

# Week 2: Literature Review

(Karan) RK Rajkumar

## Paper Information

**Title:** Predicting Employee Attrition Using Machine Learning Techniques **Authors:** Francesca Fallucchi, Marco Coladangelo, Romeo Giuliano, and Ernesto William De Luca  
**Year:** 2020 **Journal:** Computers (MDPI)

## 1. Background and Motivation

Enterprise turnover costs have long hampered organizations in translating their human capital into profitability and performance. This is particularly apparent the higher up the value chain an employee is situated or the more specialized the occupation in which they operate. Fallucchi et al. [1] stipulate that this financial liability is driven by the compounded costs of recruitment, onboarding, and retraining required in the succession process.

These observations are borne out by industry bodies such as the Society for Human Resource Management (SHRM), which estimates that total replacement costs often range between **0.5 and 2 times the annual salary** of the departing employee [2]. Furthermore, it is widely acknowledged by scholars in Organizational Behavior that engaged employees (who demonstrate positive behavioral intent) exhibit superior task performance and creativity, which in tandem drives organizational productivity [3].

Having defined the scope of the problem, the authors identify a critical gap in the conventional understanding of the motivational and relational determinants that influence an employee's decision to depart. To ameliorate this deficit, this study aims to isolate the **causal factors** behind attrition via categorical data exploration and machine learning techniques. The analysis utilizes the **IBM HR Analytics** dataset, consisting of 1,470 observations and 35 distinctive features, to construct a predictive framework for retention [1].

## 2. Methods Used

Fallucchi et al. [1] operationalize their study via the **Team Data Science Process (TDSP)**, a comprehensive five-step framework that rigorously enforces a sequence of Data Acquisition, Data Cleaning, Exploratory Data Analysis (EDA), Modeling, and Deployment. This methodology, which is commensurate with industry standards for provisioning machine learning pipelines, is designed not merely to predict outcomes but to utilize the feature importance metrics of the optimal model to isolate the specific features—or causal reasons—most strongly correlated with employee attrition.

The authors selected the **Attrition** variable of the IBM Dataset (a binary “Yes/No” indicator) as the target vector for both descriptive statistics and classification tasks. A critical precursor to the modeling phase was the execution of a correlation analysis across all predictors. This step is a sound preventative measure in statistical learning, as highly correlated or collinear features can artificially inflate variance, resulting in unstable models that are overfit to training data and violate the fundamental assumptions of algorithms such as Logistic Regression.

Finally, the study necessitated the application of **Categorical Encoding** to facilitate algorithmic ingestion. This process translates all qualitative, non-numerical columns into numerical arrays, a standard requirement when training models in scripting environments like Python or R. This transformation is essential because the underlying linear algebra of classifiers like Support Vector Machines (SVM) requires vector-based numerical inputs to function effectively.

## 3. Significance of the Work

The study offers a rigorous evaluation of various supervised learning algorithms, adjudicating their performance on structured human resources data for the specific objective of predicting attrition. Prioritized here is a targeted ensemble of classification models, with a definitive emphasis placed on **Recall** (Sensitivity) over Precision. Fallucchi et al. [1] underscore that in the domain of human capital risk, the minimization of “False Negatives” is paramount; in an enterprise context, this metric translates to the operational imperative of detecting an employee’s impending exit *before* the separation occurs.

Methodologically, the authors utilize a correlation matrix to identify pairwise dependencies between model variables. As noted by Shrestha [4], this is an industry-standard prophylactic measure against **multicollinearity**, a phenomenon in feature selection that can obscure the intelligibility of model coefficients and inflate variance.

Furthermore, the findings of this study aver that **MonthlyIncome**, **OverTime**, and **Age** are the dominant predictors of attrition. This empirical evidence is congruent with the theoretical framework of Kushwaha et al. [5], who posit that the remediation of an employee’s “Deficiency Needs” (specifically surplus work and impecunity) are necessary managerial antecedents to sustainable retention practices.

## 4. Connection to Other Work

This paper builds directly upon the foundational turnover research synthesized by Griffeth et al. [6], which codified **demographic characteristics** (e.g., **age** and **tenure**) as the baseline predictors of attrition. However, Fallucchi et al. [1] diverge from these traditional econometric paradigms by operationalizing **machine learning classifiers** to predict exit events rather than merely retrospectively explaining them.

**4.1 The Shift from Explanation to Prediction** Where prior literature focused on identifying “who leaves” using linear correlations, Fallucchi et al. confront the **stochastic reality of class imbalance** inherent in turnover data. Their benchmarking of classifiers demonstrates that **Naive Bayes** achieved the superior **Recall** for the attrition class, underscoring the strategic imperative that in HR analytics, the sensitivity to detect actual leavers is far more valuable than raw model accuracy. This finding validates the transition from *descriptive* heuristics to *predictive* risk mitigation.

**4.2 Addressing the “Black Box” Dilemma** A recurring critique in Organizational Behavior literature is the opacity of algorithmic decision-making. While Griffeth provided clear coefficients for demographic risk, Fallucchi’s work illuminates the inherent trade-off between **predictive sensitivity and interpretability**. By prioritizing the identification of “At-Risk” employees (High Recall) over the transparency of specific variables, this work situates itself not merely as a computational exercise, but as a bridge between the precision of Data Science and the interventionist mandates of Human Resource Management.

## 5. Relevance to Capstone Project

This paper serves as a direct technical benchmark for my Capstone project, “**Drivers of Attrition: An Application of the MARS Model**.” Since my project also utilizes the **IBM HR Analytics** dataset, the preprocessing pipeline detailed by Fallucchi (specifically the encoding of categorical variables and the normalization of salary data) will be adopted to ensure data quality.

**5.1 Operationalizing the MARS Model** The most critical contribution of this paper to my work is the empirical validation of variables that map to the **MARS Model** (Motivation, Ability, Role Perceptions, Situational Factors) as defined by McShane and Von Glinow [3]. While Fallucchi et al. provided a foundation, I will expand the feature selection to ensure all four pillars are represented:

- **Motivation (M):** The authors’ finding that **JobInvolvement** is a top predictor allows me to use this variable as a direct proxy for the “Motivation” pillar. Furthermore, though not emphasized in Fallucchi’s analysis, **RelationshipSatisfaction** (the interactional experience of the employee) has conventionally been understood as a critical driver of engagement and will thus feature prominently in my analysis.

- **Ability (A):** Within the MARS paradigm, Ability is defined as the transferable and non-transferable skills an employee brings to their vocation. To conceptualize this facet, I will utilize **Education**, **JobLevel**, and **TrainingTimesLastYear**, postulating that these metrics quantify the employee's technical competency and aptitude.
- **Role Perceptions (R):** Ostensibly excluded from the authors' analysis due to challenges in measurability, Role Perceptions remain a foundational plank of the MARS model. I will operationalize this construct using **YearsInCurrentRole** and **YearsWithCurrManager**, aiming to measure the clarity of behavioral expectations established through tenure and supervision.
- **Situational Factors (S):** The statistical significance of **OverTime** and **DistanceFromHome** provides the necessary justification to classify these as "Situational Factors"—external constraints that force an employee's hand regardless of their loyalty [1]. Additionally, I will leverage **EnvironmentSatisfaction** as a bellwether variable to determine an employee's dispositional state towards their physical surroundings and workplace conditions.
- **Distributive Justice:** The dominance of **MonthlyIncome** in the feature importance rankings supports the inclusion of "Distributive Justice" theory in my analysis, characterizing turnover as a rational economic response to wage inequity [3].

**5.2 Methodological Divergence** While Fallucchi et al. prioritized **Recall** to maximize the detection of quitters, my Capstone will diverge by prioritizing **Interpretability** via models such as Decision Trees, Random Forest, and Logistic Regression. My objective is not merely to flag "Flight Risks," but to explain *why* they are at risk by calculating odds ratios (e.g., "Employees working overtime are 3x more likely to leave"). This approach aims to provide actionable intelligence that transcends simple prediction, facilitating structural organizational interventions.

## References

- [1] F. Fallucchi, M. Coladangelo, R. Giuliano, and E. W. De Luca, “Predicting employee attrition using machine learning techniques,” *Computers*, vol. 9, no. 4, p. 86, 2020.
- [2] SHRM Executive Network, “The myth of replaceability: Preparing for the loss of key employees.” [Online]. Available: <https://www.shrm.org/executive-network/insights/myth-replaceability-preparing-loss-key-employees>
- [3] S. L. McShane and M. A. Von Glinow, *Organizational behavior: Emerging knowledge, global reality*, 10th ed. McGraw Hill, 2024.
- [4] N. Shrestha, “Detecting multicollinearity in regression analysis,” *American Journal of Applied Mathematics and Statistics*, vol. 8, no. 2, pp. 39–42, 2020.
- [5] B. P. Kushwaha, V. Tyagi, P. B. Sharma, and R. K. Singh, “Mediating role of growth needs and job satisfaction on talent sustainability in BPOs and call centres: An evidence from india,” *Journal of Public Affairs*, vol. 22, no. e2400, 2020.
- [6] R. W. Griffeth, P. W. Hom, and S. Gaertner, “A meta-analysis of antecedents and correlates of employee turnover: Update, moderator tests, and research implications for the next millennium,” *Journal of Management*, vol. 26, no. 3, pp. 463–488, 2000.