**National College of Ireland**

**Project Submission Sheet**

Student Name:
Nachiket Anil Mehendale
………………………………………………………………………………………………………

Student ID:
X23272473
………………………………………………………………………………………………………

Programme:
Msc AI
……………………………………………………………… Year: 2024-25………………

Module:
Data Analytics for Artificial Intelligence (MSCAI1)
………………………………………………………………………………………………………

Lecturer:
Prof. Anh Duong Trinh
………………………………………………………………………………………………………

Submission Due Date:
06/12/2024
………………………………………………………………………………………………………

Project Title:
CA-Project Report – Predicting Visa Applications with Regression Algorithms
………………………………………………………………………………………………………

Word Count:
2647
………………………………………………………………………………………………………

**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 references section.  Students are encouraged to use the Harvard Referencing Standard supplied by the Library.  To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.  Students may be required to undergo a viva (oral examination) if there is suspicion about the validity of their submitted work.**

Signature:
Nachiket Mehendale
………………………………………………………………………………………………………

Date:
06/12/2024
………………………………………………………………………………………………………

PLEASE READ THE FOLLOWING INSTRUCTIONS:
**1.Please attach a completed copy of this sheet to each project (including multiple copies).**
**2.Projects should be submitted to your Programme Coordinator.**
**3.You must ensure that you retain a HARD COPY of ALL projects, both for your own reference and in case a project is lost or mislaid.  It is not sufficient to keep a copy on computer.  Please do not bind projects or place in covers unless specifically requested.**
**4.You must ensure that all projects are submitted to your Programme Coordinator on or before the required submission date.  Late submissions will incur penalties.**
**5.All projects must be submitted and passed in order to successfully complete the year.  Any project/assignment not submitted will be marked as a fail.**

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

**Continuous Assessment Project Report**

Nachiket Mehendale (x23272473)

Data Analytics for Artificial Intelligence – H9DAI

MSCAI_SEP24

School of Computing

National College of Ireland

# Title

The title of the report is **Predicting Visa Applications with Regression Algorithms**. In this study, we are considering the **visa applications dataset** published by the Department of Justice of Ireland [1]. The report demonstrates how we addressed **5 learning outcomes** by analyzing the visa applications dataset - **(i)** Retrieve, extract, manipulate, synthesize, explore, and visualize data in preparation for data analysis and machine learning. **(ii)** Concepts and methods associated with the analysis of data using numerical and statistical techniques to assist on decision-making. **(iii)** Fundamental machine learning concepts and techniques to build and evaluate machine learning models on various problem domains. **(iv)** Evaluate and employ graphical tools for building comprehensive analytics processes and dashboards. **(v)** Critically analyze, compare, summarize, and present results to support decision making and address requirements in real-world problems.

# Background research

## 2.1 Regression analysis

Thomas C. Redman. (2023) offers this example scenario: Suppose you're a sales manager trying to predict next month's numbers. You know that dozens, perhaps even hundreds of factors — from the weather to a competitor's promotion to the rumour of a new and improved model — can impact the numbers. Perhaps people in your organization even have a theory about what will have the biggest effect on sales. "Trust me. The more rain we have, the more we sell." "Six weeks after the competitor's promotion, sales jump."
Regression analysis is a way of mathematically sorting out which of those variables does indeed have an impact. It answers the questions: Which factors matter most? Which can we ignore? How do those factors interact with one another? And, perhaps most important, how certain are we about all these factors? In regression analysis, those factors are called "variables." You have your *dependent variable* — the main factor that you're trying to understand or predict. In Redman's example above, the dependent variable is monthly sales. And then you have your *independent variables*— the factors you suspect have an impact on your dependent variable.

## 2.2 Types of regression algorithms

### 2.2.1 Linear Regression

Linear Regression is one of the simplest regression models used for prediction of results. It models the relationship between the independent variable and the dependent variable. In the case of simple linear

regression, one independent variable and one dependent variable are involved. In our dataset, however, we use Multiple Linear Regression since there are several independent variables and a single dependent variable. Multiple Linear Regression attempts to model the relationship between two or more explanatory variables and a response variable by fitting a linear equation to observed data. Every value of the independent variable x is associated with a value of the dependent variable y. (Acharya *et al*., 2019)

$$y = a_1x_1 + a_2x_2 + a_3x_3 + \ldots + a_nx_n$$

The equation represents a Multiple Linear Regression where y is the dependent variable and x1, x2,..xn are independent variables.

## 2.2.2 Support Vector Machine Regression

Support vector regression (SVR) is a type of support vector machine (SVM) that is used for regression tasks. It tries to find a function that best predicts the continuous output value for a given input value. SVR can use both linear and non-linear kernels. A linear kernel is a simple dot product between two input vectors, while a non-linear kernel is a more complex function that can capture more intricate patterns in the data. The choice of kernel depends on the data's characteristics and the task's complexity. (geeksforgeeks editorial team, 2023)

## 2.2.3 Gradient Boosting Regression

Gradient boosting is a machine learning ensemble technique that sequentially combines the predictions of multiple weak learners, typically decision trees. It aims to improve overall predictive performance by optimizing the model's weights based on the errors of previous iterations, gradually reducing prediction errors and enhancing the model's accuracy. This technique is most commonly used for linear regression. (Anshul, 2024)

## 2.2.4 Random Forest Regression

Random Forest Regression in machine learning is an ensemble technique capable of performing both regression and classification tasks with the use of multiple decision trees and a technique called Bootstrap and Aggregation, commonly known as bagging. The basic idea behind this is to combine multiple decision trees in determining the final output rather than relying on individual decision trees. Random Forest has multiple decision trees as base learning models. We randomly perform row sampling and feature sampling from the dataset forming sample datasets for every model. This part is called Bootstrap. (geeksforgeeks editorial team, 2024)


## 2.3 Visa Application Dataset – (A regression problem)

A visa is an endorsement placed within a passport that permits the holder to present to an Immigration Officer for permission to enter, leave or stay in the state for a specified time period. This dataset relates to paper and online visa applications received and decisions by year and nationality. This dataset has been generated at a specific point in time and is subject to revision. (Department of Justice / Gov of Ireland, 2024) *Fig1 shows a snapshot of the dataset after importing it into the Jupyter Notebook. Fig 2 shows how many rows and columns the dataset has or the shape of the dataset and Fig 3 shows columns and their datatypes*

| Nationality | 2017 | 2018 | 2019 | 2020 | 2021 | 2022 | 2023 |
|---|---|---|---|---|---|---|---|
| Afghanistan | 128 | 127 | 178 | 58 | 67 | 189 | 220 |
| Albania | 409 | 528 | 521 | 109 | 159 | 571 | 396 |
| Algeria | 371 | 369 | 519 | 102 | 82 | 339 | 359 |
| Angola | 78 | 122 | 79 | 10 | 8 | 88 | 87 |
| Armenia | 174 | 164 | 270 | 41 | 27 | 241 | 244 |

**Fig1. Snapshot of Visa Application Dataset**

```
print("The size of the dataset is:",df.shape)
print("The dataset contains",df.shape[0],"rows")
print("The dataset contains",df.shape[1],"columns")

The size of the dataset is: (112, 8)
The dataset contains 112 rows
The dataset contains 8 columns
```

**Fig2: Shape of the dataset**

```
[30]: print("Datatypes of the columns or attributes are as follows:\n")
      df.dtypes

Datatypes of the columns or attributes are as follows:

[30]: Nationality    object
      2017           object
      2018           object
      2019           object
      2020           object
      2021           object
      2022           object
      2023           object
      dtype: object
```

**Fig3: Datatypes of the columns**

In this experiment, we would use different regression models which would analyze the visa applications in the years 2017 through 2022 and then forecast/predict the applications for year 2023. We then compare the model predictions against the original column 2023 and evaluate the error margin.

# Data Preprocessing

When it comes to creating a Machine Learning model, data preprocessing is the first step marking the initiation of the process. Typically, real-world data is incomplete, inconsistent, inaccurate (contains errors or outliers), and often lacks specific attribute values/trends. This is where data preprocessing enters the scenario – it helps to clean, format, and organize the raw data, thereby making it ready-to-go for Machine Learning models. Let's explore various steps of data preprocessing in the case of our visa applications dataset.

### 3.1 Typecasting

Machine learning models understand the data only in the numeric format during the computation. This dataset has all the year columns in object format and we have to typecast those columns to numbers. (*refer to Fig 4*)

```
•[31]: # List of year columns
       year_columns = ['2017', '2018', '2019', '2020', '2021', '2022', '2023']

       # Convert year columns to integer type
       df[year_columns] = df[year_columns].apply(pd.to_numeric, errors='coerce', axis=0)

       df.dtypes

[31]: Nationality      object
      2017            float64
      2018            float64
      2019            float64
      2020            float64
      2021            float64
      2022            float64
      2023            float64
      dtype: object
```

**Fig 4: Datatype change in year columns**

## 3.2 Missing values

Handling missing values in machine learning is crucial because most algorithms cannot work with missing data directly. If missing values are ignored or not properly addressed, they can lead to inaccurate models and unreliable predictions. (*Fig 5 shows the missing values in our dataset*)

### 2.1 (a) Check Missing Values

```
[35]: print("Check NULL or Missing values in the dataset. \n")
      df.isnull().sum()

      Check NULL or Missing values in the dataset.

[35]: Nationality      0
      2017             6
      2018             2
      2019             4
      2020            11
      2021            16
      2022            10
      2023             8
      dtype: int64
```
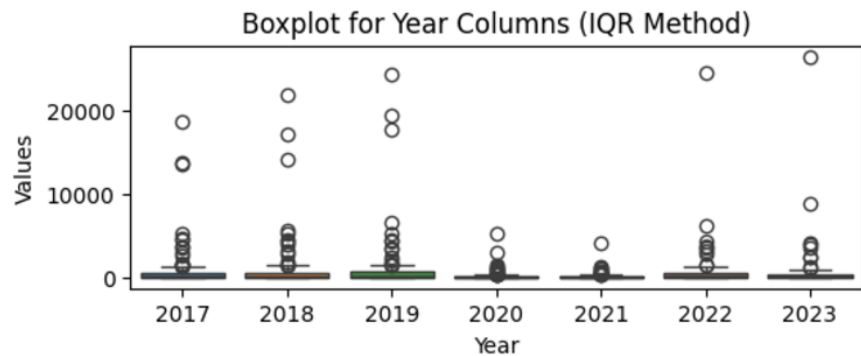
**Fig 5: Missing value count**

There is a considerable number of cells where values are missing. If we drop them, we would lose many records and therefore we decided to fill those missing values with 0. (*refer to fig 6*)

```
# Assuming 'df' is your DataFrame
df.fillna(0, inplace=True)

# Optionally, check if there are still any missing values
print("Checking for missing values after removal:")
print(df.isnull().sum())

Checking for missing values after removal:
Nationality      0
2017             0
2018             0
2019             0
2020             0
2021             0
2022             0
2023             0
dtype: int64
```

**Fig 6: Missing values filled with 0**

## 3.3 Outliers

**Outlier Detection:** Quantile-based outlier detection identifies data points that fall outside a specific range of quantiles, typically determined by the Interquartile Range (IQR). It calculates the 25th (Q1) and 75th (Q3) percentiles, and defines outliers as values below Q1−1.5×IQR or above Q3+1.5×IQR, where IQR=Q3−Q1.

A boxplot visually identifies outliers by plotting data points beyond the whiskers, which represent the range from Q1−1.5×IQR to Q3+1.5×IQR. Outliers appear as individual points. (*refer to Fig7*)

```
Outliers detected using IQR method:
2017    15
2018    14
2019    14
2020    13
2021    12
2022     9
2023     9
dtype: int64
```

**Fig 7: Outlier detection using Box plot**

**Outlier Handling**: There are quite a lot of outliers (visa applications). There is a possibility for some practical reasons a country may receive outstanding visa applications and therefore here we are not deleting or removing the outliers. Instead, we have replaced the outliers with the median value. We found that is a more reasonable way to deal with the outliers. (*refer to Fig 8*)

*# Replace outliers with the median of the respective row*
*df[year_columns] = df[year_columns].apply(lambda row: row.where(~outliers_iqr.loc[row.name],row.median(axis=0)), axis=1)*

```
Data shape after treating the outliers using IQR: (112, 8)
```
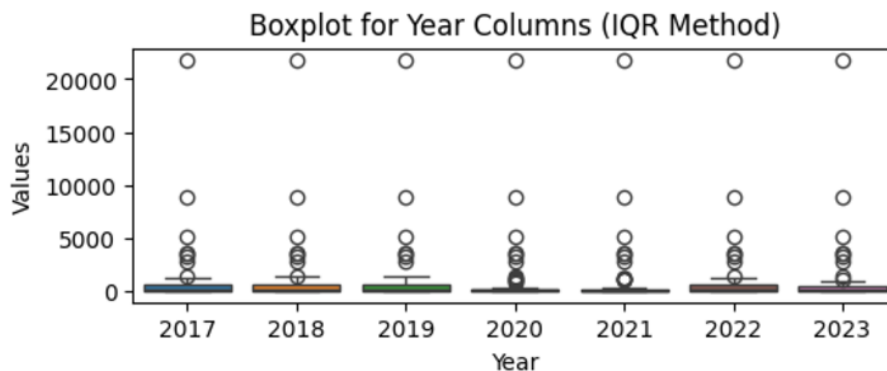
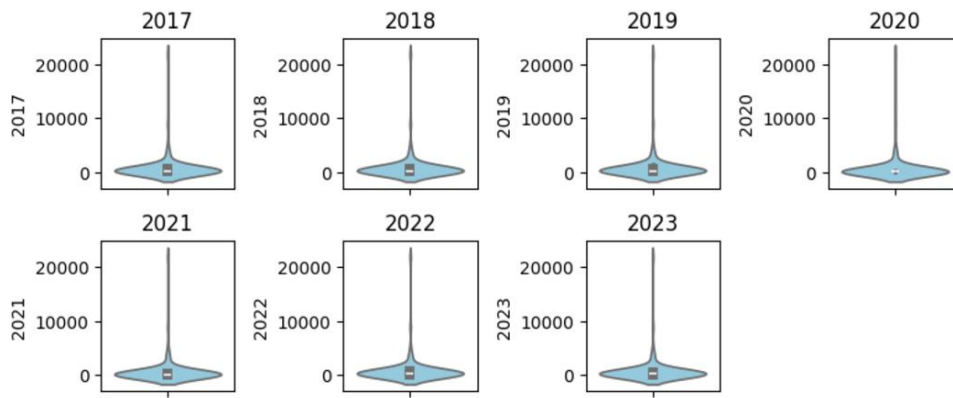**Fig 8: outliers after performing the median-replacement operation**

# Data Visualisation

Visualizing Data provides a perspective on data by showing its meaning in the larger scheme of things. It demonstrates how particular data references stand concerning the overall data picture. In the following
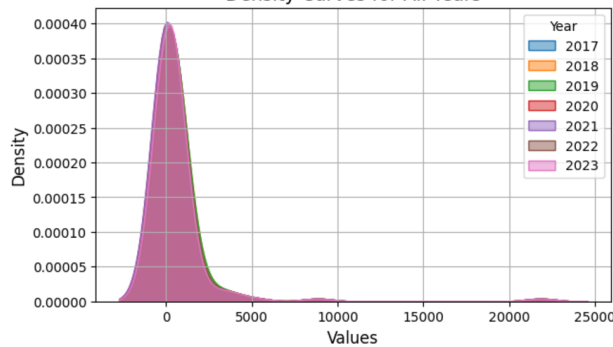
section, we have provided the different visualization techniques such univariate analysis, bivariate analysis and multivariate analysis.

## 4.1 Univariate analysis

Univariate analysis in machine learning examines the distribution and characteristics of a single variable. It helps understand data patterns, detect outliers, and summarize central tendencies and spread using statistical metrics (mean, median) or visualizations like histograms, boxplots, and violins. In the *Fig 9 below we have shown the frequency distribution of the visa applications for the all years individually using the violine charts. And Fig 10 shows kernel density plot.* It can be seen that more or less data distribution is same for all the years except for the fact that year 2023 has more applications and 2020 has more dense data.
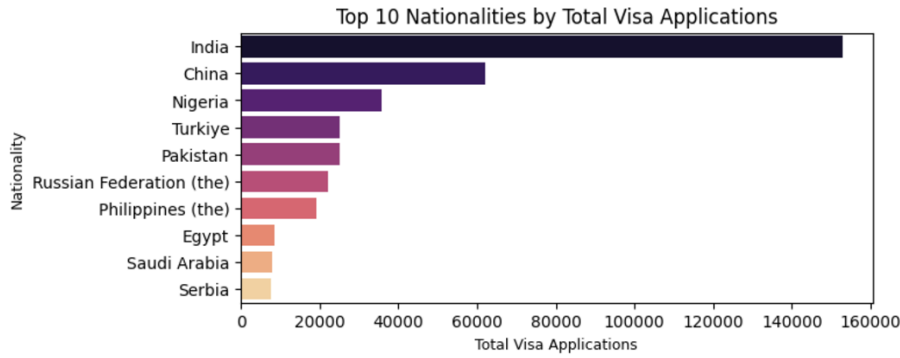


**Fig 9: Violine charts showing data distribution for all the year columns**
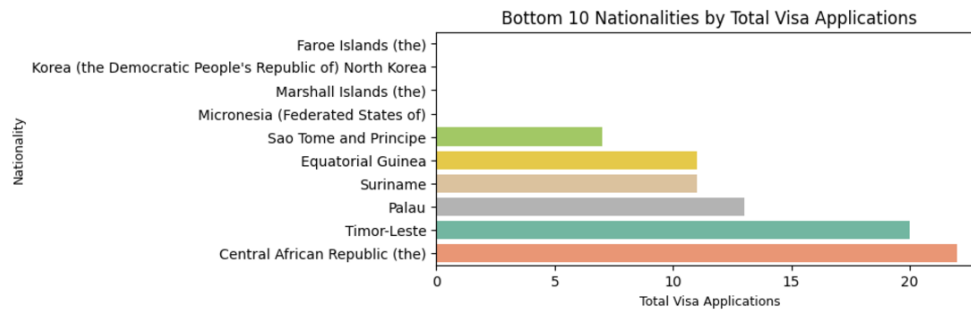


**Fig 10: K-D plot for Visa Applications for each year**

## 4.2 Bivariate analysis

Bi-variate analysis in machine learning examines the relationship between two variables. It helps identify patterns, correlations, and dependencies. *Fig 11 and Fig 12 show Top 10 and Bottom 10 countries receiving Visa Applications.* India received maximum visa applications where Faroe Island received no applications throughout the period between 2017 and 2023.
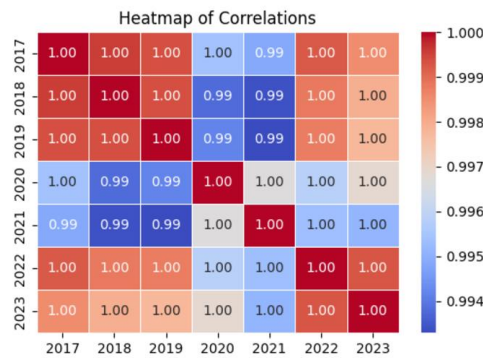
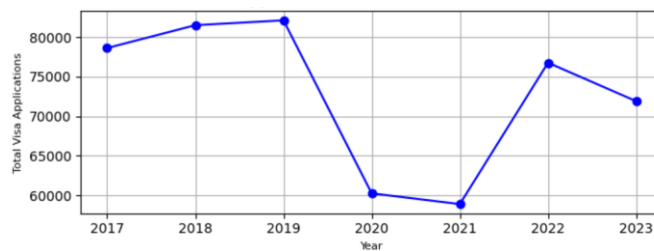**Fig11: Top10 Nationalities receiving Visa Applications**



**Fig12: Bottom10 Nationalities receiving Visa Applications**

## 4.3 Multivariate analysis

Multivariate analysis in machine learning explores relationships between multiple variables to understand their interactions. It uses techniques like correlation matrices, pair plots, and multivariate visualization to gain insights and improve predictive models. Here, we have plotted – Heatmap of correlations and Population trend from 2017 to 2023. *Fig 13 shows heatmap and Fig 14 shows the trend.*



**Fig 13. Heatmap of correlations**



**Fig 14: Trend of Total Visa Applications throughout the period**

It can be inferred that in the year 2020, 2021 – Visa Applications are reduced. This is probably due to Covid pandemic. Countries received maximum visa applications in the year 2019.

# Machine Learning Algorithms

Visa Applications is a <u>regression</u> dataset. We proposed to perform modelling using **4** types of algorithms – ***Linear Regression, Random Forest Regression, SVM Regression and Gradient Boost Regression***.

- The frequency distribution of visa applications is not normal for any of the years and therefore we decided to apply **Log Transformation**.

*year_columns = ['2017', '2018', '2019', '2020', '2021', '2022', '2023']*
*df.loc[:, year_columns] = df[year_columns].apply(lambda x: np.log(x + 1))  # Adding 1 to avoid log(0)*

- The columns are separated into X & y arrays, and **standard scaling** has been applied.
- We then called **4 regression functions**, and trained the model.

| Regression | Function |
|---|---|
| Linear | lr_model = LinearRegression() |
| Random Forest | rf_model = RandomForestRegressor(n_estimators=100, random_state=42) |
| SVM | svr_model = SVR(kernel='rbf', C=100, gamma=0.1, epsilon=0.1) |
| Gradient Boosting | gbr_model = GradientBoostingRegressor(n_estimators=100, learning_rate=0.1, max_depth=3, random_state=42) |

- We performed several permutations & combinations, and finally arrived at the best model parameters as shown in the chart above.
- **K-Fold=15** : K-Fold Cross Validation splits data into 15 parts, training on 14 and testing on 1, repeating 15 times for robust evaluation. (*refer to Fig 15*)

```python
def evaluate_model_cv(model, X, y, model_name, cv_folds=15):
    kf = KFold(n_splits=cv_folds, shuffle=True, random_state=42)
    y_pred = cross_val_predict(model, X, y, cv=kf)
```

**Fig 15- Code snippet of K-Fold Cross Validation**

# Evaluation and Discussion

Regression models used in our experiment predict how close the predicted visa applications for the year 2023 against the ground truth or the expected values (original 2023 column). The difference between their predicted and the actual value is specified in terms an "Error". There are 4 commonly used error functions we used to measure the performance of the regression models.

**Mean Squared Error (MSE)** – It is the mean or average of the squared differences between predicted and expected target values in a dataset. The formula of MSE is :
MSE = $\Sigma(y_i - p_i)2/n$, where $y_i$ is the ith observed value, $p_i$ is the corresponding predicted value for $y_i$, and n is the number of observations. The $\Sigma$ indicates that a summation is performed over all values of i. (Ken Stewart, 2023)

**Root Mean Squared Error (RMSE)** – The RMSE can be calculated as follows:
RMSE = sqrt(1 / N * sum for i to N $(y\_i - yhat\_i)^2$)
Where y_i is the i'th expected value in the dataset, yhat_i is the i'th predicted value, and sqrt() is the square root function. We can restate the RMSE in terms of the MSE as: RMSE = sqrt(MSE) (Jason Brownlee, 2021)

9

**Mean Absolute Error** (MAE) – It is a measure of the average size of the mistakes in a collection of predictions, without taking their direction into account. It is measured as the average absolute difference between the predicted values and the actual values and is used to assess the effectiveness of a regression model. The MAE loss function formula:

MAE = $(1/n) \Sigma(i=1$ to n) $|y\_i – \hat{y}\_i|$  where: n is the number of observations in the dataset. $y\_i$ is the true value. $\hat{y}\_i$ is the predicted value. (deepchecks team, 2023)

**R-squared** – It is also known as the coefficient of determination, is a statistical measure used in machine learning to evaluate the quality of a regression model. (Ejiro Onose 2023) It measures how well the model fits the data by assessing the proportion of variance in the dependent variable explained by the independent variables. R-squared is calculated mathematically by comparing the Sum of Squares of Errors (SSE) or the Sum of Squared Residuals (SSR) to the Total Sum of Squares (SST). Note: SSE and SSR can be used interchangeably. R-squared can be calculated using the following formula: $R^2 = 1 – (SSE/SST)$

Below chart (*Fig 16*) shows the values returned by the error functions or the evaluation metrics used to calculate the performance of the regression models used in this exercise.

| | Model | MSE | RMSE | MAE | $R^2$ |
|---|---|---|---|---|---|
| 0 | Linear Regression | 0.1458 | 0.3818 | 0.2649 | 0.9717 |
| 1 | Random Forest Regressor | 0.2138 | 0.4624 | 0.3142 | 0.9586 |
| 2 | Support Vector Regressor (SVM) | 0.2300 | 0.4796 | 0.3315 | 0.9554 |
| 3 | Gradient Boosting Regressor | 0.2380 | 0.4879 | 0.3285 | 0.9539 |

**Fig 16: Error Functions/Evaluation Metrics**

Out of the four models that we tested, **Linear Regressor** was found to be performing best in terms of all the error functions. Hence, the final model used for predicting the visa applications was Linear  Regression Algorithm.
We did the comparison between Ground truth and Linear Regressor Predicted values in the *Fig 18* below.

| | Actual/Ground Truth/(2023) | Predicted value(2023) by Linear Regressor |
|---|---|---|
| 0 | 2.833213 | 2.793644 |
| 1 | 5.545177 | 5.296836 |
| 2 | 5.501258 | 5.161758 |
| 3 | 5.579730 | 5.523104 |
| 4 | 6.508769 | 6.570430 |
| 5 | 3.951244 | 4.109584 |
| 6 | 1.791759 | 2.066545 |
| 7 | 5.575949 | 5.769106 |
| 8 | 0.000000 | 0.603916 |
| 9 | 5.749393 | 5.738046 |

**Fig 18: Actual Vs Prediction for Linear Regressor**

We plotted the actual vs predicted values and residuals using the graphs as shown below (*refer to Fig 19 and Fig 20*)
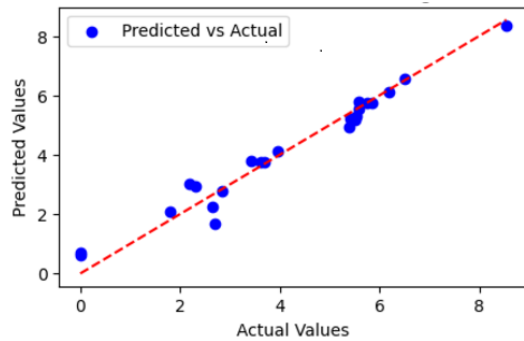


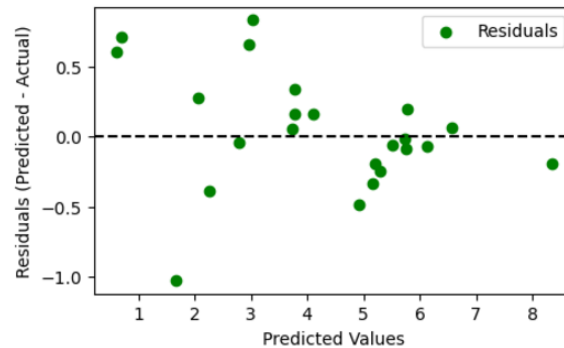**Fig 19: Actual Vs Predicted Plot for Linear Regressor**



**Fig 20: Residual Error Plot for Linear Regressor**

# Conclusion

Through experimentation on the Visa Applications dataset, we addressed all the 5 learning outcomes that we mentioned earlier in the title of this report. We learnt how to extract the data from the file, and load it into a dataframe using Python. We explored various techniques of data preprocessing such as how to handle null or missing, how to manipulate data types of the columns, and how to treat outliers in the dataset. Later we studied the theoretical concepts of **regression algorithms** in detail and also saw its different types such as **Linear Regression, Random Forest Regressor, SVM regressor, and Gradient Boost Regressor**. Then we learnt to perform several visualizations such as Violine charts, KD Plots to analyze the frequency distribution in the variables. Also, we found the Top10 and Bottom10 analysis through bar graphs and we plotted heatmap to see the correlation between multiple numerical columns of the dataset. All those visuals helped us understand the underlying pattern in the data. In the modelling section, we saw how to do split the dataset and how to apply logarithmic transformation, how to perform standard scaling, how to do cross validation and how to train the dataset using various regression algorithms. We learnt 4 error functions – Mean Squared Error, Root Mean Squared Error, Mean Absolute Error and R-Squared and found that for the Linear Regression Algorithm the error values were least. This proved that **Linear Regressor** was found to be the **best performing algorithm** in our study. Then we also compared few initial actual versus predicted values using the Linear Algorithm and found that there was a very close match. Linear regression is a simple model compared to other models. The given dataset was small and hence Linear regression was the best choice as it minimizes the risk of overfitting. Also, it had lowest MSE, RMSE, and MAE among the models and Its R-Squared value was the maximum among others.

11

# References

Acharya,Mohan & Armaan,Asfia & S Antony,Aneeta.(2019) 'A Comparison of Regression Models for Prediction of Graduate Admissions', Journel Researchgate,pp. 1-4, https://www.researchgate.net/publication/336436894_A_Comparison_of_Regression_Models_for_Prediction_of_Graduate_Admissions

geeksforgeeks editorial team (2023) 'Support Vector Regression (SVR) using Linear and Non-Linear Kernels in Scikit Learn' Available at https://www.geeksforgeeks.org/support-vector-regression-svr-using-linear-and-non-linear-kernels-in-scikit-learn/ [Accessed 04 Dec 2024]

Anshul (2024) 'Gradient Boosting Algorithm: A Complete Guide for Beginners' Available at https://www.analyticsvidhya.com/blog/2021/09/gradient-boosting-algorithm-a-complete-guide-for-beginners/ [Accessed 04 Dec 2024]

geeksforgeeks editorial team (2024) 'Random Forest Regression in Python' Available at https://www.geeksforgeeks.org/random-forest-regression-in-python/ [Accessed 04 Dec 2024]

Thomas C. Redman. (2023) People and Data: Uniting to Transform Your Business (Kogan Page)

Department of Justice / Gov of Ireland (2024) 'Visa applications and decisions - year and nationality' Available at https://data.gov.ie/dataset/visa-applications-and-decisions-year-and-nationality/resource/f4477ebf-7f24-4f2e-b222-fd661d2fbd42 [Accessed 15 Nov 2024]

Ken Stewart (2023) 'mean squared error' Available at https://www.britannica.com/science/mean-squared-error [Accessed at 30-11-2023]

Jason Brownlee (2021) 'Regression Metrics for Machine Learning' Available at https://machinelearningmastery.com/regression-metrics-for-machine-learning/[Accessed 29 Nov 2024]

deepchecks team (2023) 'Mean Absolute Error' Available at https://www.deepchecks.com/glossary/mean-absolute-error/#:~:text=Mean%20Absolute%20Error%20(MAE)%20is,effectiveness%20of%20a%20regression%20model [Accessed 02-12-2024]

Ejiro Onose (2023) 'R Squared: Understanding the Coefficient of Determination' Available at https://arize.com/blog-course/r-squared-understanding-the-coefficient-of-determination/#:~:text=R%2Dsquared%2C%20also%20known%20as,explained%20by%20the%20independent%20variables.[Accessed 03-12-2024]

Data Analytics for Artificial Intelligence (MSCAI1)

# CA Project Report Title: "Predicting Visa Applications with Regression Algorithms"

The project report explores regression analysis; its different types in detail, and finds a best model for forecasting the visa applications.

| Your Name/Student Number | Course | Date |
|---|---|---|
| **Nachiket Mehendale** | Data Analytics for Artificial Intelligence (MSCAI1) | 06-12-2024 |

# AI Acknowledgment

| Tool Name | Brief Description | Link to tool |
|---|---|---|
| | | |
| ChatGPT | For assistance in refining concepts, summarizing methodologies, and exploring key challenges related to the research project. | https://chat.openai.com/ |

# Description of AI Usage :

# AI Prompts:

1. Prompt: "Kindly Describe File Encoding Format such as UTF-8 and UTF-16. There is some issue with respect to the file not opening up in Jupyter Notebook due to encoding issue. (In 60 words)"

2. Prompt: "Can you please highlight the importance of Regression Algorithms? (In 50 words)"

3. Prompt: "Can you please provide some real-life or daily examples where we can find regression problems (In 3 points, Word limit - 100)"

4. Prompt: "Can you please tell me the Key Points in the regression algorithms - Linear, SVM, Random Forest and Gradient Boost (In 100 words)"

5. Prompt: "Which Regression algorithm suits best for the relationship - linear as well as non-linear? (In 70 words)"

6. Prompt: "Why is the linear regression is the simplest regression algorithm? and what is the history behind it? (Explain in 120 words)"

7. Prompt: "What is the systematic way to optimize the model parameters in Regression problems? (In 120 words)"

8. Prompt: "What are the significance of error functions in Regression problems? (70 words)"

9. Prompt: "In how many ways we can detect outliers? Is that always necessary to remove the outliers from the given dataset? Explain in 100 words."

10. Prompt: " What is the significance of the Violine chart in machine learning algorithm?"

11. Prompt: "What is the preferable and safest way to deal with the missing values in sensitive dataset such as visa applications (In 80 words)?"

12. Prompt: "What is the significance of Log Transformation in Machine Learning? (In 70 words)"

13. Prompt: "How K-Fold Cross-Validation works, for example If I define KFold=20 (In 70 words)"

14. Prompt: "How is MAE different from MSE, RMSE and R-Squared? (70 words)"

15. Prompt: "What is the best strategy to model the dataset when it is extremely small in size (In 100 words)"

# Evidence of AI Usage:

1.

> Kindly Describe File Encoding Format such as UTF-8 and UTF-16.
> There is some issue with respect to the file not opening up in Jupyter
> Notebook due to encoding issue. (In 60 words)

File encoding formats like UTF-8 and UTF-16 determine how text characters are represented as bytes. UTF-8 is compact and widely used, supporting all Unicode characters, while UTF-16 uses 2 or 4 bytes per character. Encoding issues in Jupyter Notebook arise when the file's encoding doesn't match the reader's expectation. Resolving this involves specifying the correct encoding when opening the file.

2.

Can you please highlight the importance of Regression Algorithms ? (In 50 words)

Regression algorithms are crucial in predicting continuous outcomes, identifying trends, and understanding relationships between variables. They enable data-driven decision-making in fields like finance, healthcare, and marketing. By modeling dependencies, regression helps optimize processes, forecast outcomes, and extract valuable insights from data, making it a cornerstone of machine learning and predictive analytics.

3.

Can you please provide some real-life or daily examples where we can find regression problems ? (In 3 points, Word limit - 100 )

1. **Housing Market Analysis**: Predicting house prices based on features like location, size, and amenities.

2. **Weather Forecasting**: Estimating future temperatures or rainfall based on historical weather data.

3. **Stock Market Prediction**: Forecasting stock prices or returns using past trends and economic indicators.

Regression is widely applied in decision-making scenarios requiring continuous outcome predictions.

Can you please provide some real-life or daily examples where we can find regression problems ? (In 3 points, Word limit - 100 )

1. **Housing Market Analysis**: Predicting house prices based on features like location, size, and amenities.

2. **Weather Forecasting**: Estimating future temperatures or rainfall based on historical weather data.

3. **Stock Market Prediction**: Forecasting stock prices or returns using past trends and economic indicators.

4.

Regression is widely applied in decision-making scenarios requiring continuous outcome predictions.

Which Regression algorithm suits best for the relationship - linear as well as non-linear ? (In 70 words)

Random Forest and Gradient Boosting regressors are well-suited for both linear and non-linear relationships. These ensemble methods work by combining multiple decision trees, effectively capturing complex patterns in data. Random Forest excels in robustness and handling overfitting, while Gradient Boosting is known for higher accuracy by iteratively optimizing errors. Both models adapt well to various data structures, making them versatile for diverse regression problems.

5.

Why is the linear regression is the simplest regression algorithm ? and what is the history behind it ? (Explain in 120 words)

Linear Regression is the simplest regression algorithm because it models the relationship between dependent and independent variables as a straight line, using a basic mathematical equation: $y = mx + c$. It assumes a linear relationship, making it computationally efficient and easy to interpret. The algorithm only requires fitting a line by minimizing the sum of squared errors (ordinary least squares), which is straightforward compared to complex models.

The history of linear regression dates back to the 19th century. Sir Francis Galton first described the concept in the context of heredity, observing "regression toward the mean." Later, Karl Pearson formalized the method mathematically, and it became a foundational tool in statistics and data science, influencing diverse fields.

6.

What is the systematic way to optimize the model parameters in Regression problems ? (In 120 words)

Optimizing model parameters in regression problems involves a systematic approach:

1. **Feature Selection and Engineering**: Start by selecting relevant features and transforming data (e.g., scaling or encoding) to improve model accuracy.

2. **Hyperparameter Tuning**: Use techniques like Grid Search or Random Search to test combinations of hyperparameters systematically. Advanced methods like Bayesian Optimization or Tree-structured Parzen Estimators (TPE) can automate and refine this process.

3. **Regularization**: Apply Lasso (L1), Ridge (L2), or Elastic Net regularization to penalize overfitting and improve generalization.

4. **Cross-Validation**: Implement k-fold cross-validation to ensure the model performs well on unseen data.

5. **Evaluation Metrics**: Compare performance using metrics like MSE, RMSE, MAE, or $R^2$. Iteratively refine based on results to balance bias-variance tradeoff effectively.

7.    ◁» ⬚ 👍 👎            ↓

What are the significance of error functions in Regression problems ? (70 words)

Error functions are crucial in regression problems as they quantify how well the model predicts outcomes. They guide optimization during training, enabling the model to minimize errors and improve accuracy. For example, **MSE** emphasizes large errors, **MAE** provides robustness against outliers, and **Huber Loss** balances both. By selecting appropriate error functions, models can align with problem-specific needs, ensuring effective performance and reliable predictions in real-world applications.

8.    ◁» ⬚ 👍 👎

In how many way we can detect outliers ? Is that always necessary to remove the outliers from the given dataset ? Explain in 100 words.

Outliers can be detected using several methods:

1. **Statistical methods** like Z-scores (values beyond a threshold of 3) or IQR (Interquartile Range) for values outside 1.5 times the IQR.

2. **Visual methods** like box plots, scatter plots, or histograms to identify data points that differ significantly.

3. **Model-based methods** using machine learning algorithms to detect unusual patterns. While removing outliers can improve model performance, it's not always necessary. In some cases, outliers represent genuine variance or critical insights. Deciding whether to remove them depends on the data context and the impact on model accuracy.

9.

What is the significance of the Violine chart in machine learning algorithm ?

The **violin chart** in machine learning is a powerful visualization tool that combines aspects of both **box plots** and **density plots**. It provides insights into the distribution of data, showing the probability density of the data at different values. The significance of the violin chart in machine learning is:

1. **Visualizes Distribution**: It shows the data's distribution, including skewness, multimodality, and outliers.

2. **Comparison of Multiple Groups**: It allows for easy comparison of the distribution of different categories or classes, helping to understand class imbalances or variations.

3. **Identifies Central Tendencies and Variability**: It reveals the median, quartiles, and overall range of the data, which is useful for feature engineering and understanding feature importance.

10.

What is the preferable and safest way to deal with the missing values in sensitive dataset such visa applications (In 80 words)?

In sensitive datasets like visa applications, the safest way to handle missing values is through **imputation** rather than removal, ensuring no loss of critical data. Common techniques include using **mean, median, or mode imputation** for numerical values, or employing **predictive modeling** (e.g., regression or k-NN) to estimate missing data based on other features. For categorical data, **mode imputation** or using a placeholder category like "Unknown" can be effective. Always ensure ethical considerations and transparency in data processing.

11.

What is the significance of Log Transformation in Machine Learning?
(In 70 words)

Log transformation in machine learning is used to stabilize variance, normalize data, and reduce skewness. It helps make relationships between variables more linear, which is crucial for algorithms like linear regression. By compressing large values, log transformation can also prevent issues related to extreme outliers, making the model more robust and improving its predictive accuracy. It's particularly useful when the data spans multiple orders of magnitude.

12.

How K-Fold Cross-Validation works, for example If I define KFold=20
(In 70 words)

If you define KFold=20, the dataset is split into 20 equally sized subsets (folds). The model is trained on 19 folds and tested on the remaining fold. This process is repeated 20 times, with each fold being used as the test set once. The model's performance is averaged over the 20 iterations, providing a robust estimate of its ability to generalize and reducing bias or variance in the results.

13.

How is MAE different from MSE, RMSE and R-Squared ? (70 words)

MAE (Mean Absolute Error) measures the average absolute difference between predicted and actual values, providing a straightforward error metric. MSE (Mean Squared Error) squares the differences, penalizing larger errors more. RMSE (Root Mean Squared Error) is the square root of MSE, making it interpretable in the same units as the data. R-Squared measures the proportion of variance explained by the model, indicating goodness-of-fit. MAE is more robust to outliers than MSE and RMSE.

14.

What is the best strategy to model the dataset when it is extremely small in size (In 100 words)

When working with an extremely small dataset, it's crucial to avoid overfitting and ensure the model generalizes well. A good strategy includes:

1. **Simplify the Model:** Use simpler models like Linear Regression or Ridge Regression to prevent complexity from causing overfitting.

2. **Regularization:** Employ techniques like L1 (Lasso) or L2 (Ridge) regularization to reduce overfitting.

3. **Cross-Validation:** Use techniques like k-fold cross-validation to ensure robust evaluation of model performance.

4. **Data Augmentation:** If possible, generate synthetic data or apply transformations to increase the dataset size.

5. **Transfer Learning:** Leverage pre-trained models or features from larger datasets when applicable.

15.