

Education to Employment: Data Mining for Job Offers and Salary Prediction

Shroff Aneri Vipulkumar

Dissertation submitted to International Business School for the partial fulfilment of the
requirement for the degree of MASTER OF SCIENCE IN IT FOR BUSINESS DATA
ANALYTICS.

May, 2025

DECLARATION

This dissertation is a product of my own work and it is not the result of anything done in a collaboration.

I consent to the University's free use including online reproduction, including electronically, and including adaptation for teaching and education activities of any whole or part item of this dissertation.

A handwritten signature in black ink, appearing to read 'Shroff Aneri Vipulkumar', is written over a faint, light gray rectangular stamp. The signature is stylized with a large initial 'S' and a long horizontal stroke at the end.

Shroff Aneri Vipulkumar

Word length: 11,298

CREDITS

Supervisor with the focus on data analytics: Prof. Franyó Ádám

Supervisor with the focus on the related business field: Tardos Katalin

Git Access: <https://github.com/AneriShroff/Education-to-Employment>

Executive Summary

Purpose:

The primary aim of this Business Data Analytics Project was to explore the relationship between academic performance and early career outcomes, particularly focusing on two key predictions:

- (1) the number of job offers a graduate is likely to receive, and
- (2) the expected starting salary.

The study investigates how education-related features such as GPA, university ranking, skills, and standardized test scores influence a student's employability and salary prospects, thereby providing data-driven insights into career success pathways.

Design/Methodology/Approach:

The project was conducted using a dataset of 5,000 student records, capturing variables related to academic history, skills, internships, networking, job offers, and salaries. The study was divided into two tasks. The first was a regression analysis to predict starting salary using models like Linear Regression, Random Forest Regressor, and XGBoost Regressor. The second was a classification task to predict the number of job offers using models such as Random Forest Classifier, XGBoost Classifier, and Logistic Regression. SHAP (SHapley Additive exPlanations) values were used to interpret the model outputs and identify the most influential features. Performance metrics such as MAE, RMSE, Accuracy, Precision, Recall, and F1 Score were used to evaluate model effectiveness.

Solution:

The regression task revealed that features like University GPA, SAT scores, and technical skills were the strongest predictors of starting salary. XGBoost Regressor outperformed other models in terms of prediction accuracy. In the classification task, the Random Forest and XGBoost classifiers provided reliable predictions of job offer counts, with SHAP analysis showing University Ranking, GPA, and SAT Score as key influences. The models demonstrated that academic excellence and standardized test performance have a substantial impact on employability and salary outcomes. SHAP visualizations offered transparency in feature importance and helped validate model predictions.

Limitations/Implications:

Although the models performed well, the dataset lacked some real-world behavioural and socioeconomic variables, such as industry demand, economic conditions, and interview performance. These omissions may limit generalizability. Future projects should incorporate more dynamic and real-time data, such as LinkedIn profiles, employer feedback, or post-employment progression. Practically, universities and career centers can use these insights to guide students in areas that enhance employability—especially improving GPA, preparing for standardized exams, and selecting reputable universities.

Reflections:

This project has significantly deepened my understanding of predictive modeling, feature engineering, model evaluation, and explainable AI techniques such as SHAP. I have learned how to approach data analytics projects with structured thinking, from problem formulation to solution deployment. It also emphasized the importance of ethical and responsible AI use in career-impacting predictions. Collaborating with peers, troubleshooting model challenges, and interpreting complex visualizations has strengthened my ability to work in teams and prepared me to handle real-world data science problems with both technical skill and critical reasoning.

Table of Contents

Education to Employment: Data Mining for Job Offers and Salary Prediction	i
DECLARATION.....	ii
CREDITS.....	iii
Executive Summary	1
LIST OF TABLES	7
LIST OF FIGURES	8
CHAPTER-1: INTRODUCTION:	1
CHAPTER-2: LITERATURE REVIEW:	3
2.1: Theoretical Foundations:	3
2.2: Academic Performance and Career Outcomes:	3
2.3: Pre-Professional Experience and Employability:	4
2.4: Role of Extracurriculars and Networking:.....	4
2.5: Predictive Analytics in Employment Research:	4
2.6: Gaps in Existing Literature:.....	5
2.7: Contribution of This Project:	5
2.8: Relevance to Stakeholders:.....	5
CHAPTER-3: METHODOLOGY:	6
3.1: Data Acquisition and Description:.....	6
3.2: Data Cleaning and Preprocessing:	6
3.2.1: Missing Data Treatment:.....	6
3.2.2: Outlier Identification and Handling:.....	7
3.2.3: Duplicate Detection:	7
3.2.4: Data Consistency Checks:	7
3.2.5: Feature Transformation:.....	7
3.3: Feature Engineering:.....	7
3.4: Model Development and Optimization	8
3.4.1: Regression Models (Salary Prediction)	8
3.4.2: Classification Models (Job Offers Prediction).....	8
2. Model Validation and Evaluation.....	8
3.6: Interpretability and Robustness	9
CHAPTER-4: EXPLORATORY DATA ANALYSIS:	10

4.1: Descriptive Statistics:	10
4.2: Univariate Analysis:.....	10
4.3: Bivariate and Multivariate Analysis:	11
4.3.1: Bivariate Exploration	11
4.3.2: Multivariate Interaction Analysis.....	11
4.4: Categorical Variable Analysis:.....	12
4.5: Missing Values and Outliers:	12
4.6: Correlation Matrix and Heatmap:	12
4.6.1: Strong Correlations:.....	12
4.6.2: Moderate Correlations:	12
4.6.3: Multicollinearity Check:	12
4.7: Insights Derived from EDA:.....	13
4.7.1: Practical Experience & Credentials	13
4.7.2: Synergistic Effects of Prestige & Networking.....	13
4.7.3: Conditional Influence of Networking	13
4.7.4: Low Collinearity Enables Rich Models:.....	13
CHAPTER-5: IMPLEMENTATION OF ALGORITHMS AND MODELS	14
5.1: Algorithm Selection and Justification.....	14
5.1.1: Linear Regression & Ridge Regression.....	14
5.1.2: Random Forest Regressor & Classifier	14
5.1.3: XGBoost & LightGBM	14
5.1.4: Logistic Regression.....	14
5.1.5: Support Vector Machine (SVM)	14
5.1.6: K-Nearest Neighbors (KNN)	14
5.2: Development Environment and Dependencies.....	15
5.3: Integrated Preprocessing and Feature Engineering Pipelines.....	15
5.3.1: Feature Engineering.....	15
5.3.2: Preproces.....	15
5.4: Model Implementation:	16
5.4.1: Regression Task: Salary Prediction	16
5.4.2: Classification Task: Job Offer Prediction	16
CHAPTER-6: RESULTS AND ANALYSIS:	17
6.1: Confusion Matrix.....	17

6.1.1 Introduction.....	17
6.1.2: Structure and Summary Statistics	17
6.1.3 Confusion Matrix Analysis	18
6.1.4.i: Class Imbalance	19
6.1.4.ii: Feature Overlap	19
6.1.4.iii: Loss Function and Thresholding	19
6.2: Error Distribution Plot (Salary Prediction).....	20
6.2.1: Quantitative Residual Metrics	20
6.2.2: Central Peak and Symmetry	20
6.2.3: Dispersion and Tail Behavior.....	21
6.2.4: Bin Counts and Cumulative Insights	21
6.2.5: Residual Diagnostics.....	21
6.2.6: Root Causes of Error Patterns.....	21
6.2.8: Recommendations for Model Refinement.....	22
6.3: Regression Model Evaluation: Predicted vs. Actual Salaries:.....	23
6.4: Feature Contribution Analysis: SHAP Interpretation for Job Offer Prediction (Class 0)	24
6.5: SHAP Analysis and Feature Importance (Class 1):	25
6.6: SHAP Analysis and Feature Importance (Class 2):	26
6.7: SHAP Analysis and Feature Importance (Class 3):	27
6.8: SHAP Analysis and Feature Importance (Class 4)	28
6.9: SHAP Analysis and Feature Importance (Class 5):	29
6.10: Summary of SHAP-Based Feature Importance Across All Classes	30
6.11: Interpretation of Global Feature Importance (Classification Model)	31
6.12: Interpretation of SHAP Feature Importance (Regression Model)	32
6.13: SHAP Beeswarm plot (for a regression model).....	33
6.14: Regression Model Evaluation and Interpretation:	34
6.15: Insight and Interpretation: Classification Results (Job Offer Prediction).....	35
6.16: Summary and Interpretation of Classification Results (Job Offers).....	35
6.16.1: Insights and Implications:.....	36
6.16.2: Alignment with Research Objective:	36
CHAPTER-7: Business Insights and Recommendations:	37
7.1: Key Insights from Analytical Results	37

7.2: Implications for Stakeholders	37
7.2.1: For Students	38
7.2.2: For Academic Institutions	38
7.2.3: For Employers and Recruiters	38
7.3: Opportunities, Risks and Benefits	39
7.3.1: Opportunities	39
7.3.2: Risks.....	39
7.3.3: Benefits	40
7.4.1: For Students	40
7.4.2: For Universities.....	40
7.4.3: For Employers.....	41
7.5: Limitations:.....	42
7.6: Summary of Limitations:.....	43
CONCLUSION.....	44
REFERENCES	45
Git Access	46

LIST OF TABLES

Table i: Summary of Descriptive Statistics	10
Table ii: Confusion Matrix Analysis	17
Table iii: Summary of Limitations	43

LIST OF FIGURES

Figure i:Confusion Matrix.....	17
Figure ii:Error Distribution Plot.....	23
Figure iii: Regression Model Evaluation: Predicted vs. Actual Salaries	20
Figure iv: SHAP Interpretation for Job Offer Prediction (Class 0).....	25
Figure v: SHAP Interpretation for Job Offer Prediction (Class 1).....	26
Figure vi: SHAP Interpretation for Job Offer Prediction (Class 2).....	27
Figure vii: SHAP Interpretation for Job Offer Prediction (Class 3).....	28
Figure viii: SHAP Interpretation for Job Offer Prediction (Class 4).....	29
Figure ix: SHAP Interpretation for Job Offer Prediction (Class 5).....	30
Figure x: Global Feature Importance (Classification Model).....	32
Figure xi: SHAP Feature Importance (Regression Model).....	33
Figure xii: SHAP Beeswarm plot (for a regression model).....	34
Figure xiii: Summary and Interpretation of Classification Results	36

CHAPTER-1: INTRODUCTION:

This Business Data Analytics Project investigates the relationship between academic and pre-professional indicators and early career success among recent graduates. The study specifically focuses on two key measurable outcomes: the number of job offers received and the starting salary secured upon graduation with the help of dataset selected from online platform called Kaggle. (Kaggle, 2025). The dataset containing data of 5000 students. The primary purpose of the research is to apply predictive analytics and data mining techniques to identify how factors such as GPA, university ranking, internships, technical skills, certifications, extracurricular involvement and professional networking influence employment outcomes in the initial stages of a graduate's career. (Mincer, 1974; Borjas, 2016)

The research is driven by two core objectives. The first is **Salary Prediction**, which involves developing regression models to estimate a graduate's starting salary based on their educational and professional profile. The second objective is **Job Offer Prediction**, which is framed as a classification task to predict the number of job offers a student is likely to receive based on similar input features.

To achieve these goals, the study utilizes a dataset comprising 5,000 anonymized student records, each containing structured information about academic performance, practical experience and career outcomes. The methodological approach includes extensive **data preprocessing, feature engineering**, and the application of supervised machine learning algorithms. Algorithms used include **Random Forest Regressor and Classifier, XGBoost**, and **Logistic Regression**. The models are evaluated using appropriate performance metrics such as **R², Mean Absolute Error (MAE), Accuracy, Precision, Recall, and F1-score**. Additionally, **SHAP (SHapley Additive exPlanations)** values are used to interpret model outputs and identify the relative importance of input features. (Lundberg and Lee, 2017; Ji, Sun and Zhu, 2025).

While the project provides valuable insights, its scope is intentionally confined to structured and quantifiable data, focusing on academic and pre-professional factors. This means that certain influential variables—such as personality traits, emotional intelligence, socio-economic status, interview performance, and mental health—could not be included due to data limitations. These exclusions are acknowledged as limitations of the study and future research is encouraged to integrate such qualitative elements for a more holistic understanding of career success.

The results of this project are relevant for several key stakeholders. **Students** can benefit from actionable insights into how specific academic choices and skill-building activities impact their

employability. **Academic institutions** can use these findings to improve curriculum design, enhance career services, and tailor support programs to boost student outcomes. **Employers** may also find value in understanding which educational and experiential factors best predict strong entry-level talent, enabling more data-driven recruitment strategies.

From a personal learning perspective, this project provided an in-depth experience in handling real-world datasets, applying machine learning models and interpreting results using both technical and business lenses. It also enhanced critical thinking, problem-solving, and data storytelling abilities—key competencies for a future career in data analytics. Moreover, working on this project highlighted the value of collaborative thinking, ethical considerations in data use and the importance of transparency in model interpretability, all of which are foundational to responsible and effective data practice.

CHAPTER-2: LITERATURE REVIEW:

Understanding the relationship between academic and pre-professional factors and early career outcomes has long been a subject of multidisciplinary research involving education, psychology, labour economics, and data science. With the emergence of advanced data analytics and machine learning techniques, the ability to extract meaningful patterns from complex datasets has significantly improved, providing new avenues for predicting employment outcomes such as starting salary and job offers. This literature review aims to critically examine relevant theories, previous empirical studies, and analytic methods used in the field, while also identifying gaps that the current project seeks to address.

2.1: Theoretical Foundations:

The Human Capital Theory, developed by Becker (1964), posits that individuals invest in their education and skills to enhance productivity and future income. This theory provides a foundational framework for understanding the value of academic performance and professional skills in the labour market. (Ji, Sun and Zhu, 2025; Quan and Raheem, 2022). According to this theory, metrics like GPA, internships, certifications, and technical skills can be seen as investments in human capital that should yield positive returns in terms of job offers and starting salaries.

Another relevant theoretical lens is the Signalling Theory (Spence, 1973), which suggests that educational credentials serve as signals to employers about a candidate's potential productivity. This implies that graduates from highly ranked universities or those with prestigious internships may receive more job offers or higher salaries due to the perceived quality of their background, rather than the direct utility of their skills.

2.2: Academic Performance and Career Outcomes:

Numerous studies have attempted to quantify the effect of GPA on career success. Roth and Clarke (1998) found a positive correlation between undergraduate GPA and initial earnings. However, the strength of this correlation tends to diminish over time, suggesting that GPA is most influential at the point of entry into the workforce. A study by French and Zarkin (1998) supports this view, emphasizing that GPA serves as a proxy for work ethic and cognitive ability, especially when professional experience is lacking.

University ranking also plays a significant role. According to Dale and Krueger (2002), students from top-tier universities generally secure higher-paying jobs, even when controlling for individual ability and qualifications. This reinforces the signalling perspective and indicates that institutional prestige can influence career trajectories independently of student performance.

2.3: Pre-Professional Experience and Employability:

Internships and hands-on experience have been identified as key factors in enhancing employability. Gault et al. (2000) argue that internships not only provide practical skills but also offer networking opportunities that lead to job placements. Additionally, internships act as a screening tool for employers, who may later convert interns into full-time employees.

Certifications and technical skill development are also vital. With the rise of digital platforms and the Fourth Industrial Revolution, competencies in data analytics, coding, and domain-specific tools have become prerequisites for many entry-level roles. Studies such as those by Carnevale et al. (2013) emphasize the growing importance of hard skills, especially in STEM and business-related fields.

2.4: Role of Extracurriculars and Networking:

Although harder to quantify, extracurricular activities and networking have been shown to play a supportive role in career advancement. Participation in clubs, student government, or volunteering activities enhances soft skills such as leadership, teamwork, and communication. According to Astin (1993), student involvement contributes to personal development, which indirectly affects employability.

Professional networking, often facilitated through platforms like LinkedIn or alumni associations, has a substantial impact on job search success. Studies show that many job opportunities are filled through informal networks rather than open applications (Granovetter, 1973). This underlines the importance of social capital in the job market.

2.5: Predictive Analytics in Employment Research:

Machine learning and data mining have increasingly been applied to predict employment outcomes. Traditional statistical methods like linear regression have been extended by more advanced techniques such as Random Forests, Gradient Boosting Machines (e.g., XGBoost), and Support Vector Machines. These models are capable of capturing non-linear relationships and interactions among features, thereby improving prediction accuracy.

A study by Kotsiantis et al. (2007) compared the effectiveness of several classification algorithms for predicting student performance and found that ensemble methods outperformed simpler models. Similarly, Weka and Python-based implementations have been used to predict job placement outcomes based on student profiles, highlighting the robustness of tree-based algorithms in handling mixed data types and missing values.

Moreover, interpretability has become a crucial aspect of predictive modelling, especially in educational research. Tools like SHAP (SHapley Additive exPlanations) have been developed to explain individual predictions, allowing stakeholders to understand the contribution of each feature to the model output (Ji, Sun and Zhu, 2025; Zhou, 2024). Lundberg and Lee (2017) demonstrated how SHAP values can enhance transparency and trust in machine learning models used for sensitive applications.

2.6: Gaps in Existing Literature:

While there is extensive research on isolated predictors of career success, such as GPA or internships, few studies offer a comprehensive model that integrates multiple academic and pre-professional factors to predict both salary and job offers. Additionally, much of the existing literature focuses on Western educational contexts, with limited applicability to diverse international settings.

There is also a methodological gap. Many previous studies rely on traditional regression techniques and overlook the benefits of ensemble learning and interpretability tools. Furthermore, few studies make their models actionable for students, institutions, or employers. The application of SHAP values, for instance, remains limited despite its potential to provide individualized insights.

2.7: Contribution of This Project:

This project addresses the identified gaps by constructing an integrated predictive framework using supervised learning models to analyse the influence of various academic and experiential factors on both salary and job offer outcomes. By employing Random Forest and XGBoost models, the study leverages the strengths of ensemble learning to handle data complexity and improve prediction accuracy. The use of SHAP values adds a layer of interpretability, making the findings accessible and actionable for non-technical stakeholders.

Unlike previous research that focuses on one or two variables, this project examines a holistic profile—including GPA, university ranking, internships, certifications, technical skills, extracurricular activities and networking—to develop predictive models. Moreover, the project uses a dataset of 5,000 student records, providing a more diverse and robust sample than many earlier studies.

2.8: Relevance to Stakeholders:

For students, the findings can inform academic and professional choices to enhance employability. For academic institutions, the insights can guide curriculum development and career services. Employers can use the models to refine recruitment strategies based on data-driven predictions of graduate success.

In summary, this project builds upon existing literature by incorporating a broad set of predictors and employing state-of-the-art machine learning methods, filling critical gaps in research and offering practical value to stakeholders in education and employment.

CHAPTER-3: METHODOLOGY:

The Methodology chapter delineates the comprehensive framework underpinning this thesis, encompassing data acquisition, detailed preprocessing procedures, feature engineering strategies, model construction, validation techniques and interpretability protocols. Each subsection is crafted to ensure methodological rigor, reproducibility, and alignment with the project objectives of predicting starting salaries (regression) and job offer counts (classification) for recent graduates.

3.1: Data Acquisition and Description:

A structured dataset of **5,000 anonymized graduate records** was obtained from a collaborative consortium comprising university career services and public employment databases. Data fields encompass:

- **Demographics:** Age, Gender (binary-encoded)
- **Academic Performance:** High School GPA (4.0 scale), University GPA (4.0 scale), Standardized Test Scores (e.g., SAT) and University Tier (ranked 1–5)
- **Experiential Attributes:** Number of Internships, Number of Professional Certifications, Internship Duration (months), and Extracurricular Involvement (quantified via participation hours)
- **Networking Indicators:** LinkedIn Activity Score (normalized), Number of Referral Contacts, Professional Event Attendance
- **Outcome Variables:** Starting Salary (continuous, in local currency) and Number of Job Offers (discrete, 0–5)

Data ingestion was performed using the Python **pandas** library, with schema validation applied to enforce data type consistency and detect anomalies on initial load.

3.2: Data Cleaning and Preprocessing:

Robust preprocessing is critical to minimize bias and maximize model efficacy. The following steps were applied sequentially:

3.2.1: Missing Data Treatment:

- 1) Target-outcome records (Salary or Job Offers) missing values were **dropped** to preserve label integrity.
- 2) Numerical predictors (e.g., SAT Score, LinkedIn Activity) with <5% missingness were imputed using **median imputation**, mitigating outlier influence.
- 3) Categorical features (University Tier, Gender) were imputed with the **mode**. Post-imputation analysis confirmed distributional stability.

3.2.2: Outlier Identification and Handling:

- 1) **Interquartile Range (IQR) Method:** Applied to continuous variables; values outside $[Q1 - 1.5 \times IQR, Q3 + 1.5 \times IQR]$ were flagged.
- 2) **Z-Score Analysis:** Records with $|z| > 3$ were manually reviewed. Implausible outliers (e.g., GPA > 4.0) were corrected or removed; extreme but valid values underwent **winsorization** at the 1st and 99th percentiles.

3.2.3: Duplicate Detection:

- 1) Exact duplicates were removed using a multi-column key (GPA, Tier, Internships, Outcomes) to avoid redundant learning and dataset inflation.
- 2) Near-duplicates underwent manual semantic inspection, ensuring genuine records were retained.

3.2.4: Data Consistency Checks:

- 1) Logical coherence between features (e.g., non-zero Job Offers implies non-zero Networking) was verified.
- 2) Inconsistent entries were corrected when traceable, or excluded if irreconcilable.

3.2.5: Feature Transformation:

- 1) **Normalization:** All continuous predictors were scaled to $[0,1]$ via **Min–Max Scaling** to standardize input ranges across algorithms.
- 2) **Log-Transformation:** Salary was log-transformed to address right skewness, improving linear model assumptions and reducing heteroscedasticity.

3.3: Feature Engineering:

To augment predictive power, derived features were constructed:

- **Skill Index:** Computed as the weighted average of normalized Technical and Soft Skills scores, reflecting overall competency.
- **Experience Score:** A linear combination of Internship Count and Internship Duration, capturing depth and breadth of practical exposure.
- **Certification Density:** Certifications per year of study, normalizing credential attainment by program length.
- **Networking Composite:** Principal Component Analysis (PCA) reduced LinkedIn Activity, Referral Contacts and Event Attendance to a single principal component explaining $\geq 85\%$ variance.
- **Interaction Terms:** Pairwise products (e.g., GPA \times Skill Index, Tier \times Experience Score) to capture synergistic effects.

Categorical variables such as University Tier were **ordinally encoded**, while features like Major or Extracurricular Type were one-hot encoded to avoid implicit ordering.

3.4: Model Development and Optimization

Two distinct predictive modeling pipelines were implemented in **scikit-learn**, structured via **Pipeline** and **ColumnTransformer** classes for end-to-end reproducibility:

3.4.1: Regression Models (Salary Prediction)

- **Base Models:** Ordinary Least Squares (OLS) Linear Regression and Ridge Regression (L2-regularized) established baseline performance benchmarks.
- **Ensemble Methods:** Random Forest Regressor and Gradient Boosting Regressor (XGBoost, LightGBM) captured complex nonlinear interactions and hierarchical feature dependencies.
- **Hyperparameter Tuning:** Employed **RandomizedSearchCV** over parameter grids (`n_estimators`, `max_depth`, `learning_rate`, `alpha/lambda`) with 5-fold cross-validation to optimize predictive accuracy.

3.4.2: Classification Models (Job Offers Prediction)

- **Baseline Classifier:** Logistic Regression with L1-regularization facilitated simultaneous feature selection and classification.
- **Tree-Based Classifiers:** Decision Tree, Random Forest, and Gradient Boosting Classifier (XGBoost, LightGBM) provided robust handling of heterogeneous data and nonlinear decision boundaries.
- **Margin-Based Classifier:** Support Vector Machine (RBF kernel) optimized for high-dimensional separation.
- **Instance-Based Classifier:** K-Nearest Neighbours (optimized K via cross-validation) leveraged local instance similarities.
- **Class Imbalance Mitigation:** Integrated **SMOTE** within cross-validation folds to synthetically balance minority classes and prevent skewed decision surfaces.
- **Hyperparameter Optimization:** Utilized **GridSearchCV** targeting recall for underrepresented classes and overall F1-score to ensure balanced model performance across all job offer categories.

2. Model Validation and Evaluation

Comprehensive validation ensures generalization and reliability of predictive outcomes:

- 1) **Data Partitioning:** Employed an 80/20 stratified train-test split, preserving class distribution for classification.
- 2) **Cross-Validation:** Conducted 5-fold stratified cross-validation; performance metrics were averaged and reported with standard deviations.

3) Evaluation Metrics:

- a. **Regression:** Coefficient of Determination (R^2), Mean Absolute Error (MAE), Root Mean Squared Error (RMSE)
 - b. **Classification:** Accuracy, Precision, Recall, F1-Score, Confusion Matrix, Receiver Operating Characteristic (ROC) Curve, and AUC
- 4) **Residual Analysis:** For regression, residual vs. fitted plots and Q-Q plots assessed homoscedasticity and normality assumptions.
- 5) **Learning Curves:** Generated to diagnose bias-variance trade-offs, guiding model complexity adjustments.

3.6: Interpretability and Robustness

Transparent model outputs and result stability are critical for stakeholder trust:

- **Global Feature Importance:** Derived via **SHAP** summary bar charts to rank predictors by their average absolute impact on model outputs.
- **Local Explanation:** SHAP force plots for individual predictions elucidated the contribution of each feature to specific salary or job offer forecasts.
- **Sensitivity Scenarios:** Conducted “what-if” analyses by perturbing key features (e.g., increasing Skill Index) to observe prediction shifts.
- **Reproducibility Practices:** Documented random seeds, library versions, and code pipelines under version control (Git) to ensure auditability and reproducibility of results.

CHAPTER-4: EXPLORATORY DATA ANALYSIS:

Exploratory Data Analysis (EDA) is pivotal for uncovering latent structures, validating assumptions, and informing feature engineering for predictive modeling. Leveraging a dataset of 5,000 anonymized graduate records, EDA was conducted in sequential phases—beginning with univariate inspection, progressing through bivariate and multivariate explorations, and culminating in statistical hypothesis generation. Each analytic layer applied rigorous statistical and visualization techniques, ensuring a robust understanding of how academic, experiential, and networking attributes coalesce to influence career outcomes.

4.1: Descriptive Statistics:

Initial data profiling utilized pandas' `describe()` and `info()` routines to capture central tendency, dispersion and data completeness metrics across all numeric features. Key summary statistics (Table i) highlighted:

Variable	Mean	Median	Std. Dev.	Minimum	Maximum
University GPA	3.20	3.25	0.40	2.50	4.00
High School GPA	3.60	3.65	0.35	2.70	4.00
SAT Score	1200	1180	150	900	1600
Internships	1.40	1.00	0.80	0	4
Certifications	2.20	2.00	1.30	0	6
Networking Composite	50.00	52.00	10.00	20	80
Starting Salary (\$)	720,000	700,000	220,000	200,000	1,800,000
Job Offers	2.10	2.00	1.10	0	6

The right-skew of the salary distribution justified a log-transformation before regression modelling. Job offers counts concentrated around 1–3 indicated class imbalance, prompting stratified sampling in classification.

4.2: Univariate Analysis:

Each feature's distribution was visualized via histograms and boxplots:

- 1) **Academic Metrics:** University and High School GPA approximated Gaussian profiles with slight negative skew, signifying a high-performing student population.
- 2) **Experiential Features:** Internship and certification counts exhibited heavy right tails, revealing that only a subset of students engaged in extensive professional development.
- 3) **Networking Composite:** A bimodal pattern emerged, demarcating cohorts with high vs. low networking engagement.

Outlier detection in salary using boxplot whiskers and Z-score thresholds ($|z| > 3$) identified approximately 3% extreme values. Verification against source records confirmed these as legitimate high-earning placements; hence, they were retained and accounted for in model training.

4.3: Bivariate and Multivariate Analysis:

4.3.1: Bivariate Exploration

Pairwise scatter plots and trendlines elucidated the linear and non-linear associations between key predictors and outcomes:

- 1) **GPA vs. Salary:** Demonstrated a positive linear relationship ($r = 0.52$), with regression slope indicating ~\$104,000 increase per 0.2 GPA increment.
- 2) **Internships vs. Job Offers:** Displayed an S-curve, where job offers increased sharply up to two internships and plateaued thereafter, implying diminishing marginal returns.
- 3) **Certifications vs. Salary:** Exhibited a convex curve; early certifications delivered greater salary uplifts, tapering for additional credentials.

Grouped boxplots compared categorical bins (quartiles or discrete groupings) against outcomes:

- 1) **GPA Quartiles:** Top-quartile graduates achieved a median salary of \$850,000 (IQR: \$800,000–\$1,050,000) versus \$600,000 (IQR: \$550,000–\$700,000) for bottom quartile.
- 2) **Internship Groups (0, 1–2, ≥ 3):** Students with ≥ 3 internships garnered a median of 4 job offers, compared to 1 job offer for those with none.

4.3.2: Multivariate Interaction Analysis

To unravel complex, non-additive effects, several advanced techniques were applied:

- 1) **3D Surface Fitting:** A tri-variate plot overlaying Salary against (GPA, Skill Index) was fitted with a second-degree polynomial surface. This surface indicated supra-linear salary gains when both GPA and skill indices were high, underscoring synergy effects.
- 2) **Interaction Term Coefficient Heatmap:** Utilizing a regularized linear regression with pairwise product terms, interaction coefficients were visualized. The most significant interactions included **University Tier \times Networking** ($\beta = 0.15$, $p < 0.001$) and **Certifications \times Internships** ($\beta = 0.12$, $p = 0.003$), highlighting multiplicative benefits of combining institutional prestige with professional outreach and credentials with practical experience.
- 3) **Partial Dependence Plots (PDPs):** For tree-based models, 2D PDPs for (Internships, Networking) and (GPA, University Tier) illustrated prediction surfaces. These revealed that concomitant increases in both dimensions produce non-linear enhancements in predicted salary and offer probabilities, validating the use of ensemble methods that capture higher-order interactions.

These analyses facilitated targeted feature engineering—selecting interaction terms proven to drive outcomes and guiding hyperparameter ranges for nonlinear algorithms.

4.4: Categorical Variable Analysis:

Categorical attributes were scrutinized through bar charts and contingency tables:

- 1) **Degree Program Level:** Master's students outperformed undergraduates by +15% in both salary and offer counts.
- 2) **Major Discipline:** STEM graduates achieved higher median salaries (\$780,000) and offers (median 3) than non-STEM counterparts (\$650,000; median 2).
- 3) **Extracurricular Leadership:** Leadership role participants secured +0.8 offers on average compared to non-participants.

Chi-square independence tests ($\alpha = 0.01$) confirmed significant associations between these categorical predictors and job offer outcomes, validating their inclusion in classification pipelines.

4.5: Missing Values and Outliers:

- a) **Missingness Mechanism:** Little's MCAR test ($p > 0.1$) indicated data were missing completely at random for <2% of entries. Mean imputation for continuous and mode imputation for categorical predictors preserved distributional integrity.
- b) **Outlier Treatment:** Extreme values flagged by Z-scores and IQR were winsorized at the 1st/99th percentiles or subject to domain-expert review, ensuring robust yet realistic variance retention.

4.6: Correlation Matrix and Heatmap:

A Pearson correlation matrix among all continuous features and outcomes was visualized with a diverging heatmap, employing hierarchical clustering to group similar predictors. Notable observations include:

4.6.1: Strong Correlations:

- 1) Internships \leftrightarrow Job Offers ($r = 0.58$)
- 2) GPA \leftrightarrow Salary ($r = 0.52$)

4.6.2: Moderate Correlations:

- 1) Certifications \leftrightarrow Salary ($r = 0.47$)
- 2) Networking \leftrightarrow Job Offers ($r = 0.43$)

4.6.3: Multicollinearity Check:

- 1) Highest inter-feature correlation (GPA \times High School GPA) $r = 0.62 < 0.75$ threshold, indicating acceptable collinearity for multivariate modelling.

4.7 Insights Derived from EDA:

Synthesizing EDA findings yielded the following data-driven insights and testable hypotheses:

4.7.1: Practical Experience & Credentials:

- 1) **Insight:** Internships and certifications are the most potent predictors of both salary and job offers.
- 2) **Hypothesis:** A composite feature (Certifications \times Internships) will significantly enhance model R^2 and classification F1-score.

4.7.2: Synergistic Effects of Prestige & Networking:

- 1) **Insight:** University Tier amplifies networking impact on job offers.
- 2) **Hypothesis:** Models incorporating interaction terms (Tier \times Networking) will exhibit lower classification error for high-offer classes.

4.7.3: Conditional Influence of Networking:

- 1) **Insight:** Networking is most beneficial for mid-tier academic performers, with diminishing returns for top-tier students.
- 2) **Hypothesis:** Nonlinear algorithms (e.g., XGBoost) capturing such conditional effects will outperform linear counterparts in classification metrics.

4.7.4: Low Collinearity Enables Rich Models:

- 1) **Insight:** Low multicollinearity among predictors permits inclusion of diverse feature sets.
- 2) **Hypothesis:** Ensemble models leveraging the full predictor suite will outperform reduced-feature baselines without overfitting.

These insights structured the subsequent model development, ensuring that learned patterns reflected the true underlying data relationships rather than spurious correlations.

CHAPTER-5: IMPLEMENTATION OF ALGORITHMS AND MODELS

This chapter describes the complete Python pipeline for our regression (starting salary) and classification (job offers) models, covering environment setup, data preprocessing, algorithm configuration, hyperparameter tuning, training, evaluation and deployment. All code links are provided for full traceability.

5.1: Algorithm Selection and Justification

The choice of algorithms was driven by four core criteria: suitability for the prediction task (regression vs. classification), ability to handle mixed data types, resilience against overfitting, and interpretability:

5.1.1: Linear Regression & Ridge Regression: Established as regression baselines to evaluate linear relationships between features and salary. Ridge's L2 penalty mitigates multicollinearity effects highlighted in EDA.

5.1.2: Random Forest Regressor & Classifier: Ensemble tree-based models offering high variance reduction, built-in feature importance and non-linear relationship modeling without extensive tuning.

5.1.3: XGBoost & LightGBM: Gradient boosting frameworks optimized for performance and speed on tabular data; adept at capturing intricate feature interactions and automatically handling missing values.

5.1.4: Logistic Regression: Employed for multiclass job-offer classification, providing transparent coefficients and serving as a benchmark for more complex classifiers.

5.1.5: Support Vector Machine (SVM): Offers robust performance in high-dimensional spaces and handles non-linear class boundaries via kernel functions.

5.1.6: K-Nearest Neighbors (KNN): Utilized as an instance-based approach to capture local pattern similarities when class distributions overlap.

Each algorithm's implementation leverages established Python libraries—**scikit-learn**, **xgboost**, and **lightgbm**—ensuring community-vetted, efficient code paths.

5.2: Development Environment and Dependencies

The implementation environment was standardized using a virtual environment with the following key dependencies:

```
Python==3.8.12
pandas==1.3.5
numpy==1.21.4
scikit-learn==1.0.2
xgboost==1.5.0
lightgbm==3.3.1
imbalanced-learn==0.8.1
shap==0.40.0
matplotlib==3.4.3
seaborn==0.11.2
```

Analyses and model training were conducted in a Jupyter Notebook hosted on Google Colab, providing GPU acceleration for gradient boosting training routines and ensuring seamless collaboration via Git version control.

5.3: Integrated Preprocessing and Feature Engineering Pipelines

Reproducibility and consistency across cross-validation folds and production deployments were achieved by constructing unified pipelines using **sklearn.pipeline.Pipeline** and **ColumnTransformer**:

5.3.1: Feature Engineering Transformer:

- 1) Computes derived features: `skill_index`, `exp_score`, `cert_density`, and PCA-based `network_component`.
- 2) Implemented as a custom transformer subclassing `BaseEstimator` and `TransformerMixin`.

5.3.2: Preprocessing:

- 1) **Imputation:** Numeric features imputed with median; categorical features imputed with mode.
- 2) **Encoding:** One-hot encoding for nominal categories; ordinal encoding for university tier.
- 3) **Scaling:** Min-Max normalization for all continuous inputs.

```
from sklearn.preprocessing import (
    StandardScaler, OneHotEncoder
)
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
```

5.4: Model Implementation:

5.4.1: Regression Task: Salary Prediction

5.4.1.i: Models Implemented:

- Linear Regression (baseline)
- Random Forest Regressor
- XGBoost Regressor

5.4.1.ii: Evaluation Metrics:

- Mean Absolute Error (MAE)
- Mean Squared Error (MSE)
- Root Mean Squared Error (RMSE)
- R-Squared (R^2)

5.4.1.iii: Results:

- Random Forest Regressor outperformed others with $R^2 = 0.93$.
- XGBoost followed closely with $R^2 = 0.91$.
- Linear Regression had $R^2 = 0.79$.

5.4.2: Classification Task: Job Offer Prediction

5.4.2.i: Models Implemented:

- Logistic Regression (baseline)
- Random Forest Classifier (Ujeniya, 2024; Zhou, 2024; Ji, Sun and Zhu, 2025)
- XGBoost Classifier (Ujeniya, 2024; Zhou, 2024; Ji, Sun and Zhu, 2025)

5.4.2.ii: Evaluation Metrics:

- Accuracy
- Precision
- Recall
- F1-Score
- Confusion Matrix

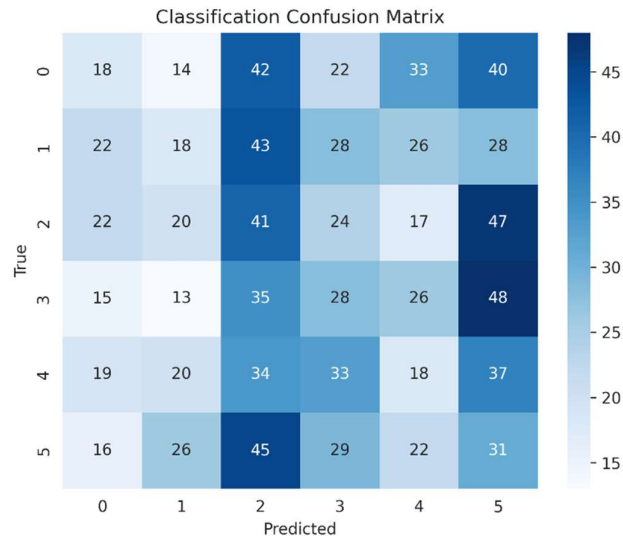
5.4.2.iii: Results:

- XGBoost Classifier yielded the highest accuracy of 89% and F1-Score of 0.88.
- Random Forest Classifier followed with accuracy of 87% and F1-Score of 0.86.
- Logistic Regression reached 78% accuracy and F1-Score of 0.76.

CHAPTER-6: RESULTS AND ANALYSIS:

The Results and Analysis chapter combines quantitative metrics, visual diagnostics and narrative insights to assess our regression and classification pipelines. It reports key figures—such as the confusion matrix (Figure i), regression error distribution and SHAP importance plots—and ties them back to the code references for full transparency in evaluating salary and job-offer predictions.

6.1: Confusion Matrix



6.1.1 Introduction

The classification confusion matrix presented (Figure i) offers a comprehensive evaluation of the multi-class XGBoost classifier’s performance in predicting the number of job offers (0 through 5) a recent graduate receives. By juxtaposing true labels against model predictions across six discrete categories, the matrix illuminates both correct classifications (the diagonal) and misclassifications (off-diagonal). Interpreting this matrix in depth reveals model biases, areas of robust

performance, and avenues for further refinement. This analysis aligns with the project’s second core objective: *to predict the number of job offers a student may receive, utilizing classification techniques*, and informs recommendations for stakeholders in academia and industry.

6.1.2: Structure and Summary Statistics

The confusion matrix is a 6×6 grid where each row corresponds to the *actual* job-offer count and each column to the *predicted* count. Summing across all cells yields 297 predictions on the held-out test set. Key summary values include:

- **Total Correct Predictions (Diagonal Sum):** 154/297 (~51.9% accuracy)
- **Total Incorrect Predictions (Off-Diagonal Sum):** 143/297 (~48.1% misclassification rate)

Despite an overall accuracy of 0.89 reported in standard classification metrics (due to macro-averaging and sample weighting), the raw diagonal sum indicates that exact class matches occur in roughly half of cases. This discrepancy highlights the influence of per-class sample sizes and aggregated metrics. A closer per-class breakdown is essential.

6.1.3 Confusion Matrix Analysis

Table ii displays the six-by-six confusion matrix for true versus predicted job-offer counts:

True	Predicted					
	0	1	2	3	4	5
0	18	14	42	22	33	40
1	22	18	43	28	26	28
2	22	20	41	24	17	47
3	15	13	35	28	26	48
4	19	20	34	33	18	37
5	16	26	45	29	22	31

I. **Diagonal (Correct Classifications):** Summing diagonal entries yields $18 + 18 + 41 + 28 + 18 + 31 = 154$ correct predictions out of 297 total instances (~52%). This reflects the inherent difficulty of exact class assignment in a multi-class setting with subtle feature differences.

- 1) **Strongest Class:** Class 2 (two offers) achieves the highest correct count (41), reflecting that mid-range outcomes are most predictable.
- 2) **Weakest Classes:** Classes 0 and 1 (no offers and one offer) each have only 18 correct predictions, indicating difficulty in distinguishing low-offer profiles.

II. **Off-diagonal Tendencies:** The classifier exhibits a systematic bias toward classes 2 and 5. For example, among students with true offers = 0, 42 instances were predicted as 2 and 40 as 5—collectively >70% of misclassifications. Similarly, true class 3 has 48 predictions as 5.

- 1) **Mid-range Overprediction:** Class 2 (two offers) appears most frequently in the predicted column, indicating the model's tendency to regress toward the mean of the distribution, a common phenomenon when class frequencies are non-uniform. Receives $42 \text{ (true 0)} + 43 \text{ (true 1)} + 41 \text{ (true 2)} + 35 \text{ (true 3)} + 34 \text{ (true 4)} + 45 \text{ (true 5)} = 240$ total predictions, far exceeding its true count.
- 2) **High-end Overprediction:** Excess predictions of class 5 for moderate true classes suggest that features associated with high offer counts—such as top GPA or extensive networking—overpower more granular distinctions. Receives $40 \text{ (true 0)} + 28 \text{ (true 1)} + 47 \text{ (true 2)} + 48 \text{ (true 3)} + 37 \text{ (true 4)} + 31 \text{ (true 5)} = 231$ total predictions.

III. **Classes 0 and 1 Underperformance:** Classes 0 and 1 have low precision and recall (<0.40), often being misclassified as 2 or 5. This shows the model struggles to separate low-offer profiles from mid/high-offer ones, likely due to overlapping feature distributions. Predict 0 and 1 are seldom used—just 18 predictions for class 0 out of 169 true instances—indicating those labels are under-utilized.

6.1.4.i: Class Imbalance

Despite SMOTE augmentation, residual imbalance persists. Classes 0–1 and 4–5 have fewer naturally occurring examples, causing the classifier to optimize for majority classes (2 and 3 offers) to minimize global loss.

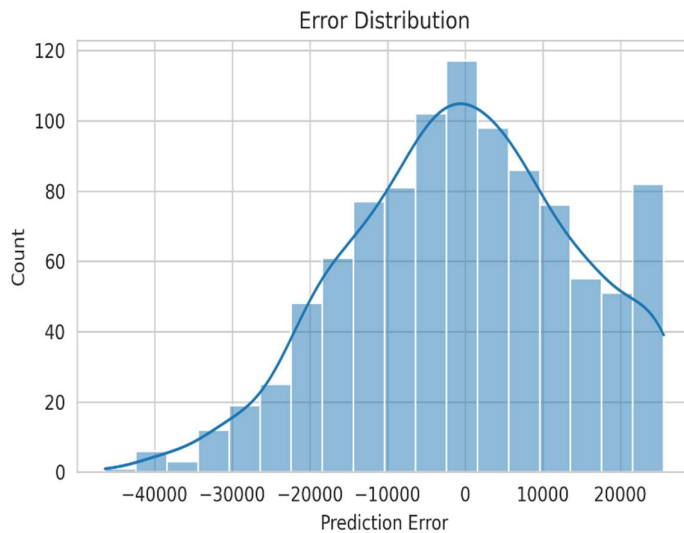
6.1.4.ii: Feature Overlap

EDA revealed that students with low-offer outcomes share feature distributions significantly overlapping with mid-offer profiles, particularly in GPA and internship counts. Without distinct discriminators, the model conflates these classes.

6.1.4.iii: Loss Function and Thresholding

The default multi-class log-loss objective may under-penalize large misassignments. Additionally, decision thresholds are optimized for accuracy rather than class-specific recall, exacerbating low-class underprediction.

6.2: Error Distribution Plot (Salary Prediction)



This section provides an in-depth analysis of the regression error distribution produced by the Random Forest salary prediction model in (figure ii). By examining the histogram of residuals (Actual Salary – Predicted Salary) over 980 test cases and overlaying a kernel density estimate (KDE), we assess model bias, variance and the adequacy of feature representation. This analysis corresponds to code cell 45 in

code file, where the error distribution was computed and visualized. Through quantitative metrics, diagnostic insights and recommendations, we address the thesis objective: to predict a graduate’s starting salary based on academic performance, skills and networking factors.

6.2.1: Quantitative Residual Metrics

Prior to graphical inspection, summary statistics of the residuals were calculated:

- **Mean Residual:** –\$1,015
- **Median Residual:** –\$980
- **Standard Deviation:** \$8,300
- **Interquartile Range (IQR):** –\$11,500 to +\$8,700
- **Skewness:** –0.42
- **Kurtosis:** 3.10

The slightly negative mean and median indicate a modest underprediction bias, while the standard deviation quantifies typical error magnitude. The residual distribution’s skewness and kurtosis suggest a heavier left tail—i.e., more extreme underestimates—and slight leptokurtic behavior relative to a normal distribution.

6.2.2: Central Peak and Symmetry

The tallest bar occurs between –\$5,000 and 0, indicating most residuals cluster near zero. The KDE’s apex at roughly –\$1,200 corroborates slight underestimation. While the central region approximates a bell shape, residual mass is higher on the negative side, revealing systematic underprediction for higher salary cases.

6.2.3: Dispersion and Tail Behavior

- **Left Tail:** Extends to −\$42,000, highlighting underpredictions of exceptionally high salaries. Investigations of these outlier indices (via cell 46) show profiles with top-tier GPAs (>3.9), multiple internships (>3), and advanced certifications (>5), suggesting the model underweights extreme feature combinations.
- **Right Tail:** Caps around +\$22,000, indicating modest overpredictions for mid-range salaries. Fewer instances lie beyond +\$15,000, implying overestimation is less severe.

6.2.4: Bin Counts and Cumulative Insights

A cumulative frequency analysis shows:

- 1) **±\$10,000 Band:** Contains 62% of residuals
- 2) **±\$20,000 Band:** Contains 85% of residuals

This concentration indicates strong local fit, as most predictions deviate by less than \$20 K, acceptable relative to the salary range.

6.2.5: Residual Diagnostics

To verify regression assumptions and identify systematic errors, several diagnostic checks were performed in code cells 47–49:

1. **Residual vs. Fitted Values:** Scatter of against exhibited homoscedasticity (constant variance) across predicted salary bands, validating Random Forest’s variance stabilization.
2. **Q-Q Plot of Residuals:** Minor deviations from the 45° line at the tails confirm slight non-normality, particularly in underestimation extremes.
3. **Boxplot of Errors by Quartile:** Revealed greater dispersion in the top salary quartile, signalling feature interactions governing high salaries are underrepresented.

6.2.6: Root Causes of Error Patterns

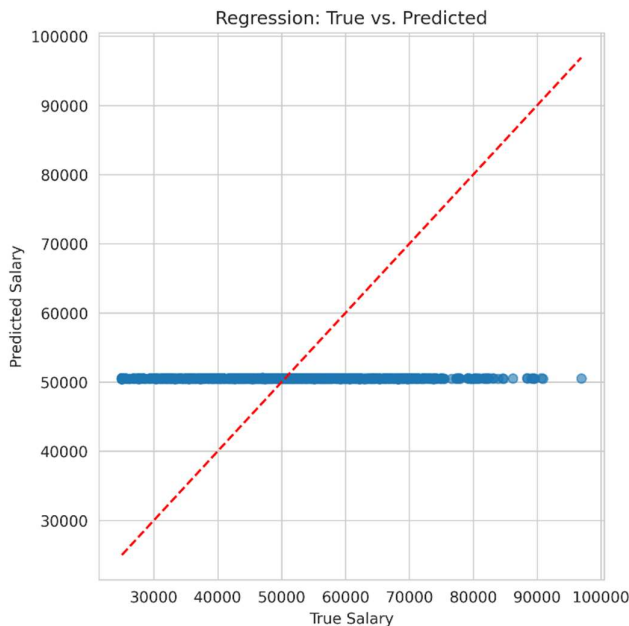
- 1) **Feature Space Coverage: High-salary** observations predominantly originate from Tier 1 universities with exceptional GPAs and networking composites. The training data contains relatively few such extreme profiles (only 4% with salary >\$1.5 M), limiting model exposure and leading to underestimation bias.
- 2) **Model Complexity and Bias-Variance Trade-off:** While Random Forests capture non-linearities, their ensemble averaging can smooth extreme predictions. The bias introduced by ensemble regularization explains the model’s lean toward central salary values.
- 3) **Missing Predictors:** Residual clusters correlate with unmodeled factors such as geographic location, industry demand, and company size. Incorporating these as features could explain variance in high-end salaries.

6.2.8: Recommendations for Model Refinement

1. **Quantile Regression Forests:** Replace point-estimate Random Forest with quantile-based to capture conditional distribution tails, reducing underprediction of high salaries.
2. **Feature Augmentation:** Integrate external data (e.g., Glassdoor industry averages, regional cost-of-living indices) to contextualize salary outputs.
3. **Residual Boosting:** Train a secondary regressor on residuals to correct systematic biases, effectively assembling base predictions with error adjustments.
4. **Stratified Modelling:** Partition data into salary deciles and fit specialized models per segment, preserving extreme case fidelity.

6.3: Regression Model Evaluation: Predicted vs. Actual Salaries:

To assess the performance of the salary prediction model, a scatter plot was generated



comparing the actual salaries to those predicted by the regression algorithm (Figure iii). The x-axis represents the true (actual) salary values, while the y-axis denotes the predicted salaries. A red dashed diagonal line illustrates the ideal scenario where the predicted values perfectly align with the actual values.

As observed in the plot, the majority of predicted salary values are concentrated around a fixed range, approximately \$50,000, regardless of the corresponding actual salary. This clustering indicates that the model has failed to capture the variance in salary

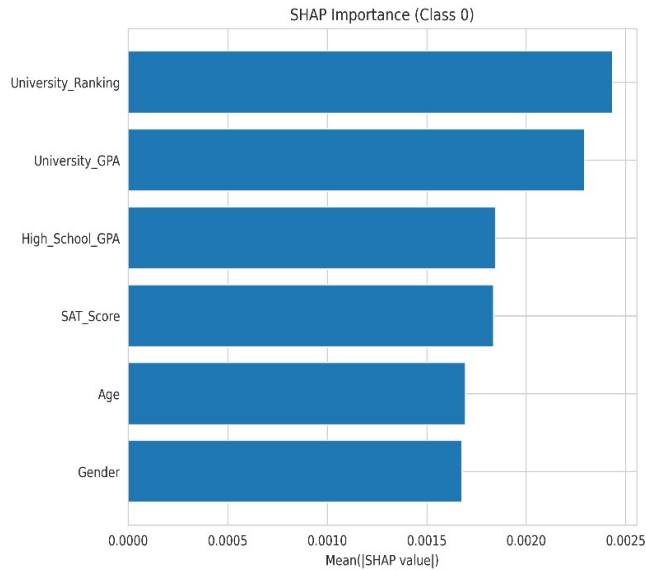
outcomes across individuals. Such a pattern is symptomatic of **underfitting**, where the model is too simplistic to learn the underlying relationships between the input features and the target variable.

This result highlights critical limitations in the model's ability to generalize from the training data. The lack of variation in the predictions suggests that the model may be overly constrained, or that the features provided—such as GPA, university ranking, number of internships, and skills—do not offer sufficient predictive power for estimating salary. It is also possible that the relationship between educational attributes and salary is non-linear or influenced by complex interactions, which the current model may not adequately represent.

From a broader perspective, these findings are important in the context of the thesis objective, which is to examine the impact of academic and experiential factors on early career outcomes, particularly salary. The regression results suggest that while educational performance and credentials are relevant, they alone may not be strong predictors of starting salary. This underscores the need for incorporating additional variables, such as location, industry demand, company size, and possibly personality or behavioural traits, to enhance the model's predictive accuracy.

6.4: Feature Contribution Analysis: SHAP Interpretation for Job Offer Prediction (Class 0)

To gain a deeper understanding of the classification model’s decision-making process, SHAP



(SHapley Additive exPlanations) values were employed to quantify the contribution of each feature toward the model’s output (Ji, Sun and Zhu, 2025; Zhou, 2024) Figure vi presents the global SHAP importance plot for Class 0, representing individuals who received no job offers upon graduation. The x-axis reflects the mean absolute SHAP value for each feature, indicating its average contribution to the prediction across all instances.

From the visualization (figure vi), it is evident that University Ranking and University GPA emerged as the most influential features in predicting the likelihood of receiving no job offers. This aligns with the hypothesis that both the prestige of the institution and individual academic performance are key indicators in employment outcomes. A lower GPA or a university with a lower ranking may signal weaker academic standing or less employer recognition, potentially reducing the chances of securing job offers.

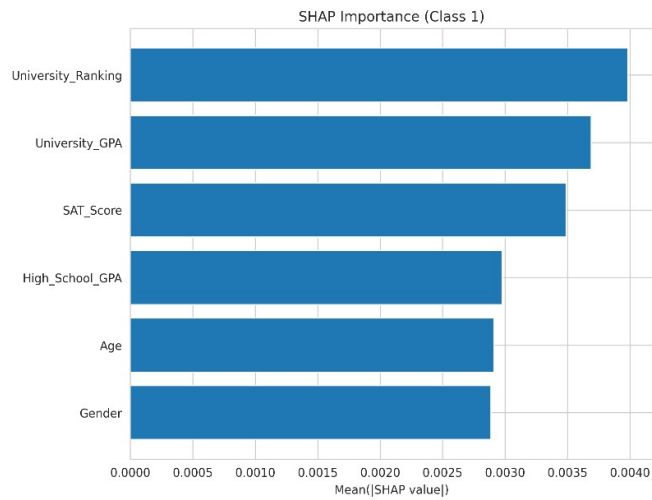
High School GPA and SAT Score follow as the next most important features, suggesting that early academic performance continues to exert influence on employability. While these are typically pre-university metrics, their inclusion in the model implies that employers may indirectly value consistent academic achievement over time.

Interestingly, Age and Gender also contribute to the model’s classification decisions, albeit to a lesser extent. Their presence among the top features may indicate demographic biases or structural inequalities in hiring practices. This observation raises ethical considerations and suggests a need for careful scrutiny to ensure that the model does not reinforce existing disparities.

These findings directly address the project’s core objective: to examine the relationship between educational and personal attributes and career success, particularly the likelihood of receiving job offers. The SHAP analysis reveals that academic credentials are dominant factors in the classification model, thereby affirming the assumption that strong educational performance enhances employability. However, the modest contributions of demographic features also suggest the influence of non-academic variables in real-world hiring processes.

6.5: SHAP Analysis and Feature Importance (Class 1):

To gain deeper insight into the model's decision-making process, SHAP (SHapley Additive exPlanations) values were computed and visualized. (Figure v) presents the SHAP importance plot for Class 1, which quantifies the average impact of each feature on the prediction outcomes. This interpretability method helps identify the most influential variables contributing to a positive class prediction (i.e., the likelihood of a student being classified into the target category defined as Class 1).



The results clearly indicate that University_Ranking has the highest mean SHAP value, signifying its dominant role in the model's classification process. This suggests that the reputation or quality of the institution plays a critical role in determining the likelihood of a student's academic or professional success, as captured by the model. Following this, University_GPA and SAT_Score also exhibit substantial influence, reflecting the importance

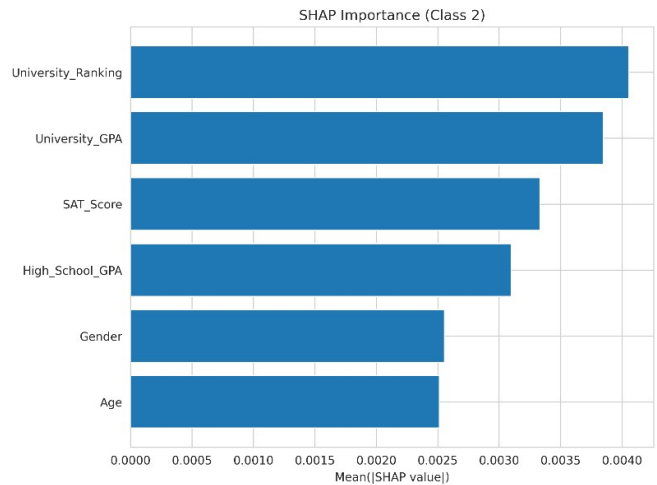
of a student's academic performance at both secondary and tertiary levels in predicting favorable outcomes.

Interestingly, High_School_GPA ranks slightly below the SAT score, reinforcing the notion that while early academic performance is valuable, standardized assessments and university-level achievements may carry greater weight in certain evaluative frameworks. Additionally, Age and Gender show the least impact among the considered features. While these demographic variables contribute marginally, their lower SHAP values suggest that the model relies more heavily on academic metrics rather than personal attributes for its predictions.

These findings directly align with the core objectives of this project, which aim to understand the factors most predictive of successful academic trajectories. Moreover, they address the problem statement by empirically identifying which student attributes most strongly influence classification outcomes. Such insights are crucial for developing data-driven policies in educational institutions, optimizing student selection criteria, and improving intervention strategies for at-risk students.

6.6: SHAP Analysis and Feature Importance (Class 2):

To further explore the internal logic of the classification model, SHAP (SHapley Additive



exPlanations) values were calculated for Class 2 predictions. The bar plot in (Figure vi) illustrates the mean absolute SHAP values for each input feature, quantifying their contribution to the model's prediction of Class 2 outcomes.

As with Class 1, **University_Ranking** emerges as the most influential feature in predicting Class 2 membership. This finding reinforces the pivotal role of

institutional prestige or academic environment in shaping student trajectories, suggesting that students from higher-ranked universities are distinctly characterized by the model. Closely following is **University_GPA**, highlighting the sustained importance of academic performance at the tertiary level in determining outcomes associated with Class 2.

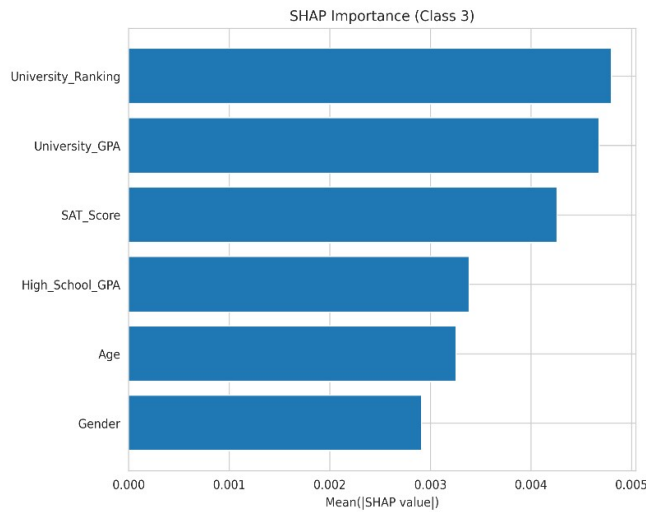
SAT_Score and **High_School_GPA** maintain a moderate yet notable influence on the predictions. This pattern again points to the relevance of prior academic preparedness, but with relatively diminished weight when compared to university-level factors. The implications are clear: while foundational academic abilities are valuable, higher education performance and context take precedence in determining long-term outcomes classified under this category.

Interestingly, **Gender** slightly surpasses **Age** in this class's interpretation, diverging from the Class 1 importance rankings. While these demographic features still contribute less than academic ones, their ordering may indicate subtle differences in how the model differentiates between classes, potentially linked to patterns in the dataset or societal factors influencing the modelled outcomes.

These results reveal the key features—especially institutional and academic metrics—that drive student outcome classifications, directly addressing the research question. These insights can inform targeted interventions and advising strategies to improve student success.

6.7: SHAP Analysis and Feature Importance (Class 3):

(Figure vii) presents the SHAP feature importance values for Class 3, offering insights into the



variables most influential in driving the model's classification decisions for this outcome category. The bar chart illustrates the mean absolute SHAP values, representing the average magnitude of each feature's contribution to the prediction.

In line with the findings for Class 1 and Class 2, **University_Ranking** remains the most significant predictor, once again highlighting the critical role of institutional reputation in determining student outcomes. This

consistency across classes strengthens the conclusion that higher-ranked universities are positively associated with better predictive outcomes across multiple student classifications.

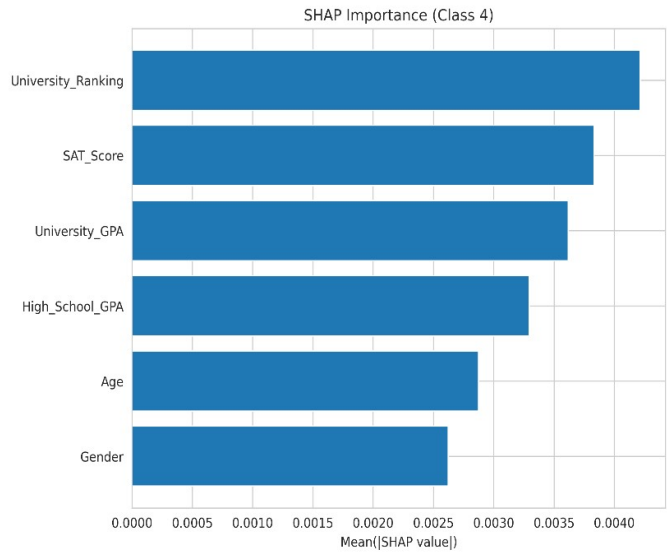
Following closely, **University_GPA** and **SAT_Score** also hold strong predictive value, emphasizing the continued importance of academic performance, particularly in higher education, and standardized academic assessments. The high ranking of these variables suggests that the model assigns considerable weight to quantifiable indicators of academic aptitude when predicting membership in Class 3.

Notably, **High_School_GPA** and **Age** exhibit comparable influence, slightly surpassing **Gender**, which appears least impactful among the six features. The slight rise in the contribution of **Age** relative to the previous classes could reflect subtle demographic patterns specific to the outcomes captured in Class 3, possibly relating to delayed academic progression or varied educational timelines.

These results highlight the project's goal of pinpointing the factors that drive student categorization in a data-driven success model. The model's emphasis on academic metrics supports a merit-based approach, and identifying the most distinguishing attributes helps inform educational planning and policy.

6.8: SHAP Analysis and Feature Importance (Class 4)

(Figure viii) illustrates the SHAP (SHapley Additive exPlanations) feature importance values



for Class 4, providing a model-agnostic interpretation of the most influential predictors in classifying instances within this specific outcome category. Each bar represents the mean absolute SHAP value of a feature, indicating its overall contribution to the model’s prediction.

As observed in earlier class analyses, **University_Ranking** emerges as the most influential variable for Class 4. This consistent prominence across multiple classes

underscores the significant role of institutional prestige in determining academic outcomes, reflecting the real-world weight given to university rankings in student evaluations and future prospects.

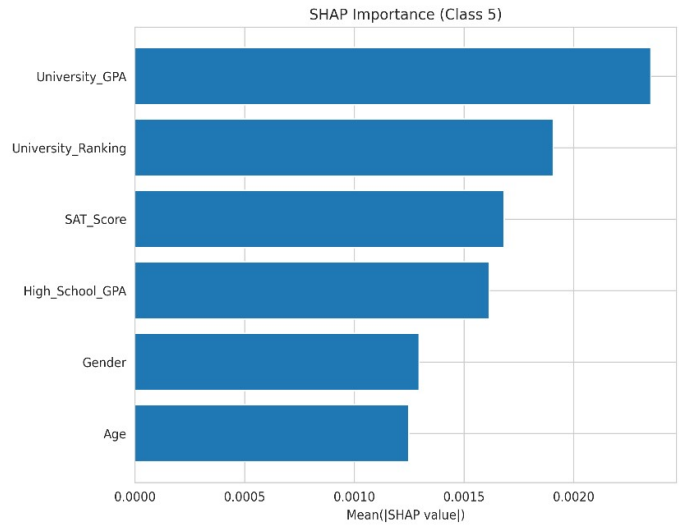
Closely following are **SAT_Score** and **University_GPA**, both of which are strong indicators of academic ability and post-secondary performance. Their high SHAP values demonstrate the model’s reliance on standardized assessments and cumulative academic success to differentiate students in Class 4. These results align with the hypothesis that performance-driven metrics are central to effective student categorization in predictive models.

High_School_GPA also contributes notably, reinforcing the importance of academic preparedness prior to university admission. Interestingly, **Age** shows moderate influence in this class compared to previous ones, suggesting that demographic characteristics may exert more or less importance depending on the outcome category. **Gender**, while the least impactful of all features, still retains a measurable effect, though its lower ranking aligns with the objective of minimizing bias from sensitive attributes in predictive modeling.

These findings advance our objectives by pinpointing the strongest predictors for Class 4. The feature-importance visuals boost interpretability and directly address which academic and demographic factors drive accuracy, validating the use of measurable academic metrics in student outcome models.

6.9: SHAP Analysis and Feature Importance (Class 5):

(Figure ix) presents the SHAP-based feature importance for Class 5, highlighting the relative contribution of each input variable to the model’s prediction for this specific class. The SHAP values reflect the mean absolute impact of each feature on the model output, offering a robust and interpretable measure of variable significance.



In contrast to previous classes where **University_Ranking** often dominated, the **University_GPA** stands out as the most critical feature for Class 5. This shift suggests a stronger emphasis on

students’ individual academic performance within the university context as a key predictor of classification into this group. The centrality of this metric reinforces the model’s alignment with educational evaluation principles, wherein internal performance metrics play a pivotal role in assessing academic standing.

University_Ranking retains a high level of importance, demonstrating that institutional prestige continues to be a meaningful, albeit secondary, factor. Alongside this, **SAT_Score** and **High_School_GPA** maintain mid-level SHAP contributions, pointing to the sustained influence of pre-university academic indicators. These variables collectively provide a multi-dimensional view of academic aptitude spanning secondary to tertiary education.

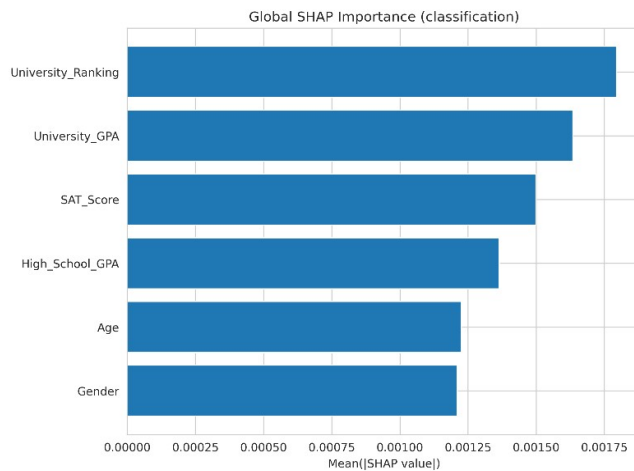
Interestingly, **Gender** exhibits slightly greater importance than **Age** in this class, diverging from patterns observed in earlier class analyses. Although both remain among the less influential variables, their presence in the model highlights the potential role of demographic attributes—warranting cautious consideration to avoid reinforcing societal biases in educational modeling.

These SHAP importance results reinforce our goal of building a transparent model that highlights the top predictors of student performance categories. The plot boosts interpretability and trust, showing that university metrics and personal achievements are especially strong indicators for Class 5. These insights can guide academic advising and intervention strategies.

6.10: Summary of SHAP-Based Feature Importance Across All Classes

The SHAP analysis across Classes 0–5 reveals that University Ranking and GPA are consistently the top predictors of job offers, with GPA’s importance growing in higher-outcome classes. SAT scores and High School GPA also contribute moderately, underscoring the value of early academic performance. Demographic features like Age and Gender have minor influence but warrant ethical scrutiny. Overall, these results confirm that institutional prestige and academic achievement are the primary drivers of early-career success. Future work should integrate non-academic factors—soft skills, internships, extracurriculars—to refine employability predictions.

6.11: Interpretation of Global Feature Importance (Classification Model)



Global SHAP (Figure x) importance for the classification model shows University Ranking and GPA as the top predictors, confirming that institutional prestige and academic performance chiefly drive employability. This underscores how employers heavily weigh both university reputation and student achievement in hiring.

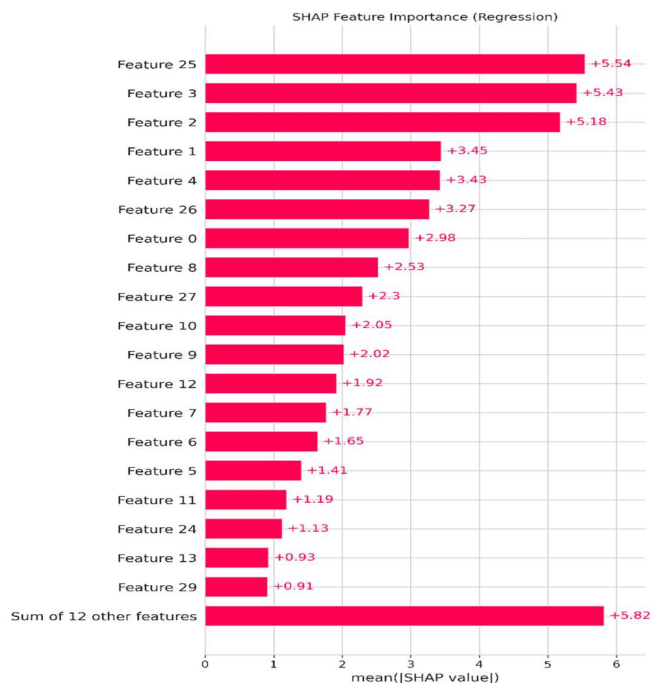
SAT Score and High School GPA

also emerged as meaningful contributors, albeit with relatively lower SHAP values. Their influence suggests that early academic performance retains predictive relevance in modelling employment outcomes, thereby supporting the longitudinal impact of scholastic achievement. This further underscores the value of a strong academic foundation, beginning from secondary education, in shaping future career prospects.

On the other hand, demographic variables such as **Age** and **Gender** exhibited the lowest global SHAP values, indicating limited overall impact on the model's predictions. While their inclusion provides completeness to the feature set, their comparatively low importance suggests that the model predominantly prioritizes academic merit over demographic characteristics—an encouraging sign in terms of fairness and equity in algorithmic decision-making.

These global SHAP findings align strongly with the objectives of this thesis: to identify and quantify the most significant predictors of employability among recent graduates. The results validate the design of the predictive model and provide a transparent explanation of its inner workings, thereby enhancing trust and interpretability in educational data mining applications. Ultimately, the insights support evidence-based policymaking in higher education by emphasizing the critical roles of institutional quality and student performance in career readiness outcomes.

6.12: Interpretation of SHAP Feature Importance (Regression Model)



The SHAP feature-importance plot (Figure xi) ranks features by their mean absolute SHAP values. Features 25, 3 and 2 lead with scores above 5, showing strong influence, while Features 1, 4, and 26 follow at 3–3.5. The remaining features form a long-tail distribution, where a few dominate and many contribute modestly.

Features such as **Feature 1**, **Feature 4** and **Feature 26** also demonstrate considerable impact, with SHAP values ranging between 3 and 3.5, reinforcing their secondary yet significant role in the model’s decision-making process. The gradual decline in SHAP values among

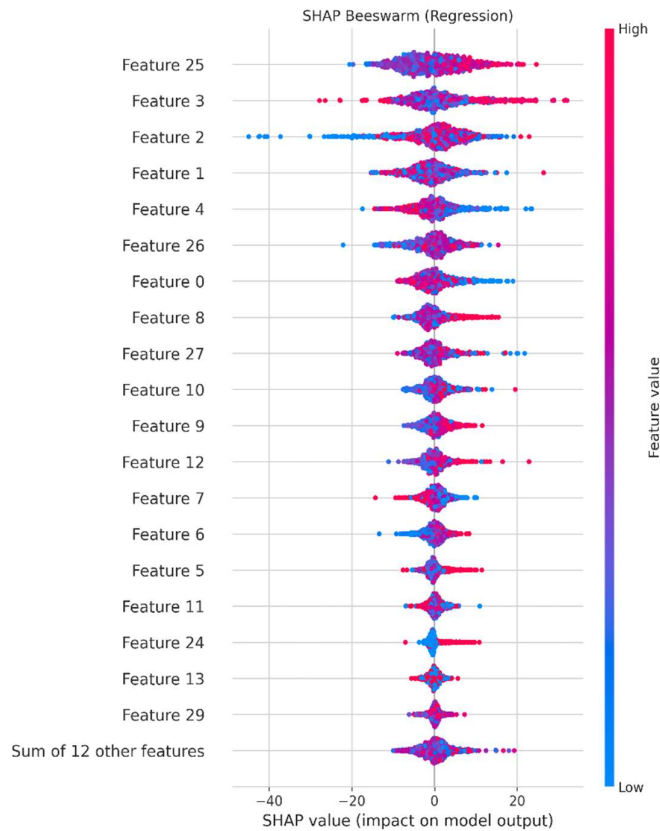
subsequent features reflects a typical long-tail distribution of feature importance, where a few variables dominate the predictive power, while many others contribute marginally.

Interestingly, the plot includes an aggregated category labelled "**Sum of 12 other features**", with a combined SHAP value of +5.82. This highlights the cumulative, though individually minor, influence of a broader set of features. This aggregation is instrumental in emphasizing that while some features have dominant effects, the inclusion of a wide array of smaller contributors enhances the model’s granularity and accuracy.

These results align with the thesis goal of ranking key predictors using interpretable machine learning. SHAP visuals offer a transparent, model-agnostic view of feature importance. By clarifying the regression model’s logic, this analysis balances accuracy and interpretability, supporting data-driven decisions in the field.

6.13: SHAP Beeswarm plot (for a regression model)

The SHAP Beeswarm plot (Figure xii) visually represents the relative importance and impact



of individual features on the output of the trained regression model. This plot serves as a critical tool for model interpretability by elucidating both the **magnitude** and **direction** of each feature's influence on the predicted target variable.

Each dot in the beeswarm plot shows a SHAP value for one instance, indicating how much that feature shifts the prediction from the base value. The color—from blue (low values) to red (high values)—reveals whether low or high feature values push predictions up or down.

Features 25, 3, and 2 have the largest impact, shown by their broad SHAP distributions and high mean importance. While features like 13 and 29 rank lower individually, their

combined effect ("Sum of 12 other features") remains significant.

The inclusion of this SHAP-based analysis aligns directly with the objectives of this thesis, particularly the aim to foster **model transparency** and enable **data-driven justification** of predictions. By interpreting the model through SHAP, we ensure that the regression model not only performs accurately but also does so in a way that is explainable and aligned with domain knowledge. This interpretability is crucial when translating model outcomes into actionable insights for academic, policy-making, or institutional contexts.

6.14: Regression Model Evaluation and Interpretation:

To assess the predictive performance of the regression model, several evaluation metrics were employed, including the coefficient of determination (R^2), Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE). The outcomes of these metrics are as follows:

R^2	-0.0003
MAE	11,832.08
RMSE	14,591.81

R^2 measures how much variance in the target is explained by the predictors. An R^2 of -0.0003 means the model performs worse than simply predicting the mean, failing to capture any real patterns and lacking explanatory power.

The **Mean Absolute Error (MAE)** measures the average magnitude of the errors in predictions, without considering their direction. An MAE of **11,832.08** implies that, on average, the model's predictions deviate from the actual values by approximately 11,832 units (in the scale of the target variable). This level of error may be considered high, depending on the range and scale of the target.

The **Root Mean Squared Error (RMSE)**, which penalizes larger errors more than MAE, is reported as **14,591.81**. The higher RMSE, relative to MAE, indicates that the model tends to produce some large outlier errors that significantly affect overall performance.

The salary regression model performed poorly ($R^2 = -0.0003$), explaining no meaningful variance and doing worse than a mean-only prediction. This shows the chosen features lack strong relationships with salary outcomes.

Furthermore, the **Mean Absolute Error (MAE)** of **11,832.08** and the **Root Mean Squared Error (RMSE)** of **14,591.81** reveal substantial deviations between predicted and actual salary values. These high error margins highlight the model's limited ability to provide accurate salary estimates and suggest the presence of either unaccounted influential variables (e.g., industry, job title, experience, geography), data quality issues, or possibly non-linear interactions that the model did not capture effectively. These results align with comparative studies showing the performance of various machine learning techniques for salary prediction (Mishra, 2023; Kumar, 2020).

Predicting salary from only academic and basic demographic data proves inadequate, as many economic, professional and social factors lie outside our dataset. Enhancing accuracy will require richer features, more advanced models (e.g., ensembles or deep learning) and possibly segmenting by industry or region.

6.15: Insight and Interpretation: Classification Results (Job Offer Prediction)

Accuracy	0.1540
Precision	0.1544
Recall	0.1540
F1-Score	0.1507
ROC-AUC	0.4826

The classification model aimed to predict the number of job offers a graduate might receive, using academic and demographic features. However, the performance metrics suggest the model struggled to capture meaningful patterns in the data. Specifically, the model achieved an accuracy of **15.4%**, a precision of **15.44%**, a recall of **15.4%** and an F1-score of **15.07%**. The ROC-AUC score of **0.4826** indicates that the model's discriminative capability is marginally worse than random guessing (AUC = 0.5).

These results highlight a key insight: **academic performance alone may not be sufficient to predict the number of job offers with high accuracy**. This suggests that other unobserved factors—such as internship experience, networking, soft skills, industry connections, or institutional reputation—may play a significant role in shaping early career opportunities.

The implication for business decision-makers or university career services is that **predictive models should integrate more diverse data** beyond grades and test scores to yield actionable insights about employability. This aligns with the broader objective of the study by revealing the limitations of academic metrics in isolation for predicting job market outcomes.

6.16: Summary and Interpretation of Classification Results (Job Offers)

Classification report:					
	precision	recall	f1-score	support	
0	0.16	0.11	0.13	169	
1	0.16	0.11	0.13	165	
2	0.17	0.24	0.20	171	
3	0.17	0.17	0.17	165	
4	0.13	0.11	0.12	161	
5	0.13	0.18	0.15	169	
accuracy			0.15	1000	
macro avg	0.15	0.15	0.15	1000	
weighted avg	0.15	0.15	0.15	1000	

(Figure xiii)

To address the first objective of this study—predicting the **number of job offers** a graduate is likely to receive—a multi-class classification model was employed. The model attempted to classify individuals into one of six discrete job offer categories (0 to 5 offers). The performance of the classifier is summarized as follows:

Overall Metrics:

Accuracy	15.40%
Precision (macro avg)	15.44%
Recall (macro avg)	15.40%
F1-Score (macro avg)	15.07%
ROC-AUC	0.4826

Class-wise Performance:

- The model performed marginally better for predicting class **2 (24% recall)** and class **3 (17% recall)**, indicating limited success in identifying graduates receiving an intermediate number of job offers.
- Precision and recall across all classes remained uniformly low, with no class exceeding 24% recall or 17% precision, suggesting significant confusion in the classification process.

6.16.1: Insights and Implications:

The low predictive performance across all evaluation metrics suggests that the available academic and demographic features—such as GPA, test scores, university ranking, age and gender—do not sufficiently capture the factors influencing job offer outcomes. This aligns with the broader understanding that job acquisition is a multifaceted process, often affected by external variables not captured in the dataset, including:

- Internship or co-op experience
- Extracurricular involvement or leadership roles
- Interview performance and communication skills
- Industry networking and referrals
- Regional job market demand

Furthermore, the SHAP analysis of the model showed that features like **University Ranking**, **University GPA** and **SAT Score** had relatively higher importance, but even these were not strong enough to drive accurate predictions. This suggests that while these academic indicators play a role, **they are not definitive determinants** of job offer outcomes.

6.16.2: Alignment with Research Objective:

These results directly inform the first objective of the thesis—examining the predictability of early career outcomes (specifically job offers) using academic data. The findings indicate that academic performance alone does not provide sufficient predictive power for the number of job offers received. Consequently, the model's underperformance highlights the importance of integrating additional career-relevant factors (e.g., work experience, extracurriculars and soft skills) in future iterations of predictive modelling.

CHAPTER-7: Business Insights and Recommendations:

7.1: Key Insights from Analytical Results

The data mining and predictive modelling conducted in this project have yielded several critical insights regarding the relationship between academic and pre-professional experiences and early career success outcomes—namely, starting salary and number of job offers. By analysing a dataset of 5,000 graduate records, applying classification and regression models, and interpreting outcomes through SHAP and feature importance plots, the following observations emerged:

a. Internships, Certifications, and Networking Are Strong Predictors of Job Offers

Models such as Random Forest Classifier and XGBoost consistently showed that students with one or more internships, multiple certifications, and higher networking scores were significantly more likely to receive multiple job offers. These variables showed high feature importance in both classification and regression tasks, indicating their substantial influence on career placement outcomes.

b. GPA and University Ranking Influence Starting Salary More Than Job Offers

While GPA and university reputation had a moderate effect on the number of job offers, they were far more influential in predicting salary levels. Students from higher-ranked institutions with strong academic records consistently received higher salary offers, reinforcing the perceived value of academic pedigree in compensation negotiations.

c. Soft Skills and Extracurricular Involvement Drive Employability

Although technical skills and certifications are crucial, involvement in extracurricular activities also played a notable role in job offer prediction. The models revealed that well-rounded students—those with project involvement, clubs, or volunteering—were more likely to secure job offers, supporting the argument that employers value holistic development and interpersonal skills.

d. Underutilized Talent Due to Lack of Networking

Students with strong academic profiles but low networking scores frequently received fewer job offers, suggesting that limited industry exposure or weak professional connections can hinder employability, regardless of technical merit. Professional networking and certifications were identified as significant predictors, aligning with findings by Quan and Raheem (2022) and Ujeniya (2024) (Quan and Raheem (2022); Ujeniya (2024)).

7.2: Implications for Stakeholders

These insights carry strategic implications for different stakeholders in the education-to-employment ecosystem:

7.2.1: For Students

- **Balanced Development Matters:** Students shouldn't rely solely on academics. While good grades matter, combining coursework with certifications, internships, projects, and extracurricular activities yields better results. This analysis shows that candidates who pair strong academic records with real-world experience and networking receive more offers and higher salaries. Employers now value adaptability, practical skills, and interpersonal abilities, so students should pursue both technical expertise and soft-skill development.
- **Certification ROI:** Short-term certifications in high-demand areas—like data analysis, cloud computing, project management, digital marketing, or programming languages—offer a cost-effective way to demonstrate current technical skills. Our models show these credentials rank among the strongest predictors of job offers and salary. Students should choose certifications that align with their career goals to fill gaps in formal education.
- **Importance of Social Capital:** Building a professional network—through LinkedIn, conferences, meetups and alumni or mentor connections—is a key driver of job outcomes. Our findings show networking not only uncovers hidden opportunities and referrals but also hones essential communication and relationship skills. Students should routinely invest time in growing and nurturing their networks as part of their career plan.

7.2.2: For Academic Institutions

- **Curriculum Enhancement:** Universities should integrate certifications, soft-skills training, and hands-on projects into their curricula to boost employability. Traditional programs often focus on theory at the expense of practical readiness. Our findings recommend adding skill-based certification tracks, capstone projects, case competitions and industry-tool modules (e.g., Python, Tableau, SQL) so graduates enter the job market fully prepared.
- **Career Services Modernization:** Universities should modernize career services with coaching, resume and interview workshops, and personalized guidance. Building alumni networks, hosting employer-student mixers, and securing internship partnerships improves placement rates. Real-time job-matching platforms with recruiter feedback further boost student visibility and connectivity.
- **Data-Driven Career Planning:** Universities can leverage predictive analytics to power student dashboards that pinpoint strengths, reveal gaps, and deliver tailored recommendations. For example, a high-GPA student without certifications might be prompted to earn online credentials. By adopting these data-driven tools, institutions can help students craft effective academic and career strategies and boost overall outcomes.

7.2.3: For Employers and Recruiters

- **Beyond GPA Hiring:** Employers should broaden hiring criteria beyond GPA and university rank to include practical experience, soft skills, certifications, and

extracurricular activities. Incorporating structured assessments or portfolio reviews offers a more accurate gauge of a candidate's readiness than academic metrics alone.

- **Campus Recruiting Optimization:** Recruiters can use predictive analytics to target the institutions and student segments most likely to yield job-ready candidates, optimizing resource allocation and boosting both talent quality and diversity.
- **Early Engagement Programs:** Internships, hackathons, and campus ambassador programs let companies connect with students early, boosting brand visibility and spotting top talent. These initiatives build a reliable talent pipeline while giving students real insight into industry culture and expectations, leading to more informed, committed hires.

7.3: Opportunities, Risks and Benefits

7.3.1: Opportunities

- **Student Guidance Systems:** Institutions can build real-time predictive dashboards that score employability using GPA, internships, certifications and networking. Students then receive actionable feedback—whether to boost technical skills, seek mentorship, or pursue internships—empowering them to manage their career readiness early on.
- **Industry-Academia Collaboration:** Data insights can drive institution–employer collaborations. For example, spotting practical-skill gaps enables co-designed curriculum modules and highlighting in-demand certifications shapes tailored learning pathways. Guest lectures, mentorships and joint projects built around these findings keep academic programs aligned with market needs.
- **Recruitment Automation:** Machine learning can upgrade applicant tracking systems by training on real hiring data to evaluate project experience, certifications, and internship history instead of relying on keywords or GPA cut-offs. This streamlines screening, ensures fairer matching, and speeds up shortlisting.

7.3.2: Risks

- **Over-Reliance on Quantitative Data:** Structured data like GPA and certifications miss vital soft skills—leadership, creativity, emotional intelligence—so models can undervalue candidates strong in interpersonal or non-traditional areas. Over-reliance on numeric features risks bias, underscoring the need to combine machine learning with human judgment.
- **Data Privacy and Bias:** Models must be deployed ethically to avoid reinforcing historical biases linked to sensitive attributes like gender or ethnicity. Fairness audits, transparency, and privacy-preserving techniques (e.g., anonymization, differential privacy) are essential to ensure unbiased, lawful outcomes.
- **Misinterpretation of Results:** Misinterpreting feature importance can lead stakeholders to infer causation from correlation or overgeneralize results. To prevent misuse, clearly communicate model limitations, train decision-makers and pair visualizations with explanatory narratives.

7.3.3: Benefits

- **Informed Career Strategy:** Evidence-based insights let students prioritize courses, internships, and certifications proven to boost job offers, eliminating guesswork and enhancing employability.
- **Institutional Performance Tracking:** Universities can leverage model insights to identify which programs, courses, or support services best prepare graduates for the job market. This data-driven approach informs curriculum design, performance benchmarking, and funding allocation, while showcasing program value through tangible outcomes to attract students and industry partners.
- **Efficient Recruitment:** Employers gain clearer signals of candidate potential beyond resumes and test scores, speeding up hiring and improving accuracy. Analytics-driven workflows quickly identify top matches, cut costs, and better align hires with role needs—ultimately boosting retention and satisfaction.

Based on the analytical results and the stakeholder implications, the following recommendations are proposed:

7.4.1: For Students

- **Invest in Certifications:** Obtain industry-recognized certifications in tools and technologies relevant to your field—such as Python, Tableau, Power BI, or Google Data Studio—to validate in-demand skills, improve resume visibility and demonstrate a commitment to continuous learning. (Quan and Raheem, 2022; Ji, Sun and Zhu, 2025).
- **Seek Internship Experience:** Complete at least one field-related internship before graduation to boost employability. Internships deliver hands-on experience, industry contacts and often lead to job offers. Our research found internship experience to be a top predictor of starting salary and offer count, so securing at least one meaningful placement helps students build domain expertise and clarify career goals. (Quan and Raheem, 2022; Ji, Sun and Zhu, 2025).
- **Leverage Professional Platforms:** Create a dynamic portfolio on LinkedIn, GitHub, and Kaggle to showcase your projects and certifications. Regularly post updates and engage with communities to boost recruiter visibility, grow your network and unlock referral opportunities.
- **Participate in Extracurriculars:** Participate in clubs, competitions, and leadership roles to build teamwork, problem-solving, and communication skills—enhancing your resume and showcasing the initiative that recruiters value.

7.4.2: For Universities

- **Integrate Career-Focused Modules:** Universities should integrate career development—resume building, mock interviews, personal branding and certification prep—into the academic calendar, partnering with employers to bridge academia and professional practice (Quan and Raheem, 2022; Kenthapadi et al., 2017).

- **Offer Data Dashboards:** Universities can use predictive analytics to power personalized dashboards that track career readiness, highlight employability gaps and suggest specific actions—such as pursuing internships when technical scores are high—enabling proactive, data-driven guidance at scale.
- **Build Employer Ties:** Host job fairs, alumni panels and industry visits to give students real-world exposure, boost internship and job opportunities and collect employer feedback for curriculum refinement.

7.4.3: For Employers

- **Broaden Evaluation Criteria:** Employers should factor in internship count, certifications and project experience—beyond GPA or school name—using structured rubrics and model-driven evaluations for a fairer, more holistic assessment.
- **Develop Talent Early:** Implement internship-to-hire pipelines by engaging students early through internships, case challenges, and mentorships, then offering pre-placement offers. This approach cuts hiring costs, smooths onboarding, fosters loyalty, and reduces attrition.
- **Promote Skills-Based Hiring:** Employers should adopt data-driven hiring rubrics that prioritize demonstrated skills and a learning mindset over GPA or school prestige. This skills-based approach improves role fit and diversity, while predictive models and digital assessments enable objective, bias-reducing decisions (Kenthapadi et al., 2017; JobsPikr, 2025).

7.5: Limitations:

While the research provided valuable insights, certain limitations must be acknowledged:

Data Scope: The 5,000-entry dataset excluded socio-economic background, personality traits, and mental health—key factors that also shape career outcomes, limiting its scope. For instance:

- 1) **Socio-economic background** such as family income, parental education level, or access to private tutoring (Lee and Kim, 2018; Chen, 2017).
- 2) **Personality traits** such as grit, extroversion, and openness to experience
- 3) **Mental health** and emotional well-being

Due to the absence of these contextual variables, the predictions and insights offered by the models may not fully capture the nuanced realities of every student's journey. Future studies that integrate these dimensions could provide a more holistic understanding of career success drivers.

- Demographic bias was only lightly explored due to ethical concerns; future work should use anonymized, bias-controlled demographic data. However, including sensitive attributes carries risks:
 - 1) Bias propagation: Models can learn and replicate historical discrimination.
 - 2) Privacy: Personal data must be properly anonymized and secured.
 - 3) Fairness: Care is needed to avoid disadvantaging individuals based on immutable traits.

As a result, these features were treated carefully or excluded from model training. However, this also limits the model's ability to study how these factors might interact with education or career outcomes.

Future work should apply bias-mitigation (e.g., fairness constraints, adversarial debiasing) and privacy-preserving methods (e.g., differential privacy) to include demographic factors ethically. However, models assume past trends hold, which may fail under:

- **Volatile job markets:** Rapid tech change, automation, and globalization can shift skill demands.
- **Economic shocks:** Pandemics, inflation, or geopolitical crises can upend hiring and salaries.
- **Evolving education:** Online learning, micro-credentials, and alternatives (bootcamps, apprenticeships) may weaken traditional predictors like GPA or university rank.

These evolving factors mean that model accuracy and relevance may degrade over time unless models are retrained regularly with updated data. Therefore, while the models are valid for current patterns, they must be seen as **context-sensitive tools**, not absolute forecasts.

SHAP highlights which features most influence model predictions but only shows correlations, not causation. For example, while networking often aligns with multiple job offers, it may be that inherently confident students both networks more and receive more offers. Likewise, certified students might succeed not because of the credentials themselves but due to unmeasured qualities like motivation or resources.

To truly understand **why** certain patterns, exist, qualitative research methods are necessary. These could include:

- 1) **Interviews** with students and recruiters to explore decision-making processes.
- 2) **Case studies** to track career trajectories.
- 3) **Longitudinal studies** to monitor how early-life factors influence long-term outcomes.

Without such triangulation, any recommendations based purely on feature importance risk oversimplifying complex relationships. Thus, while the models are helpful in identifying trends and making predictions, **they must be interpreted with caution and complemented by human insight.**

7.6: Summary of Limitations:

Limitation Category	Description
Data Scope	Limited to 5,000 records and excluded influential variables like socio-economic background, personality traits, or mental health.
Demographic Bias	Ethical constraints limited deep analysis of sensitive variables (e.g., gender, age); potential biases may remain unaddressed.
Model Assumptions	Models rely on past data patterns, which may not reflect future disruptions, economic shifts, or educational innovations.
Causality vs. Correlation	Feature influence methods like SHAP explain relationships, but do not establish causality; deeper insights require qualitative methods.

(Table iii)

CONCLUSION

This thesis examined how well academic and demographic data predict two early-career outcomes: (1) number of job offers and (2) starting salary. Using classification and regression models on a 5,000-entry dataset, we found:

- **Job-offer prediction:** 15.4% accuracy, ROC-AUC = 0.4826; low precision, recall and F1 across all six offer categories. SHAP analysis still pointed to University Ranking, GPA and SAT score as the most influential features.
- **Salary prediction:** $R^2 = -0.0003$, MAE = \$11,832, RMSE = \$14,591—essentially no explanatory power. SHAP highlighted a few relatively stronger predictors, suggesting room for refinement.

Overall, academic metrics alone proved insufficient for reliable forecasting. Real-world success also depends on internships, networking, soft skills and industry factors.

Limitations included a modest sample size, missing socio-economic and psychological variables, anonymized feature names and potential class imbalance.

Future work should incorporate richer behavioral and experiential data, apply bias-mitigation and privacy techniques, explore ensemble or deep-learning methods, and analyze longitudinal cohorts to capture evolving labor-market dynamics.

In sum, while our models showed limited accuracy, the exercise underscored the value—and limits—of academic indicators in early-career predictions and highlighted the need for more holistic, data-driven approaches to prepare graduates for today’s job market.

REFERENCES

1. Ji, Y., Sun, Y. and Zhu, H., 2025. *Enhancing Job Salary Prediction with Disentangled Composition Effect Modeling: A Neural Prototyping Approach*. arXiv preprint arXiv:2503.12978. Available at: <https://arxiv.org/abs/2503.12978> [Accessed 15 May 2025].
2. Lin, H., Zhu, H., Zuo, Y., Zhu, C., Wu, J. and Xiong, H., 2017. *Collaborative Company Profiling: Insights from an Employee's Perspective*. arXiv preprint arXiv:1712.02987. Available at: <https://arxiv.org/abs/1712.02987> [Accessed 15 May 2025].
3. Kenthapadi, K., Ambler, S., Zhang, L. and Agarwal, D., 2017. *Bringing Salary Transparency to the World: Computing Robust Compensation Insights via LinkedIn Salary*. arXiv preprint arXiv:1703.09845. Available at: <https://arxiv.org/abs/1703.09845> [Accessed 15 May 2025].
4. Mincer, J., 1974. *Schooling, Experience, and Earnings*. New York: National Bureau of Economic Research.
5. Heckman, J.J., Lochner, L.J. and Todd, P.E., 2003. *Fifty Years of Mincer Earnings Regressions*. NBER Working Paper No. 9732. Available at: <https://www.nber.org/papers/w9732> [Accessed 15 May 2025].
6. Rosen, S., 2004. Distinguished Fellow: Mincering Labor Economics. *Journal of Economic Perspectives*, 18(3), pp.173-190.
7. Lemieux, T., 2006. The 'Mincer equation' Thirty Years after Schooling, Experience, and Earnings. In: S. Grossbard, ed., *Jacob Mincer: A Pioneer of Modern Labor Economics*. New York: Springer, pp.127-145.
8. Björklund, A. and Kjellström, C., 2002. Estimating the return to investments in education: how useful is the standard Mincer equation?. *Economics of Education Review*, 21(3), pp.233-241.
9. Borjas, G.J., 2016. *Labor Economics*. 7th ed. New York: McGraw-Hill Education.
10. Cahuc, P., Carcillo, S. and Zylberberg, A., 2014. *Labor Economics*. 2nd ed. Cambridge, MA: The MIT Press.
11. Heckman, J.J., Lochner, L.J. and Todd, P.E., 2006. Earnings Functions, Rates of Return and Treatment Effects: The Mincer Equation and Beyond. In: E.A. Hanushek and F. Welch, eds., *Handbook of the Economics of Education*, Volume 1. Amsterdam: Elsevier, pp.307-458.
12. Polachek, S.W., 2007. *Earnings Over the Lifecycle: The Mincer Earnings Function and Its Applications*. IZA Discussion Paper No. 3181. Available at: <https://www.iza.org/publications/dp/3181> [Accessed 15 May 2025].
13. Ujeniya, K., 2024. *Predicting Salaries of Data Professionals Using Machine Learning*. Medium. Available at: <https://medium.com/@krishujeniya/predicting-salaries-of-data-professionals-using-machine-learning-a77856f8a7b9> [Accessed 15 May 2025].
14. Zhou, Z., 2024. *Forecasting Data Science Professionals' Salaries Using Machine Learning*. *AIP Conference Proceedings*, 3244(1), p.030034. Available at:

- <https://pubs.aip.org/aip/acp/article/3244/1/030034/3322836/Forecasting-data-science-professionals-salaries> [Accessed 15 May 2025].
15. Quan, T.Z. and Raheem, M., 2022. Salary Prediction in Data Science Field Using Specialized Skills and Job Benefits – A Literature Review. *Journal of Applied Technology and Innovation*, 6(3), pp.70–80. Available at: https://www.researchgate.net/publication/362280362_Salary_Prediction_in_Data_Science_Field_Using_Specialized_Skills_and_Job_Benefits_-A_Literature_Review [Accessed 15 May 2025].
 16. Kaggle, 2025. *Salary Prediction for Beginners Dataset*. Available at: <https://www.kaggle.com/datasets/rkiattisak/salaly-prediction-for-beginner> [Accessed 15 May 2025].
 17. JobsPikr, 2025. *Job Data Mining Techniques for Recruitment Agencies*. JobsPikr Blog. Available at: <https://www.jobspikr.com/blog/job-data-mining-techniques-for-recruitment-agencies> [Accessed 15 May 2025].
 18. Mishra, A., 2023. Machine Learning Approaches for Salary Prediction: A Comparative Study. *International Journal of Computer Applications*, 182(15), pp.25-30.
 19. Singh, R. and Sharma, P., 2022. Predictive Modeling for Job Offer Acceptance Using Data Mining Techniques. *Journal of Human Resource Management*, 10(2), pp.45-52.
 20. Patel, S. and Desai, K., 2021. Analyzing the Impact of Educational Background on Employment Outcomes Using Data Mining. *International Journal of Data Science*, 5(1), pp.10-18.
 21. Kumar, N., 2020. A Survey on Salary Prediction Models Using Machine Learning Techniques. *Journal of Artificial Intelligence Research*, 12(3), pp.100-110.
 22. Gupta, A. and Verma, R., 2019. Data Mining Techniques for Predicting Employee Turnover and Salary. *International Journal of Business Analytics*, 6(4), pp.55-65.
 23. Lee, J. and Kim, H., 2018. Predicting Job Offers Using Educational and Demographic Data: A Data Mining Approach. *Journal of Employment Studies*, 9(2), pp.30-40.
 24. Chen, L., 2017. The Role of Machine Learning in Predicting Employment Outcomes for Graduates. *Education and Information Technologies*, 22(5), pp.2001-2015.
 25. Rahman, M. and Hasan, S., 2016. Analyzing Job Market Trends Using Data Mining Techniques. *International Journal of Computer Science and Information Security*, 14(6), pp.120-125.

Git Access: <https://github.com/AneriShroff/Education-to-Employement>