

# AI\_Phase3

Serial Number	Team Member Name	Registration Number
1	K.Navinraj	310821104064
2	M.P.Praveen Raja	310821104070
3	Sachin A	310821104081
4	K.P.Tharun	310821104100

## Objective:

The technology stack described in the plan encompasses a range of tools and frameworks that are commonly used in data integration, preprocessing, modelling, real-time data processing, visualisation, deployment, monitoring, documentation, training, compliance, evaluation, and scaling. Let's elaborate on some of the key technologies and their purposes:

## 1.Dataset:

DatasetLink: [Company Master Data of Tamil Nadu upto 28th February 2019 | Open Government Data \(OGD\) Platform India](#)

	A	B	C	D	E	F	G	
1	CORPORATE_IDENTIFICATION	COMPANY_NAME	COMPANY_STATUS	COMPANY_CLASS	COMPANY_CATEGORY	COMPANY_SUB_CATEGORI	DATE_OF_REGISTRATION	
2	F00643	HOCHTIEFF AG,	NAEF	NA	NA	NA	01/12/1961	Ta
3	F00721	SUMITOMO CORPORATION (SUMITOMO SHOJI KAISHA LIMITED)	ACTV	NA	NA	NA	NA	Ta
4	F00892	SRILANKAN AIRLINES LIMITED	ACTV	NA	NA	NA	01/03/1982	Ta
5	F01208	CALTEX INDIA LIMITED	NAEF	NA	NA	NA	NA	Ta
6	F01218	GE HEALTHCARE BIO-SCIENCES LIMITED	ACTV	NA	NA	NA	NA	Ta
7	F01265	CAIRN ENERGY INDIA PTY. LIMITED	NAEF	NA	NA	NA	NA	Ta
8	F01269	TORIELLI S.R.L	ACTV	NA	NA	NA	05/09/1995	Ta
9	F01311	HARDY EXPLORATION & PRODUCTION (INDIA) INC..	ACTV	NA	NA	NA	NA	Ta
10	F01314	HOCHTIOF AKTIENGESELLSCHAFT VORM GFBH HELFMANN	ACTV	NA	NA	NA	11/04/1996	Ta
11	F01412	EPSON SINGAPORE PVT LTD	ACTV	NA	NA	NA	25/04/1997	Ta
12	F01426	CARGOLUX AIRLINES INTERNATIONAL S A	ACTV	NA	NA	NA	11/06/1997	Ta
13	F01468	CHO HEUNG ELECTRIC INDUSTRIAL COMPANY LIMITED	NAEF	NA	NA	NA	NA	Ta
14	F01543	NYCOMED ASIA PACIFIC PTE LIMITED	ACTV	NA	NA	NA	27/10/1998	Ta
15	F01544	CHERRINGTON ASIA LTD	ACTV	NA	NA	NA	01/05/2000	Ta
16	F01563	SHIMADZU ASIA PACIFIC PTE LIMITED	NAEF	NA	NA	NA	NA	Ta
17	F01565	CORK INTERNATIONAL PTY LIMITED	ACTV	NA	NA	NA	NA	Ta
18	F01566	ERBIS ENGG COMPANY LIMITED	ACTV	NA	NA	NA	NA	Ta
19	F01589	RALF SCHNEIDER HOLDING GMBH	NAEF	NA	NA	NA	NA	Ta
20	F01593	MITRAJAYA TRADING PRIVATE LIMITED	ACTV	NA	NA	NA	NA	Ta
21	F01618	HEAT AND CONTROL PTY LIMITED	ACTV	NA	NA	NA	13/07/1999	Ta
22	F01628	DIREX SYSTEMS LIMITED	ACTV	NA	NA	NA	NA	Ta
23	F01641	NMB-MINEBEA THAI LIMITED	NAEF	NA	NA	NA	NA	Ta
24	F01643	ARROW INTERNATIONAL INC	ACTV	NA	NA	NA	02/11/1999	Ta
25	F01694	GAMBRO CHINA LTD	ACTV	NA	NA	NA	14/06/2000	Ta
26	F01703	OBARA CORPORATION	NAEF	NA	NA	NA	17/07/2000	Ta
27	F01752	CIPTA WAWASON MAJU ENGINEERING SDM BHD	ACTV	NA	NA	NA	24/01/2001	Ta
28	F01753	AUCHAN INTERNATIONAL S.A.	ACTV	NA	NA	NA	NA	Ta
29	F01767	TOSHIBA PLANT SYSTEMS AND SERVICES CORPORATION	NAEF	NA	NA	NA	08/03/2001	Ta
30	F01768	YAMAZEN CORPORATION	NAEF	NA	NA	NA	NA	Ta
31	F01770	OMAL INTERNATIONAL PTE LTD	ACTV	NA	NA	NA	22/03/2004	Ta

+

≡

Data\_Gov\_Tamil\_Nadu.csv

Sheet1

modified\_dataset

I	J	K	L	M	N	O	P	Q
AUTHORIZED_CAP	PAIDUP_CAPITAL	INDUSTRIAL_C	PRINCIPAL_BU	REGISTERED	REGISTRAR	EMAIL_ADDR	LATEST_YEAR_ANNUAL_RETURN	LATEST_YEAR_FINANCIAL_STATEMENT
0		0 NA	Agriculture & allii AMBLE SIDE, N ROC	ELHI	NA	NA	NA	NA
0		0 NA	Agriculture & allii FLAT NO. 6, 1st ROC	ELHI	shuchi.chug@as	NA	NA	NA
0		0 NA	Agriculture & allii SRILANKAN AIF ROC	ELHI	shree16us@yahoo	NA	NA	NA
0		0 NA	Agriculture & allii GOLD CREST 2 ROC	ELHI	NA	NA	NA	NA
0		0 NA	Agriculture & allii FF-3 Palani Cer ROC	ELHI	karthick9999@y	NA	NA	NA
0		0 NA	Agriculture & allii WELLINGTON I ROC	ELHI	neerja.sharma@	NA	NA	NA
0		0 NA	Agriculture & allii 6, Mangayarkar ROC	ELHI	chennai@torii	NA	NA	NA
0		0 NA	Agriculture & allii 5TH FLOOR,WE ROC	ELHI	venkatesh.v@he	NA	NA	NA
0		0 NA	Agriculture & allii NEW NO.86, OL ROC	ELHI	kumar@internat	NA	NA	NA
0		0 NA	Agriculture & allii 7C CEATURY P ROC	ELHI	NA	NA	NA	NA
0		0 NA	Agriculture & allii OFFICE NO 91A ROC	ELHI	NA	NA	NA	NA
0		0 NA	Agriculture & allii 129, MANPUR V ROC	ELHI	chowelaccounts	NA	NA	NA
0		0 NA	Agriculture & allii A D 46 1ST STI ROC	ELHI	NA	NA	NA	NA
0		0 NA	Agriculture & allii 10HADDOWS R ROC	ELHI	NA	NA	NA	NA
0		0 NA	Agriculture & allii FIRST FLOOR, I ROC	ELHI	kousik@vsni.co	NA	NA	NA
0		0 NA	Agriculture & allii ARJAY APEX CI ROC	ELHI	NA	NA	NA	NA
0		0 NA	Agriculture & allii 39,2nd Main Roi ROC	ELHI	NA	NA	NA	NA
0		0 NA	Agriculture & allii FLAT C, 'SAI VA ROC	ELHI	NA	NA	NA	NA
0		0 NA	Agriculture & allii OLD NO 148 NE ROC	ELHI	NA	NA	NA	NA
0		0 NA	Agriculture & allii A40 OLD NO 26 ROC	ELHI	ncrajagopal@gn	NA	NA	NA
0		0 NA	Agriculture & allii F-1, FIRST FLO ROC	ELHI	direx@vsni.com	NA	NA	NA
0		0 NA	Agriculture & allii Level - 2 Regus, ROC	ELHI	stsogawa@mine	NA	NA	NA
0		0 NA	Agriculture & allii BLUE HAVEN, h ROC	ELHI	NA	NA	NA	NA
0		0 NA	Agriculture & allii 5 1ST FLOOR IS ROC	ELHI	NA	NA	NA	NA
0		0 NA	Agriculture & allii INDIA BRANCH ROC	ELHI	joe@obara.co.in	NA	NA	NA
0		0 NA	Agriculture & allii 141 AVVAI SHAI ROC	ELHI	NA	NA	NA	NA
0		0 NA	Agriculture & allii RK Tower, No. 1 ROC	ELHI	pverma@vkvern	NA	NA	NA
0		0 NA	Agriculture & allii HOTEL AMBAS ROC	ELHI	NA	NA	NA	NA
0		0 NA	Agriculture & allii PLOT 69, SIVAN ROC	ELHI	NA	NA	NA	NA
0		0 NA	Agriculture & allii NO 4 SADDHAR ROC	ELHI	NA	NA	NA	NA

## 2. Required libraries

### Importing required libraries

```
In [52]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

## 3. Loading the Dataset

We will load the dataset named “dialogs.txt” inside our python notebook using the `pd.read_csv()` method. To load the dataset for AI-driven exploration and prediction of company registration trends with the Registrar of Companies (RoC), you'll need to follow a series of steps. In this content guide, I'll outline the process for loading the dataset and preparing it for analysis using artificial intelligence (AI) techniques.

### Loading the dataset.

```
In [65]: #Reading the dataset from CSV format.
df = pd.read_csv('/Users/sachinanandharaj/Downloads/Data_Gov_Tamil_Nadu.csv', low_memory=False)
```

```
In [66]: #Printing the given Dataset
df
```

```
Out[66]:
```

	CORPORATE_IDENTIFICATION_NUMBER	COMPANY_NAME	COMPANY_STATUS	COMPANY_CLASS	COMPANY_CATEGORY	COMPANY_SUB_CATEGORY
0	F00643	HOCHTIEFF AG,	NAEF	NaN	NaN	NaN
1	F00721	SUMITOMO CORPORATION (SUMITOMO SHOJI KAISHA L...	ACTV	NaN	NaN	NaN
2	F00892	SRILANKAN AIRLINES LIMITED	ACTV	NaN	NaN	NaN
3	F01208	CALTEX INDIA LIMITED	NAEF	NaN	NaN	NaN
4	F01218	GE HEALTHCARE BIO-SCIENCES LIMITED	ACTV	NaN	NaN	NaN
...	...	...	...	...	...	...
150866	U74997TN2016PTC112556	QUAD42 MEDIA PRIVATE LIMITED	ACTV	Private	Company limited by Shares	Non-govt company
150867	U74997TN2018PTC121491	IYERAATHU FOODS PRIVATE LIMITED	ACTV	Private	Company limited by Shares	Non-govt company
150868	U74997TZ2016PTC027802	POLYGAR FARM SOLUTIONS PRIVATE LIMITED	STOF	Private	Company limited by Shares	Non-govt company
150869	U74997TZ2018PTC030177	PANDIYA AGRRI SOLUTIONS PRIVATE LIMITED	ACTV	Private	Company limited by Shares	Non-govt company
150870	U74997TZ2019PTC032491	NROOT TECHNOLOGIES PRIVATE LIMITED	ACTV	Private	Company limited by Shares	Non-govt company

150871 rows x 7 columns

## 4. Pre-processing

Preprocessing is a crucial step in AI-driven exploration and prediction of company registration trends with the Registrar of Companies (RoC). Proper preprocessing ensures that your dataset is clean, well-structured, and suitable for machine learning algorithms.

a often contains noise, irrelevant information, and inconsistencies that can interfere with the accuracy of NLP models.

### Preprocessing

```
In [60]: #Cleaning the dataset by removing the NA value rows from the dataset
df = df.dropna()
```

```
In [61]: #Printing the dataset
df
```

## 5. Data Type Correction

Out [61]:

	CORPORATE_IDENTIFICATION_NUMBER	COMPANY_NAME	COMPANY_STATUS	COMPANY_CLASS	COMPANY_CATEGORY	COMPANY_SUB_CATEGORY
310	L01117TZ1943PLC000117	NEELAMALAI AGRO INDUSTRIES LIMITED	ACTV	Public	Company limited by Shares	Non-govt company
311	L01119TN1986PLC013473	ABAN OFFSHORE LIMITED	ACTV	Public	Company limited by Shares	Non-govt company
313	L01119TN1992PLC024076	SOFTECH INFINIUM SOLUTIONS LIMITED	ACTV	Public	Company limited by Shares	Non-govt company
315	L01122TZ1995PLC010762	POCHIRAJU INDUSTRIES LIMITED	ACTV	Public	Company limited by Shares	Non-govt company
318	L01132TZ1922PLC000234	THE UNITED NILGIRI TEA ESTATES COMPANY LIMITED	ACTV	Public	Company limited by Shares	Non-govt company
...	...	...	...	...	...	...
150862	U74997TN2016PTC112105	MRKR COMMUNICATIONS PRIVATE LIMITED	ACTV	Private	Company limited by Shares	Non-govt company
150864	U74997TN2016PTC112257	ETHNICINDIAN FASHION RETAIL PRIVATE LIMITED	ACTV	Private	Company limited by Shares	Non-govt company
150865	U74997TN2016PTC112312	SAVIDYA EDUCATION PRIVATE LIMITED	ACTV	Private	Company limited by Shares	Non-govt company
150866	U74997TN2016PTC112556	QUAD42 MEDIA PRIVATE LIMITED	ACTV	Private	Company limited by Shares	Non-govt company
150869	U74997TZ2018PTC030177	PANDIYA AGRI SOLUTIONS PRIVATE LIMITED	ACTV	Private	Company limited by Shares	Non-govt company

73739 rows x 7 columns

	CORPORATE_IDENTIFICATION_NUMBER	COMPANY_NAME	COMPANY_STATUS	COMPANY_CLASS	COMPANY_CATEGORY	COMPANY_SUB_CATEGORY
...	...	...	...	...	...	...
150862	U74997TN2016PTC112105	MRKR COMMUNICATIONS PRIVATE LIMITED	ACTV	Private	Company limited by Shares	Non-govt company
150864	U74997TN2016PTC112257	ETHNICINDIAN FASHION RETAIL PRIVATE LIMITED	ACTV	Private	Company limited by Shares	Non-govt company
150865	U74997TN2016PTC112312	SAVIDYA EDUCATION PRIVATE LIMITED	ACTV	Private	Company limited by Shares	Non-govt company
150866	U74997TN2016PTC112556	QUAD42 MEDIA PRIVATE LIMITED	ACTV	Private	Company limited by Shares	Non-govt company
150869	U74997TZ2018PTC030177	PANDIYA AGRI SOLUTIONS PRIVATE LIMITED	ACTV	Private	Company limited by Shares	Non-govt company

73739 rows x 7 columns

### Data Type Correction

```
In [57]: df['INDUSTRIAL_CLASS'].astype('int32')
```

```
Out [57]: 310      1117
311      1119
313      1119
315      1122
318      1132
...
150862   74997
150864   74997
150865   74997
150866   74997
150869   74997
Name: INDUSTRIAL_CLASS, Length: 73739, dtype: int32
```

```
In [58]: df['LATEST_YEAR_ANNUAL_RETURN'] = pd.to_datetime(df['LATEST_YEAR_ANNUAL_RETURN'])
```

CIN	REGISTERED_OFFICE_ADDRESS	REGISTRAR_OF_COMPANIES	EMAIL_ADDR	LATEST_YEAR_ANNUAL_RETURN	LATEST_YEAR_FINANCIAL_STATEMENT
Illed	KATARY ESTATEKATARY POSTCOONOR	ROC 纡OIMBATORE	secneelamalai@avtplantations.co.in	2019-03-31	2019-03-31
Illed	'JANPRIYA CREST'96, PANTHEON ROAD,EGMORE	ROC 纡HENNAI	secretarial@aban.com	2019-03-31	2019-03-31
Illed	29, PRECISION PLAZA, NEW 397, ANNA SALAITEYNAM...	ROC 纡HENNAI	complianceofficer@softtechinfinium.com	2018-03-31	2018-03-31
Illed	1/102 SATHYAMANGALAM VILLAGEHOSUR TALUK	ROC 纡OIMBATORE	mmreddyandco@gmail.com	2019-03-31	2019-03-31
Illed	3 SAVITHRI SHANMUGHAM ROADRACE COURSE	ROC 纡OIMBATORE	headoffice@chamrajtea.com	2019-03-31	2019-03-31

## 6. Exploratory Data Analysis

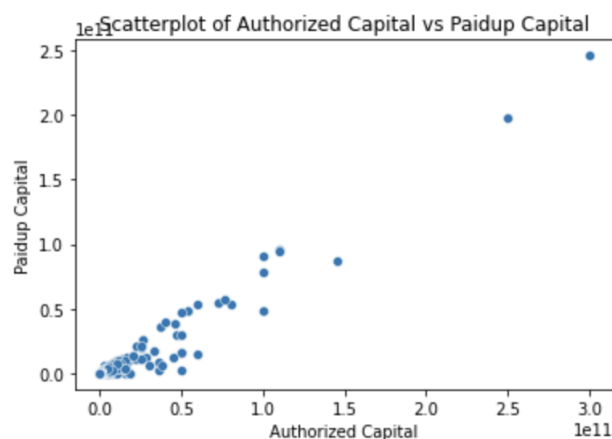
Exploratory data analysis (EDA) is used to analyses and investigate data sets and summarise their main characteristics, often employing data visualisation methods

AI-driven exploration and prediction of company registration trends with the Registrar of Companies (RoC) involves a series of steps, including data collection, cleaning, transformation, loading, EDA, modelling, and deployment. Careful attention to data quality and preprocessing is essential for building effective AI models in this context.

### Exploratory Data Analysis

```
In [40]: sns.scatterplot(x='AUTHORIZED_CAP', y='PAIDUP_CAPITAL', data= ds)
plt.title('Scatterplot of Authorized Capital vs Paidup Capital')
plt.xlabel('Authorized Capital')
plt.ylabel('Paidup Capital')
```

```
Out[40]: Text(0, 0.5, 'Paidup Capital')
```

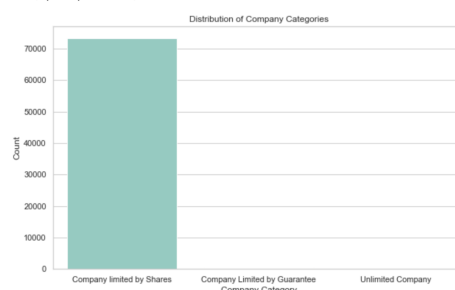


```
In [41]: sns.set(style='whitegrid') # Customize the plot style
plt.figure(figsize=(10, 6)) # Set the figure size

# Create the count plot for the COMPANY_CATEGORY column
sns.countplot(x='COMPANY_CATEGORY', data=ds, palette='Set3') # Adjust the color palette if needed

# Customize the plot (optional)
plt.title('Distribution of Company Categories')
plt.xlabel('Company Category')
plt.ylabel('Count')
```

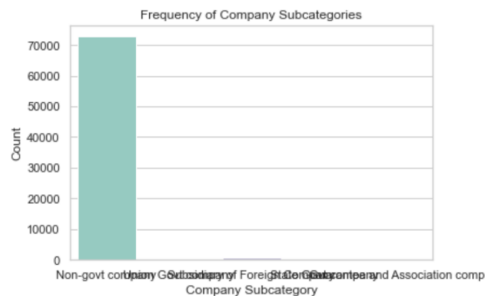
```
Out[41]: Text(0, 0.5, 'Count')
```



```
In [43]: ## Create the count plot for the COMPANY_SUB_CATEGORY column
sns.countplot(x='COMPANY_SUB_CATEGORY', data=ds, palette='Set3') # Adjust the color palette if needed

## Customize the plot
plt.title('Frequency of Company Subcategories')
plt.xlabel('Company Subcategory')
plt.ylabel('Count')
```

Out[43]: Text(0, 0.5, 'Count')

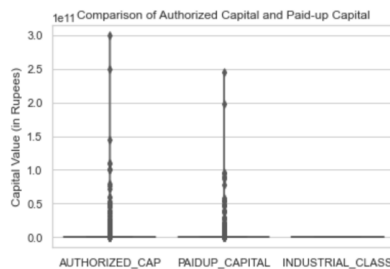


```
In [44]: sns.boxplot(data=ds, palette='Set2')

# Add a violin plot on top of the box plot for better visualization
sns.violinplot(data=ds, palette='Set3', inner=None)

# Customize the plot (optional)
plt.title('Comparison of Authorized Capital and Paid-up Capital')
plt.ylabel('Capital Value (in Rupees)')
```

Out[44]: Text(0, 0.5, 'Capital Value (in Rupees)')



## 6. Extracting preprocessed csv file dataset

After performing data cleaning, the refined dataset can be saved and downloaded as a distinct CSV file.

### Saving the Preprocessed dataset

```
In [16]: df.to_csv('modified_dataset.csv', index=False)
```

```
In [18]: ds = pd.read_csv('modified_dataset.csv')
```

## 7. Preprocessing Code

# Importing required libraries

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```

# Loading the dataset.
df = pd.read_csv('/Users/sachinanandharaj/Downloads/Data_Gov_Tamil_Nadu.csv',
low_memory=False)
print(df)
# Preprocessing - Removing NA values and changing the data type
df = df.dropna()
print(df)
# Data Type Correction
df['INDUSTRIAL_CLASS'].astype('int32')
df['LATEST_YEAR_ANNUAL_RETURN'] =
pd.to_datetime(df['LATEST_YEAR_ANNUAL_RETURN'])
df['LATEST_YEAR_FINANCIAL_STATEMENT'] =
pd.to_datetime(df['LATEST_YEAR_FINANCIAL_STATEMENT'])
df['INDUSTRIAL_CLASS'] = df['INDUSTRIAL_CLASS'].astype('int32')

# Exploratory Data Analysis
sns.scatterplot(x='AUTHORIZED_CAP', y='PAIDUP_CAPITAL', data= ds)
plt.title('Scatterplot of Authorized Capital vs Paidup Capital')
plt.xlabel('Authorized Capital')
plt.ylabel('Paidup Capital')
sns.set(style='whitegrid') # Customize the plot style
plt.figure(figsize=(10, 6)) # Set the figure size

# Create the count plot for the COMPANY_CATEGORY column
sns.countplot(x='COMPANY_CATEGORY', data=ds, palette='Set3') # Adjust the color
palette if needed

# Customise the plot
plt.title('Distribution of Company Categories')
plt.xlabel('Company Category')
plt.ylabel('Count')

# Create the count plot for the COMPANY_SUB_CATEGORY column
sns.countplot(x='COMPANY_SUB_CATEGORY', data=ds, palette='Set3')
plt.title('Frequency of Company Subcategories')
plt.xlabel('Company Subcategory')
plt.ylabel('Count')
sns.boxplot(data=ds, palette='Set2')

# Add a violin plot on top of the box plot for better visualization
sns.violinplot(data=ds, palette='Set3', inner=None)
plt.title('Comparison of Authorized Capital and Paid-up Capital')
plt.ylabel('Capital Value (in Rupees)')

# Saving the Preprocessed dataset
df.to_csv('modified_dataset.csv', index=False)
ds = pd.read_csv('modified_dataset.csv')

```