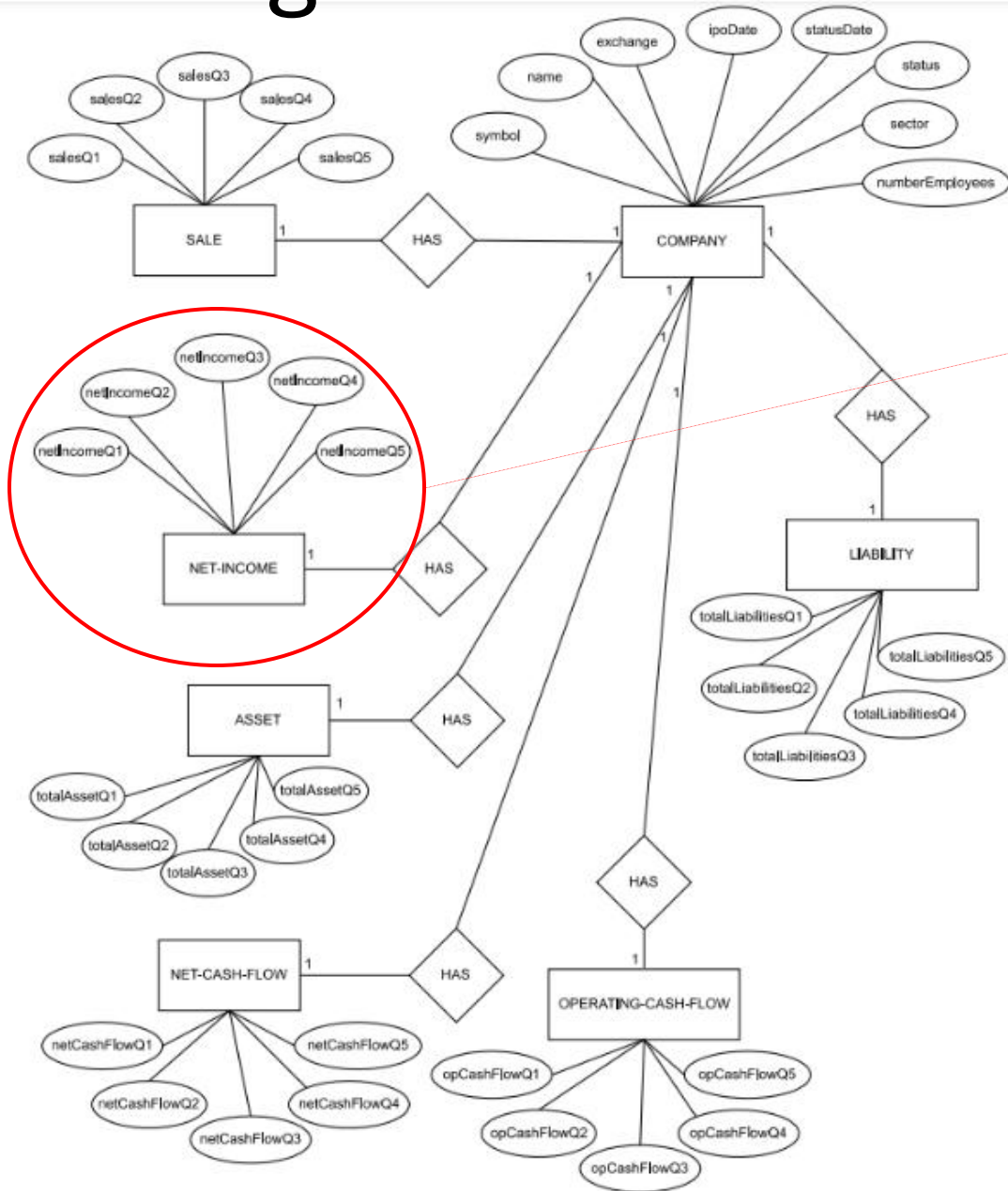


# Schema Integration 3 and Mapping :

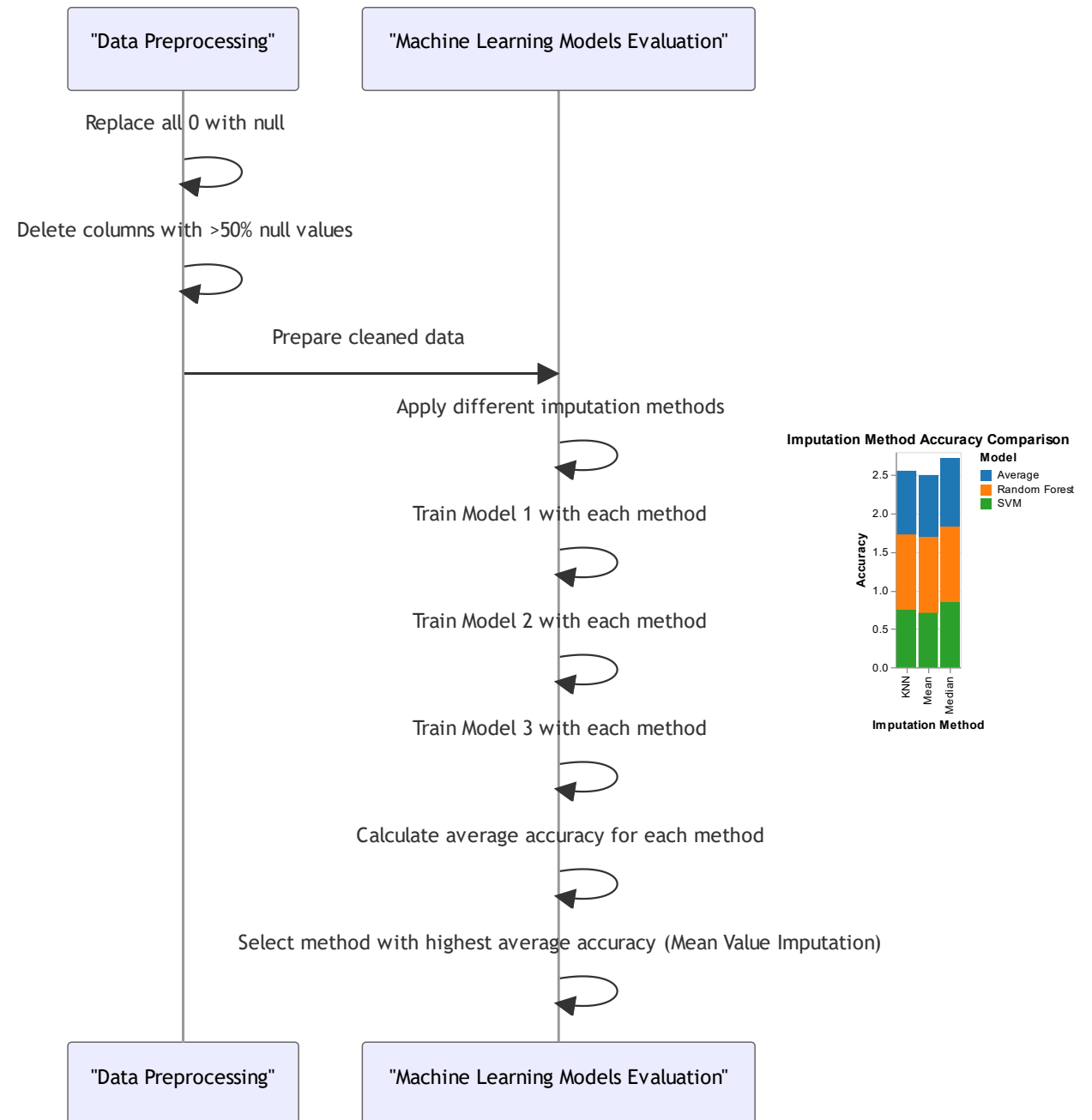
## Listed and Delisted IPO Data



# ER Diagram

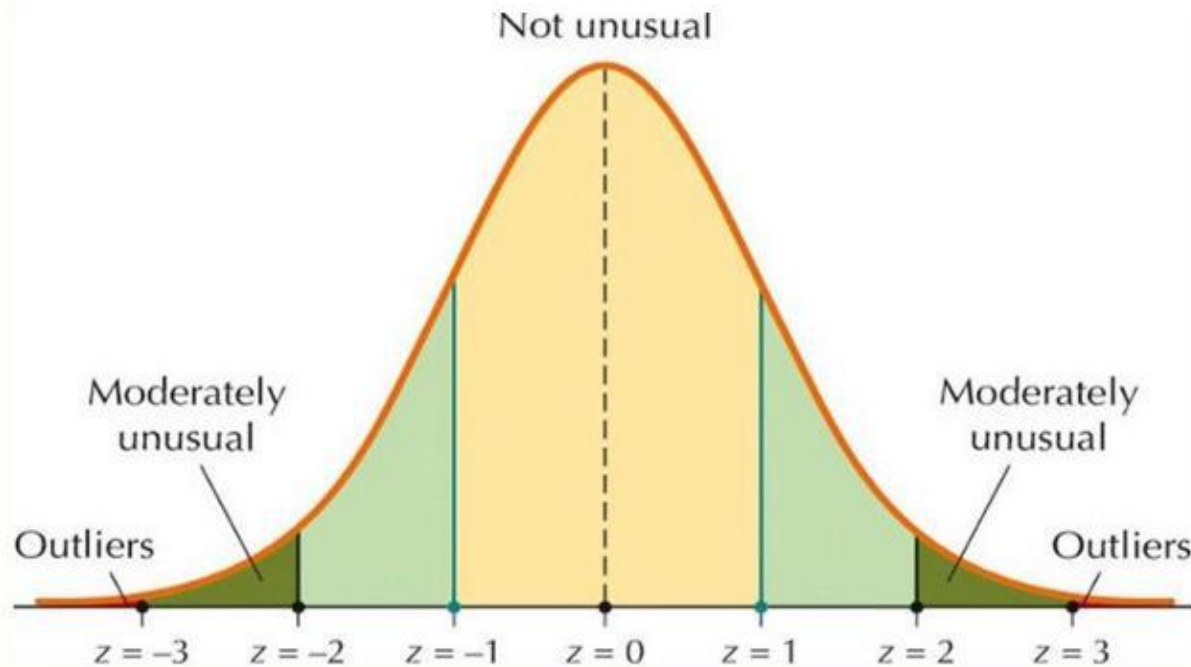


### 3. Data Cleaning: Support DC on Null values



### 3. Data Cleaning: Support DC for outliers

#### Detecting Outliers with z-Scores



- After visualize our data, we notice it follows a **Gaussian normal distribution**.
- data points with a Z-Score greater than a **threshold=3** are considered outliers



Model	Description	Accuracy Score
SVM	Before Data Cleaning	0.5969
SVM	After Data Cleaning	0.7286

### 3. Data Cleaning: ML model Improvement



# 4. Data Transformation

# Data Processing – Normalization

- API: `sklearn.preprocessing.Normalizer`
- This API preserves unit norm, since our ML model is SVM, it relies on the distance between support vectors and boundary that maximizes the margin

```
df_scaled.head()
```

	salesQ5	salesQ4	salesQ3	salesQ2	salesQ1	netIncomeQ5	netIncomeQ4	netIncomeQ3	netIncomeQ2	netIncomeQ1	...
0	0.004745	0.004261	0.005316	0.004351	0.004593	-0.000097	-8.839976e-05	-5.628035e-05	-1.165012e-04	-1.133295e-03	...
1	0.454581	0.408193	0.509275	0.416830	0.440023	0.000005	9.235170e-06	1.633772e-05	3.597637e-06	1.088563e-05	...
2	0.454582	0.408194	0.509276	0.416832	0.440024	0.000001	-1.854457e-07	-5.563370e-07	-2.781685e-07	-4.079804e-07	...
3	0.004745	0.004261	0.005316	0.004351	0.004593	-0.000097	-8.839976e-05	-5.628035e-05	-1.165012e-04	-1.133295e-03	...
4	0.004745	0.004261	0.005316	0.004351	0.004593	-0.000097	-8.839976e-05	-5.628035e-05	-1.165012e-04	-1.133295e-03	...

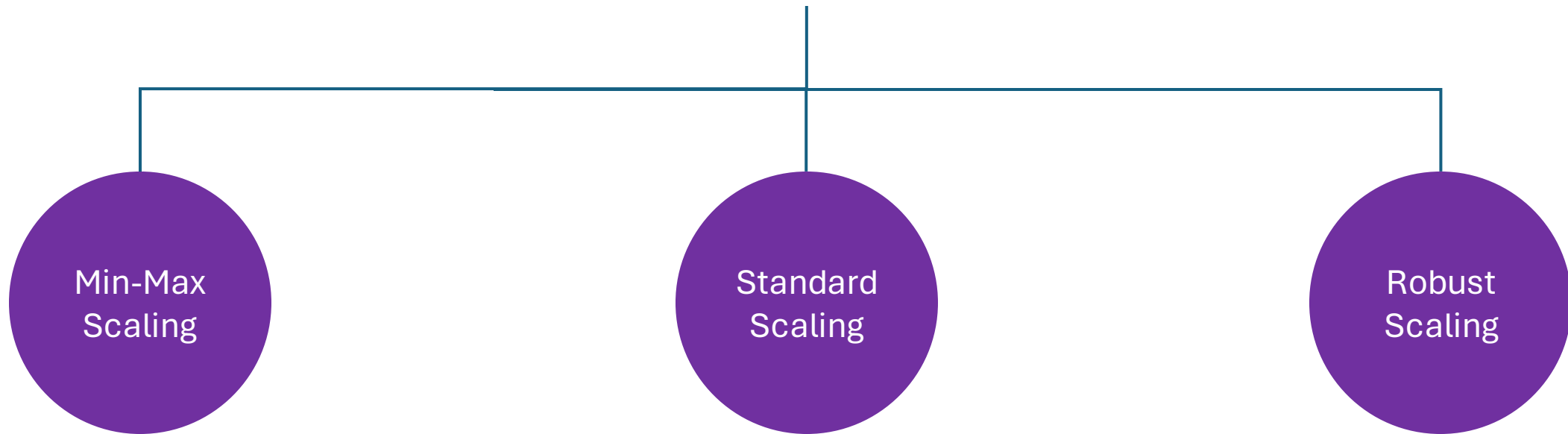
rows × 31 columns

# 4. Data Transformation – Applying different scaling functions after normalization

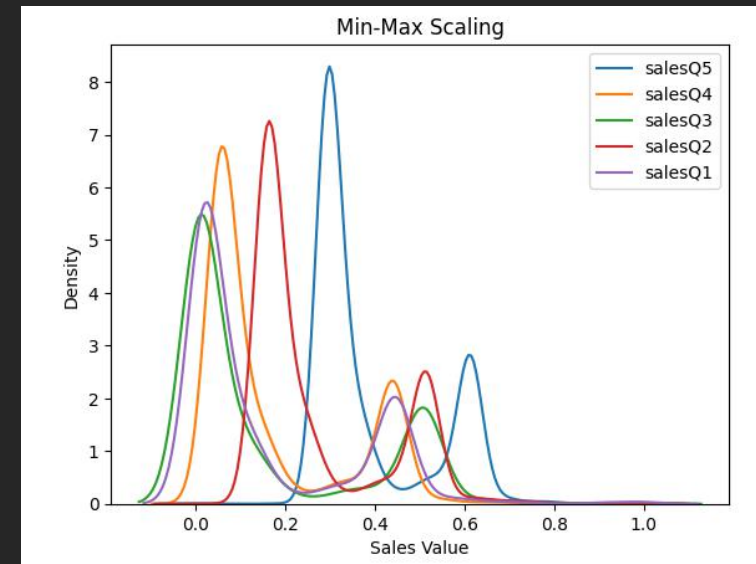
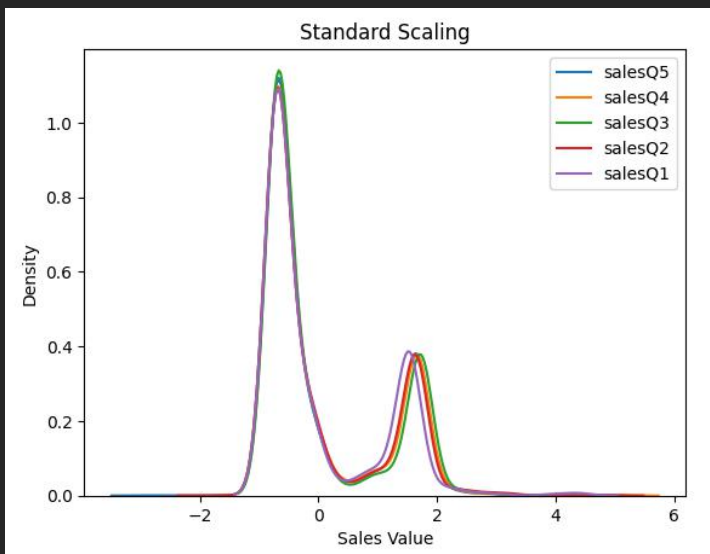
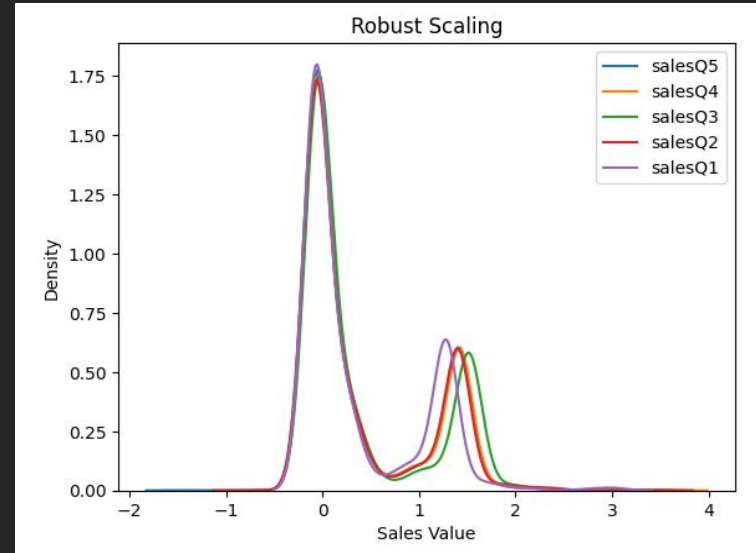
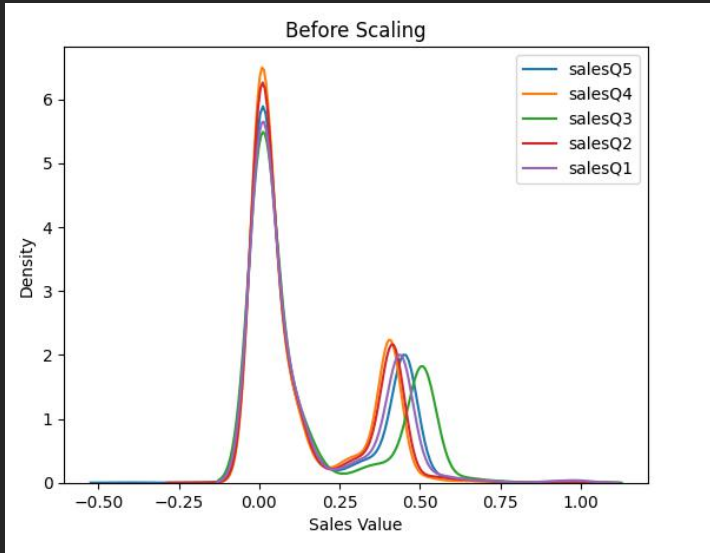
- We considered different scaling functions in our scaling process



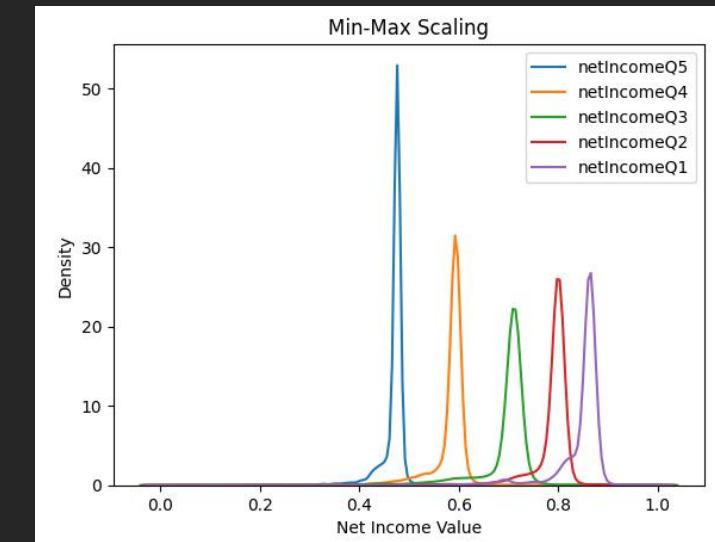
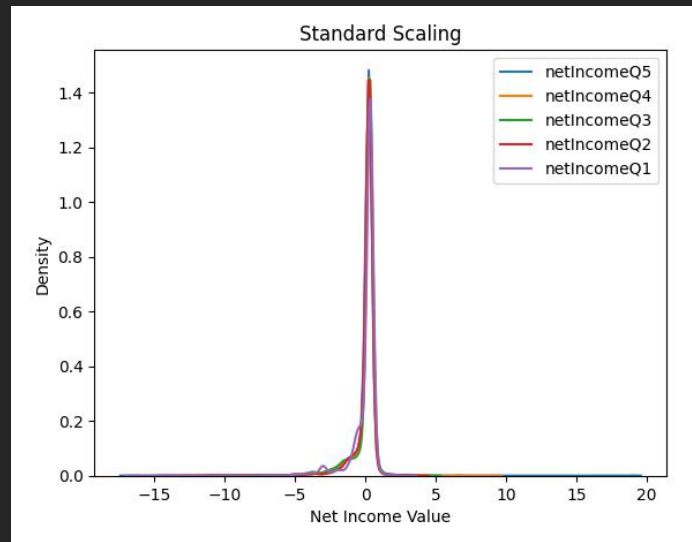
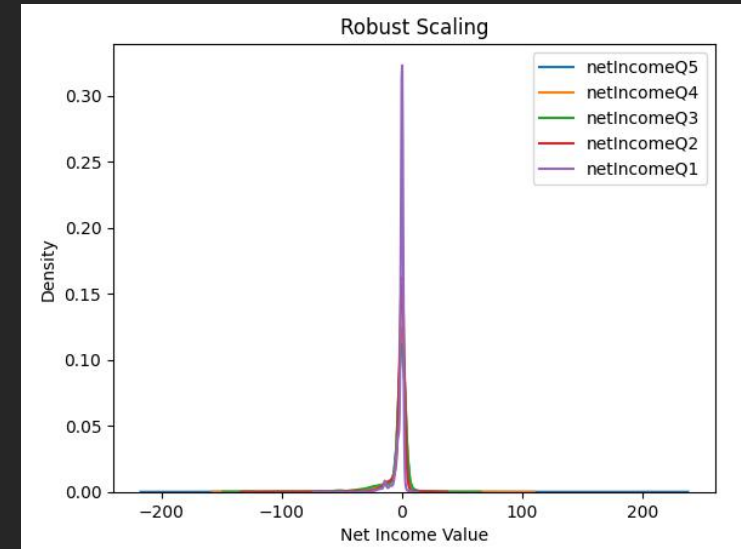
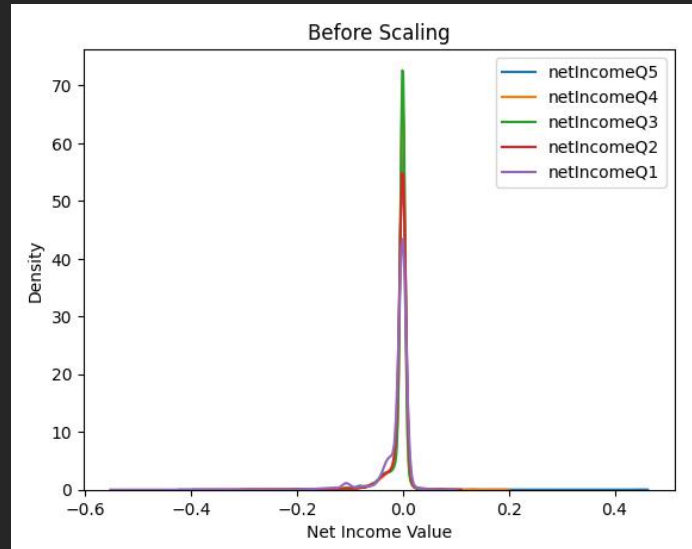
DT-normalized



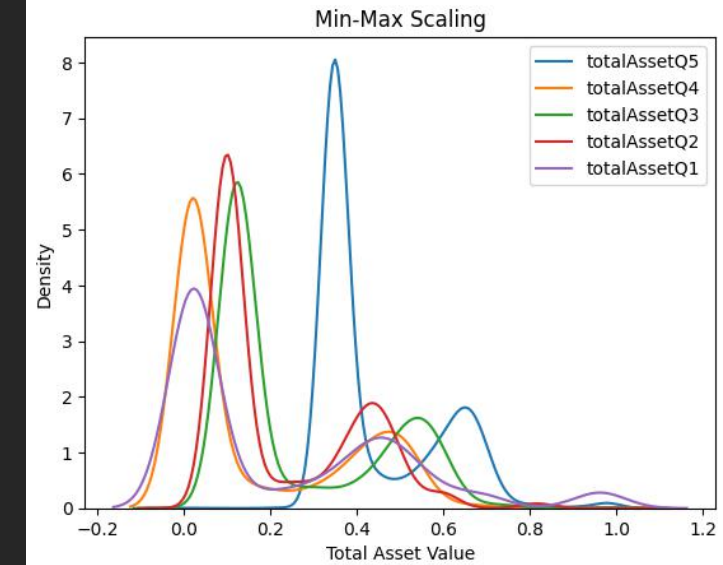
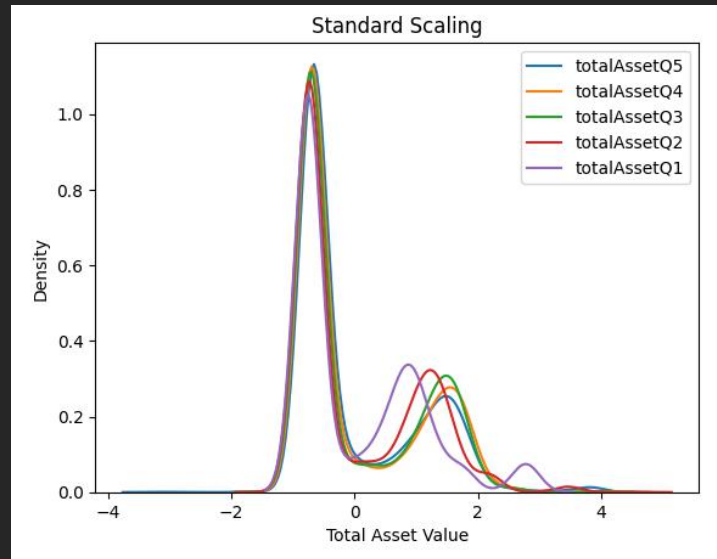
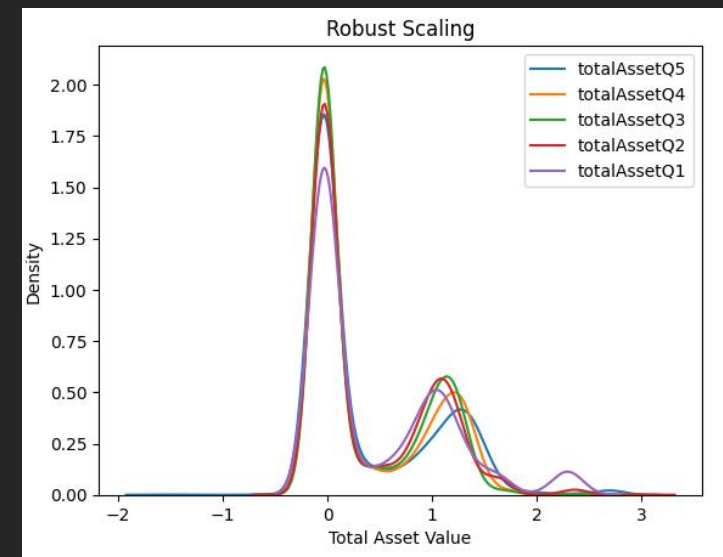
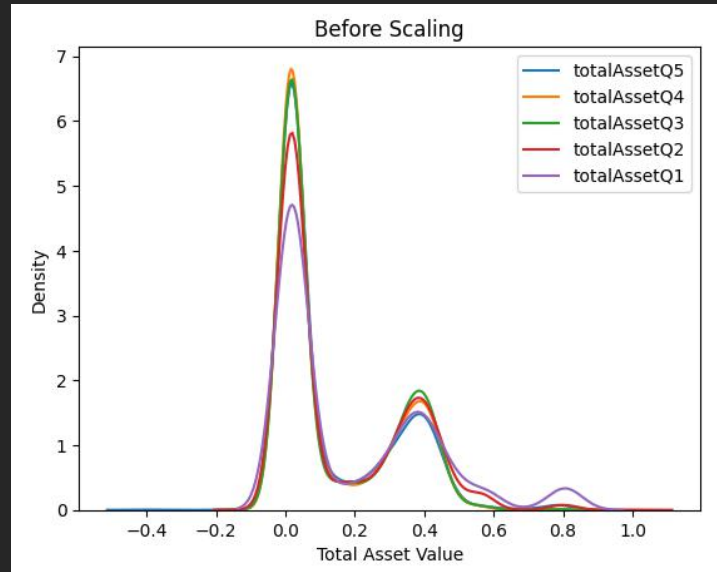
# Sales



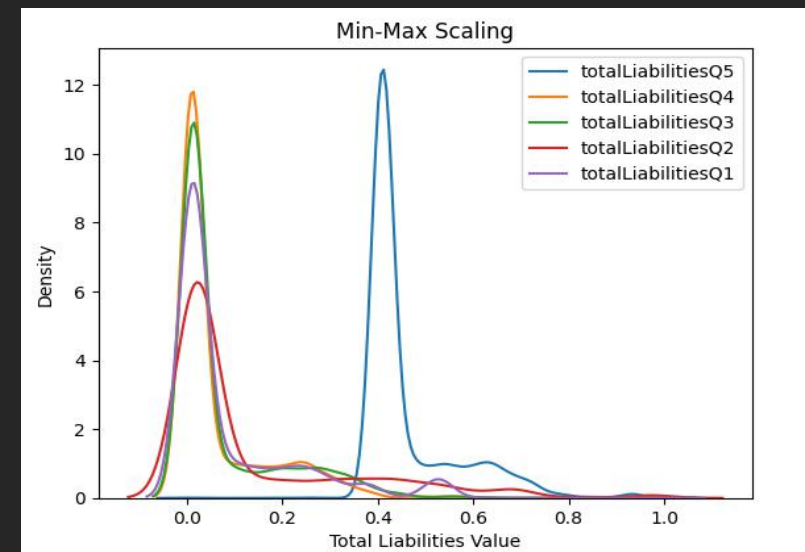
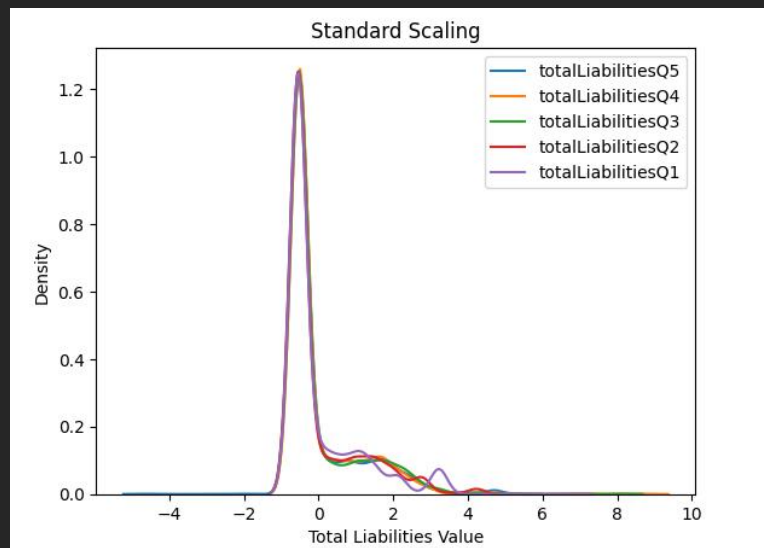
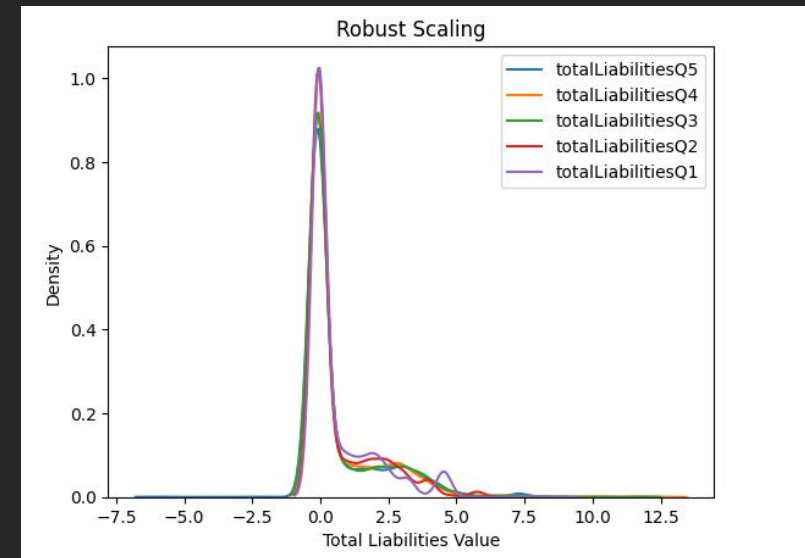
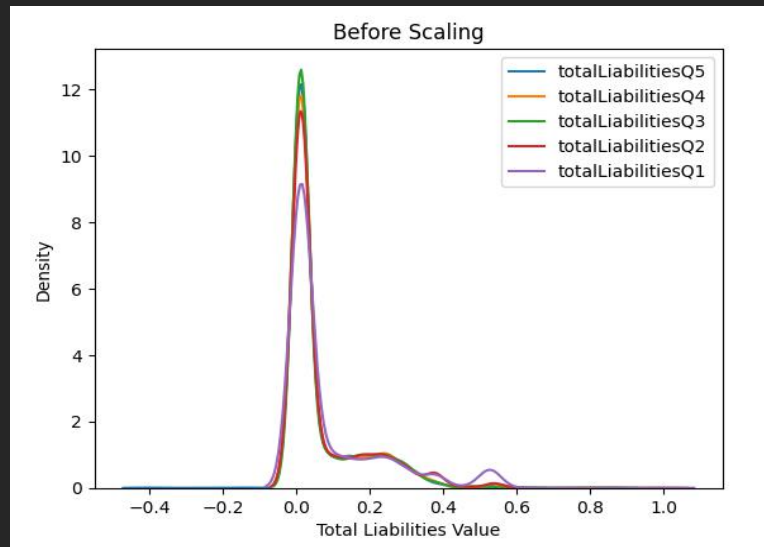
# Net Income



# Total Asset

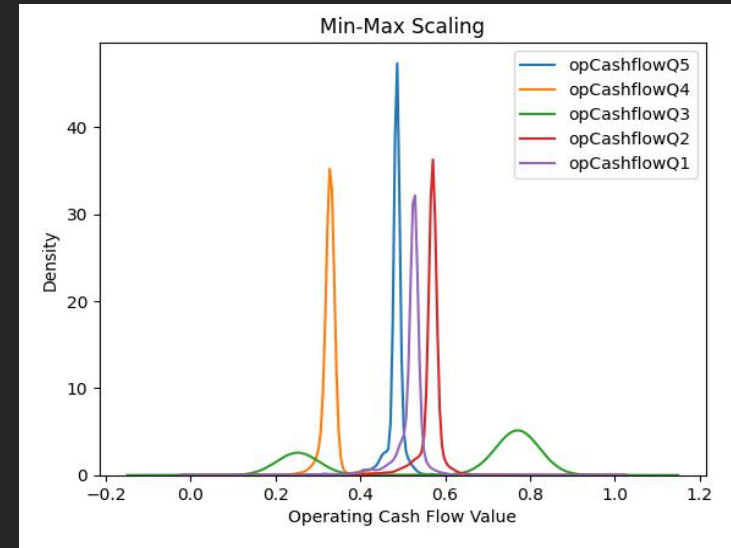
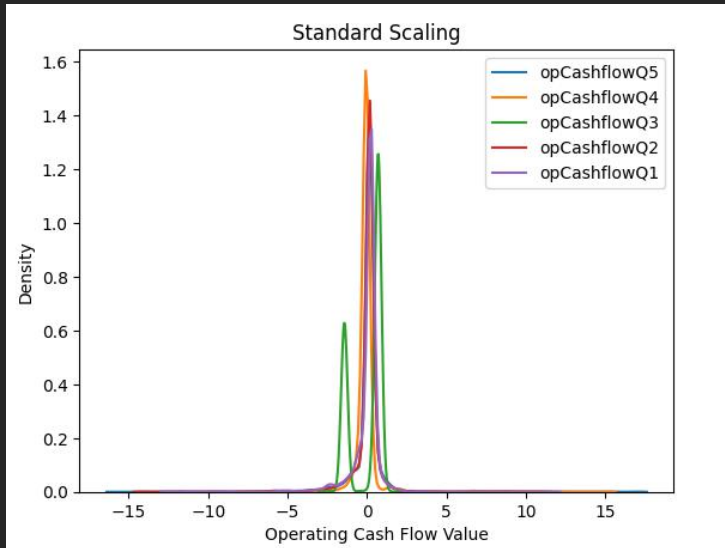
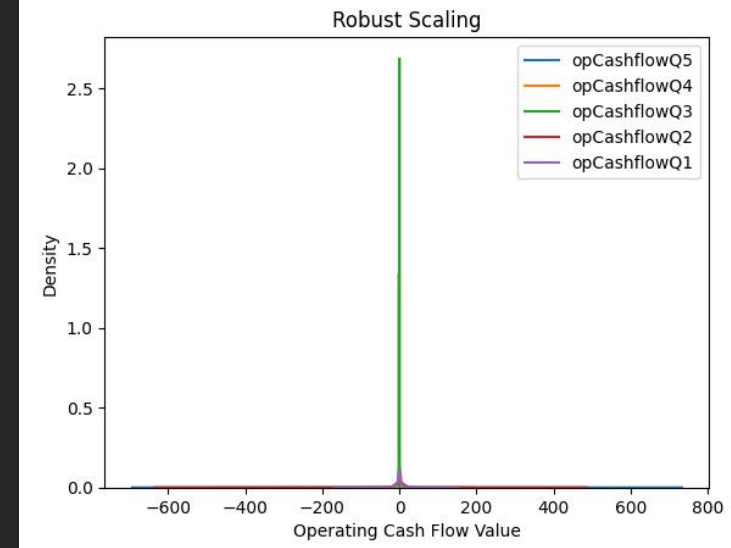
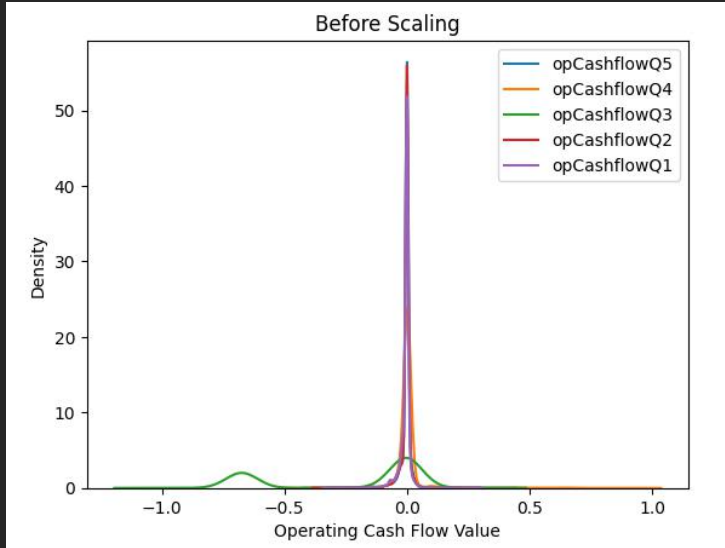


# Total Liabilities

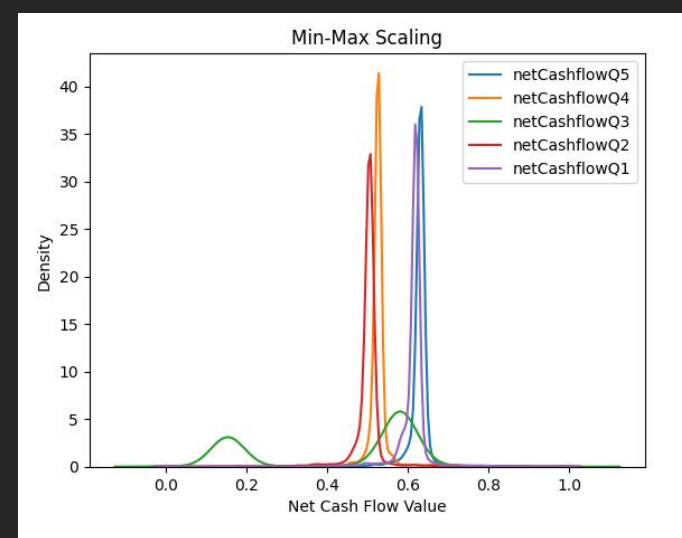
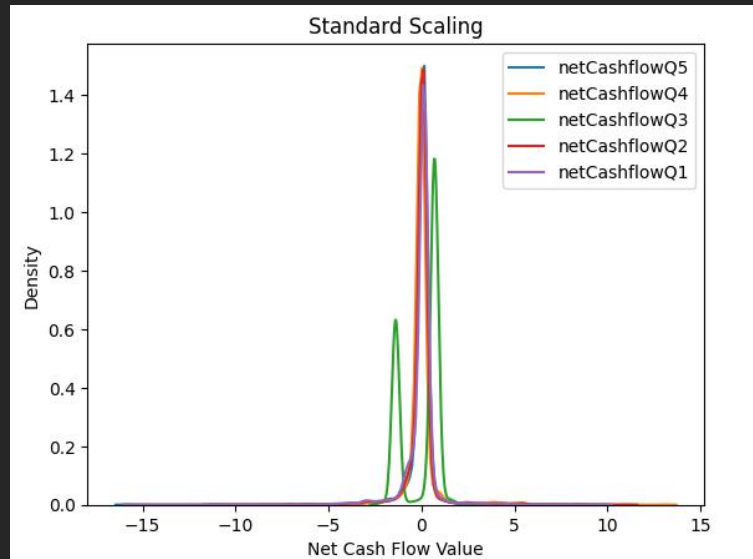
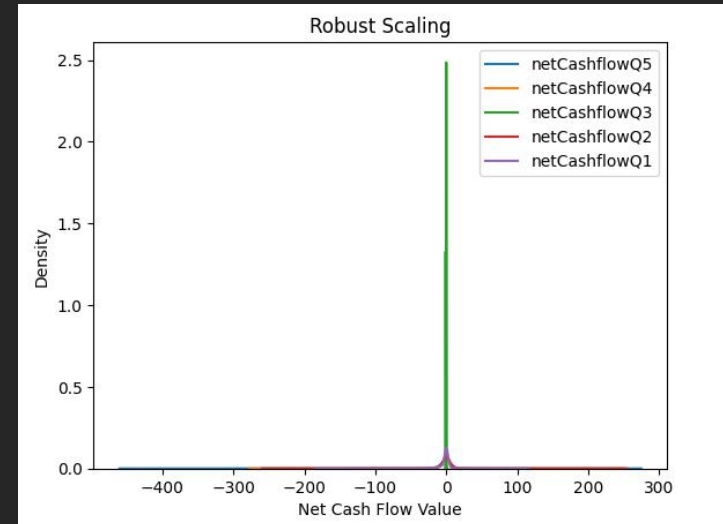
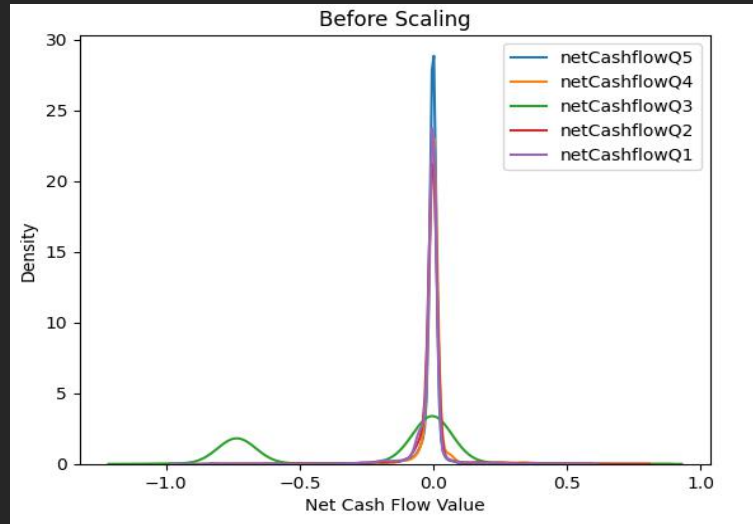




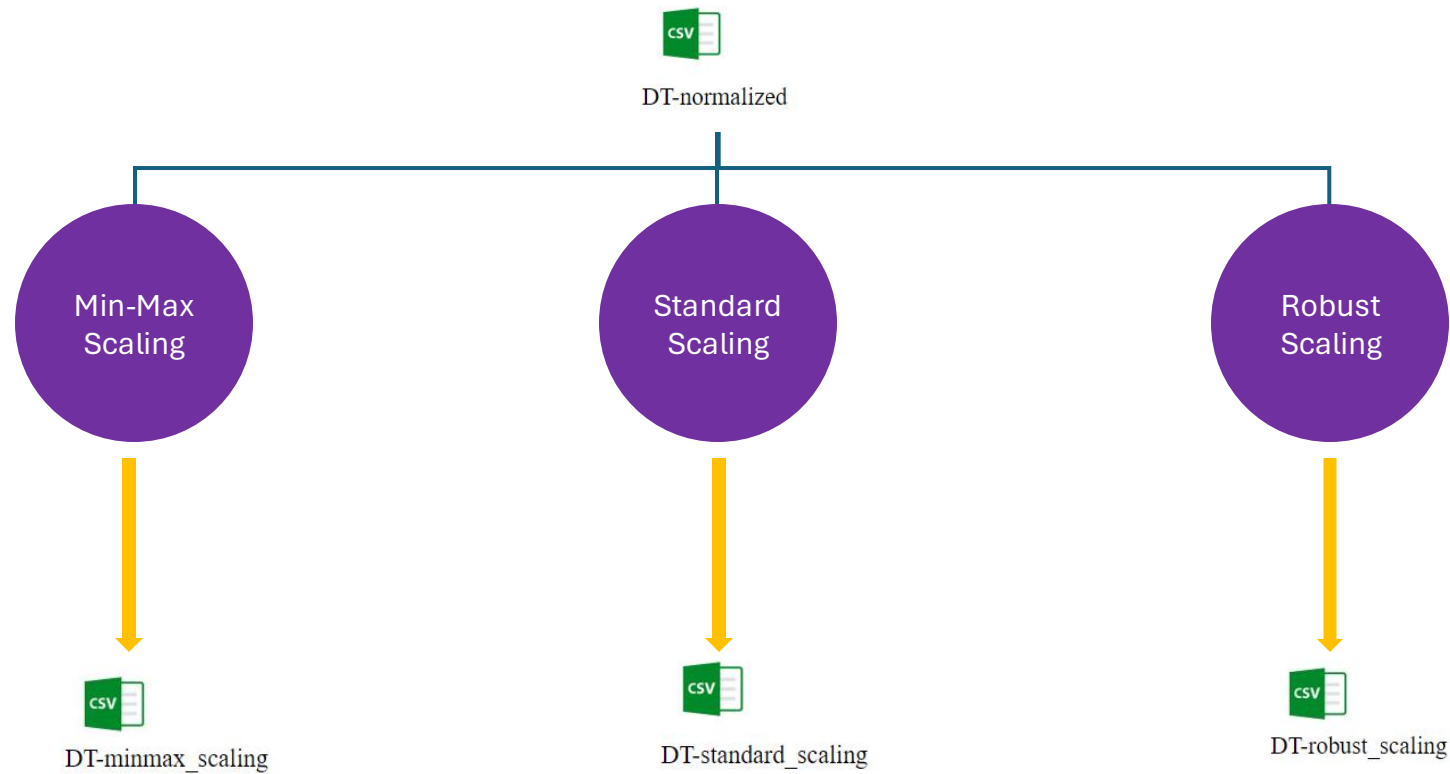
# Operating Cash Flow



# Net Cash Flow



## 4. Data Transformation – Resulting output data



# 4. Data Transformation – Data Encoding

- We used the following to do our data encoding API: sklearn.preprocessing.**LabelEncoder**
- We transform status columns from "listed", "delisted" to 0 and 1

## ✓ Data Encoding

```
encoder = LabelEncoder()  
df3_cleaned['status'] = encoder.fit_transform(df3_cleaned['status'])  
df3_cleaned
```

	status	salesQ5	salesQ4	salesQ3	salesQ2	salesQ1	netIncomeQ5	netIncomeQ4	netIncomeQ3	netIncomeQ2	...
0	0	4.744901e-03	4.260704e-03	5.315800e-03	0.004351	0.004593	-0.000097	-8.839976e-05	-5.628035e-05	-1.165012e-04	...
1	0	4.545807e-01	4.081927e-01	5.092752e-01	0.416830	0.440023	0.000005	9.235170e-06	1.633772e-05	3.597637e-06	...
2	0	4.545819e-01	4.081938e-01	5.092765e-01	0.416832	0.440024	0.000001	-1.854457e-07	-5.563370e-07	-2.781685e-07	...
3	0	4.744901e-03	4.260704e-03	5.315800e-03	0.004351	0.004593	-0.000097	-8.839976e-05	-5.628035e-05	-1.165012e-04	...
4	0	4.744901e-03	4.260704e-03	5.315800e-03	0.004351	0.004593	-0.000097	-8.839976e-05	-5.628035e-05	-1.165012e-04	...
...	...	...	...	...	...	...	...	...	...	...	...
2795	1	5.586174e-01	5.016129e-01	1.473739e-03	0.001158	0.540728	-0.001010	3.249664e-04	6.626467e-04	5.744568e-05	...
2796	1	1.895937e-05	1.004333e-06	7.850978e-08	0.332960	0.351486	0.004391	3.952301e-03	3.445666e-03	2.209753e-03	...
2797	1	2.108767e-03	2.479783e-03	1.443247e-03	0.000999	0.911613	0.000973	1.195231e-03	5.832293e-04	6.583507e-04	...
2798	1	7.998089e-09	2.132824e-08	9.331104e-08	0.599251	0.632594	-0.000194	3.176116e-03	1.404976e-03	-7.768811e-06	...
2799	1	1.532547e-03	2.321270e-03	2.031616e-03	0.001533	0.888784	-0.000010	7.094075e-04	-8.396840e-04	5.323506e-04	...

## 4. Data Transformation – ML Result

Condition	SVM Accuracy
Before Data Transformation	0.7976
After Data Transformation:	
MinMaxScaler	0.8619
StandardScaler	0.8810
RobustScaler	0.7976

# 5. Data Visualization and ETL-Use of Different Tools

- Python (Pandas, SQLite, BeautifulSoup, csv, requests, selenium), API, Google Collab.
- AWS Lambda, AWS SQS for creating multi-threading web scrapping tool.
- Delta Lake From Databricks (<https://docs.databricks.com/en/delta/index.html>).
- AWS Glue for ETL pipeline.

## 5.ETL to find Features-Feature Selection

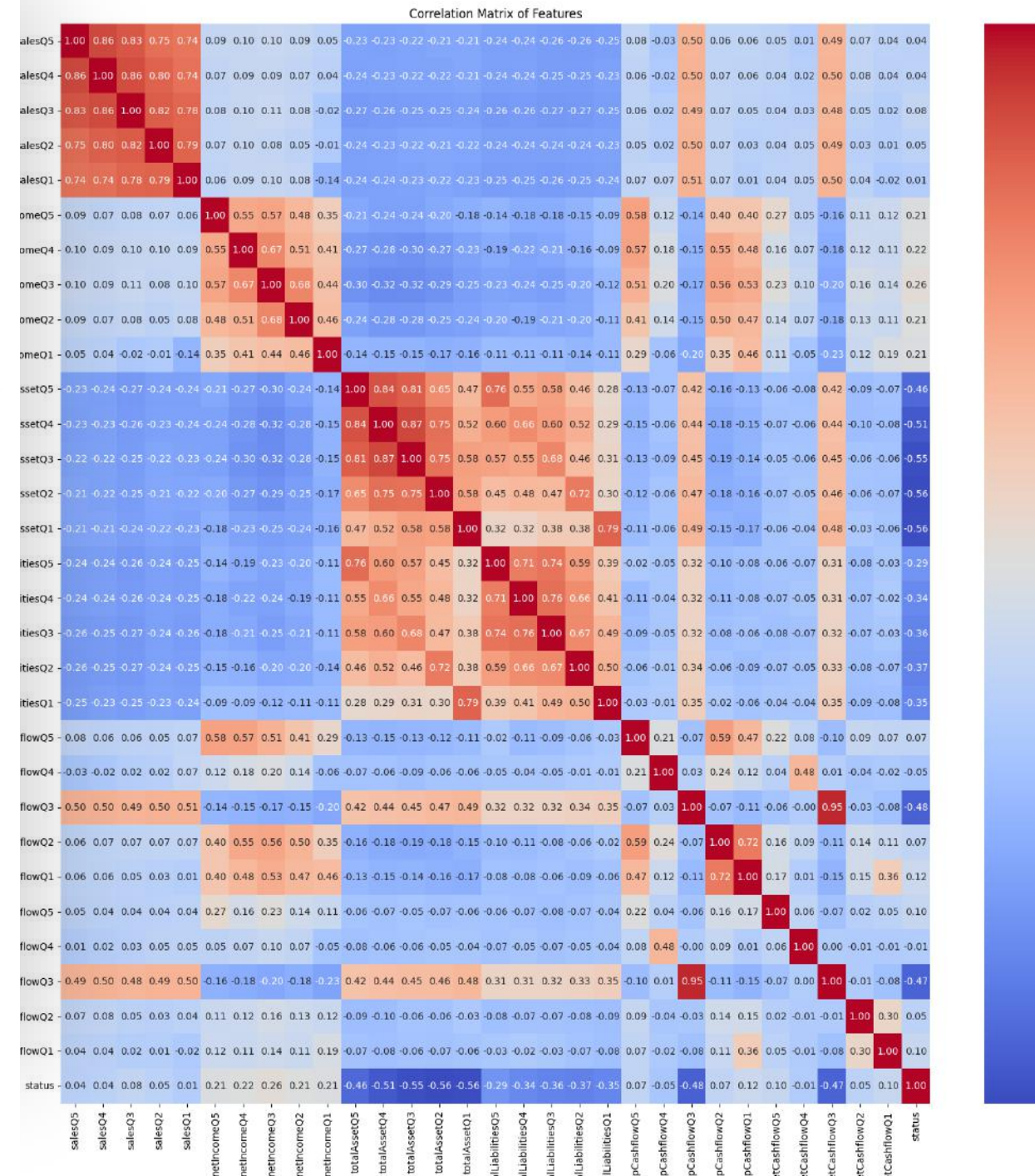
Removing irrelevant Features

We only keep Times Series data

- ~~Stock Symbol~~
- ~~Company Name~~
- ~~Exchange Center~~
- ~~IPO Delisted Date~~(data leakage, in a real world, you won't know when your company is delisted)

# 5.ETL to find Features- Feature Selection

- Set correlation threshold=0.85
- Drop columns:  
 'salesQ3', 'totalAssetQ1',  
 'totalLiabilitiesQ4', 'opCashflowQ3',  
 'salesQ5', 'totalAssetQ5', 'totalAssetQ4',  
 'salesQ4', 'salesQ2'





## 5.ETL to find Features- Feature Selection ML accuracy

Description	SVM Accuracy
Before Feature Selection	0.8821
After Feature Selection	0.8833



# 5.ETL to find Features- Features Transformation

We perform PCA for features Transformation

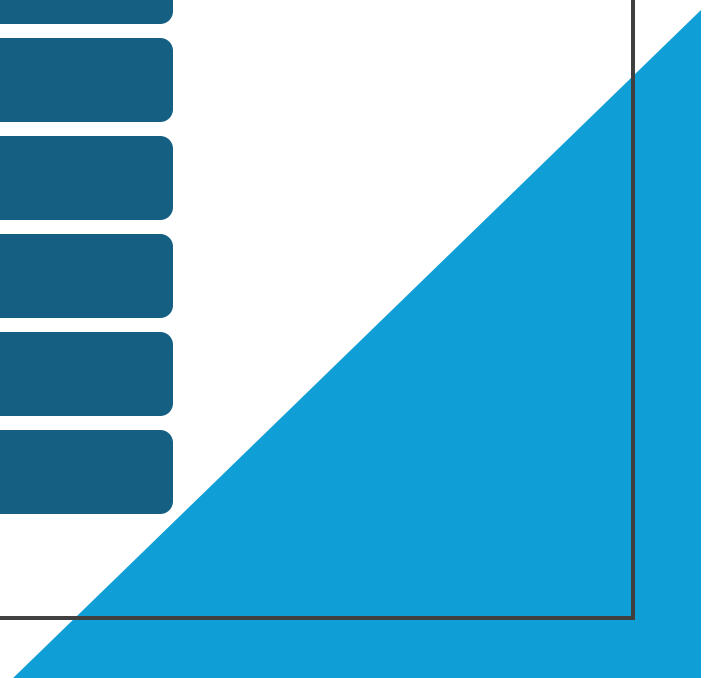
Why?

- It maximize the variance of the dataset.

- filter out noise and less significant details

- decrease the feature correlation

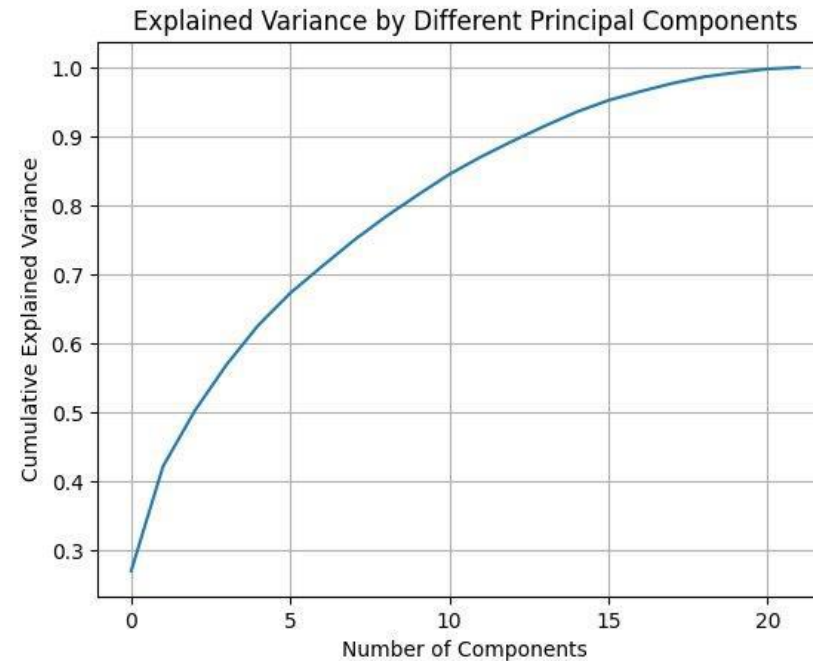
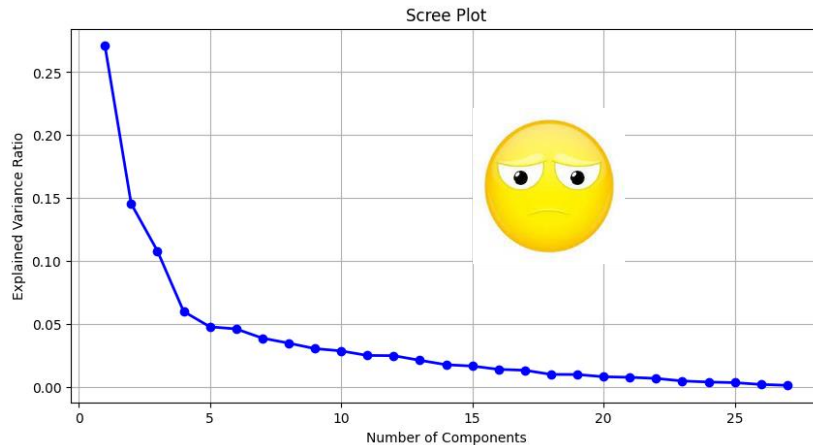
- easy for visualize as less columns remains



# 5.ETL to find Features-Hyperparameter Tuning

Objective: Let  $N = n$  of components, Select  $N$  to cover 95% of variance,  $1 \leq N \leq 31$

Plot the variance for each values:



Best hyperparameter  $N$   
=18

# 5.ETL to find Features-ML Accuracy

A screenshot of a Jupyter Notebook interface. The top bar shows the command 'principal\_df.head()' and a play button. Below it, a table displays the first five rows of principal components (PCA1 to PCA18). The data is as follows:

	PCA1	PCA2	PCA3	PCA4	PCA5	PCA6	PCA7	PCA8	PCA9	PCA10	PCA11	PCA12	PCA13	PCA14	PCA15	PCA16	PCA17	PCA18
0	-1.676521	-0.257513	2.018692	-0.043753	-0.208860	-0.133962	-0.007171	-0.106446	0.047713	0.058022	-0.003011	-0.013878	0.019516	-0.026920	-0.007651	-0.019073	0.010205	-0.047502
1	-2.640882	1.392012	-2.090201	0.535550	-0.188258	-0.299712	-0.233052	-0.198329	-0.171183	-0.268602	-0.006329	0.018682	0.179673	0.095123	0.066136	0.057468	-0.029552	-0.073766
2	-2.648911	1.397961	-2.087499	0.535232	-0.187373	-0.298361	-0.231196	-0.196152	-0.171992	-0.266810	-0.006455	0.018430	0.179822	0.095117	0.066745	0.057600	-0.029385	-0.072753
3	-1.676521	-0.257513	2.018692	-0.043753	-0.208860	-0.133962	-0.007171	-0.106446	0.047713	0.058022	-0.003011	-0.013878	0.019516	-0.026920	-0.007651	-0.019073	0.010205	-0.047502
4	-1.676521	-0.257513	2.018692	-0.043753	-0.208860	-0.133962	-0.007171	-0.106446	0.047713	0.058022	-0.003011	-0.013878	0.019516	-0.026920	-0.007651	-0.019073	0.010205	-0.047502

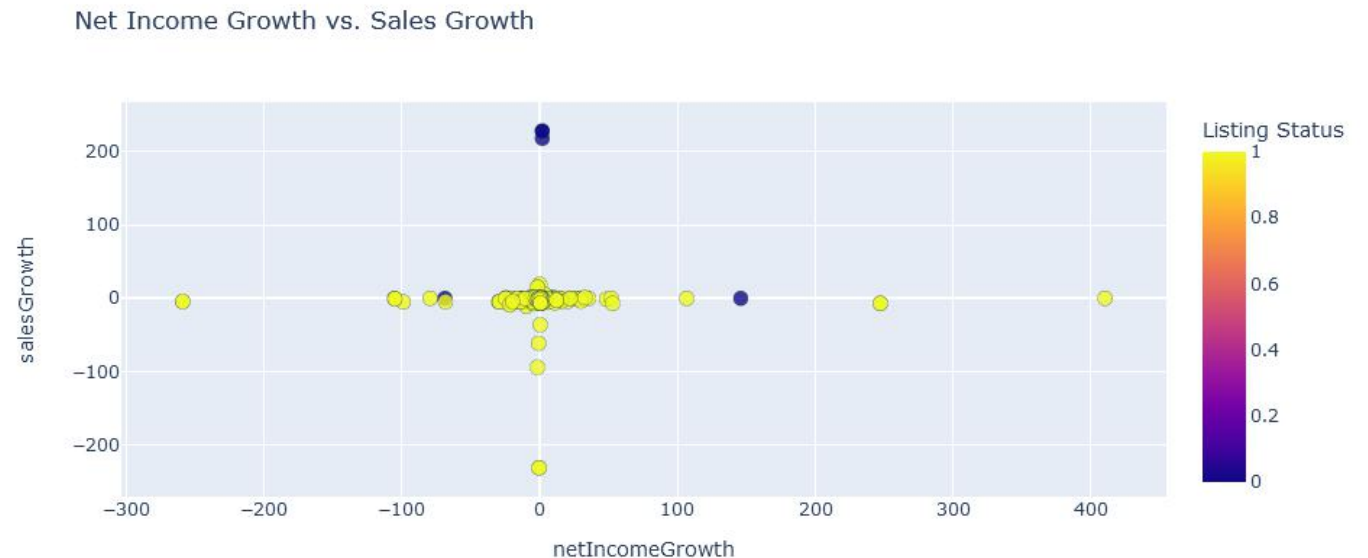
Description	SVM Accuracy
Before Feature Selection and Transformation	0.8821
After Feature Selection and Transformation	0.9000



# 5.1. Data Visualization/Dashboard

- Used Python Dash + Plotly for visualization
- Some problems:  
Timing across time-series data only has 5 points(Q1-Q5).
- Making a few types of plots less reliable.

Net Income Growth vs. Sales Growth vs. Status

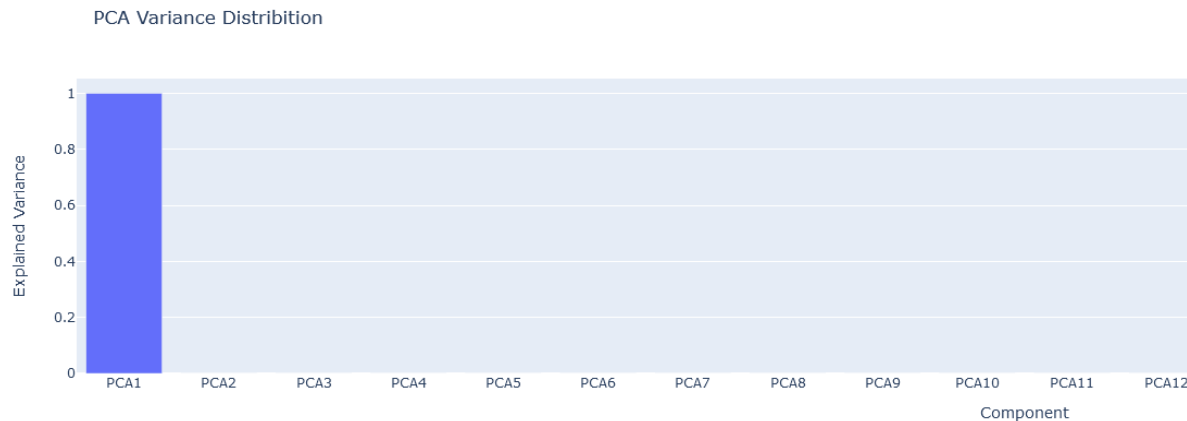


Scatter plot be looking like this 🤔

# 5.1. Data Visualization/Dashboard

- Principal component analysis is used to help with this problem.

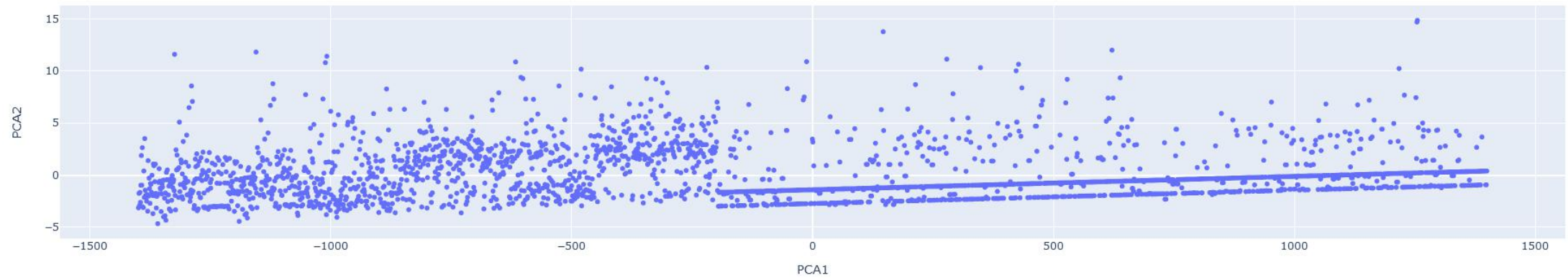
## PCA Visualization Dashboard



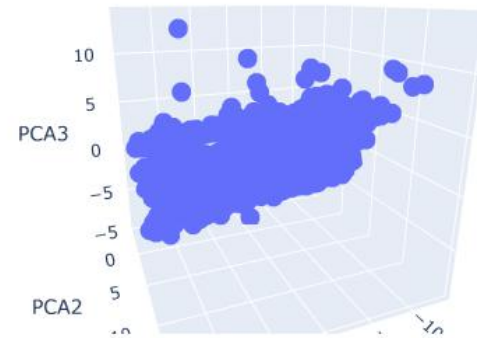
# 5.1. Data Visualization/Dashboard

📷 🔍 + 📏 🗨️ + 📄 🗑️

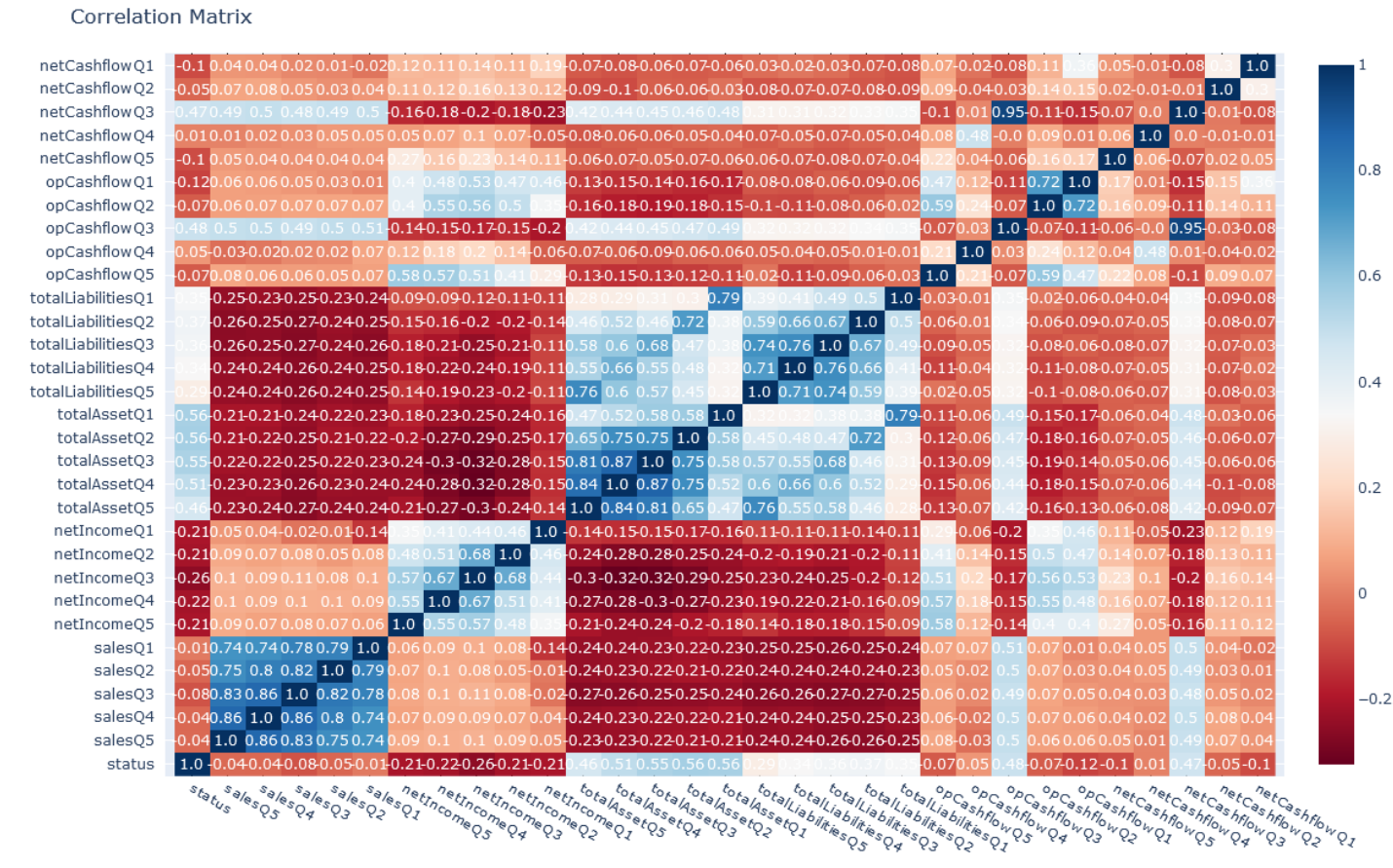
PCA Biplot (PCA1 vs PCA2)



3D PCA Biplot (PCA1 vs PCA2 vs PCA3)



# 5.1. Data Visualization/Dashboard



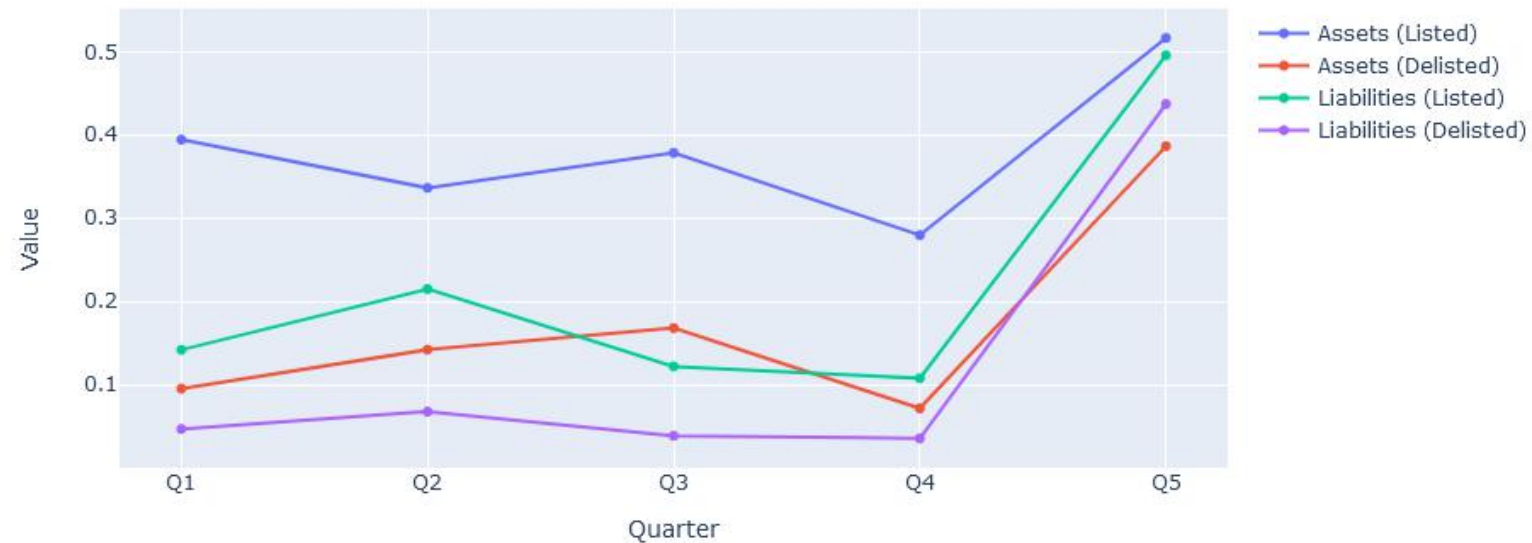
Correlation Matrix



# 5.1. Data Visualization/Dashboard

Select Metric for Time Series Analysis:

Assets/Liabilities Over Quarters for Listed vs Delisted Companies



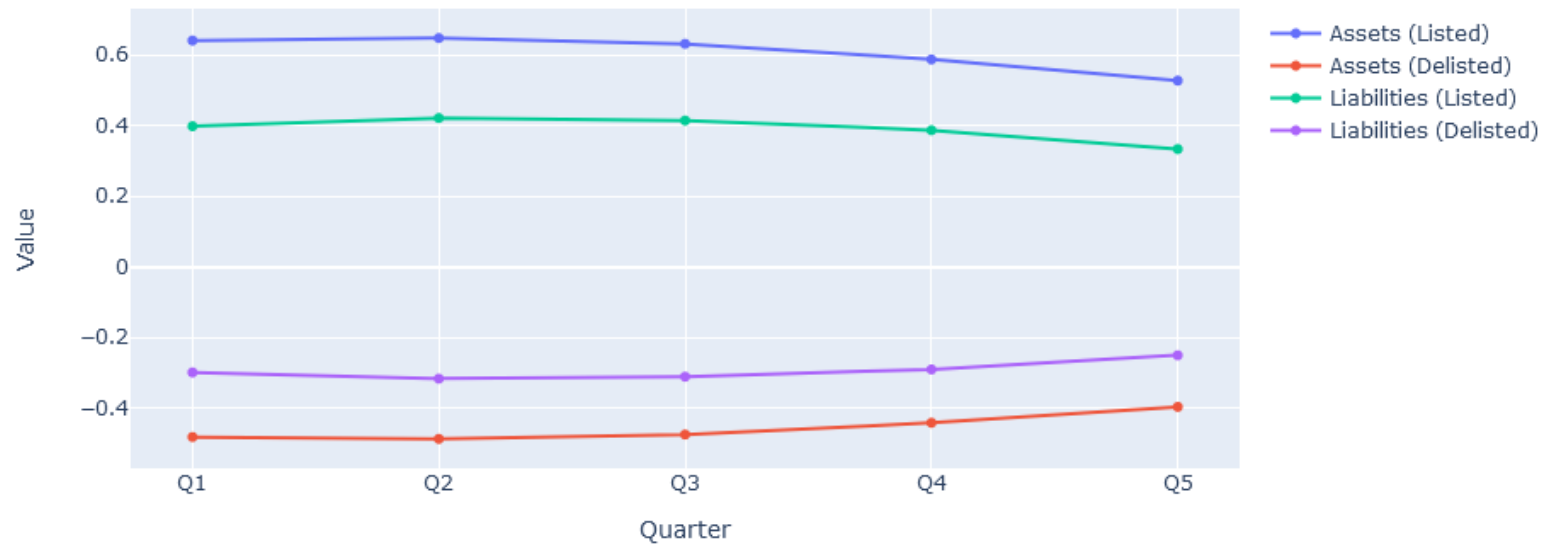
Minmax-scaled Dataset

# 5.1. Data Visualization/Dashboard

Select Metric for Time Series Analysis:



Assets/Liabilities Over Quarters for Listed vs Delisted Companies



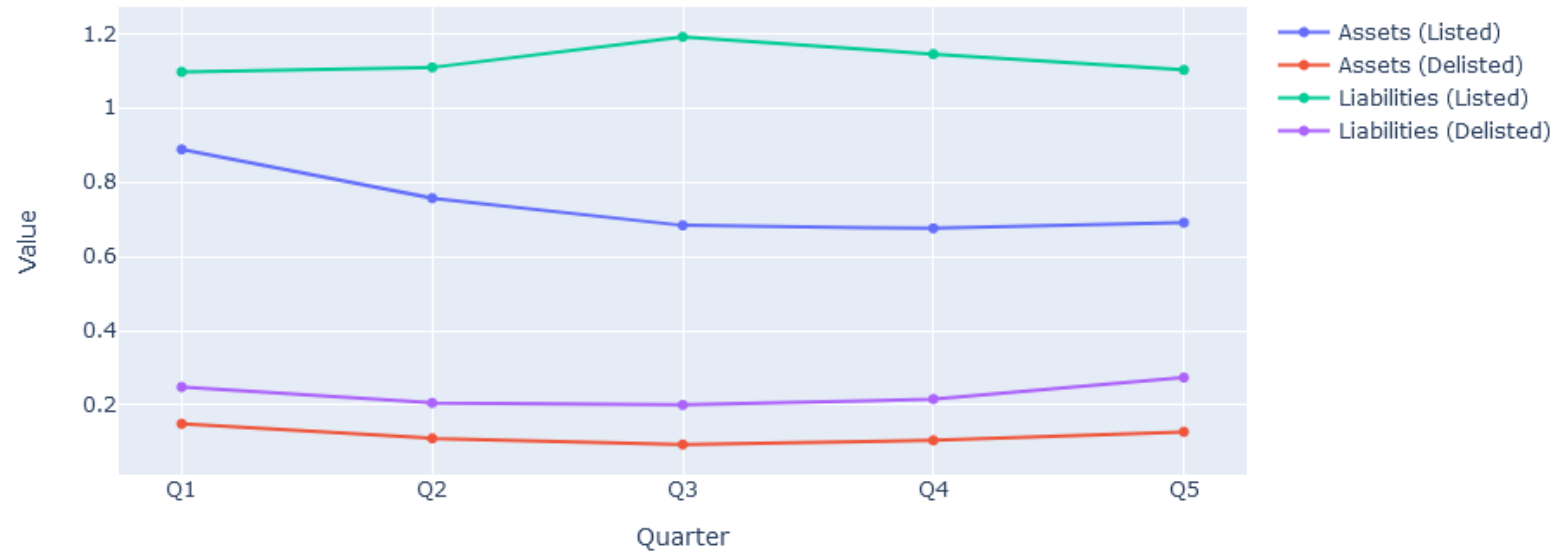
Standard-scaled Dataset

# 5.1. Data Visualization/Dashboard

Select Metric for Time Series Analysis:

Total Assets  Total Liabilities  ▼

Assets/Liabilities Over Quarters for Listed vs Delisted Companies



Robust-scaled Dataset