National College of
Ireland

# Configuration Manual

MSc Research Project
Data Analytics

# Chandrashekar Gettam Rajgopal

Student ID: X21226075

School of Computing
National College of Ireland

Supervisor:     Mayank Jain

| Student Name: | Chandrashekar Gettam Rajgopal |
|---|---|
| Student ID: | X21226075 |
| Programme: | Data Analytics |
| Year: | 2023 |
| Module: | MSc Research Project |
| Supervisor: | Mayank Jain |
| Submission Due Date: | 14/12/2023 |
| Project Title: | Configuration Manual |
| Word Count: | 1468 |
| Page Count: | 10 |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

**ALL** internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

| Signature: | Chandrashekar Gettam Rajgopal |
|---|---|
| Date: | 29th January 2024 |

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:**

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies). | ☑ |
| **Attach a Moodle submission receipt of the online project submission**, to each project (including multiple copies). | ☑ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | ☑ |

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Configuration Manual

## Chandrashekar Gettam Rajgopal
### X21226075

# 1 Introduction

This configuration manual contains step-by-step instructions for configuring the system and the manual's main goal is to provide comprehensive instructions on how to carry out the research study. It also details the machine setup needed to compile and execute the models, model evaluations, code snippets, and visualizations from exploratory data research. The processes involve installing the necessary programs and packages in addition to the minimal configuration that is advised for a project to succeed.

# 2 System Configuration

## 2.1 Hardware Requirements

- Operating System: Google Colab will run on a modern web browser and all Operating systems that can run a modern web browser can be used.

- RAM: Minimum 4GB but recommended 8GB for better performance.

- Processor: Modern processor Intel, AMD or Apple Silicon which can handle multiple tasks and run multiple browser tabs smoothly.

- Storage: Minimum of 256 GB Hard disk or Solid state drive to install the browser and to run.

## 2.2 Software Requirements

- Web Browser: Mozilla Firefox, Google Chrome, Safari, or any other Chromium browser for a good experience.

- Google Drive: Store the dataset in Google Drive and mount it to Google Colab.

- Microsoft Excel: To do initial exploration.

- Tableau: Used for initial exploration and exploratory data analysis.

- Python: Any version of Python 3.x

# 3 Code snippets

## 3.1 Data Collection

Data is downloaded from Kaggle using this link Investments VC Dataset Arindam235 (2023). The dataset contains an investments_VC.csv file that has information on various companies investment details.

## 3.2 Setup Google Colab

Store the downloaded investments_VC.csv file in Google Drive and mount the Google Drive into Google Colab as shown in Figure 1.
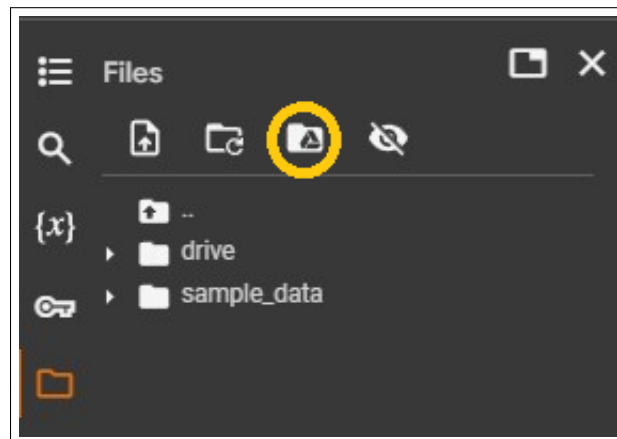


Figure 1: Google Drive mount in Google Colab

Choosing the run time of Google Colab Figure 2. Choose Python for Runtime type and CPU can be chosen by default and click on save.



Figure 2: Google Colab Runtime

## 3.3 Install necessary Packages and Libraries

The next step is installing Pandas, Scikit-learn, Deap, gplearn, geopandas and plotly. packages as shown in Figure 3.



Figure 3: Packages

Next step is importing all necessary libraries as shown in Figure 4

- `matplotlib.pyplot as plt`: Matplotlib is a popular library for creating static, animated, and interactive visualizations in Python.

- `numpy as np`: NumPy is used for numerical operations and provides support for arrays and matrices.

- `pandas as pd`: Pandas is used for data manipulation and analysis, particularly for working with structured data in tabular format.

- `plotly.express as px`: Plotly is a library for creating interactive and visually appealing plots and charts.

- `seaborn as sns`: Seaborn is a high-level interface for creating informative and attractive statistical graphics.

- `matplotlib.dates as mdates`: A module within Matplotlib that provides tools for working with date and time data in plots.

- `deap`: A library for evolutionary algorithms and genetic programming.

3

Figure 4: Importing Libraries

- `gplearn.genetic`: A library for symbolic regression using genetic programming.

- `sklearn.compose`: Part of scikit-learn, used for composing transformers and estimators into pipelines.

- `sklearn.ensemble`: Scikit-learn's ensemble methods library, which includes techniques like random forests.

- `sklearn.impute`: Provides tools for imputing missing data in scikit-learn pipelines.

- `sklearn.linear_model`: Scikit-learn's linear modeling tools, including linear regression.

- `sklearn.manifold`: Part of scikit-learn, used for manifold learning techniques like t-SNE.

- `sklearn.metrics`: Scikit-learn's library for evaluating model performance with various metrics.

- `sklearn.model_selection`: Provides tools for model selection, including cross-validation and train-test splitting.

- `sklearn.neighbors`: Scikit-learn's k-neighbors-based methods, including k-nearest neighbors regression.

- `sklearn.pipeline`: Part of scikit-learn, used for creating machine learning pipelines.

- `sklearn.preprocessing`: Scikit-learn's library for data preprocessing, including label encoding, standard scaling, and one-hot encoding.

- **xgboost**: A library for gradient boosting, often used for predictive modeling.

- **scipy.stats**: Part of SciPy, provides statistical functions and distributions.

- **wordcloud**: A library for creating word clouds, often used for visualizing word frequencies in text data.

## 3.4 Loading the dataset

Loading the dataset and printing the dataset to understand the features present. Figure 5



Figure 5: Loading Dataset

## 3.5 Data preprocessing

Removing features that are not important for predicting the total funding USD. Figure 6 shows these features are removed from the dataset as it is not required for total funding prediction "permalink", "name", "homepage_url" and a few features like "category_list, "founded_month", "founded_quarter", "founded_at" are redundant and repetitive in other features present in the dataset. After this many other preprocessing were performed on the dataset and performed in-depth analysis of the dataset through extensive exploratory data analysis.

Figure 6: Removing unnecessary columns

### 3.5.1 Grouping the market into market group bins to understand as sectors

There is a lot of market present in the dataset and it would be easier to analyze and understand different markets when grouped into the market sector each market sector contains a list of markets that is more relevant to that market sector as shown in figure 7. Market Groups are as below:

- Telecommunications

- Travel and Hospitality

- Others

- Technology and Software

- Marketing and Advertising

- Commerce and Retail

- Health and Wellness

- Real Estate and Construction

- Education and Training

- Food and Beverages

- Media and Entertainment

- Transportation and Mobility

Figure 7: Market to Market Group

- Financial Services

- Energy and Environment

- Manufacturing and Industrial

- Other

### 3.5.2 Handling state city missing data using external mapping

Several state code values were empty but the city value was present. So by using an external mapper that maps the city to its state code values are imputed as shown in figure 8. This sample of code shows several stages in the preprocessing and imputation of data for a dataset. It loads a DataFrame with investment data and a CSV file with state and city information first. Common suffixes like "City" and "County" are removed from the city names to make them more pristine. After that, it updates the investment dataset's missing state codes and builds a mapping from city to state code using the state-city dataset that was loaded. Furthermore, it uses the 'first_funding_at_datetime' column to impute missing values in the 'founded_year' column. 'state_code' and 'city' values that are lacking are further imputed by the code by mapping them to the most prevalent values in individual groupings. Lastly, it reports the count of null or missing values for each column in the dataset and imputes missing "status" values using the mode of the "status" column. These procedures guarantee that the dataset is full and cleaner for the analysis that comes after.

Figure 8: Handling missing data using imputing

### 3.5.3 TSNE

Code as shown in Figure 9 to generate the TSNE which is used for understanding the dataset for dimensionality reduction and viewing the dataset in 2D graph.

- 'city'

- 'founded_year'

- 'venture'

- 'market_group'

- 'state_code'

- 'funding_total_usd'

Above are among the company data subsets that are first chosen and stored in data subset funding bins. After that, this subset is divided into bins according to the 'funding_total_usd' using predetermined ranges; a new column called 'funding_bins' is then created with names for these ranges. The high-dimensional data is then reduced to two dimensions using a t-SNE analysis on the subset, which helps to visualize the data's patterns and clusters. The financing bins that correspond to the two-dimensional characteristics that are obtained using t-SNE are saved in a DataFrame called tsne_df. Ultimately, a scatter plot is made to show the clusters of data points, with each point colored according to its funding bin. This gives an overview of the distribution of funding amounts among the companies and how they relate to each other.

```
data_subset_funding_bins = data_tsne_selected_columns[['city', 'founded_year', 'venture','market_group','state_code', 'funding_total_usd']]

# Define bins for 'funding_total_usd' and create a binned column
bins = [0, 10000, 50000, 100000, 500000, 1000000, 5000000, 10000000, 50000000, 100000000, 500000000, 1000000000, 5000000000, 10000000000, 50000000000, 1000
bin_labels = ['<10K', '10K-50K', '50K-100K', '100K-500K', '500K-1M', '1M-5M', '5M-10M', '10M-50M', '50M-100M','100M-500M', '500M-1B', '1B-5B', '5B-10B', '1
data_subset_funding_bins['funding_bins'] = pd.cut(data_subset_funding_bins['funding_total_usd'], bins=bins, labels=bin_labels)

# Run t-SNE
tsne = TSNE(n_components=2, random_state=42)
tsne_results = tsne.fit_transform(data_subset_funding_bins[['city', 'founded_year', 'venture', 'market_group', 'state_code', 'funding_total_usd']])

# Create a new DataFrame for the t-SNE results
tsne_df = pd.DataFrame(tsne_results, columns=['TSNE1', 'TSNE2'])
tsne_df['funding_bins'] = data_subset_funding_bins['funding_bins'].values

# Plotting
plt.figure(figsize=(10, 6))
for i, label in enumerate(bin_labels):
    subset = tsne_df[tsne_df['funding_bins'] == label]
    plt.scatter(subset['TSNE1'], subset['TSNE2'], label=label, alpha=0.5)

plt.xlabel('t-SNE Feature 1')
plt.ylabel('t-SNE Feature 2')
plt.title('t-SNE of Funding Data')
plt.legend(title='Funding Bin')
plt.show()
```

Figure 9: TSNE

### 3.5.4 Model building

A code sample of one of the model buildings is given in Figure 10. The provided code describes how to pick features for a linear regression model using a genetic algorithm (GA) (Dasgupta et al.; 2013). Given that the goal is to minimize mistakes, a unique 'FitnessMin' class with negative weights has been developed. Furthermore defined is an 'Individual' class that represents a solution (a set of features). The toolbox comes pre-configured with ways to build individuals and populations, evaluate fitness using a custom function called evalModel, and produce binary attributes (to determine whether to include each characteristic). Using cross-validation on the training data, this function chooses features that match the ones in each person's binary string. If no features are selected, a large error is issued.

The population size, crossover probability, mutation probability, and number of generations are established, and the genetic operators for mating, mutating, and selection are defined. The algorithm determines the best individuals based on their fitness and runs for the designated generations, estimating the time to completion. Lastly, it publishes the total time as well as the most attractive and fit person. To reduce the prediction error, this procedure automatically chooses the most important features for the linear regression model.

```
def evalModel(individual):
    # Select features based on individual
    features = [i for i, bit in enumerate(individual) if bit == 1]
    if not features:
        return 1e6,  # Huge error if no features selected

    # Create a linear regression model
    model = LinearRegression()

    # Use cross-validation to evaluate the model
    scores = cross_val_score(model, X_train[:, features], y_train, cv=10, scoring='neg_mean_squared_error')
    return (np.mean(scores),)

# Register the genetic operators
toolbox.register("evaluate", evalModel)
toolbox.register("mate", tools.cxTwoPoint)
toolbox.register("mutate", tools.mutFlipBit, indpb=0.05)
toolbox.register("select", tools.selTournament, tournsize=3)

# Genetic Algorithm parameters
population_size = 30
crossover_probability = 0.7
mutation_probability = 0.2
number_of_generations = 20

# Create initial population
population = toolbox.population(n=population_size)

# Run the genetic algorithm
start_time = time.time()
for gen in range(number_of_generations):
    start_gen_time = time.time()
    offspring = algorithms.varAnd(population, toolbox, cxpb=crossover_probability, mutpb=mutation_probability)
    fits = toolbox.map(toolbox.evaluate, offspring)
    for fit, ind in zip(fits, offspring):
        ind.fitness.values = fit
    population = toolbox.select(offspring, k=len(population))
    end_gen_time = time.time()

    # Print generation info
    best_ind = tools.selBest(population, 1)[0]
    print(f"Generation {gen+1}/{number_of_generations} - Best: {best_ind}, Fitness: {best_ind.fitness.values[0]}")
    remaining_gens = number_of_generations - (gen + 1)
    gen_time = end_gen_time - start_gen_time
    eta = gen_time * remaining_gens
    print(f"Remaining Generations: {remaining_gens}, ETA: {eta:.2f} seconds")

# Total time taken
end_time = time.time()
print(f"Total time taken: {end_time - start_time:.2f} seconds")
```

Figure 10: Model building

# References

Arindam235 (2023). Startup investments (crunchbase), https://www.kaggle.com/datasets/arindam235/startup-investments-crunchbase. Accessed: 25/10/2023.

Dasgupta, K., Mandal, B., Dutta, P., Mandal, J. K. and Dam, S. (2013). A genetic algorithm (ga) based load balancing strategy for cloud computing, *Procedia Technology* **10**: 340–347. First International Conference on Computational Intelligence: Modeling Techniques and Applications (CIMTA) 2013.