**Research Article**

# PCA-SWF: A Principal Component Analysis-Based Stacked Model Approach for Weather Forecasting

Shimaila[1], Dr. Sifatullah Siddiqi[2]

[1, 2] *Department of Computer Science, Integral University, Lucknow, India,*

*shimailaphd@gmail.com, sifatullah.siddiqi@gmail.com*

*\*Corresponding Author: Shimaila; Email: shimailaphd@gmail.com*

| ARTICLE INFO | ABSTRACT |
|---|---|
| | Weather forecasting plays a crucial role across various sectors of society, enabling timely predictions of severe weather events such as hurricanes, tornadoes, storms, floods, and heatwaves. Accurate forecasts provide essential data for issuing public warnings, allowing individuals and authorities to take necessary precautions to safeguard lives and property. In agriculture, farmers rely on weather forecasts to optimize irrigation, planting, and harvesting schedules, ensuring efficient water resource management, maximizing crop yields, and minimizing damage caused by extreme weather. Additionally, weather predictions significantly impact transportation systems by providing insights into potential hazards such as ice, snow, poor visibility, wind speeds, and road conditions. This information helps railways, airlines, and shipping industries adjust schedules and ensure safe and efficient transport operations. In this study, we propose a weather prediction approach using a stacked model, which outperforms single classifiers such as Decision Trees, Support Vector Machines, K-Nearest Neighbors, and Multilayer Perceptrons. Principal Component Analysis (PCA) is employed for feature extraction to reduce dimensionality and improve prediction accuracy. The performance of the proposed model is evaluated using metrics such as accuracy, Matthews Correlation Coefficient (MCC), F1-score, and the Receiver Operating Characteristic (ROC) curve, along with the Area Under the Curve (AUC). The results demonstrate the effectiveness of the stacked model in achieving robust and reliable weather forecasting.<br><br> |

## INTRODUCTION

Weather forecasts are crucial for making decisions in a variety of industries, such as transportation, agriculture, and disaster management. Accurate forecasts have the power to differentiate between a successful and unsuccessful agricultural output, the smooth operation of transportation networks, and the effective management of natural disasters. With the increasing unpredictability of climate trends, accurate weather forecasting is more crucial than ever [1]. The incorporation of machine learning into conventional weather forecast techniques has resulted in a revolutionary change in recent years. The integration of sophisticated data analytics and meteorology has the potential to greatly improve the precision and dependability of weather predictions [2]. By breaking down the procedures and giving you the tools to explore the enormous potential of machine learning in meteorology, this paper will walk you through the process of creating a weather prediction model [3]. AI has greatly improved data processing skills, revolutionizing weather forecasting. AI is superior at real-time data integration and analysis because it can handle the vast amounts and variety of meteorological data that traditional approaches find difficult to handle [4] [5]. Artificial Intelligence (AI) can make more precise and timely predictions by using machine learning algorithms to find intricate patterns in both history and present data. More accuracy in forecasting is made possible by meteorologists because to this sophisticated data management, especially for short-term weather changes. The potential of AI to increase the precision of extreme weather forecasts is one of its most important contributions. AI models can anticipate catastrophic weather phenomena like hurricanes, tornadoes, and floods more accurately and with longer lead periods thanks to deep learning techniques and ensemble methodologies [6]

[7]. These models' ability to provide early warnings is essential for disaster planning and mitigation as it has the potential to save lives and reduce financial harm. Further improving forecasting efficacy, AI's pattern recognition skills also reveal minor signs of severe weather that conventional models could overlook. Long-term climate models and customized weather services are made possible by AI [8] [9]. Artificial intelligence (AI)-driven solutions enhance user experience and accessibility to vital meteorological information by offering hyper-local predictions customized to individual requirements. By evaluating vast amounts of historical data and modeling several future scenarios, artificial intelligence (AI) helps climate scientists create longer-term models that are more accurate [10] [11]. These developments highlight the critical role AI plays in both short-term weather forecasting and long-term environmental planning by assisting scientists and policymakers in understanding climate patterns and developing effective policies to mitigate climate change.

## LITERATURE REVIEW

Grönquist et al. [12] propose a mixed model using a subset of weather trajectories combined with deep neural networks for post-processing ensemble forecasts. This approach significantly improves forecast skill, particularly for extreme weather events. Guan and Zhu [13] introduce verification methodologies for extreme weather forecasts, specifically focusing on extreme cold temperatures and precipitation events. They compare ANF and EFI algorithms and their performance in forecasting extreme events. Li et al. [14] review progress in ensemble forecasting of extreme weather based on numerical models. They emphasize the dominance of dynamical models in extreme weather forecasting and discuss approaches for improving ensemble probabilistic forecasts. Xu et al. [15] introduce ExtremeCast, a method for improving extreme value prediction in global weather forecasts using machine learning. They address the issue of biased predictions for extreme events and demonstrate superior performance compared to traditional methods. Bouallègue et al. [16] assess the rise of data-driven weather forecasting using machine learning in an operational context. They compare ML-generated forecasts with standard NWP-based forecasts and highlight the potential of ML methods for improving forecast accuracy. Lopes et al. [17] evaluate daily temperature extremes in ECMWF operational forecasts and ERA5 reanalysis, highlighting improvements in ERA5 accuracy over the past decade. Xu et al. [18] propose an AI-driven regional weather model to improve disastrous extreme precipitation forecasting in North China, demonstrating superior performance over traditional models. Bouallègue et al. [19] introduce PoET, a post-processing approach using hierarchical transformers to improve ensemble weather forecasts, achieving significant skill improvement globally. Beimel et al. [20] focus on improving wind forecasts for sailing events using a combination of numerical modeling and machine learning post-processing, showcasing the potential of ML models to enhance forecast accuracy. Wang and Ikegaya [21] predict annual extreme winds in Iran using numerical weather forecasting, meteorological observation, and statistical models, providing insights into predicting extreme wind speeds. Varshney et al. [22] explore the role of Graph Neural Networks (GNNs) in weather prediction, aiming to enhance prediction accuracy and reliability by capturing complex relationships within meteorological data. Obisesan [23] compares six machine learning models for predicting meteorological variables in a tropical location, identifying Random Forest as the best-performing model. Nguyen et al. [24] introduce ClimateLearn, an open-source library for benchmarking machine learning models in weather and climate modeling, promoting reproducibility and collaboration in the field.

## RESEARCH METHODOLOGY

When investigating the impact of weather prediction towards the environment study, it is crucial to tailor the research philosophy, approach, data collection methods, analysis, and sampling structures accordingly. Here we use Stacking, also known as stacked machine learning, is an ensemble learning method [25][26] that enhances prediction performance by merging many machine learning models. The goal is to combine the advantages of several models or classifiers to produce a strong meta-model that outperforms any individual model or classifier. When dealing with complicated datasets, where no one technique exhibits optimal performance across all characteristics, this approach is quite helpful. In this article, we present a layered machine learning model (Figure 1) that aims to improve task-specific forecasting accuracy.
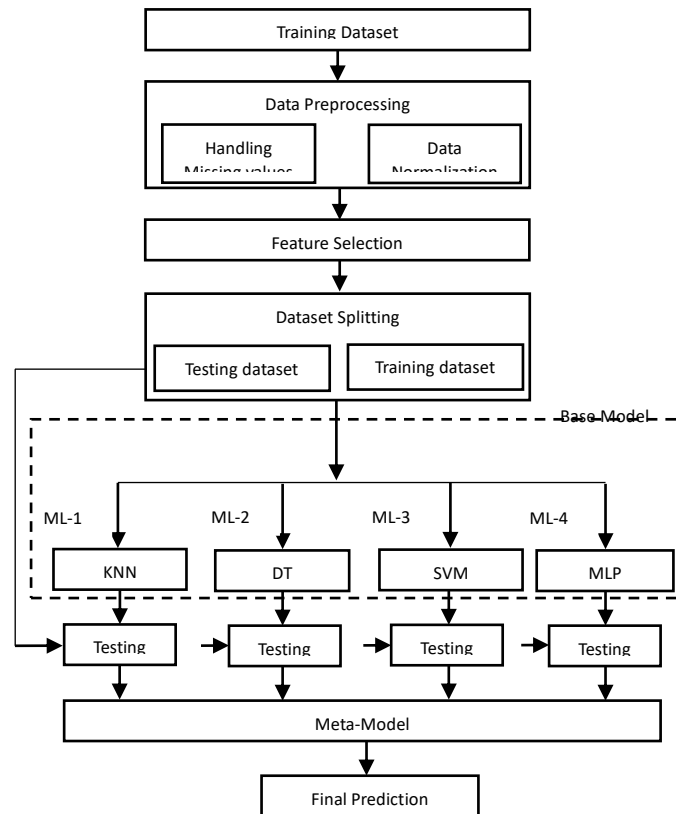
Figure 1: Proposed model work flow

Two layers make up the Stacked model: Base learners and Meta-learner. Below is a summary of every layer:

**Base Learners (Level 0 Models):** These are the first layer models that analyze the input data on their own. For identifying various patterns in the data, the diversity of base learners is essential. We suggest making use of the base learners listed as: KNN, DT, SVM and MLP.

**Meta-Learner (Level-1 Model):** The output of the Meta Layer is generated by Logistic regression which takes the predictions of the base learners as input and generate Final prediction as output.

**IMPLEMENTATION AND RESULT**

**Dataset :**

Kaggle, an online company that serves as a repository for various types of data, was the source of the dataset. The file format was a Comma-Separated Values (CSV) file, and the file contains 25000 rows and 25 columns. The information collected for this study contains 10 years data.

**Prediction by Classifiers :**

The data must be pre-processed to enable high performance of the algorithm and precise forecasts. We cleaned the dataset by removing unnecessary features and handling missing values. Figure 2 displays the correlation values between the variables as a correlation matrix. In every table cell, the correlation between two variables is shown. The value is in the range of -1 to 1. If the correlation coefficient between two variables is 1, then there is a complete positive linear relationship between them.
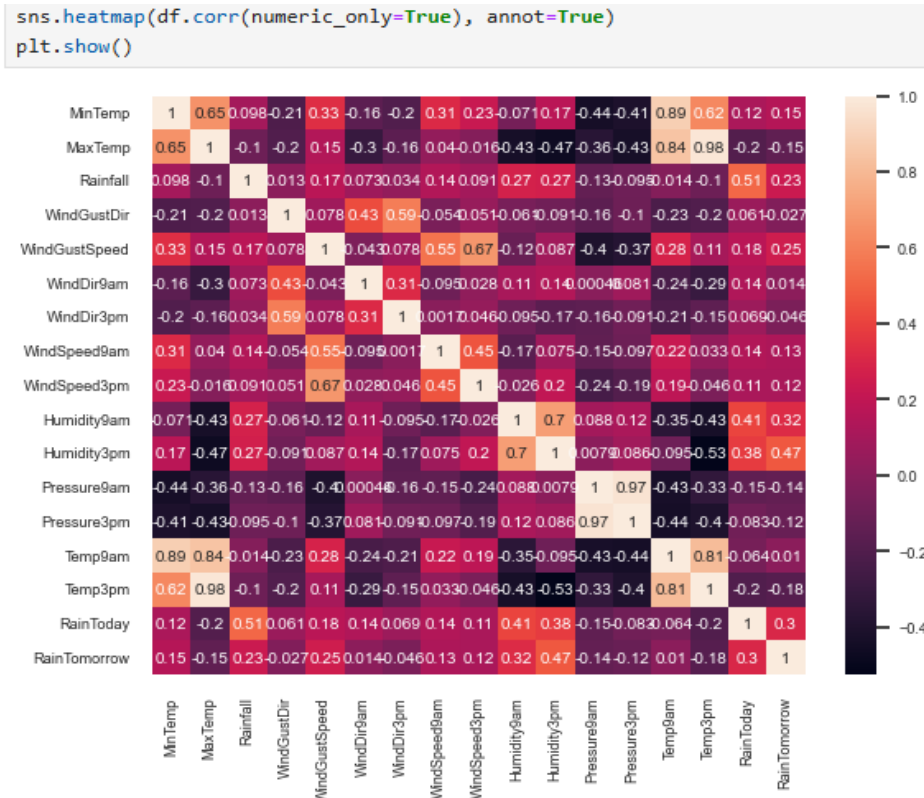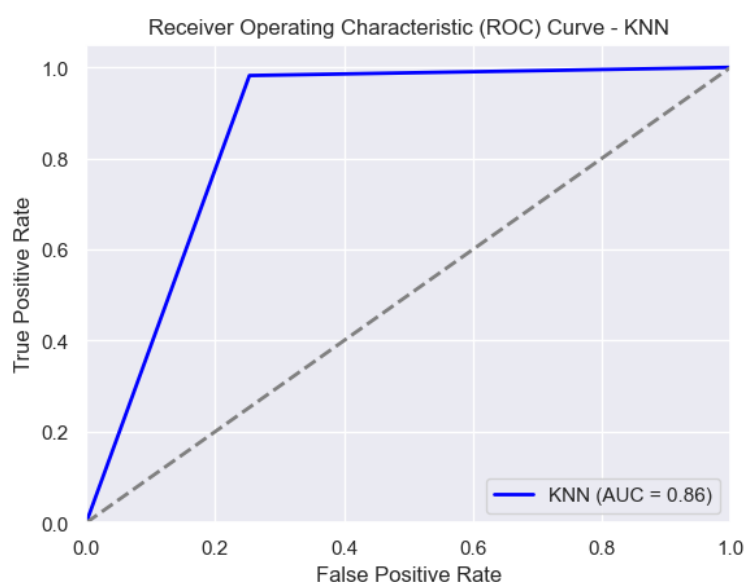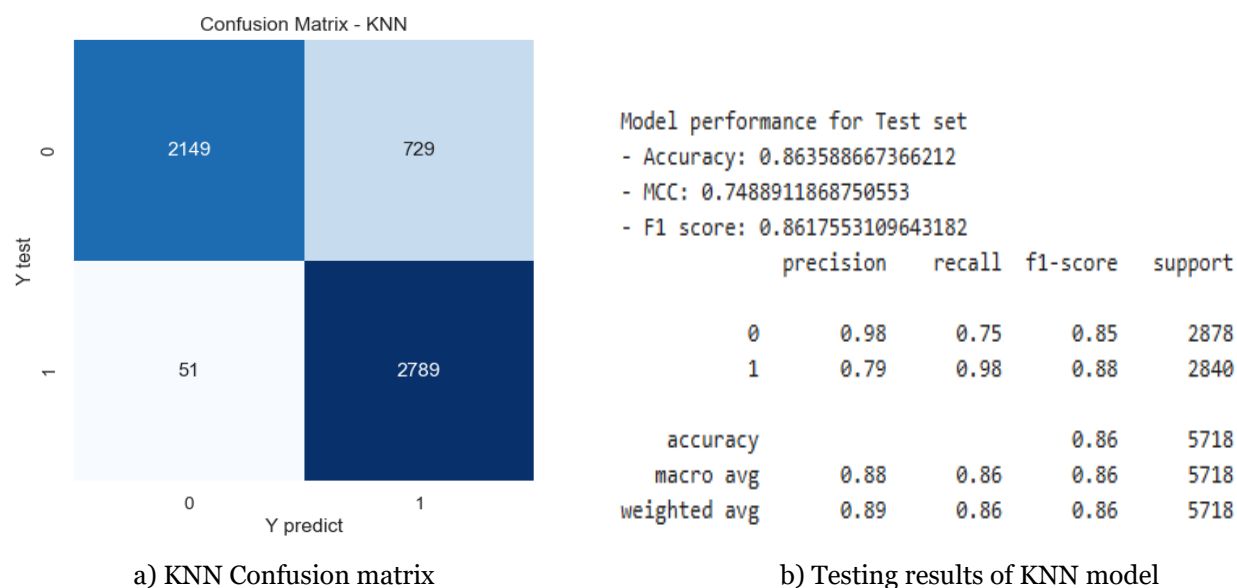
```
sns.heatmap(df.corr(numeric_only=True), annot=True)
plt.show()
```



**Figure 2:** Correlation matrix

Principal component analysis a feature extraction technique is used to select 10 Principal Components PC1, PC2, PC3, PC4, PC5, PC6, PC7, PC8, PC9, PC10 out of the features present in the dataset. Now the four classifiers based on DT, SVM, KNN and MLP [27] are applied on 10 Principal Components after which the outputs are send to the Meta-Model which applies Logistic Regression for prediction. The four Classifiers and the Proposed Model are as follows:
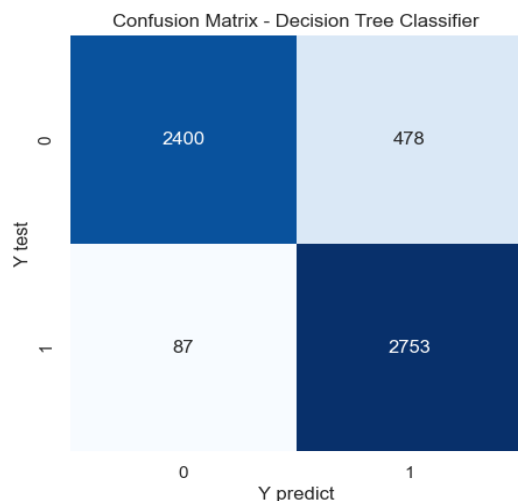
**KNN :**

KNN is a simple yet powerful machine learning method for regression and classification issues. The "K" in KNN refers to the number of closest neighbors to consider while making a prediction for a new data point [28]. The confusion matrix score is shown in Figure 3 (a). Figure 3 (b) shows the test dataset results having accuracy of 0.8635, MCC of 0.7488, F1-score of 0.8617 indicates a significant correlation between predicted and actual values. Performance is also shown in graphical form Figure 3 (c), by the ROC Curve (AUC) of 0.86.

a) KNN Confusion matrix

b) Testing results of KNN model



c) ROC curve for KNN

**Figure 3:** Overall KNN results

**DT:**

A DT is a hierarchical structure in which the result or class label is represented by each leaf node and the inside nodes, each representing a choice based on a characteristic [29]. Figure 4(a) shows the confusion matrix score. With an accuracy of 0.9011, MCC of 0.8101, F1-score of 0.9007 and AUC of 0.90 DT performed well, as shown in Figure 4(b) and Figure 4(c) respectively. The DT model was successful in identifying the underlying patterns in the data and producing precise predictions.
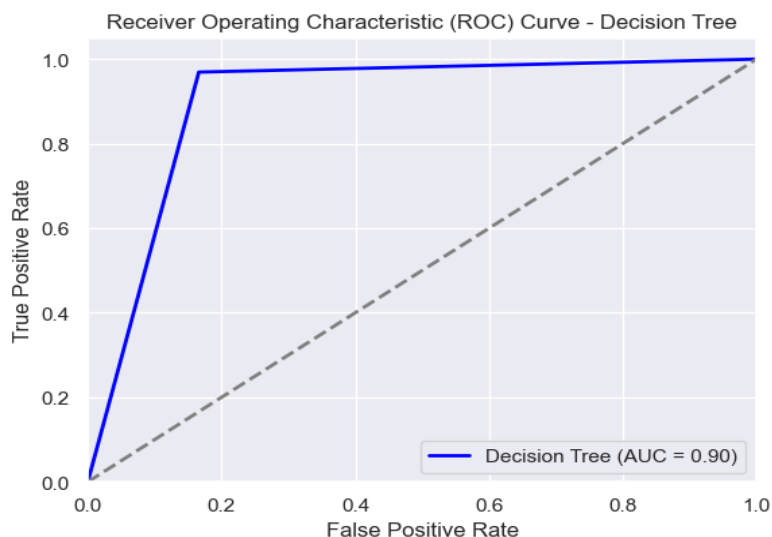
Confusion Matrix - Decision Tree Classifier

a) DT Confusion matrix

```
Model performance for Test set
- Accuracy: 0.9011892270024484
- MCC: 0.8101480854406503
- F1 score: 0.900770504982467
                 precision    recall  f1-score   support

             0       0.97      0.83      0.89      2878
             1       0.85      0.97      0.91      2840

      accuracy                           0.90      5718
     macro avg       0.91      0.90      0.90      5718
  weighted avg       0.91      0.90      0.90      5718
```
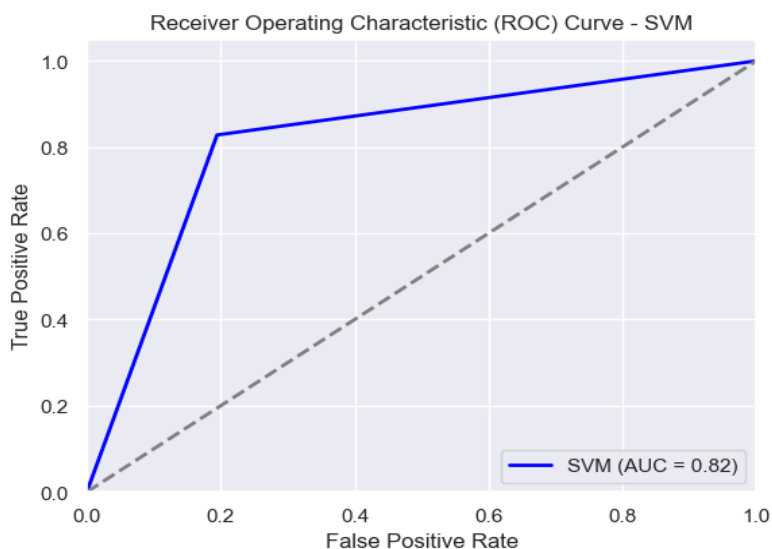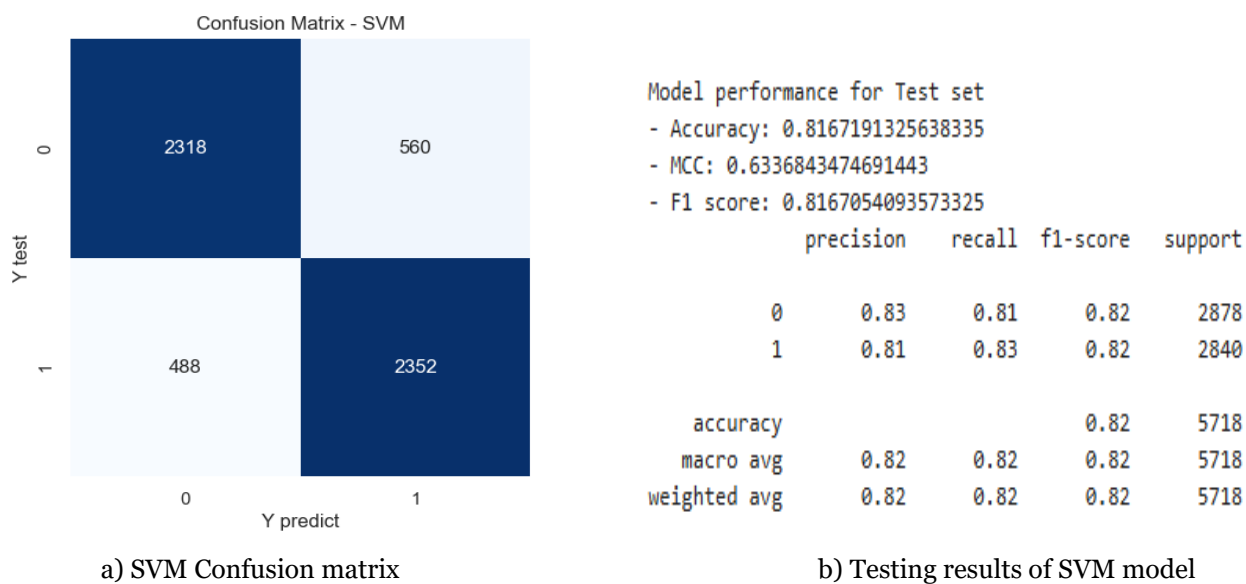
b) Testing results of DT model



c) ROC curve for DT

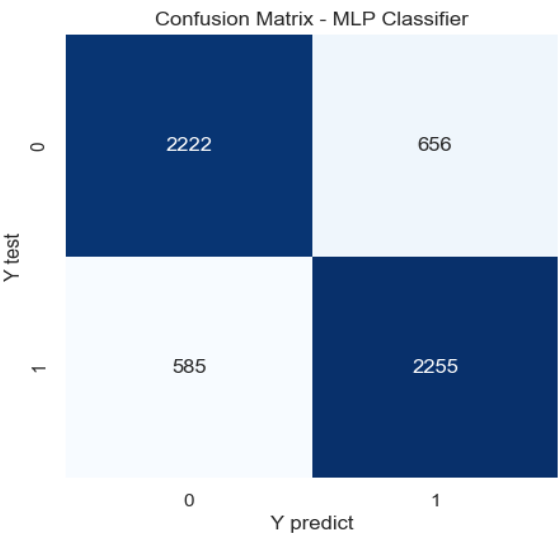**Figure 4:** Overall DT results

**SVM:**

SVM, or the data points that are closest to the decision boundary, are crucial for figuring out the margin. Because SVMs only utilize a fraction of training data points that are near the decision boundary, they are memory-efficient for large datasets [30]. Figure 5(a) shows the confusion matrix score. Figure 5(b) shows the test dataset result, which has an accuracy of 0.8167, MCC of 0.6336, and an F1-score of 0.8167. Performance is also shown in graphical form in Figure 5(c), by ROC curve (AUC) of 0.82.

a) SVM Confusion matrix

b) Testing results of SVM model



c) ROC curve for SVM

**Figure 5:** Overall SVM results

**MLP:**

MLP is an artificial neural network that consists of multiple layers of nodes, or neurons, including an input layer, one or more hidden layers, and an output layer [31]. The confusion matrix score is shown in Figure 6(a). The test dataset result and the performance in graphical form are shown in Figure 6 (b) and Figure 6 (c). The accuracy, MCC, F1 score and AUC are 0.7829, 0.5661, 0.7829 and 0.78 respectively. Even though MLP is a strong model, it does not perform well on this dataset.
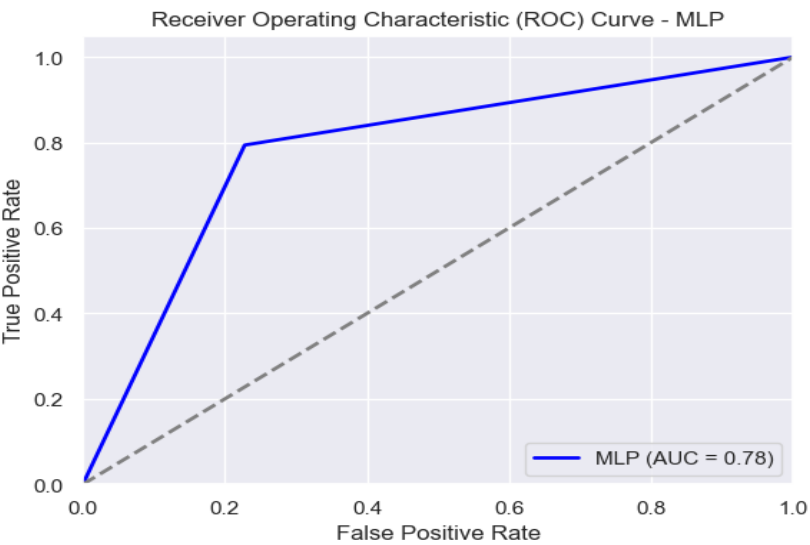
a) MLP Confusion matrix

b) Testing results of MLP model



c) ROC curve for MLP

**Figure 6:** Overall MLP results

**Proposed Stacked Model:**

In machine learning ensembles, stacked models are a potent strategy that provide improved predictive performance by utilizing the advantages of several base models and skillfully merging their predictions via a meta-model. The confusion matrix score is shown in Figure 7(a). Figure 7(b) shows the test dataset results with an excellent performance across all metrics, with an accuracy of 0.9309, MCC of 0.8643 and F1-score of 0.9308. Performance in graphical form is shown in Figure 7(c) with AUC of 0.93. This shows better classification, a high correlation between expected and actual values, and test dataset performance that is better than all other models in the given performance metrics, suggesting that ensemble approaches or layered models are useful for enhancing prediction accuracy, MCC, F1-score and AUC.
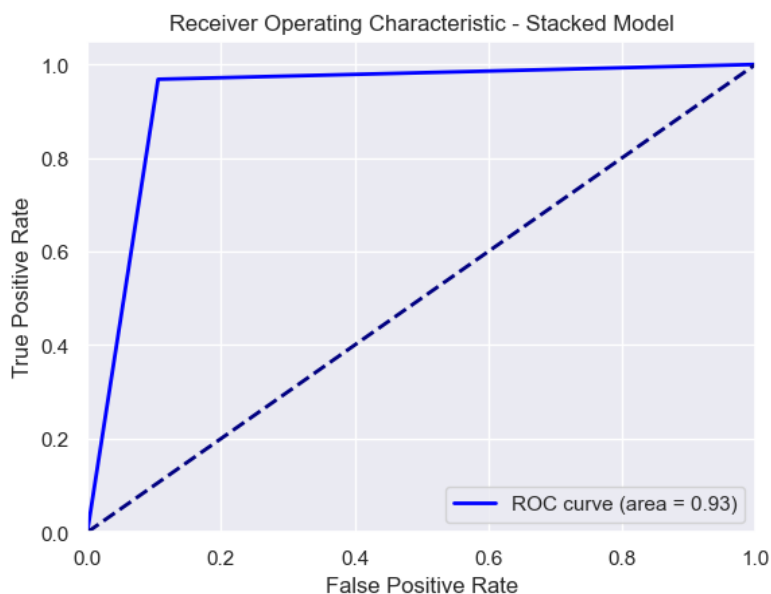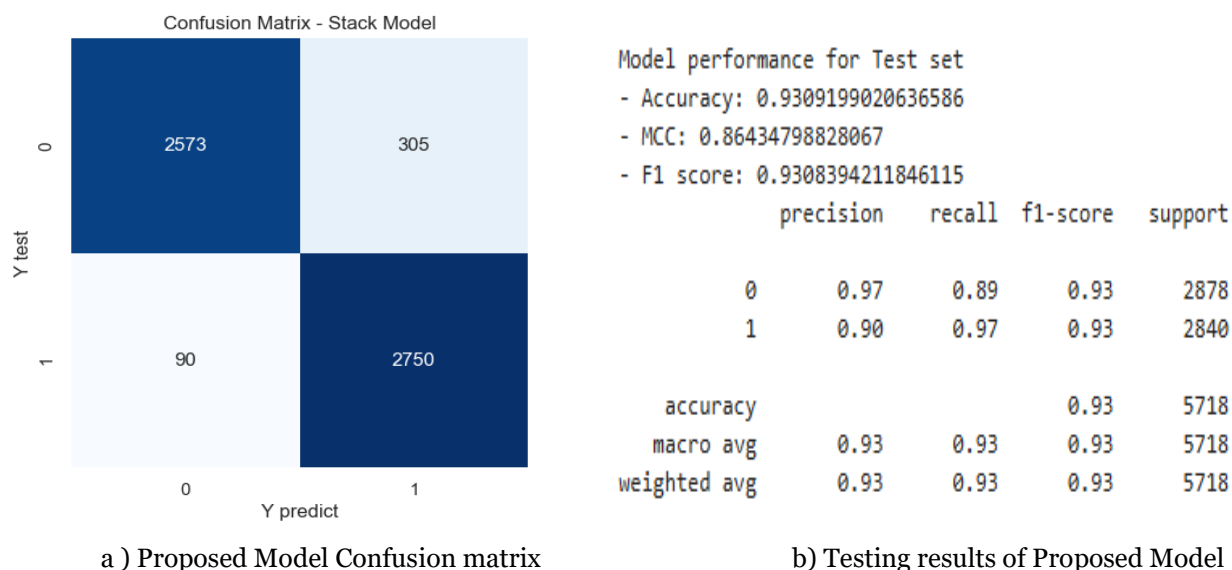
Confusion Matrix - Stack Model

|  | 0 | 1 |
|---|---|---|
| 0 | 2573 | 305 |
| 1 | 90 | 2750 |

```
Model performance for Test set
- Accuracy: 0.9309199020636586
- MCC: 0.86434798828067
- F1 score: 0.9308394211846115
                precision   recall  f1-score   support

           0       0.97      0.89      0.93      2878
           1       0.90      0.97      0.93      2840

    accuracy                           0.93      5718
   macro avg       0.93      0.93      0.93      5718
weighted avg       0.93      0.93      0.93      5718
```

a ) Proposed Model Confusion matrix            b) Testing results of Proposed Model

Receiver Operating Characteristic - Stacked Model

ROC curve (area = 0.93)

c) ROC curve for Proposed Model

**Figure 7:** Overall Proposed Model results

The comparative analysis of Test Results in Figure 8, shows that the proposed model significantly outperforms individual machine learning models in terms of accuracy, MCC and F1 score. The proposed model achieved an outstanding accuracy of 0.9309, MCC of 0.8643, and F1-score of 0.9308, indicating better classification and a strong correlation between predicted and actual values. These results show the effectiveness of ensemble methods, particularly the hybrid stacked model approach, in enhancing predictive accuracy and robustness in Weather prediction. The comparative analysis of Test Results by ROC Curve (AUC) is shown in Figure 9, in which the Stacked Model outperforms all other individual Classifiers.

**CONCLUSION**

The results analysis, which are displayed in Figure 8 and Figure 9, shows the metrics for the performance of several machine learning models that have been assessed using the dataset. With accuracy of 0.8635, MCC of 0.7488 and F1-score 0.8617, the KNN model demonstrated strong performance. The test dataset result of SVM is the least with accuracy of 0.84672, MCC 0.507379, and an F1-score of 0.827682. The MLP classifier, shows the least performance with accuracy,
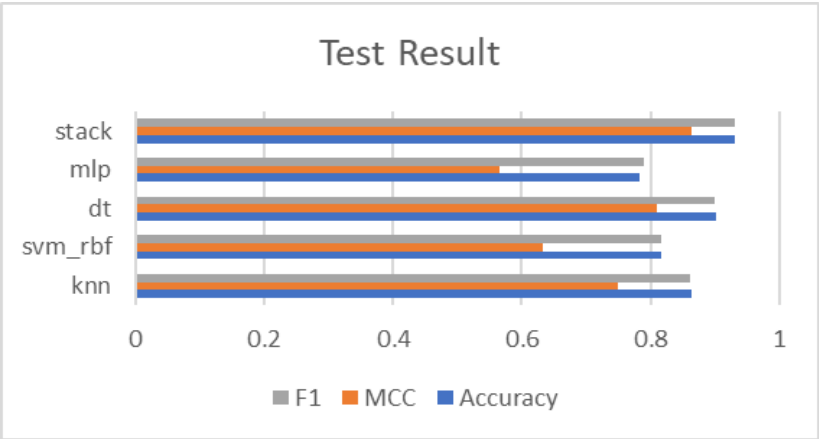
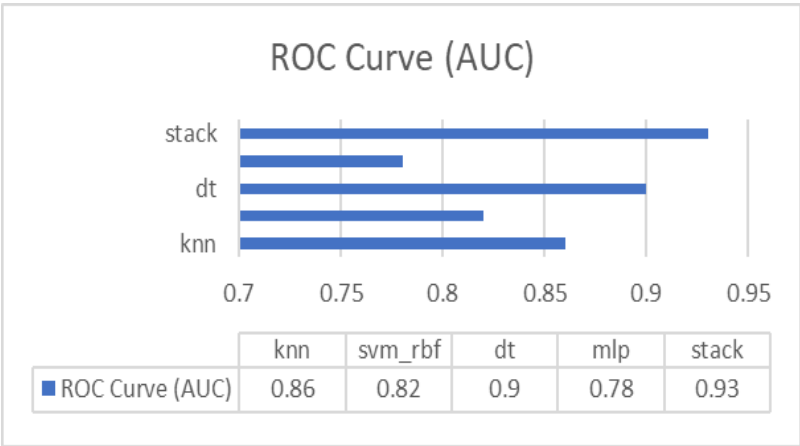**Figure 8:** Comparative analysis of Test Results



**Figure 9:** AUC Test Result analysis

MCC and F1 score of 0.7829, 0.5661 and 0.7829. DT performed better than KNN, with accuracy of 0.9011, MCC of 0.8101, and an F1-score of 0.9007. With remarkable metrics, accuracy of 0.9309, MCC of 0.8643, and F1-score of 0.9308, the stacked model significantly outperformed all individual models, indicating its better capacity to handle the dataset across all important performance measures.

## CONFLICT OF INTEREST

The authors declare no conflict of interest.

## REFERENCES

[1]    Fathi, M., Haghi Kashani, M., Jameii, S. M., & Mahdipour, E. (2022). Big data analytics in weather forecasting: A systematic review. Archives of Computational Methods in Engineering, 29(2), 1247-1275.

[2]    Lam, R., Sanchez-Gonzalez, A., Willson, M., Wirnsberger, P., Fortunato, M., Alet, F., ... & Battaglia, P. (2023). Learning skillful medium-range global weather forecasting. Science, 382(6677), 1416-1421.

[3]    Hewage, P., Trovati, M., Pereira, E., & Behera, A. (2021). Deep learning-based effective fine-grained weather forecasting model. Pattern Analysis and Applications, 24(1), 343-366.

[4]    Mall, P. K., & Singh, P. K. (2023). Credence-Net: a semi-supervised deep learning approach for medical images. International Journal of Nanotechnology, 20(5-10), 897-914.

[5]    Wu, J. M. T., Zhan, J., & Lin, J. C. W. (2017). An ACO-based approach to mine high-utility itemsets. Knowledge-Based Systems, 116, 102-113.

[6]   Tu, C. J., Chuang, L. Y., Chang, J. Y., & Yang, C. H. (2007). Feature Selection using PSO-SVM. IAENG International journal of computer science, 33(1).

[7]   Amoozegar, M., & Minaei-Bidgoli, B. (2018). Optimizing multi-objective PSO based feature selection method using a feature elitism mechanism. Expert Systems with Applications, 113, 499-514.

[8]   Bhattacharya, A., Goswami, R. T., & Mukherjee, K. (2019). A feature selection technique based on rough set and improvised PSO algorithm (PSORS-FS) for permission based detection of Android malwares. International journal of machine learning and cybernetics, 10, 1893-1907.

[9]   Wei, B., Zhang, W., Xia, X., Zhang, Y., Yu, F., & Zhu, Z. (2019). Efficient feature selection algorithm based on particle swarm optimization with learning memory. IEEE Access, 7, 166066-166078.

[10]  Shehadeh, H. A., Jebril, I. H., Jaradat, G. M., Ibrahim, D., Sihwail, R., Al Hamad, H., ... & Alia, M. A. (2023). Intelligent Diagnostic Prediction and Classification System for Parkinson's Disease by Incorporating Sperm Swarm Optimization (SSO) and Density-Based Feature Selection Methods. International Journal of Advances in Soft Computing & Its Applications, 15(3).

[11]  Abukhodair, F., Alsaggaf, W., Jamal, A. T., Abdel-Khalek, S., & Mansour, R. F. (2021). An intelligent metaheuristic binary pigeon optimization-based feature selection and big data classification in a MapReduce environment. Mathematics, 9(20), 2627.

[12]  Grönquist, P., Yao, C., Ben-Nun, T., Dryden, N., Dueben, P., Li, S., & Hoefler, T. (2021). Deep learning for post-processing ensemble weather forecasts. Philosophical Transactions of the Royal Society A, 379(2194), 20200092.

[13]  Guan, H., & Zhu, Y. (2017). Development of verification methodology for extreme weather forecasts. Weather and forecasting, 32(2), 479-491.

[14]  Li, G., Jing, C., Jiawen, Z., & Quanliang, C. (2019). Progress in researches on ensemble forecasting of extreme weather based on numerical models. Advances in Earth Science, 34(7), 706.

[15]  Xu, W., Chen, K., Han, T., Chen, H., Ouyang, W., & Bai, L. (2024). ExtremeCast: Boosting Extreme Value Prediction for Global Weather Forecast. arXiv preprint arXiv:2402.01295.

[16]  Ben Bouallègue, Z., Clare, M. C., Magnusson, L., Gascon, E., Maier-Gerber, M., Janoušek, M., ... & Pappenberger, F. (2024). The rise of data-driven weather forecasting: A first statistical assessment of machine learning-based weather forecasts in an operational-like context. Bulletin of the American Meteorological Society.

[17]  Lopes, F. M., Dutra, E., & Boussetta, S. (2024). Evaluation of Daily Temperature Extremes in the ECMWF Operational Weather Forecasts and ERA5 Reanalysis. Atmosphere, 15(1), 93.

[18]  Xu, H., Zhao, Y., Zhao, D., Duan, Y., & Xu, X. (2024). Improvement of disastrous extreme precipitation forecasting in North China by pangu-weather AI-Driven regional WRF model. Environmental Research Letters.

[19]  Bouallègue, Z. B., Weyn, J. A., Clare, M. C., Dramsch, J., Dueben, P., & Chantry, M. (2024). Improving medium-range ensemble weather forecasts with hierarchical ensemble transformers. Artificial Intelligence for the Earth Systems, 3(1), e230027.

[20]  Beimel, S., Suari, Y., & Gabbay, F. (2024). Improving Weather Forecasts for Sailing Events Using a Combination of a Numerical Forecast Model and Machine Learning Postprocessing. Applied Sciences, 14(7), 2950.

[21]  Wang, W., & Ikegaya, N. Predicting Annual Extreme Winds in Iran Using Numerical Weather Forecasting, Meteorological Observation, and Statistical Model.

[22]  Varshney, Y., Kumar, V., Dubey, D. K., & Sharma, S. Forecasting Precision: The Role of Graph Neural Networks and Dynamic GNNs in Weather Prediction.

[23]  Obisesan, O. E. (2024). Machine Learning Models for Prediction of Meteorological Variables for Weather Forecasting. International Journal of Environment and Climate Change, 14(1), 234-252.

[24]  Nguyen, T., Jewik, J., Bansal, H., Sharma, P., & Grover, A. (2024). Climatelearn: Benchmarking machine learning for weather and climate modeling. Advances in Neural Information Processing Systems, 36.

[25]  Fatima, N., & Siddiqi, S (2024). Acute Myocardial Infarction: Prediction and Patient Assessment through Different ML Techniques. International Journal of Intelligent Systems and Applications in Engineering, 12(13s), 106–121.

[26]  Fatima, N., & Siddiqi, S (2024).Enhanced Myocardial Infarction Prediction Using Machine Learning Algorithms and Gender-Specific Insights. Journal of Electrical Systems,  Vol. 20 No. 7s (2024), 973-988.

[27] Agrawal, U., Arora, J., Singh, R., Gupta, D., Khanna, A., & Khamparia, A. (2020). Hybrid wolf-bat algorithm for optimization of connection weights in multi-layer perceptron. ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM), 16(1s), 1-20.

[28] Taunk, K., S. De, Verma, S., and Swetapadma, A., "A Brief Review of Nearest Neighbor Algorithm for Learning and Classification," 2019 International Conference on Intelligent Computing and Control Systems (ICCS), Madurai, India, 2019, pp. 1255-1260, doi: 10.1109/ICCS45141.2019.9065747.

[29] Bahzad, J., Abdulazeez M., Adnan. (2021). Classification Based on Decision Tree Algorithm for Machine Learning. Journal of Applied Science and Technology Trends. 2. 20-28.

[30] Cervantes, J., Lamont, F. G., Mazahua, L.G., Lopez, A., A comprehensive survey on support vector machine classification: Applications, challenges and trends, Neurocomputing, Volume 408, 2020, Pages 189-215.

[31] Hassan, R., & Amine, Mohammed & Idrissi, Janati & Ghanou, Youssef & Ettaouil, Mohamed. (2016). Multilayer Perceptron: Architecture Optimization and Training. International Journal of Interactive Multimedia and Artificial Inteligence. 4. 26-30. 10.9781/ijimai.2016.415.