TEXAS
The University of Texas at Austin

# PREDICTING MOST LUCRATIVE DATA SCIENCE JOBS

Data Science Programming – Group Project

Group 19

The University of Texas at Austin

# TEAMMATES

Jiarui Chang

Santhosh Kumar

Vi Tran

Meenal Gaba

## Background

- Data science is one of the hottest job profiles in the market right now and it is experiencing remarkable growth over the years
- In 2022, the data and analytics industry is worth *USD ~250 billion and it is expected to grow to USD ~330 billion (⬆35% increase)*
- The U.S. Bureau of Labor Statistics predicts there will be *~11.6 million data science-related jobs in 2026*
- Irrespective of their educational backgrounds, people are shifting towards data science due to its relevance and attractive salary
- Understanding the variation in this market due to multiple factors like *job profile, size of the company, location of company, experience, mode of work* etc., will be key to understanding the earning potential
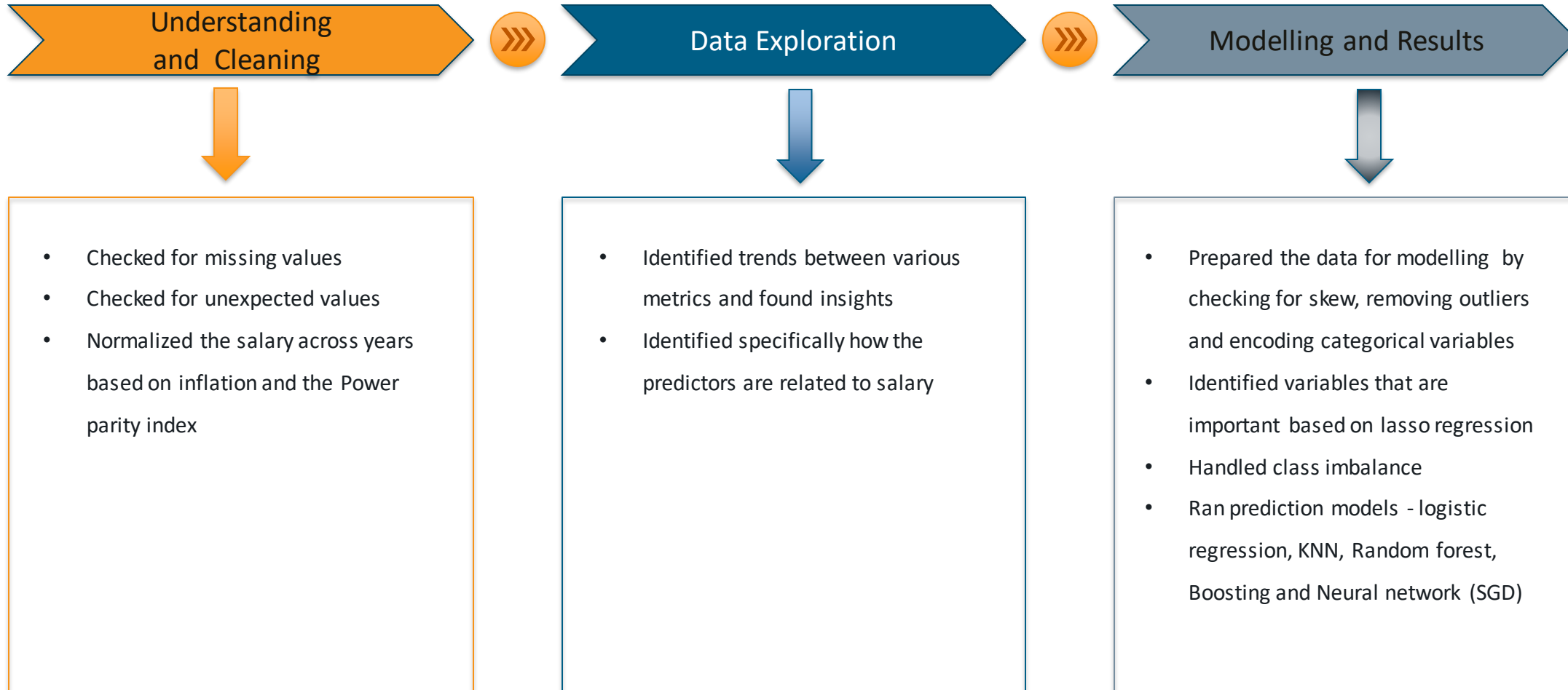
## Objectives

- Understanding the lay of the land in terms of how different metrics influence each other in the field of data science jobs
- Figuring out the metrics that have an impact on the salary
- Predicting the most lucrative jobs within the data science domain - *Who are the data scientists that earn more the 75% of people in this domain?*

# We have followed a three-step process to arrive at the solution

## Understanding and Cleaning

## Data Exploration

## Modelling and Results

**Understanding and Cleaning**
- Checked for missing values
- Checked for unexpected values
- Normalized the salary across years based on inflation and the Power parity index

**Data Exploration**
- Identified trends between various metrics and found insights
- Identified specifically how the predictors are related to salary

**Modelling and Results**
- Prepared the data for modelling by checking for skew, removing outliers and encoding categorical variables
- Identified variables that are important based on lasso regression
- Handled class imbalance
- Ran prediction models - logistic regression, KNN, Random forest, Boosting and Neural network (SGD)

## Data Snapshot

| Work Year | Job Category | Experience Level | Employment Type | Employee Location | Mode of Work | Company Location | Company Size | Salary | Salary Currency | Salary in USD |
|---|---|---|---|---|---|---|---|---|---|---|
| 2023 | Other | Senior | Full-time | Spain | Full-Remote | Spain | LARGE | 80000 | EUR | 85847 |
| 2023 | Machine Learning | Mid/Intermediate level | Contractor | India | Full-Remote | United States | SMALL | 1000000 | USD | 13000 |
| 2023 | Machine Learning | Mid/Intermediate level | Contractor | United States | Full-Remote | United States | SMALL | 25500 | USD | 25500 |
| 2023 | Data Science | Senior | Full-time | Canada | Hybrid | Canada | MEDIUM | 175000 | USD | 175000 |
| 2022 | Data Science | Senior | Full-time | Canada | Full-Remote | Canada | MEDIUM | 120000 | USD | 120000 |
| 2021 | Data Science | Senior | Full-time | United States | On-Site | United States | LARGE | 222200 | USD | 222200 |
| 2020 | Data Science | Senior | Full-time | United States | On-Site | United States | LARGE | 136000 | USD | 136000 |
| 2023 | Data Engineering | Senior | Full-time | United States | Full-Remote | United States | MEDIUM | 130000 | USD | 130000 |
| 2023 | Data Science | Senior | Full-time | Canada | On-Site | Canada | MEDIUM | 141000 | USD | 141000 |
| 2023 | Data Science | Senior | Full-time | United States | On-Site | United States | MEDIUM | 147100 | USD | 147100 |
| 2023 | Data Science | Senior | Full-time | United States | On-Site | United States | MEDIUM | 90700 | USD | 90700 |
| 2023 | Data Engineering | Senior | Full-time | United States | Full-Remote | United States | MEDIUM | 130000 | USD | 130000 |
| 2023 | Data Engineering | Senior | Full-time | United States | Full-Remote | United States | MEDIUM | 100000 | USD | 100000 |

## Data Dictionary

| Variable | Description | Variable | Description |
|---|---|---|---|
| Work Year | The year the salary was paid (2020 - 2023) | Employee Location | Employee's country of residence during the work year |
| Job Category | Data Science | ML Engineer | Data Architect | Data Engineer | Others | Company Location | The country of the employer's office |
| Experience Level | Senior | Mid | Entry | Mode of work | On-Site | Hybrid | Remote |
| Employment Type | Full time | Contractor | Part time | Freelance | Company Size | The median number of people that worked for the company during the year ( Large | Medium | Small) |
| Salary | Salary Currency | The total gross salary amount paid and local currency | Salary in US | The salary in USD |

# Normalized salaries based on inflation rates and purchasing power parity to ensure fair comparisons

## 1 Adjusting for Inflation

Global Yearly Inflation Rates (Actual 2019-2023 & Forecast 2024-2028)



- As we have salaries across 2020-2023, it is important to bring them to one scale to make fair comparison
- We have scaled up the salary data from 2020 – 2022 to 2023 based on the *inflation rates of each country*

## Normalized salary

## 2 Adjusting for Purchasing Power Parity

**Burgernomics: The Price of a Big Mac in Comparison**
Price of a Big Mac in selected countries (in U.S. dollars)

| | |
|---|---|
| Switzerland | 6.71 |
| United States | 5.67 |
| Brazil | 4.80 |
| United Kingdom | 4.41 |
| South Korea | 3.89 |
| Japan | 3.54 |
| China | 3.12 |
| India | 2.65 |
| Russia | 2.20 |
| South Africa | 2.15 |

As of January 2020
Source: The Economist

statista

| Country | Avg. DS Salary | Absolute Salary (USD) | PPP Adjusted Salary |
|---------|---------------|----------------------|---------------------|
| India | ₹ 1,300,000 | $15.6K | **$56K** |
| China | ¥293,000 | $41.1K | **$70K** |

**AGENDA**

- Background and Objectives

- Understanding the dataset and cleaning

- Exploratory data analysis

- Predictive Modeling

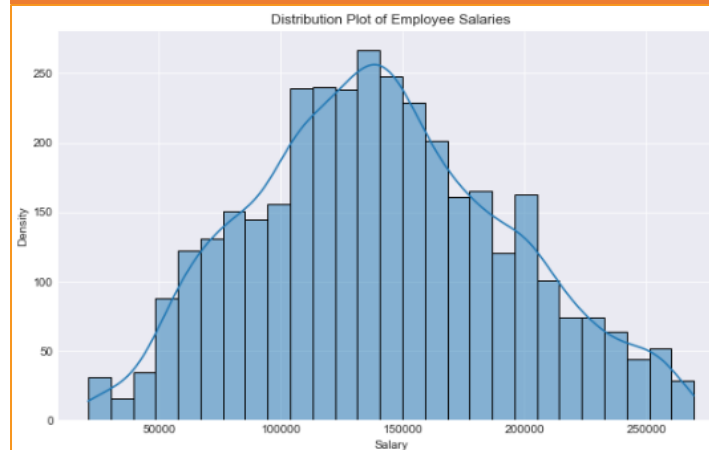- Recommendations and conclusion

# EDA – Opportunity Growth

## Global Impact



Over the years, data science opportunities have seen a remarkable surge across the globe.

## Mode of Work



➤ Full-remote jobs decline from 2020 to 2023, shifting towards on-site opportunities indicating the **COVID impact**

# EDA – Geographical Distribution



Map of Business concentration

### Key Takeaways

➢ Maximum data science roles lie around north American region.

# EDA – Pay scale

## Company Size



Salary Distribution Across company size

➢ Medium-sized companies offer higher average compensation than large or small-sized companies.

## Employment Type



Salary Distribution Across Different Employment Types

➢ Full-time employees enjoy a higher pay scale in comparison to other employee types.

# EDA – Pay scale

**Experience Level**

Salary Distribution Across Different Experience levels

**Job Category**

Salary Distribution across Job category

➢ As experience grows, so does the offered salary

➢ Job roles do not significantly impact the Pay scale

➢ Among the most prevalent roles, data engineering positions exhibit comparatively lower salaries

# EDA – Geographical Distribution



Average Salary by Company Location

## Key Takeaways

➢ Technologically developed countries such as the United States, Singapore, and Canada has the highest average salary.

➢ Less developed areas such as the majority of South America earned less.

*By analyzing these patterns, it is clear that the normalized salary metric has some relationship with the predictors ( experience level, employment type, remote ratio, company size, company location and employee location)*

# Data Preparation for Modelling

## Identifying And Excluding Outliers

Distribution Plot of Employee Salaries

It is clear that there are a **few outliers** in the target variable. We will **remove records that fall outside of mean +/- 2 SD confidence interval** for the purpose of modelling (**~170 records**)

## No Skew in salary

Distribution Plot of Employee Salaries

Salary is *normally distributed*

## Converting salary into category

➢ As our objective is to predict if the salary is going to fall above 75% of the people in this domain, we will be converting salary into a **binary categorical variable**

➢ **1 --> Above 75th percentile**

➢ **0 --> Below 75th percentile**

## One Hot Encoding of Categorical Variables

In one hot encoding, we convert each categorical value of a variable into a **new categorical column and assign a binary value of 1 or 0 to** those columns. *See example to the right:*

| Mode of work | | On-Site | Hybrid | Full-Reote |
|---|---|---|---|---|
| On-Site | One Hot Encoding → | 1 | 0 | 0 |
| Hybrid | | 0 | 1 | 0 |
| Full-Remote | | 0 | 0 | 1 |

# Variable selection and handling class imbalance

| | Step 1 | Step 2 | Step 3 | Step 4 |
|---|---|---|---|---|
| **Model →** | *Base Model* | *Lasso Regression* | *Model - After variable selection* | *Model - Class imbalance fixed* |
| **Description** | A logistic regression model run with all the response variables | A Lasso regression model to identify the variable importance | Logistic regression model to see if there is improvement in prediction after variable selection | Updated the sampling technique based on *RandomUnderSampler* to fix the class imbalance |
| **Results** | ➤ Accuracy: 76%<br>➤ Precision: 73%<br>➤ Recall: 76% | ➤ Accuracy: 75%<br>➤ Precision: 70%<br>➤ Recall: 75% | ➤ Accuracy: 76%<br>➤ Precision: 73%<br>➤ Recall: 76% | ➤ Accuracy: 62%<br>➤ Precision: 73%<br>➤ Recall: 62% |
| **Confusion Matrix** | *Training Data Split*<br>Below 75th → 1,891<br>Above 75th → 618 | *Training Data Split*<br>Below 75th → 1,891<br>Above 75th → 618 | *Training Data Split*<br>Below 75th → 1,891<br>Above 75th → 618 | *Training Data Split*<br>Below 75th → 618<br>Above 75th → 618 |
| **Takeaway** | ➤ *Very poor* Recall for Above 75th percentile due to low samples size for training – **High** Class Imbalance | ➤ **Removed** all the **variables** that were **penalized to 0 based on L1 penalty** | ➤ **No significant difference** before and after variable selection<br>➤ **Still High** Class imbalance | ➤ Though overall accuracy has dropped, the model does a **decent job in predicting both categories** – Class imbalance treated |

Confusion matrices (Predicted: Below 75th / Above 75th; Actual rows):

Step 1:
| Actual | Below 75th (Predicted) | Above 75th (Predicted) |
|---|---|---|
| Below 75th | 99% | 1% |
| Above 75th | 93% | 7% |

Step 2:
| Actual | Below 75th (Predicted) | Above 75th (Predicted) |
|---|---|---|
| Below 75th | 99% | 1% |
| Above 75th | 97% | 3% |

Step 3:
| Actual | Below 75th (Predicted) | Above 75th (Predicted) |
|---|---|---|
| Below 75th | 99% | 1% |
| Above 75th | 93% | 7% |

Step 4:
| Actual | Below 75th (Predicted) | Above 75th (Predicted) |
|---|---|---|
| Below 75th | 60% | 40% |
| Above 75th | 30% | 70% |

# Trying out different models to further improve prediction

| Model → | KNN Classification | Random Forest | Boosting (GBM) | Neural Network (SGD) |
|---|---|---|---|---|
| **Description** | Identified the best K value using cross-validation and built a KNN classifier | Using 5 fold cross validation and Grid search, identified the best hyperparameters for RF | Using 5 fold cross validation and Grid search, identified the best hyperparameters for boosting | Using SGD built a neural network model |
| **Results** | ➢ Best K: 27<br>➢ Accuracy: 76%<br>➢ Precision: 74%<br>➢ Recall: 76% | ➢ Depth: 5 \| Feature: 3 \| Trees: 1000<br>➢ Accuracy: 62%<br>➢ Precision: 74%<br>➢ Recall: 62% | ➢ $\lambda$=0.01\| Depth: 3 \| Trees: 500<br>➢ Accuracy: 62%<br>➢ Precision: 73%<br>➢ Recall: 62% | ➢ Accuracy: 70%<br>➢ Precision: 73%<br>➢ Recall: 70% |
| **Confusion Matrix** | Predicted / Below 75th / Above 75th<br>Actual Below 75th: 98% \| 2%<br>Actual Above 75th: 87% \| 13% | Predicted / Below 75th / Above 75th<br>Actual Below 75th: 60% \| 40%<br>Actual Above 75th: 30% \| 70% | Predicted / Below 75th / Above 75th<br>Actual Below 75th: 60% \| 40%<br>Actual Above 75th: 30% \| 70% | Predicted / Below 75th / Above 75th<br>Actual Below 75th: 75% \| 25%<br>Actual Above 75th: 45% \| 55% |
| **Takeaway** | ➢ **Worse** than Logistic Regression | ➢ Same as Logistic Regression | ➢ Same as Logistic Regression | ➢ **Slightly better than logistic regression / RF/ Boosting** due to better accuracy and similar recall% |

# Recommendations and conclusion

| Model → | Logistic Regression / RF / Boosting | Neural Network (SGD) |
|---|---|---|
| **Results** | ➢ Accuracy: 62%<br>➢ Precision: 73%<br>➢ Recall: 62% | ➢ Accuracy: 70%<br>➢ Precision: 73%<br>➢ Recall: 70% |
| **Confusion Matrix** | | |

**Logistic Regression / RF / Boosting — Confusion Matrix**

| | | Predicted | |
|---|---|---|---|
| | | Below 75th | Above 75th |
| Actual | Below 75th | 482 (60%) | 322 (40%) |
| | Above 75th | 83 (30%) | 189 (70%) |

**Neural Network (SGD) — Confusion Matrix**

| | | Predicted | |
|---|---|---|---|
| | | Below 75th | Above 75th |
| Actual | Below 75th | 603 (75%) | 201 (25%) |
| | Above 75th | 122 (45%) | 150 (55%) |

## Recommendations and conclusions

➢ While the neural network model performed better than other models, the prediction accuracy, precision and recall are not so high

➢ Maybe advanced feature engineering techniques and a bigger sample size can help predict the salary better


Feature Importance