# Deliverable 3

### 1. Team

**a) Members**
    **1. Vishnu Vardhan Darimidi**
    **2. Bharath Kumar Gompa**
    **3. Koosha Sharifani**
    **4. Tharun Abhinav Suraneni**
    **5. Alex Miller**

**b) Communication plan to include project artifact repository:**

Methods of communication :- e-mail, Whatsapp(once a week)
Communication response times:- Same day

Meeting attendance:- After looking into everyone's schedule we decided to meet weekly once virtually to discuss the progress of the project and complete the required things and meet the deliverable requirements.

Running meetings:- Virtual.

Git Link:-  https://github.com/Vishnu-Vardhan1809/Big-Data-Project

### 2. Selection of data to analyze from:

We looked at the information and data for our project. We'll be using a dataset called "2021 Kaggle Machine Learning and Data Science survey" from a website called Kaggle\competitions. We'll study this dataset to see how it's set up and what it contains, and then use that information to figure out some results.

Dataset -Kaggle link: https://www.kaggle.com/competitions/kaggle-survey-2021/data

### 3. Business Problem or Opportunity, Domain Knowledge (link to information on domain relative to data, problem or opportunity):

With the growing rate of data, the demand to hire data scientists in the market also increases. So to get a better understanding of the people that break into the field of data and the requirement to meet today's market's expectation to deliver the project as early as possible with high accuracy, we need predictive analysis of people who pursuing data science and machine learning job and what are the certain features that some of them have that make them as more significant potential in the job market.

The dataset contains a complete industry-wide survey that represents information about job seekers including the age, level of their education, skill sets, carrier, and origin. It has around 10,000 rows and 369 columns.

 **Business Problem or Opportunity:**

1. How does the level of education of job seekers in the field of data science and machine learning correlate with their salary and job opportunities
2. What are the most common career backgrounds of people who transition into data science and machine learning jobs
3. Are there any specific skills or programming languages that are more desirable for data science and machine learning jobs in the current job market.

# Input Dataset:

There are 369 columns in this dataset. Below shown are few important columns in the dataset.

1. Age
2. Gender
3. Country
4. Education level
5. Job role
6. Years of coding Experience
7. Programming language used most often
8. Computing platform used most often
9. Type of most data used often
10. Annual Compensation

## 4. Research Objectives and Question(s)

1. Analyzing the efficacy of a machine learning algorithm in forecasting job roles and selecting the algorithm that produces the most optimal results for the model.
2. Examine patterns and trends within the data science industry, and explore how various factors such as job roles, education, and compensation may be influencing these trends over time.
3. Explore the potential effects of emerging technologies and advancements in the data science industry on job roles, skills, and employment opportunities. Understand how the landscape of data science may shift in response to these changes over time.
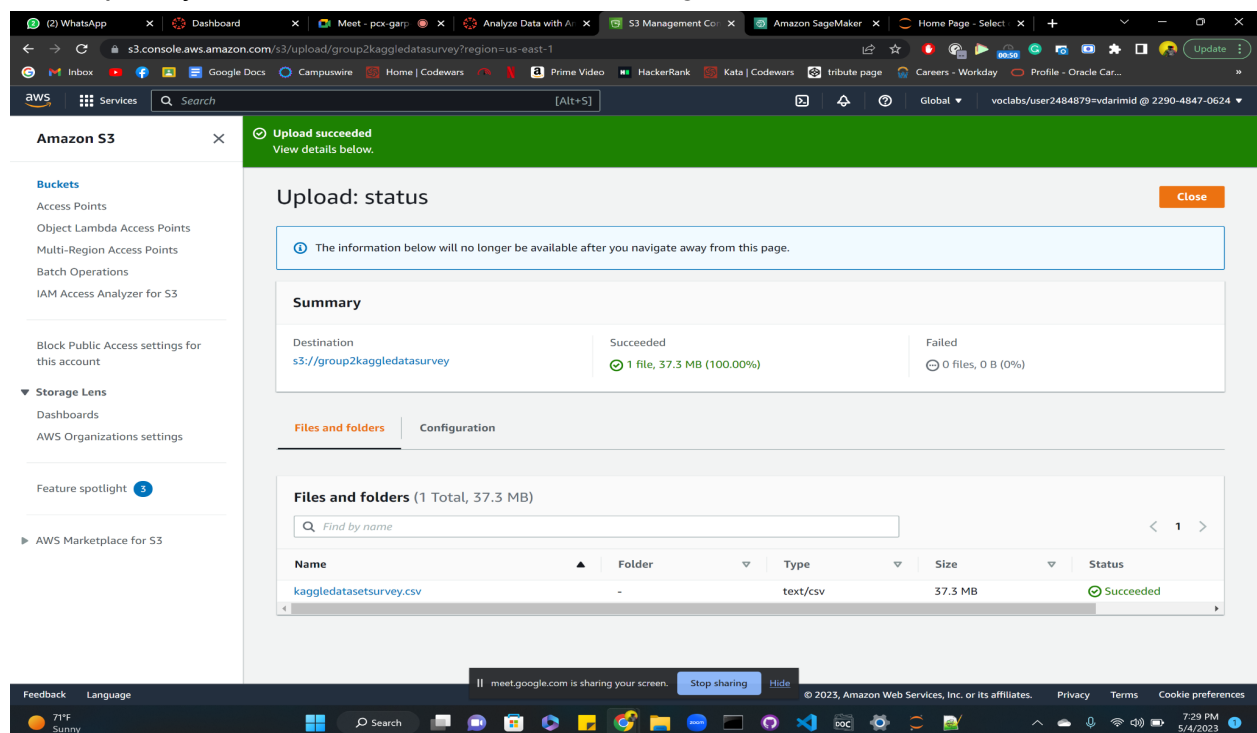
## Analytics and Machine Learning

- Prior to developing any analytical model, the dataset was cleaned to remove missing values and special characters.
- We employed a classification technique to categorize the output column and, therefore, did not take outliers into account as they were deemed irrelevant.
- The dataset was split into training and testing sets with 60% and 40% of the data, respectively, to create two models for predicting job roles using random forest classification algorithm.
- Standard metrics such as accuracy and confusion matrix were used to evaluate the performance of these models.
- The  random forest classification algorithm. had an accuracy score of 0.5495.
- The statistical approaches utilized by these machine learning algorithms enable them to generate predictions based on the patterns and correlations they uncover within the data.

## Implementation:

The data was uploaded to an S3 bucket named group2kaggledatasurvey and was then loaded into a Jupyter notebook instance using the AWS S3 copy command. Subsequently, the data was combined into a single dataframe.

We have used Random Forest Classifier model:

We have calculated accuracy and Confusion matrix for the above model.



## Evaluation and Optimization:

- Accuracy scores and confusion matrix were used to evaluate the performance of the models. The accuracy metric assesses the proportion of correct predictions, while the confusion matrix presents a summary of the model's prediction results.

- For 40% of the testing data, the model's accuracy score of 0.4943 indicates that it is not accurately identifying the cover type.

## Results:

Accuracy:0.4943 and Confusion Matrix.

# Future work, Comments:

**1) What was unique about the data? Did you have to deal with imbalance? What data cleaning did you do? Outlier treatment? Imputation?**

The dataset that we are dealing with is conducted an industry-wide survey that presents a truly comprehensive view of the state of data science and machine learning. What we did as a research team was utilize these data in order to do both descriptive and predictive analysis. As we already know, descriptive analytics tells what happened in your business in the past week, month, or year. On the other hand, Predictive analytics determines the potential outcomes of present and past actions and trends. Some of the features in your dataset, such as age, gender, country, education level, and current role, are descriptive because they reflect on things past consumer behavior. Some of the features, such as programming languages and years of coding experience, are predictive because they can help assess the likelihood of a certain event happening in the future, such as getting a job role as a data scientist.

The unique thing about the dataset is that it is relevant to a wide variety of industries, including genetics, health sciences, economics, chemometrics, sociological surveys, environmental sciences, and finance. Furthermore, The features are flexible for both descriptive and predictive analysis depending on our goal. For example, you can use descriptive analytics to compare the demographics and skills of different job roles in your dataset. You can also use predictive analytics to forecast the demand and supply of different job roles based on the trends and patterns in your dataset. Here we are trying to predict the job role based on the remaining factors like languages they know and degree they achieved and age and gender.

After cleaning the data and preparing it for preprocessing, one of the things that we considered was whether the data is imbalanced or not. Imbalanced data refers to a situation, where one target class represents a significant portion of observations. Also, There are different ways to deal with imbalanced data, depending on the context such as Down sampling, Up sampling, and Using different metrics. For our dataset, we were dealing with imbalanced data which crate problems when we try to predict any field in the data frame and specifically the columns related to program languages that our users know had some issues with imbalanced data. Balanced data will help the model to predict accurately. Therefore, we used a resampling method called Synthetic Minority Oversampling Technique to address this issue. Our dataset mostly consisted of object data types. But one column contains integer data but it shows the object type in its

property. So processed the particular field by changing the type of the column. Null values are replaced by the mode of the column. The rows which consist of maximum null values are deleted.

## 2). Did you create any new additional features / variables?

One of the ways to deal with imbalanced data was to create additional variables and columns but it would make multicollinearity or redundancy in your data or an overfitting problem with the final prediction. In the same way, we did not apply any outlier treatment because our target column is a categorized column so outlier treatment is not needed. Overall, it is important to carefully consider data preprocessing and handling techniques when dealing with imbalanced datasets in order to ensure accurate model predictions.

## 3) What was the process you used for evaluation? What was the best result?

For the project, the process we used for evaluation using was the random forest classifier which starts by splitting the dataset into three subsets: train, validation, and test. For the ratio, we used an 80/20/10 split ratio for the train/validation/test sets. The reason that we choose a random forest classifier is that it uses averaging to improve the predictive accuracy and control over-fitting. It can handle both categorical and numerical features, and it is robust to imbalanced data and outliers. It performs well with imbalanced data and can handle non-linear relationships between the features. In the end, after creating and setting the model we evaluated the model performance on the test set using various metrics. I used accuracy, precision, recall, F1-score, ROC curve, and PR curve to measure how well the model can predict the job role of new data points.

## 4) What were the problems you faced? How did you solve them?

During the project we all faced a few challenges in the preprocessing phase. Deciding the threshold value to delete rows and fitting data with the right estimates is a challenging task. Furthermore, handling imbalanced data can be difficult as it requires careful consideration of various techniques such as resampling, boosting, and bagging.

**5) What future work would you like to do?**

for future work, we could explore other machine learning algorithms and techniques specifically designed to handle imbalanced data or use different features to predict other variables related to job roles or entrepreneurial competency

**6) Instructions for individuals that may want to use your work**
In order to obtain comparable outcomes, individuals interested in utilizing our work should consider the following guidelines:

Install the essential Python libraries, including Numpy, Pandas, Matplotlib, Seaborn, and Scikit-learn, that are widely used in this project
Ensure that the dataset is thoroughly cleaned, transformed, and prepared to eliminate the majority of the outliers. This will result in the development of accurate machine learning models and visualizations