

2021

Big Data Architecture & Governance



Northeastern University

Assignment Name:

Group Assignment 1 – Team 4

Student Names:

Apoorva Mishra

Atulya Sharma

Dhruv Panchal



Contents

Assignment	2
1.1. Case	2
1.2. Assignment Goals	2
1.2.1. Visualization Deliverables	2
1.2.2. Other deliverables	3
Documentation	4
1. Vision Diagram.....	4
2. Velero Screenshots.....	5
3. Data Profiling, Wrangling and Cleansing	18
4. Database Installation	31
5. Data Mapping and Integration.....	32
6. Data Validation and Data Visualization.....	35
7. System Integration and User Acceptance Testing.....	43
8. Challenges Encountered.....	45
9. End User Instructions.....	46



Assignment

1.1. Case

Each team should select a dataset to analyze and build an analytical dashboard as a Proof-of-concept to illustrate the value of data driven analytics. You need to present your dataset.

1.2. Assignment Goals

To work with datasets, Perform/Create:

- Group Assignment/Project on Velero with below mentioned activities:
 - Tasks/ Project Short/ Long Form/ Group Allocation/ Timesheet/ Issues & Risks.
- Data Profiling – Using Python profiling library, describe your understanding of the data.
- Data Wrangling and Cleansing - Pandas/Alteryx/XSV
 - Filtering and Aggregating if needed.
 - Missing value handling.
 - Deriving additional columns from existing datasets if needed.
 - Cleaning (removing blank spaces, formatting dates, Capitalizing etc.) .
- Database Installation: Install NEO4J database .
- Data Mapping and Integration to your Database for the Entire Dataset.
- Business and Technical Metadata – develop business term list describing all the data elements available in the file.
- Data Validation and Data Visualization using Python – Validate the data using python and provide a dashboard using python visualization libraries.
- System Integration and User Acceptance Testing - Test Cases – describe your validation & testing process.
- Risks/Issues of project.
- Describe challenges encountered and how you resolved them.
- End User Instructions (Steps to run your Dashboard) – provide a full description how to run your process:
 - Database Creation and load.
 - Visualization interpretation - describe information regarding your findings.

1.2.1. VISUALIZATION DELIVERABLES

Once you wrangle/clean/join/integrate the data, import the data into **NEO4J** and illustrate how to use the appropriate graph to illustrate various aspects of analysis.

Questions to consider:

- Columns used for dimensions, and columns that are used for measurement.
- How would you generate new dimensions?
- Who would use this dashboard and how they benefit from your dashboard?
- What value would be generated using this dashboard?



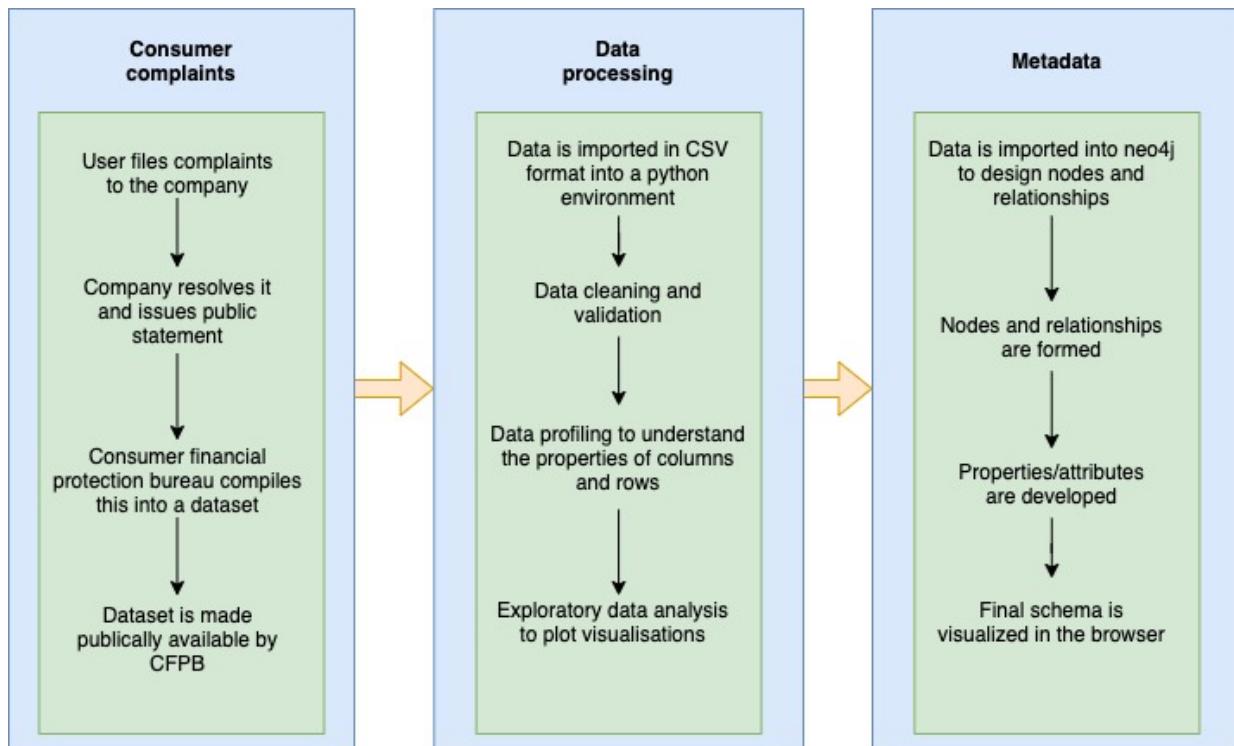
1.2.2. OTHER DELIVERABLES

- Presentation of the entire work from the first step till the dashboards including the Velero screenshots.
- Business and technical metadata presentation – Identifying all available business terms and extracting related technical metadata.
- Complete explanation of the dashboard and usability.
- Complete instruction as how to implement and run the database load, technical meta data extraction, and dashboard.



Documentation

1. Vision Diagram





2. Velero Screenshots

- Velero's Enterprise Transparency Platform consists of enterprise portfolio planning and execution.
- Velero ETP enable executives quickly view the calculated estimated costs and full resource allocation information. Velero ETP includes fully integrated time, resource and budget management modules.
- Velero's integrated solution provides a fast and easy access to all aspect of portfolio planning and execution such as defining strategic objectives, estimated resources, budget, risk/issue, mandate, status and activities, with the ability to attach and upload related documents directly within the product

Below are several screenshots for different use cases in Velero.

Issues/Risks

We have identified two risks that we have accounted in Velero.

1. Issue while loading data in Neo4j

The screenshot shows the 'Risk/Issue Add' dialog box. The form contains the following fields:

Field	Value
Issue Date	04/08/2021
Date Expected	04/10/2021
Designation	--Issue--
Probability	75
Date Assigned	04/08/2021
Date Closed	04/15/2021
Impact	9
Exposure Recipient	Company
Assigned to	External Resource
Department Employee	Aboorva.Mishra 1A99001
Print	Include
Action	Escalate
Description*	Issues while loading the data on Neo4j. Data takes a lot of time to load and the script seems to run endlessly.
Mitigation*	[Empty text area]



2. Data Incorrectness

Task is updated successfully.

Risk/Issue Add

Save

Issue Date 04/12/2021	Designation -Issue--	Probability 50	Impact 7	Print Include
Date Expected 04/15/2021	Date Assigned 04/12/2021	Date Closed 04/16/2021		
Exposure Recipient Company	Assigned to Apoorva Mishra	Department Employee Apoorva.Mishra (A9900)	Action Consolidate	
Description* Data incorrectness. We have null values for multiple columns				
mitigation* Perform data cleaning and data wrangling and data validation to have consistent data				

3. Trouble connecting data with Neo4j

Risk/Issue Add

Save

Issue Date 04/08/2021	Designation -Issue--	Probability 89	Impact 4	Print Include
Date Expected 04/09/2021	Date Assigned 04/08/2021	Date Closed 04/09/2021		
Exposure Recipient External Resource	Assigned to Apoorva.Mishra (A9900)	Department Employee Apoorva.Mishra (A9900)	Action Escalate	
Description* Trouble connecting Neo4j with CSV data				
mitigation*				



Heat Map Status

List Project*'s By Include Options Portfolio Functions

Extended Status Maintenance

==> Record is Updated!

Heatmap Status	Accomplishments	Not Accomplished	Next Week Focus	Go Green Actions	Save
Schedule	Green	We are following the project timeline and adhering to the schedule			
Resource	Green	We have divided the resources per our requirements			
\$ Cost	Green	We have optimized the budget utilization			
Scope	Green	To get a visual representation of the consumer complaints data			

Resource Planning

Resource Management for: Group 4 - Consumer Complaints Analysis (Start Planning year: 2021)

Category/Name	Year	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
Application Engineer	2021	0.00	0.00	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Data Analyst	2021	0.00	0.00	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Data Engineer	2021	0.00	0.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Data visualizers	2021	0.00	0.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Database Admin	2021	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Students	2021	0.00	0.00	0.80	1.20	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
[444] Apoorva, Mishra	2021	0.00	0.00	40.00	60.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
[454] Atulya, Sharma	2021	0.00	0.00	40.00	60.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
[446] Dhruv, Panchal	2021	0.00	0.00	40.00	60.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Systems Architect	2021	0.00	0.00	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00



Milestone/Tasks with Allotted Times per Activity

Milestone/TaskList1

		PLC	Order	Type	Milestone/Task Description	%Complete	Est Hours	Est HIC	Assigned To	Start Date	End Date	Status
<input checked="" type="checkbox"/>	Complete	1-Initiation	1	Analysis	Project Initiation	100.00%	2.00	0.00	Apoorva, Mishra	04/05/2021	04/09/2021	Complete
<input checked="" type="checkbox"/>	Complete	1-Initiation	1	Analysis	Project Initiation	100.00%	2.00	0.00	Atulya, Sharma	04/05/2021	04/09/2021	Complete
<input checked="" type="checkbox"/>	Complete	1-Initiation	2	Analysis	Architecture Design	100.00%	10.00	0.00	Atulya, Sharma	04/15/2021	04/24/2021	Complete
<input checked="" type="checkbox"/>	Complete	1-Initiation	3	Analysis	Non-functional Requirements	100.00%	15.00	0.00	Dhruv, Panchal	04/21/2021	04/24/2021	Complete
<input checked="" type="checkbox"/>	Green	2-Planning	1	Analysis	Talent Management & Optimization	90.00%	4.00	0.40	Dhruv, Panchal	04/17/2021	04/24/2021	Inprocess
<input checked="" type="checkbox"/>	Complete	2-Planning	1	Analysis	Project Planning	100.00%	5.00	0.00	Atulya, Sharma	04/08/2021	04/13/2021	Complete
<input checked="" type="checkbox"/>	Complete	3-Execution	1	Analysis	Architecture review & Approval	100.00%	16.00	0.00	Apoorva, Mishra	04/15/2021	04/22/2021	Complete
<input checked="" type="checkbox"/>	Complete	3-Execution	2	Development	Data Cleaning	100.00%	14.00	0.00	Atulya, Sharma	04/12/2021	04/16/2021	Complete
<input checked="" type="checkbox"/>	Complete	3-Execution	2	Development	Functional Requirements	100.00%	8.00	0.00	Atulya, Sharma	04/17/2021	04/24/2021	Complete
<input checked="" type="checkbox"/>	Complete	3-Execution	3	Development	Data Extraction	100.00%	5.00	0.00	Atulya, Sharma	04/19/2021	04/21/2021	Complete
<input checked="" type="checkbox"/>	Complete	3-Execution	3	Development	Neo4j Script Development	100.00%	19.00	0.00	Apoorva, Mishra	04/12/2021	04/16/2021	Complete

Milestone/TaskList2

Project* Milestones/Tasks: [Group 4 - Consumer Complaints Analysis]											GanttChart	Burndown	New	Back
		Copy	Excel	Column visibility	Show 11 entries	Search:								
		PLC	Order	Type	Milestone/Task Description	%Complete	Est Hours	Est HIC	Assigned To	Start Date	End Date	Status		
<input checked="" type="checkbox"/>	Complete	3-Execution	3	Development	Vision Diagram	100.00%	9.00	0.00	Dhruv, Panchal	04/12/2021	04/16/2021	Complete		
<input checked="" type="checkbox"/>	Complete	3-Execution	3	Development	Visualization Dashboard	100.00%	15.00	0.00	Atulya, Sharma	04/12/2021	04/16/2021	Complete		
<input checked="" type="checkbox"/>	Complete	3-Execution	4	Other	Software Installation	100.00%	2.00	0.00	Apoorva, Mishra	04/07/2021	04/08/2021	Complete		
<input checked="" type="checkbox"/>	Complete	3-Execution	5	QA	Write Master Test Plan	100.00%	2.00	0.00	Dhruv, Panchal	04/21/2021	04/23/2021	Complete		
<input checked="" type="checkbox"/>	Complete	3-Execution	7	Development	Transformation & Load Process	100.00%	10.00	0.00	Atulya, Sharma	04/13/2021	04/16/2021	Complete		
<input checked="" type="checkbox"/>	Complete	3-Execution	10	QA	User Acceptance Testing	100.00%	3.00	0.00	Dhruv, Panchal	04/22/2021	04/22/2021	Complete		
<input checked="" type="checkbox"/>	Complete	4-Controlling	1	Analysis	Vendor/Tool Selection	100.00%	8.00	0.00	Dhruv, Panchal	04/17/2021	04/24/2021	Complete		
<input checked="" type="checkbox"/>	Complete	5-Closing	2	Systems	Implementation	100.00%	7.00	0.00	Apoorva, Mishra	04/22/2021	04/22/2021	Complete		
<input checked="" type="checkbox"/>	Complete	5-Closing	4	Other	Closing	100.00%	2.00	0.00	Atulya, Sharma	04/23/2021	04/23/2021	Complete		
Showing 19 to 20 of 20 entries											Previous	1	2	Next



Project Initiation

ES/Tasks, [Group 4 - Consumer Complaints Analysis]

Milestone/Task Maintenance

Order	PLC Category	Type	Milestone/Task Description
1	1-Initiation	Analysis	Project Initiation

Assigned To: Atulva. Sharma (A99001) | Start Date: 04/05/2021 | End Date: 04/09/2021 | Heatmap: 1-Green | Status: Complete

% Completed: 100 | Est Hours: 2.00 | Rem Hours: 0.0 | Print [Y]: Include

Save

Description*: Discussed project roadmap divided tasks

Note(1)*: Note(2)*

Project Initiation

ES/Tasks, [Group 4 - Consumer Complaints Analysis]

Milestone/Task Maintenance

Order	PLC Category	Type	Milestone/Task Description
1	1-Initiation	Analysis	Project Initiation

Assigned To: Aboorva. Mishra (A99001) | Start Date: 04/05/2021 | End Date: 04/09/2021 | Heatmap: 4-Complete | Status: Complete

% Completed: 100.00 | Est Hours: 2.00 | Rem Hours: 0.0 | Print [Y]: Include | **Delete** | **Save**

Description*: Started Project on Velero for Project Management Outlined the requirements Divided all tasks

Note(1)*: Note(2)*



Data Extraction

Milestone/Task Maintenance

Order	PLC Category	Type	Milestone/Task Description
3	3-Execution	Development	Data Extraction

Assigned To	Start Date	End Date	Heatmap	Status
Atulva. Sharma [A99001]	04/19/2021	04/21/2021	1-Green	Complete

% Completed	Est Hours	Rem Hours	Print [Y]	Save
100	5	0.0	Include	Save

Description*
Used python script to extract the technical metadata

Note(1)*

Note(2)*

Data Cleaning

Milestone/Task Maintenance

Order	PLC Category	Type	Milestone/Task Description
2	3-Execution	Development	Data Cleaning

Assigned To	Start Date	End Date	Heatmap	Status
Atulva. Sharma [A99001]	04/12/2021	04/16/2021	1-Green	Complete

% Completed	Est Hours	Rem Hours	Print [N]	Delete	Save
100	14.00	0.0	Include	Delete	Save

Description*
Written Neo

Note(1)*

Note(2)*



Architecture Design

ES TASKS: [Group 4 - Consumer Complaints Analysis]

Milestone/Task Maintenance

Order	PLC Category	Type	Milestone/Task Description
2	1-Initiation	Analysis	Architecture Design

Assigned To: Atulva. Sharma (A99001) Start Date: 04/15/2021 End Date: 04/24/2021 Heatmap: 4-Complete Status: Complete

% Completed: 100.00 Est Hours: 10.00 Rem Hours: 0.0 Print [Y]: Include

Delete **Save**

Description*
Worked on the architecture of the data model based on the given dataset
figured out the relationships between columns

Note(1)*

Note(2)*

Non-Functional Requirement

Milestone/Task Maintenance

Order	PLC Category	Type	Milestone/Task Description
3	1-Initiation	Analysis	Non-functional Requirements

Assigned To: Dhruv. Panchal (A99001) Start Date: 04/21/2021 End Date: 04/24/2021 Heatmap: 4-Complete Status: Complete

% Completed: 100.00 Est Hours: 15.00 Rem Hours: 0.0 Print [Y]: Include

Delete **Save**

Description*
Prepared Report based on all inputs
Prepared Presentation for stakeholders

Note(1)*

Note(2)*



Neo4j Script Development

Milestone/Tasks, [Group 4 - Consumer Complaints Analysis]

Milestone/Task Maintenance

Order	PLC Category	Type	Milestone/Task Description
3	3-Execution	Development	Neo4j Script Development

Assigned To: Aboorva. Mishra (A99001) | Start Date: 04/12/2021 | End Date: 04/16/2021 | Heatmap: 1-Green | Status: Complete

% Completed: 100 | Est Hours: 19 | Rem Hours: 0.0 | Print [Y]: Include

Description*: Scripted queries for pulling data into Neo4j
Developed a graph based on the data model.

Note(1)*

Note(2)*

Save

Vision Diagram

Milestone/Tasks, [Group 4 - Consumer Complaints Analysis]

Milestone/Task Maintenance

Order	PLC Category	Type	Milestone/Task Description
3	3-Execution	Development	Vision Diagram

Assigned To: Dhruv. Panchal (A99001) | Start Date: 04/12/2021 | End Date: 04/16/2021 | Heatmap: 1-Green | Status: Complete

% Completed: 100 | Est Hours: 9 | Rem Hours: 0.0 | Print [Y]: Include

Description*: Prepared vision diagram based on the requirement.

Note(1)*

Note(2)*

Save



Visualization Dashboards

Project Tasks: [Group 4 - Consumer Complaints Analysis]

Milestone/Task Maintenance

Order	PLC Category	Type	Milestone/Task Description
3	3-Execution	Development	Visualization Dashboard

Assigned To: Atulva. Sharma (A99001) Start Date: 04/12/2021 End Date: 04/16/2021 Heatmap: 1-Green Status: Complete

% Completed: 100 Est Hours: 15 Rem Hours: 0.0 Print [Y]: Include

Save

Description*: Prepared Dashboards in python for visual representation and understanding of stakeholders

Note(1)* Note(2)*

Project Planning

Project Phases: 1-Initiation 1 Analysis Project Initiation 100.00% 2.00 0.00 Apoorva, 04/05/2021 04/09/2021

Milestone/Task Maintenance

Order	PLC Category	Type	Milestone/Task Description
1	2-Planning	Analysis	Project Planning

Assigned To: Atulva. Sharma (A99001) Start Date: 04/08/2021 End Date: 04/13/2021 Heatmap: 4-Complete Status: Complete

% Completed: 100.00 Est Hours: 5.00 Rem Hours: 0.0 Print [Y]: Include

Delete Save

Description*: Designed a roadmap about the tasks to be completed
Divided the tasks
Scheduled calls for weekly update

Note(1)* Note(2)*



Talent Management & Optimization

Milestone/Task Maintenance

Order	PLC Category	Type	Milestone/Task Description	
1	2-Planning	Analysis	Talent Management & Optimization	
Assigned To	Start Date	End Date	Heatmap	Status
Dhruv. Panchal (A99001)	04/17/2021	04/24/2021	1-Green	Inprocess
% Completed	Est Hours	Rem Hours	Print [Y]	
90.00	4.00	0.4	Include	<button>Delete</button> <button>Save</button>
Description*			Note(1)*	
Dedicated resources for all tasks Integrated all pieces together for readiness of deliverables Followed up on the tasks assigned via weekly updates				
			Note(2)*	

Architecture Review & Approval

Milestone/Task Maintenance

Order	PLC Category	Type	Milestone/Task Description	
1	3-Execution	Analysis	Architecture review & Approval	
Assigned To	Start Date	End Date	Heatmap	Status
Aboorva. Mishra (A99001)	04/15/2021	04/22/2021	4-Complete	Complete
% Completed	Est Hours	Rem Hours	Print [Y]	
100.00	16.00	0.0	Include	<button>Delete</button> <button>Save</button>
Description*			Note(1)*	
Got the data model reviewed Collected feedback and reported to the team Worked with the team on correcting data per the feedback				
			Note(2)*	



Write Master Test Plan

The screenshot shows the 'Milestone/Task Maintenance' form. The 'Milestone/Task Description' field contains 'Write Master Test Plan'. The 'Assigned To' field is set to 'Dhruv. Panchal [A99001]'. The 'Start Date' is '04/21/2021' and the 'End Date' is '04/23/2021'. The 'Heatmap' status is '1-Green' and the 'Status' is 'Complete'. The 'Description*' field is empty. There are two notes: 'Note(1)*' and 'Note(2)*', both of which are also empty.

Transformation and Load Processing

The screenshot shows the 'Milestone/Task Maintenance' form. The 'Milestone/Task Description' field contains 'Transformation & Load Process'. The 'Assigned To' field is set to 'Atulva. Sharma [A99001]'. The 'Start Date' is '04/13/2021' and the 'End Date' is '04/16/2021'. The 'Heatmap' status is '1-Green' and the 'Status' is 'Complete'. The 'Description*' field contains 'Performed ETL on data'. There are two notes: 'Note(1)*' and 'Note(2)*', both of which are empty.



User Acceptance Testing

Milestone/Task Maintenance

Order	PLC Category	Type	Milestone/Task Description
1C	3-Execution	QA	User Acceptance Testing

Assigned To	Start Date	End Date	Heatmap	Status
Dhruv. Panchal (A99001)	04/22/2021	04/22/2021	1-Green	Complete

% Completed	Est Hours	Rem Hours	Print [N]
100	3.00	0.0	Include

Delete **Save**

Description*
Note(1)*
Note(2)*

Software Installation

Milestone/Task Maintenance

Order	PLC Category	Type	Milestone/Task Description
4	3-Execution	Other	Software Installation

Assigned To	Start Date	End Date	Heatmap	Status
Apoorva. Mishra (A99001)	04/07/2021	04/08/2021	1-Green	Complete

% Completed	Est Hours	Rem Hours	Print [Y]
100	2	0.0	Include

Save

Description*
Installed all required
Note(1)*
Note(2)*



Closing

Milestone/Task Maintenance

Order	PLC Category	Type	Milestone/Task Description
4	5-Closna	Other	Closing

Assigned To	Start Date	End Date	Heatmap	Status
Atulva. Sharma [A99001]	04/23/2021	04/23/2021	1--Green	Complete

% Completed	Est Hours	Rem Hours	Print [N]	Delete	Save
100.00	2.00	0.0	Include	<button>Delete</button>	<button>Save</button>

Description*
Completed all tasks
Verified pre production checklist

Note(1)*

Note(2)*



3. Data Profiling, Wrangling and Cleansing

The dataset we are using is called “Consumer Complaints”. This data is a collection of complaints about consumer financial products and services that was sent to the companies for response. Complaints are published after the company responds, confirming a commercial relationship with the consumer, or after 15 days, whichever comes first.

The dataset comprises of Consumer Complaints on financial products, and we'll see how to classify consumer complaints text into these categories: Debt collection, Consumer Loan, Mortgage, Credit card, Credit reporting, Student loan, Bank account or service, Payday loan, Money transfers, Other financial service, Prepaid card.

Data Profiling has been done using the pandas profiling library on python. It gives us information that is an extension of the info provided after using `df.describe()` function of pandas. A detailed report can be generated in the form of an interactive HTML page or can be integrated directly into the python environment.

Initial Installation of Pandas Profiling Library

Pandas Profiling can be installed by running the following on the Command Line Interface:

```
pip install pandas-profiling
```

Getting Started

To use pandas-profiling, we need to import it to our notebook using the following code:

```
Import numpy as np  
Import pandas as pd  
From pandas_profiling import ProfileReport
```

This is followed by loading the dataset on the python notebook.

```
df1 = pd.read_csv("Consumer_Complaints.csv")
```

To view the dataset, type the following code:



df1.head()

Date received	Product	Sub-product	Issue	Sub-issue	Consumer complaint narrative	Company public response	Company	State	ZIP code	Tags	Consumer consent provided?	Submitted via	Date sent to company
07/29/2013	Consumer Loan	Vehicle loan	Managing the loan or lease	NaN	NaN	NaN	Wells Fargo & Company	VA	24540	NaN	NaN	Phone	07/30/2013
07/29/2013	Bank account or service	Checking account	Using a debit or ATM card	NaN	NaN	NaN	Wells Fargo & Company	CA	95992	Older American	NaN	Web	07/31/2013
07/29/2013	Bank account or service	Checking account	Account opening, closing, or management	NaN	NaN	NaN	Santander Bank US	NY	10065	NaN	NaN	Fax	07/31/2013
07/29/2013	Bank account or service	Checking account	Deposits and withdrawals	NaN	NaN	NaN	Wells Fargo & Company	GA	30084	NaN	NaN	Web	07/30/2013
07/29/2013	Mortgage	Conventional fixed mortgage	Loan servicing, payments, escrow account	NaN	NaN	NaN	Franklin Credit Management	CT	06106	NaN	NaN	Web	07/30/2013

By using the function df1.describe(), we get the following outcome:

Complaint ID

```
count    6.705980e+05
mean     1.145473e+06
std      6.468630e+05
min      1.000000e+00
25%      5.930575e+05
50%      1.174558e+06
75%      1.718337e+06
max      2.218987e+06
```

This obviously does not give a lot of information about the dataset that can be used to do analysis. Thus, profiling is needed.

Profiling Report Before Data Cleansing/Wrangling

We have used the following code to get the profiling report:

```
profile = ProfileReport(df1, title = 'Consumer Complaints Report')
```

We then have the option to either view the report within Jupyter Notebook or in the form of a separate document by using:



```
profile_eda.to_notebook_iframe()
```

OR

```
profile_eda.to_file("Post_EDA_Consumer_Complaints_Report.html")
```

Understanding of the Profiling Report

The data has 18 columns with 670,598 **rows**. **23%** of the cells have missing data. Out of these, 15 columns are Categorical (String), 2 are Boolean (true/false) and 1 is Numeric.

Dataset statistics

Number of variables	18
Number of observations	670598
Missing cells	2770590
Missing cells (%)	23.0%
Duplicate rows	0
Duplicate rows (%)	0.0%
Total size in memory	92.1 MiB
Average record size in memory	144.0 B

Variable types

Categorical	15
Boolean	2
Numeric	1



The Dataset has multiple columns with high cardinality. This basically means that these columns have many distinct values, but they are not necessarily unique throughout the dataset.

Another Observation that can be made is that multiple columns also have null or NaN values.

“Complaint ID” is the only column with unique values where none of the rows have are missing values.

Warnings

Date received	has a high cardinality: 1818 distinct values	High cardinality
Issue	has a high cardinality: 95 distinct values	High cardinality
Sub-issue	has a high cardinality: 68 distinct values	High cardinality
Consumer complaint narrative	has a high cardinality: 112690 distinct values	High cardinality
Company	has a high cardinality: 3933 distinct values	High cardinality
State	has a high cardinality: 62 distinct values	High cardinality
ZIP code	has a high cardinality: 27889 distinct values	High cardinality
Date sent to company	has a high cardinality: 1767 distinct values	High cardinality
Submitted via	is highly correlated with Consumer consent provided?	High correlation
Consumer consent provided?	is highly correlated with Submitted via	High correlation
Product	is highly correlated with Sub-issue and 1 other fields	High correlation
Sub-issue	is highly correlated with Product and 1 other fields	High correlation
Issue	is highly correlated with Product and 1 other fields	High correlation
Sub-product	has 198202 (29.6%) missing values	Missing
Sub-issue	has 400730 (59.8%) missing values	Missing
Consumer complaint narrative	has 555894 (82.9%) missing values	Missing
Company public response	has 525401 (78.3%) missing values	Missing
Tags	has 575868 (85.9%) missing values	Missing
Consumer consent provided?	has 462447 (69.0%) missing values	Missing
Consumer disputed?	has 41419 (6.2%) missing values	Missing
Consumer complaint narrative	is uniformly distributed	Uniform
Complaint ID	has unique values	Unique

Analysis of each column:

- Date received: This column defines the date at which the user complaint was filed. This has no missing data but has high cardinality.

Date received	Distinct	1818	08/27/2015	963
Categorical	Distinct (%)	0.3%	06/26/2014	916
HIGH CARDINALITY	Missing	0	07/06/2016	914
	Missing (%)	0.0%	08/26/2015	912
	Memory size	5.1 MiB	09/20/2016	909
			Other values (1813)	665984



- Product: This column defines the product the complaint was filed. This has no missing data. It contains 12 distinct values and has high correlation with Sub-issue and one other field.

Product	Distinct	12	Mortgage	210324
Categorical	Distinct (%)	< 0.1%	Debt collection	124236
HIGH CORRELATION	Missing	0	Credit reporting	119195
	Missing (%)	0.0%	Credit card	79007
	Memory size	5.1 MiB	Bank account or service	76084
			Other values (7)	61752

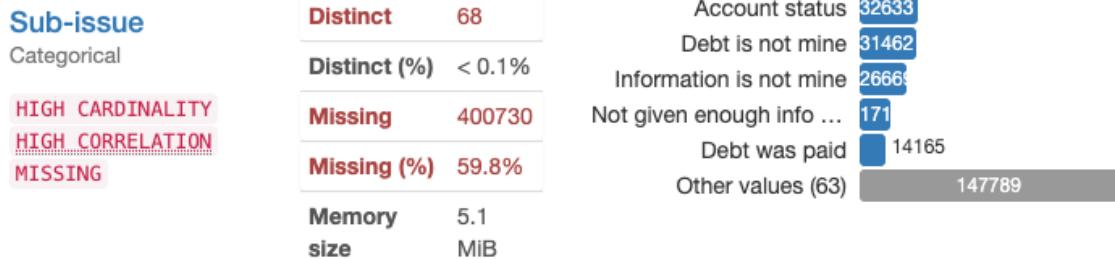
- Sub-product: This column is a sub column for the product category and defines the product on a granular level.

Sub-product	Distinct	47	Other mortgage	81715
Categorical	Distinct (%)	< 0.1%	Conventional fixed mo...	65036
MISSING	Missing	198202	Checking account	53016
	Missing (%)	29.6%	Other (i.e. phone, healt...	3717
	Memory size	5.1 MiB	I do not know	260
			Other values (42)	209435

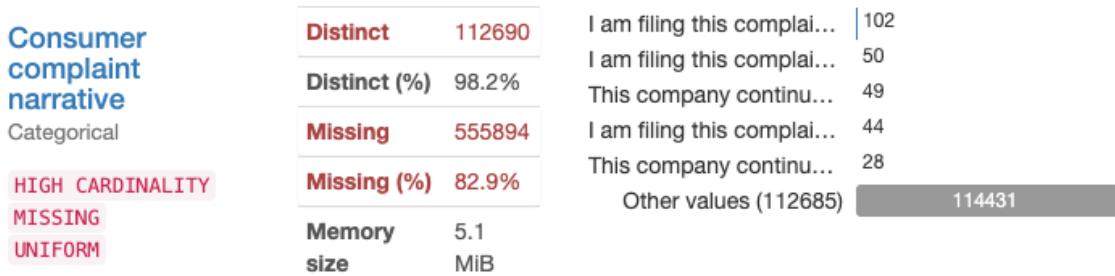
- Issue: This column defines the issue for which the complaint was filed. This has no missing data but has multiple distinct values and it also high correlated to 2 columns.

Issue	Distinct	95	Loan modification,coll...	106455
Categorical	Distinct (%)	< 0.1%	Incorrect information o...	86904
HIGH CARDINALITY	Missing	0	Loan servicing, payme...	70166
HIGH CORRELATION	Missing (%)	0.0%	Cont'd attempts collec...	5159
	Memory size	5.1 MiB	Account opening, closi...	333
			Other values (90)	322136

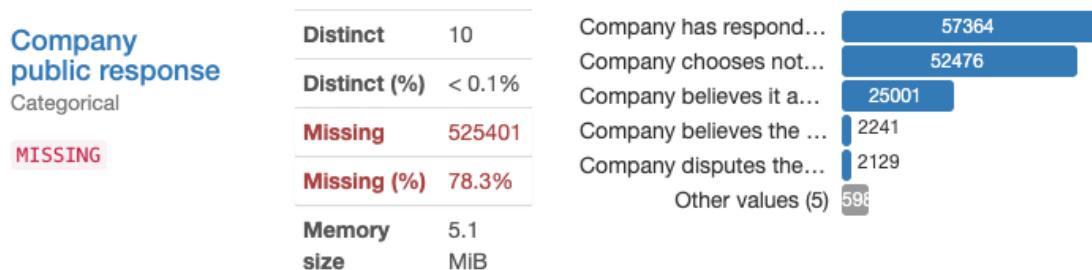
- Sub-issue: This column is a secondary column to the issue column and defines the issue on a granular level.



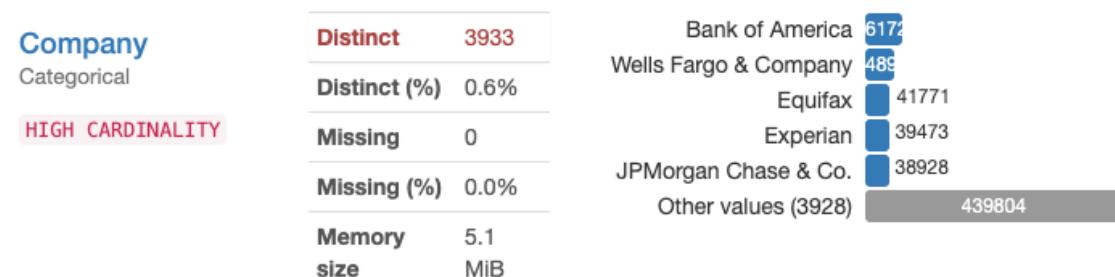
- Consumer complaint narrative: This column is the textual representation of the customers issue and has majority of the rows missing but has a uniform distribution of values.



- Company public response: This column defines the public response the company gave for the complaint filed. It has 78.3% values missing.

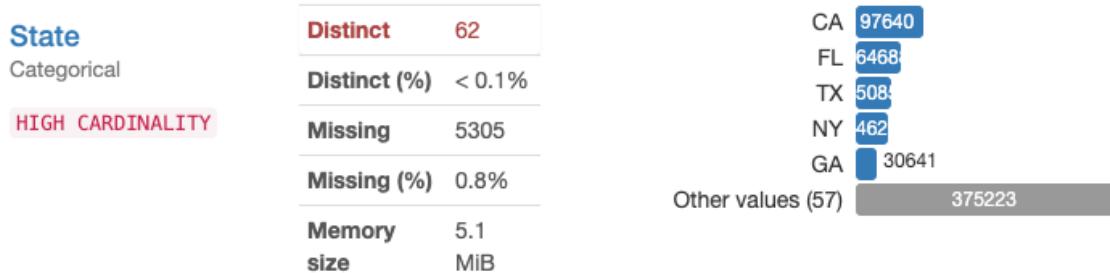


- Company: This column gives the name of the company against whom the complaint is filed. This a no missing values but more than 3900 distinct values.

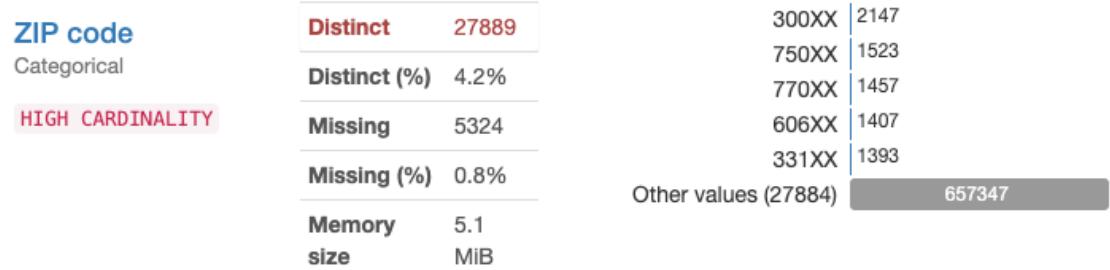




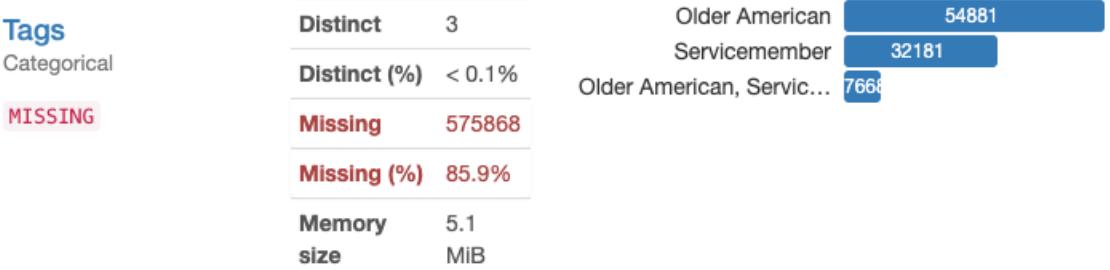
- State: Defines the US state the complaint was filed in. This column has less than 1% of the values missing.



- Zip code: This column defines the zipcode of the complaint. This is a sub column of the state and defines it on a granular level. This should have been a numeric but due to certain columns having incorrect values, this column is not the most consistent.

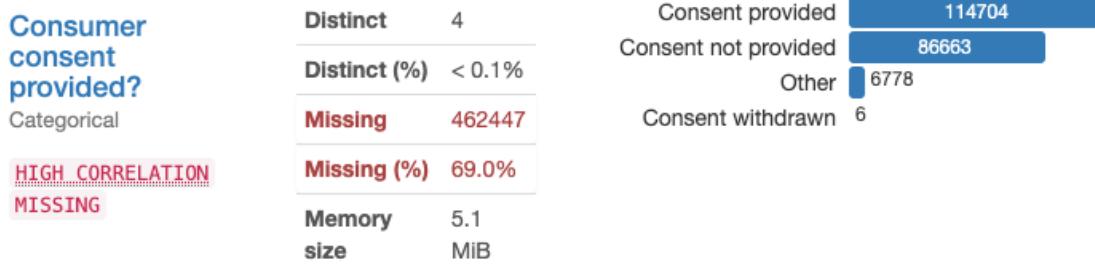


- Tags: This column defines the tags given to the complaint if any. Majority of rows have no value.

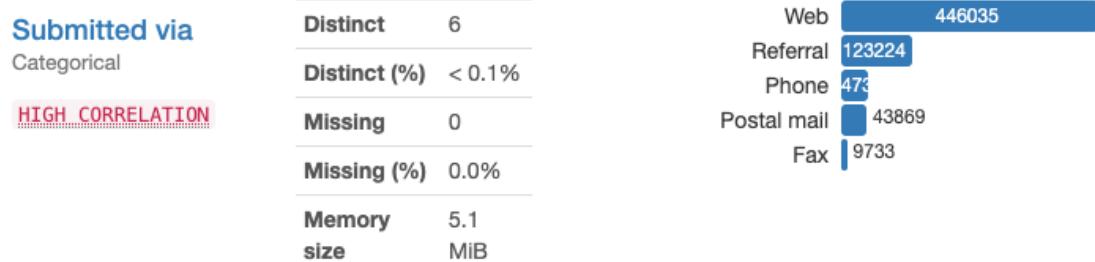




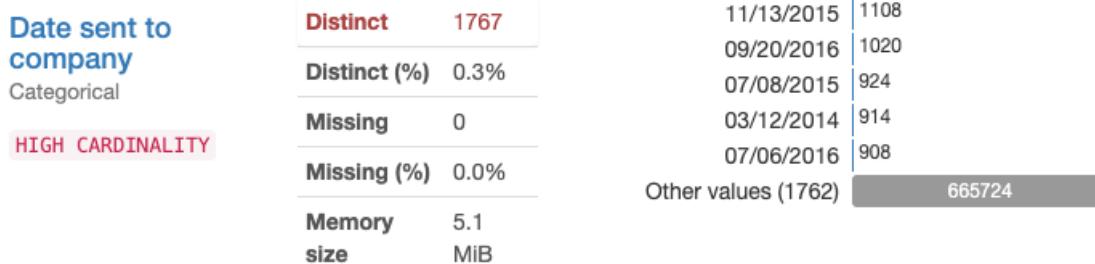
- Consumer consent provided: This column says whether consent was given to access information.



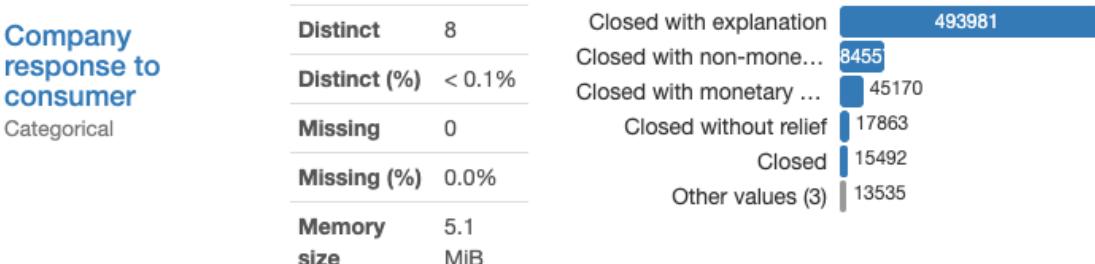
- Submitted via: This column defines the way the complaint was filed.



- Date sent to company: This tells the date the complaint was sent to the company.



- Company response to consumer: This column defines the response that the company gave to the consumer.





- **Timely_response:** This column shows whether the response was timely.

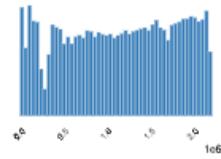
Timely response?	Distinct	2	True	652848
Boolean	Distinct (%)	< 0.1%	False	17750
	Missing	0		
	Missing (%)	0.0%		
	Memory size	655.0 KiB		

- **Consumer disputed:** This column describes the outcome of the resolution or the complaint.

Consumer disputed?	Distinct	2	False	496466
Boolean	Distinct (%)	< 0.1%	True	132713
	Missing	41419	(Missing)	41419
	Missing (%)	6.2%		
	Memory size	1.3 MiB		

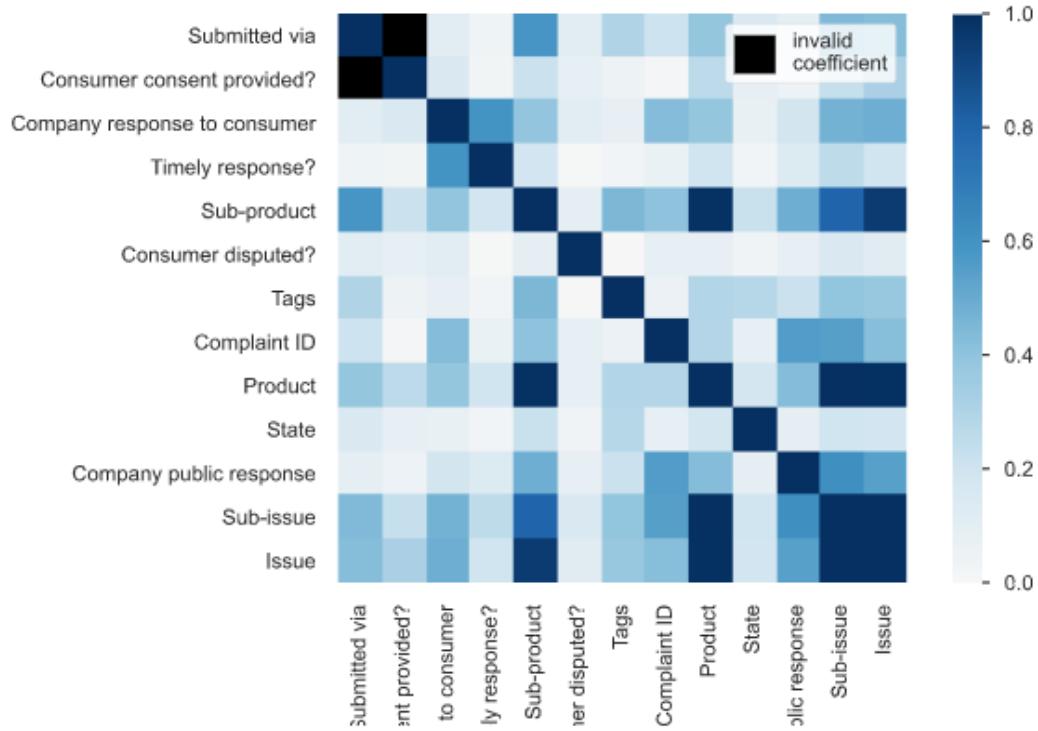
- **Complaint ID:** This column defines the compliant id if the compliant.

Complaint ID	Distinct	670598	Mean	1145473.241
Real number ($\mathbb{R}_{\geq 0}$)	Distinct (%)	100.0%	Minimum	1
	Missing	0	Maximum	2218987
	Missing (%)	0.0%	Zeros	0
	Infinite	0	Zeros (%)	0.0%
	Infinite (%)	0.0%	Memory size	5.1 MiB

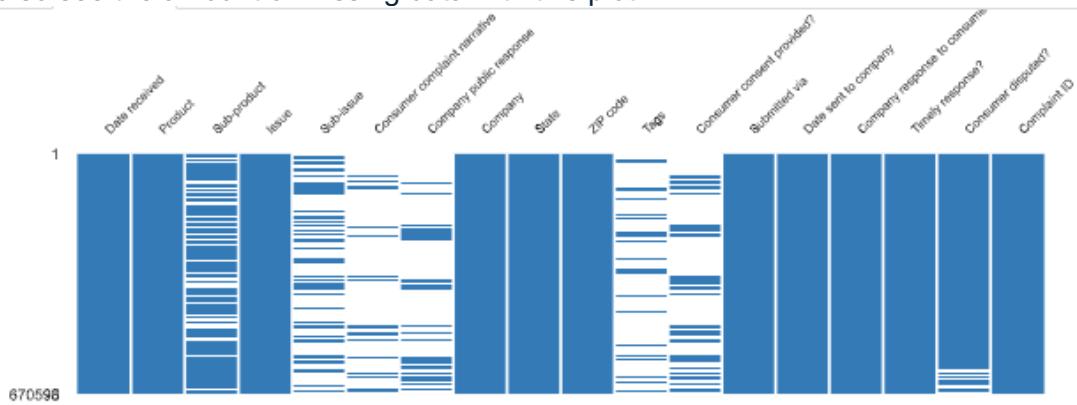




The data was then checked for correlation and this the matrix that was formed. This shows the strong dependency between some columns but also the invalid coefficient error. We try to solve this by cleansing this data.



We also see the amount of missing data with this plot.





We then perform some EDA and data cleansing to improve this data. The following steps were performed to fix the data.

- Renaming Columns to remove spaces between words.

```
# Renaming Columns

df2.rename(columns = {'Date received' : 'Date_Received', 'Product' : 'Product', 'Sub-product' : 'Sub_Product',
                     'Issue' : 'Issue', 'Sub-issue' : 'Sub_Issue',
                     'Consumer complaint narrative' : 'Consumer_Complaint_Narrative',
                     'Company public response' : 'Company_Public_Response', 'Company' : 'Company',
                     'State' : 'State', 'ZIP code' : 'Zip_Code', 'Tags' : 'Tags',
                     'Consumer consent provided?' : 'Consumer_Consent_Provided',
                     'Submitted via' : 'Submitted_Via', 'Date sent to company' : 'Data_Sent_To_Company',
                     'Company response to consumer' : 'Company_Response_To_Consumer',
                     'Timely response?' : 'Timely_Response', 'Consumer disputed?' : 'Consumer_Disputed',
                     'Complaint ID' : 'Complaint_ID'}, inplace = True)
```

- Filling empty cells with 'unknown'

```
1 df3.Sub_Product = df2.Sub_Product.fillna('Unknown')
2 df3.Sub_Issue = df2.Sub_Issue.fillna('Unknown')
3 df3.Consumer_Complaint_Narrative = df2.Consumer_Complaint_Narrative.fillna('Unknown')
4 df3.Company_Public_Response = df2.Company_Public_Response.fillna('Unknown')
5 df3.Tags = df2.Tags.fillna('No_Tags')
6 df3.Consumer_Consent_Provided = df2.Consumer_Consent_Provided.fillna('Unknown')
7 df3.Consumer_Disputed = df2.Consumer_Disputed.fillna('Unknown')
8 df3.State = df2.State.fillna('Unknown')
9 df3.Zip_Code = df2.Zip_Code.fillna('00000')
```

- Date_Received and Data_Sent_To_Company are converted to datetime datatype. Multiple values are replaced in different columns to make the dataset as consistent as possible.

```
1 df3['Date_Received'] = pd.to_datetime(df3['Date_Received'])
2 df3['Data_Sent_To_Company'] = pd.to_datetime(df3['Data_Sent_To_Company'])
3 df3['Zip_Code'] = df3['Zip_Code'].str.replace('X', '0')
4 df3['Zip_Code'] = df3['Zip_Code'].str.replace('-', '0')
5 df3['Zip_Code'] = df3['Zip_Code'].str.replace('~', '0')
6 df3['Zip_Code'] = df3['Zip_Code'].str.replace('[', '0')
7 df3['Zip_Code'] = df3['Zip_Code'].str.replace(']', '0')
8 df3['Zip_Code'] = df3['Zip_Code'].str.replace('\'', '0')
9 df3['Zip_Code'] = df3['Zip_Code'].str.replace('+', '0')
10 df3['Zip_Code'] = df3['Zip_Code'].str.replace('.','0')
11 df3['Zip_Code'] = df3['Zip_Code'].str.replace('/','0')
12 df3['Zip_Code'] = df3['Zip_Code'].str.replace('$','0')
13 df3['Zip_Code'] = df3['Zip_Code'].str.replace('!','0')
14 df3['Zip_Code'] = df3['Zip_Code'].str.replace('*','0')
15 df3['Product'] = df3['Product'].str.replace(' ', '_')
16 df3['Sub_Product'] = df3['Sub_Product'].str.replace(' ', '_')
17 df3['Issue'] = df3['Issue'].str.replace(' ', '_')
18 df3['Sub_Issue'] = df3['Sub_Issue'].str.replace(' ', '_')
19 df3['Consumer_Complaint_Narrative'] = df3['Consumer_Complaint_Narrative'].str.replace(' ', '_')
20 df3['Company_Public_Response'] = df3['Company_Public_Response'].str.replace(' ', '_')
21 df3['Company'] = df3['Company'].str.replace(' ', '_')
22 df3['Tags'] = df3['Tags'].str.replace(' ', '_')
23 df3['Consumer_Consent_Provided'] = df3['Consumer_Consent_Provided'].str.replace(' ', '_')
24 df3['Company_Response_To_Consumer'] = df3['Company_Response_To_Consumer'].str.replace(' ', '_')
25 df3
```



- Reordering the columns.

```
1 df4 = df3[['Complaint_ID', 'Company', 'Date_Received', 'Data_Sent_To_Company', 'Product',
2   'Sub_Product', 'Issue', 'Sub_Issue', 'Submitted_Via', 'Consumer_Complaint_Narrative',
3   'State', 'Zip_Code', 'Tags', 'Company_Public_Response', 'Company_Response_To_Consumer',
4   'Timely_Response', 'Consumer_Disputed', 'Consumer_Consent_Provided']]
```

- Zip code was converted into int datatype.

```
1 df4['Zip_Code'] = df4['Zip_Code'].astype(int)
2 df4
```

- Finally, the new dataset is downloaded as a csv.

```
df4.to_csv('Post_EDA.csv', index = False)
```

`df4.shape`

(670598, 18)

- Our final dataset after cleansing.

Complaint_ID	Company	Date_Received	Data_Sent_To_Company	Product	Sub_Product
468882	Wells_Fargo_&_Company	2013-07-29	2013-07-30	Consumer_Loan	Vehicle_loan
468889	Wells_Fargo_&_Company	2013-07-29	2013-07-31	Bank_account_or_service	Checking_account
468879	Santander_Bank_US	2013-07-29	2013-07-31	Bank_account_or_service	Checking_account
468949	Wells_Fargo_&_Company	2013-07-29	2013-07-30	Bank_account_or_service	Checking_account
475823	Franklin_Credit_Management	2013-07-29	2013-07-30	Mortgage	Conventional_fixed_mortgage
...
2211891	Equifax	2016-11-16	2016-11-16	Credit_reporting	Unknown
2179768	Citibank	2016-10-26	2016-10-26	Mortgage	Conventional_adjustable_mortgage_(ARM)
2212094	Equifax	2016-11-17	2016-11-17	Credit_reporting	Unknown
2126003	Amex	2016-09-22	2016-09-22	Credit_card	Unknown
2062419	BancorpSouth_Bank	2016-08-13	2016-09-07	Consumer_Loan	Installment_loan



Post EDA the report was taken, and this is the result of the process.

- Missing values have been fixed.

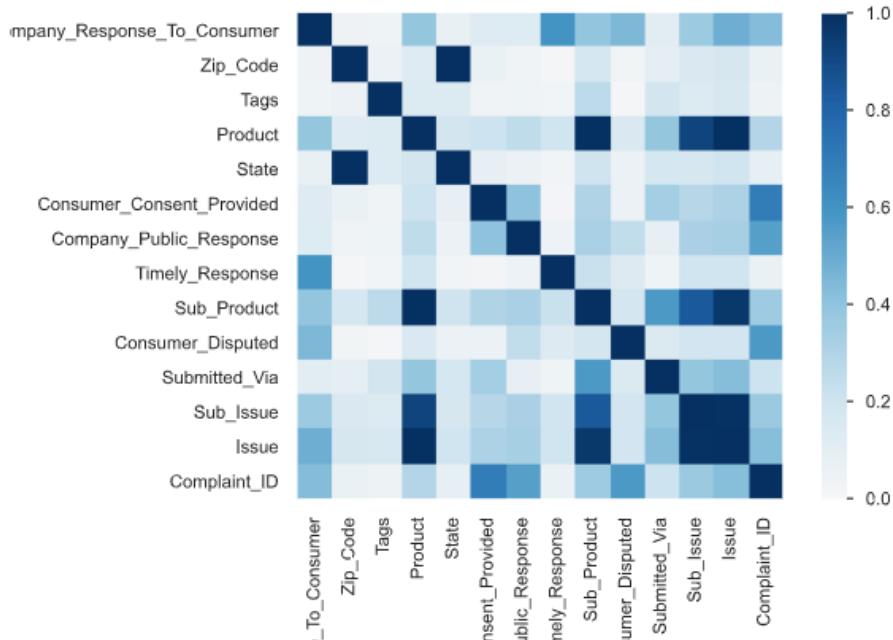
Dataset statistics

Number of variables	18
Number of observations	670598
Missing cells	0
Missing cells (%)	0.0%
Duplicate rows	0
Duplicate rows (%)	0.0%
Total size in memory	92.1 MiB
Average record size in memory	144.0 B

Variable types

Numeric	2
Categorical	15
Boolean	1

- The correlation has improved greatly. We do not have any invalid correlations anymore.





4. Database Installation

Neo4j

Neo4j is a NoSQL Graph database that uses Cypher Query Language to perform operations on the data that is being loaded.

The first step in Neo4j is to create a new Project where the Database can be added. We created a project called “CSYE_PROJ”.

Adding the database is the next step. We have Created a local DBMS “Consumer_Complaints” and then the dataset is added to the import folder of Neo4j.

The screenshot shows the Neo4j Neo4jDB interface. At the top, there is a header bar with the project name "CSYE_PROJ" on the left and a "Add" button with a plus sign on the right. Below the header, there is a list of databases. One database, "Consumer_Complaints 4.2.5", is highlighted and marked as "ACTIVE". To the right of the database list are three buttons: "Stop" (red), "Open" (blue), and "...". Below the database list, there is a sidebar with a "File" section containing a file icon and the name "Post_EDA.csv". To the right of the sidebar, there is a "Filename" dropdown menu.



5. Data Mapping and Integration

Data Mapping begins with the creation of constraints.

```
1 CREATE CONSTRAINT ON (complaint:Complaint) ASSERT complaint.ComplaintId IS UNIQUE;
2 CREATE CONSTRAINT ON (company:Company) ASSERT company.name IS UNIQUE;
3 CREATE CONSTRAINT ON (response:Response) ASSERT response.companyResponse IS UNIQUE;
4 CREATE CONSTRAINT ON (product:Product) ASSERT product.ProductName IS UNIQUE;
5 CREATE CONSTRAINT ON (state:State) ASSERT state.StateName IS UNIQUE;
6 CREATE CONSTRAINT ON (zipcode:Zipcode) ASSERT zipcode.ZipId IS UNIQUE;
7 CREATE CONSTRAINT ON (subProduct:SubProduct) ASSERT subProduct.SubProductName IS UNIQUE;
8 CREATE CONSTRAINT ON (issue:Issue) ASSERT issue.IssueName IS UNIQUE;
9 CREATE CONSTRAINT ON (subIssue:SubIssue) ASSERT subIssue.SubIssueName IS UNIQUE;
```

Creating Nodes Complaint and Company with their Relationship.

```
1 //Create Node Complaint
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MERGE (complaint:Complaint {ComplaintId:toInteger(row.Complaint_ID)})
5 ON CREATE SET complaint.DateReceived = row.Date_Received;

1 //Create Node Company
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MERGE (company:Company {name:row.Company});

1 //Create Relationship complaint-against-company
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MATCH (complaint:Complaint {ComplaintId:toInteger(row.Complaint_ID)})
5 MATCH (company:Company {name:row.Company})
6 MERGE (complaint)-[:AGAINST]→(company);
```

Creating Node Response and its Relationship with Company.

```
1 //Create Node Response
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MERGE (response:Response {companyResponse:row.Company_Response_To_Consumer});

1 //Create Relationship response-from-company
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MATCH (company:Company {name:row.Company})
5 MATCH (response:Response {companyResponse:row.Company_Response_To_Consumer})
6 MERGE (response)-[:FROM]→(company);
```

Creating Nodes Product and Sub-Product and their Relationships with Complaint and each other.

```
1 //Create Node Product
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MERGE (product:Product {ProductName:row.Product});

1 //Create Relationship complaint-about-product
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MATCH (complaint:Complaint {ComplaintId:toInteger(row.Complaint_ID)})
5 MATCH (product:Product {ProductName:row.Product})
6 MERGE (complaint)-[:ABOUT]→(product);

1 //Create Node SubProduct
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MERGE (subProduct:SubProduct {SubProductName:row.Sub_Product});

1 //Create Relationship complaint-about-subproduct
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MATCH (complaint:Complaint {ComplaintId:toInteger(row.Complaint_ID)})
5 MATCH (subProduct:SubProduct {SubProductName:row.Sub_Product})
6 MERGE (complaint)-[:ABOUT]→(subProduct);

1 //Create Relationship product-contains-subproduct
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MATCH (product:Product {ProductName:row.Product})
5 MATCH (subProduct:SubProduct {SubProductName:row.Sub_Product})
6 MERGE (product)-[:CONTAINS]→(subProduct);
```



Creating Node State and Zip Code.

```
1 //Create Node State
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MERGE (state:State {StateName:row.State});
1 //Create Node Zip
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MERGE (zipcode:Zipcode {ZipId:toInteger(row.Zip_Code)});
```

Creating Relationship between State & Complaint and Zip Code & Complaint.

```
1 //Create Relationship complaint-originatedfrom-state
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MATCH (complaint:Complaint {ComplaintId:toInteger(row.Complaint_ID)})
5 MATCH (state:State {StateName:row.State})
6 MERGE (complaint)-[:ORIGINATED_FROM]→(state);
1 //Create Relationship complaint-originatedfrom-zip
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MATCH (complaint:Complaint {ComplaintId:toInteger(row.Complaint_ID)})
5 MATCH (zipcode:Zipcode {ZipId:toInteger(row.Zip_Code)})
6 MERGE (complaint)-[:ORIGINATED_FROM]→(zipcode);
```

Creating Nodes Issue & Sub Issue and their Relationship with Complaint.

```
1 //Create Node Issue
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MERGE (issue:Issue {IssueName:row.Issue});
1 //Create Relationship complaint-with-issue
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MATCH (complaint:Complaint {ComplaintId:toInteger(row.Complaint_ID)})
5 MATCH (issue:Issue {IssueName:row.Issue})
6 MERGE (complaint)-[:WITH]→(issue);
1 //Create Node SubIssue
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MERGE (subIssue:SubIssue {SubIssueName:row.Sub_Issue});
1 //Create Relationship complaint-with-subissue
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MATCH (complaint:Complaint {ComplaintId:toInteger(row.Complaint_ID)})
5 MATCH (subIssue:SubIssue {SubIssueName:row.Sub_Issue})
6 MERGE (complaint)-[:WITH]→(subIssue);
```

Creating Relationship between Issue and Sub Issue.

```
1 //Create Relationship subissue-incategory-issue
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MATCH (issue:Issue {IssueName:row.Issue})
5 MATCH (subIssue:SubIssue {SubIssueName:row.Sub_Issue})
6 MERGE (subIssue)-[:IN_CATEGORY]→(issue);
```

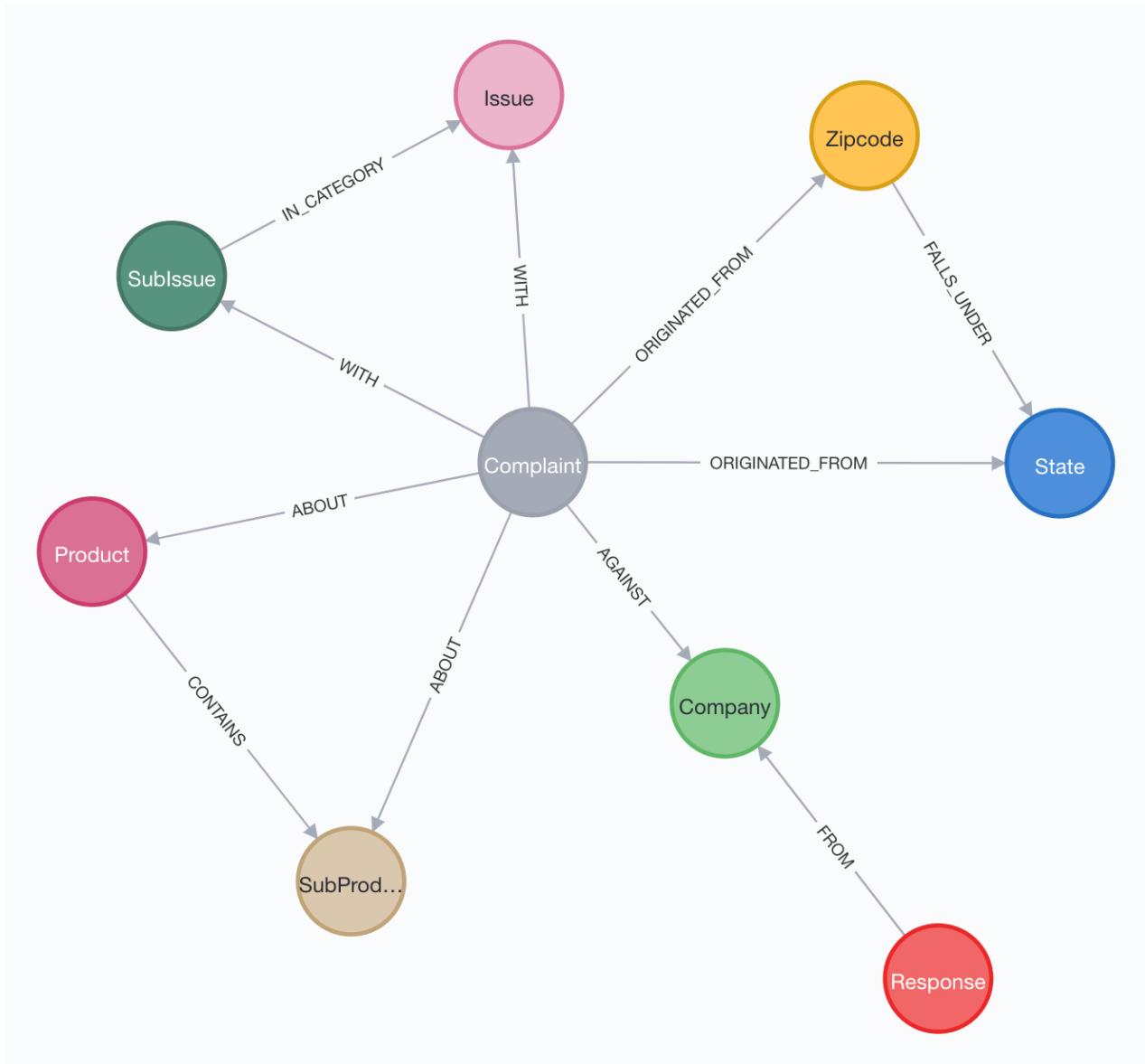
Creating Relationship between Zip Code and State.

```
1 //Create Relationship Zip-fallsunder-State
2 :auto USING PERIODIC COMMIT 500
3 LOAD CSV WITH HEADERS FROM 'file:///Post_EDA.csv' AS row
4 MATCH (state:State {StateName:row.State})
5 MATCH (zipcode:Zipcode {ZipId:toInteger(row.Zip_Code)})
6 MERGE (zipcode)-[:FALLS_UNDER]→(state);
```



Schema Visualization.

```
call db.schema.visualization();
```

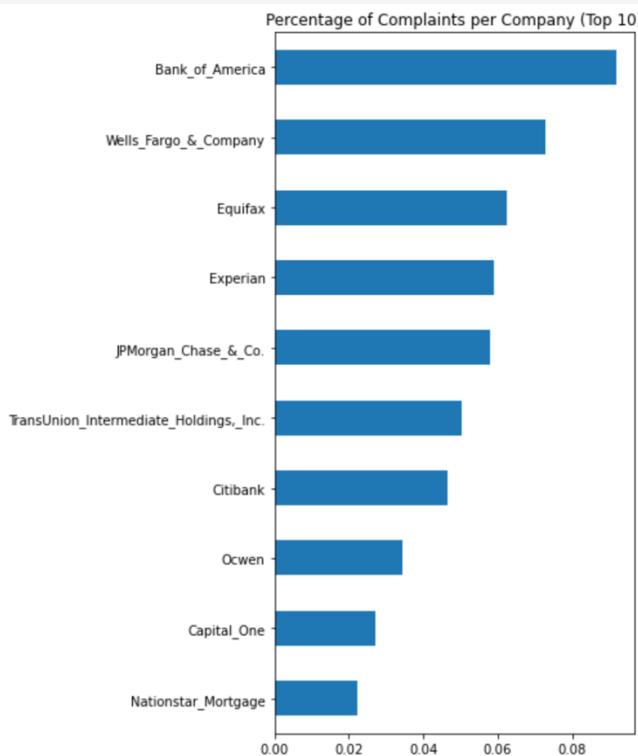




6. Data Validation and Data Visualization

The bar chart shows the top 10 companies that have received the maximum amount of complaints

```
fig_dims = (5, 10)
fig, ax = plt.subplots(figsize=fig_dims)
df1['Company'].value_counts(normalize=True)[-10:].plot(ax=ax, kind='barh')
ax.set_title('Percentage of Complaints per Company (Top 10)')
```

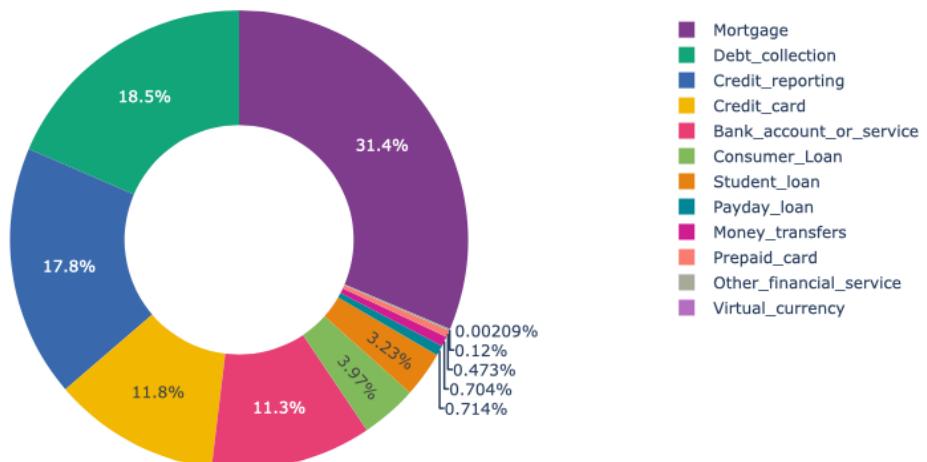




This donut plot displays the product the complaint is for.

```
pie_chart_1 = px.pie(data_frame= dfdict, values = 'Number_Of_Complaints', names = 'Product',
                      title = 'Percentage of Products that receive complaints',
                      color_discrete_sequence=px.colors.qualitative.Bold, hole=0.5)
pie_chart_1.update_layout(showlegend=True)
```

Percentage of Products that receive complaints





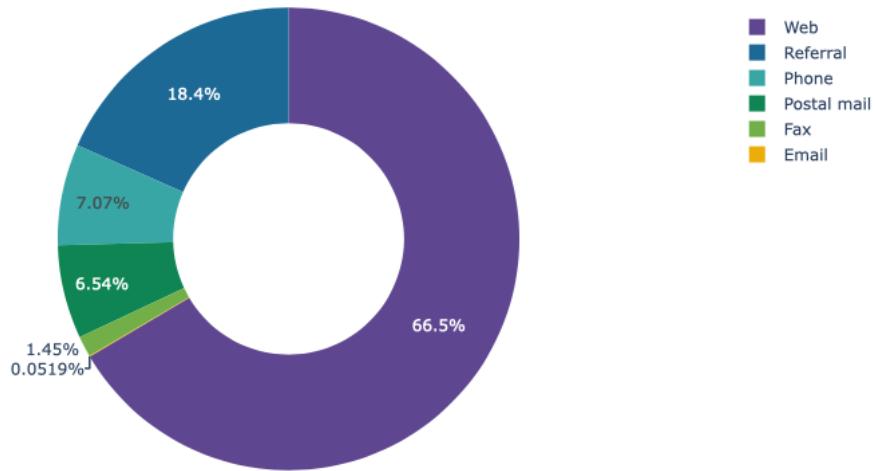
This plot shows the way the complaint was filed.

```
df_via = df1.groupby(["Submitted_Via"]).size().reset_index(name="Number_of_Complaints")
df_via.head()
```

Submitted_Via	Number_of_Complaints	
0	Email	348
1	Fax	9733
2	Phone	47389
3	Postal mail	43869
4	Referral	123224

```
pie_chart_2 = px.pie(data_frame= df_via, values = 'Number_of_Complaints', names = 'Submitted_Via',
                      title = 'Medium of Submitting Complaints',color_discrete_sequence=px.colors.qualitative.Prism,
                      hole=0.5)
pie_chart_2.update_layout(showlegend=True)
```

Medium of Submitting Complaints





This bar chart displays the timely vs untimely response graph.

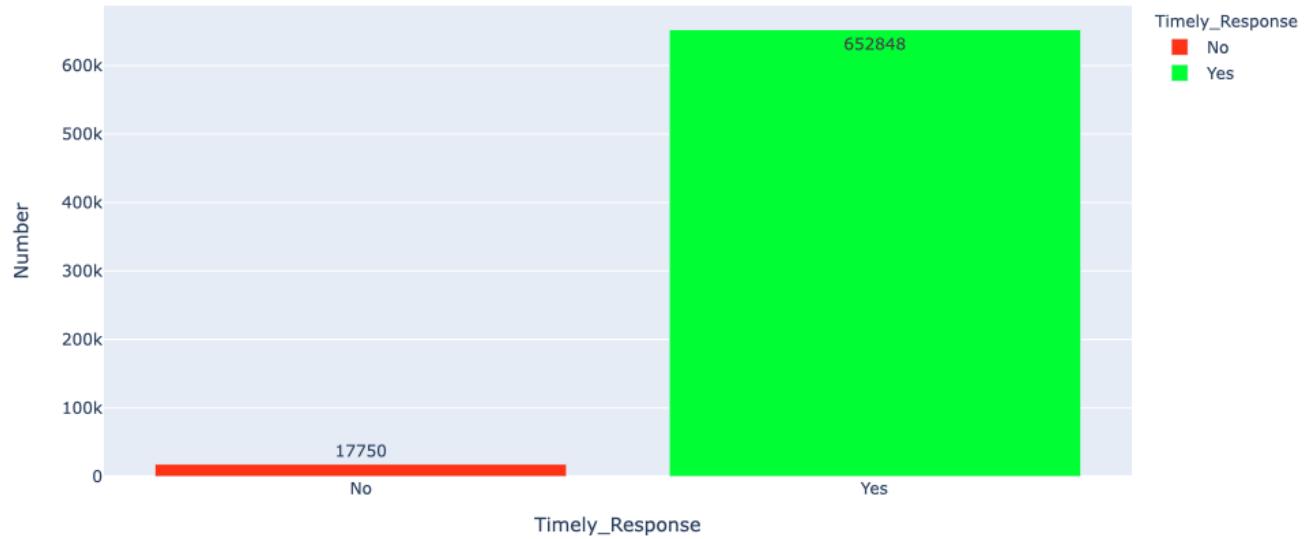
```
timely_resp_dict = df1.groupby(["Timely_Response"]).size().reset_index(name="Number")
timely_resp_dict.head()
```

Timely_Response	Number
0	No 17750
1	Yes 652848

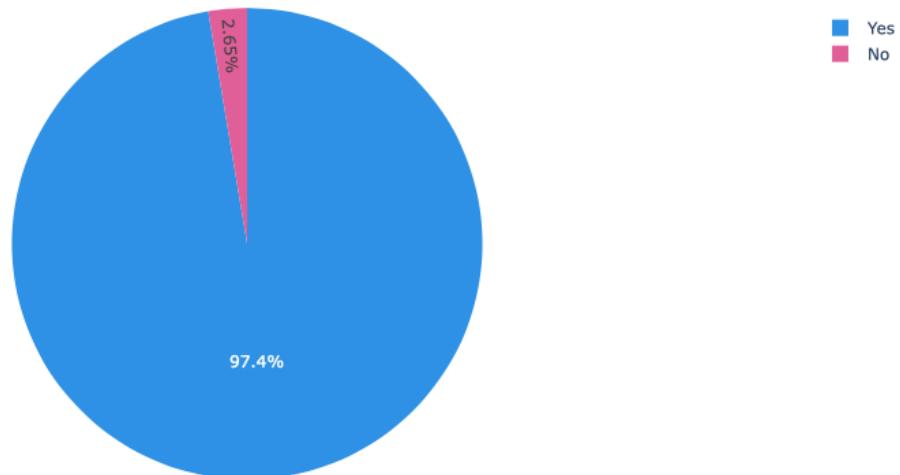


```
bar_chart_1 = px.bar(data_frame = timely_resp_dict, x = 'Timely_Response', y = 'Number', color = 'Timely_Response',
                      title = 'Number Of Times A Timely Response Was Given',text = 'Number',
                      color_discrete_sequence=px.colors.qualitative.Light24)
bar_chart_1
```

Number Of Times A Timely Response Was Given



Comparing Percentage of Timely Responses





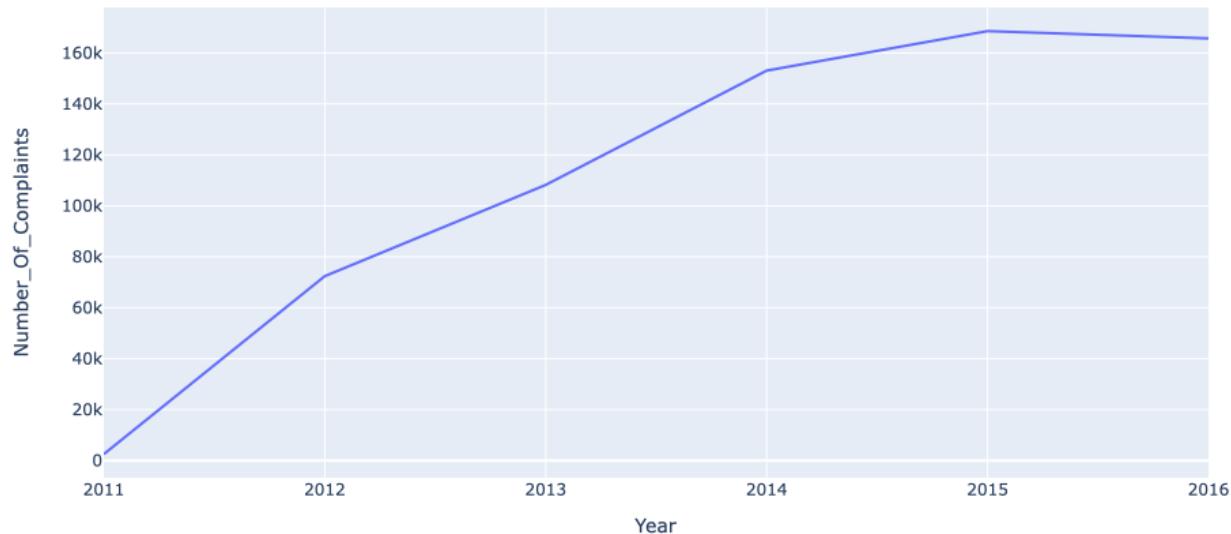
This table show the cumulative complaints over the years.

```
df_date = df1.groupby(["Year"]).size().reset_index(name="Number_Of_Complaints")
df_date.head()
```

Year	Number_Of_Complaints
0	2540
1	72403
2	108226
3	153088
4	168600

```
lineplot = px.line(data_frame= df_date, x ='Year', y = 'Number_Of_Complaints',
                     title = 'Number Of Complaints Received Every Year')
lineplot.update_layout(showlegend=True)
```

Number Of Complaints Received Every Year





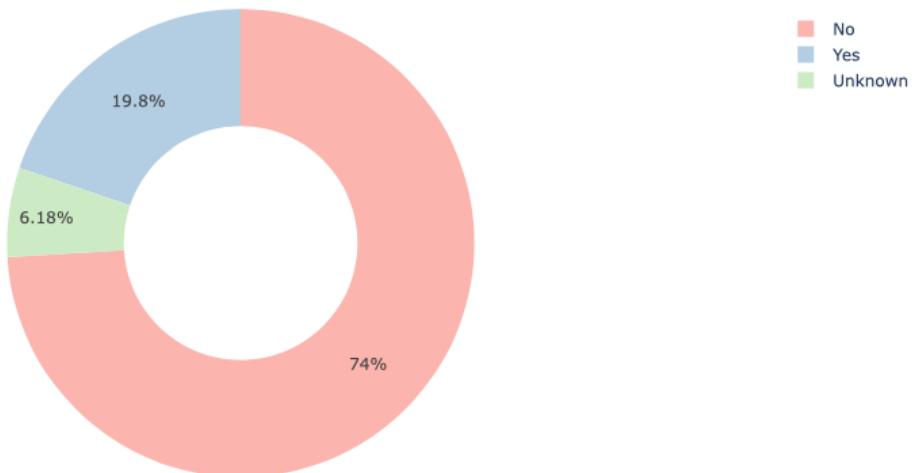
This plot shows the dispute denial percentage

```
df_disputed = df1.groupby(["Consumer_Disputed"]).size().reset_index(name="Number_Of_Complaints")
df_disputed.head()
```

Consumer_Disputed	Number_Of_Complaints
0	No 496466
1	Unknown 41419
2	Yes 132713

```
pie_chart_4 = px.pie(data_frame=df_disputed, values='Number_Of_Complaints', names='Consumer_Disputed',
                      title='Dispute Denial Percentage', color_discrete_sequence=px.colors.qualitative.Pastel1,
                      hole=0.5)
pie_chart_4.update_layout(showlegend=True)
```

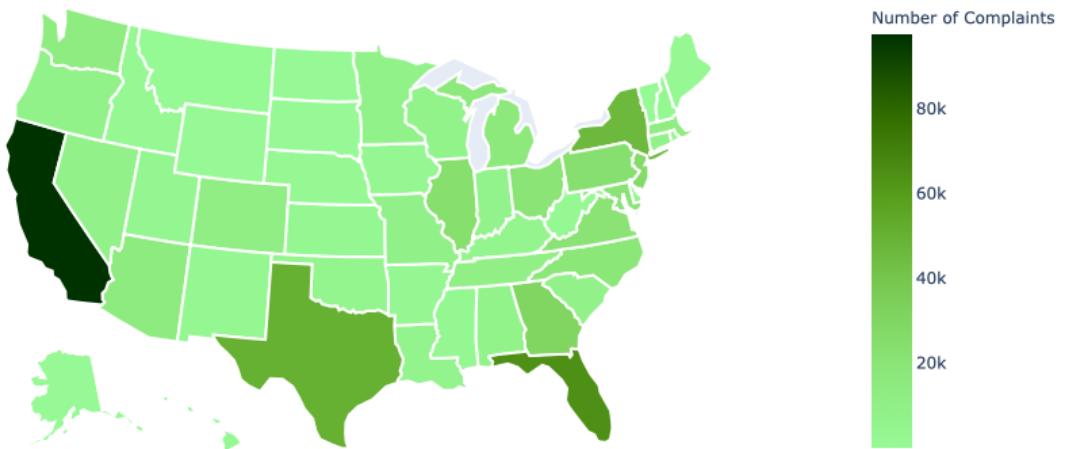
Dispute Denial Percentage





This plot shows the map for most complaints filed from state.

Number of Complaints by State





7. System Integration and User Acceptance Testing

System integration was tested prior to ingesting data in the database.

The following piece of code was run to validate system connection. The py2neo is the package. It is open source and available to download using the PIP package manager.

The neo4j is running on localhost and the authentication is passed as parameter.

```
from py2neo import Graph
graph = Graph("bolt://localhost:7687", auth=("neo4j", "neo4j"))
try:
    graph.run("Match () Return 1 Limit 1")
    print('Connection successful')
except Exception:
    print('Connection unsuccessful')
```

We connect to the database and send a query which if runs without error means the database is connected.

Upon no error we print **Connection Successful**

If the code returns some error, we say the **Connection Unsuccessful**

The output of the above code concludes system integration testing.

For integration testing, we test the contents of both the input file and the data inserted in the database.

For Data validation we check the count of the rows of data in the neo4j database and the python script used for EDA.

The following is the output for the neo4j row count query.

The match statement is used to match all the rows in the Complaint column, and we return the count of rows.

```
MATCH (:Complaint)
RETURN count(*) as count
```

count
1 670598



```
1 df1= pd.read_csv ("Post_EDA.csv")
2 df1.shape[0]
```

670598

Count of all rows was cross verified with the input csv file using Python.

This concludes integration testing.

User acceptance testing:

User Acceptance Testing, also known as UAT testing means **A process of verifying that a solution works for the user.**

The ideal workflow will ideally be a user entering the data and it will successfully get inserted into the database.

The issue that could arise in this is the data is either invalid or has inconsistencies.

To validate the user can check the data that is to be inserted in the database using the python script to count the number of rows. Then we use a second script to check the amount of rows in the neo4j.

If the output of both the query and the python code match then we can conclude the application is working perfectly and this concludes user acceptance testing.



8. Challenges Encountered

The following are the challenges we encountered

- Data mapping
 - Here we had to collectively understand what the data provided to us contained.
 - We had to understand the relationship between the columns.
- Data visualization
 - For visualization, a new package had to be installed to build interactive graphs and plots
 - Zip code and State relationship had to be built to plot the data on an interactive map
- Neo4j
 - The concepts and inner working of this new database had a learning curve.
 - The working knowledge of how nodes and relations are built in Neo4j had to be figured out.
 - New syntax had to be understood to efficiently query the database and insert data into it.



9. End User Instructions

The end user for our project will be the **Consumer Financial Protection Bureau or CFPB**. The dataset and the dashboard created by us will be helpful for them to protect people from getting into trouble due to the financial decision that they make.

Database Creation and loading

- Open Neo4j and create a project, followed by creating a local database.
- Load the Post_EDA.csv file to the import folder.
- Start and then open the database.
- Create constraints followed by nodes and relationships.
- Visualize the schema.

Visualization interpretation

Following is the dashboard that can be used by the user:



The dashboard can be interpreted as:

- A personal from CFPB can look at the dashboard to identify the companies that have been receiving the maximum number of complaints and set up an investigation against them to improve the quality of services provided.



- One can also see the products that need careful consideration before investing the money into. Since, high number of complaints is not a good sign.
- The states from where, the complaints filed are maximum can be evaluated using the map. Further, the state government rules can be modified if the companies are exploiting the consumers.
- There are a few companies that have denied any dispute on the complaints being filed. The CFPB can get in touch with them to identify the problems and act if necessary.
- The dashboard shows that the number of complaints has been increasing every subsequent year. This can be an indication that there might be something wrong with the policies that allow companies to make rules that the consumer is not happy with.