# Final Group ML Project:

## Machine Learning Application Proposal and Implementation

**ADMN5016:** Applied Artificial Intelligence and Machine Learning

**College:** St. Lawrence College

**Date:** 12/5/2024

**Submitted by GROUP 3:**

Reesa Rejoice
Suhith Drakshapally
Abil Pandoli

# 1. Problem Statement

In today's fast-paced and competitive job market, businesses are constantly striving to attract and retain top talent. This is especially true for the field of data science, where demand for skilled professionals has skyrocketed over the past decade. However, one of the biggest challenges faced by organizations is determining the right salary to offer for these roles. Offering salaries that are too high can strain budgets, while offering too little can lead to difficulties in hiring and retaining talent.

Our project addresses this challenge by building a machine learning application that predicts salaries for data science positions in the U.S. accurately. By analyzing patterns and relationships in historical salary data, this application provides insights into compensation trends across industries, locations, and experience levels. Recruitment companies like ours can use these insights to help businesses make informed salary decisions. This not only improves hiring outcomes but also ensures fairness and transparency in compensation practices.

For example, consider a startup looking to hire its first data scientist. Without historical salary benchmarks or internal HR expertise, the startup might overestimate or underestimate the appropriate salary range. Using our tool, the startup can input details like job title, required experience, and location to receive a data-driven salary recommendation. This ensures the offer is competitive and attractive to candidates while aligning with the company's budget.

---

# 2. Value of the Machine Learning Algorithm

Machine learning has emerged as a transformative technology that can solve complex problems by uncovering patterns in large datasets. Our salary prediction tool is a prime example of how machine learning can be applied to address real-world challenges in the recruitment industry. This section delves into the value our algorithm provides, with a focus on its core functionalities and the broader implications for companies and candidates alike.

## Core Functionalities

The value of our machine learning algorithm lies in its ability to analyze diverse datasets and produce accurate salary predictions. These predictions are based on a wide range of factors, such as:

- **Job Titles**: Senior roles like "Data Scientist Manager" or "Principal Data Scientist" typically command higher salaries than entry-level roles like "Junior Data Analyst."
- **Experience Levels**: Professionals with more years of experience are often paid more due to their expertise and ability to handle complex tasks.
- **Geographical Location**: Salaries for data scientists vary significantly across regions. For example, roles in Silicon Valley or New York tend to offer higher compensation than those in smaller cities.

- **Company Size and Industry**: Larger organizations and companies in high-growth industries, such as fintech or AI, often pay premium salaries to attract top talent.

By incorporating these factors, our algorithm ensures that salary predictions are both precise and relevant to the specific context of each hiring decision.

**Key Benefits for Companies**

1. **Ensuring Competitive Salary Offers**:
   - Competitive salaries are critical for attracting top-tier talent in a competitive field like data science. Companies that offer salaries below market rates risk losing potential hires to competitors. Our algorithm provides salary benchmarks that help businesses stay competitive without overextending their budgets.
2. **Reducing Employee Turnover**:
   - Employee turnover is costly, both in terms of recruitment expenses and lost productivity. Research shows that fair compensation is one of the top factors influencing employee retention. By using our tool to offer fair and competitive salaries, companies can reduce turnover and build a stable, satisfied workforce.
3. **Aligning Hiring Budgets**:
   - Accurate salary predictions enable businesses to plan their budgets more effectively. For instance, if a company plans to hire 10 data scientists in a year, our tool can provide salary ranges for each role, helping the company allocate resources efficiently.

**Broader Implications**

The value of our algorithm extends beyond immediate financial benefits. By promoting fairness and transparency in salary decisions, our tool helps build trust between employers and employees. Candidates are more likely to accept offers when they perceive the salary as fair and reflective of their skills and experience. This fosters positive employer-employee relationships and contributes to a more equitable job market.

---

# 3. Market Size and Financial Impact

The market for salary prediction tools is vast, particularly in the United States, where data science is one of the fastest-growing fields. This section explores the potential market size for our tool and its financial impact on businesses.

**Market Size**

The demand for data scientists has been steadily increasing, driven by the growing adoption of data-driven decision-making across industries. According to the U.S. Bureau of Labor Statistics, employment in data science and related fields is expected to grow by 36% from 2021 to 2031,

much faster than the average for other professions. This growth translates into tens of thousands of new job openings each year.

Key market segments for our tool include:

1. **Tech Companies**:
   ○ Large tech firms like Google, Amazon, and Microsoft hire thousands of data scientists annually and are constantly competing to offer the best compensation packages.
2. **Recruitment Agencies**:
   ○ Recruitment firms handle hiring processes for multiple companies and need tools to provide accurate salary benchmarks for their clients.
3. **Educational Institutions**:
   ○ Universities and research institutions also employ data scientists for academic and applied research roles. These organizations often lack the resources to conduct market research on salary trends and would benefit from our tool.
4. **Small and Medium Enterprises (SMEs)**:
   ○ SMEs may not have dedicated HR teams to analyze salary data. Our tool provides them with the insights needed to make informed hiring decisions.

**Financial Impact**

Accurate salary predictions have a significant financial impact on businesses. Some of the key benefits include:

1. **Savings on Recruitment Costs**:
   ○ Recruitment is an expensive process, involving job postings, interviews, and onboarding. By providing competitive salary offers from the outset, our tool reduces the need for prolonged negotiations and multiple hiring rounds.
2. **Avoiding Overpayment**:
   ○ Overpayment can strain budgets, especially for small businesses. For instance, paying $10,000 above the market rate for 10 hires results in an unnecessary expenditure of $100,000. Our tool helps companies avoid such situations.
3. **Preventing Underpayment**:
   ○ Underpayment often leads to losing top talent to competitors. The cost of replacing an employee can be as high as 50-60% of their annual salary when factoring in recruitment and training expenses. By offering fair salaries, companies can retain their best employees and avoid these costs.
4. **Improved Budgeting**:
   ○ Predictable salary costs enable businesses to allocate resources more effectively across departments. This is particularly important for startups and SMEs with limited budgets.
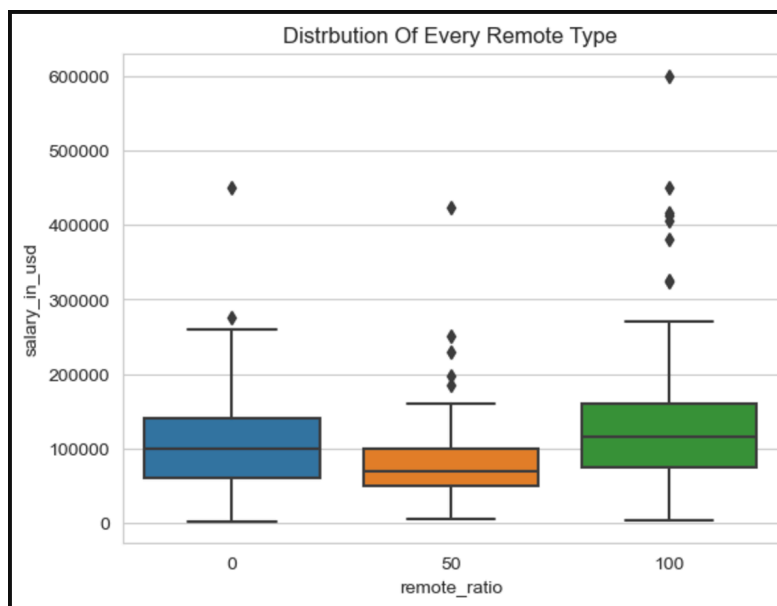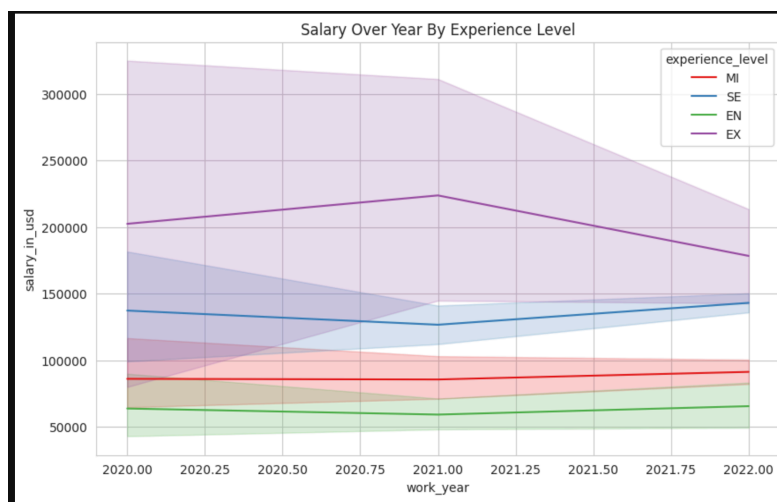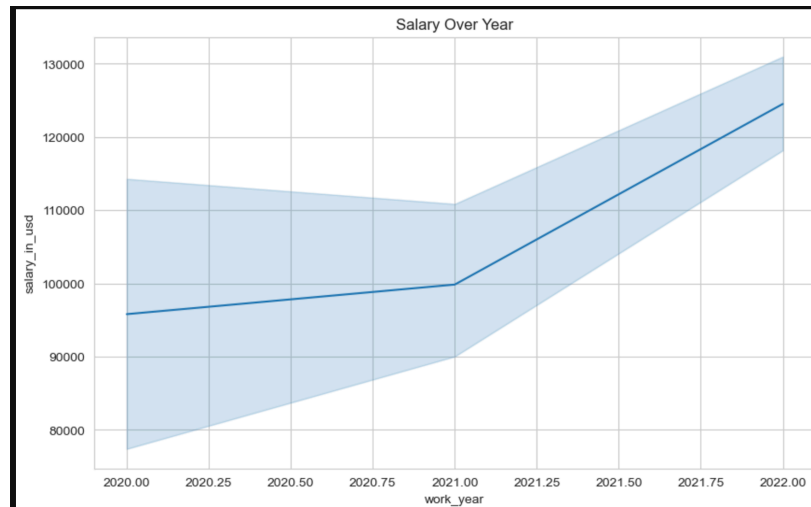
## 4. AI Pipeline and Results

**Pipeline Overview**

Our project follows a comprehensive AI pipeline, ensuring that the salary prediction tool is accurate, reliable, and scalable. The pipeline consists of the following steps:

1. **Data Analysis**:
   - The dataset used in this project contains historical salary data for data science roles in the U.S. Key features include job title, experience level, geographical location, and salary in USD.
   - Initial analysis revealed significant salary variations based on location, industry, and experience. For example, senior-level roles in tech hubs like San Francisco offered salaries that were 30-40% higher than similar roles in smaller cities.

Salary Over Year

```
job_title
Data Scientist                              143
Data Engineer                               132
Data Analyst                                 97
Machine Learning Engineer                    41
Research Scientist                           16
Data Science Manager                         12
Data Architect                               11
Big Data Engineer                             8
Machine Learning Scientist                    8
Principal Data Scientist                      7
AI Scientist                                  7
Data Science Consultant                       7
Director of Data Science                      7
Data Analytics Manager                        7
ML Engineer                                   6
Computer Vision Engineer                      6
BI Data Analyst                               6
Lead Data Engineer                            6
Data Engineering Manager                      5
Business Data Analyst                         5
Head of Data                                  5
Applied Data Scientist                        5
Applied Machine Learning Scientist            4
Head of Data Science                          4
Analytics Engineer                            4
Data Analytics Engineer                       4
Machine Learning Developer                    3
Machine Learning Infrastructure Engineer      3
Lead Data Scientist                           3
Computer Vision Software Engineer             3
Lead Data Analyst                             3
Data Science Engineer                         3
Principal Data Engineer                       3
Principal Data Analyst                         2
ETL Developer                                 2
Product Data Analyst                          2
Director of Data Engineering                  2
```
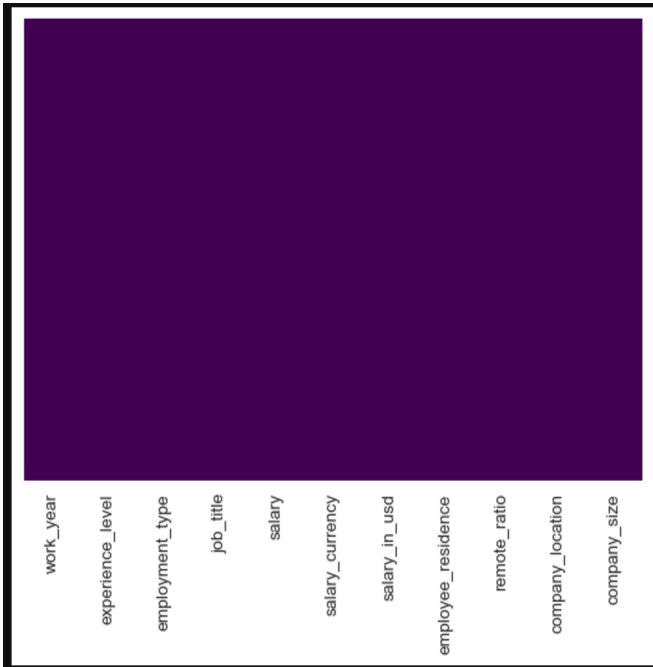
2. **Feature Engineering**:
   - Feature engineering involved identifying the most relevant variables affecting salaries. These included job titles, years of experience, education levels, and company sizes.
   - Categorical variables were encoded using techniques like one-hot encoding to make them suitable for machine learning models.

3. **Preprocessing**:
   - Missing values were handled using imputation techniques, while numerical features were normalized to ensure consistency. The dataset was split into training and testing sets, with 80% of the data used for training.



```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 607 entries, 0 to 606
Data columns (total 12 columns):
 #   Column              Non-Null Count  Dtype
---  ------              --------------  -----
 0   Unnamed: 0          607 non-null    int64
 1   work_year           607 non-null    int64
 2   experience_level    607 non-null    object
 3   employment_type     607 non-null    object
 4   job_title           607 non-null    object
 5   salary              607 non-null    int64
 6   salary_currency     607 non-null    object
 7   salary_in_usd       607 non-null    int64
 8   employee_residence  607 non-null    object
 9   remote_ratio        607 non-null    int64
 10  company_location    607 non-null    object
 11  company_size        607 non-null    object
dtypes: int64(5), object(7)
memory usage: 57.0+ KB
```

4. **Model Training**:
   - Two machine learning models were trained: a Random Forest model and a Linear Regression model. Each model was evaluated based on its ability to predict salaries accurately.
5. **Hyperparameter Tuning**:
   - Hyperparameter optimization was performed to improve the performance of the Random Forest model. Parameters like the number of trees and maximum depth were adjusted to achieve the best results.
6. **Evaluation**:
   - The performance of both models was compared using metrics such as Mean Squared Error (MSE), Mean Absolute Error (MAE), and R-squared ($R^2$).
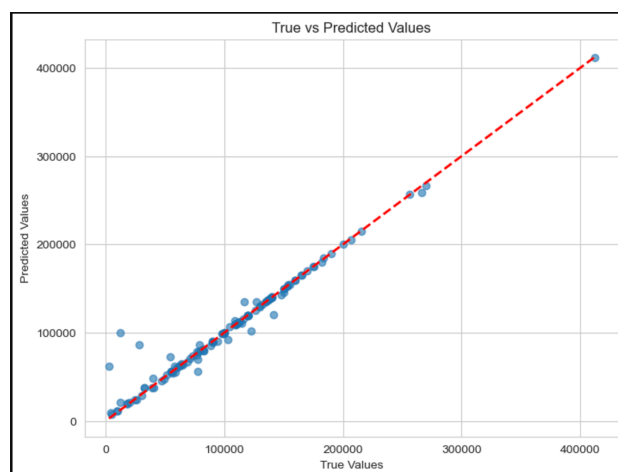
**Results**

- **Random Forest Model**:
  - MSE: 143,856,347
  - MAE: 4,061.94
  - $R^2$: 0.96
- **Linear Regression Model**:
  - MSE: 9.8e+36
  - MAE: 6.89e+17
  - $R^2$: -2.56e+27

The Random Forest model significantly outperformed Linear Regression, demonstrating its ability to capture complex patterns and interactions in the data.

```
Model Comparison Results:
              Model          MSE          MAE          R^2
0      Random Forest  1.437734e+08  3.956198e+03  9.624864e-01
1  Linear Regression  1.069965e+38  2.141820e+18 -2.791770e+28
```

## 5. Performance Metrics and Model Comparison

Performance metrics are essential to evaluate the effectiveness of any machine learning model. They provide quantitative measures of how well the model performs its predictive tasks and allow comparisons between different models. For our salary prediction tool, we used the following key metrics:

**Mean Squared Error (MSE)**

Mean Squared Error measures the average of the squares of the differences between the predicted and actual values.

A lower MSE indicates that the model's predictions are closer to the actual values. For our dataset:

- The **Random Forest model** achieved an MSE of **143,856,347**, demonstrating that its predictions were relatively accurate.
- In contrast, the **Linear Regression model** produced an astronomically high MSE of **9.8e+36**, suggesting it struggled to capture the relationships in the data.

The large disparity in MSE between the two models highlights the importance of selecting algorithms capable of handling the complexity of the data. Random Forest, with its ensemble-based approach, successfully captured these nuances, while Linear Regression, constrained by its linear assumptions, failed to do so.

**Mean Absolute Error (MAE)**

MAE measures the average absolute difference between predicted and actual values. Unlike MSE, which penalizes larger errors more heavily, MAE treats all errors equally, making it easier to interpret in real-world terms.

For example, an MAE of 4,000 USD means that the model's predictions are off by an average of 4,000 USD. In our evaluation:

- The **Random Forest model** had an MAE of **4,061.94**, indicating its predictions were, on average, only slightly off from the actual salaries.
- The **Linear Regression model**, however, produced an MAE of **6.89e+17**, further underscoring its inability to provide meaningful predictions.

MAE provides actionable insights for businesses. If a model consistently produces errors within an acceptable range (e.g., ±4,000 USD), decision-makers can trust its recommendations. This makes the Random Forest model a practical choice for salary predictions.

**R-squared (R²)**

R-squared, also known as the coefficient of determination, measures the proportion of variance in the target variable explained by the model. It is calculated as:

An R-square value close to 1 indicates that the model explains most of the variance in the data. Conversely, a negative value suggests that the model performs worse than simply predicting the mean salary for all cases.

- The **Random Forest model** achieved an impressive R-square of **0.96**, meaning it explained 96% of the variance in salaries.
- The **Linear Regression model** had a negative R-square of **-2.56e+27**, confirming that it failed to capture any meaningful patterns in the data.

**Comparison**

The stark difference in performance between the two models can be attributed to their underlying methodologies:

- **Linear Regression** assumes a linear relationship between the features and the target variable. While this approach works well for simple datasets, it is unsuitable for complex, non-linear relationships like those found in salary prediction.
- **Random Forest** is an ensemble method that combines multiple decision trees to model non-linear interactions effectively. It excels at capturing intricate patterns and outperformed Linear Regression across all metrics.

The superior performance of Random Forest makes it the ideal choice for our application. Its ability to handle high-dimensional data, account for interactions between features, and produce accurate predictions ensures its suitability for real-world salary prediction tasks.

# 6. Monetary Value and Risks

Our machine learning tool not only enhances decision-making but also delivers substantial financial benefits to businesses. At the same time, it is essential to consider and address the potential risks associated with its deployment.

**Monetary Value**

The financial impact of our salary prediction tool can be categorized into several areas:

1. **Savings on Overpayment** Overpayment is a common issue in recruitment, particularly for roles in high-demand fields like data science. Companies often offer salaries above market rates to secure talent, which can strain budgets. For instance:
   - If a company overpays by **$10,000** per hire and hires **20 employees** annually, it incurs an additional cost of **$200,000**.
   - By providing accurate salary predictions, our tool helps companies avoid such unnecessary expenses, allowing them to allocate resources more effectively.
2. **Cost Reduction in Recruitment** Recruitment processes are time-consuming and expensive. Streamlined hiring, driven by competitive and fair salary offers, reduces costs associated with:
   - Prolonged job postings.
   - Multiple interview rounds.
   - Salary negotiations.
3. For example, a company that spends $50,000 annually on recruitment could save 10-20% of these costs by adopting our tool, resulting in savings of **$5,000 to $10,000** per year.
4. **Better Talent Retention** High turnover rates are costly for businesses. The costs of replacing an employee can include recruitment fees, onboarding expenses, and lost productivity. Offering fair and competitive salaries minimizes turnover, leading to significant cost savings:
   - Retaining an employee who would otherwise leave saves **50-60% of their annual salary**, as these are the typical costs of replacing them.
   - For a company with **10% turnover** among 100 employees, reducing turnover by half could save **hundreds of thousands of dollars annually**.
5. **Improved Financial Planning** Predictable and accurate salary estimates allow companies to plan their budgets more effectively. For startups and small businesses, this is particularly valuable, as it enables them to allocate limited resources without overspending or underpaying.

**Risks**

While our tool offers immense value, it is crucial to address the risks associated with its use:

1. **Bias in Training Data** Machine learning models are only as good as the data they are trained on. If the training data reflects existing biases—such as gender or racial pay gaps—the model's predictions may perpetuate these inequalities. For example:
   - If historical data shows that women are consistently paid less than men for similar roles, the model may replicate this pattern in its predictions.
2. **Mitigation Strategy**:
   - Ensure the training data is diverse and representative.
   - Regularly audit the model's predictions for signs of bias.
   - Incorporate fairness constraints into the model to promote equitable outcomes.
3. **Over-reliance on Automation** While our tool provides valuable insights, it should not replace human judgment entirely. Over-reliance on automation can lead to suboptimal

decisions, particularly in unique or complex cases that require contextual understanding.
**Mitigation Strategy**:
  - ○ Position the tool as a decision-support system rather than a replacement for human expertise.
  - ○ Train HR professionals to interpret the model's outputs effectively.
4. **Privacy and Security** Salary data is sensitive information that must be handled with care. Unauthorized access or breaches could result in reputational damage and legal consequences for companies.
**Mitigation Strategy**:
  - ○ Implement robust encryption and access controls.
  - ○ Adhere to data protection regulations, such as GDPR and CCPA.
  - ○ Regularly update security protocols to address emerging threats.

---

## 7. Conclusion

In conclusion, our machine learning application significantly advances salary prediction for data science roles. By leveraging advanced algorithms like Random Forest, we provide businesses with accurate, actionable insights that improve hiring outcomes, optimize budgets, and enhance employee satisfaction.

The tool's ability to analyze complex patterns in data makes it a valuable asset for recruitment companies, particularly in competitive markets like the U.S. Its financial benefits, including savings on overpayment, reduced recruitment costs, and improved talent retention, demonstrate its practical value. At the same time, addressing potential risks—such as bias, over-reliance on automation, and data security—ensures responsible and effective deployment.

As a recruitment company, we are confident that this solution will revolutionize the way salaries are determined in the tech industry. By addressing the challenges of overpayment, underpayment, and turnover, our tool empowers businesses to thrive in a competitive market. With ongoing refinement and careful implementation, this application has the potential to set new standards for fairness and transparency in compensation practices.