# HR Analytics — Predicting Employee Attrition

## Abstract

This report presents a practical machine learning workflow to predict employee attrition using an HR analytics dataset. We analyze the data pipeline, modeling choices, and results extracted from the provided Jupyter notebook. The solution compares baseline linear and tree-based classifiers and demonstrates how data preprocessing (feature scaling and a train–test split) affects performance. The best model achieved approximately **0.891** accuracy on the held-out test data. We discuss operational insights for HR teams—such as identifying risk drivers and prioritizing retention actions—and outline recommendations to improve model robustness (handling class imbalance, richer evaluation metrics, and feature engineering). The goal is to support proactive, data-driven talent management while keeping the approach reproducible and business-friendly.

## 1. Introduction

Organizations face significant costs when employees leave—both direct replacement costs and indirect loss of knowledge and productivity. Predictive attrition modeling helps HR teams flag at-risk employees early so managers can intervene with targeted actions (career growth, workload balancing, compensation adjustments). This project applies supervised learning to historical employee records to estimate attrition likelihood and surface drivers correlated with churn.

## 2. Data & Preprocessing

The notebook loads a standard HR attrition dataset and performs essential preprocessing steps. Numerical features are scaled (Standard/MinMax scaling), and the data is split into training and testing partitions to estimate generalization. In a production setting, categorical variables should be encoded (e.g., One−Hot) and pipeline components chained to avoid leakage. We recommend adding checks for missing values, outliers, and class imbalance (attrition is typically a minority class in HR data). Where imbalance is present, techniques such as class weights or SMOTE can improve recall for the minority (leaving) class.

## 3. Modeling Approach

The notebook trains multiple classifiers, including **DecisionTreeClassifier, LogisticRegression**. Linear models like Logistic Regression provide well-calibrated probabilities and interpretability, while tree-based models (Decision Trees, and potentially Random Forests/Gradient Boosting) capture nonlinear feature interactions. Hyperparameters may be tuned with Grid/Randomized Search; for a balanced evaluation, we advise stratified splits and cross−validation. Feature importance (from trees or coefficients from logistic regression) can guide HR policy by highlighting risk factors (e.g., overtime, tenure, job role, environment satisfaction).

## 4. Evaluation

Model performance was evaluated on a held out test set. The top model reached an accuracy of **0.891**. While accuracy is useful, attrition use cases benefit from metrics sensitive to minority class detection: precision, recall, F1−score, ROC AUC, and Precision Recall AUC. Confusion matrices should be reviewed to quantify False Negatives (missed at risk employees). Where feasible, threshold tuning (operating point selection) can trade precision vs. recall to match business priorities.

## 5. Insights & Business Impact

Even a modestly accurate model can add value when embedded in HR workflows. High−risk segments often share characteristics such as low engagement scores, longer commute, low compensation relative to peers, frequent overtime, or limited internal mobility. Recommended actions include scheduled career conversations, mentoring, workload balancing, targeted training, and compensation reviews. Importantly, models should be used ethically: ensure transparency, avoid proxy bias, and complement predictions with human judgment.

## 6. Recommendations

• Add robust categorical encoding and feature engineering (e.g., tenure buckets, overtime rate, internal transfers, manager changes).
• Address class imbalance with class_weight='balanced' or SMOTE; monitor recall and PR-AUC.
• Use cross−validation and hyperparameter tuning for Logistic Regression (C, penalty) and tree ensembles (depth, estimators).
• Log model versioning and build a reproducible pipeline (train/serve parity).
• Calibrate probabilities (Platt/isotonic) and choose decision thresholds based on business cost of False Negatives vs. False Positives.
• Track drift in production; refresh the model with new data on a schedule.

## 7. Conclusion

This HR analytics workflow demonstrates how supervised learning can predict employee attrition and inform proactive retention strategies. By strengthening preprocessing, addressing imbalance, expanding models to calibrated tree ensembles or gradient boosting, and adopting richer evaluation, the solution can move from proof of concept to production ready. The approach should remain transparent and fair, assisting—rather than replacing—people leaders in decision making.