

**ANALYSIS OF MACHINE LEARNING MODELS FOR AIRLINE
PASSENGERS' SATISFACTION USING PYSPARK**

**COURSEWORK:
BIG DATA ANALYTICS AND DATA VISUALISATION
7153CEM**

BY

**AYOBAMI ENIOLA AGBOOLA
15620073**

**DATA SCIENCE AND COMPUTATIONAL INTELLIGENCE
COVENTRY UNIVERSITY,
COVENTRY, UK.**

DECEMBER, 2024

ANALYSIS OF MACHINE LEARNING MODELS FOR AIRLINE PASSENGERS' SATISFACTION USING PYSPARK

ABSTRACT

The rapid growth of airline travel and large volume of data in the industry has necessitated advanced tools to analyse and predict passenger satisfaction, a critical factor in maintaining competitive advantage. This study employs PySpark to implement machine learning models for predicting airline passenger satisfaction. The dataset, consisting more than 120,000 records of passenger feedback across various categories, was preprocessed to handle missing values, encode categorical variables, and assemble feature vectors. Four machine learning models, logistic regression, support vector machine (SVM), random forest, and gradient boosted trees (GBT) were trained, evaluated, and also optimized through hyperparameter tuning. The results revealed that ensemble models outperformed linear models in the predictive performance, achieving accuracy of 96% after tuning. Feature importance analysis highlighted technological service quality, travel type, and customer type as significant contributors to satisfaction predictions. Additionally, the study observed consistent evaluation metrics across models, underscoring the robustness of the PySpark framework in handling large, complex datasets. This research demonstrates the applicability of big data analytics in customer satisfaction analysis and provides a pathway for future studies to enhance airline customer experience through data driven insights. The findings emphasize the critical role of machine learning in transforming passenger feedback into actionable strategies for airlines.

INTRODUCTION

The aviation industry plays a pivotal role in global transportation and connects millions of passengers daily across various destinations. In order to retain customer patronage, increase market-share and, ultimately profitability, airlines need to prioritize high-quality service to their passengers, thus ensuring customers' satisfaction has become a critical factor for the sustainability and success of airlines in a highly competitive market (Liou et al., 2011, p. 131).

Satisfaction levels are influenced by a range of factors, including flight schedules, in-flight services, customer service, and ticket pricing (Clemes et al., 2008, p. 49). Understanding and predicting passenger satisfaction is essential for airlines to improve service quality, retain customers, and maintain a competitive edge (Shah et al, 2020, p. 159). Therefore, it is very important for airlines to periodically evaluate customer satisfaction, perceived service quality and loyalty intentions of the business (Liou et al., 2011).

Data driven approaches leveraging machine learning models have emerged as powerful tools to uncover insights into customer satisfaction patterns and predict outcomes effectively. The analysis, processing, and storing of big data can be challenging due to data imperfection, massive data size, computational difficulty, and lengthy evaluation time because of the sheer volume, velocity, and variety. Big Data frameworks like Hadoop MapReduce and Spark are potent tools that provide an effective way to analyse huge datasets and Apache Spark is one of the most widely used large-scale data processing engines due to its speed, low latency in-memory computing, and powerful analytics (Malik et al., 2024, p. 151785).

This study investigates the ability of four classification models to predict airline passengers' satisfaction while employing PySpark to implement scalable machine learning pipelines due to its suitability in handling large datasets.

RELATED WORKS

Customer satisfaction has become a key intermediary objective in service operations due to the benefits it brings to organizations. Previous research has demonstrated that satisfaction is strongly associated with re-purchase. Chen (2008, p. 709) highlighted the need for organizations to provide their passengers with quality services in order to gain a competitive advantage and earn a high profit with sustained development in the airline industry.

Chang and Yeh (2002) used fuzzy multicriteria analysis modelling to formulate the service quality of airlines, applying fuzzy set theory to describe the ambiguity between the criteria weight and performance ratings of each airline, and they found that flight safety is the most important factor among the 15 attributes used to measure the quality of service. In their work, Chen and Chang (2005) examined the gap between the actual service rendered by airlines and passengers' service expectations in order to determine airlines' service quality. They reported that gaps did exist and that passengers were most interested in the responsiveness and assurance dimensions from airline frontline staff.

Feng and Park (2023) utilized applied machine learning based classifiers on a public data from the national information society of South Korea to predict users' satisfaction and identified 10 essential features that act as indicators in the prediction. Anand and Supriya (2023) utilized PySpark framework approach for airline delay prediction, implementing multiple algorithms and reported that logistic regression and random forest classifiers performed better than others.

DATASET DESCRIPTION

The airline passenger satisfaction dataset used in this study was obtained from Kaggle (<https://www.kaggle.com/api/v1/datasets/download/mysarahmadbhat/airline-passenger-satisfaction>). The dataset contains 24 features, the information about more than 120,000 passengers and their evaluation of different factors such as cleanliness, comfort, service and overall experience. The dataset has 5 categorical columns, 18 integer columns and 1 column of type float. All columns have no null values, except 'Arrival Delay' which has 393 null values. The target variable is the 'Satisfaction' column with 71858 'Neutral/Dissatisfied' instances and 55088 instances of 'Satisfied' passengers. Table 1 shows the description of each variable column.

METHODOLOGY

The growing demand for processing large-scale datasets has necessitated the development of efficient tools like Apache Spark and to leverage the features Hadoop and Spark offers in big data analysis, therefore it is imperative to install and configure the programs.

Hadoop

In order to install the Hadoop, Linux was firstly installed on VirtualBox using Ubuntu. After that, a "temp" folder was created on Desktop from the Terminal and the openjdk version of java installed. Using 'wget', hadoop 3.4.0 version was installed and uncompressed with 'tar -xzf'. The location for Java "export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64" and Hadoop "\$ export HADOOP_HOME=/home/Ayobami/Desktop/temp/hadoop-3.4.0" were set and added to path, and the hadoop was installed. The installation confirmation of Hadoop and some of its usage for map reduce are shown in "Appendix A".

Spark

After installing Hadoop 3.4.0, Spark was downloaded and uncompressed with the ‘tar -xzf’. The location for the Spark was set “\$ export SPARK_HOME=/home/Ayobami/Desktop/temp/spark-3.5.3-bin-hadoop3 and added to path. Java and Hadoop location were also set and added to path prior. The installation confirmation for Spark is shown in “Appendix B”.

Data Preprocessing

PySpark, the Python API for Apache Spark, offers robust distributed computing capabilities, making it ideal for handling large datasets and performing advanced machine learning tasks. The PySpark’s architecture supports iterative computations and in-memory processing, which significantly accelerates data analysis and machine learning workflows. PySpark’s MLlib library includes tools for classification, regression, clustering, and feature transformation, making it a versatile solution for diverse data science tasks.

The dataset was uploaded into PySpark, and an exploratory data analysis process was conducted to assess its structure (“Figure 1”), identify missing values, and evaluate data quality. It was observed that the ‘Arrival Delay’ column contained 393 missing values, representing about 3% of the total dataset. Given the relatively small proportion of missing data, these rows were removed following best practices in data preprocessing (Kaiser, 2014, p. 44). As revealed by the box plot (“Figure 2”), outliers were removed from the ‘Arrival Delay’, ‘Departure Delay’ and ‘Flight Distance’ columns.

The PySpark session was initialized to enable distributed data processing. This session served as the foundation for managing data transformations and model training. The dataset’s Numerical features were directly used in model training, while categorical features were preprocessed for compatibility with machine learning algorithms. The Categorical features were transformation using the PySpark’s StringIndexer. The ‘Satisfaction’ column was designated as the target variable (label) and other columns were combined into a single feature vector using PySpark’s VectorAssembler.

The preprocessed dataset was split into training and testing subsets using PySpark’s randomSplit method. The training set Contained 80% of the dataset for model training, ensuring sufficient data diversity for learning patterns and the testing set has 20% of the data for unbiased evaluation of model performance.

Model Development

This study employed four popular machine learning models, logistic regression, support vector machines (SVM), random forest, and gradient boosted trees (GBT) to predict customer satisfaction using the Airline Passenger Satisfaction dataset. After splitting the dataset into training and test sets for each model, the models were defined with ‘features’ for ‘featuresCol’ and ‘label’ for ‘labelCol’. The training process involved fitting the models to the training set and predictions made with the test set, using the transform method of each model. The hyperparameter optimization was performed by defining grids using ParamGridBuilder and cross validation with the CrossValidator, which utilized the training set to identify the best-performing hypaparameters and predictions from the test set generated using the best model.

The performance of each model was evaluated with the accuracy, precision, recall, and F1-score metrics. Additionally, feature importance analysis was conducted, the absolute coefficients values were ranked to determine the important features for logistic regression and SVM models and the built-in importance scores used to rank the important features of random forest and GBT.

Logistic regression

The logistic function is a mathematical function which converts the output of linear regression into a probability score between 0 and 1, where the score represents the likelihood of the event being predicted (Anand and Supriya, 2023, p. 3). Logistic Regression is a statistical approach to analyse methods of predicting a data value that is entirely based on the observations and is suitable for interpreting the influence of input features on the output (Sen et al., 2023).

SVM

Support vector machine (SVM) is a statistical learning theory that seeks to identify the hyperplane that most effectively divides instances. It is excellent at defining ideal decision borders because it can find data points that affect where the boundary between classes is placed. SVM's performance might vary depending on the parameters selected and the features of the dataset (Haque and Hassan 2023, p. 5)

Random Forest

Random forest classification is based on concept of ensemble learning, a technique for integrating many classifiers to handle tough problems. It contains a number of decision trees on various subsets of a given dataset and takes the average to improve the predictive accuracy of that dataset. The outcome depends on majority votes of projections rather than relying solely on one decision tree (Viswanatha et al., 2023, p. 9).

Gradient Boosting Trees (GBT)

Gradient boosting tree, known for its ability to handle complex relationships and large feature spaces is a method where trees are built in series and compared to each other based on mathematical scores of splits which combines the principles of gradient boosting with decision trees. An ensemble of decision trees are built sequentially and each subsequent tree corrects the mistakes of the previous tree which iteratively improves the model's predictions by minimizing a loss function using gradient descent optimization. Therefore, the GBT creates a strong predictive model by combining multiple weak decision trees to classify observations accurately (Anand and Supriya, 2023, p. 3).

RESULTS

Model Performance Metrics

The models' performances were analysed using accuracy, precision, recall, and F1-score as performance metrics. Each model underwent both baseline and hyperparameter-tuned evaluations to enhance performance.

All the models achieved high accuracy ("Figure 3") and were excellent at predicting both classes ("Figure 4"). The GBT model (95%) particularly outperformed others, followed by random forest (92%). The SVM and logistic regression models had identical performance metrics (87%), indicating similar behaviour in predicting customers' satisfaction from the dataset. Hyperparameter tuning improved the performance of the random forest and GBT models to 96%, demonstrating the impact of optimized parameter selection on ensemble methods. While logistic regression and SVM demonstrated comparable performance, they were however significantly outperformed by random forest and GBT. This performance gap highlights the advantage of ensemble learning techniques in capturing complex feature interactions and their robustness in handling diverse datasets (Wei et al, 2023, p. 349).

The models were uniformly effective across the evaluation metrics, as it was observed that the values of accuracy, precision, recall, and F1-score were identical in each model, highlighting the reliability of the modelling process and the suitability of the dataset for the predictive tasks. This comprehensive analysis demonstrates the effectiveness of machine learning models in predicting customer satisfaction while identifying actionable insights through feature importance.

Feature Importance Analysis

Feature importance was analysed (“Figure 5”) for each model to understand the factors contributing most significantly to customer satisfaction predictions. Features like ‘On board Wifi Service’ and ‘Class’ emerged as critical predictors across all models, Logistic Regression and SVM gave utmost priority to ‘Type of Travel’ and ‘Customer Type’ while the ensemble methods (Random Forest and Gradient Boosted Trees) gave highest importance to ‘Online Boarding’, which indicate passenger satisfaction with technological services.

Conclusion

In this study, machine learning models were applied using PySpark, to predict customer satisfaction from airline passenger data. The dataset, comprising various features that reflected customers’ experiences, which enabled the exploration of predictive performance with logistic regression, support vector machine (SVM), random forest, and gradient boosted trees (GBT) models. Comprehensive preprocessing steps which include handling missing values, indexing categorical features, and vectorizing the predictors ensured that the data was suitably prepared for analysis.

The results demonstrated that ensemble methods, particularly GBT, outperformed other models, whereas logistic regression and SVM maintained consistent metrics across all evaluations. The consistent performance of logistic regression and SVM can be attributed to their linear structures, which suited the dataset's characteristics.

Feature importance analysis revealed that customer centric attributes, such as ‘travel type’, ‘customer type’, and ‘onboard services’, were significant predictors of satisfaction across all the models. Ensemble methods particularly emphasized features linked to technological engagement, such as ‘online boarding’ and inflight entertainment, indicating these elements play a critical role in modern customer satisfaction.

Overall, this study has shown the effectiveness of machine learning models using PySpark for large datasets to provide robust, scalable solutions for predictive analytics. While ensemble methods exhibited superior performance, simpler models like logistic regression and SVM remain valuable for interpretability and computational efficiency in less complex scenarios. This study has therefore shown that businesses can understand and anticipate customer needs better by leveraging big data analytics, which would help foster improved customer satisfaction and in return, loyalty.

Future work could explore other algorithms and include additional metrics to deepen insights into model behaviour and performance nuances. Also, more granular analysis is recommended to further refine and explore subtle differences between the models’ performances.

REFERENCES

- Anand, R. N., & Supriya, M. (2023). Enhancing Airline Operations by Flight Delay Prediction - A PySpark Framework Approach. *2023 International Conference on Ambient Intelligence, Knowledge Informatics and Industrial Electronics*, 1–4. <https://doi.org/10.1109/aikiie60097.2023.10389960>
- Chang, W.-J., Chen, G.-J., & Yeh, Y.-L. (2002). Fuzzy Control of Dynamic Positioning Systems for Ships. *Journal of Marine Science and Technology*, 10(1). <https://doi.org/10.51400/2709-6998.2300>
- Chen, C.-F. (2008). Investigating structural relationships between service quality, perceived value, satisfaction, and behavioral intentions for air passengers: Evidence from Taiwan. *Transportation Research Part A: Policy and Practice*, 42(4), 709–717. <https://doi.org/10.1016/j.tra.2008.01.007>
- Chen, F.-Y., & Chang, Y.-H. (2005). Examining airline service quality from a process perspective. *Journal of Air Transport Management*, 11(2), 79–87. <https://doi.org/10.1016/j.jairtraman.2004.09.002>
- Clemes, M. D., Gan, C., Kao, T.-H., & Choong, M. (2008). An empirical analysis of customer satisfaction in international air travel. *Innovative Marketing*, 4(2), 48–62.
- Feng, Y., & Park, J. (2023). Using machine learning-based binary classifiers for predicting organizational members' user satisfaction with collaboration software. *PeerJ Computer Science*, 9, e1481.
- Haque, F. M., & Hassan, M. M. (n.d.). Bank Loan Prediction Using Machine Learning Techniques. *ArXiv Preprint ArXiv:2410.08886*.
- Kaiser, J. (2014). Dealing with Missing Values in Data. *Journal of Systems Integration*, 5(1), 42–51. <https://doi.org/10.20470/jsi.v5i1.178>
- Liou, J. J. H., & Tzeng, G.-H. (2007). A non-additive model for evaluating airline service quality. *Journal of Air Transport Management*, 13(3), 131–138. <https://doi.org/10.1016/j.jairtraman.2006.12.002>
- Malik, H., Sronach, R. L., & Witzig, S. B. (n.d.). Nature photography: a spark for discussion and action on environmental issues. *Journal of Outdoor and Environmental Education*, 1-23.
- Shah, S. S. A., & Khan, Z. (2019). Corporate social responsibility: a pathway to sustainable competitive advantage? *International Journal of Bank Marketing*, ahead-of-print(ahead-of-print), 159–174. <https://doi.org/10.1108/ijbm-01-2019-0037>
- Viswanatha, V., Ramachandra, A. C., Vishwas, K. N., & Adithya, G. (2023). LOAN APPROVAL PREDICTION USING MACHINE LEARNING. *International Research Journal of Modernization in Engineering Technology and Science*, 13(4), 7–19. <https://doi.org/10.56726/irjmets39658>
- Wei, Y., Zhang, D., Gao, M., Tian, Y., He, Y., Huang, B., & Zheng, C. (2023). Breast cancer prediction based on machine learning. *Journal of Software Engineering and Applications*, 16(18), 348–360.

Table 1. Data Description

Variable Name	Description
ID	Unique passenger identifier
Gender	Gender of the passenger (Female/Male)
Age	Age of the passenger
Customer Type	Type of airline customer (First-time/Returning)
Type of Travel	Purpose of the flight (Business/Personal)
Class	Travel class in the airplane for the passenger seat
Flight Distance	Flight distance in miles
Departure Delay	Flight departure delay in minutes
Arrival Delay	Flight arrival delay in minutes
Departure and Arrival Time Convenience	Satisfaction level with the convenience of the flight departure and arrival times from 1 (lowest) to 5 (highest) - 0 means "not applicable"
Ease of Online Booking	Satisfaction level with the online booking experience from 1 (lowest) to 5 (highest) - 0 means "not applicable"
Check-in Service	Satisfaction level with the check-in service from 1 (lowest) to 5 (highest) - 0 means "not applicable"
Online Boarding	Satisfaction level with the online boarding experience from 1 (lowest) to 5 (highest) - 0 means "not applicable"
Gate Location	Satisfaction level with the gate location in the airport from 1 (lowest) to 5 (highest) - 0 means "not applicable"
On-board Service	Satisfaction level with the on-boarding service in the airport from 1 (lowest) to 5 (highest) - 0 means "not applicable"
Seat Comfort	Satisfaction level with the comfort of the airplane seat from 1 (lowest) to 5 (highest) - 0 means "not applicable"
Leg Room Service	Satisfaction level with the leg room of the airplane seat from 1 (lowest) to 5 (highest) - 0 means "not applicable"
Cleanliness	Satisfaction level with the cleanliness of the airplane from 1 (lowest) to 5 (highest) - 0 means "not applicable"
Food and Drink	Satisfaction level with the food and drinks on the airplane from 1 (lowest) to 5 (highest) - 0 means "not applicable"
In-flight Service	Satisfaction level with the in-flight service from 1 (lowest) to 5 (highest) - 0 means "not applicable"
In-flight Wifi Service	Satisfaction level with the in-flight Wifi service from 1 (lowest) to 5 (highest) - 0 means "not applicable"
In-flight Entertainment	Satisfaction level with the in-flight entertainment from 1 (lowest) to 5 (highest) - 0 means "not applicable"
Baggage Handling	Satisfaction level with the baggage handling from the airline from 1 (lowest) to 5 (highest) - 0 means "not applicable"
Satisfaction	Overall satisfaction level with the airline (Satisfied/Neutral or unsatisfied)

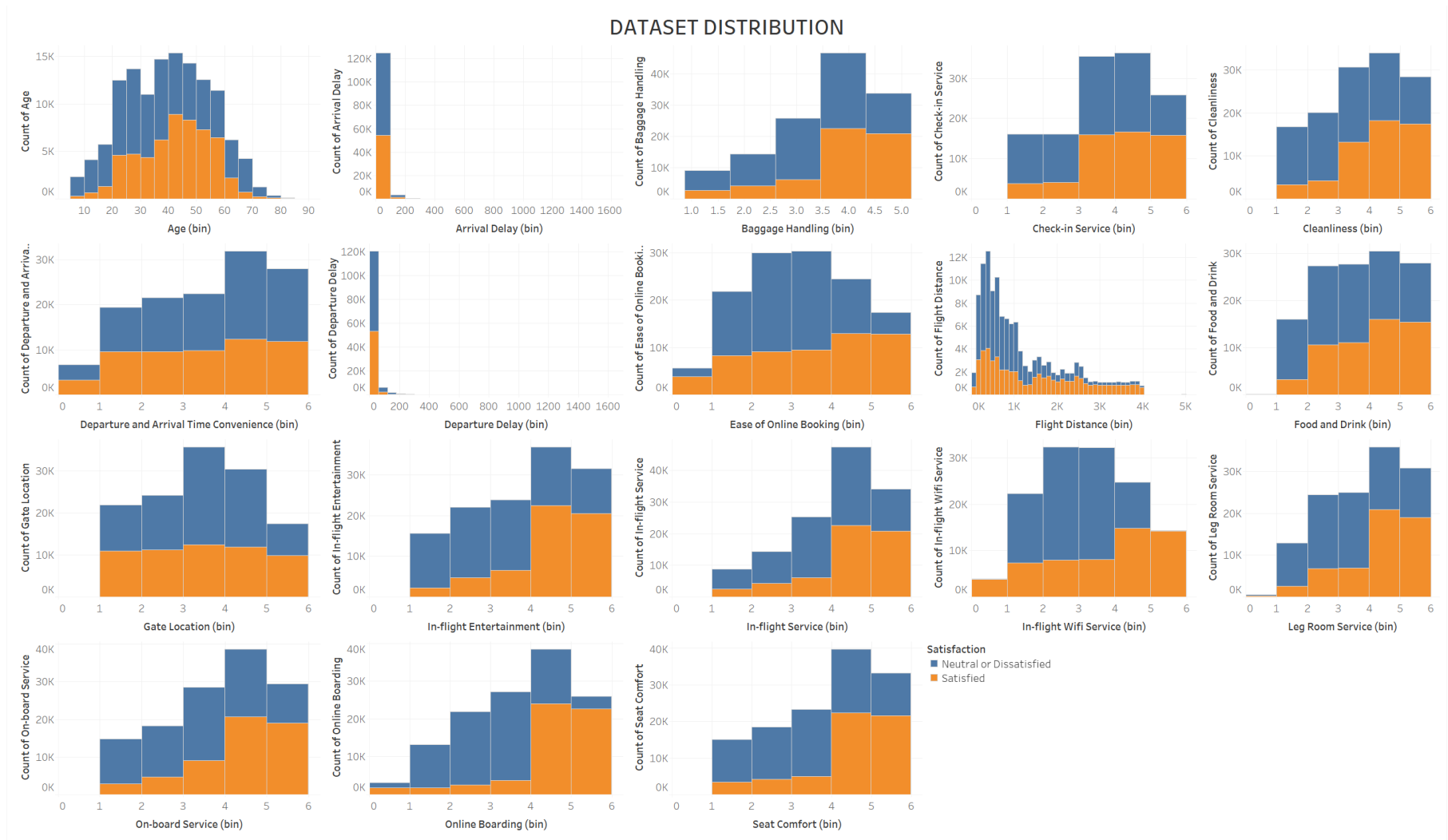


Fig. 1 Histogram showing distribution of the numerical features by ‘Satisfaction’

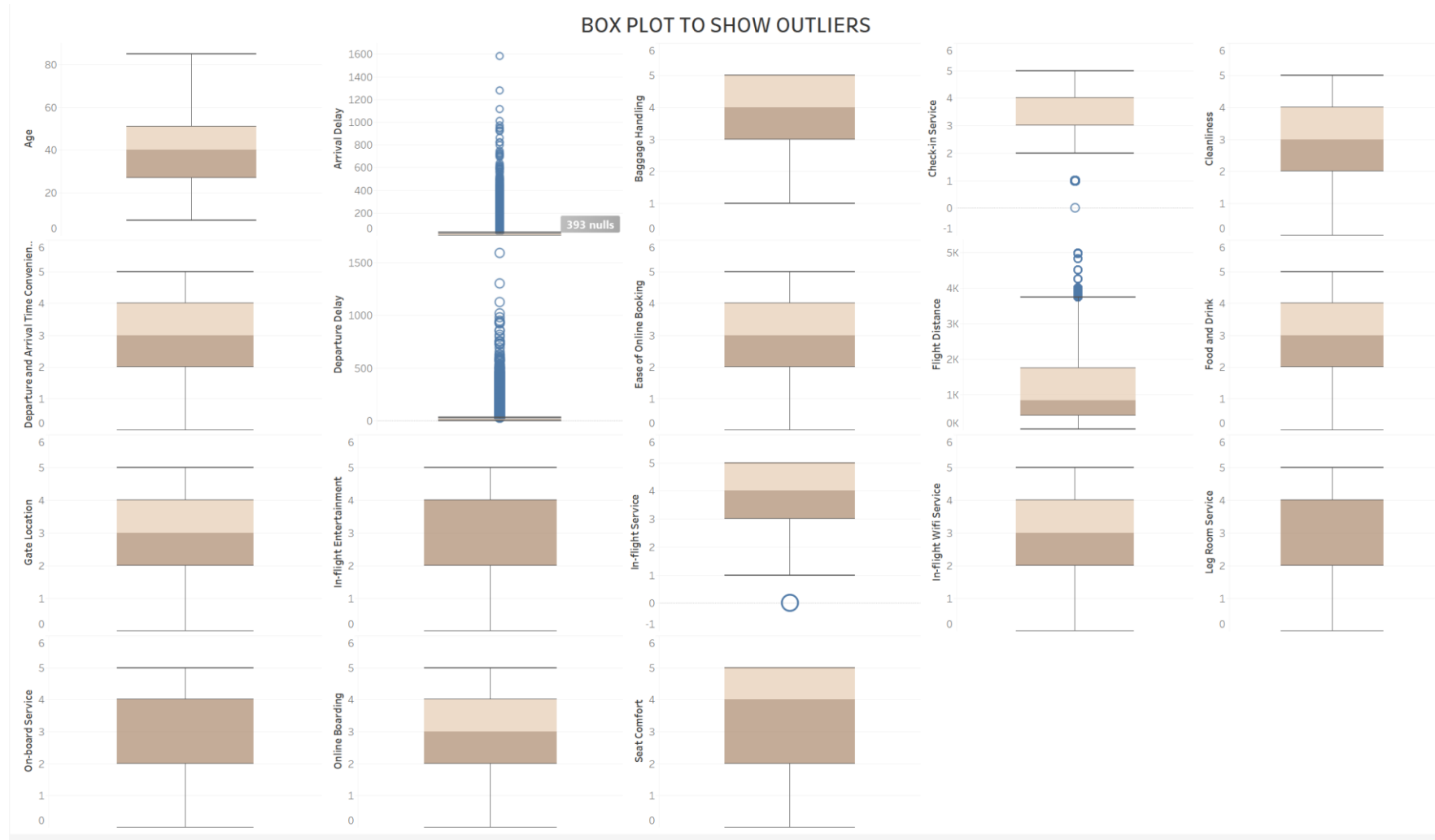


Fig. 2. Box plot showing outliers

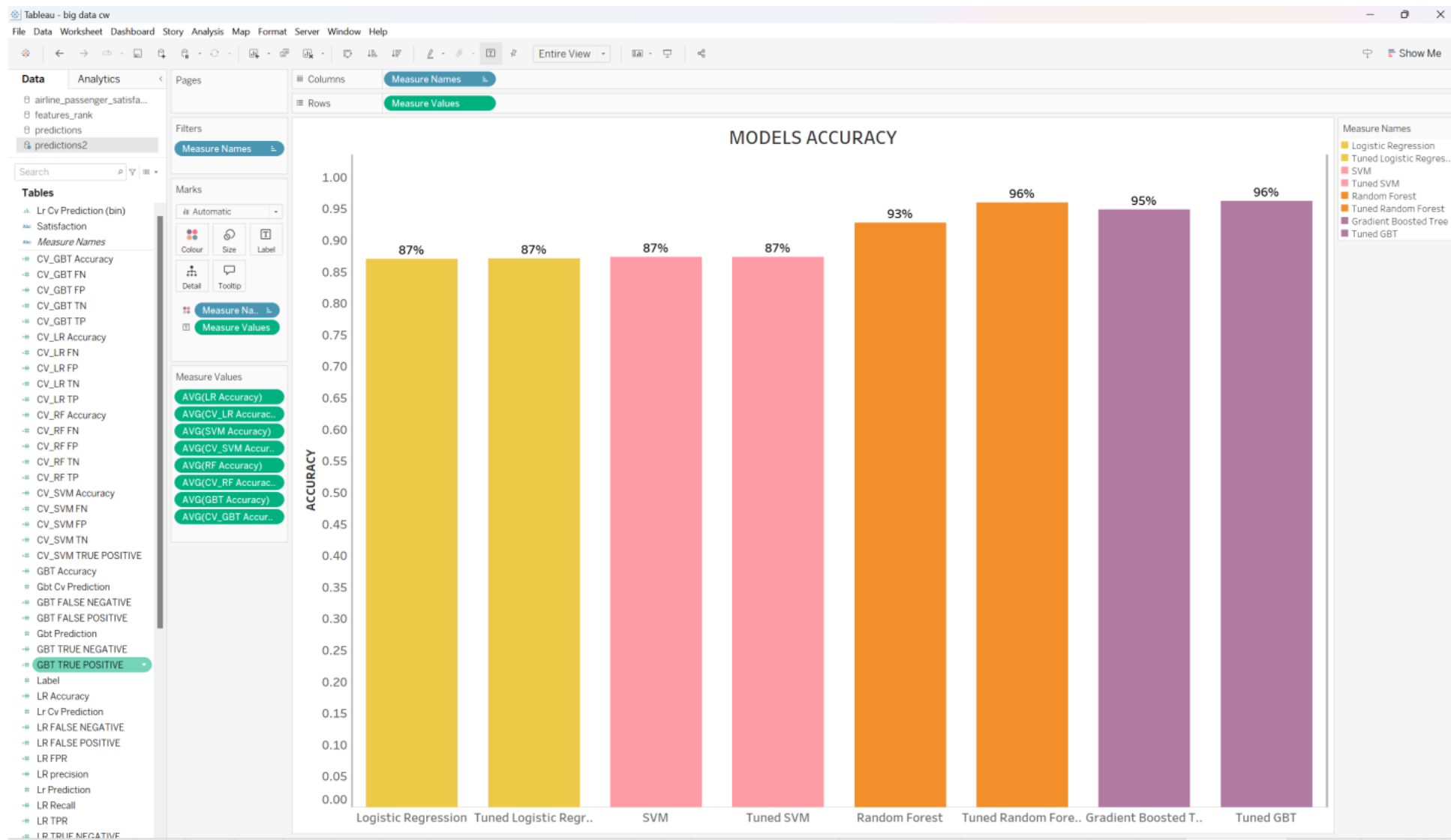


Fig 3. Models' Accuracy

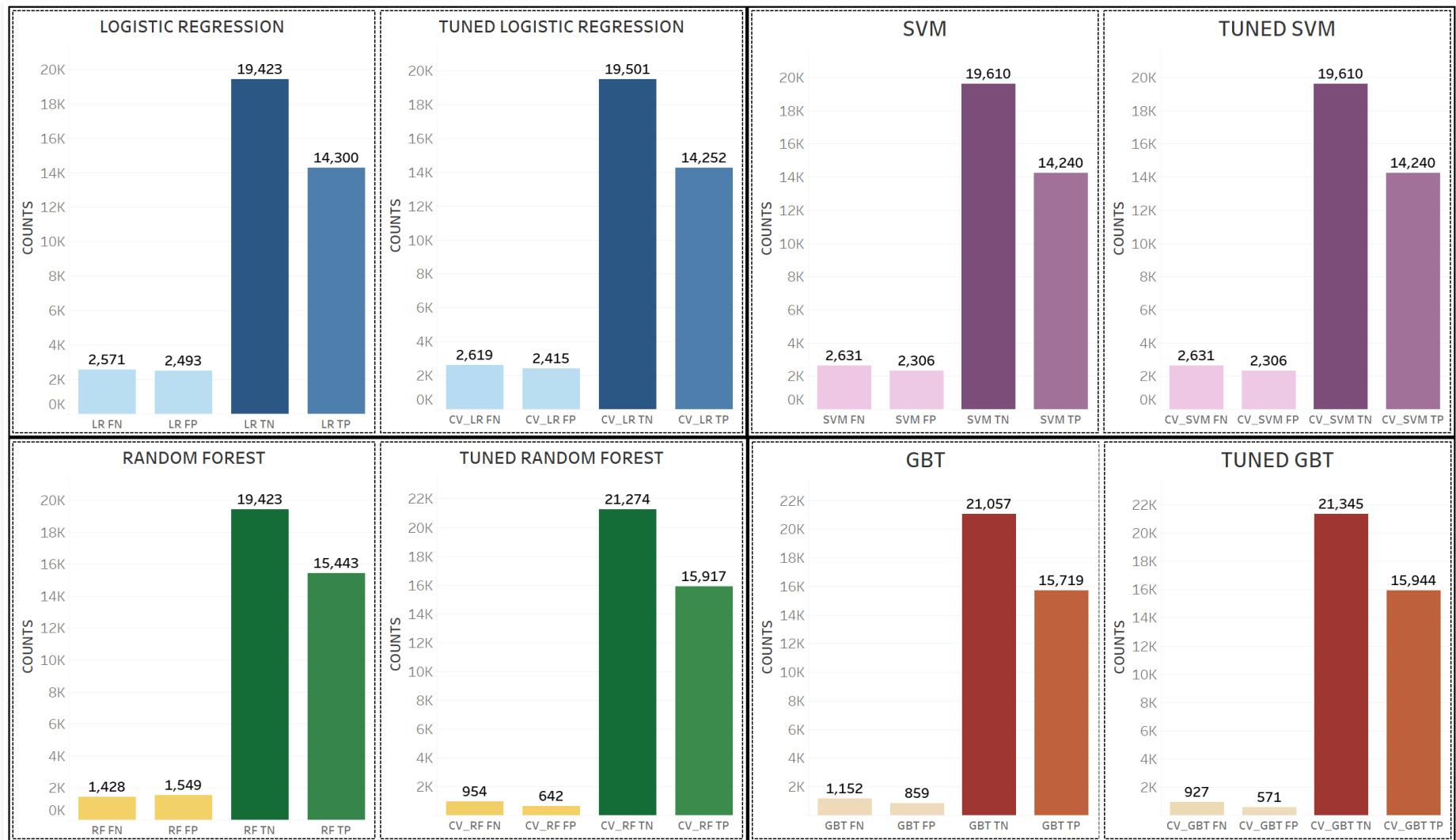
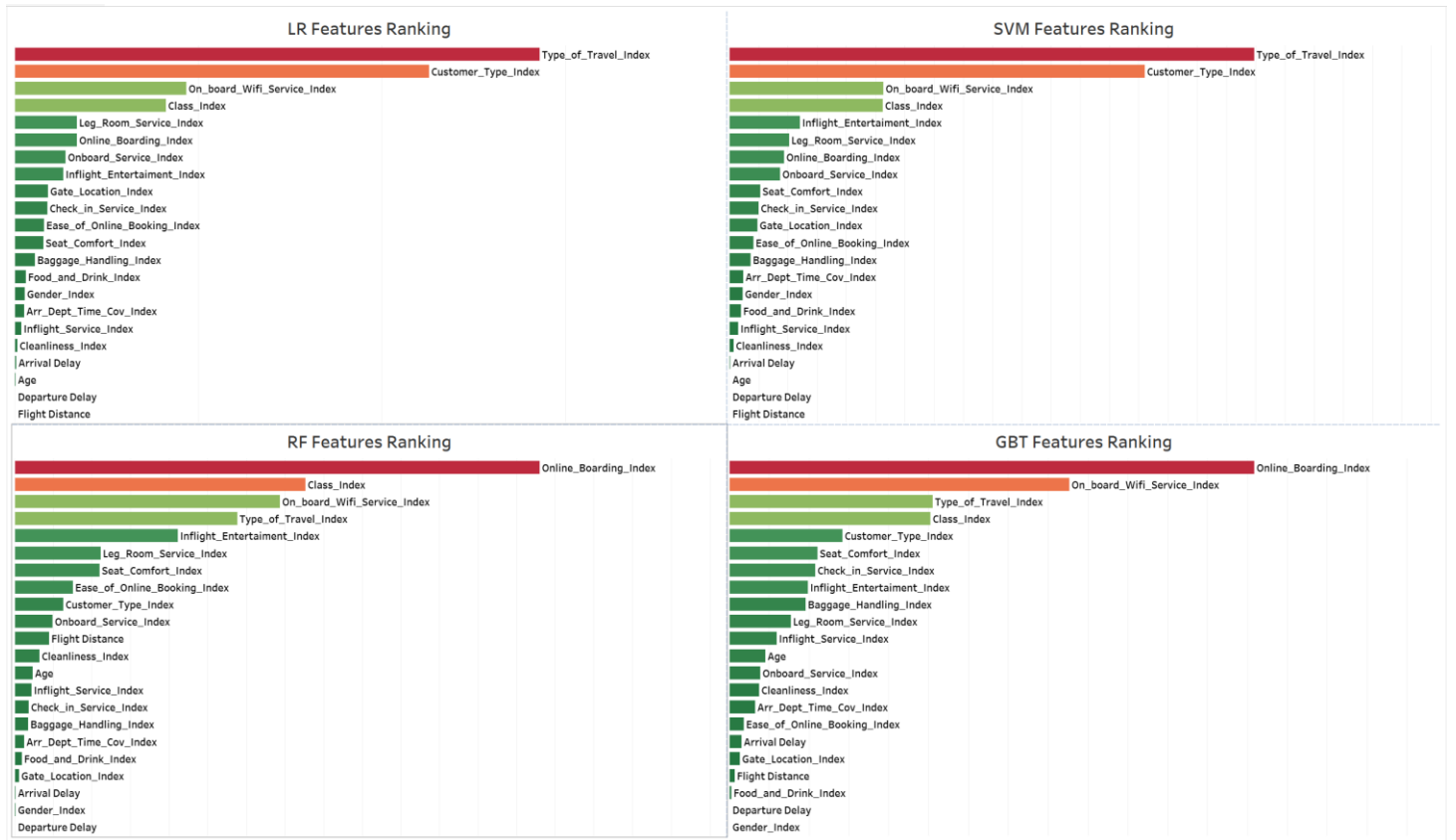
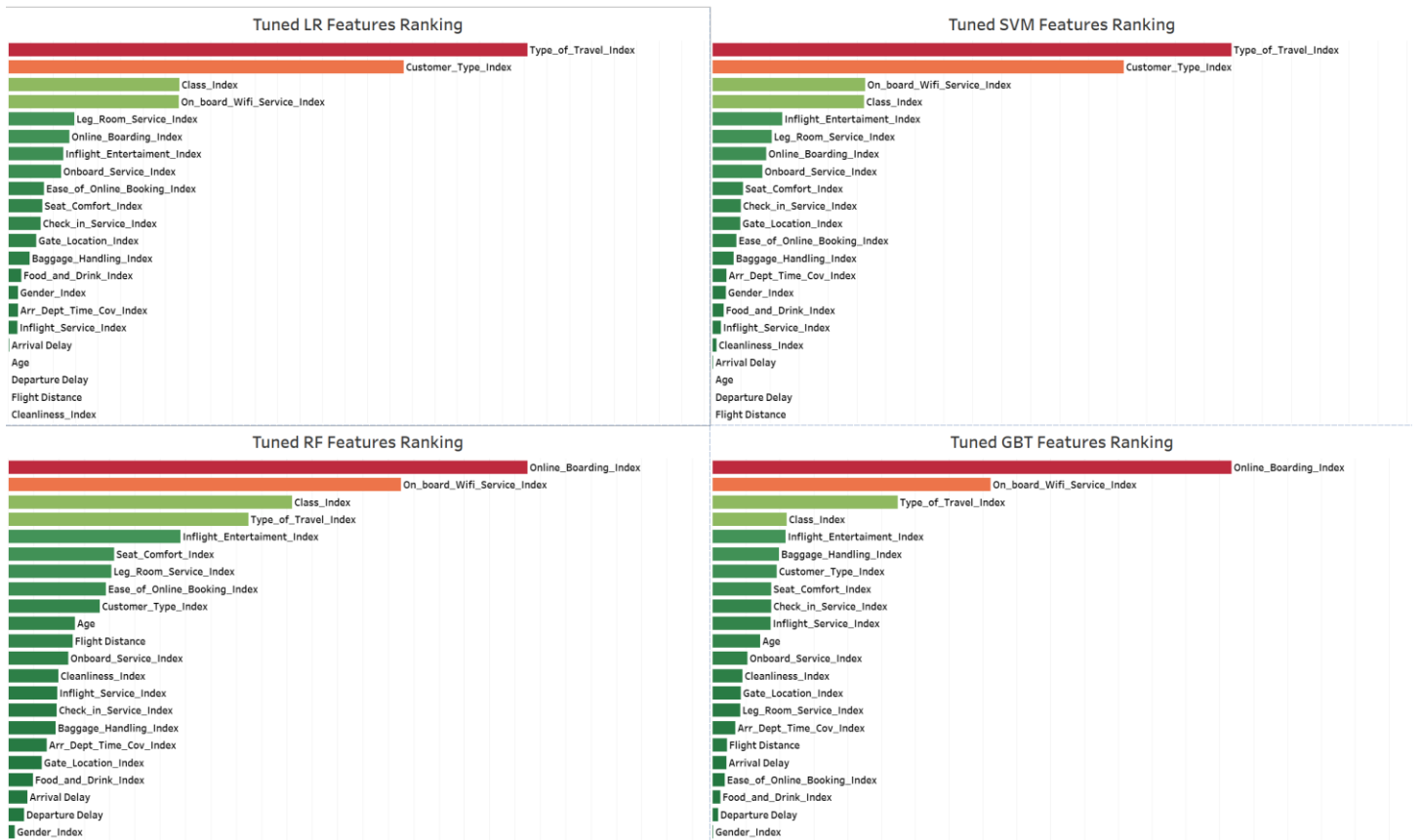


Fig. 4. Models' Predictions



(a) Normal models

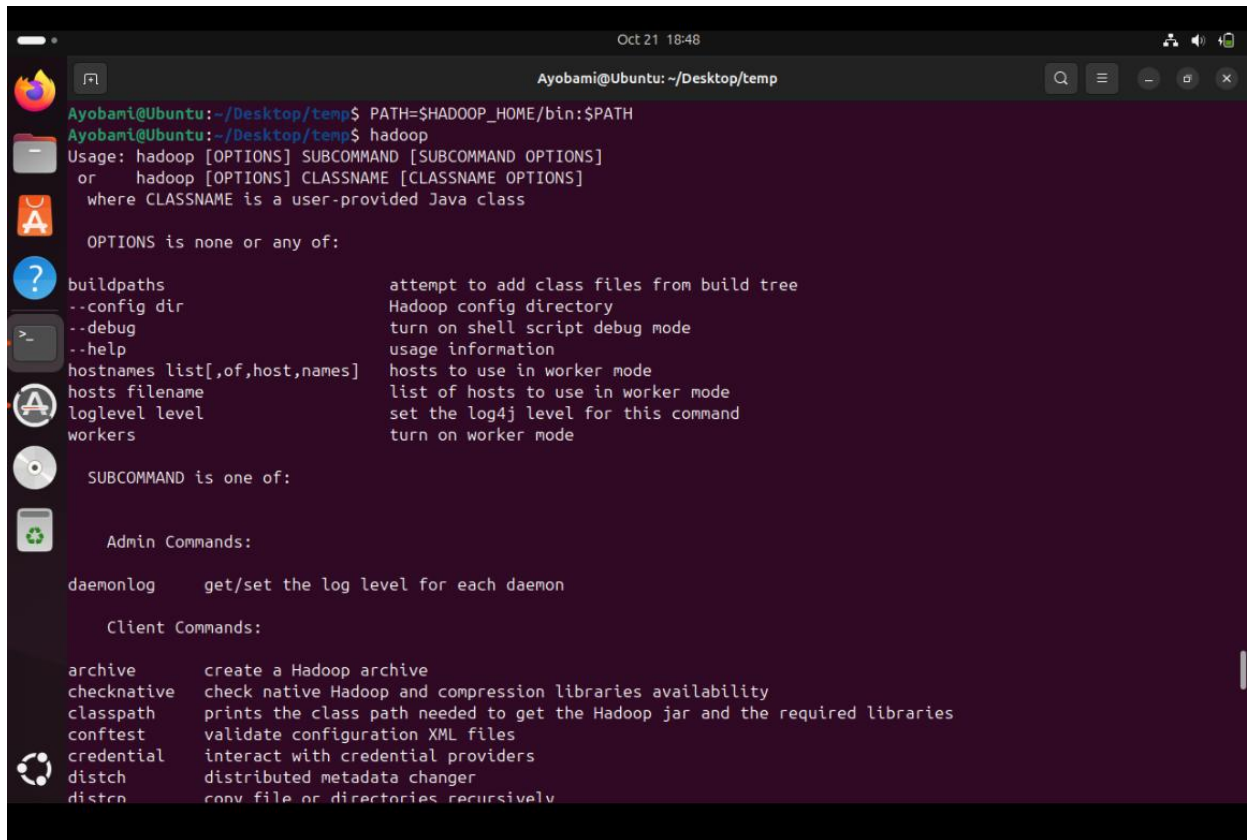


(b) Optimized models

Fig. 5. Important Features

APPENDIX A

HADOOP



A terminal window titled "Ayobami@Ubuntu: ~/Desktop/temp" showing the output of the `hadoop` command. The output displays the usage of `hadoop` with various options and subcommands. The window has a dark theme and a sidebar with application icons on the left.

```
Ayobami@Ubuntu:~/Desktop/temp$ PATH=$HADOOP_HOME/bin:$PATH
Ayobami@Ubuntu:~/Desktop/temp$ hadoop
Usage: hadoop [OPTIONS] SUBCOMMAND [SUBCOMMAND OPTIONS]
or      hadoop [OPTIONS] CLASSNAME [CLASSNAME OPTIONS]
       where CLASSNAME is a user-provided Java class

       OPTIONS is none or any of:

       buildpaths          attempt to add class files from build tree
       --config dir        Hadoop config directory
       --debug             turn on shell script debug mode
       --help             usage information
       hostnames list[,of,host,names] hosts to use in worker mode
       hosts filename      list of hosts to use in worker mode
       loglevel level      set the log4j level for this command
       workers            turn on worker mode

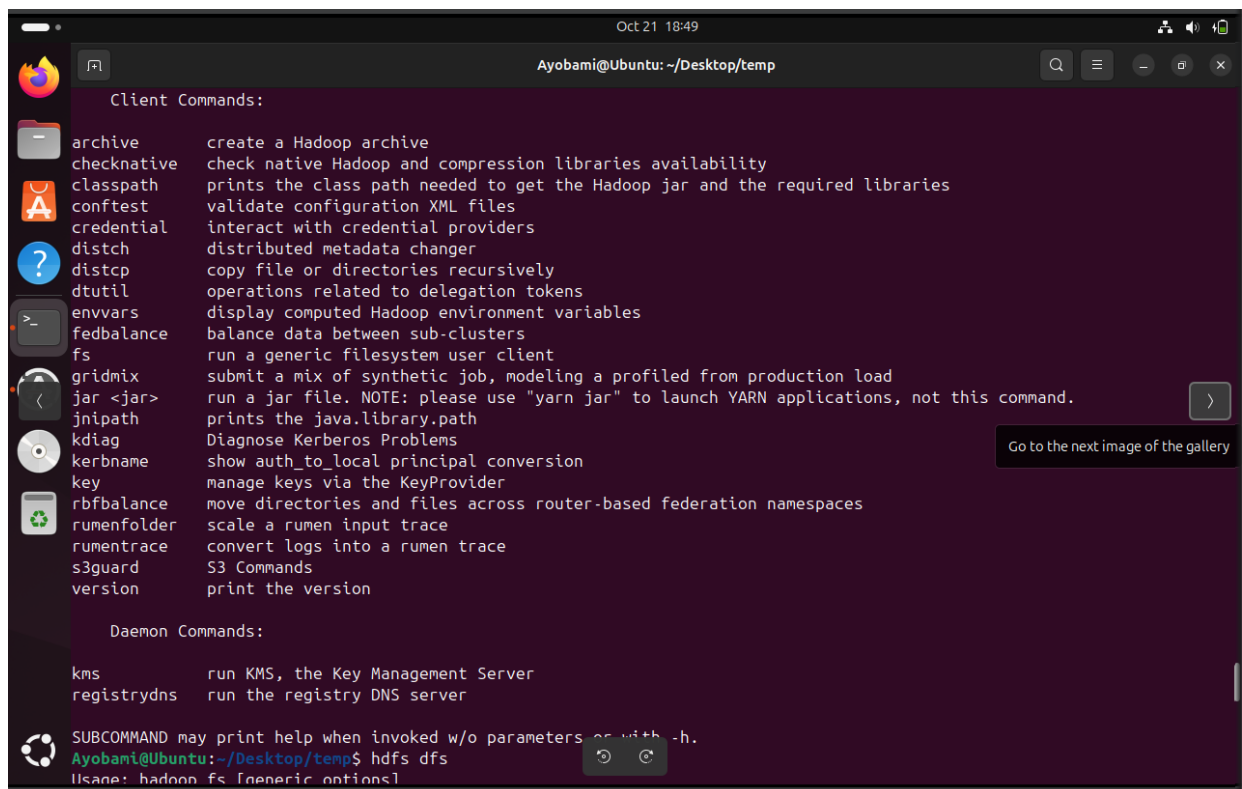
       SUBCOMMAND is one of:

       Admin Commands:

       daemonlog          get/set the log level for each daemon

       Client Commands:

       archive            create a Hadoop archive
       checknative        check native Hadoop and compression libraries availability
       classpath          prints the class path needed to get the Hadoop jar and the required libraries
       conftest           validate configuration XML files
       credential         interact with credential providers
       distch             distributed metadata changer
       distcp            copy file or directories recursively
```



A terminal window titled "Ayobami@Ubuntu: ~/Desktop/temp" showing the output of the `hadoop` command. The output displays the usage of `hadoop` with various options and subcommands. The window has a dark theme and a sidebar with application icons on the left.

```
Ayobami@Ubuntu:~/Desktop/temp$ hadoop
Usage: hadoop [OPTIONS] SUBCOMMAND [SUBCOMMAND OPTIONS]
or      hadoop [OPTIONS] CLASSNAME [CLASSNAME OPTIONS]
       where CLASSNAME is a user-provided Java class

       OPTIONS is none or any of:

       buildpaths          attempt to add class files from build tree
       --config dir        Hadoop config directory
       --debug             turn on shell script debug mode
       --help             usage information
       hostnames list[,of,host,names] hosts to use in worker mode
       hosts filename      list of hosts to use in worker mode
       loglevel level      set the log4j level for this command
       workers            turn on worker mode

       SUBCOMMAND is one of:

       Admin Commands:

       daemonlog          get/set the log level for each daemon

       Client Commands:

       archive            create a Hadoop archive
       checknative        check native Hadoop and compression libraries availability
       classpath          prints the class path needed to get the Hadoop jar and the required libraries
       conftest           validate configuration XML files
       credential         interact with credential providers
       distch             distributed metadata changer
       distcp            copy file or directories recursively
       dtutil            operations related to delegation tokens
       envvars           display computed Hadoop environment variables
       fedbalance        balance data between sub-clusters
       fs                run a generic filesystem user client
       gridmix           submit a mix of synthetic job, modeling a profiled from production load
       jar <jar>         run a jar file. NOTE: please use "yarn jar" to launch YARN applications, not this command.
       jnipath           prints the java.library.path
       kdiag            Diagnose Kerberos Problems
       kerbname         show auth to local principal conversion
       key              manage keys via the KeyProvider
       rfbalance        move directories and files across router-based federation namespaces
       rumenfolder       scale a rumen input trace
       rumentrace       convert logs into a rumen trace
       s3guard          S3 Commands
       version          print the version

       Daemon Commands:

       kms              run KMS, the Key Management Server
       registrydns     run the registry DNS server

       SUBCOMMAND may print help when invoked w/o parameters or with -h.
Ayobami@Ubuntu:~/Desktop/temp$ hadoop fs dfs
Usage: hadoop fs [generic options]
```

```
Oct 21 18:49
Ayobami@Ubuntu: ~/Desktop/temp

registrydns run the registry DNS server

SUBCOMMAND may print help when invoked w/o parameters or with -h.
Ayobami@Ubuntu:~/Desktop/temp$ hdfs dfs
Usage: hadoop fs [generic options]
    [-appendToFile [-n] <localsrc> ... <dst>]
    [-cat [-ignoreCrc] <src> ...]
    [-checksum [-v] <src> ...]
    [-chgrp [-R] GROUP PATH...]
    [-chmod [-R] <MODE[,MODE]... | OCTALMODE> PATH...]
    [-chown [-R] [OWNER][:[GROUP]] PATH...]
    [-concat <target path> <src path> <src path> ...]
    [-copyFromLocal [-f] [-p] [-l] [-d] [-t <thread count>] [-q <thread pool queue size>] <localsrc> ... <dst>]
    [-copyToLocal [-f] [-p] [-crc] [-ignoreCrc] [-t <thread count>] [-q <thread pool queue size>] <src> ... <localdst>]
    [-count [-q] [-h] [-v] [-t <storage type>] [-u] [-x] [-e] [-s] <path> ...]
    [-cp [-f] [-p | -p[topax]] [-d] [-t <thread count>] [-q <thread pool queue size>] <src> ... <dst>]
    [-createSnapshot <snapshotDir> [<snapshotName>]]
    [-deleteSnapshot <snapshotDir> <snapshotName>]
    [-df [-h] [<path> ...]]
    [-du [-s] [-h] [-v] [-x] <path> ...]
    [-expunge [-immediate] [-fs <path>]]
    [-find <path> ... <expression> ...]
    [-get [-f] [-p] [-crc] [-ignoreCrc] [-t <thread count>] [-q <thread pool queue size>] <src> ... <localdst>]
    [-getfacl [-R] <path>]
    [-getfattr [-R] {-n name | -d} [-e en] <path>]
    [-getmerge [-nl] [-skip-empty-file] <src> <localdst>]
    [-head <file>]
    [-help [cmd ...]]
    [-ls [-C] [-d] [-h] [-q] [-R] [-t] [-S] [-r] [-u] [-e] [<path> ...]]
    [-mkdir [-p] <path> ...]
    [-moveFromLocal [-f] [-p] [-l] [-d] <localsrc> <localdst>]
    [-moveToLocal <src> <localdst>]
```

```
Oct 21 18:49
Ayobami@Ubuntu: ~/Desktop/temp

[-moveToLocal <src> <localdst>]
[-mv <src> ... <dst>]
[-put [-f] [-p] [-l] [-d] [-t <thread count>] [-q <thread pool queue size>] <localsrc> ... <dst>]
[-renameSnapshot <snapshotDir> <oldName> <newName>]
[-rm [-f] [-r] [-R] [-skipTrash] [-safely] <src> ...]
[-rmdir [-ignore-fail-on-non-empty] <dir> ...]
[-setfacl [-R] [{-b|-k} {-m|-x <acl_spec>} <path>][--set <acl_spec> <path>]]
[-setfattr {-n name [-v value] | -x name} <path>]
[-setrep [-R] [-w] <rep> <path> ...]
[-stat [format] <path> ...]
[-tail [-f] [-s <sleep interval>] <file>]
[-test [-defswrz] <path>]
[-text [-ignoreCrc] <src> ...]
[-touch [-a] [-m] [-t TIMESTAMP (yyyyMMdd:HHmmss) ] [-c] <path> ...]
[-touchz <path> ...]
[-truncate [-w] <length> <path> ...]
[-usage [cmd ...]]

Generic options supported are:
    -conf <configuration file>          specify an application configuration file
    -D <property=value>                  define a value for a given property
    -fs <file:///|hdfs://namenode:port> specify default filesystem URL to use, overrides 'fs.defaultFS' property from configurations.
    -jt <local|resourcemanager:port>     specify a ResourceManager
    -files <file1,...>                   specify a comma-separated list of files to be copied to the map reduce cluster
    -libjars <jar1,...>                  specify a comma-separated list of jar files to be included in the classpath
    -archives <archive1,...>             specify a comma-separated list of archives to be unarchived on the compute machines

The general command line syntax is:
command [genericOptions] [commandOptions]

Ayobami@Ubuntu:~/Desktop/temp$
```

```
Oct 21 18:59
Ayobami@Ubuntu: ~/Desktop/temp

Paris, France
Berlin, Germany
London, England
Copenhagen, Denmark
Manchester, England
Ayobami@Ubuntu:~/Desktop/temp$ ls input
cities.csv
Ayobami@Ubuntu:~/Desktop/temp$ hadoop jar $HADOOP_HOME/share/hadoop/mapreduce/hadoop-mapreduce-examples-3.4.0.jar wordcount input output
2024-10-21 18:59:13,714 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2024-10-21 18:59:13,939 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2024-10-21 18:59:13,947 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2024-10-21 18:59:14,196 INFO input.FileInputFormat: Total input files to process : 1
2024-10-21 18:59:14,259 INFO mapreduce.JobSubmitter: number of splits:1
2024-10-21 18:59:14,497 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local774138703_0001
2024-10-21 18:59:14,498 INFO mapreduce.JobSubmitter: Executing with tokens: []
2024-10-21 18:59:14,674 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2024-10-21 18:59:14,684 INFO mapreduce.Job: Running job: job_local774138703_0001
2024-10-21 18:59:14,706 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2024-10-21 18:59:14,728 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitterFactory
2024-10-21 18:59:14,729 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2024-10-21 18:59:14,729 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2024-10-21 18:59:14,729 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
2024-10-21 18:59:14,827 INFO mapred.LocalJobRunner: Waiting for map tasks
2024-10-21 18:59:14,834 INFO mapred.LocalJobRunner: Starting task: attempt_local774138703_0001_m_000000_0
2024-10-21 18:59:14,944 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitterFactory
2024-10-21 18:59:14,944 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2024-10-21 18:59:14,944 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
```

```
Oct 21 19:00
Ayobami@Ubuntu: ~/Desktop/temp

2024-10-21 18:59:15,170 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
2024-10-21 18:59:15,242 INFO mapred.LocalJobRunner:
2024-10-21 18:59:15,243 INFO mapred.MapTask: Starting flush of map output
2024-10-21 18:59:15,243 INFO mapred.MapTask: Spilling map output
2024-10-21 18:59:15,243 INFO mapred.MapTask: bufstart = 0; bufend = 126; bufvoid = 104857600
2024-10-21 18:59:15,243 INFO mapred.MapTask: kvstart = 26214396(104857584); kvend = 26214360(104857440); length = 37/6553600
2024-10-21 18:59:15,261 INFO mapred.MapTask: Finished spill 0
2024-10-21 18:59:15,318 INFO mapred.Task: Task:attempt_local774138703_0001_m_000000_0 is done. And is in the process of committing
2024-10-21 18:59:15,325 INFO mapred.LocalJobRunner: map
2024-10-21 18:59:15,326 INFO mapred.Task: Task 'attempt_local774138703_0001_m_000000_0' done.
2024-10-21 18:59:15,345 INFO mapred.Task: Final Counters for attempt_local774138703_0001_m_000000_0: Counters: 18
File System Counters
  FILE: Number of bytes read=281878
  FILE: Number of bytes written=997479
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
Map-Reduce Framework
  Map input records=5
  Map output records=10
  Map output bytes=126
  Map output materialized bytes=138
  Input split bytes=113
  Combine input records=10
  Combine output records=9
  Spilled Records=9
  Failed Shuffles=0
  Merged Map outputs=0
  GC time elapsed (ms)=24
Total committed bytes=15214360
```



```
Oct 21 19:00
Ayobami@Ubuntu: ~/Desktop/temp

2024-10-21 18:59:15,345 INFO mapred.LocalJobRunner: map task executor complete.
2024-10-21 18:59:15,347 INFO mapred.LocalJobRunner: Starting task: attempt_local774138703_0001_r_000000_0
2024-10-21 18:59:15,347 INFO mapred.LocalJobRunner: Waiting for reduce tasks
2024-10-21 18:59:15,363 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitterFactory
2024-10-21 18:59:15,363 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2024-10-21 18:59:15,363 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2024-10-21 18:59:15,364 INFO mapred.Task: Using ResourceCalculatorProcessTree : [ ]
2024-10-21 18:59:15,386 INFO mapred.ReduceTask: Using ShuffleConsumerPlugin: org.apache.hadoop.mapreduce.task.reduce.ShuffleConsumer
2024-10-21 18:59:15,396 WARN impl.MetricsSystemImpl: JobTracker metrics system already initialized!
2024-10-21 18:59:15,449 INFO reduce.MergeManagerImpl: MergerManager: memoryLimit=350532384, maxSingleShuffleLimit=87633096, mergeThreshold=231351376, ioSortFactor=10, memToMemMergeOutputsThreshold=10
2024-10-21 18:59:15,458 INFO reduce.EventFetcher: attempt_local774138703_0001_r_000000_0 Thread started: EventFetcher for fetching Map Completion Events
2024-10-21 18:59:15,500 INFO reduce.LocalFetcher: localfetcher#1 about to shuffle output of map attempt_local774138703_0001_m_000000_0 decomp: 134 len: 138 to MEMORY
2024-10-21 18:59:15,514 INFO reduce.InMemoryMapOutput: Read 134 bytes from map-output for attempt_local774138703_0001_m_000000_0
2024-10-21 18:59:15,521 INFO reduce.MergeManagerImpl: closeInMemoryFile -> map-output of size: 134, inMemoryMapOutputs.size() -> 1, commitMemory -> 0, usedMemory -> 134
2024-10-21 18:59:15,522 INFO reduce.EventFetcher: EventFetcher is interrupted.. Returning
2024-10-21 18:59:15,523 INFO mapred.LocalJobRunner: 1 / 1 copied.
2024-10-21 18:59:15,524 INFO reduce.MergeManagerImpl: finalMerge called with 1 in-memory map-outputs and 0 on-disk map-outputs
2024-10-21 18:59:15,542 INFO mapred.Merger: Merging 1 sorted segments
2024-10-21 18:59:15,542 INFO mapred.Merger: Down to the last merge-pass, with 1 segments left of total size: 124 bytes
2024-10-21 18:59:15,544 INFO reduce.MergeManagerImpl: Merged 1 segments, 134 bytes to disk to satisfy reduce memory limit
2024-10-21 18:59:15,551 INFO reduce.MergeManagerImpl: Merging 1 files, 138 bytes from disk
2024-10-21 18:59:15,551 INFO reduce.MergeManagerImpl: Merging 0 segments, 0 bytes from memory into reduce
2024-10-21 18:59:15,551 INFO mapred.Merger: Merging 1 sorted segments
```

```
Oct 21 19:00
Ayobami@Ubuntu: ~/Desktop/temp

Map-Reduce Framework
  Combine input records=0
  Combine output records=0
  Reduce input groups=9
  Reduce shuffle bytes=138
  Reduce input records=9
  Reduce output records=9
  Spilled Records=9
  Shuffled Maps =1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=0
  Total committed heap usage (bytes)=151531520

Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0

File Output Format Counters
  Bytes Written=108
2024-10-21 18:59:15,607 INFO mapred.LocalJobRunner: Finishing task: attempt_local774138703_0001_r_000000_0
2024-10-21 18:59:15,608 INFO mapred.LocalJobRunner: reduce task executor complete.
2024-10-21 18:59:15,727 INFO mapreduce.Job: Job job_local774138703_0001 running in uber mode : false
2024-10-21 18:59:15,736 INFO mapreduce.Job: map 100% reduce 100%
2024-10-21 18:59:15,738 INFO mapreduce.Job: Job job_local774138703_0001 completed successfully
2024-10-21 18:59:15,767 INFO mapreduce.Job: Counters: 30

File System Counters
  FILE: Number of bytes read=564064
  FILE: Number of bytes written=1995204
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
```

```
Oct 21 19:00
Ayobami@Ubuntu: ~/Desktop/temp

FILE: Number of large read operations=0
FILE: Number of write operations=0
Map-Reduce Framework
  Map input records=5
  Map output records=10
  Map output bytes=126
  Map output materialized bytes=138
  Input split bytes=113
  Combine input records=10
  Combine output records=9
  Reduce input groups=9
  Reduce shuffle bytes=138
  Reduce input records=9
  Reduce output records=9
  Spilled Records=18
  Shuffled Maps =1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=24
  Total committed heap usage (bytes)=303063040
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
File Input Format Counters
  Bytes Read=102
File Output Format Counters
  Bytes Written=108
Ayobami@Ubuntu:~/Desktop/temp$
```

```
Oct 21 19:02
Ayobami@Ubuntu: ~/Desktop/temp

  Shuffled Maps =1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=24
  Total committed heap usage (bytes)=303063040
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
File Input Format Counters
  Bytes Read=102
File Output Format Counters
  Bytes Written=108
Ayobami@Ubuntu:~/Desktop/temp$ hdfs dfs -ls output
Found 2 items
-rw-r--r--  1 Ayobami Ayobami          0 2024-10-21 18:59 output/_SUCCESS
-rw-r--r--  1 Ayobami Ayobami       96 2024-10-21 18:59 output/part-r-00000
Ayobami@Ubuntu:~/Desktop/temp$ hdfs dfs -getmerge output /tmp/result
Ayobami@Ubuntu:~/Desktop/temp$ cat /tmp/result
Berlin, 1
Copenhagen, 1
Denmark 1
England 2
France 1
Germany 1
London, 1
Manchester, 1
Paris, 1
Ayobami@Ubuntu:~/Desktop/temp$
```

SPARK



APPENDIX C

CODE SYNTAX

```
1. # **ANALYSIS OF MACHINE LEARNING MODELS FOR AIRLINE PASSENGER'S
   SATISFACTION USING PYSPARK**
2.
3. ---
4. ---
5. We wanted to understand what factors make airline passengers more or
   less satisfied with their flight
6. experience. By studying this, airlines could improve services and
   prioritize what passengers care most about.
7.
8. To do this, we employed four machine learning models (logistic
   regression, Support Vector Machine,
9. Random Forest and Gradient Boosted Tree) to help predict whether a
   passenger is satisfied or not based
10. on other related information such as the comfort of the seats,
   inflight Wi-Fi service, distance
11. traveled, etc. We were also able to determine which factors
   are most important for passenger satisfaction.
12.
13.
14.
15.     #Importing Data Exploratory libraries
16.     import pandas as pd
17.     import numpy as np
18.     import matplotlib.pyplot as plt
19.     import seaborn as sns
20.
21.     #Loading the dataset
22.
23.     data =
pd.read_csv('/content/sample_data/airline_passenger_satisfaction.csv
')
24.     data.head()
25.
26.     #Checking the number of rows and columns
27.
28.     data.shape
29.
30.     #Counts of the label
31.     data['Satisfaction'].value_counts()
32.
33.     #Checking information about the dataset
34.
35.     data.info()
36.
```

```

37.     #Checking the columns with null values
38.
39.     data.isnull().sum()
40.
41.     #Drop the null values since they are less than 10% of the data
42.
43.     data.dropna(inplace=True)
44.     data.shape
45.
46.     data['Satisfaction'].value_counts()
47.
48.     #Description statistics of the dataset
49.
50.     data.describe()
51.
52.     #Checking the number of unique values each column has
53.
54.     data.nunique()
55.
56.     #Dropping outlier data
57.
58.     arr_delay_outlier = data['Arrival Delay'] < 650
59.     dep_delay_outlier = data['Departure Delay'] < 650
60.     flight_dist_outlier = data['Flight Distance'] < 4000
61.     Data = data[arr_delay_outlier & dep_delay_outlier &
62.     flight_dist_outlier]
63.     Data.shape
64.
65.     **PYSPARK ENVIRONMENT**
66.
67.     #Install pyspark environment
68.     !pip install pyspark
69.
70.     #Importing the required spark libraries
71.     from pyspark.sql import SparkSession
72.     from pyspark.ml.linalg import Vectors
73.     from pyspark.ml.feature import VectorAssembler, StringIndexer
74.     from pyspark.ml.classification import LogisticRegression,
75.     LinearSVC, RandomForestClassifier, GBTCClassifier
76.     from pyspark.ml.evaluation import
77.     MulticlassClassificationEvaluator
78.     from pyspark.mllib.evaluation import MulticlassMetrics
79.     from pyspark.ml.tuning import CrossValidator, ParamGridBuilder
80.
81.     #Initialising the Spark Session
82.     spark =
83.     SparkSession.builder.appName('AirlineSatisfaction_RF').getOrCreate()
84.
85.     #Load the RF dataset

```

```

82.
83.     New_Data = spark.createDataFrame(Data, verifySchema=True)
84.     New_Data.printSchema()
85.
86.     #Indexing categorical columns
87.
88.     #Satisfaction column
89.     satisfaction_indexer = StringIndexer(inputCol='Satisfaction',
90.     outputCol='label')
91.     New_Data = satisfaction_indexer.fit(New_Data).transform(New_Data)
92.
93.     #Gender column
94.     gender_indexer = StringIndexer(inputCol='Gender',
95.     outputCol='Gender_Index')
96.     New_Data = gender_indexer.fit(New_Data).transform(New_Data)
97.
98.     #Customer Type column
99.     customer_type_indexer = StringIndexer(inputCol='Customer
100.     Type', outputCol='Customer_Type_Index')
101.     New_Data = customer_type_indexer.fit(New_Data).transform(New_Data)
102.
103.     #Type of Travel column
104.     type_of_travel_indexer = StringIndexer(inputCol='Type of
105.     Travel', outputCol='Type_of_Travel_Index')
106.     New_Data = type_of_travel_indexer.fit(New_Data).transform(New_Data)
107.
108.     #Class column
109.     class_indexer = StringIndexer(inputCol='Class',
110.     outputCol='Class_Index')
111.     New_Data = class_indexer.fit(New_Data).transform(New_Data)
112.
113.     #Departure and Arrival Time Convenience column
114.     Arrival_Depart_Conv_indexer = StringIndexer(inputCol='Departure and Arrival Time Convenience',
115.     outputCol='Arr_Dept_Time_Cov_Index')
116.     New_Data = Arrival_Depart_Conv_indexer.fit(New_Data).transform(New_Data)
117.
118.     #Ease of Online Booking column
119.     Ease_of_Online_Booking_indexer = StringIndexer(inputCol='Ease
120.     of Online Booking', outputCol='Ease_of_Online_Booking_Index')
121.     New_Data = Ease_of_Online_Booking_indexer.fit(New_Data).transform(New_Data)
122.
123.     #Check-in Service column
124.     Check_in_Service_indexer = StringIndexer(inputCol='Check-in
125.     Service', outputCol='Check_in_Service_Index')

```

```

117.     New_Data                                                    =
        Check_in_Service_indexer.fit(New_Data).transform(New_Data)
118.
119.     #Online Boarding column
120.     Online_Boarding_indexer = StringIndexer(inputCol='Online
        Boarding', outputCol='Online_Boarding_Index')
121.     New_Data                                                    =
        Online_Boarding_indexer.fit(New_Data).transform(New_Data)
122.
123.     #Gate Location column
124.     Gate_Location_indexer = StringIndexer(inputCol='Gate
        Location', outputCol='Gate_Location_Index')
125.     New_Data                                                    =
        Gate_Location_indexer.fit(New_Data).transform(New_Data)
126.
127.     #On-board Service column
128.     Onboard_Service_indexer = StringIndexer(inputCol='On-board
        Service', outputCol='Onboard_Service_Index')
129.     New_Data                                                    =
        Onboard_Service_indexer.fit(New_Data).transform(New_Data)
130.
131.     #Seat Comfort column
132.     Seat_Comfort_indexer = StringIndexer(inputCol='Seat Comfort',
        outputCol='Seat_Comfort_Index')
133.     New_Data                                                    =
        Seat_Comfort_indexer.fit(New_Data).transform(New_Data)
134.
135.     #Leg Rom Service column
136.     Leg_Room_Service_indexer = StringIndexer(inputCol='Leg Room
        Service', outputCol='Leg_Room_Service_Index')
137.     New_Data                                                    =
        Leg_Room_Service_indexer.fit(New_Data).transform(New_Data)
138.
139.     #Cleanliness column
140.     Cleanliness_indexer = StringIndexer(inputCol='Cleanliness',
        outputCol='Cleanliness_Index')
141.     New_Data                                                    =
        Cleanliness_indexer.fit(New_Data).transform(New_Data)
142.
143.     #Food and Drink column
144.     Food_and_Drink_indexer = StringIndexer(inputCol='Food and
        Drink', outputCol='Food_and_Drink_Index')
145.     New_Data                                                    =
        Food_and_Drink_indexer.fit(New_Data).transform(New_Data)
146.
147.     #In-flight Service column
148.     Inflight_Service_indexer = StringIndexer(inputCol='In-flight
        Service', outputCol='Inflight_Service_Index')

```

```

149.     New_Data =
    Inflight_Service_indexer.fit(New_Data).transform(New_Data)
150.
151.     #In-flight Entertainment column
152.     Inflight_Entertainment_indexer = StringIndexer(inputCol='In-
    flight Entertainment', outputCol='Inflight_Entertainment_Index')
153.     New_Data =
    Inflight_Entertainment_indexer.fit(New_Data).transform(New_Data)
154.
155.     #In-flight Wifi Service column
156.     On_board_Wifi_Service_indexer = StringIndexer(inputCol='In-
    flight Wifi Service', outputCol='On_board_Wifi_Service_Index')
157.     New_Data =
    On_board_Wifi_Service_indexer.fit(New_Data).transform(New_Data)
158.
159.     #Baggage Handling column
160.     Baggage_Handling_indexer = StringIndexer(inputCol='Baggage
    Handling', outputCol='Baggage_Handling_Index')
161.     New_Data =
    Baggage_Handling_indexer.fit(New_Data).transform(New_Data)
162.
163.     #A peep into the new columns
164.
165.     New_Data.columns
166.
167.     #Extracting the indexed categorical columns with the numerical
    columns as relevant features
168.
169.     feature_columns = ['Age', 'Flight Distance', 'Departure
    Delay', 'Arrival Delay', 'Gender_Index',
170.                        'Customer_Type_Index',
    'Type_of_Travel_Index', 'Class_Index', 'Arr_Dept_Time_Cov_Index',
171.                        'Ease_of_Online_Booking_Index', 'Check_in_Se
    rvice_Index', 'Online_Boarding_Index',
172.                        'Gate_Location_Index',
    'Onboard_Service_Index', 'Seat_Comfort_Index',
    'Leg_Room_Service_Index',
173.                        'Cleanliness_Index',
    'Food_and_Drink_Index', 'Inflight_Service_Index',
174.                        'Inflight_Entertainment_Index',
    'On_board_Wifi_Service_Index', 'Baggage_Handling_Index']
175.
176.     #A peep into the first 5 rows of the feature columns
177.
178.     New_Data.select(feature_columns).show(5)
179.
180.     #Transforming the feature columns using Vector Assembler
181.

```



```

182.     assembler      =      VectorAssembler(inputCols=feature_columns,
        outputCol="features")
183.     output = assembler.transform(New_Data)
184.     output.select("features").show(5)
185.
186.     #Extracting the transformed features column and label column
        for the classification
187.
188.     final_data = output.select("features", "label")
189.     final_data.printSchema()
190.
191.     #Defining Evaluation Metrics for the Models
192.
193.     accuracy_evaluator                                     =
        MulticlassClassificationEvaluator(labelCol="label",
        predictionCol="prediction", metricName="accuracy")
194.     precision_evaluator                                   =
        MulticlassClassificationEvaluator(labelCol="label",
        predictionCol="prediction", metricName="weightedPrecision")
195.     recall_evaluator                                     =
        MulticlassClassificationEvaluator(labelCol="label",
        predictionCol="prediction", metricName="weightedRecall")
196.     f1_evaluator                                         =
        MulticlassClassificationEvaluator(labelCol="label",
        predictionCol="prediction", metricName="f1")
197.
198.
199.     **LOGISTIC REGRESSION MODEL**
200.
201.     #Copying transformed dataset for the for the Logistic
        Regression
202.
203.     lr_data = final_data
204.     lr_data.printSchema()
205.
206.     #Splitting the lr_data into training and test sets
207.
208.     train_data_lr, test_data_lr = lr_data.randomSplit([0.7, 0.3],
        seed=20)
209.
210.     #Initialising the Logistic Regression classifier
211.
212.     lr              =      LogisticRegression(featuresCol="features",
        labelCol="label")
213.
214.     #Trainning the lr model
215.
216.     lr_model = lr.fit(train_data_lr)
217.

```

```

218.     #Predictions with the lr model on the test set
219.
220.     lr_predictions = lr_model.transform(test_data_lr)
221.
222.     #Evaluating the lr model
223.
224.     lr_accuracy = accuracy_evaluator.evaluate(lr_predictions)
225.     lr_weightedPrecision =
precision_evaluator.evaluate(lr_predictions)
226.     lr_weightedRecall = recall_evaluator.evaluate(lr_predictions)
227.     lr_f1 = f1_evaluator.evaluate(lr_predictions)
228.
229.     print(f"Logistic Regression Model Accuracy: {lr_accuracy:
.2f}")
230.     print(f"Logistic Regression Model Weighted Precision:
{lr_weightedPrecision: .2f}")
231.     print(f"Logistic Regression Model Weighted Recall:
{lr_weightedRecall: .2f}")
232.     print(f"Logistic Regression Model f1: {lr_f1: .2f}")
233.
234.     #Preview first three rows of the lr prediction
235.
236.     lr_predicted= lr_predictions.select("label", "prediction",
"features")
237.     lr_predicted.show(3)
238.
239.     #Extract the logistic coefficients
240.
241.     lr_coefficients = lr_model.coefficients.toArray()
242.
243.     #Matching features to their absolute coefficients
244.
245.     lr_feature_importance = list(zip(feature_columns,
abs(lr_coefficients)))
246.
247.     #Sort the feature importances in descending order
248.
249.     sorted_lr_importances = sorted(lr_feature_importance,
key=lambda x: x[1], reverse=True)
250.
251.
252.     lr_features_ranking = pd.DataFrame(sorted_lr_importances,
columns=['Feature', 'LR_coefficient'])
253.     lr_features_ranking
254.
255.     ***Applying Hyperparameter tuning to the Logistic Regression
Classifier***
256.

```

```

257.     #Parameter grid for the logistic regression Hyperparameter
        tuning
258.
259.     lr_paramGrid = (ParamGridBuilder()
260.                     .addGrid(lr.regParam, [0.1, 0.01])
261.                     .addGrid(lr.maxIter, [50, 100])
262.                     .build())
263.
264.     # Defining the Cross-validator for the logistic regression
265.
266.     lr_crossval = CrossValidator(estimator=lr,
267.                                  estimatorParamMaps=lr_paramGrid,
268.                                  evaluator=accuracy_evaluator,
        numFolds=5)
269.
270.     #Cross-validation fitting for the lr model
271.
272.     lr_cv_model = lr_crossval.fit(train_data_lr)
273.
274.     # Extracting lr_cv Best model
275.
276.     lr_best_model = lr_cv_model.bestModel
277.
278.     #Predictions with the SVM CV on the test data
279.
280.     lr_cv_predictions = lr_best_model.transform(test_data_lr)
281.
282.     #Evaluating the lr_cv model
283.
284.     lr_cv_accuracy = accuracy_evaluator.evaluate(lr_cv_predictions)
285.     lr_cv_weightedPrecision = precision_evaluator.evaluate(lr_cv_predictions)
286.     lr_cv_weightedRecall = recall_evaluator.evaluate(lr_cv_predictions)
287.     lr_cv_f1 = f1_evaluator.evaluate(lr_cv_predictions)
288.
289.     print(f"CV      Logistic      Regression      Model      Accuracy:
        {lr_cv_accuracy: .2f}")
290.     print(f"CV Logistic Regression Model Weighted Precision:
        {lr_cv_weightedPrecision: .2f}")
291.     print(f"CV Logistic Regression Model Weighted Recall:
        {lr_cv_weightedRecall: .2f}")
292.     print(f"CV Logistic Regression Model f1: {lr_cv_f1: .2f}")
293.
294.     #Preview first three rows of the lr_Cv prediction
295.
296.     cv_lr_predicted= lr_predictions.select("label", "prediction",
        "features")

```

```

297.     cv_lr_predicted.show(3)
298.
299.     #Extract the logistic coefficients
300.
301.     lr_cv_coefficients = lr_best_model.coefficients.toArray()
302.
303.     #Matching features to their absolute coefficients
304.
305.     lr_cv_feature_importance      =      list(zip(feature_columns,
abs(lr_cv_coefficients)))
306.
307.     #Sort the feature importances in descending order
308.
309.     sorted_lr_cv_importances      =      sorted(lr_cv_feature_importance,
key=lambda x: x[1], reverse=True)
310.
311.     cv_lr_feature_ranking = pd.DataFrame(sorted_lr_cv_importances,
columns=['Feature', 'CV_LR_coefficient'])
312.     cv_lr_feature_ranking
313.
314.     **LINEAR SUPPORT VECTOR MACHINE (SVM) MODEL**
315.
316.     #Copying transformed dataset for the SVM model
317.
318.     svm_data = final_data
319.     svm_data.printSchema()
320.
321.     #Splitting the svm_data into training and test sets
322.
323.     train_data_svm,  test_data_svm  =  svm_data.randomSplit([0.7,
0.3], seed=20)
324.
325.     #Initialising the Linear SVM classifier
326.
327.     svm      =  LinearSVC(featuresCol="features",  labelCol="label",
maxIter=100, regParam=0.01)
328.
329.     #Training the SVM model
330.
331.     svm_model = svm.fit(train_data_svm)
332.
333.     #Predictions with the SVM model on the test set
334.
335.     svm_predictions = svm_model.transform(test_data_svm)
336.
337.     #Evaluating the SVM model
338.
339.     svm_accuracy = accuracy_evaluator.evaluate(svm_predictions)

```

```

340.     svm_weightedPrecision                                     =
precision_evaluator.evaluate(svm_predictions)
341.     svm_weightedRecall                                       =
recall_evaluator.evaluate(svm_predictions)
342.     svm_f1 = f1_evaluator.evaluate(svm_predictions)
343.
344.     print(f"SVM Model Accuracy: {svm_accuracy: .2f}")
345.     print(f"SVM Model Weighted Precision: {svm_weightedPrecision:
.2f}")
346.     print(f"SVM Model Weighted Recall: {svm_weightedRecall: .2f}")
347.     print(f"SVM Model f1: {svm_f1: .2f}")
348.
349.     #Preview first three rows of the svm prediction
350.
351.     svm_predicted= svm_predictions.select("label", "prediction",
"features")
352.     svm_predicted.show(3)
353.
354.     #Extract the SVM coefficients
355.
356.     svm_coefficients = svm_model.coefficients.toArray()
357.
358.     #Matching features to their absolute coefficients
359.
360.     svm_feature_importance = list(zip(feature_columns,
abs(svm_coefficients)))
361.
362.     #Sort the feature importances in descending order
363.
364.     sorted_svm_importances = sorted(svm_feature_importance,
key=lambda x: x[1], reverse=True)
365.
366.     svm_feature_ranking = pd.DataFrame(sorted_svm_importances,
columns=['Feature', 'SVM_coefficient'])
367.     svm_feature_ranking
368.
369.     ***Applying hyperparameter tuning to the Linear SVM
Classifier***
370.
371.     #Parameter grid for the SVM Hyperparameter tuning
372.
373.     svm_param_grid = (ParamGridBuilder()
374.                        .addGrid(svm.maxIter, [50, 100, 200])
375.                        .addGrid(svm.regParam, [0.01, 0.1, 1.0])
376.                        .build())
377.
378.     # Defining the Cross-validator for the SVM
379.     svm_crossval = CrossValidator(estimator=svm,
380.                                   estimatorParamMaps=svm_param_grid,

```

```

381.                                     evaluator=accuracy_evaluator,
382.                                     numFolds=3)
383.
384.     #Cross-validation fitting for the svm model
385.
386.     svm_cv_model = svm_crossval.fit(train_data_svm)
387.
388.     # Extracting SVM Best model
389.
390.     svm_best_model = svm_cv_model.bestModel
391.
392.     #Predictions with the SVM CV on the test data
393.
394.     svm_cv_predictions = svm_best_model.transform(test_data_svm)
395.
396.     #Evaluating the cv SVM model
397.
398.     svm_cv_accuracy                                     =
accuracy_evaluator.evaluate(svm_cv_predictions)
399.     svm_cv_weightedPrecision                             =
precision_evaluator.evaluate(svm_cv_predictions)
400.     svm_cv_weightedRecall                               =
recall_evaluator.evaluate(svm_cv_predictions)
401.     svm_cv_f1 = f1_evaluator.evaluate(svm_cv_predictions)
402.
403.     print(f"CV SVM Model Accuracy: {svm_cv_accuracy: .2f}")
404.     print(f"CV          SVM          Model          Weighted          Precision:
{svm_cv_weightedPrecision: .2f}")
405.     print(f"CV SVM Model Weighted Recall: {svm_cv_weightedRecall:
.2f}")
406.     print(f"CV SVM Model f1: {svm_cv_f1: .2f}")
407.
408.     #Preview first three rows of the SVM_cv prediction
409.
410.     cv_svm_predicted=                                svm_predictions.select("label",
"prediction", "features")
411.     cv_svm_predicted.show(3)
412.
413.     #Extract the CV_SVM coefficients
414.
415.     svm_cv_coefficients = svm_best_model.coefficients.toArray()
416.
417.     #Matching features to their absolute coefficients
418.
419.     svm_cv_feature_importance      =      list(zip(feature_columns,
abs(svm_cv_coefficients)))
420.
421.     #Sort the feature importances in descending order
422.

```

```

423.     sorted_svm_cv_importances = sorted(svm_cv_feature_importance,
key=lambda x: x[1], reverse=True)
424.
425.     cv_svm_feature_ranking =
pd.DataFrame(sorted_svm_cv_importances, columns=['Feature',
'CV_SVM_coefficient'])
426.     cv_svm_feature_ranking
427.
428.     **RANDOM FOREST MODEL**
429.
430.     #Copying transformed dataset for the Logistic Regression
431.
432.     rf_data = final_data
433.     rf_data.printSchema()
434.
435.     #Split the data into training and test sets
436.
437.     train_data_rf, test_data_rf = rf_data.randomSplit([0.7, 0.3],
seed=20)
438.
439.     #Initialising Random Forest Model
440.
441.     rf = RandomForestClassifier(featuresCol='features',
labelCol='label', numTrees=100)
442.
443.     #Training the Model
444.
445.     rf_model = rf.fit(train_data_rf)
446.
447.     #The prediction on the test data
448.
449.     rf_predictions = rf_model.transform(test_data_rf)
450.
451.     #Evaluating the random forest model
452.
453.     rf_accuracy = accuracy_evaluator.evaluate(rf_predictions)
454.     rf_weightedPrecision =
precision_evaluator.evaluate(rf_predictions)
455.     rf_weightedRecall = recall_evaluator.evaluate(rf_predictions)
456.     rf_f1 = f1_evaluator.evaluate(rf_predictions)
457.
458.     print(f"Random Forest Model Accuracy: {rf_accuracy: .2f}")
459.     print(f"Random Forest Model Weighted Precision:
{rf_weightedPrecision: .2f}")
460.     print(f"Random Forest Model Weighted Recall:
{rf_weightedRecall: .2f}")
461.     print(f"Random Forest Model f1: {rf_f1: .2f}")
462.
463.     #Preview first three rows of the rf prediction

```

```

464.
465.     rf_predicted= rf_predictions.select("label", "prediction",
      "features")
466.     rf_predicted.show(3)
467.
468.     #Extract the feature importances
469.
470.     rf_importances = rf_model.featureImportances
471.
472.     #Matching features to their importances
473.
474.     rf_feature_importances = list(zip(feature_columns,
      rf_importances))
475.
476.     #Sort the feature importances in descending order
477.
478.     sorted_rf_importances = sorted(rf_feature_importances,
      key=lambda x: x[1], reverse=True)
479.
480.     rf_feature_ranking = pd.DataFrame(sorted_rf_importances,
      columns=['Feature', 'RF_Importance'])
481.     rf_feature_ranking
482.
483.     ***Applying Hyperparameter tuning to the Random Forest***
484.
485.     #Defining the
486.
487.     rf_param_grid = (ParamGridBuilder()
488.                      .addGrid(rf.numTrees, [50, 100, 150])
489.                      .addGrid(rf.maxDepth, [5, 10, 15])
490.                      .addGrid(rf.maxBins, [32, 64])
491.                      .build())
492.
493.     #The cross-validator
494.
495.     rf_crossval = CrossValidator(estimator=rf,
496.                                  estimatorParamMaps=rf_param_grid,
497.                                  evaluator=accuracy_evaluator,
498.                                  numFolds=3)
499.
500.     #Training the model
501.
502.     rf_cv_model = rf_crossval.fit(train_data_rf)
503.
504.     #The best model
505.
506.     best_rf_model = rf_cv_model.bestModel
507.
508.     #Predictions on the test data

```



```

509.
510.     rf_cv_predictions = best_rf_model.transform(test_data_rf)
511.
512.     #Evaluating the cv random forest model
513.
514.     rf_cv_accuracy =
accuracy_evaluator.evaluate(rf_cv_predictions)
515.     rf_cv_weightedPrecision =
precision_evaluator.evaluate(rf_cv_predictions)
516.     rf_cv_weightedRecall =
recall_evaluator.evaluate(rf_cv_predictions)
517.     rf_cv_f1 = f1_evaluator.evaluate(rf_cv_predictions)
518.
519.     print(f"CV RF Model Accuracy: {rf_cv_accuracy: .2f}")
520.     print(f"CV RF Model Weighted Precision:
{rf_cv_weightedPrecision: .2f}")
521.     print(f"CV RF Model Weighted Recall: {rf_cv_weightedRecall:
.2f}")
522.     print(f"CV RF Model f1: {rf_cv_f1: .2f}")
523.
524.     #Preview first three rows of the svm prediction
525.
526.     cv_rf_predicted= rf_cv_predictions.select("label",
"prediction", "features")
527.     cv_rf_predicted.show(3)
528.
529.     #Extract feature importances
530.     cv_rf_feature_importances = best_rf_model.featureImportances
531.
532.     #Matching features to their importances
533.     cv_rf_importances = list(zip(feature_columns,
cv_rf_feature_importances))
534.
535.     #Sort the feature importances in descending order
536.     sorted_cv_rf_importances = sorted(cv_rf_importances,
key=lambda x: x[1], reverse=True)
537.
538.     cv_rf_features_ranking =
pd.DataFrame(sorted_cv_rf_importances, columns=['Feature',
'CV_RF_Importance'])
539.     cv_rf_features_ranking
540.
541.     **GRADIENT-BOOSTED TREE (GBT) CLASSIFIER MODEL**
542.
543.     #Copying transformed dataset for the GBT model
544.
545.     gbt_data = final_data
546.     gbt_data.printSchema()
547.

```

```

548.     #Split the data into training and test sets
549.
550.     train_data_gbt,  test_data_gbt  =  gbt_data.randomSplit([0.7,
551.     0.3], seed=20)
552.     # Initialize Gradient-Boosted Trees classifier
553.
554.     gbt = GBTClassifier(featuresCol="features", labelCol="label",
555.     maxIter=50)
556.     # Training the GBT model
557.
558.     gbt_model = gbt.fit(train_data_gbt)
559.
560.     #Prediction on the test set
561.
562.     gbt_predictions = gbt_model.transform(test_data_gbt)
563.
564.     #Evaluating the GBT model
565.
566.     gbt_accuracy = accuracy_evaluator.evaluate(gbt_predictions)
567.     gbt_weightedPrecision =
568.     precision_evaluator.evaluate(gbt_predictions)
569.     gbt_weightedRecall =
570.     recall_evaluator.evaluate(gbt_predictions)
571.     gbt_f1 = f1_evaluator.evaluate(gbt_predictions)
572.
573.     print(f"Gradient Boost Tree Model Accuracy: {gbt_accuracy:
574.     .2f}")
575.     print(f"Gradient Boost Tree Model Weighted Precision:
576.     {gbt_weightedPrecision: .2f}")
577.     print(f"Gradient Boost Tree Model Weighted Recall:
578.     {gbt_weightedRecall: .2f}")
579.     print(f"Gradient Boost Tree Model f1: {gbt_f1: .2f}")
580.
581.     #Preview first three rows of the GBT prediction
582.
583.     gbt_predicted= gbt_predictions.select("label", "prediction",
584.     "features")
585.     gbt_predicted.show(3)
586.
587.     # Extracting feature importances from GBT model
588.     gbt_importances = gbt_model.featureImportances.toArray()
589.
590.     # Combine with feature names
591.     gbt_feature_importance = list(zip(feature_columns,
592.     gbt_importances))
593.
594.     #Sorting by importances

```

```

588.     gbt_sorted_importance = sorted(gbt_feature_importance,
    key=lambda x: x[1], reverse=True)
589.
590.     gbt_feature_ranking = pd.DataFrame(gbt_sorted_importance,
    columns=['Feature', 'GBT_Importance'])
591.     gbt_feature_ranking
592.
593.     *Hyperparameter tuning for the Gradient Boosted Tree
    Classifier*
594.
595.     # Parameter grid for the GBT hyperparameter tuning
596.
597.     gbt_param_grid = (ParamGridBuilder()
598.                        .addGrid(gbt.maxIter, [10, 50, 100])
599.                        .addGrid(gbt.maxDepth, [3, 5, 7])
600.                        .addGrid(gbt.stepSize, [0.1, 0.2])
601.                        .build())
602.
603.     # Defining the Cross-validator for the GBT
604.
605.     gbt_crossval = CrossValidator(estimator=gbt,
606.                                   estimatorParamMaps=gbt_param_grid,
607.                                   evaluator=accuracy_evaluator,
608.                                   numFolds=3)
609.
610.     #Performing the GBT cross-validation
611.
612.     gbt_cv_model = gbt_crossval.fit(train_data_gbt)
613.
614.     # Extracting GBT Best model
615.
616.     gbt_best_model = gbt_cv_model.bestModel
617.
618.     #Predictions on the test data
619.
620.     gbt_cv_predictions = gbt_best_model.transform(test_data_gbt)
621.
622.     #Evaluating the cv GBT model
623.
624.     gbt_cv_accuracy =
    accuracy_evaluator.evaluate(gbt_cv_predictions)
625.     gbt_cv_weightedPrecision =
    precision_evaluator.evaluate(gbt_cv_predictions)
626.     gbt_cv_weightedRecall =
    recall_evaluator.evaluate(gbt_cv_predictions)
627.     gbt_cv_f1 = f1_evaluator.evaluate(gbt_cv_predictions)
628.
629.     print(f"CV GBT Model Accuracy: {gbt_cv_accuracy: .2f}")

```

```

630.     print(f"CV      GBT      Model      Weighted      Precision:
        {gbt_cv_weightedPrecision: .2f}")
631.     print(f"CV GBT Model Weighted Recall: {gbt_cv_weightedRecall:
        .2f}")
632.     print(f"CV GBT Model f1: {gbt_cv_f1: .2f}")
633.
634.     #Preview first three rows of the GBT prediction
635.
636.     cv_gbt_predicted=                gbt_cv_predictions.select("label",
        "prediction", "features")
637.     cv_gbt_predicted.show(3)
638.
639.     # Extracting feature importances from tuned GBT model
640.     cv_gbt_importances                =
        gbt_best_model.featureImportances.toArray()
641.
642.     # Combine with feature names
643.     cv_gbt_feature_importance        =      list(zip(feature_columns,
        cv_gbt_importances))
644.
645.     #Sorting by importances
646.     cv_gbt_sorted_importance        =      sorted(cv_gbt_feature_importance,
        key=lambda x: x[1], reverse=True)
647.
648.     cv_gbt_feature_ranking            =
        pd.DataFrame(cv_gbt_sorted_importance,                columns=['Feature',
        'GBT_CV_Importance'])
649.     cv_gbt_feature_ranking
650.
651.     features_rank                    =      pd.concat([lr_features_ranking,
        cv_lr_feature_ranking, svm_feature_ranking, cv_svm_feature_ranking,
        rf_feature_ranking, cv_rf_features_ranking, gbt_feature_ranking,
        cv_gbt_feature_ranking], axis=1)
652.     features_rank
653.
654.     Satisfaction = Data['Satisfaction']
655.     logistic_prediction                =
        pd.DataFrame(lr_predicted.select("label",      "prediction").collect(),
        columns=["Label", "lr_prediction"])
656.     tuned_logistic_prediction          =
        pd.DataFrame(lr_cv_predictions.select("prediction").collect(),
        columns=["lr_cv_prediction"])
657.     svm_prediction                    =
        pd.DataFrame(svm_predicted.select("prediction").collect(),
        columns=["svm_prediction"])
658.     tuned_svm_prediction                =
        pd.DataFrame(svm_cv_predictions.select("prediction").collect(),
        columns=["svm_cv_prediction"])

```

```

659.     random_forest_prediction                                =
        pd.DataFrame(rf_predicted.select("prediction").collect(),
                      columns=["rf_prediction"])
660.     tuned_random_forest_prediction                          =
        pd.DataFrame(rf_cv_predictions.select("prediction").collect(),
                      columns=["rf_cv_prediction"])
661.     gbt_prediction                                           =
        pd.DataFrame(gbt_predicted.select("prediction").collect(),
                      columns=["gbt_prediction"])
662.     tuned_gbt_prediction                                     =
        pd.DataFrame(gbt_cv_predictions.select("prediction").collect(),
                      columns=["gbt_cv_prediction"])
663.
664.     all_predictions = pd.concat([Satisfaction,
        logistic_prediction,    tuned_logistic_prediction,    svm_prediction,
        tuned_svm_prediction,    random_forest_prediction,
        tuned_random_forest_prediction,    gbt_prediction,
        tuned_gbt_prediction], axis=1)
665.     all_predictions.head(50)
666.     all_predictions.to_csv('predictions2.csv', index=False)
667.
668.     features_rank.to_csv('features_rank.csv', index=False)
669.     all_predictions.to_csv('predictions.csv', index=False)

```