Contents lists available at https://citedness.com/index.php/jdsi

# Data Science Insights

Journal Page is available at https://citedness.com/index.php/jdsi

Research articles

# Comprehensive Approach to Weather Prediction with the Random Forest Algorithm

*Pedro Joyarieb [1] , Vian Candra Silalahi [2] , Vallencia Anggelica [3] , Khatrina Kelly Ongso[4]*

Faculty of Computer Science, Pelita Indonesia Institute of Business and Technology, Pekanbaru, Indonesia
email: [1] pedro.joyarieb@student.pelitaindonesia.ac.id , [2] vian.candra@student.pelitaindonesia.ac.id ,
[3] vallencia.anggelica@student.pelitaindonesia.ac.id , [4] khatrina@student.pelitaindonesia.ac.id

## ARTICLE INFO

## ABSTRACT

Weather is an air condition that is very important in everyday life. Accurate weather predictions can help people anticipate and deal with weather changes that can have an impact on daily activities. This research aims to develop an effective weather prediction model using machine learning algorithms. In this research, we use three popular machine learning algorithms, namely Random Forest, Support Vector Machine (SVM), and Decision Tree. The data used consists of historical weather data, including air temperature, air humidity, rainfall, wind direction, air pressure, wind speed, and solar radiation. The research results show that the Random Forest algorithm has the highest accuracy, with a prediction rate of 83%. The SVM algorithm is next, with a prediction rate of 78%, while the Decision Tree algorithm has a prediction rate of 72%. These findings show that Random Forest is the most effective algorithm in predicting weather, especially in predicting air temperature and rainfall. This research has significant practical implications in increasing the accuracy of weather predictions, which can help society anticipate and deal with weather changes that can impact in daily activities. In the future, this research can be used as a basis for developing more accurate and reliable weather prediction systems.

Correspondence:
Pedro Joyarieb, Faculty of Computer Science, Pelita Indonesia Institute of Business and Technology, Pekanbaru, Indonesia
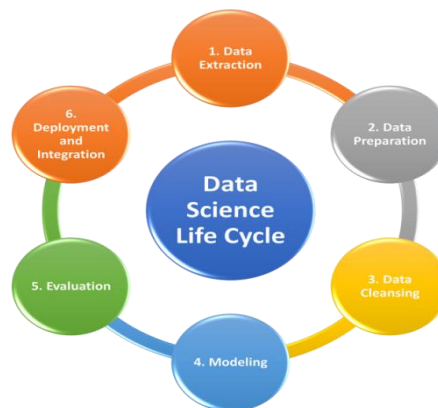
## 1. Introduction

Weather is an important factor in all aspects of life. Significant weather changes affect human activities. Several activities are related to weather such as agricultural activities, plantations and aviation. [1] Weather is an important factor in all aspects of life. Significant weather changes affect human activities. Several activities are related to weather such as agricultural activities, plantations and aviation. It is important to monitor the weather so that potential solutions such as heavy rain and lightning can be prepared. [2] Weather prediction is an important element in daily life that affects various sectors such as agriculture, aviation, urban planning , and natural disaster mitigation. Forecasting is a tool used to predict a value in the future by paying attention to relevant data or features, both based on past data and current data. [3] Over the decades, weather prediction methods have evolved from traditional numerical models based on physics equations and atmospheric dynamics, to more sophisticated approaches. ML is a branch of Artificial Intelligence (AI) that allows computers to develop behavior based on empirical data or data obtained from observations or experiments. [4] On the other hand , machine learning algorithms have shown great potential in improving the accuracy of weather predictions. One algorithm that stands out in this regard is Random Forest. [5] Random Forest is an ensemble learning technique that combines predictions from several decision trees to produce more accurate and stable results. This algorithm is known for its ability to handle data with a large number of features and overcome the problem of overfitting, which is often a challenge in machine learning models. [6] Furthermore, there is a decision tree algorithm which helps researchers to obtain information that is easy to understand, analyze and interpret accurately based on the weather dataset in the city of Jakarta in January 2018. [7] and the last algorithm used, Support Vector Machines (SVM) is a classification model in machine learning with a binary or discriminative model, working on two differentiation classes. SVM is used to find the best hyperplane by maximizing the distance between classes [5].

The advantage of Random Forest in weather prediction lies in its ability to handle the complexity and non-linearity of meteorological data. Weather data usually consists of various variables such as temperature, humidity, atmospheric pressure, wind speed, and historical weather data, all of which have complex relationships with each other. Research shows that Random Forest has the greatest influence on weather prediction results with accuracy. [8] This research aims to explore and implement the Random Forest algorithm in weather prediction, with a focus on increasing accuracy and efficiency. Using a comprehensive meteorological dataset and rigorous validation methods, this research will analyze Random Forest's performance in various weather conditions and compare it with other prediction models. In the process of creating weather forecast information, there are several obstacles. First, it is difficult to create forecast information because it involves many data sources such as observation data, weather application model data, cloud condition image data from satellites, cloud condition data from radar. Second, maritime weather forecasts generally rely on the abilities of forecasters, so that the resulting interpretations can differ from one forecaster to another because it depends on their individual experience. Differences in interpretation can confuse users, which ultimately has the potential to reduce the quality of the information conveyed. [9]

This research shows that the Random Forest model with the Cross-Validation technique can be used to predict rain with high accuracy. The research results show that the Random Forest algorithm has the highest accuracy, with a prediction rate of 83%. The SVM algorithm is next, with a prediction rate of 78%, while the Decision Tree algorithm has a prediction rate of 72%. In this way, it is hoped that it can provide accurate rainy weather forecasts and provide precise information by applying machine learning models.

## 2. Research Methods

In this research, researchers succeeded in developing a Random Forest model with 83% accuracy to predict weather types based on prepared data. Apart from that, researchers also evaluated the SVM model as a comparison and found that Random Forest had higher accuracy (83%) than SVM (78%) and Decision Tree (72%). These steps provide an overview of how the data science lifecycle is applied in solving business problems using weather data. Random Forest is an ensemble learning that is built from decision trees. [2] and Classification cases can be found in various domains. [2] Several studies implemented machine learning algorithms for weather classification and prediction [2] , carrying out weather classification with several classification algorithms such as SVM, Decision Tree and Random Forest. [15] By following these stages, researchers can ensure that the data-based solutions they develop are relevant, accurate, and reliable for decision making.



1st Picture. Lifcycle science data used

Based on the diagram in Picture 1, the Data Science Lifecycle scheme consists of six main stages:

### 2.1. Data Extraction (Data Extraction)

This stage is the first step in the Data Science Lifecycle. At this stage, the data required for the project is collected from various sources, such as databases, text files, APIs, and sensors. This data can be structured data, semi-structured data, or unstructured data. In collecting this data, there are several things that researchers do to obtain This data, namely: The data used comes from Kaggle, with a dataset named "S eattle - Weather - Dataset". The data was downloaded on May 2 4 2024 , on 22 o'clock . 00 WIB . Extraction datasets from Kaggle. Data cleaning to remove irrelevant data, fill in values missing data, and handling invalid data. Pre-processing data For prepare it For analysis, including changing formats and merging tables if necessary. Researchers will use Visual Studio Code software to explore and serve aspects Which can give outlook Which valuable about Weather

Prediction with method visualize data. By using this tool, researchers can create visualizations interesting and informative to understand more deeply predicting the weather .

### 2.2. Data Preparation (Data Preparation)

Once the data is extracted, it needs to be prepared for analysis. This stage includes data cleaning, data transformation, and data integration. Data cleaning aims to overcome data quality problems, such as missing data, inconsistent data, and outlier data. Data transformation aims to change the data format, such as changing from CSV format to JSON format, or changing the data scale. Data integration aims to combine data from various sources into one coherent dataset.

### 2.3. Data Cleansing (Data Cleaning)

This stage is part of Data Preparation which focuses on cleaning data from errors, inconsistencies and missing values. Dirty data can lead to inaccurate and misleading analysis results.

### 2.4. Modeling (Modeling)

At this stage, statistical or machine learning models are built to analyze the data and generate new predictions or insights. This model can be a regression, classification, or clustering model.

### 2.5. Evaluation

Once the model is built, it needs to be evaluated to ensure that it is accurate and reliable. Model evaluation is carried out using various metrics, such as accuracy, precision, and recall.
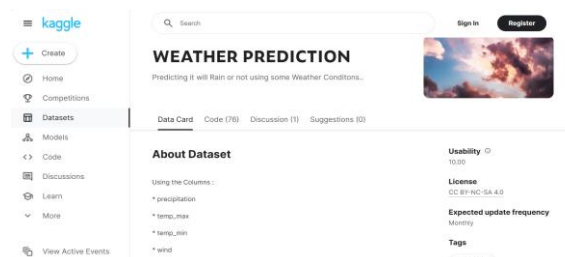
### 2.6. Deployments

The final stage in the Data Science Lifecycle is deploying the model to a production environment. These models can be integrated into web applications, mobile applications, or other systems to provide predictions or new insights in real-time. The Data Science Lifecycle schema is not always linear. The stages in this scheme can be repeated and adjusted according to project needs. It is important to note that the Data Science Lifecycle schema is just a framework. This framework can be adapted and modified to meet specific project needs.

## 3. Results and Discussion

### 3.1. Data Extraction Results

Data Which used writer obtained from site web Kaggle ( www.kaggle.com ) like on image 2 below :



2nd Picture. Data source

Following is data Which has researcher get:



3rd Picture. Data extraction results

Each data in the table has a different data type and there are also types the same data. Data type In data there are data types , Numerical , Object , Date time . Date ( Date ) = Date time , Precipitation (Rainfall) = Numerical , Maximum Temperature = Numerical , Minimum Temperature = Numerical , Wind ( Wind pressure ) = Numerical , Weather (Approx weather) = Object t

### 3.2. Data Preparation Results



4th Picture. Data preparation results

### 3.3. Data Cleaning Results



5th Picture. Data have been data cleaned

From Figure 5 above you can see the weather collected from 2012, where the author recorded weather data, maximum and minimum temperatures, rainfall and wind. The following is link from part data study, Look in https://docs.google.com/spreadsheets/d/1My0pJqEZbbhm_5QA6-3dYVE-e2OVcnpYfFh5ySO_KFU/edit?gid=0#gid=0 for versionin full.

### 3.4. Data Modeling and Evaluation

*Random Forest Model*

In weather dataset analysis using the Random Forest algorithm, this model achieved 83% accuracy. This shows that the model performs well in predicting weather conditions based on features such as precipitation, maximum temperature, minimum temperature, and wind speed. Random Forest is a strong choice due to its ability to handle complex and large datasets, and provide accurate and stable results.

```
# Menghitung akurasi dan Laporan klasifikasi
accuracy = accuracy_score(y_test, y_pred)
print(f'Akurasi Random Forest: {accuracy*100:.2f}%')
print(classification_report(y_test, y_pred))

✓ 0.0s

Akurasi Random Forest: 83.14%
              precision    recall  f1-score   support

     drizzle       0.17      0.07      0.10        14
         fog       0.33      0.16      0.21        32
        rain       0.96      0.92      0.94       192
        snow       0.50      0.25      0.33         8
         sun       0.78      0.94      0.85       193

    accuracy                           0.83       439
   macro avg       0.55      0.47      0.49       439
weighted avg       0.80      0.83      0.81       439
```

6th Picture. Prediction Model Using the Random Forest Algorithm

*Support Vector Machine (SVM) Models*

In comparison, the SVM algorithm was applied to the same dataset and achieved 78% accuracy. Even though it is lower than Random Forest, SVM still provides quite good results. However, SVM has disadvantages in terms of computing time and sensitivity to features of different scales.

```
# Evaluasi akurasi
accuracy = accuracy_score(y_test, y_pred)
print(f'Akurasi SVM: {accuracy*100:.2f}&')
# Laporan klasifikasi
print(classification_report(y_test, y_pred))


Akurasi SVM: 77.82&
              precision    recall  f1-score   support

     drizzle       0.00      0.00      0.00         9
         fog       0.00      0.00      0.00        25
        rain       0.93      0.81      0.87       120
        snow       0.00      0.00      0.00         8
         sun       0.69      1.00      0.82       131

    accuracy                           0.78       293
   macro avg       0.33      0.36      0.34       293
weighted avg       0.69      0.78      0.72       293
```

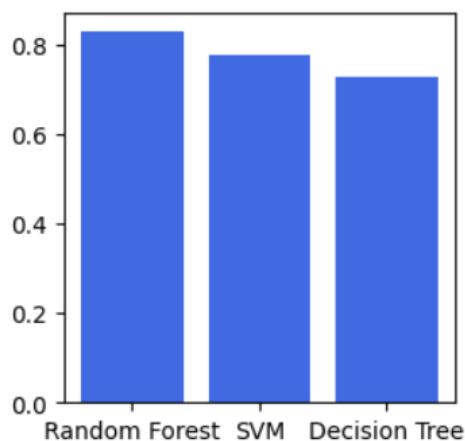7th Picture. Prediction Model Using the SVM Algorithm

*Decision Tree Model*

The Decision Tree model succeeded in achieving 73.3% accuracy in classifying weather types based on other weather parameters such as rainfall, maximum temperature, minimum temperature and wind speed. This accuracy shows that the model is able to recognize patterns in the data quite well, but there is still room for improvement.

```
# Menghitung Akurasi
accuracy = accuracy_score(y_test, y_pred)
print(f"Akurasi Decision Tree: {accuracy * 100:.2f}%")
# Laporan Klasifikasi
print(classification_report(y_test, y_pred))
✓ 0.0s

Akurasi Decision Tree: 72.70%
              precision    recall  f1-score   support

           0       0.00      0.00      0.00         9
           1       0.23      0.20      0.21        25
           2       0.88      0.89      0.88       120
           3       0.38      0.38      0.38         8
           4       0.77      0.75      0.76       131

    accuracy                           0.73       293
   macro avg       0.45      0.44      0.45       293
weighted avg       0.73      0.73      0.73       293
```
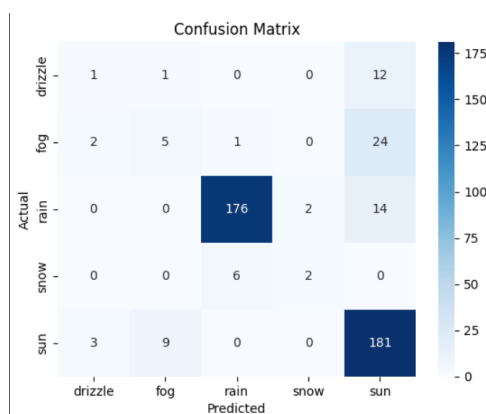
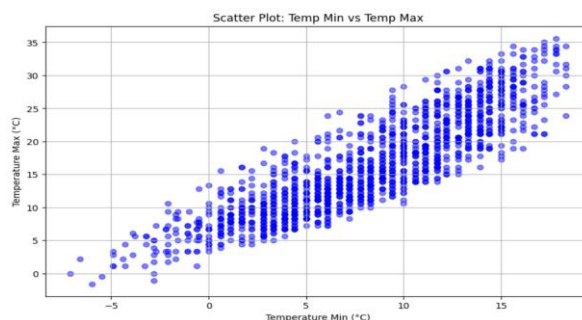8th Picture. Prediction Model Using the Decision Tree Algorithm

9th Picture. Visualization of comparative accuracy of the models used

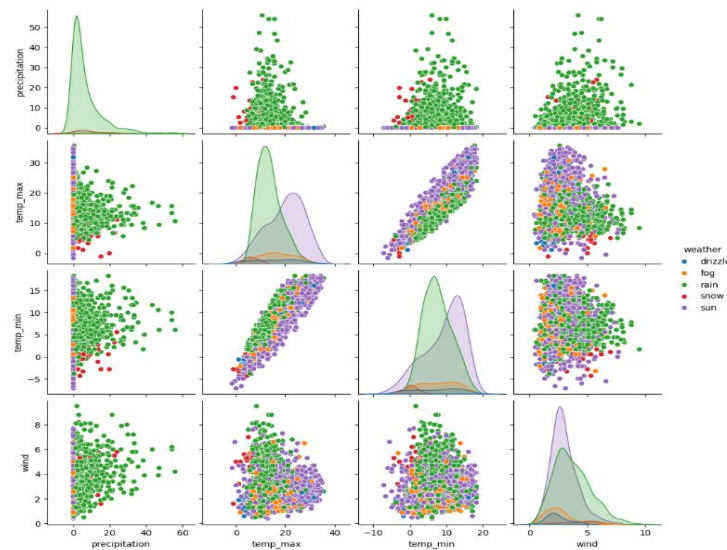Here are some data visualizations that can be concluded:



10th Picture. Confusion Matrix

The most likely weather is sun: This can be seen from the value 175 in the "sun" cell in the "sun" row. This means that there are 175 events where the actual weather condition is the sun and the predicted weather condition is also the sun. The least likely weather is snow: This can be seen from the value 0 in the "snow" cell in the "snow" row. This means that there are no instances where the actual weather condition is snow and the predicted weather condition is also snow. Rainy weather predictions are quite accurate: This can be seen from the value 150 in the "rain" cell in the "rain" row. This means that there are 150 events where the actual weather conditions are rain and the predicted weather conditions are also rain. The drizzle weather prediction is relatively accurate: This can be seen from the value 176 in the "drizzle" cell in the "drizzle" row. This means that there were 176 events where the actual weather condition was drizzle and the predicted weather condition was also drizzle. The fog weather prediction is less accurate: This can be seen from the value 75 in the "fog" cell in the "fog" row. This means that there were 75 events where the actual weather condition was fog and the predicted weather condition was also fog.



11th Picture. Scatter Plot Visualization

This plot depicts the relationship between minimum and maximum temperatures. Each dot represents a specific day. The points are distributed with an upward trend, indicating that when the minimum temperature is higher, the maximum temperature also tends to be higher.



12th Picture. Visualization Using Pair Plot

This plot shows the relationship between all the numerical variables in the dataset. The diagonal plot shows the distribution of each variable. The off-diagonal scatter plot shows the relationship of pairs of variables, with coloring based on weather type.

## 4. Conclusion

Leveraging machine learning algorithms to predict weather is a very important step because this technique is able to capture the complexity and hidden patterns in huge and diverse weather data, which is often difficult to identify by traditional prediction methods. In the analysis carried out, the Random Forest algorithm provided the highest accuracy of 83%. This shows its superiority in handling multiple features effectively and reducing the risk of overfitting through the use of ensemble learning, where multiple decision trees are combined to make more robust and stable predictions. The SVM (Support Vector Machine) algorithm with an accuracy of 77.8% also shows excellent capabilities in classifying data by separating classes optimally in a high feature space, even though the complexity of this model is higher compared to Decision Tree. Meanwhile, Decision Tree provides an accuracy of 73.3%, and despite its lower accuracy, the model offers very high interpretability, allowing us to understand and explain how decisions are made based on various weather features. This accuracy advantage is very important because more precise and reliable weather predictions can help in various aspects of daily life such as activity planning, agriculture, natural resource management and disaster mitigation. With increased accuracy of weather predictions, we can be better prepared to deal with extreme weather conditions and make more informed decisions.

## References

[1]     I. Intan, S. Aminah Dinayati Ghani, AT Koswara, U. Dipa Makassar, K. National Archives of the Republic of Indonesia, and JP Independen, "Performance Analysis of Weather Forecasting Using Machine Learning Algorithms Performance Analysis of Weather Forecasting using Machine Learning Algorithms," *Jurnal_Pekommas_Vol._6_No* , vol. 2, pp. 1–8, 2021, doi: 10.30818/jpkm.2021.2060221.

[2]     F. Hamami and IA Dahlan, "Weather Classification for DKI Jakarta Province Using the Random Forest Algorithm with Oversampling Techniques," *J. Teknoinfo* , vol. 16, no. 1, p. 87, 2022, doi: 10.33365/jti.v16i1.1533.

[3]     AA Karim, MA Prasetyo, and MR Saputro, "Comparison of Random Forest, K-Nearest Neighbor, and SVM Methods in Accuracy Prediction of Italian League Matches," *Semin. Nas. Technol. Science* , vol. 2, no. 1, pp. 377–382, 2023, [Online]. Available: https://proceeding.unpkediri.ac.id/index.php/stains/article/view/2877

[4]     B. Suma and U. Pasundan, "Bandung September 2020," no. January, 2021, doi: 10.13140/RG.2.2.16086.47680.

[5]     R. Risanti, "Analysis of Weather Prediction Models Using Support Vector Machine, Gradient Boosting, Random Forest, and Decision Tree," vol. XII, pp. 119–128, 2024, doi: 10.21009/03.1201.fa18.

[6]     ZA Dwiyanti and C. Prianto, "Weather Prediction for Jakarta City Using the Random Forest Method," *J. Tekno Incentive* , vol. 17, no. 2, pp. 127–137, 2023, doi: 10.36787/jti.v17i2.1136.

[7]     F. Indriaharti Harida and N. Khazizah, "Weather analysis in Jakarta City in January 2018 using the Decision Tree Algorithm," *J. Poros Tek.* , vol. 14, no. 1, pp. 33–37, 2022, [Online]. Available: https://www.kaggle.com/datasets/msf1203/pr

[8]     M. Arrieta-Prieto and KR Schell, "Spatially transferable machine learning wind power prediction models: v−logit random forests," *Renew. Energy* , vol. 223, no. June 2023, p. 120066, 2024, doi: 10.1016/j.renene.2024.120066.

[9]     S. Mujiasih, "Utilization of Data Mining for Weather Forecasts," *J. Meteorol. and Geofis.* , vol. 12, no. 2, 2011, doi: 10.31172/jmg.v12i2.100.

[10]    MN Rifqi and RT Aldisa, "Application of the Support Vector Machine Method in Predicting Weather Predictions," *J. Comput. Syst. Informatics* , vol. 5, no. 2, pp. 368–379, 2024, doi: 10.47065/josyc.v5i2.4961.

[11]    ARI Pratama, SA Latipah, and BN Sari, "Optimization of Rainfall Classification Using Support Vector Machine (Svm) and Recursive Feature Elimination (Rfe)," *JIPI (Journal of Scientific Research and Information Learning* , vol. 7, no. 2, pp. 314–324, 2022, doi: 10.29100/jipi.v7i2.2675.

[12]    AM Siregar, "Classification for Weather Prediction Using Esemble Learning," *Petir* , vol. 13, no. 2, pp. 138–147, 2020, doi: 10.33322/petir.v13i2.998.

[13]    D. Husen *et al.* , "Forest Fire Prediction Analysis Using the Random Forest Classifier Algorithm Forest and land fires in Indonesia have become an international concern, especially since the forest fires that occurred in the 80s [2]. Causes of immortality," vol. 16, pp. 150–155, 2022.

[14]    C. Dewi, "The Influence of Anfis Architecture on Weather Forecasting," *J. Environmental Eng. Sustain. Technol.* , vol. 2, no. 1, pp. 12–19, 2015, doi: 10.21776/ub.jeest.2015.002.01.3.

[15]    S. Shalsabilla, P. Rachmawati, K. Vidya Prakusa, and S. Rihastuti, "Application of Data Mining with the Decision Tree Method for Weather Prediction in the City of Seattle using the Weka Application," *Semin. Nas. Amikom Surakarta* , no. November, pp. 93–100, 2023.