Predicting Employee Attrition in the Saudi Private Sector Using Machine Learning

Haya Alqahtani, Hana Almagrabi, Amal Alharbi Department of Information Systems-Faculty of Computing and Information Technology, King Abdulaziz University, Jeddah, Saudi Arabia

Abstract-Employee attrition represents a prominent issue facing organizations, as human capital represents one of the most valuable resources. Attrition refers to the voluntary or involuntary reduction in the number of employees, which can negatively affect profitability, reputation, and overall organizational performance. Therefore, a comprehensive understanding of this phenomenon, its causal factors, and the development of retention strategies is crucial for mitigating employee turnover. The purpose of this work is to predict employee attrition in the Saudi private sector and identify the key factors contributing to employee turnover using machine learning approaches. in addition, the research structurally evaluates the performance of multiple Machine Learning (ML) algorithms within the proposed framework to determine the most effective predictive model for employee attrition. This study utilized a training dataset obtained from an online survey targeting employees in the Saudi private sector in order to investigate employee attrition and identify its most prominent causes within this context. Thus, various Machine Learning (ML) algorithms, including Logistic Regression (LR), Support Vector Machine (SVM), K-Nearest Neighbors (KNN), Decision Tree (DT), Random Forest (RF), Extreme Gradient Boosting (XGBoost), Bagging ensemble, and Voting Classifier (VC) were evaluated. The results demonstrate that the Voting Classifier (VC) yielded the highest accuracy at 90%. Moreover, the analysis identified job opportunities and job titles as some of the most influential factors driving employee turnover.

Keywords—Employee attrition; attrition prediction; predictive models; machine learning; voting classifier; ensemble methods; Saudi private sector; employee turnover; employee retention; feature importance

I. INTRODUCTION

Employee attrition is a critical challenge for organizations due to its impact on productivity, operational costs, and long-term business sustainability. Attrition occurs when employees leave their organizations due to resignation, retirement, or involuntary termination. There are various key factors influencing employee attrition, including job satisfaction, work—life balance, compensation, work environment, and age. Employees who perceive themselves as undervalued or undercompensated are more likely to leave their organizations in pursuit of better opportunities [1]. Employee attrition can be measured using metrics such as sales growth, return on equity, customer service quality, and profitability [2].

As a matter of fact, voluntary turnover is particularly disruptive, resulting in the loss of institutional knowledge and interrupting ongoing projects [3]. Moreover, new hires require significant time to achieve the performance level of their predecessors, which reduces organizational efficiency [4], [5]. Studies have shown that voluntary employee turnover creates

significant costs due to the expenses of recruitment, hiring, training, and employee development, which can be more than 1.5 times an employee's annual salary [6].

Given the significant impact of attrition, predictive analytics has become an essential tool for forecasting employee turnover and formulating strategies to mitigate it [7]. Identifying employees prone to attrition enables organizations to proactively improve retention, enhance employee satisfaction, and reduce turnover-related costs [8]. HR analytics, known as talent analytics [9], plays a vital role in quantifying workforce factors that influence company outcomes, thereby supporting better decision-making [10]. Consequently, HR professionals must develop an in-depth understanding of employee profiles and the critical factors influencing turnover to accurately identify individuals at risk of departure [11], [12], [13].

Data is a crucial element for the success of a people analytics team [14]. However, in HR analytics, the focus should shift from simply collecting large volumes of data to ensuring the data provides meaningful insights and is used effectively to create organizational value [10]. Previous research on attrition prediction has predominantly utilized publicly available datasets, such as the IBM Employee Attrition dataset, as well as datasets obtained from platforms like Kaggle. Nevertheless, such datasets do not capture the complexity of real-world employee behavior and organizational dynamics, in addition to being characteristically imbalanced, causing predictive bias if not properly handled [8]. In order to address existing limitations, the primary objective of this study was to develop a realistic dataset and propose a robust ML framework for predicting employee attrition.

The main contributions of this paper are threefold:

- 1) Development of a contextually-grounded dataset for the Saudi private sector: This research introduces an empirically-collected dataset from 1,191 real employees across diverse Saudi industries, addressing the critical limitation of prior studies that rely predominantly on synthetic datasets (e.g., IBM HR dataset). The dataset uniquely captures region-specific factors shaped by Saudization policies and cultural dynamics, while incorporating previously underexplored variables such as job opportunities, emotional commitment, health issues, and sector diversity. This contribution fills a significant gap in attrition research for emerging markets and provides a reusable resource for future studies in similar contexts [15].
- 2) A comprehensive comparative ML framework with systematic feature selection: This study systematically evaluates eight ML algorithms—including conventional methods (LR, SVM, KNN, DT) and ensemble techniques (RF, XGBoost,

Bagging, Voting Classifier)—across seven feature selection configurations with SMOTE-based class balancing [16]. The Voting Classifier with RFE-selected features achieves 90% accuracy, demonstrating that ensemble methods significantly outperform conventional approaches while maintaining interpretability. This design provides generalizable methodological guidance for predictive modeling in real organizational settings where data imbalance and feature redundancy are prevalent. This framework provides HR practitioners with evidence-based guidance for selecting optimal modeling strategies for imbalanced attrition datasets.

3) Actionable insights revealing a paradigm shift in attrition drivers: The feature importance analysis reveals that dynamic developmental factors (job opportunities, career advancement, job titles) significantly outweigh traditional demographic factors (age, Academic degree) in predicting attrition. The study provides data-driven visions for targeted interventions—including career development programs and advancement pathways—that enable organizations to proactively identify at-risk employees.

Experimental results demonstrate the effectiveness of the proposed approach, with promising accuracy and generalization capabilities of the models. The remainder of the paper is structured as follows: Section II presents a review of previous studies on employee attrition prediction. Section III details the methodology adopted in this study, including dataset development, modeling framework, and evaluation strategy. Section IV reports the experimental results and feature analysis. Section V discusses the findings and their implications. Section VI offers concluding remarks and potential directions for future research.

II. LITERATURE REVIEW

A. Employee Attrition Overview

Employee attrition is defined as the departure of employees from an organization for various reasons, including resignation, retirement, or death [17]. Attrition is also described as a reduction in the workforce resulting from these departures and is sometimes referred to as total turnover or wastage [18], [19]. While attrition reflects a decline in the number of employees, employee turnover specifically denotes the replacement of departing employees with new hires. These terms are often used interchangeably in workforce analyses due to their conceptual overlap. Turnover can be classified as voluntary or involuntary [20]. Even though dysfunctional turnover, characterized by the loss of high-performing employees, is detrimental to an organization, the turnover of employees with easily replaceable skills tends to have less impact [21].

In fact, employee attrition is a critical issue impacting various sectors, particularly those where employee roles are essential, creating barriers to profitability, competitiveness, and productivity by causing overtime, delayed project delivery, customer dissatisfaction, and reduced morale among remaining staff who must fill vacancies [19]. Apart from this, high attrition rates contribute to increased employee turnover, resulting in substantial expenses related to recruitment, training, and development of new hires [20], [22]. Therefore, retaining employees is financially crucial as it directly influences revenue growth through customer acquisition and retention

[23]. Furthermore, voluntary turnover results in workforce movement across companies and roles, transferring valuable knowledge to competitors, and consequently posing a major risk to organizational interests [24].

B. Factors Influencing Employee Attrition

While human resources represent one of the most valuable assets contributing to a company's success, one of the most prominent challenges faced by HR management is employee attrition [22]. The initial step towards addressing this issue involves identifying the underlying causes and examining the factors associated with attrition. In this sub-section, several studies that have identified or investigated these factors are discussed.

For instance, the study in [17] reportes that employee attrition occurs for various reasons, including insufficient support for employees in lower organizational positions, lack of direct communication lines between these employees and their supervisors, monotonous tasks leading to a loss of creativity in repetitive jobs, limited career advancement prospects, and a lack of opportunities. As well as that, authors in [25] confirm that employees tend to leave their jobs when they perceive inadequate compensation, experience disagreements with their managers, or confront issues related to work hours, career development, and family obligations.

Moreover, a study conducted on 40 private financial institutions—including banks, investment fund organizations, and insurance companies—in the Kottayam district of India concluded that a strained relationship between management and employees, employee dissatisfaction, lack of job security, and job-related pressure significantly contribute to employee attrition [19]. Additionally, extensive research highlights a range of factors influencing turnover rates, which can broadly be categorized into organizational factors, job-related factors, work environment factors, and personal factors [26].

C. Employee Attrition and Retention in Saudi Private Sector

In this subsection, studies investigating the factors that impact employee turnover in the Saudi private sector are presented to highlight the significance of this issue and the necessity for continued research. The following studies explore some organizational and job-related factors closely related to employee attrition, such as job satisfaction, work participation, work pressures, rewards, effective training, and recognition.

For example, the study in [27] examined the impact of training and development opportunities following Saudization on employees' intentions to leave their organizations. The findings indicated that providing employee training contributed to retention exceeding five years, with training and development opportunities facilitating promotions. This was identified as a critical factor in reducing turnover rates. Besides that, a study conducted in [28] aimed to identify the causes of employee turnover and its effects on the Saudi business environment. The study focused on factors such as job satisfaction, organizational commitment, support, feedback, rewards, and recognition. The results revealed that the primary factors influencing employees' decisions to leave included organizational support, rewards, appreciation, and commitment, whereas job satisfaction and

supervisor support had a comparatively less impact on turnover intention.

In contrast, the authors in [29] conducted a study to identify the employer retention methods and the factors affecting retention from the employee's perspective. The study was carried out in private businesses in three Saudi Arabian cities: Jeddah, Rabigh, and Yanbu. The findings indicated that increased rewards and privileges are key retention strategies for businesses aiming to maintain their staff long-term. Another study [30] demonstrated a negative relationship between employee experience and turnover intention, suggesting that as employees gain more experience, their likelihood of leaving decreases. Additionally, age plays a significant role in turnover motivation, with older workers exhibiting a lower propensity to leave compared to younger employees.

D. Overview of Predictive Models for Employee Attrition

Predictive ML involves training algorithms on historical data to generate predictions for new, unseen data. By leveraging various factors—such as job performance, length of service, and other relevant indicators— HR managers can evaluate an employee's intention to leave their position. Consequently, ML models have been employed extensively for the prediction of employee attrition. Inspired by previous studies, combinations of several of these models and techniques have been applied in this study to the newly created dataset, and the results have been analyzed in order to obtain robust predictions.

For instance, this study [13] suggested that people analytics be used to predict employee attrition, with a focus on data quality over quantity. The study employed three datasets: the IBM HR Analytics Employee Attrition dataset, a simulated HR dataset based on Kaggle data, and a survey-based dataset on the reasons for attrition, and applied a variety of algorithms, including deep learning methods, traditional ML models: DT, SVM, LR, RF, and ensemble techniques like XGBoost, VC, and stacking ANNs. The results indicated that the VC outperformed other models across all datasets. In the same way, the study in [31] evaluated other algorithms including KNN, Naive Bayes,RF, and LR on IBM HR Analytics Employee Attrition dataset, reporting the highest accuracy of 89% for Naive Bayes.

In addition, the study in [32] used various ML and ensemble learning methods to predict employee attrition, comparing four ML techniques and three ensemble methods, implementing grid search for hyperparameter tuning. The results reveal that the RF outperformed other models with an accuracy of 95%, identifying job satisfaction and stock options as primary attrition drivers. Similarly, Garg et al in [33] trained RF, AdaBoost, XGBoost, and ensemble stacking models on the IBM dataset, reporting the RF classifier as the top performing with 87.41% accuracy and 88% precision, and indicating that monthly pay significantly influences employee turnover.

In study [34], the IBM HR Analytics Employee Attrition and Performance dataset was used to train multiple models, including LR, RF, SVM, KNN, and DT. The LR model achieved the highest accuracy, with 87% precision and an F1 score of 83%. The findings identified monthly income, overtime, years with the current manager, age, distance from home, total working years, and years at the company as key

attrition indicators. As well as previous research, [35] has used the IBM dataset available to evaluate LR, DT, RF, GBM, XGBoost, SVM, and KNN for attrition prediction. As a result, RF and GBM demonstrated superior predictive performance for employee retention. However, the SVM model achieved an accuracy of 88.78% and maintained a balanced performance with respect to accuracy and recall for the minority class (employees who departed), highlighting its usefulness in attrition prediction.

E. Research Gap and Study Positioning

Despite the growing body of research on employee attrition prediction using ML, several critical gaps remain that limit the practical applicability and contextual relevance of existing studies.

- 1) First: The majority of prior research has relied heavily on synthetic datasets, particularly the IBM HR Analytics Employee Attrition dataset [31], [33], [34], [35], or simulated data from platforms like Kaggle [13], [32]. While these datasets provide standardized benchmarks for model comparison, they fail to capture the complexity and nuances of real-world organizational dynamics, cultural contexts, and region-specific employment practices. This dependency on generic datasets raises questions about the generalizability and practical utility of the developed models in diverse organizational settings.
- 2) Second: There is a notable scarcity of research examining employee attrition within the Saudi Arabian context, despite the Kingdom's unique labor market characteristics shaped by Saudization policies, cultural factors, and rapid economic transformation [27], [28], [29]. The limited studies that do focus on Saudi Arabia primarily employ traditional statistical methods or exploratory analyses rather than comprehensive machine learning frameworks, leaving a gap in predictive analytics tailored to this specific workforce.
- 3) Third: Existing attrition prediction studies often utilize incomplete or inconsistent feature sets, omitting important variables identified in recent retention literature. Variables such as job opportunities, emotional commitment, health issues, job stability, and sector-specific factors are frequently overlooked in predictive models [27], [28], [31], despite evidence of their significance in employee turnover decisions.
- 4) Fourth: While ensemble methods have shown promise in attrition prediction [31], [32], [33], systematic comparisons across multiple feature selection techniques and their impact on model performance remain limited. Most studies employ a single feature selection approach without exploring how different methods influence predictive accuracy and model interpretability, particularly in the context of imbalanced datasets typical of attrition scenarios.
- 5) Finally: There is insufficient emphasis on deriving actionable insights for HR management from predictive models. Many studies focus primarily on achieving high prediction accuracy without adequately analyzing which factors most strongly influence attrition or how organizations can leverage these insights to develop targeted retention strategies [7], [14].

This study addresses these gaps by: 1) introducing a novel, empirically-collected dataset from real employees in the Saudi private sector that reflects authentic organizational contexts and cultural dynamics; 2) incorporating a comprehensive set of attrition factors derived from both international literature and Saudi-specific studies, including previously underexplored variables, making this dataset valuable for addressing a range of factors not available in other attrition datasets, such as health issues, benefits, medical insurance, job opportunities, and others; 3) systematically comparing multiple machine learning algorithms and ensemble methods across various feature selection strategies to identify optimal prediction configurations; 4) providing detailed feature importance analysis to generate actionable insights for HR practitioners; and 5) demonstrating the practical applicability of predictive analytics in addressing attrition challenges within the evolving Saudi labor market. By filling these gaps, this research contributes both methodologically and contextually to the field of employee attrition prediction, offering a robust framework adaptable to other regional and organizational contexts.

III. METHODOLOGY

The methodology is structured to address three key objectives: 1) collecting authentic employee data that captures Saudi-specific factors, 2) preparing and balancing the dataset to mitigate prediction bias, and 3) systematically evaluating multiple ML models to identify the optimal prediction approach. Each phase builds upon the previous one to ensure robust and reliable attrition predictions. The data collection phase establishes the empirical foundation through a targeted survey of 1,191 Saudi private sector employees. The preprocessing phase addresses data quality through outlier removal, categorical encoding, and SMOTE-based class balancing. The modeling phase implements a comprehensive evaluation framework comparing eight algorithms across seven feature selection strategies, with performance assessed using accuracy and F1 score metrics optimized for imbalanced classification tasks. Fig. 1 shows the basic phases of the methodology followed in this study.

A. Data Collection

Employee data related to the attrition factors examined in this study were collected through an online questionnaire, specifically designed and implemented using the Google Forms platform. This survey specifically targeted real employees working within the Saudi private sector, thereby addressing the limitations of relying on synthetic datasets often used in prior research. In total, 1200 responses were gathered. The responses were thoroughly filtered to exclude individuals who did not belong to the Saudi private sector, resulting in a final sample of 1191 participants. Additionally, all respondents were employees from a wide range of organizations spread across the various regions of the Kingdom of Saudi Arabia, thereby ensuring diverse geographic representation within the dataset.

Moreover, the questionnaire was conducted anonymously to maintain participant confidentiality. With regard to demographics, male respondents constituted 43%, while female respondents made up 57% of the total sample. Additionally, 515 respondents indicated that they had left their jobs, providing direct insight into employee attrition.

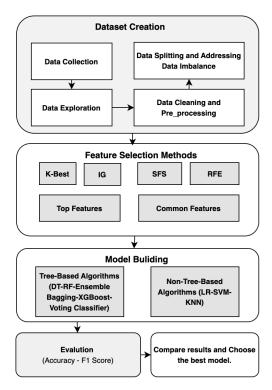


Fig. 1. Methodology phases.

TABLE I. DATASET FEATURES

Category	Features		
Personal Features	Gender, Age, Marital Status, Academic Degree, Health Issues, Years of Experience, Years of Ex- perience in last organization, Job Opportunities.		
Job-related Features	Sector, Department, Job Title, Monthly Salary, Allowances, Medical Insurance, Bonus, Overtime, Payment Overtime, Rewards and Wages Satisfaction, Get Deserved Promotion, Training Programs During the Last Three Years, Useful Training Programs, Business Travel, Job Support, Recognition, Emotional Commitment, Job Engagement, Distance to Work.		
Psychological Features	Physical Stress, Psychological Exhaustion, Job Stability, Environment Satisfaction, Job Satisfac- tion, Work-Life Balance.		
Target Feature	Attrition		

B. Data Exploration

The dataset used in this study consists of 34 attributes, which are categorized into three primary groups: personality attributes, job-related attributes, and psychological and satisfaction attributes. First, personality attributes include individual characteristics such as gender, age, and social status. Second, job-related attributes encompass factors related to the work itself, such as job title, salary, sector, and department. Finally, psychological and satisfaction attributes relate to how employees perceive their work environment, including job stability, environmental satisfaction, and overall job satisfaction. The dataset also features the target variable of attrition, as detailed in Table I.

The variables included in this study were primarily selected based on insights from previous research on employee attrition prediction. In addition to these established factors, the dataset integrates a set of novel variables that have not been widely applied in attrition prediction models. These novel variables originate from studies specifically focusing on the Saudi private sector, as well as recommended attrition predictors in the literature. For instance, key examples of these additional variables include job opportunities, allowances, job support, medical insurance, deserved promotion [27], [28], emotional commitment, health issues, and job stability [13]. Correspondingly, the sector variable was included to reflect the complex and diverse composition of the Saudi private sector, which encompasses a wide range of industries and organizational forms.

In this study, data exploration involved visualizing the distributions of key features, identifying correlations among variables, and detecting potential outliers or anomalies. To better understand the relationships between various factors and employee attrition, a thorough exploratory analysis was performed. The detailed distributions of the main features related to attrition are presented and further evaluated in the discussion section.

C. Data Cleaning and Pre-processing

The dataset contains both numerical and categorical features. Numerical variables included features such as employee allowances, whereas the remaining attributes were categorical in nature. In this section, the steps taken to preprocess the dataset are outlined to build ML models for employee attrition prediction.

- 1) Handling missing values: The dataset was examined for the presence of missing entries using the <code>isnull().sum()</code> function in Pandas. As no missing values were detected, neither imputation procedures nor the elimination of incomplete records was required.
- 2) Outlier detection and removal: To mitigate the influence of anomalous data points on model training, numerical features were evaluated for outliers using the Z-score technique. Instances with an absolute Z-score greater than three were classified as outliers and excluded from the dataset. This step ensured that extreme values did not distort the learning process of the predictive models.
- 3) Addressing categorical features categorical attributes processed using a two-stage approach:
- a) Standardization: Preliminary exploration of the dataset revealed inconsistencies in the representation of categorical variables due to irregular spacing and capitalization. These issues were rectified by removing extraneous spaces and enforcing uniform formatting across feature values. This standardization process was applied to attributes such as Business_Travel, Years_Experience, and MonthlySalary.
- b) Grouping rare categories: Several categorical variables, including MonthlySalary, Sector, Years_experience_lastorganization, Years_Experience, Age, and Department exhibited long-tailed distributions characterized by a large number of infrequent categories. To reduce model complexity and mitigate the risk of overfitting, these rare categories were consolidated into broader, more meaningful groups.

4) Encoding categorical variables: As a preliminary step, the LabelEncoder and OrdinalEncoder were employed to encode nominal and ordinal features, respectively. Due to the overfitting of tree-based models, in particular, the dataset was prepared in two distinct versions, with the key distinction being the encoding of categorical variables. In both dataset versions, ordinal features were encoded using Ordinal Encoding to preserve the intrinsic order among categories. Similarly, Label Encoding was consistently applied to binary features across both versions to ensure uniform representation.

For the tree-based dataset, nominal features with multiple categories, including Maritalstatus, Sector, Department, and JobTitle were encoded using the M-Estimator method. This encoder assigns numerical representations to categories, considering both the global distribution of the target variable and the distribution within each category [36]. By providing smooth and target-aware encodings, M-Estimator encoding enables tree-based models to create more effective and informative splits.

On the other hand, for the same nominal features, one-hot encoding was applied in the dataset version for non-tree-based models. This encoder generates binary dummy variables for each category, where the presence of a category is represented as 1 and its absence as 0, providing an orthogonal and non-ordinal representation suitable for these algorithms [37].

5) Scaling numerical features: Numerical and encoded features were normalized to ensure that the models could converge more quickly and avoid biasing the models towards features with larger ranges.

D. Data Splitting and Addressing Data Imbalance

The target variable, Attrition, exhibited significant class imbalance, with a notably higher proportion of instances labeled as "No" (employees who remained) compared to "Yes" (employees who left), as illustrated in Fig. 2. The dataset was partitioned into training and testing subsets, with an 80 %-20% split respectively, while maintaining the class distribution within each subset as detailed in Table II.

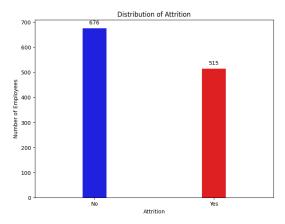


Fig. 2. Attrition distribution.

To address the class imbalance in the training dataset, the Synthetic Minority Oversampling Technique (SMOTE) was applied exclusively to the training set, targeting the minority

TABLE II. DATASET SPLIT AND ATTRITION DISTRIBUTION

#	Total Records	Non-Attrition	Attrition
Total records	1174	676	498
Training set 80%	939	534	405
SMOTE set	1068	534	534
Testing set 20%	235	142	93

class labeled "Yes". SMOTE generates synthetic samples by interpolating between existing minority class instances, thereby balancing the class distribution. This approach mitigates bias towards the majority class, enhances model performance on the minority class, and exposes the model to a more diverse and balanced set of examples, improving its generalization capability.

E. Feature Selection Methods

Feature selection was performed to identify the most informative features for predicting employee attrition in both tree-based and non-tree-based datasets. Four methods—SelectKBest, Mutual Information Gain (IG), Recursive Feature Elimination (RFE), and Sequential Feature Selection (SFS)—were applied, testing subsets of the top 10, 20, 25, and 30 features. The model performance was best with the top 25 features, consistent with prior studies in [38]. Therefore, this subset was selected for model training and evaluation.

F. Models Training and Validation

The modeling process aimed to predict employee attrition using a combination of tree-based and non-tree-based ML algorithms. This subsection details the models employed, evaluation metrics, and the overall procedures for training and testing

- 1) Feature sets: Seven feature sets were used to evaluate each model, reflecting various feature selection strategies. Each of these feature sets was evaluated to determine the optimal set of features for the models. These sets included: 1) all features without selection, 2) features from the SelectKBest method, 3) Information Gain-selected features, 4) Recursive Feature Elimination (RFE) features, 5) Sequential Feature Selection (SFS) features, 6) the union of features selected by all methods (common features), and 7) the intersection of features selected by any method (top features). This comprehensive evaluation allowed assessing model sensitivity to various feature subsets and selecting the most predictive features.
- 2) Model building: In this section, we describe the ML models used for predicting employee attrition, including both non-tree-based and tree-based algorithms. The non-tree-based algorithms employed in this study are listed in Table III, which includes LR, SVM, and KNN. These models are well-suited for datasets with one-hot encoded categorical features. Additionally, tree-based algorithms, as presented in Table IV, were utilized to leverage their ability to handle complex feature interactions and non-linear relationships. The tree-based models include DT, RF, Bagging ensemble, XGBoost, and a Voting Classifier combining Random Forest and XGBoost.

TABLE III. NON-TREE-BASED ALGORITHMS

No.	Model
1	Logistic Regression
2	Support Vector Machine (SVM)
3	K-Nearest Neighbors (KNN)

TABLE IV. TREE-BASED ALGORITHMS

No.	Model
1	Decision Tree
2	Random Forest
3	Bagging
4	XGBoost (Extreme Gradient Boosting)
5	Voting Classifier (Random Forest + XGBoost)

3) Evaluation metrics: Accuracy measures the proportion of correctly classified instances out of the total instances. It is calculated as:

$$Accuracy = \frac{True\ Positives + True\ Negatives}{Total\ Instances}$$

Accuracy provides a simple and intuitive measure of the model's overall correctness. However, it can be misleading when dealing with imbalanced datasets.

The F1 Score is defined as the harmonic mean of precision and recall, providing a quantitative measure that balances these two metrics. This metric is formally expressed as:

F1 Score =
$$2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

The F1 score is especially advantageous when it is necessary to account for both false positives and false negatives. A higher F1 score indicates robust model performance across both dimensions. Notably, the F1 score often provides a more informative evaluation than accuracy in scenarios characterized by class imbalance, due to reflecting the model's proficiency with respect to the minority.

IV. RESULTS

A. AL Models with Different Features Selection Methods for the Top 25 Features

The performance of various ML models using different feature selection strategies was evaluated. The models evaluated include LR, DT, RF, XGBoost, Ensemble Bagging, SVM, KNN, and VC. Each model was assessed using seven feature selection strategies: All Features, SelectKBest, IG, RFE, SFS, Common Feature, and Top Feature. These strategies utilized the top 25 ranked features based on their respective scores. After training under all feature selection scenarios, results were compared to identify the optimal configuration per model. For each model, only the best result and corresponding feature set are reported as shown in Table V, which summarizes the best performance of each model with feature sets. Additionally, the highest accuracy observed in this study was 90%, and the highest F1 score was 85%, both achieved by the Voting Classifier leveraging the feature set selected by RFE.

TABLE V. MODELS PERFORMANCE

Model	Fatures selection	Accuracy	F1 Score
LR	All	77.00	74.00
SVM	RFE	80.34	71.00
K-NN	Common	78.63	75.25
DT	Common	84.62	76.32
RF	SFS	89.00	84.71
Ensemble Bagging	SFS	87.18	82.35
XGBoost	All	83.76	75.95
Voting Classifier	RFE	90.00	85.00

B. Comparative Analysis with Existing Methods

To validate the effectiveness of the proposed approach and assess its importance relative to existing methods, we compare our best-performing model (Voting Classifier with RFE features) against results reported in recent studies. Table VI presents this comparison based on the methodologies and results explicitly reported in the referenced studies.

The comparative analysis demonstrates that the proposed method achieves competitive performance, with the Voting Classifier achieving 90% accuracy and 85% F1 score. This performance is comparable to the best results reported in existing literature, with several important distinctions.

First, unlike prior studies that predominantly utilize the IBM HR Analytics Employee Attrition dataset [31], [33], [34], [35], our approach employs a novel, empirically-collected dataset from the Saudi private sector, which introduces additional complexity due to real-world data characteristics and regional-specific factors. Second, our method achieves this performance while addressing class imbalance through SMOTE and systematic feature selection across multiple methods (SelectKBest, IG, RFE, SFS), providing a more comprehensive evaluation framework than many existing approaches. Third, the high F1 score (85%) indicates strong performance on both precision and recall, which is particularly important for imbalanced attrition datasets where the minority class (employees who leave) is of primary interest.

The study by Yahia et al. [13] also reported that Voting Classifier outperformed other models across multiple datasets, which aligns with our findings and validates the effectiveness of ensemble methods for attrition prediction. However, our study extends this by demonstrating the applicability of these methods to region-specific, real-world data and by providing detailed feature importance analysis that reveals job opportunities and job titles as primary attrition drivers in the Saudi context.

The effectiveness of hybrid ensemble approaches is further demonstrated by Govindarajan et al. [39], whose Hybrid Model integrating multiple algorithms achieved 95% accuracy on the IBM HR dataset. Their comparative analysis revealed that individual models achieved lower performance (SVM: 88.6%, Random Forest and Gradient Boosting: 87.3%, Decision Tree: 80.5%), while the hybrid approach effectively harnessed complementary strengths while mitigating individual weaknesses. This finding directly supports our methodology of employing a Voting Classifier that combines Random Forest and XG-Boost, demonstrating that ensemble integration consistently outperforms individual algorithms across different studies and

datasets. The parallel between their Hybrid Model (95%) and our Voting Classifier (90%) performance, despite our use of more challenging real-world data, validates the robustness of ensemble-based approaches for attrition prediction.

V. DISCUSSION

This study addresses the pressing issue of employee attrition within the Saudi private sector, with a focus on developing predictive models using ML techniques. The research aims to determine how well ML models can predict employee attrition, and to identify the feature configurations that most effectively enhance prediction accuracy. The subsequent discussion synthesizes the key findings in relation to the research questions, objectives, and methodology introduced earlier.

A. Discussion Related to Dataset

This study utilized a new, real-world dataset that stands out for its relevance, depth, and contextual specificity to the Saudi Arabian private sector workforce. Unlike many prior works relying on synthetic or general-purpose datasets such as IBM's attrition data [31], [33], [34], and [35], this research provides a more realistic and localized perspective. The dataset comprises responses from 1,191 participants gathered through an online questionnaire targeting Saudi private-sector employees across various industries. The dataset's contextual specificity enables the extraction of actionable insights tailored to the distinct work culture and employment practices prevalent in Saudi Arabia, which often diverge from Western contexts. A significant insight from this dataset is the class imbalance in the target variable: approximately 43% of participants had left a job, while 57% had not. To address this imbalance, the Synthetic Minority Over-sampling Technique (SMOTE) was applied during training, which is critical for mitigating prediction bias in high-stakes human resource analytics. However, the geographic and cultural specificity of the dataset may constrain the generalizability of the findings to other labor markets.

B. Discussion Related to Employee Attrition Factors

A variety of feature selection techniques were employed in this study. Due to the multiplicity and potential conflict among selected feature groups, a general, model-agnostic approach—binary logistic regression—was utilized to assess the predictive significance of variables. Logistic regression offers interpretable correlation measures in the form of odds ratios, which provide straightforward insights compared to the often complex interpretability of ML outputs [40]. In addition, a subset of common variables was identified as the intersection of all applied feature sets, representing features consistently deemed significant across methods.

The results of the logistic regression analysis revealed that the most significant factors influencing employee attrition were job opportunities, job title, deserved promotion, training programs, emotional commitment, job engagement, benefits, and medical insurance. The common set of attributes selected from all attribute selection methods included the following variables: years of experience, job opportunities, sector, job engagement, job support, years of experience in the previous organization, job stability, training programs, appreciation, emotional commitment, health issues, and deserved promotion.

Study	Dataset	Best Model	Accuracy	F1 Score
Rana et al. [31]	IBM HR Dataset	Naive Bayes	89%	-
Chung et al. [32]	IBM HR Dataset	Random Forest	95%	-
Garg et al. [33]	IBM Dataset	Random Forest	87.41%	88% (Precision)
Solomon et al. [34]	IBM HR Dataset	Logistic Regression	87%	83%
Amzad et al. [35]	IBM Dataset	SVM	88.78%	-
Yahia et al. [13]	Large, medium simulated HR datasets, small (450 real employees)	Deep and ensemble learning	96%, 98%, 99% respectively	-
Proposed Method	Saudi Private Sector (1,191 real employees)	Voting Classifier (RFE)	90%	85%

TABLE VI. COMPARISON WITH EXISTING EMPLOYEE ATTRITION PREDICTION METHODS

The results align with prior studies investigating employee attrition determinants in the Saudi private sector. Study [27] highlighted the significant influence of job opportunities, training programs, and promotion on attrition rates. In the same way, the study in [13] underscored the impact of job title and job engagement, emphasizing additional factors such as job stability, health conditions, and emotional commitment. Moreover, study [28] affirms the role of job support and recognition in reducing attrition. Furthermore, study [29] identified job support, recognition, financial allowances, and bonuses as critical factors in employee retention consistent with. The importance of years of experience in predicting employee attrition was also highlighted in [34]. On the other hand, several variables were not statistically significant in predicting workforce attrition, including age, academic degree, monthly salary, overtime hours, bonuses, job satisfaction, environmental satisfaction, business travel, and distance to work, which contradicts previous studies [13], [34]. These contradictions reflect contextual differences in the Saudi private sector. The non-significance of salary and bonuses suggests competitive compensation is a baseline expectation rather than a differentiator, shifting focus to career growth. Age's reduced importance aligns with Saudi Arabia's young workforce in Vision 2030 transformation [30], where career prospects uniformly influence turnover. Job satisfaction's non-significance indicates employees may be satisfied yet leave for better opportunities, suggesting satisfaction alone is insufficient without development prospects. The non-significance of business travel contradicts Yahia et al. [13], likely reflecting cultural and worklife balance differences in the Saudi context. These findings emphasize that retention strategies from Western datasets require adaptation to Middle Eastern labor markets.

C. Discussion Related to ML Models

1) ML Models: Based on the study results, the voting classifier with the feature set selected by Recursive Feature Elimination (RFE) achieved the best performance, with 90% accuracy and an 85% F1 score. This finding aligns with the results reported in study [13]. Ensemble learning methods, such as voting classifiers, combine multiple weak learners to leverage their complementary strengths, resulting in more accurate predictions. Combining outputs from multiple models is a well-established technique to improve overall model performance [41]. Apart from this, RF achieved the second-best performance, consistent with its superiority in previous studies [32],[33], and [35]. Conversely, the LR model did not outperform the other models, contrasting with findings reported in [38].

2) Impact of Feature Sets: The "All Features" set includes all available features for comparison, showing good overall

performance. Specifically, the LR achieved the best result with all features with 77% accuracy and 74% F1 score. The evaluation of "Common Features" across models revealed varying impact: while the accuracy of VC and RF improved, the performance of LR, SVM, and Bagging declined significantly. Besides that, the "Top Features" method consistently improved model performance, although it was less effective compared to SFS and RFE.

In contrast, most models demonstrated improved performance with feature selection methods. Models such as SVM, RF, and VC benefited from feature selection by reducing noise and focusing on only relevant variables. SFS emerged as the most effective feature selection method for most models. For instance, with ensemble bagging, SFS provided the highest accuracy at 87.18%. SFS emerged as the most effective technique overall, providing the highest accuracy of 87.18% for the ensemble bagging model. Additionally, RFE ranked as the second-best feature selection method after SFS, improving model performance, particularly for RF, Bagging, and VC models.

While SelectKBest and IG effectively identified important features, their impact on improving model performance was less significant compared to SFS and RFE. For example, Logistic Regression (LR) achieved an F1 score of 73.58% using KBest, which increased to 79.31% with SFS. Similarly, Random Forest (RF) obtained an F1 score of 73.58% with IG, which improved to 84.71% using SFS and reached 80% with RFE as shown in detail in Fig. 3, which presents a comparison of accuracy across various models and feature selection methods.

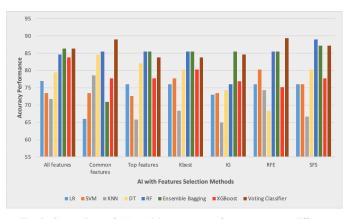


Fig. 3. Comparison of AI model accuracy performance across different feature sets and feature selection methods.

VI. CONCLUSION

Employee attrition remains a critical challenge for the private sector, particularly in monitoring employee intentions and proactively addressing turnover. This study aimed to analyze key factors predicting employee departure and the underlying reasons for turnover, thereby enabling organizations to implement effective retention strategies. The primary objective of this study was to develop and evaluate a predictive model using ML and ensemble methods to address attrition challenges tailored to the Saudi private sector workforce.

The study has four key contributions. First, it introduces a new dataset derived from previously identified causes of employee attrition. Next, it presents the design, development, and evaluation of predictive models using ML and ensemble methods. Additionally, the study compares the performance of various ML models across different feature configurations. Lastly, it provides a thorough analysis of factors affecting turnover intentions among employees in the Saudi private sector, contributing valuable insights to inform retention policy development.

The study results indicate the superiority of ensemble models over individual algorithms in predicting employee attrition. Furthermore, demographic and overall satisfaction variables showed limited predictive power within this dataset, while dynamic and developmental factors had a more pronounced influence on attrition. These findings reveal an important shift in employee attrition patterns: traditional demographicbased retention strategies may be insufficient in contemporary organizational contexts. Instead, the prominence of dynamic factors—such as job opportunities, career development prospects, and job titles—suggests that employees in the Saudi private sector prioritize growth-oriented and opportunity-driven workplace environments. This insight challenges conventional HR assumptions that focus primarily on compensation and job satisfaction, highlighting the need for organizations to invest in career development programs, clear advancement pathways, and competitive positioning within the job market. in addition to, need to adopt proactive strategies focused on development and career advancement. The proposed model can help identify at-risk employees early and support data-driven retention decisions in the Saudi private sector.

From a strategic perspective, the high predictive accuracy of ensemble models (90% with Voting Classifier) demonstrates the viability of data-driven early warning systems for attrition risk. Organizations can leverage such models to identify atrisk employees proactively, enabling targeted interventions before turnover intentions crystallize into actual departures. This predictive capability is particularly valuable for retaining high-performing employees and reducing the substantial costs associated with talent replacement and organizational knowledge loss.

However, the study faced limitations, including potential bias inherent to questionnaire-based data collection, while relying on actual employee records in private sector companies may provide more reliable data. Additionally, the focus on the Saudi private sector as a whole suggests that sector-specific studies—for example, healthcare or education—could yield more precise insights into attrition drivers.

Future research should pursue several promising direc-

tions to advance employee attrition prediction and understanding. First, addressing the challenge of imbalanced data typical of turnover datasets through advanced techniques such as SMOTE (Synthetic Minority Over-sampling Technique), ADASYN (Adaptive Synthetic Sampling), or cost-sensitive learning algorithms could improve model performance and generalizability. Second, exploring deep learning architectures, including recurrent neural networks (RNNs) and long shortterm memory (LSTM) networks, may capture temporal patterns in employee behavior and career trajectories, potentially revealing early warning signals that precede turnover decisions. Third, incorporating explainable AI (XAI) techniques such as SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-agnostic Explanations) would enhance model interpretability, enabling HR practitioners to understand not only predictions but also the reasoning behind them, thereby facilitating more informed decision-making. Fourth, conducting longitudinal studies that track employees over extended periods could reveal the dynamic evolution of attrition risk factors and validate the temporal stability of predictive models. Finally, sector-specific investigations within healthcare, education, technology, and finance sectors would uncover industry-unique attrition drivers and enable the development of tailored retention strategies.

Ultimately, this study offers a foundational step toward enhancing employee retention strategies within Saudi Arabia's evolving private sector landscape, demonstrating that the integration of machine learning techniques with empirical workforce data can transform HR practices from reactive problem-solving to proactive talent management.

REFERENCES

- [1] K. M. Mitravinda and S. Shetty, "Employee attrition: Prediction, analysis of contributory factors and recommendations for employee retention," in 2022 IEEE Int. Conf. Women Innovation, Technol. Entrepreneurship (ICWITE), Dec. 2022, pp. 1–6.
- [2] F. Mozaffari, M. Rahimi, H. Yazdani, and B. Sohrabi, "Employee attrition prediction in a pharmaceutical company using both machine learning approach and qualitative data," *Benchmarking: An Int. J.*, 2022, ahead-of-print.
- [3] R. Joseph, S. Udupa, S. Jangale, K. Kotkar, and P. Pawar, "Employee attrition using machine learning and depression analysis," in 2021 5th Int. Conf. Intell. Comput. Control Syst. (ICICCS). IEEE, 2021, pp. 1000–1005.
- [4] S. Gupta, G. Bhardwaj, M. Arora, R. Rani, P. Bansal, and R. Kumar, "Employee attrition prediction in industries using machine learning algorithms," in 2023 10th Int. Conf. Comput. Sustainable Global Develop. (INDIACom). IEEE, 2023, pp. 945–950.
- [5] G. R. Rajeswari, R. Murugesan, R. Aruna, B. Jayakrishnan, and K. Nilavathy, "Predicting employee attrition through machine learning," in 2022 3rd Int. Conf. Smart Electron. Commun. (ICOSEC). IEEE, 2022, pp. 1022–1027.
- [6] E. Rombaut and M.-A. Guerry, "Predicting voluntary turnover through human resources database analysis," *Manage. Res. Rev.*, vol. 41, no. 1, pp. 96–112, 2018.
- [7] N. B. Yahia, J. Hlel, and R. Colomo-Palacios, "From big data to deep data to support people analytics for employee attrition prediction," *IEEE Access*, vol. 9, pp. 60447–60458, 2021.
- [8] A. I. Al-Alawi and Y. A. Ghanem, "Predicting employee attrition using machine learning: A systematic literature review," in 2024 ASU Int. Conf. Emerging Technol. Sustainability Intell. Syst. (ICETSIS). IEEE, Jan. 2024.
- [9] A. Tursunbayeva, S. D. Lauro, and C. Pagliari, "People analytics—A scoping review of conceptual boundaries and value propositions," *Int. J. Inf. Manage.*, vol. 43, pp. 224–247, 2018.

- [10] D. Angrave, A. Charlwood, I. Kirkpatrick, M. Lawrence, and M. Stuart, "HR and analytics: Why HR is set to fail the big data challenge," *Human Resource Manage*. J., vol. 26, no. 1, pp. 1–11, 2016.
- [11] T. Pape, "Prioritising data items for business analytics: Framework and application to human resources," *Eur. J. Oper. Res.*, vol. 252, no. 2, pp. 687–698, 2016.
- [12] P. K. Jain, M. Jain, and R. Pamula, "Explaining and predicting employees' attrition: A machine learning approach," SN Appl. Sci., vol. 2, pp. 1–11, 2020.
- [13] N. B. Yahia, J. Hlel, and R. Colomo-Palacios, "From big data to deep data to support people analytics for employee attrition prediction," *IEEE Access*, vol. 9, pp. 60 447–60 458, 2021.
- [14] T. Peeters, J. Paauwe, and K. V. D. Voorde, "People analytics effectiveness: Developing a framework," *J. Organizational Effectiveness: People Performance*, vol. 7, no. 2, pp. 203–219, 2020.
- [15] H. Alqahtani, H. Almagrabi, and A. Alharbi, "Employee attrition prediction using machine learning models: A review paper," *Int. J. Artif. Intell. Appl.*, vol. 15, no. 2, 2024.
- [16] A. Raza, K. Munir, M. Almutairi, F. Younas, and M. M. S. Fareed, "Predicting employee attrition using machine learning approaches," *Appl. Sci.*, vol. 12, no. 13, p. 6424, 2022.
- [17] A. Mhatre, A. Mahalingam, M. Narayanan, A. Nair, and S. Jaju, "Predicting employee attrition along with identifying high risk employees using big data and machine learning," in 2020 2nd Int. Conf. Adv. Comput., Commun. Control Netw. (ICACCCN). IEEE, 2020, pp. 269–276.
- [18] S. Dutta, S. K. Bandyopadhyay, and S. K. Bandyopadhyay, "Employee attrition prediction using neural network cross validation method," *Int. J. Commerce Manage. Res.*, vol. 6, no. 3, pp. 80–85, 2020.
- [19] K. Ashokkumar, S. Jacob, and A. Joseph, "A study on attrition management in private sector financial institutions—A special reference to Kottayam District in Kerala," 2020, technical report.
- [20] D. S. Sisodia, S. Vishwakarma, and A. Pujahari, "Evaluation of machine learning models for employee churn prediction," in 2017 Int. Conf. Inventive Comput. Inform. (ICICI). IEEE, 2017, pp. 1016–1020.
- [21] D. Alao and A. B. Adeyemo, "Analyzing employee attrition using decision tree algorithms," *Comput., Inf. Syst., Develop. Inform. Allied Res. J.*, vol. 4, no. 1, pp. 17–28, 2013.
- [22] F. H. Wardhani and K. M. Lhaksmana, "Predicting employee attrition using logistic regression with feature selection," Sinkron: Jurnal dan Penelitian Teknik Inform., vol. 7, no. 4, pp. 2214–2222, 2022.
- [23] J. T. Johnson, R. W. Griffeth, and M. Griffin, "Factors discriminating functional and dysfunctional salesforce turnover," J. Bus. Ind. Marketing, 2000.
- [24] J. Scott, S. Waite, and D. Reede, "Voluntary employee turnover: A literature review and evidence-based, user-centered strategies to improve retention," J. Amer. College Radiol., vol. 18, no. 3, pp. 442–450, 2021.
- [25] R. Basariya and R. R. Ahmed, "A study on attrition—turnover intentions of employees," *Int. J. Civil Eng. Technol.*, vol. 10, no. 1, pp. 2594–2601, 2010
- [26] P. Ghosh, R. Satyawadi, J. P. Joshi, and M. Shadman, "Who stays with you? Factors predicting employees' intention to stay," *Int. J. Organizational Anal.*, vol. 21, no. 3, pp. 288–312, 2013.

- [27] S. Muzaffar and U. Javed, "Training and development opportunities and turnover intentions post saudization," *PalArch's J. Archaeology Egypt/Egyptology*, vol. 18, no. 14, pp. 521–531, 2021.
- [28] W. Bashraheel and M. Abunar, "The factors of employee turnover intention in private sectors organization in Jeddah, Saudi Arabia," PalArch's J. Archaeology Egypt/Egyptology, vol. 18, no. 12, pp. 87– 94, 2021.
- [29] R. H. Al-Gharbi and U. Javed, "Factors impacting employee retention: A case of private companies in Saudi Arabia," *PalArch's J. Archaeology Egypt/Egyptology*, vol. 18, no. 14, pp. 771–778, 2021.
- [30] S. Faisal, M. Naushad, and M. Faridi, "A study on the level and relationship of job embeddedness and turnover intentions among Saudi Arabian working-class," *Manage. Sci. Lett.*, vol. 10, no. 13, pp. 3167– 3172, 2020.
- [31] A. Rana, S. Malik, and M. Chauhan, "Employee attrition prediction using ML techniques," *J. Kufa Math. Comput.*, vol. 11, no. 2, pp. 88– 100, Aug. 2024.
- [32] D. Chung, J. Yun, J. Lee, and Y. Jeon, "Predictive model of employee attrition based on stacking ensemble learning," *Expert Syst. Appl.*, vol. 215, p. 119364, 2023.
- [33] U. Garg, N. Gupta, M. Manchanda, and K. Purohit, "Classification and prediction of employee attrition rate using ML classifiers," in 2024 Int. Conf. Inventive Comput. Technol. (ICICT). IEEE, Apr. 2024, pp. 608– 613.
- [34] C. H. Solomon, D. Mohankumar, and C. Sivanandam, "Employee attrition analysis using ML," in 2024 Second Int. Conf. Emerging Trends Inf. Technol. Eng. (ICETITE). IEEE, Feb. 2024, pp. 1–8.
- [35] M. S. A. Basha, O. Varikunta, A. U. Devi, and S. Raja, "ML in HR analytics: A comparative study on the predictive accuracy of attrition models," in 2024 2nd Int. Conf. Device Intell., Comput. Commun. Technol. (DICCT). IEEE, Mar. 2024, pp. 475–480.
- [36] D. Micci-Barreca, "A preprocessing scheme for high-cardinality categorical attributes in classification and prediction problems," SIGKDD Explor., vol. 3, pp. 27–32, Jul. 2001.
- [37] Various Authors, "Encoding categorical data: Is there yet anything 'hotter' than one-hot encoding?" *arXiv preprint arXiv:2312.16930*, 2023, available: https://ar5iv.org/abs/2312.16930.
- [38] S. F. Sari and K. M. Lhaksmana, "Employee attrition prediction using feature selection with information gain and random forest classification," *J. Comput. Syst. Inform. (JoSYC)*, vol. 3, no. 4, pp. 410–419, 2022
- [39] R. Govindarajan, N. K. Kumar, P. S. Reddy, E. S. Pravallika, B. Dhatri, and G. P. Kumar, "Predicting employee attrition: A comparative analysis of machine learning models using the IBM human resource analytics dataset," *Procedia Comput. Sci.*, vol. [volume], p. [pages], 2025.
- [40] H. S. R. Rajula, G. Verlato, M. Manchia, N. Antonucci, and V. Fanos, "Comparison of conventional statistical methods with ML in medicine: Diagnosis, drug development, and treatment," *Medicina*, vol. 56, no. 9, p. 455, Sep. 2020, pMID: 32911665; PMCID: PMC7560135.
- [41] M. A. Ganaie, M. Hu, A. K. Malik, M. Tanveer, and P. N. Suganthan, "Ensemble deep learning: A review," Eng. Appl. Artif. Intell., vol. 115, p. 105151, 2022.