# Data Analysis Project Report

## Summary:

This large-scale data analysis project's main goal was to do a thorough investigation of the data analysis employment market. We methodically gathered job listings from a well-known job search portal by using web scraping techniques. The job listings included important information including job titles, salary, firm information, addresses, and necessary abilities. Web scraping, data collecting, data cleaning, data processing, and data visualization were among the stages of the project. Every stage was carried out with great care in order to get deep understanding of the subtleties and complexity of the data analysis work environment.

## Project Phases:

1. **Web Scraping:** Using BeautifulSoup and Selenium for web crawling was the first step. These tools were very helpful in browsing the job search portal, processing critical information including addresses, job titles, salary, business data, and necessary abilities, and extracting HTML text.

2. **Data Collection:** A carefully selected list of search phrases relevant to data analysis locations was used to guarantee an exhaustive dataset. A representative sample of the data analysis job market was obtained by methodically visiting the produced URLs in order to collect a variety of job ads.

3. **Data Cleaning:** In order to guarantee the accuracy and excellence of the dataset, data cleaning procedures were essential. Formatting mistakes, missing values, and inconsistencies were all fixed by thorough cleaning. Crucial data, including job descriptions, addresses, firm names, job titles, and wages, were painstakingly retrieved and arranged.

4. **Data Processing:** The pandas software was used to convert the processed data into a structured DataFrame. To improve the dataset's richness, more characteristics were computed, such as the number of talents needed for each job posting, and special emphasis was paid to standardizing pay values.

5. **Visualization:** The goal of the visualization phase was to convey findings in an understandable and educational way. Among the noteworthy visuals were:

   a. **City-wise Job Distribution:** geographical distribution of job postings among various cities, providing information on the characteristics of the local labor market. Bar Plot

b. **Correlation Analysis:** investigation of relationships between the quantity of skills needed and a derived metric (Points_of_Job), offering a more complex view of the dynamics of skills in job postings. Scatter Plot and Correlation Matrix

c. **Job Class Distribution:** Using a weighted score method, job advertisements are categorized into three different groups (1, 2, or 3) to help with job hierarchies. Bar Plot

d. **Salary Distribution Across Top Cities:** Exploration of salary distributions across the top 15 cities, highlighting variations in compensation in key job markets. Bar Plot

e. **Skills Frequency Analysis:** Finding the most often needed abilities across all job postings provides insightful information about the skill requirements in the market for data analysis jobs. Bar Plot

f. **Training Offered Analysis:** investigation of job postings that provide employee training, looking at how pay are split in these situations. Bar Plot and Box Plot

g. **Company-wise Job Listings:** An examination of the businesses with the most job postings offers valuable information about the top employers in the data analysis industry. Bar Plot

6. **Machine Learning:** We added machine learning to our research later on to improve our analysis. A popular classification approach called logistic regression was used to forecast job types according to particular criteria. Here's a condensed explanation:

a. **Feature Selection:** We considered four key features: Salary, Experience, Job Count, and Skills Count. These features were chosen for their relevance in influencing job classifications.

b. **Data Splitting:** To evaluate the model's effectiveness, we divided our dataset into two subsets: 80% for training and 20% for testing.

c. **Model Training:** Our Logistic Regression model was trained using the training subset, where it learned patterns and relationships between the selected features and job classes.

d. **Model Evaluation:** The model's performance was assessed using standard metrics like Accuracy Score and Confusion Matrix, helping us understand how well it predicted job classes.

    **e.** **Interpretation of Results:** The coefficients from the trained model provided insights into the relative importance of each feature in determining job classes.

## Conclusion:

    This extensive data analysis project offers stakeholders unmatched insights into the data analysis job market through a thorough and methodical methodology. The project adeptly tackles the complexities involved in gathering, cleansing, processing, and visualizing data. The comprehensive findings may be used by stakeholders, such as recruiters, job seekers, and industry analysts, to inform strategic decision-making. Together, the complex visualizations, statistical analysis, and machine learning forecasts provide a thorough grasp of the ever-changing field of data analysis jobs. This project goes above and above to provide a thorough and comprehensive examination of the data analysis job market in addition to accomplishing its main purpose of collecting, cleaning, and analyzing pertinent data.