

# A HYBRID MACHINE LEARNING-BASED PREDICTION APPROACH: THE ACCOUNTANT BEHAVIOR

\*<sup>1</sup>Udeagha E.O., <sup>1</sup>Choji N.D., <sup>2</sup>Mbanaso U.O., <sup>2</sup>Aimufua G.I.O.

<sup>1</sup>Department of Computer Science, University of Jos, Nigeria

<sup>2</sup>Department of Computer Science, Nasarawa State University, Keffi, Nigeria

\*Corresponding Author Email Address: [emmanjab@gmail.com](mailto:emmanjab@gmail.com)

## ABSTRACT

The integration of technology in accounting roles raises questions about the adaptability and skills of accountants in utilizing these tools effectively. Understanding how accountants' behavior is influenced by technology is crucial for their professional development and the accounting industry's future. This study focused on the development of a predictive model, leveraging both Naive Bayes and K-Nearest Neighbors (KNN) models. The research methodology involved the use of Pandas DataFrame to establish a robust framework for the dataset, incorporating both established and innovative features as input variables. These datasets were then utilized as the training data for the predictive model, with the primary objective of extracting valuable insights for decision-making and forecasting accountant behavior. The key findings of the study shed light on the performance of the different models employed. The Naïve Bayes model emerged as a standout performer, achieving an accuracy rate of 63% and an exceptional recall rate of 97%. This underscores its effectiveness in predicting accountant behavior, especially in identifying positive instances. On the other hand, the K-Nearest Neighbors model displayed a balanced trade-off between precision and recall, achieving an accuracy rate of 52% and an F1 score of 64%. This suggests that the model provides a reasonable compromise between accurately identifying positive cases and overall performance. Furthermore, the hybrid KNN-NB model, which amalgamates elements from both approaches, also achieved an accuracy rate of 52%. This finding indicates that the hybrid model has the potential to harness the strengths of both algorithms, offering a versatile approach to predicting accountant behavior.

**Keywords:** Machine Learning, Accountants Behaviour, Pandas Data Frame, Artificial Intelligence, Prediction

## INTRODUCTION

The digitalization of the world continues and new innovations within accounting and finance will affect every day work tasks. Without insights about what is currently happening on the field and what will most likely change in the foreseeable future, various professions will be put at risk. The gap between information technology and the traditional accounting and finance roles is predicted to rapidly diminish (Huttunen.etal.,2019). The understanding of concepts such as Big Data, Accounting Systems, Machine learning, Cloud Computing, and Data Science Applications will be crucial for succeeding in the coming job market as a decision maker in financial sector. The study focuses on the efficient use of data for business decision-making in the financial industry. The Big Data revolution promises to transform how we live, work, and think by enabling process optimization, empowering insight discovery and

improving decision making. The realization of this grand potential relies on the ability to extract value from such massive data through data analytics; machine learning is at its core because of its ability to learn from data and provide data driven insights, decisions, and predictions. However, traditional machine learning approaches were developed in a different era, and thus are based upon multiple assumptions, such as the data set fitting entirely into memory, what unfortunately no longer holds true in this new context (Abdualgalil and Abraham, 2020). These broken assumptions, together with the Big Data characteristics, are creating obstacles for the traditional techniques. Consequently, this thesis compiles, summarizes, and organizes machine learning challenges with Big Data.

The threat that big data analytics will replace many of the tasks traditionally played by accountants may be particularly salient in the auditing context. For example, instead of relying on traditional sampling techniques to perform tests of details, automated processes could examine entire populations for unusual patterns and anomalies. In place of auditors sending out manual confirmations, automatic confirmations could be achieved by a block chain type of technology. A consensus has been emerging among academics that once access to big data analytic techniques becomes ubiquitous in business, users of financial statements will expect audited financial statements on demand, necessitating a shift from traditional sample based auditing to continuous 'auditing by exception', where data analytic techniques direct auditor attention on a real-time basis to "instances where the data does not match the auditor's expectations based on his or her knowledge of the client's business" (Earley et al., 2018).

Furthermore, it is possible that in the future, public accounting firms may face greater competition for the provision of audit services from non-audit firms. As further advances in data analytic and data visualization techniques become available, it will become easier for non-accountants with competencies in data analysis to obtain audit evidence and complete financial statement audits by applying data analytic techniques to big data (Brown-Liburd et al. 2017; Earley 2018). The use of indigenous developed software for computation of Big Data will greatly enhance the work of accountants, and by so doing give both individuals, corporate organizations, governments at different levels and the public the required data for development which will enrich the society.

Accounting and Finance are no exception to this premise, especially when it comes to predicting the outcome of any business decision. Although decision-making always involves future outcomes, however, the effect of any financial decision could be exceptionally phenomenal on the future of any business. Financial decisions involve day-to-day business operations on the one hand and long-term strategy on the other. A business needs to decide about the pricing of the product on the one hand and the pricing of

securities on the other. A business needs to review the business risk as well as the financial risk. All these situations involve financial predictions involving a complex interaction between various datasets. Unfortunately, the existing financial theories are not able to handle such complex decision situations, though they may give some idea. (Heaton et al., 2018). According to (Dixon & Halperin 2019) claim that even though machine learning has been prevalent in financial services for over four decades now, only in the past few years, its impact has been felt in investment management and trading. Both computational and theoretical developments in machine learning have resulted in the increased practice of machine learning in the finance area.

The following is a list of the contribution of this paper:

- Proposed an Optimal Feature selection strategy using Pandas Data Frame to establish a framework for the dataset.
- Proposed a new Hybrid ML model by combining KNN & Optimized NB.
- Proposed an innovative features as input variables such as big data use, positive effect, threat and skill.

The following portion of the paper is carried out as follows: The work is divided into four 4 main sections. The paper examines previously completed work in its second section. Section 3 provides a description of the proposed accounting behavior prediction. A comprehensive discussion of the outcomes obtained with the proposed model is provided in Section 4. The conclusion is found in 5th section.

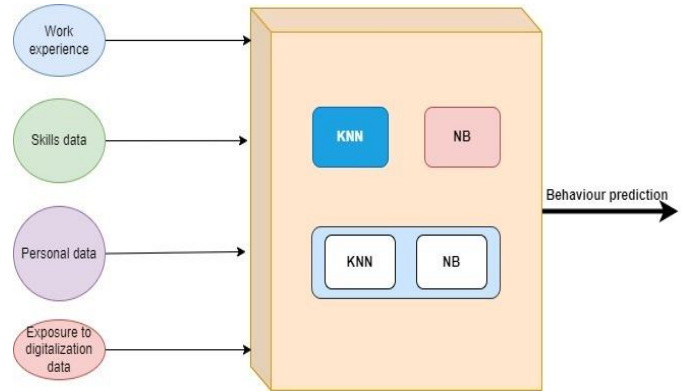
**MATERIALS AND METHODS**

This study utilized a research approach termed as design and creation research, a variation of Design Science Research (DSR) commonly applied in the realm of information systems and computing research for addressing issues through a problem-solving framework (Oates, 2006). DSR, grounded in engineering and scientific principles, primarily centers around the development of artifacts (Oates, 2006). Nonetheless, there remains ongoing debate regarding the definition, concepts, terminology, and supporting theories of DSR as a research paradigm (Johannesson et al., 2013). While the tenets of DSR lack a systematic delineation due to the intricate nature of science, information systems, and computing research, it constitutes a field often involving the convergence of organizational, human, and technological elements within problem-solving contexts. Hence, it is particularly well-suited for tackling challenges that encompass organizational, human, and technological dimensions.

**Framework for the Model for Predicting Accountant Behavior**

The framework is a guiding structure that directs the design and implementation of a system or process. It establishes a set of principles, standards, and best practices to ensure coherence and consistency in achieving desired objectives. The system framework is shown in Figure 1, and the left-hand side of the framework represents the data required for analysis, which is categorized into distinct group. The first category is "work experience," which includes information regarding experiences about the accountant, the years spent in service. The second category is "skill data," which comprises data related to the accountant's skills acquired over the years. The third category, "personal data," encompasses personal information about the accountants, such as name, age, gender, contact information, and any other relevant identifiers that help identify the accountant. The final category, "exposure to

digitalization data," pertains to information specific to the digitalization for which the account is being made. Once the data is categorized, it is fed into the model for processing and analysis. The model utilizes this input to predict the behavior of accountants. By employing suitable Naïve Bayes and KNN techniques, the model assesses the provided data and determines the behavioral activity of the accountants in Nigeria.



**Figure 1:** Framework for the Model for Predicting Accountant Behavior

The behavior predictions when analyzed helps to determine the effectiveness or otherwise of an accountant in the use of big data to analyzed relevant information required for an organization or company to make progress in their day to day running of their establishment. This information can help organization to predict both their future input and results output for years to come, helping them to establish a trend.

The pseudocode provided in Algorithm 1 for Predicting Accountant Behavior that utilizes Naïve Bayes and KNN techniques. It describes the basic steps that the algorithm would follow to detect account behavior.

**Algorithm 1**

**Algorithm for Predicting Accountant Behaviour**

1. Import required packages
2. Read dataset
3. Describe data
4. Fill in NaN values
5. Print correlation
6. Extract features
7. Normalization
8. Data splitting
9. Model initializations
10. Model combination
11. Models' evaluation
12. Plots

**Formulation of The Dataset**

The dataset used in this study was obtained via the administration of questionnaires. The gathered data, as shown in Table 1, underwent a process of cleaning and preprocessing. This dataset provides the researcher with the opportunity to train algorithms using the whole dataset.

**Table 1:** Brief Demography of Participants

A	B	C	D	E	F	G	H	I
S/N	AGE	NO OF YEAR OF PRACTICE	YEARS OF VIEW	YEAR OF WORK	AWARENESS	YEARS OF AWARENESS	% OF COMPUTER LITERACY	NO OF Y
1	35	5	2.5	8	1	23	60	
2	35	5	2.5	2.5	1	2.5	40	
3	45	10	13	13	0	8	60	
4	35	5	8	2.5	1	2.5	40	
5	45	10	13	13	0	8	20	
6	55	20	13	13	1	8	60	
7	55	20	18	18	0	13	60	
8	45	10	13	18	0	8	60	
9	45	10	13	8	0	13	40	
10	45	10	13	13	0	13	60	
11	45	10	8	8	0	8	60	
12	55	10	18	18	0	18	60	
13	45	20	13	13	1	2.5	60	
14	45	20	8	8	0	13	20	
15	55	10	18	13	0	13	80	
16	45	10	8	8	0	8	80	
17	45	10	2.5	13	0	13	80	
18	45	5	8	8	1	13	60	
19	35	10	2.5	2.5	1	2.5	80	

From Table 1, the data collected from the survey was then used to train a model, which was used to predict the behavior of accountants using naïve bayes and KNN algorithms. The dataset includes demographic variables (e.g., age, gender, education level), variables related to the use of machine learning in accounting (e.g., familiarity with machine learning, frequency of use), and variables related to attitudes towards machine learning (e.g., perceived usefulness, perceived ease of use). A total of 46 accountants participated in the survey, and the majority of the respondents were males (30%) compared to females (15%). The age distribution of the respondents was between 25 and 55 years, with an average age of 34 years. Furthermore, the highest proportion of respondents had a bachelor's degree (25%), while 10% had a master's degree, and 5% had a PhD. Perceived Benefits of Machine Learning The survey results indicated that the majority of the accountants (35%) believed that machine learning can enhance their job performance. The respondents stated that machine learning could improve their accuracy in financial reporting, reduce the time taken to complete tasks, and increase efficiency in their work processes. Moreover, 10% of the respondents believed that machine learning could help them to develop new skills, while 5% stated that machine learning could enhance their decision-making abilities. The results also revealed that 30% of the accountants perceived challenges in adopting machine learning in their profession. The respondents identified the lack of knowledge and skills as the most significant barrier to the adoption of machine learning in their work. They also cited the cost of acquiring and maintaining machine learning technology and the fear of job displacement as other significant challenges. Furthermore, 10% of the respondents believed that machine learning could pose a threat to the privacy and security of financial data. During this phase it was found that attributes such as no of years of practice, years of view, year of work, awareness, sex and age have a great impact on the mental education, whereas others, like Marital Status and Computer literacy, are considered redundant. The best features according to the results of the previous step were selected and the learning algorithms that had the best accuracy results were run on them.

**Performance Evaluation Matrices**

Diverse metrics are often used to assess the effectiveness of a categorization model. The measurements include Accuracy, Confusion Matrix, Precision, Recall, F1-score, Error Rate, and Training time.

**Classification Accuracy:** The term "accuracy" refers to the ratio of correct predictions produced by the model to the total number of

predictions made which is shown in equation 1.

$$\text{Accuracy (\%)} = \frac{TP+TN}{TP+TN+FP+FN} \times 100 \quad 1)$$

**Precision:** Precision is determined by the ratio of true positives to the total of true positives and false positives.

The formula for precision is given in equation 2:

$$\text{Precision} = \frac{TP}{TP+FP} \quad 2)$$

**Sensitivity;** Sensitivity is a numerical metric used to calculate the proportion of correctly identified positive situations that were incorrectly labeled as negative by the model. It is sometimes denoted as recall or true positive rate. Mathematically, it is defined as the ratio of the number of true positive (TP) occurrences to the sum of true positive and false negative (FN) cases which is shown in equation 3.

$$\text{Sensitivity} = \frac{TP}{TP+FN} \quad 3)$$

**Specificity**

Specificity is a synonym for the actual negative rate. The theoretical definition of the term involves the calculation of the ratio between the number of true negative (TN) instances and the sum of true negative and false positive (FP) cases.

Mathematically, the expression is given in equation 4:

$$\text{Specificity} = \frac{TN}{TN+FP} \quad 4)$$

**F-Score**

The F-score is a statistical metric used to evaluate the effectiveness of a binary classification model by measuring its capacity to reliably anticipate occurrences of the positive class. The computation employs the metrics of accuracy and recall. It is mathematically calculated using equation 5:

$$F1\text{-score} = 2 * \frac{\text{precision}+\text{recall}}{\text{precision}+\text{recall}} \quad 5)$$

**Error Rate (EER)**

The error rate may be computed by dividing the total count of wrong predictions made on the test set by the total count of predictions made on the test set.

Mathematically, it is expressed in equation 6:

$$\text{Error Rate} = \frac{\text{Incorrect Predictions}}{\text{Total Predictions}} \quad 6)$$

**RESULTS AND DISCUSSION**

**Numerical Experimental Performance of the Proposed Model**

This step evaluates the proposed model's quantitative experimental performance and refines the sample size and feature collection. The findings are thoroughly documented, using tables and figures.

**Importing the Libraries**

The Python Libraries Pandas, Numpy, CSV, Matplotlib, Seaborn, and Joblib were loaded into Jupyter Notebook. Numpy enables fast computation and broadcasting over multi-dimensional arrays by vectorization (Stančin & Jović, 2019). Pandas is a sophisticated and intuitive open-source program designed for the analysis and manipulation of data. It is constructed using the Python programming language (Subasi, 2020). Matplotlib and Seaborn are Python libraries especially tailored for data visualization. They

provide an intuitive interface for generating visually captivating and practical graphs. Seaborn is constructed upon the framework of Matplotlib and provides a slightly smaller range of functionalities in comparison to Matplotlib (Pintor et al., 2019). Figure 2 depicts the inclusion of Pandas, Numpy, CSV, Matplotlib, Seaborn, and Joblib Python libraries into the Jupyter Notebook.

```
import pandas as pd # For data manipulation
import numpy as np # For scientific computing
import csv
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.graph_objects as go
import joblib
```

Figure 1: Importing Python Libraries

### Loading the Dataset

Table 2 depicts the imported dataset. The 'hvplot.pandas' was employed for interactive data visualization, loads a dataset from 'Dataset.csv' using 'pd.read\_csv()', and provides initial insights through 'df.head()' and 'df.describe()'. It sets the stage for comprehensive data analysis and decision-making.

Table 2: Imported dataset and headers

	AGE	NO OF YEAR OF PRACTICE	YEARS OF VIEW	YEAR OF WORK	AWARENESS	YEARS OF AWARENESS	PER_OF_COMPUTER_LITERACY	NO_OF_YEAR	TRAINING
count	999.000000	999.000000	999.000000	999.000000	999.000000	999.000000	999.000000	999.000000	999.000000
mean	40.362362	12.333333	10.267611	9.799209	6.496496	5.055455	14.867868	70.454454	
std	8.829020	4.563952	5.006791	4.861046	5.500238	6.454079	17.267806	6.088990	18.133542
min	25.000000	5.000000	2.000000	2.000000	0.000000	2.000000	0.000000	5.000000	40.000000
25%	33.000000	8.000000	6.000000	6.000000	0.000000	7.000000	36.000000	10.000000	55.000000
50%	41.000000	12.000000	10.000000	10.000000	6.000000	12.000000	51.000000	15.000000	70.000000
75%	48.000000	16.000000	15.000000	14.000000	1.000000	18.000000	65.000000	20.000000	86.000000
max	55.000000	20.000000	18.000000	18.000000	1.000000	23.000000	80.000000	25.000000	100.000000

	MENTAL_EDUCATION	BIG_DATA_USE	POSITIVE_EFFECT	THREAT	SKILLS	SKILLS_2	LEVEL_OF_WORK	TIME_OF_WORK	RESULT	LEVEL
count	999.000000	999.000000	999.000000	999.000000	999.000000	999.000000	999.000000	999.000000	999.000000	999.000000
mean	38.179179	79.568569	51.487487	49.496496	68.840841	51.438438	49.870871	3.479980	6.086386	59.736732
std	23.737225	12.504631	24.259486	24.065453	18.062120	13.093389	11.855017	1.116285	2.046194	17.816133
min	0.000000	20.000000	10.000000	10.000000	10.000000	30.000000	30.000000	2.000000	1.500000	10.000000
25%	18.000000	69.000000	30.000000	29.000000	54.000000	40.000000	40.000000	2.750000	4.000000	44.000000
50%	40.000000	80.000000	52.000000	50.000000	69.000000	52.000000	50.000000	3.000000	6.000000	61.000000
75%	60.000000	90.000000	73.000000	70.000000	84.000000	61.000000	60.000000	4.000000	8.000000	75.000000
max	80.000000	100.000000	100.000000	90.000000	100.000000	100.000000	70.000000	5.000000	5.500000	90.000000

### Replacing Unknown Values in the Dataset

In the data pre-processing step, the system replaces all columns that contain question marks with non-null values. This process aims to handle missing or unknown values in the dataset. After replacing the question marks with non-null values, the system proceeds to print the updated data. This allows for a visual inspection of the dataset to verify the changes and ensure that the non-null values have been properly substituted. By replacing the question marks with non-null values and printing the data, the system addresses missing or unknown values in the dataset, ensuring that subsequent analysis or modeling can be performed accurately and reliably.

```
Replace missing values if any and replace non integral values.
In [ ]: # Positive: 1
df["BEHAVIOUR"].replace("Positive", 1, inplace=True)
# Negative: 0
df["BEHAVIOUR"].replace("Negative", 0, inplace=True)
# print
print(df["BEHAVIOUR"])

# df["BEHAVIOUR"] = df["BEHAVIOUR"].astype(float)
# print(df["BEHAVIOUR"])

# df["AWARENESS"] = df["AWARENESS"].astype(float)
# print(df["AWARENESS"])

0      0
1      0
2      1
3      1
4      1
...
994    1
995    1
996    1
997    0
998    0
Name: BEHAVIOUR, length: 999, dtype: int64
```

Figure 3: Replacing Unknown Values in the Dataset

### Printing of the Correlation

Table 3 shows the dataset correlation analysis. Eliminating "behavior" and "awareness" columns from the dataset streamlined the prediction process and improved efficiency. After being deemed unnecessary, these columns were removed. Note that deleting these columns made the dataset more concentrated and meaningful (Table 3). This intentional dataset dimensionality reduction has several benefits. First, it makes data analysis easier to handle and analyze. Second, it decreases computing overhead by processing fewer data points, which is especially useful for huge datasets. After removing unneeded columns, the algorithm reassesses variable correlation. After-processing correlation analysis is crucial for understanding the linkages and interdependencies between retained variables. Understanding these connections is crucial for understanding how factors interact and affect prediction. This rigorous technique to data reduction and correlation reevaluation streamlines and refines the dataset. It optimizes the dataset for modelling or analysis, ensuring that the variables are useful to prediction. This increases modeling efficiency, prediction accuracy, and interpretability, resulting in more informed and accurate results.

Table 3: Reduced Data after Dropping Irrelevant Columns

	AGE	NO OF YEAR OF PRACTICE	YEARS OF VIEW	YEAR OF WORK	AWARENESS	YEARS OF AWARENESS	PER_OF_COMPUTER_LITERACY	NO_OF_YEAR	TRAINING
AGE	1.000000	0.041774	0.008811	-0.014866	-0.041698	0.035110	-0.010366	0.077300	0.009832
NO OF YEAR OF PRACTICE	0.041774	1.000000	0.022052	-0.008164	-0.020360	0.023404	-0.001401	0.008112	0.007912
YEARS OF VIEW	0.008811	0.022052	1.000000	0.030335	0.015268	0.026275	-0.011987	0.023265	-0.017164
YEAR OF WORK	-0.014866	-0.008164	0.030335	1.000000	-0.012342	0.021947	0.042396	0.027303	-0.048027
AWARENESS	-0.041698	-0.020360	0.015268	-0.012342	1.000000	-0.022115	0.025245	0.000177	-0.017609
YEARS OF AWARENESS	0.035110	0.023404	0.026275	0.021947	-0.022115	1.000000	-0.007147	0.030234	-0.009110
PER_OF_COMPUTER_LITERACY	-0.010366	-0.001401	-0.011987	0.042396	0.025245	-0.007147	1.000000	0.028878	0.028316
NO_OF_YEAR	0.077300	0.008112	0.026265	0.027303	0.000177	0.030234	0.028878	1.000000	0.014821
TRAINING	0.009832	0.007912	-0.014744	-0.024927	-0.017609	0.009110	0.028316	0.014821	1.000000

### Feature Extraction

In this step of the data preparation process, the dataset is divided into its fundamental components: features (X) and the target variable (y). This is achieved by utilizing the loaded Pandas DataFrame 'df'. The 'X' variable is formed by excluding two specific columns, "BEHAVIOUR" and "AWARENESS," from the dataset, representing the input attributes for modeling. Simultaneously, the 'y' variable is created by isolating the "BEHAVIOUR" column, representing the target variable that signifies the behavior of interest. This step depicted in depicted in Table 4 is pivotal in organizing the dataset for subsequent data preprocessing, model development, and predictive performance evaluation.

**Table 4: Features Extraction**

	MENTAL_EDUCATION	BIG_DATA_USE	POSITIVE_EFFECT	THREAT	SKILLS	\
0	20	100	100	90	90	
1	0	60	80	50	50	
2	40	100	100	10	50	
3	40	80	80	10	30	
4	40	40	60	50	70	
..	...	...	...	...	...	...
994	24	98	25	44	69	
995	33	75	87	69	41	
996	27	72	13	69	65	
997	65	98	32	89	78	
998	67	61	69	34	63	

	SKILLS_2	LEVEL_OF_WORK	TIME_OF_WORK	RESULT	LEVEL	50
0	100	60	2.5	9.5	50	
1	40	30	5.5	7.5	30	
2	40	50	5.5	5.5	50	
3	60	30	5.5	7.5	30	
4	40	30	5.5	5.5	50	
..	...	...	...	...	...	...
994	43	59	2.0	4.0	74	
995	66	35	2.0	3.0	82	
996	56	58	3.0	4.0	88	
997	66	65	2.0	3.0	54	
998	44	47	2.0	6.0	72	

**Data Normalization**

In this step of the data preparation process, normalization is performed using the 'MinMaxScaler' from the 'scikit-learn' library. This process transforms the feature values to a standardized scale, typically ranging between 0 and 1. Normalization ensures that all features have the same scale, preventing any particular feature from dominating the modeling process due to its larger numerical range. This step enhances the model's ability to effectively learn from the data. The dataset is split into training and testing sets. This is achieved with the 'train\_test\_split' function, which partitions the data into two subsets. The 'X\_train' and 'y\_train' sets are designated for training the predictive model, while the 'X\_test' and 'y\_test' sets are reserved for evaluating the model's performance. The 'test\_size' parameter, set at 0.2, indicates that 20% of the data is allocated for testing, leaving 80% for training. Additionally, a 'random\_state' is set to ensure reproducibility of the split.

**Scale and split**

```
In [ ]: # Step 3: Normalize the features to have values between 0 and 1
scaler = MinMaxScaler()
X_normalized = scaler.fit_transform(X)

# Step 4: Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X_normalized, y, test_size=0.2, random_state=42)
```

**Figure 4: Scaling and Splitting of Dataset**

**RESULTS**

This section will present the data analysis results based on the model computation, test and comparison of the models and plotting the graph of the computed data set.

**Model Creation**

The Hybrid KNN-NB model proposed in this study combines two distinct models: Naive Bayes (NB) and k-Nearest Neighbors (KNN). This fusion aims to leverage the advantages of both models, resulting in a hybrid approach. This model is aimed at predicting the behavior of accountants focused on.

**K-Nearest Neighbors (KNN) Model**

The model employed in this study is a k-Nearest Neighbors (k-NN), implemented using the 'KNeighborsClassifier' imported from the 'scikit-learn' library. Figure 5 illustrates the model creation process. Subsequently, the k-NN model is initialized with 'n\_neighbors' set to 3, signifying that it considers the three nearest data points when making predictions based on features. Following initialization, the 'knn\_model' is trained with the provided training data ('X\_train' and

'y\_train'), enabling it to learn from the patterns in the training dataset. Afterwards, the trained k-NN model is applied to make predictions on the testing data ('X\_test'), and these predictions are stored in the 'knn\_predictions' variable. Lastly, the 'knn\_predictions' are printed, revealing the model's computed predictions based on the input features. This step is pivotal for evaluating the model's performance and its capacity to make accurate predictions on previously unseen test data, serving as a critical aspect of model assessment.

**KNN**

```
In [ ]: from sklearn.neighbors import KNeighborsClassifier

# Create k-NN model
knn_model = KNeighborsClassifier(n_neighbors=3)
knn_model.fit(X_train, y_train)

# Make predictions
knn_predictions = knn_model.predict(X_test)
print(knn_predictions)
```

**Figure 5: Creation of K-Nearest Neighbours (KNN) Model**

**KNN Prediction**

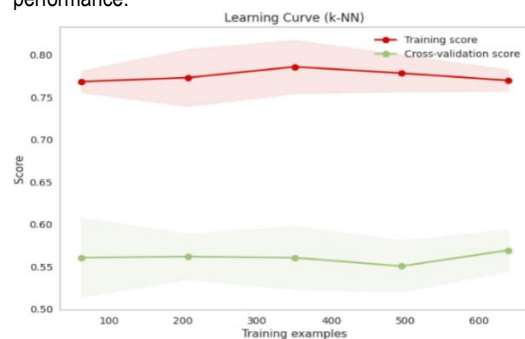
Following training, the model is applied to make predictions on the testing data ('X\_test'). The predicted values are stored in the 'knn\_predictions' variable. Finally, the 'knn\_predictions' are printed, showing the model's computed predictions. The output is a binary sequence of 0s and 1s depicted in Figure 6, representing the model's predictions for each corresponding data point in the testing dataset. These predictions are fundamental for evaluating the model's performance and its ability to accurately classify data points into the appropriate categories.

```
[0 0 1 1 1 1 1 1 0 1 1 0 1 1 1 0 1 0 1 1 1 1 1 1 1 1 0 0 1 1 0 0 1 0 1 1 1 0 1
 1 1 1 1 1 1 1 1 1 1 0 1 0 1 1 1 1 1 1 1 1 1 0 1 1 0 1 1 1 1 1 1 1 0 1 0 0 1 0 1
 0 1 1 0 0 0 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 0 1 0 0 1 1 1 1
 1 0 1 1 0 0 1 0 1 1 1 1 1 1 0 0 0 1 1 0 1 1 0 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 0
 1 0 0 1 1 1 0 0 1 1 0 1 0 1 0 1 1 1 0 0 1 0 1 0 1 1 0 1 0 1 1 1 1 1 1 1 0 1 1 0 1
 0 1 1 1 1 1 1 0 0 1 1 1 1 1 0]
```

**Figure 6: Computed Predictions for K-Nearest Neighbors (KNN) Model**

**Learning Curve for KNN Model**

In Figure 7, a learning curve is presented for the k-Nearest Neighbors (k-NN) model, achieved through the utilization of the 'plot\_learning\_curve' function. The learning curve aids in making informed decisions about model training, assisting in the diagnosis of issues like overfitting or underfitting, and guiding the optimization of model complexity and data size for enhanced predictive performance.



**Figure 7: Learning Curve for KNN Model**





## Conclusion

This paper undertakes a thorough investigation of the dynamic terrain where technology, accounting, and human behavior converge. The research primarily focuses on examining the complex interplay between accountants and the disruptive influences of information technology. It extensively investigates the effects of many topics, such as Big Data, Accounting Systems, Machine Learning, Cloud Computing, and Data Science, on the field of accounting. The study utilizes a methodological framework that integrates elements of positivism with interpretivism, enabling the researcher to conduct empirical investigations while also gaining a comprehensive grasp of the research topic.

A crucial element of this research is on the formulation and implementation of prognostic models aimed at comprehending and foretelling the conduct of accountants within the framework of contemporary technology. The concept utilizes Python, a highly adaptable programming language well recognized for its proficiency in data analysis and Machine Learning (ML). The use of Python is a deliberate decision due to its strategic advantages, as it provides users the ability to utilize robust libraries such as NumPy and Pandas for efficient data processing, as well as Matplotlib and Seaborn for effective data visualization. The use of a strong technical infrastructure guarantees that the study is firmly rooted in rigorous data analysis methodologies.

The core of this research is on the development and execution of predictive models, specifically emphasizing the K-Nearest Neighbors (KNN) and Naïve Bayes (NB) algorithms. These models are used both alone and in combination to form a hybrid model known as KNN-NB. The objective of this study is to evaluate and contrast the efficacy of various models in forecasting the behavior of accountants. The study utilizes assessment measures, including accuracy, precision, recall, and F1 score, to assess the efficacy of each model. The results of this research provide significant perspectives on the changing responsibilities of accountants in light of technology progress. The Naïve Bayes model has great performance, obtaining an accuracy rate of 63% and displaying a notable capability in accurately identifying positive cases, as shown by its high recall of 97%. In contrast, the K-Nearest Neighbors model demonstrates a balanced precision and recall, yielding an F1 score of 64%, while obtaining an accuracy rate of 52%. The hybrid KNN-NB model, which integrates components from both the K-Nearest Neighbors (KNN) and Naive Bayes (NB) techniques, has a 52% accuracy rate, suggesting its capability in leveraging the respective advantages of both algorithms. This study serves as a fundamental basis for future inquiries into the modeling and simulation of accountant behavior within the context of artificial intelligence and data analytics. The thesis recommended that accountants and accounting firms should continually improve their knowledge regarding machine learning/ artificial intelligence as this will enhance the performance of accounting functions, thereby eliminating certain accounting cost.

**Conflict of Interest:** The corresponding author, representing all authors, confirms the absence of any conflict of interest.

## REFERENCE

- Abd-Mutalib, H., Muhammad Jamil, C. Z., Mohamed, R., & Ismail, S. N. A. (2023). The determinants of environmental knowledge sharing behavior among accounting educators: a modified theory of planned behavior. *International Journal of Sustainability in Higher Education*, 24(5), 1105-1135.
- Abdualgalil and Abraham (2020): Naive Bayes Modest Probabilistic Classifier Based on the Bayesian: A Review. *International Journal of Advanced Research in Computer and Communication Engineering*, 2(1), 17-21
- Abdural, K. K., Jha, R., & Afroz, S. (2020). Data Mining Techniques for Intrusion Detection: A Review. *International Journal of Advanced Research in Computer and Communication Engineering*, 3(6), 19-23
- Agarwal, R., Joshi, M.V.:(2017): PNrule: A New Framework for Learning Classifier Models in Data Mining. Tech. Report, Dept. of Computer Science, University of Minnesota
- Aksoy, S.: k-Nearest Neighbor Classifier and Distance Functions. Technical Report, Department of Computer Engineering, Bilkent University (February 2016)
- Alcazar, et al. (2020). Classical Versus Quantum Models in Machine Learning: Insights from a Finance Application. *Machine Learning: Science and Technology*, 1(3), 035003 <https://doi.org/10.1088/2632-2153/ab9009>  
*American Journal of Industrial and Business Management*, 8, 1817-1824. <https://doi.org/10.4236/ajibm.2018.88123>
- Arnott, R., Harvey, C. R., & Markowitz, H. (2019). A Backtesting Protocol in the Era of Machine Learning. *The Journal of Financial Data Science*, 1(1), 64-74.
- Aziz, S., Dowling, M. M., Hammami, H., & Piepenbrink, A. (2019). Machine Learning in Finance: A Topic Modeling Approach. <http://dx.doi.org/10.2139/ssrn.3327277>
- Azizam, N. A., Dzulklipli, M. R., Shamimi, N. I., Maon, S. N., John, D., Belawing, J., ... & Yahaya, N. (2020). Applicability of theory of planned behavior and protection motivation theory in predicting intention to purchase health insurance. *Advances in Business Research International Journal*, 6(1), 1-9.
- Azman, N. L. A., & Vaicondam, Y. (2020). Behavioral intention in forensic accounting services. *International Journal of Psychosocial Rehabilitation*, 24(2), 1837-1846.
- Bhargava, N., Sharma, G., Bhargava, R., & Mathuria, M. (2018). Decision tree analysis on j48 algorithm for data mining. *Proceedings of International Journal of Advanced Research in Computer Science and Software Engineering*, 3(6).
- Bhavsar, Y. B., & Waghmare, K. C. (2019). Intrusion detection system using data mining technique: Support vector machine. *International Journal of Emerging Technology and Advanced Engineering*, 3(3), 581-586
- Bracke, P., & Datta, et al. (2019). Machine Learning Explain ability in Finance: An Application to Default Risk Analysis. Bank of England Working Paper No. 816. <http://dx.doi.org/10.2139/ssrn.3435104>
- Brown, J. O., Hays, J., & Stuebs Jr, M. T. (2016). Modeling accountant whistleblowing intentions: Applying the theory of planned behavior and the fraud triangle. *Accounting and the Public Interest*, 16(1), 28-56.
- Chukwuani, V. N., & Egiyi, M. A. (2020). Automation of Accounting Processes: Impact of Artificial Intelligence. *International Journal of Research and Innovation in Social Science (IJRISS)*, 4, 444-449.



- Chukwudi, O., Echefu, S., Boniface, U., & Victoria, C. (2018). Effect of Artificial Intelligence on the Performance of Accounting Operations among Accounting Firms in South East Nigeria. *Asian Journal of Economics, Business and Accounting*, 7, 1-11. <https://doi.org/10.9734/AJEB/2018/41641>
- Dixon, M., & Halperin, I. (2019). The Four Horsemen of Machine Learning in Finance. *Economia*. <https://www.icaew.com/technical/technology/artificial-intelligence/artificial-intelligence-articles/how-artificial-intelligence-will-impact-accounting>
- Frey, G. P. C., Yu, J. X., Yu, P. S., & Lu, H. (2016, August). Parameter free bursty events detection in text streams. In *Proceedings of the 31st international conference on Very large data bases* (pp. 181-192). VLDB Endowment
- Ganti, M., Prasad, P., & Wankhade, L. K. (2021, September). A new data mining based network intrusion detection model. In *Computer and Communication Technology (ICCT), 2021 International Conference on* (pp. 731-735). IEEE.
- Grano, C., Singh Solorzano, C., & Di Pucchio, A. (2022). Predictors of protective behaviours during the Italian COVID-19 pandemic: An application of protection motivation theory. *Psychology & Health*, 37(12), 1584-1604.
- Griffin, O. (2019, October 6). How Artificial Intelligence Will Impact Accounting.
- Grindley, E. J., Zizzi, S. J., & Nasypany, A. M. (2008). Use of protection motivation theory, affect, and barriers to understand and predict adherence to outpatient rehabilitation. *Physical therapy*, 88(12), 1529-1540.
- Heaton, J. B., Polson, N. G., & Witte, J. H. (2017). Deep Learning in Finance: Deep Portfolios. *Applied Stochastic Models in Business and Industry*, 3(1), 3-12. <https://doi.org/10.1002/asmb.2209>
- Huang, et al. (2020). Deep learning in finance and banking: A literature Review and classification. *Frontiers of Business Research in China*, 14, 2-24. <https://doi.org/10.1186/s11782-020-00082-6>
- Huttunen. Jeong, K., Hong, T., Chae, M. and Kim, J. (2019) Development of a Decision Support Model for Determining the Target Multi-Family Housing Complex for Green Remodeling Using Data Mining Techniques. *Energy and Buildings*, 202, Article ID: 109401. <https://doi.org/10.1016/j.enbuild.2019.109401>
- ICEAW (2018). Artificial intelligence and the future of accountancy. ISBN 978-1-78363-816-1. Available: [www.icaew.com/itfac](http://www.icaew.com/itfac). (Accessed 25 July 2021)
- Israel, et al. (2020). Can Machines 'Learn' Finance? *Journal of Investment Management*. Available at <https://doi.org/10.37200/IJIM/RP200125>
- Jalota, C. and Agrawal, R. (2019) Analysis of Educational Data Mining using Classification. *International Conference on Machine Learning, Big Data, Cloud and Parallel Computing (COMITCon)*, 14-16:
- Kaur N., Sahdev S. L., Sharma M. & Siddiqui L. (2020). Banking 4.0: "the influence of artificial intelligence on the banking industry & how AI is changing the face of modern-day banks". *International Journal of Management*. Vol. 11 No. 6. pp. 577-585.
- Kaur, B., Ahuja, L. and Kumar, V. (2019) Crime against Women: Analysis and Prediction Using Data Mining Techniques. *International Conference on Machine Learning, Big Data, Cloud and Parallel Computing (COMITCon)*, 14-16 February 2019, Faridabad. <https://doi.org/10.1109/COMITCon.2019.8862195>
- Najjar D. (2019). Is artificial intelligence (AI) the future of accounting? *The Balance Small Business*. Available: <https://www.thebalancesmb.com/is-artificial-intelligence-the-future-of-accounting-4083182> (accessed 13 July 2021)
- Olaiya, F and Adeyemo, A and Barnabas (2012). "Application of Data Mining Techniques in Weather Prediction and Climate Change Studies". *International Journal of Information Engineering and Electronic Business*, 2012, 1, 51-59
- Özer, G., & Yilmaz, E. (2011). Comparison of the theory of reasoned action and the theory of planned behavior: An application on accountants' information technology usage. *African Journal of Business Management*, 5(1), 50-58.
- Quoc, T. N., Bao, Q. P. T., Huu, B. N., & Bao, A. N. P. (2023, May). Motivating Accounting Information Systems Security Policy Compliance: Insight from the Protection Motivation Theory and the Theory of Reasoned Action. In *International Conference on Emerging Challenges: Strategic Adaptation in The World of Uncertainties (ICECH 2022)* (pp. 342-359). Atlantis Press.
- Reddy, P. S., Yayaswi, K. R. K., & Kumar, B. K. (2019). Accounting Intelligence—The New Era in Accounting. *Journal of Information and*