# For A Better or Worse: Evaluating the Hybrid Feature Selections in Predicting Mobile Network Performance

Azman Ab Malik[a], Noormadinah Allias[b,c*], Mohd Nazri Ismail[d], Roziyani Rawie[e] & Irni Hamiza Hamzah[f]

[a]School of Computer Science,11800 USM Pulau Pinang

[b]Faculty of Information Science and Technology, Multimedia University, Jalan Ayer Keroh Lama, 75450 Melaka, Malaysia.

[c]Faculty of Computer and Mathematical Sciences, Universiti Teknologi MARA, 40450 Shah Alam, Selangor Darul Ehsan

[d]National Defence University of Malaysia, Kem Sungai Besi, 57000 Kuala Lumpur, Malaysia

[e]Computer Engineering Technology, Universiti Kuala Lumpur,  1016, Jln Sultan Ismail, Bandar Wawasan, 50250 Kuala Lumpur, Wilayah Persekutuan Kuala Lumpur

[f]Faculty of Electrical Engineering, Universiti Teknologi MARA, Cawangan Pulau Pinang, Malaysia

*Corresponding author: noormadinah@uitm.edu.my

ABSTRACT

*Today, the proliferation of smart devices and mobile networks, alongside activities like social networking, online gaming, and video streaming, has led to the generation of vast amounts of data. This surge in data consumption has placed significant pressure on mobile service providers to deliver higher data throughput to meet growing demands. As a result, mobile operators require efficient feature selection strategies to optimize throughput while ensuring the effective use of network resources. Feature selection is critical in improving network performance by identifying and prioritizing key parameters that significantly influence throughput. This paper introduces a hybrid feature selection approach that combines mutual information as a filter-based method with Recursive Feature Elimination using an Extra Tree Regressor as a wrapper-based method. The selected features are evaluated using three machine learning algorithms: Extra Tree Regressor, Random Forest, and Extreme Gradient Boosting. Experimental results indicate that the proposed feature selection method, when paired with the Extra Tree Regressor, outperforms both Random Forest and Extreme Gradient Boosting in terms of Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and the R-squared ($R^2$) metric.*

*Keywords:  Hybrid feature selection; filter; wrapper; downlink throughput prediction; mobile network*

## INTRODUCTION

The development of smart devices has equipped people with powerful computing tools that are lightweight enough to fit in the palm of their hands, enabling a variety of uses beyond simple communication (Fourati, Maaloul, and Chaari 2021; Oughton et al. 2021). These events consequently lead towards the advancement of mobile communications, resulting in the birth of the Fourth, Fifth, and Sixth Generations of mobile networks (4G, 5G, and 6G) (Oughton et al. 2021; Kshirsagar et al. 2022). As a result, mobile users can enjoy an enormous amount of entertainment at their fingertips, such as playing online games, streaming towards high definition of videos,

interacting with social media platforms anytime and everywhere with network coverage (Kshirsagar et al. 2022). Conversely, to have a smooth delivery of the activities mentioned above, the network throughput plays a major role.

These throughputs, however, have a strong correlation with the success rate of data transmission over a communication channel between a source and a destination (El-Saleh et al. 2022; Supriya 2022). This is because it was a critical factor in defining the user experience in mobile networks (Imoize et al. 2020; Walelgne et al. 2018; Jha & Vijayarajan 2020). Faster download times, more responsive online interactions, and better streaming are all made

possible by higher throughput. As for mobile service providers, ensuring a good Quality of Services with high throughput is very critical in maximizing the user experiences and satisfactions. As a result, many works have been done including implementing predicting methodology.

However, throughput prediction is not trivial as it is dependent on not only the received signal strength, but also other measures including context information such as geographical location, activities performed and the mobility (Al-Thaedan et al. 2023). Therefore, it is important to identify the most important features that can lead towards a good prediction result. Furthermore, by removing the irrelevant features can reduce the difficulties of model learning task (Zhao & Liu 2023). Besides, by retaining the importance features, it will leads towards a more intuitive understanding of the underlying patterns in the dataset (Zhao & Liu 2023). Focusing on the wireless mobile networks, many methods have been proposed in selecting the important features that can contribute towards a high prediction of throughput result, including Pearson Correlation (Mostafa et al. 2022), Deep Learning (Lee & Lee 2021), and Random Forest (Raca et al. 2019; Zhao & Liu 2023).

Meanwhile some of the researchers did not discuss the decision of selecting the important features in their study (Minovski et al. 2021; Eliviani & Bandung 2022). Unfortunately, the Pearson Correlation can be highly sensitive to outliers, which can skew the correlation values and affect feature selection. Meanwhile, the utilization of Deep Learning as feature selection is typically more computationally expensive. In the meantime, the Random Forest is prone to overfitting on the training data, especially if the forest is very large or the dataset is noisy.

Hence, based on the works above, our contributions highlighted the implementation of a hybrid feature selection to predict the downlink throughput of mobile networks. We proposed a hybrid feature selection that utilizes mutual information as a filter based and machine learning models as the wrapper-based feature selection.

## LITERATURE SURVEY

This section will discuss previous research works that used feature selection in mobile wireless networks. Feature selection has become one of the critical aspects of enhancing the performance of wireless mobile networks. Previously, various studies have been conducted for selecting relevant features to improve network efficiency, especially the throughput. Basically,

feature selection methods can be classified into three groups, namely, filter, wrapper and embedded (Manikandan & Abirami 2021). Filter based method select the features independently by performing scoring criterion for each of the features for target variables (Bommert et al. 2020; Tadist et al. 2019). Mutual Information, Pearson correlation, and Information Gain are among the most common filter based feature selection (Shen & Xu 2022; Gong et al. 2024).

Instead of filter-based method, the wrapper method evaluates a subset of features based on the accuracy of a prediction model trained with them. The wrapper technique finds a subset of features by using a classifier and learning methods. The learning model in these approaches is responsible for identifying a subset of candidate features by exploring the space of primary features. It then assesses the performance of the selected subset using a classifier. Examples of wrapper methods include particle swarm optimization, simulated annealing, and genetic algorithms (Naik & Kiran 2021).

In the embedded methodology, the selection of features is rendered a fundamental component of the model development process. The procedure for feature selection is assimilated within the training of the classifier, thereby facilitating the model to concurrently acquire knowledge and discern the most advantageous subset of features. This dual functionality permits the processes of feature selection and model training to transpire simultaneously, thereby enhancing both efficiency and overall performance. (Alyasiri et al. 2022; Rostami et al. 2021).

Meanwhile, Table 1 displayed the pass studies that have been conducted on the mobile wireless network based on the feature selection classes. From the table, there are various types of feature selection methods selected by the researchers based on classes, and each of them have their own strength and limitations. Each of the methods has their own unique advantages and limitations that influenced their effectiveness. For example, Pearson correlation was simple to compute and interpret, however, it only became effective with the linear dataset. As a result, it failed to capture the non-linear associations, especially in a complex dataset (Murari et al. 2020).

In addition, Mutual Information able to capture the relationship between linear and non-linear relationship, thus providing a more comprehensive measure of dependency (Murari et al. 2020). As for the Information gain, this method will works well on the complex dataset, especially in genetics, but it can be computationally expensive and may not perform well with sparse data (Fan et al. 2011). In a wrapper class, Sequential Backward Selection (SBS) and Sequential

TABLE 1. Past surveys of feature selection in wireless mobile network

| No | Feature selection classes | Past surveys of feature selection in wireless mobile network |
|---|---|---|
| 1 | Filter | Pearson Correlation (Mostafa et al. 2022); Chi-Square (Dangi & Lalwani, 2023); Mutual Information (Dangi & Lalwani, 2023); Information Gain (Izadi et al. 2023) |
| 2 | Wrapper | Sequential Backward Selection, Sequential forward floating selection , Sequential backward floating selection (Pasyuk et al. 2019); Grey Wolf Optimisation (Sagar & Saidireddy, 2023) |
| 3 | Embedded | Random Forest (Raca et al. 2019; Zhao & Liu 2023; Boruta et al. 2023); Support Vector Machine (Izadi et al. 2023) |

Forward Floating Selection (SFFS) are simple and intuitive. In addition, both methods are easy to implement and understand. They effectively reduce dimensionality by systematically adding or removing features based on their contribution to model performance. However, these methods can be computationally intensive and may not work well for highly correlated features, resulting in sub-optimal feature subsets.

As for the Grey Wolf Optimisation (GWO), it exhibits fast convergence and requires fewer parameters, making it efficient for high-dimensional feature selection. (Wang et al. 2022; Kurniadi et al. 2023). Despite its advantages, GWO can also stagnate in local optima, particularly in complex datasets (Wang et al. 2022). Additionally, in Embedded classes, the Random Forest can be considered robust with the overfitting, but can be computationally expensive, especially if having many trees and applied on a large dataset, thus leads towards a longer training time, thus make it less practical. For the Support Vector Machine (SVM), the SVM with non-linear kernels can be complex and computationally intensive, making it harder to scale to large datasets.

# METHODOLOGY

## DATASET

As shown in Figure 1, the experiment was conducted using an open production dataset created by Raca et al. (Raca et al. 2018). The dataset consists of five mobility patterns, which is the static, pedestrian, bus, car, and train. However, this experiment used a train dataset which consisted of 38976 samples with 17 features. The train dataset consists of the route taken between Cork to Dublin for 240km and from Cork to Farranfore for about 75km.

## DATA CLEANING

Inside the dataset, there are some irrelevant columns that are not needed for prediction. For example, the 'Timestamp', 'CellID', and the 'Operatorname'. As a result, all the three features above have been dropped from the dataset. Meanwhile, the dataset also consists of the placeholder value ('-') value. The values have been replaced with the NaN for numerical columns and then these missing values were imputed by using the median.

## FEATURE ENGINEERING

The features used in this experiment consist of the categorical and numerical formats. Therefore, the OneHotEncoding is used to transform the categorical variables of 'NetworkMode' and 'State' into a format that is suitable for modelling. Meanwhile, the numerical features are standardised by using StandardScaler to ensure they contribute equally to the model performance.

## DATASET SPLITTING

The dataset is split into a set of features (X), and the target variable (y), which is the DL_bitrate. It is further divided into 70% of the training set and 30% of testing sets to validate the model performances.

## HYBRID FEATURE SELECTIONS

In this hybrid feature selection, the features are first selected from the original features of the dataset. The Mutual Information has been used to perform the action by dynamically selecting top relevant features to the target variables via SelectKBest method. After finding the $k$ best features, the Recursive Feature Elimination with Cross-Validation (RFECV) is used to further refine the feature selection. These RFECV incorporated the machine learning models; namely Random Forest (RF), Extreme Gradient Boosting (XGBoost) and Extra Tree Regressor, that later iteratively removes features and evaluates model performance using a cross-validated framework to remove the features that are less useful. This will result in future reductions in the number of features that have been evaluated based on their importance through the training process.

## MODEL EVALUATIONS

To assess the model's performance, three distinct machine learning models were employed, specifically the Extra Tree Regressor, Random Forest, and Support Vector Machine, for the purpose of predicting the DL bitrate on the testing dataset. The hyperparameters utilized for each model in this study were established based on their default values. Subsequently, the evaluation of the models was conducted using the Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and Coefficient of Determination ($R^2$) as they represent the most commonly employed metrics for evaluating model performance (Hodson, 2022).

## RESULTS AND DISCUSSION

### PERFORMANCE COMPARISON WITH CLASSICAL FEATURE SELECTIONS

Table 1 illustrates the comparative performance of the proposed method against conventional feature selection methods. The findings presented in the Table 2 illustrates the comparative performance of the proposed method against conventional feature selection methods. The findings presented in the table show that all models exhibit competitive performance metrics in terms of mean absolute error (MAE), with the wrapper method achieving a value of 0.0555, while the filter and hybrid methods recorded a value of 0.0556. Nevertheless, the wrapper method for feature selection showed the least favorable result in terms of Root Mean Square Error (RMSE); 0.1940, compared to the filter and hybrid methods. Whereas both filter and hybrid methods consistently achieved 0.1931. The sub-optimal performance of the wrapper can be attributed to its inherent nature, which evaluates subsets of features using the learning model. Consequently, this may result in the inclusion of certain irrelevant features that do not improve the overall performance of the model (Sahu & Dash 2023).

In contrast, the filter method used the statistical test to eliminate the irrelevant features before the model training (Thakkar & Lohiya 2023), hence enhance the performances. This approach has also been implemented in our hybrid model, which combines the strength of filter and wrapper. In our hybrid model, the filter was able to efficiently select the relevant features, whereby the wrapper method acted by fine tuning the selection, resulted towards the most informative and discriminative feature set, thus leading towards the improvement of the RMSE.
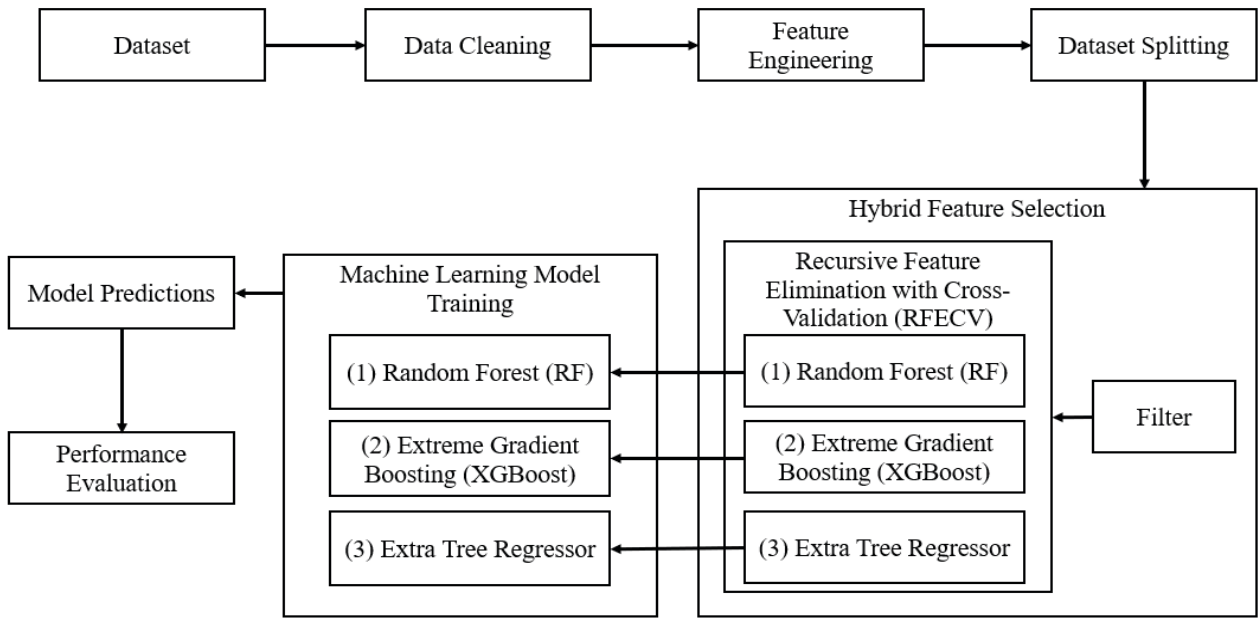


FIGURE 1. A proposed hybrid feature selection

In Table 3, the hybrid methodology employing the Extra Tree Regressor demonstrated superior performance relative to both the Random Forest and Extreme Gradient Boosting algorithms as indicated by the metrics of Mean Absolute Error (MAE) and Root Mean Square Error (RMSE). The hybrid approach utilizing the Extra Tree Regressor surpassed the Random Forest and Extreme Gradient Boosting techniques in terms of both MAE and RMSE metrics.

As illustrated in the table, the Extra Tree Regressor recorded an MAE of 0.0556 and an RMSE of 0.1931. Conversely, the hybrid approach utilizing Random Forest

attained a commendable second position, yielding an MAE of 0.0568 and an RMSE of 0.1967, trailed by the Extreme Gradient Boosting method which produced an MAE of 0.0622 and an RMSE of 0.1983.

The superior performances of Extra Tree Regressor were influenced by the ability of the model to split the nodes by choosing cut points fully at random (Geurts et al. 2006; Saeed et al. 2021).This will allow the model to choose split points without any specific criterion or optimize the process. This randomness in node splitting can lead to diverse and uncorrelated trees. In addition, as Extra Tree Regressor uses the whole learning sample rather than a bootstrap replica to grow the trees (Geurts et al. 2006), it enables each of the decision trees in the ensemble to be trained on the full dataset without any resampling or bootstrapping. Therefore, the model can capture more information, thus leading towards better generalisation and predictive accuracy.

TABLE 2. A comparison table between the proposed hybrid feature selection and classical feature selection

| No | Feature selection | Performance Evaluation | |
|----|-------------------|------|------|
| | | MAE | RMSE |
| 1 | Filter (Mutual Information) | 0.0556 | 0.1931 |
| 2 | Wrapper (Extra Tree Regressor) | 0.0555 | 0.1940 |
| 3 | Hybrid (Filter + RFECV (Extra Tree Regressor)) | 0.0556 | 0.1931 |

TABLE 3. A comparison table between the proposed hybrid feature selection and different hybrid models

| No | Feature selection | Performance Evaluation | |
|----|-------------------|------|------|
| | | MAE | RMSE |
| 1 | Hybrid (Filter + RFECV (Random Forest)) | 0.0568 | 0.1967 |
| 2 | Hybrid (Filter + RFECV (Extreme Gradient Boosting)) | 0.0622 | 0.1983 |
| 3 | Hybrid (Filter + RFECV (Extra Tree Regressor)) | 0.0556 | 0.1931 |

## FEATURES ANALYSIS

The results obtained in Table 2 and Table 3 are basically have been influenced by the features chosen by each of the methods as displayed in Table 4. Table 4 displayed the selected features that were determined using the classical and hybrid feature selection models. The data illustrates that while all models predominantly identified a comparable set of features, each model exhibited a tendency to select slightly distinct subsets of these features. For instance, the filter-based feature selection method identified nearly equivalent features but with a reduced quantity in comparison to the wrapper-based feature selection approach. The variation in the identified features stems from the fact that the wrapper method selected. The variation in the identified features stems from the fact that the wrapper method selected; 'NRxRSRP','NRxRSRQ', and 'NetworkMode_LTE', whereas the filter method opted for 'NetworkMode_HSPA+' and 'NetworkMode_LTE'.

In the meantime, the hybrid methodologies employing Random Forest and Extreme Gradient Boosting have identified the minimal set of features, comprising twelve distinct attributes. The differentiation between these two methodologies lies in the fact that the hybrid approach utilizing Random Forest selected the Received Signal Strength Indicator (RSSI), whereas the hybrid technique incorporating Extreme Gradient Boosting opted for State_I as its unique attribute. Nonetheless, the situations differ between the filter method and the proposed hybrid feature selection approaches, because both models identify the same features. The similar features were chosen because of their strong association with the target variable, emphasizing not only their importance in result prediction, but also showed that the proposed method is fused to achieve feature subset with high quality.

TABLE 4. Features selected based on feature selection models

| Feature selection | Number of features selected | Features selected |
|---|---|---|
| Filter (Mutual Information) | 15 | 'Longitude', 'Latitude', 'Speed', 'RSRP', 'RSRQ', 'SNR', 'CQI', 'RSSI', 'UL_bitrate', 'ServingCell_Lon', 'ServingCell_Lat', 'ServingCell_Distance', 'NetworkMode_HSPA+', 'NetworkMode_LTE', 'State_I' |
| Wrapper (Extra Tree Regressor) | 16 | 'Longitude', 'Latitude', 'Speed', 'RSRP', 'RSRQ', 'SNR', 'CQI', 'RSSI', 'UL_bitrate', 'NRxRSRP', 'NRxRSRQ', 'ServingCell_Lon', 'ServingCell_Lat', 'ServingCell_Distance', 'NetworkMode_LTE', 'State_I' |
| Hybrid (Filter + RFECV(Extra Tree Regressor)) | 15 | 'Longitude', 'Latitude', 'Speed', 'RSRP', 'RSRQ', 'SNR', 'CQI', 'RSSI', 'UL_bitrate', 'ServingCell_Lon', 'ServingCell_Lat', 'ServingCell_Distance', 'NetworkMode_HSPA+', 'NetworkMode_LTE', 'State_I' |
| Hybrid (Filter + RFECV (Random Forest)) | 12 | 'Longitude', 'Latitude', 'Speed', 'RSRP', 'RSRQ', 'SNR', 'CQI', 'RSSI', 'UL_bitrate', 'ServingCell_Lon', 'ServingCell_Lat', 'ServingCell_Distance' |
| Hybrid (Filter + RFECV (Extreme Gradient Boosting)) | 12 | 'Longitude', 'Latitude', 'Speed', 'RSRP', 'RSRQ', 'SNR', 'CQI', 'UL_bitrate', 'ServingCell_Lon', 'ServingCell_Lat', 'ServingCell_Distance', 'State_I' |

The SelectKBest evaluates the individual predictive efficacy of each feature, whereas RFECV uses the Extra Tree regressor to evaluate the collective contribution of a subset of features using cross-validation. Both techniques consistently select the same features, implying that these traits are highly relevant to the predictive model and provide essential insights for making exact predictions. Furthermore, the selected features from each of the models play a vital role in predicting the throughput.

For example, the latitude and longitude, are used to define the spatial location of mobile devices. Mobile network operators are now able to understand geographical distributions, which informs infrastructure planning and the development of traffic management techniques, traffic patterns, and network load evaluations (Nema & Jaafar 2020). As a result, this can help mobile network operators achieve network traffic steadiness, hence improving Quality of Service (QoS).

In addition, the RSRP and RSRQ are required to make a migration choice during intercellular movement. These values indicate the quality of communication in dBm and dB, respectively. These variables are crucial in determining the frequency band in which the communication will occur. While the communication is in the mobile state, the RSRP and RSRQ decide which base station and cell to use for data transmission and handover (Kurnaz et al. 2023) . The Channel Quality Information (CQI) is crucial in LTE-A for the base station to decide the corresponding modulation and encoding scheme. In addition, the mobile devices will estimate the CQI to maximize the data rate (Kurnaz et al.

2023). The RSSI is the key metrics in assessing signal quality. Poor RSSI values can lead to a severe issue in service quality and also signal coverage, thus neglecting the mobile user in achieving quality of experience (QoE).

## COMPARISONS BETWEEN COEFFICIENT OF DETERMINATIONS ($R^2$)

Table 2 displayed the result for the Coefficient of Determination ($R^2$). The filter that utilised Mutual Information and the hybrid Extra Tree Regressor achieved the highest $R^2$ with 0.9619, followed by the wrapper that utilized Extra Tree Regressor obtained the $R^2$ value of 0.9616, followed by the hybrid feature selection that utilized Random Forest and Extreme Gradient Boosting, both 0.9605 and 0.9598 respectively. This indicates that approximately 96.19% of the variance in the dependent variable can be explained by the independent variables selected by the hybrid Extra Tree Regressor model and filter-based model. A higher $R^2$ value suggests that the model fits the data well and provides a good representation of the relationship between the independent and dependent variables. Furthermore, according to the authors (Chicco et al. 2021) , they suggested that the $R^2$ is more informative and does not have the interpretability limitations of MAE and also RMSE. Nevertheless, it is a standard metric for evaluating regression analyses in any scientific fields.

In addition, the Random Forest selection method also performs well in capturing the relationship between the independent and dependent variables, albeit slightly lower

than the Extra Tree Regressor. The Extreme Gradient Boosting indicates that approximately 95.98% of the variance in the dependent variable can be explained by the independent variables selected by the XGBoost model. While still high, the $R^2$ value for XGBoost is slightly lower than that of the Extra Tree Regressor and Random Forest Regressor, suggesting a slightly weaker fit to the data. From

the results above, while MAE and RMSE are commonly used metrics for model evaluation, $R^2$ is suggested as a more informative metric for regression analyses, providing insights into the goodness of fit of the model without the interpretability limitations of of other metrics like MAE and RMSE (Chicco et al. 2021).
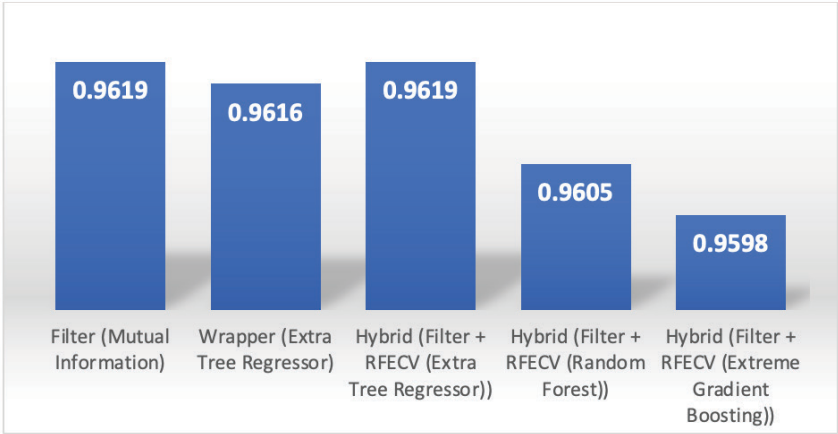


FIGURE 2. A comparison between coefficient of determinations ($R^2$)

## COMPARISONS WITH DIFFERENT DATASET

In Figure 3, we evaluated our hybrid feature selection methodologies utilizing the Fifth Generation (5G) dataset sourced from the authors (Raca et al. 2020). We selected this dataset, which encompasses driving mobility patterns with a total of 25,212 samples that engaged in streaming video content on the Netflix platform. Nevertheless, a comprehensive discussion regarding the dataset can be found in the publications authored by the researchers. All experimental procedures were conducted in accordance with the methodologies delineated in Figure 1. Nonetheless, we excluded the features such as 'Timestamp', 'CellID', 'Operatorname', 'pingavg', 'pingmin', 'pingmax', 'pingstdev', 'pingloss', 'cellhex', 'nodehex', 'lachex', and 'rawcellid', as they were inappropriate to the objectives of our research.

Based on Figure 3, the integration of filter methodologies and Recursive Feature Elimination with Cross-Validation (RFECV) utilizing the Extra Tree Regressor yields the most superior performance across all evaluated metrics. It exhibits the minimal Mean Absolute Error (MAE) and Root Mean Square Error (RMSE), indicating that the predictive outputs of the model are in closest proximity to the actual values. The coefficient of determination, represented by an $R^2$ score of 0.7816, implies that this model accounts for approximately 78.16% of the variance observed in the target variable.

Meanwhile, the implementation of RandomForest in conjunction with the hybrid feature selection approach demonstrates marginally inferior performance compared to the Extra Trees Regressor. Both the Mean Absolute Error (MAE) and the Root Mean Square Error (RMSE) exhibit elevated values, suggesting a reduction in predictive accuracy, while the coefficient of determination ($R^2$) is recorded at 0.7337, accounting for approximately 73.37% of the observed variance. Nevertheless, this model exhibits a relatively commendable performance; however, it is unequivocally surpassed by the Extra Trees methodology.

Ultimately, Extreme Gradient Boosting (XGBoost) employing this feature selection methodology demonstrates the most inferior performance. It has the highest MAE and RMSE, meaning its predictions are the least accurate. The $R^2$ score of 0.6935 shows that it explains about 69.35% of the variance, the lowest among the three models.

## CONCLUSION

In conclusion, feature selection models play a vital role in the identification of essential features that impact the results of throughput prediction. In the present investigation, a hybrid feature selection technique was introduced, which involved a combination of filter and wrapper methods, established on the Extra Tree Regressor model. This strategy delivered promising outcomes in terms of crucial

performance measures like Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and Coefficient of Determination ($R^2$). Moving forward, there exist numerous avenues for potential exploration and enhancement.

Primarily, our objective is to optimize both the predictive models and the hyperparameter values associated with the proposed feature selection strategy as a future work. By adjusting these elements, it is possible to further enhance performance and forecast accuracy. Moreover, our intention is to validate the efficacy of the proposed feature selection technique across various datasets, not focusing only on mobile network dataset. The experimentation using diverse datasets will facilitate the assessment of its adaptability and resilience in varied fields and circumstances. This comprehensive assessment will offer valuable insights into the suitability and efficacy of the feature selection methodology in diverse settings.

As a significant contribution to the field, the application of mutual information in conjunction with Recursive Feature Elimination with Cross-Validation (RFECV) during the feature selection process will guarantee that only the most pertinent and impactful features are preserved, thereby enhancing the overall efficacy of a machine learning model. Given that Mutual Information serves as an effective initial methodology by identifying features that exhibit the highest correlation with the target variable, it ensures that the model is developed using the most informative features available.
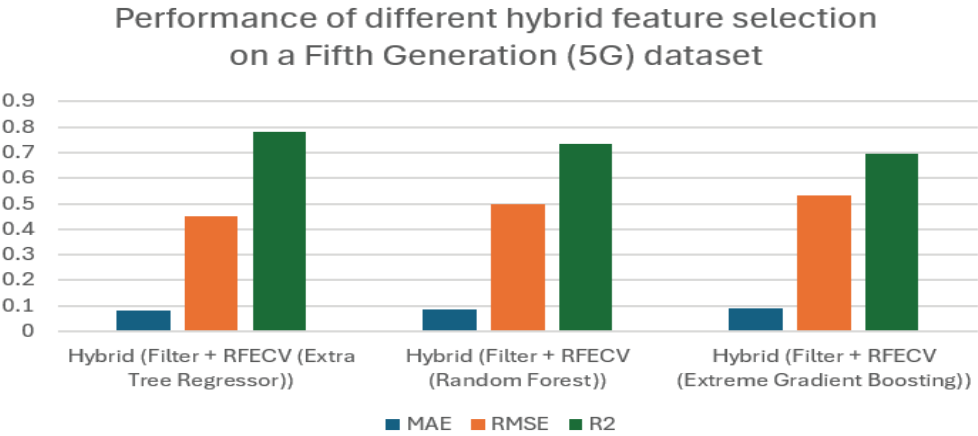


FIGURE 3. Hybrid feature selection models tested on different dataset

## ACKNOWLEDGEMENT

## DECLARATION OF COMPETING INTEREST

None.

## REFERENCES

Al-Thaedan, A., Shakir, Z., Mjhool, A. Y., Alsabah, R., Al-Sabbagh, A., Salah, M., & Zec, J. 2023. Downlink throughput prediction using machine learning models on 4G-LTE networks. *International Journal of Information Technology* 15(6): 2987-2993.

Alyasiri, O. M., Cheah, Y. N., Abasi, A. K., & Al-Janabi, O. M. 2022. Wrapper and hybrid feature selection methods using metaheuristic algorithms for English text classification: A systematic review. *IEEE Access* 10: 39833-39852.

Bommert, A., Sun, X., Bischl, B., Rahnenführer, J., & Lang, M. 2020. Benchmark for filter methods for feature selection in high-dimensional classification data. *Computational Statistics & Data Analysis* 143: 106839.

Chicco, D., Warrens, M. J., & Jurman, G. 2021. The coefficient of determination R-squared is more informative than SMAPE, MAE, MAPE, MSE and RMSE in regression analysis evaluation. *Peerj Computer Science 7*: e623.

Dangi, R., & Lalwani, P. 2023. Feature selection based machine learning models for 5G network slicing approximation. *Computer Networks* 237: 110093.

El-Saleh, A. A., Alhammadi, A., Shayea, I., Alsharif, N., Alzahrani, N. M., Khalaf, O. I., & Aldhyani, T. H. 2022. Measuring and assessing performance of mobile broadband networks and future 5G trends. *Sustainability* 14(2): 829.

Eliviani, R. 2022. Throughput prediction for multimedia IoT in wireless network. Conference Proceedings - *2022 International Symposium on Electronics and Smart Devices (ISESD)* 1-6.

Fan, R., Zhong, M., Wang, S., Zhang, Y., Andrew, A., Karagas, M., Chen, H., Amos, C. I., Xiong, M., & Moore, J. H. 2011. Entropy-based information gain approaches to detect and to characterize gene-gene and gene-environment interactions/correlations of complex diseases. *Genetic Epidemiology* 35(7): 706-721.

Fourati, H., Maaloul, R., & Chaari, L. 2021. A survey of 5G network systems: challenges and machine learning approaches. *International Journal of Machine Learning and Cybernetics* 12(2): 385-431.

Geurts, P., Ernst, D., & Wehenkel, L. 2006. Extremely randomized trees. *Machine Learning* 63: 3-42.

Gong, H., Li, Y., Zhang, J., Zhang, B., & Wang, X. 2024. A new filter feature selection algorithm for classification task by ensembling pearson correlation coefficient and mutual information. *Engineering Applications of Artificial Intelligence* 131: 107865.

Hodson, T. O. 2022. Root mean square error (RMSE) or mean absolute error (MAE): when to use them or not. *Geoscientific Model Development Discussions* 2022: 1-10.

Imoize, A. L., Orolu, K., & Atayero, A. A. A. 2020. Analysis of key performance indicators of a 4G LTE network based on experimental data obtained from a densely populated smart city. *Data in Brief* 29: 105304.

Izadi, S., Mahmood, A., & Rojia, N. 2023. Analysis of feature selection methods for network traffic classification BT - *Proceedings of the 8th International Conference on Advanced Intelligent Systems and Informatics* 65–77.

Jha, H., & Vijayarajan, V. (2020). Mobile internet throughput prediction using machine learning techniques. *Conference Proceedings - 2020 International Conference on Smart Electronics and Communication* 253-257.

Kshirsagar, P. R., Reddy, D. H., Dhingra, M., Dhabliya, D., & Gupta, A. 2022. A Review on comparative study of 4g, 5g and 6g networks. *Proceedings of 5th International Conference on Contemporary Computing and Informatics* 1830-1833.

Kurnaz, C., Kola, A. F., & Esenalp, M. O. 2023. Performance analysis and modeling based on LTE-A field measurements: a city center example. *International Journal of Information Technology* 15(4) 1919-1925.

Kurniadi, F. I., Wulandari, A., & Arifin, S. 2023. Feature selection using grey wolf optimization algorithm on light gradient boosting machine. *Heart* 13: 476.

Lee, D., & Lee, J. 2021. Machine learning and deep learning for throughput prediction. *Proceedings of Twelfth International Conference on Ubiquitous and Future Networks* 452-454.

Manikandan, G., & Abirami, S. 2021. Feature selection is important: State-of-the-art methods and application domains of feature selection on high-dimensional data. *Applications in Ubiquitous Computing* 177-196.

Minovski, D., Ögren, N., Mitra, K., & Åhlund, C. (2021). Throughput prediction using machine learning in LTE and 5G networks. *IEEE Transactions on Mobile Computing* 22(3): 1825-1840.

Mostafa, A., Elattar, M. A., & Ismail, T. 2022. Downlink throughput prediction in LTE cellular networks using time series forecasting. *Conference Proceedings - 2022 International Conference on Broadband Communications for Next Generation Networks and Multimedia Applications* 1-4.

Murari, A., Rossi, R., Lungaroni, M., Gaudio, P. & Gelfusa, M. 2020. Quantifying total influence between variables with information theoretic and machine learning techniques. *Entropy* 22(2): 141.

Naik, D. L., & Kiran, R. 2021. A novel sensitivity-based method for feature selection. *Journal of Big Data* 8(1): 128.

Nema, B. M., & Jaafar, A. N. 2020. Geo location of mobile device. *Recent Trends in Communication Networks* 55.

Oughton, E. J., Lehr, W., Katsaros, K., Selinis, I., Bubley, D., & Kusuma, J. 2021. Revisiting wireless internet connectivity: 5G vs Wi-Fi 6. *Telecommunications Policy* 45(5): 102127.

Supriya, M. 2022. Throughput analysis with effect of dimensionality reduction on 5g dataset using machine learning and deep learning models. *Conference Proceedings - 2022 International Conference on Industry 4.0 Technology* 1-7.

Pasyuk, A., Semenov, E., & Tyuhtyaev, D. 2019. Feature selection in the classification of network traffic flows. Conference Proceedings - *2019 International Multi-Conference on Industrial Engineering and Modern Technologies* 1-5.

Raca, D., Dylan, L., Cormac, J. S., & Jason, J. Q., 2020. MMSys '20. New York, NY, USA: Association for Computing Machinery.

Raca, D., Quinlan, J. J., Zahran, A. H., & Sreenan, C. J. 2018. Beyond throughput: A 4G LTE dataset with channel and context metrics. *Proceedings of 9th ACM Multimedia Systems Conference* 460-465.

Raca, D., Zahran, A. H., Sreenan, C. J., Sinha, R. K., Halepovic, E., Jana, R., Gopalakrishnan, V., Bathula, B., & Varvello, M. 2019. Empowering video players in cellular: Throughput prediction from radio network measurements. *Proceedings of 10th ACM Multimedia Systems Conference* 201-212.

Rostami, M., Berahmand, K., Nasiri, E., & Forouzandeh, S. 2021. Review of swarm intelligence-based feature selection methods. *Engineering Applications of Artificial Intelligence* 100: 104210.

180

Saeed, U., Jan, S. U., Lee, Y. D., & Koo, I. 2021. Fault diagnosis based on extremely randomized trees in wireless sensor networks. *Reliability engineering & System Safety* 205: 107284.

Sagar, D., & Saidireddy, M. 2023. Security measurement in LTE/LTE-A network based on zs-lr feature selection technique and um-tgan attack detection technique. *Expert Systems with Applications* 231: 120703.

Sahu, B., & Dash, S. 2022. Multi-filter wrapper enhanced machine learning model for cancer. *International Conference on Intelligent Systems and Machine Learning* 64-78.

Shen, J., & Xu, F. 2022. Method of fault feature selection and fusion based on poll mode and optimized weighted KPCA for bearings. *Measurement* 194: 110950.

Tadist, K., Najah, S., Nikolov, N. S., Mrabti, F., & Zahi, A. 2019. Feature selection methods and genomic big data: a systematic review. *Journal of Big Data* 6(1): 1-24.

Thakkar, A., & Lohiya, R. 2023. Fusion of statistical importance for feature selection in deep neural network-based intrusion detection system. *Information Fusion* 90: 353-363.

Walelgne, E. A., Manner, J., Bajpai, V., & Ott, J. 2018. Analyzing throughput and stability in cellular networks. *Proceedings of NOMS 2018-2018 IEEE/IFIP Network Operations and Management Symposium* 1-9.

Wang, J., Lin, D., Zhang, Y., & Huang, S. 2022. An adaptively balanced grey wolf optimization algorithm for feature selection on high-dimensional classification. *Engineering Applications of Artificial Intelligence* 114: 105088.

Zhao, X., & Liu, W. (2023, November). Artificial intelligence based feature selection and prediction of downlink IP throughput. *Conference Proceedings - 2023 International Conference on Wireless Communications and Signal Processing* 116-121.

Zhao, X., & Liu, W. 2023. Artificial Intelligence Based Feature Selection and Prediction of Downlink IP Throughput. *Conference Proceedings - 2023 International Conference on Wireless Communications and Signal Processing* 116-121