

## Course Seven

### Google Advanced Data Analytics Capstone



#### Instructions

Use this PACE strategy document to record your decisions and reflections as a data professional as you work through the capstone project. As a reminder, this document is a resource guide that you can reference in the future and a space to help guide your responses and reflections posed at various points throughout the project.

#### Portfolio Project Recap

Many of the goals you accomplished in your individual course portfolio projects are incorporated into the Advanced Data Analytics capstone project, including:

- Create a project proposal
- Demonstrate understanding of the form and function of Python
- Show how data professionals leverage Python to load, explore, extract, and organize information through custom functions
- Demonstrate understanding of how to organize and analyze a dataset to find the “story.”
- Create a Jupyter notebook for exploratory data analysis (EDA)
- Create visualization(s) using Tableau
- Use Python to compute descriptive statistics and conduct a hypothesis test
- Build a multiple linear regression model with ANOVA testing
- Evaluate the model
- Demonstrate the ability to use a notebook environment to create a series of machine learning models on a dataset to solve a problem
- Articulate findings in an executive summary for external stakeholders



## Project proposal

# Employee Attrition Prediction for Salifort Motors Project Proposal

## Overview

*This project aims to help Salifort Motors reduce employee turnover by analyzing historical HR data and predicting which employees are at risk of leaving. Using exploratory data analysis and machine learning models like Logistic Regression, Random Forest, and XGBoost, we derived actionable insights to support HR strategies.*

Milestones	Tasks	PACE stages
Project Planning	Identify business problem, audience, and key questions.	Plan
Data Cleaning and Exploration	Load data, clean missing/outlier values, explore variables (EDA).	Analyze
Statistical Testing and Hypothesis	Conduct descriptive stats and hypothesis tests to validate assumptions.	Analyze
Data Modeling	Build and evaluate Logistic Regression, Random Forest, and XGBoost models.	Construct
Insights and Recommendations	Interpret model results and generate business recommendations.	Construct
Final Report and Visualization	Summarize findings, create Tableau dashboards, compile executive summary.	Execute
Reflection and Submission	Reflect on project process, finalize deliverables, and submit portfolio.	Execute



## Data Project Questions & Considerations



### PACE: Plan Stage

#### Foundations of data science

- Who is your audience for this project?
  - The audience includes HR managers, business leaders, and decision-makers at Salifort Motors who are responsible for employee retention and organizational performance.
- What are you trying to solve or accomplish? And, what do you anticipate the impact of this work will be on the larger business need?
  - The goal is to develop a predictive model that identifies employees at risk of leaving the company. By doing so, the business can proactively address factors contributing to attrition, improve retention, and make informed strategic HR decisions. This project can reduce costly turnover, improve employee satisfaction, and enable data-driven retention strategies, ultimately saving time and resources and supporting business continuity.
- What questions need to be asked or answered?
  - What factors contribute most to employee attrition?
  - Can we accurately predict who is likely to leave?
  - What interventions could reduce employee churn?
  - Are there department or salary-based patterns?
- What resources are required to complete this project?
  - Kaggle dataset: [HR Analytics](#)
  - sklearn, pandas, seaborn, matplotlib
  - Course materials
- What are the deliverables that will need to be created over the course of this project?
  - Cleaned dataset and exploratory data analysis (EDA)
  - Visualizations (boxplots, histograms, heatmaps, etc.)
  - Predictive models (Logistic Regression, Random Forest, XGBoost)
  - Model evaluation metrics and insights
  - PACE Strategy document with reflections
  - Final summary of recommendations and insights for stakeholders



## Get Started with Python

- How can you best prepare to understand and organize the provided information?
  - By importing the dataset and reviewing its structure using `.info()` and `.describe()`.
  - Checking for missing or duplicate values.
  - Exploring variable types and distributions to guide preprocessing and modeling steps.
- What follow-along and self-review codebooks will help you perform this work?
  - Course notebooks from previous lessons on EDA, classification models, and data preprocessing.
  - Kaggle kernels or documentation on working with HR datasets.
  - Scikit-learn's official documentation: [https://scikit-learn.org/stable/user\\_guide.html](https://scikit-learn.org/stable/user_guide.html)
- What are a couple of additional activities a resourceful learner would perform before starting to code?
  - Research best practices in employee attrition prediction models.
  - Review the business context of attrition to frame analysis goals.
  - Make a backup copy of the dataset to avoid irreversible changes.
  - Create a checklist of EDA and preprocessing tasks to ensure consistency.

## Go Beyond the Numbers: Translate Data into Insights

- What are the data columns and variables, and which ones are most relevant to your deliverable?
  - The dataset includes the following variables: `satisfaction_level` (numeric, float64), `last_evaluation` (numeric, float64), `number_project` (numeric, int64), `average_monthly_hours` (numeric, int64), `time_spend_company` (numeric, int64), `work_accident` (binary, int64), `promotion_last_5years` (binary, int64), `department` (categorical, object), `salary` (categorical, object), `left` (binary target variable: 0 = stayed, 1 = left)
  - **Most relevant variables** for predicting employee attrition include: `satisfaction_level`, `last_evaluation`, `number_project`, `average_monthly_hours`, `promotion_last_5years`, `salary`, and `time_spend_company`.



- What units are your variables in?
  - `satisfaction_level` and `last_evaluation`: ratio between 0 and 1.
  - `average_monthly_hours`: number of hours. `number_project`, `time_spend_company`: count.
  - `work_accident`, `promotion_last_5years`, `left`: binary (0 or 1).
  - `department` and `salary`: categorical (string labels)
- What are your initial presumptions about the data that can inform your EDA, knowing you will need to confirm or deny with your future findings?
  - Employees with low satisfaction are more likely to leave. High working hours may correlate with burnout and turnover. Lack of promotion in recent years might drive attrition. Salary level could strongly influence retention. Certain departments may have higher turnover rates.
- Is there any missing or incomplete data?
  - No missing values were found in the dataset after initial inspection using `.isnull().sum()`.
- Are all pieces of this dataset in the same format?
  - Mostly yes. All numeric variables are correctly typed (`int64` or `float64`). Categorical variables (`department`, `salary`) are of object type and need one-hot encoding before modeling.
- Which EDA practices will be required to begin this project?
  - Check for and remove duplicate rows. Use `.describe()` and `.info()` to get summary statistics. Visualize distributions using histograms and box plots. Use scatter plots and correlation heatmaps to explore variable relationships. Identify and handle outliers, especially in `time_spend_company`



## The Power of Statistics

- What is the main purpose of this project?
  - The purpose is to analyze employee attrition data at Salifort Motors to identify key factors contributing to employee turnover and build predictive models. This analysis can help the company improve employee retention by taking data-driven action.
- What is your research question for this project?
  - **What are the most significant factors that influence whether an employee will leave Salifort Motors, and how can we use this information to predict attrition?**
- What is the importance of random sampling? In this case, what is an example of sampling bias that might occur if you didn't use random sampling?
  - Random sampling ensures that the training and testing datasets represent the overall population without bias, making the model generalizable.
  - **Sampling bias example:** If we only use recent hires or only one department's employees, the model might overfit to patterns in that group and fail to predict turnover accurately across the organization.

## Regression Analysis: Simplify Complex Data Relationships

- Who are your stakeholders for this project?
  - The stakeholders are HR managers, executive leadership at Salifort Motors, and team leads who are responsible for improving employee satisfaction, reducing turnover, and maintaining team performance.
- What are you trying to solve or accomplish?
  - We aim to identify the key factors influencing employee attrition and build accurate predictive models that can help the company proactively retain valuable employees.
- What are your initial observations when you explore the data?
  - Employees with lower satisfaction levels and higher working hours tended to leave. Employees who did not receive a promotion in the last 5 years or had high project counts showed higher attrition. Certain departments (e.g., sales, technical) and low-salary groups experienced more turnover. Outliers in time spent at the company also correlated with higher attrition.



- What resources do you find yourself using as you complete this stage? (Make sure to include the links.)
  - pandas documentation, scikit-learn documentation, Seaborn visualization guide, [Google's Advanced Data Analytics Certificate resources](#)
- Do you have any ethical considerations in this stage?
  - Yes, it's important to ensure that no discriminatory bias exists in model decisions (for example, based on salary or department). The model is not used to unfairly target or penalize employees. Employee data privacy is respected throughout analysis and reporting.

## The Nuts and Bolts of Machine Learning

- What am I trying to solve?
  - We are trying to predict whether an employee is likely to leave the company based on various features such as satisfaction level, workload, promotions, salary, department, etc. The goal is to help HR proactively address attrition risk.
- What resources do you find yourself using as you complete this stage?
  - [XGBoost documentation](#), scikit-learn ML guide, [Kaggle dataset for HR attrition](#), Google's capstone project template, and lab instructions
- Is my data reliable?
  - Yes, after removing duplicates, checking for and handling outliers, encoding categorical variables, and verifying no missing values, the dataset is cleaned and reliable for modeling.
- Do you have any additional ethical considerations in this stage?
  - Yes, ensure that model predictions are not biased against specific departments or salary groups. Prevent the misuse of predictions to unfairly treat or exclude employees. Transparently communicate how the model is used and ensure human oversight in decision-making.
- What data do I need/would I like to see in a perfect world, to answer this question?
  - Employee satisfaction over time, not just as a snapshot. Manager feedback scores or 360-degree reviews. Detailed work-life balance indicators (e.g., overtime hours, workload spikes). Exit interview reasons or employee sentiment data
- What data do I have/can I get?
  - We have: Satisfaction level, Evaluation score, Number of projects, Average monthly hours, Years at the company, Promotion and accident records, Department and salary, Whether the employee left or stayed



- What metric should I use to evaluate the success of my business objective? Why?
  - We should focus on:
    - **F1-Score** (balances false positives and false negatives)
    - **Recall** (to minimize missing potential leavers)
    - **Accuracy** (to check overall performance)
    - **AUC** (for a more balanced view of performance)

These help ensure the model doesn't just predict the majority class (those who stayed) but accurately identifies those likely to leave.





## Data Project Questions & Considerations



### PACE: Analyze Stage

#### Get Started with Python

- Will the available information be sufficient to achieve the goal based on your intuition and the analysis of the variables?
  - Yes. The dataset contains key features such as satisfaction level, evaluation score, project load, average monthly hours, tenure, salary, and department. These provide strong indicators of attrition risk, making the data sufficient to build meaningful models.

#### Go Beyond the Numbers: Translate Data into Insights

- What steps need to be taken to perform EDA in the most effective way to achieve the project goal?
  - Check data types and null values, Identify and remove duplicate entries, Detect and treat outliers, Explore class imbalance (left vs. stayed), Visualize relationships using boxplots, histograms, heatmaps, and scatter plots, Convert categorical features to numerical via encoding, Interpret trends (example, satisfaction vs. attrition, salary level vs. attrition).
- Do you need to add more data using the EDA practice of joining? What type of structuring needs to be done to this dataset, such as filtering, sorting, etc.?
  - No additional data is needed for this capstone project.  
Structuring includes: Encoding categorical columns (department, salary), removing duplicates and irrelevant columns, filtering out extreme outliers, and sorting during visualization for clarity.
- What initial assumptions do you have about the types of visualizations that might best be suited for the intended audience?
  - **Boxplots** for visualizing outliers and distributions of numeric features, **Histograms** to understand frequency distributions, **Scatter plots** to explore correlation between variables like satisfaction vs. evaluation, **Heatmaps** to understand overall correlation structure between numeric variables  
These help business stakeholders visually connect data trends to potential HR issues.



## The Power of Statistics

- Why are descriptive statistics useful?
  - They summarize the central tendencies, variability, and shape of the data distribution. This helps in identifying patterns, spotting anomalies, and guiding model selection.
- What is the difference between the null hypothesis and the alternative hypothesis?
  - **Null Hypothesis ( $H_0$ ):** There is no relationship between the feature and attrition.
  - **Alternative Hypothesis ( $H_1$ ):** There is a significant relationship between the feature and attrition.

These hypotheses guide statistical testing and model construction.

## Regression Analysis: Simplify Complex Data Relationships

- What are some purposes of EDA before constructing a multiple linear regression model?
  - Identify multicollinearity among features, detect outliers and skewed distributions, confirm linearity between independent and dependent variables, choose relevant predictors, and Address missing values or data quality issues.
- Do you have any ethical considerations in this stage?
  - Yes: Avoid reinforcing systemic biases through variables like salary or department, ensure insights do not lead to discriminatory policies, maintain employee anonymity and data privacy during analysis.

## The Nuts and Bolts of Machine Learning

- What am I trying to solve? Does it still work? Does the plan need revising?
  - We are predicting employee attrition to reduce turnover. The plan is still valid as the initial EDA supports the goal and confirms that the features correlate well with the target variable.
- Does the data break the assumptions of the model? Is that ok, or unacceptable?
  - Some models, like logistic regression, assume: No multicollinearity, Linearity in logit. These assumptions may be slightly violated, but robust tree-based models like Random Forest and XGBoost don't rely on these assumptions, making them more suitable.
- Why did you select the X variables you did?
  - Selected variables (like satisfaction level, average hours, promotions, salary, etc.) were chosen based on EDA findings, correlation strength, and business relevance. These are all meaningful predictors of whether an employee stays or leaves.



- What are some purposes of EDA before constructing a model?
  - **Exploratory Data Analysis (EDA)** helps ensure the data is clean, meaningful, and ready for modeling. Key purposes include: **Understanding data structure and distributions:** EDA reveals the shape, spread, and summary statistics of each variable. **Detecting missing values and duplicates:** Helps identify and address any incomplete or redundant data. **Identifying outliers:** Outliers can distort model training, especially in algorithms sensitive to extreme values. **Assessing feature relevance:** EDA can highlight which features correlate with the target variable and should be considered for the model. **Uncovering relationships and patterns:** Visualization tools like scatter plots and heatmaps reveal trends, correlations, and dependencies between variables. **Checking class imbalance:** Important for classification tasks, especially to determine if resampling or special evaluation metrics are needed. **Informing preprocessing decisions:** Guides encoding, scaling, or transformation strategies for categorical and numerical features.
- What has the EDA told you?
  - Low satisfaction, no promotions, and longer tenure are linked to attrition. Employees with more projects or high monthly hours are more likely to leave. Lower salary tiers see higher exit rates. Departments like sales and technical show higher attrition patterns.
- What resources do you find yourself using as you complete this stage?
  - Pandas Documentation, Matplotlib & Seaborn for visualization, [Kaggle Discussion Threads on HR Analytics](#), Course content on EDA and model selection
- Do you have any ethical considerations in this stage?
  - Yes. Interpretation of results should not stigmatize departments or salary groups. Also, machine learning models must be applied with fairness and transparency in HR decisions.



## Data Project Questions & Considerations



### **PACE: Construct Stage**

#### **Get Started with Python**

- Do any data variable averages look unusual?
  - Upon conducting descriptive analysis and visualization, most of the variables in the dataset appeared within reasonable ranges. However, one noticeable outlier was the `time_spend_company` variable, where a small group of employees had significantly longer tenures compared to the rest, which pulled the average slightly higher. While this could indicate loyalty or seniority, further inspection revealed that these outliers were often associated with higher rates of attrition, suggesting that unusually high tenure might correspond to burnout or stagnation for some employees.
- How many vendors, organizations, or groupings are included in this total data?
  - The dataset includes employees from several departments, which act as internal organizational groupings within Salifort Motors. Specifically, the `department` column captures groupings such as sales, technical, support, and management, among others. Additionally, the `salary` column includes groupings based on employee compensation: low, medium, and high. Together, these provide a layered structure to analyze patterns of attrition and employee behavior across different segments of the organization.

#### **Go Beyond the Numbers: Translate Data into Insights**

- What data visualizations, machine learning algorithms, or other data outputs will need to be built in order to complete the project goals?
  - To fulfill the project goals, a series of visualizations and predictive models were created. Visualizations such as histograms, boxplots, and heatmaps helped explore feature distributions, detect outliers, and understand variable correlations. For predictive modeling, Logistic Regression was first constructed to establish a baseline, followed by more robust classifiers such as Random Forest and XGBoost. These models allowed for a deeper understanding of which features most influenced employee attrition and how accurately we could predict it.



- What processes need to be performed in order to build the necessary data visualizations?
  - To build effective visualizations, the dataset first needed to be cleaned and preprocessed. This included removing duplicates, handling outliers, checking for null values, and encoding categorical variables for compatibility with machine learning tools. Once the data was prepared, visualization libraries like `matplotlib` and `seaborn` were used to generate plots. The boxplots visualized distribution and outliers; histograms illustrated frequencies of features; and the heatmap provided an overall view of feature correlations. These visualizations supported both exploration and communication of key findings.
- Which variables are most applicable for the visualizations in this data project?
  - The variables that stood out most during visualization were `satisfaction_level`, `last_evaluation`, `average_monthly_hours`, `number_project`, and `time_spend_company`, as they exhibited the strongest signals when comparing employees who stayed versus those who left. Additionally, the salary and department categorical variables revealed important trends when visualized using countplots and bar charts, particularly in relation to attrition rates.
- Going back to the Plan stage, how do you plan to deal with the missing data (if any)?
  - Fortunately, this dataset had no missing values, as confirmed during the cleaning phase. If there had been missing data, I would have considered imputation methods such as mean or median substitution for numerical features, or the mode for categorical ones. Alternatively, if a large portion of a feature was missing, I might have opted to drop that column entirely, depending on its importance and correlation with the target variable.

## The Power of Statistics

- How did you formulate your null hypothesis and alternative hypothesis?
  - The hypothesis formulation stemmed from the core business question: What factors contribute to employee attrition? The null hypothesis ( $H_0$ ) posited that no relationship exists between the selected features (for example, satisfaction level, promotion history, number of projects) and employee attrition. The alternative hypothesis ( $H_1$ ) suggested that at least one of these features significantly influences whether an employee leaves the company. This guided both the exploratory and predictive phases of the analysis.



- What conclusion can be drawn from the hypothesis test?
  - The models and EDA both indicated that several features have a statistically significant association with employee attrition. For example, low satisfaction levels and a lack of promotion over five years were strong indicators that an employee might leave. Therefore, we rejected the null hypothesis and accepted the alternative: that certain features do meaningfully influence attrition, justifying the need for targeted interventions based on these factors.

### **Regression Analysis: Simplify Complex Data Relationships**

- Do you notice anything odd?
  - One odd observation came from the Logistic Regression model, which showed a relatively low recall for predicting employee exits, meaning it often failed to catch those who were likely to leave. This indicated that the model wasn't capturing the nuances of employee behavior well. As a result, more complex models like Random Forest and XGBoost were implemented, both of which showed significant improvements in precision, recall, and overall accuracy. These ensemble models better captured non-linear relationships and variable interactions. If time and computational resources permitted, hyperparameter tuning and additional feature engineering could further enhance these models.
- Can you improve it? Is there anything you would change about the model?
  - Answered in the previous question.



## The Nuts and Bolts of Machine Learning

### Is there a problem? Can it be fixed? If so, how?

- One of the problems encountered during model construction was the limited performance of the baseline Logistic Regression model, particularly in identifying employees who left the company (low recall for class 1). This indicated that the model was not effectively capturing complex patterns in the data, possibly due to its linear nature. To fix this, I applied more sophisticated ensemble models, Random Forest and XGBoost, which could handle nonlinear relationships and better model interactions between features. These models significantly improved both precision and recall, especially for predicting employee attrition. Additionally, data imbalance between classes was an underlying issue, and techniques like class weighting or SMOTE (Synthetic Minority Over-sampling Technique) could further help if performance still needs improvement.
- Which independent variables did you choose for the model, and why?
  - I selected variables that showed strong correlations with the target variable (left) during exploratory data analysis. These included `satisfaction_level`, `last_evaluation`, `number_project`, `average_monthly_hours`, and `time_spend_company`, as each of these features revealed meaningful trends in employee behavior. Additionally, binary features such as `work_accident` and `promotion_last_5years` were included due to their logical relevance to employee satisfaction and motivation. To capture organizational context, I one-hot encoded the `department` and `salary` columns. These features were chosen based on both domain understanding and empirical patterns observed during data exploration, making them good candidates for predictive modeling.
- How well does your model fit the data? (What is my model's validation score?)
  - The model fit was evaluated using multiple performance metrics. The Random Forest model achieved an overall accuracy of 97.8%, while the XGBoost model slightly outperformed it with an accuracy of 97.95%. Both models demonstrated excellent precision and recall, especially for the minority class of employees who left the company. These validation scores, along with the high F1-scores and strong confusion matrix outcomes, indicate that the models generalize well and have strong predictive capabilities on unseen data.



- Can you improve it? Is there anything you would change about the model?
  - While the current results are promising, there are still opportunities for refinement. Future improvements could include hyperparameter tuning using GridSearchCV or RandomizedSearchCV to optimize the models further. Incorporating techniques to handle class imbalance, such as SMOTE or adjusting class weights, might enhance the recall of class 1 even more. Additionally, feature engineering, such as creating interaction terms or aggregating engagement trends over time, could lead to better insight and model performance. Moreover, a deeper look into feature importance could help eliminate redundant variables, simplifying the model without sacrificing accuracy.
- Do you have any ethical considerations in this stage?
  - Yes, ethical concerns are central when developing machine learning models involving human subjects, especially in employment contexts. It's essential to ensure that the model does not reinforce biases, such as penalizing certain departments or salary groups unfairly, due to historical trends embedded in the data. All data must be anonymized to protect individual privacy. Moreover, predictions should not be used for punitive actions, but rather as supportive tools to enhance employee retention and satisfaction through targeted interventions. Transparency in how predictions are used and the fairness of outcomes must be continuously evaluated.



## Data Project Questions & Considerations



### PACE: Execute Stage

#### Get Started with Python

- Given your current knowledge of the data, what would you initially recommend to your manager to investigate further prior to performing an exploratory data analysis?
  - Before starting the EDA, I would recommend that the manager investigate the completeness, cleanliness, and reliability of the employee dataset. It's crucial to ensure that there are no duplicate records, missing values, or inconsistent formats, especially across important variables such as department, satisfaction level, and time spent at the company. Verifying that the dataset represents the full spectrum of employees across tenure, departments, and job roles is also important to avoid sampling bias. Lastly, I would suggest validating that the target variable, "left," is accurate and reflects actual employee attrition, as all further modeling decisions will rely on it.
- What data initially presents as containing anomalies?
  - During the EDA phase, one variable that presented anomalies was `time_spend_company`, which had several outliers. This was evident from the boxplot visualization, where a small number of employees had significantly higher years at the company compared to the majority. These high values might represent executive-level employees or edge cases and should be further investigated to determine if they distort the model. Other variables like `average_monthly_hours` were surprisingly clean and didn't contain visible outliers after review.
- What additional types of data could strengthen this dataset?
  - Adding more detailed employee-level data would significantly enhance the analysis. For example, including employee job titles, performance scores, engagement survey results, exit interview summaries, or even manager ratings could provide deeper insight into why employees choose to stay or leave. Additionally, data related to training history, benefits utilization, or workplace satisfaction feedback could introduce new dimensions to better model employee attrition.

#### Go Beyond the Numbers: Translate Data into Insights

- What key insights emerged from your EDA and visualizations(s)?
  - The EDA revealed that employees with lower satisfaction levels and higher working hours tended to leave the company more often. Attrition was also higher among those who had not been promoted in the past five years and those with either very few or too many projects. Departments such as sales and technical experienced higher turnover compared to others, and low salary was another strong indicator of attrition. These trends were consistently visualized through histograms, boxplots, and correlation heatmaps.



- What business recommendations do you propose based on the visualization(s) built?
  - Based on the data visualizations, I would recommend that management focus on improving employee satisfaction through regular feedback mechanisms and work-life balance initiatives. Employees with high workload or long tenure without promotion could benefit from development opportunities or role changes to reduce burnout. Moreover, departments with higher turnover, such as sales, should be reviewed for cultural or workload-related issues. Lastly, compensation adjustments should be considered for low-salary bands to reduce voluntary exits.
- Given what you know about the data and the visualizations you were using, what other questions could you research for the team?
  - Given the data, it would be valuable to further investigate the interaction effects between variables. For example, does salary influence satisfaction differently across departments? Or does high workload only lead to attrition when satisfaction is also low? These layered insights could be explored through advanced modeling or segmentation analysis. Additionally, examining attrition trends over time or in response to specific internal events (e.g., policy changes, management shifts) would add depth to the findings.
- How might you share these visualizations with different audiences?
  - To share these insights with leadership, I would prepare a concise dashboard in Tableau that highlights the key metrics: attrition rates by department, salary bands, and satisfaction level distributions. For HR analysts, I would provide interactive filters so they can explore patterns within subgroups. For general staff, simplified infographics or summary slides could help communicate trends and promote transparency about initiatives being taken in response.

## The Power of Statistics

- What key business insight(s) emerged from your A/B test?
  - Although a formal A/B test wasn't conducted in this project, the insights derived from the EDA and modeling process provide similar value. For instance, comparing groups like those who received a promotion vs. those who didn't, or employees with high vs. low satisfaction, served as a quasi-A/B comparison. These comparisons revealed clear differences in attrition likelihood, suggesting that workplace policies such as promotion frequency and workload balance significantly affect retention.
- What business recommendations do you propose based on your results?
  - Based on the findings, I recommend that Salifort Motors implement targeted employee retention programs. This includes focusing on increasing employee satisfaction through internal surveys and timely responses, reducing excessive work hours where possible, and developing clear paths for career advancement. Departments with high turnover should undergo a deeper qualitative assessment, and compensation structures should be reevaluated, especially for those in the low salary bracket, as low pay was a consistent predictor of attrition.



## Regression Analysis: Simplify Complex Data Relationships

- To interpret model results, why is it important to interpret the beta coefficients?
  - Interpreting beta coefficients in logistic regression allows stakeholders to understand the direction and strength of influence each independent variable has on the likelihood of employee attrition. For example, a negative coefficient for `satisfaction_level` suggests that as satisfaction increases, the likelihood of leaving decreases. This interpretability helps decision-makers prioritize interventions. Without understanding the meaning of these coefficients, the model would function as a black box, making it harder to drive actionable business decisions.
- What potential recommendations would you make to your manager/company?
  - I would recommend that the company invest in continuous performance reviews and recognition programs to keep satisfaction high, especially among employees who have not recently been promoted or those handling multiple projects. Using the model to flag high-risk individuals early could allow HR to intervene proactively. Creating transparent growth paths and balancing workloads might help reduce the turnover rate significantly.
- Do you think your model could be improved? Why or why not? How?
  - Yes, the model could be improved by incorporating more granular and real-time employee data. While the current dataset performs well, adding time-series features like month-to-month satisfaction, performance trends, or event-based flags (e.g., manager change) could increase predictive accuracy. Additionally, using ensemble techniques or fine-tuning hyperparameters in XGBoost could improve precision and recall for the minority class (`left = 1`).
- What business recommendations do you propose based on the models built?
  - The models confirm that high workload, long tenure without promotion, low salary, and dissatisfaction are strong predictors of employee attrition. Based on this, I recommend a data-driven early warning system to flag high-risk employees and prompt timely interventions. Salary adjustments, career growth frameworks, and employee engagement initiatives should be prioritized, particularly in departments like sales and technical, where attrition is higher.
- What key insights emerged from your model(s)?
  - Key insights include the finding that employee satisfaction is the most influential factor in retention. Also, departments and salary categories have a strong influence on attrition, and workload measured by `number_project` and `average_monthly_hours` plays a significant role. Notably, the advanced models like XGBoost achieved over 97% accuracy, showing high predictive capability when these features are used effectively.



- Do you have any ethical considerations at this stage?
  - Yes, ethical considerations include ensuring the model does not reinforce existing biases, such as penalizing employees from specific departments or salary groups without considering context. It's also important to maintain transparency with employees about how data is being used. Predictive models should be used as supportive tools for decision-making, not as the sole basis for HR actions, to avoid unfair profiling.

## The Nuts and Bolts of Machine Learning

- What key insights emerged from your model(s)?
  - The most advanced model, XGBoost, provided high precision and recall, especially for predicting employees who are likely to leave. It reinforced the insight that factors such as low satisfaction, lack of promotion, and high workload are reliable indicators of attrition. The model also demonstrated that machine learning can effectively support HR decision-making when trained on relevant employee data.
- What are the criteria for model selection?
  - Model selection was based on accuracy, precision, recall, and interpretability. While logistic regression was easier to interpret and explain, its recall for identifying employees who left was relatively low. Random Forest and XGBoost provided significantly better overall performance, with XGBoost slightly outperforming Random Forest. These models were preferred for their balance between performance and generalization capability.
- Does my model make sense? Are my final results acceptable?
  - Yes, the model makes sense and aligns with known business logic. For example, dissatisfied employees or those overworked and underpaid are more likely to leave, which the model correctly identified. The final accuracy above 97% in Random Forest and XGBoost is not only acceptable but excellent for predictive modeling in HR analytics.
- Were there any features that were not important at all? What if you take them out?
  - Some of the one-hot encoded department variables had less influence on the model's predictions compared to variables like satisfaction or salary. Removing these features slightly reduced performance, but not significantly. However, their inclusion helps capture department-level trends, so keeping them is helpful for business insights.
- Given what you know about the data and the models you were using, what other questions could you address for the team?
  - Other questions that could be explored include forecasting future attrition rates over time, identifying the ideal combination of workload and satisfaction to maximize retention, and evaluating which managers have higher retention rates under their supervision. Segmenting the model by tenure or department could also provide more tailored insights.



- What resources do you find yourself using as you complete this stage?
  - Key resources included scikit-learn documentation, XGBoost library documentation, Google's course materials, Kaggle notebooks for attrition analysis, and official Python pandas and matplotlib documentation. These helped in preprocessing, visualization, and model building.
- Is my model ethical?
  - The model is ethical to the extent that it is trained on relevant and non-sensitive attributes and used for supportive decision-making rather than punitive action. It is essential that the model's use is communicated clearly to employees and that it does not replace human judgment in employment decisions.
- When my model makes a mistake, what is happening? How does that translate to my use case?
  - When the model misclassifies an employee, it typically happens when that individual exhibits a mixed profile, for example, low satisfaction but high salary, or long tenure but recent promotion. In the business context, this could result in either missing an at-risk employee or falsely flagging someone as high-risk. Both cases require caution: HR should use the model as an early alert system, not a final decision-maker.