

Application of Data Mining for Rainfall Prediction Classification in Australia with Decision Tree Algorithm and C5.0 Algorithm

Penerapan *Data Mining* Untuk Klasifikasi Prediksi Hujan di Australia dengan Algoritme *Decision Tree* dan Algoritme C5.0

Irwansyah Saputra¹, Dinar Ajeng Kristiyanti²

^{1,2} Ilmu Komputer, IPB University, Indonesia

¹ Sistem Informasi, Universitas Nusa Mandiri, Indonesia

² Teknologi Informasi, Universitas Bina Sarana Informatika, Indonesia

^{1*}92irwansyah@apps.ipb.ac.id, ^{2*}dinarajengkristiyanti@apps.ipb.ac.id

Keywords: Classification; Decision Tree; C5.0; Rstudio; Rain Australia

Abstract

Purpose: This study aims to predict rain in Australia with a machine learning classification approach. Precise and accurate rain prediction is very important for planning and management of water resources, flood warning, construction activities and aviation operations, and others.

Design/methodology/approach: The methods or stages applied in classifying rain predictions in Australia are through several stages including Data Collection, Data Pre-processing (including Missing Value handling in it), Classification Modeling by applying and comparing Decision Tree and C5.0 algorithms, Result validation using Dataset Partition and 10-Cross Fold Validation and Model Evaluation using Confusion Matrix.

Findings/result: The methods or stages applied in classifying rain predictions in Australia are through several stages including Data Collection, Data Pre-processing (including Missing Value handling in it), Classification Modeling by applying and comparing Decision Tree and C5.0 algorithms, Result of validation using Dataset Partition and 10-Cross Fold Validation and Model Evaluation using Confusion Matrix.

Originality/value/state of the art: In addition to the classification model used, dataset validation, either by partitioning the dataset or by 10-cross fold validation, can also affects the accuracy of the prediction results.

Kata kunci: Klasifikasi; *Decision Tree*; C5.0; RStudio; Hujan Australia

Abstrak

Tujuan: Penelitian ini bertujuan untuk memprediksi hujan di Australia dengan pendekatan klasifikasi *machine learning*. Prediksi hujan yang tepat dan akurat sangat penting untuk perencanaan dan pengelolaan sumber daya air, peringatan banjir, kegiatan konstruksi dan operasi penerbangan serta yang lainnya.

Perancangan/metode/pendekatan: Metode atau tahapan yang diterapkan dalam melakukan klasifikasi prediksi hujan di Australia yaitu melalui beberapa tahapan diantaranya Pengumpulan Data, *Data Pre-processing* (termasuk dilakukan penanganan *Missing Value* didalamnya), Pemodelan Klasifikasi dengan menerapkan dan membandingkan algoritme *Decision Tree* dan C5.0, Validasi Hasil menggunakan Partisi *Dataset* dan *k-Cross Fold Validation* serta Evaluasi Model menggunakan *Confusion Matrix*.

Hasil: Berdasarkan hasil yang diperoleh, evaluasi menggunakan *10-Cross Fold Validation* lebih unggul yang memiliki akurasi paling tinggi sebesar 87.35% untuk algoritme *Decision Tree* dan akurasi sebesar 86.85% untuk algoritme *C5.0 Rule-Based Model*, dibandingkan dengan metode Split 80:20 pada kasus prediksi hujan di Australia.

Keaslian/state of the art: Selain model klasifikasi yang digunakan, validasi *dataset* baik itu dengan partisi *dataset* atau *k-Cross Fold Validation* juga dapat mempengaruhi akurasi hasil prediksi.

1. Pendahuluan

Australia adalah sebuah benua dan negara yang terdiri dari beberapa negara bagian. Iklim di delapan negara bagian Australia sangat bervariasi. Sebagian besar wilayah Australia memiliki empat musim yaitu musim panas, musim dingin, musim gugur, dan musim semi [1]. Curah hujan di Australia sangat beragam, karena benua ini memiliki iklim yang beragam, dan Australia bagian utara memiliki iklim tropis. Kemudian, bagian barat daya dan pesisir selatan beriklim subtropis. Australia barat penuh dengan gurun, iklimnya kering, sangat panas di siang hari, tetapi sangat dingin di malam hari. Terakhir, pantai timur Australia memiliki iklim laut, sehingga daerah tersebut memiliki curah hujan yang cukup tinggi sepanjang tahun [2].

Prediksi hujan yang akurat merupakan salah satu tugas yang paling menantang dan penting di dunia saat ini [3], tidak terkecuali Australia karena memiliki iklim yang beragam [4]. Hal yang membuat menantang dikarenakan hujan bersifat dinamis dari fenomena iklim dan fluktuasi acak yang terlibat dalam proses fisik [5]. Biasanya prediksi hujan dibuat untuk beberapa periode waktu yang meliputi mingguan, bulanan dan prediksi musiman. Prediksi hujan yang tepat dan akurat sangat penting untuk perencanaan dan pengelolaan sumber daya air, peringatan banjir, kegiatan konstruksi dan operasi penerbangan serta yang lainnya [6]. Agar hasil prediksi hujan

menjadi optimal, berbagai kompleksitas perlu ditangani [7], seperti data statistik cuaca yang memiliki banyak fitur diantaranya kelembaban, tekanan, kecepatan angin, polutan, konsentrasi, dan lainnya.

Pendekatan *machine learning* telah banyak dilakukan dalam memprediksi hujan. Penerapannya yaitu dengan menggunakan metode komputasi dan memprediksi hujan dengan mengambil dan mengintegrasikan pengetahuan tersembunyi dari pola *linear* dan *non-linear* dari data cuaca masa lalu [8]. Teknik klasifikasi akan menghasilkan seperangkat aturan yang disebut aturan yang akan digunakan sebagai indikator untuk dapat memprediksi kelas data yang ingin diprediksi [9]. Beberapa penelitian telah banyak dilakukan untuk memprediksi hujan dengan pendekatan *machine learning* dengan menerapkan beragam algoritme klasifikasi diantaranya *Random Forest* [8][10], *Logistic Regression* [8][11] dan *Support Vector Regression* [11]. Bahkan pada penelitian [12], prediksi hujan dilakukan dengan membandingkan beragam *classifier* seperti *Naïve Bayes*, *Support Vector Machine*, *Decision Tree*, *Neural Network* dan *Random Forest*. Pada penelitian lainnya [13][14][15][16][17][18], *Decision Tree* telah banyak digunakan dalam prediksi hujan. Algoritme *Decision Tree* merupakan alat pengambilan keputusan terbaik yang bersifat sederhana dalam struktur, namun memiliki performa yang baik dalam analisis beberapa variabel [14]. Algoritme C5.0 juga telah diterapkan untuk memprediksi hujan [19]. Algoritme C5.0 adalah algoritme klasifikasi yang dapat diterapkan pada kumpulan data yang besar [19]. Algoritme C5.0 memiliki efisiensi dan memori yang lebih baik daripada algoritme C4.5 [20].

Data hujan di Australia dirilis oleh seorang ilmuwan data bernama Joe Young, dan telah diunggah ke situs kaggle dengan judul *dataset Rain in Australia* [21]. Kaggle adalah komunitas *online* ilmuwan data dan praktisi *machine learning* yang diakuisisi oleh *Google*. Selain menjadi komunitas, situs ini juga berisi kumpulan data (*dataset*) berbagai kasus yang ditemukan dari seluruh dunia [22]. Penelitian ini akan berfokus untuk memprediksi hujan di Australia dengan membandingkan akurasi kedua algoritme *Decision Tree* dan algoritme C5.0. Verifikasi hasil dibagi menjadi dua bagian, pertama, kumpulan data dibagi sebanyak 80% sebagai data latih dan 20% sebagai data uji. Kedua, *K-Fold Cross-Validation* digunakan sebanyak 10 kali untuk validasi.

2. Metode/Perancangan

Metode atau tahapan yang diterapkan dalam melakukan klasifikasi prediksi hujan di Australia yaitu melalui beberapa tahapan diantaranya Pengumpulan Data, *Data Pre-processing* (termasuk dilakukan penanganan *Missing Value* didalamnya), Pemodelan Klasifikasi dengan menerapkan dan membandingkan algoritme *Decision Tree* dan C5.0, Validasi Hasil menggunakan Partisi *Dataset* dan *10-Cross Fold Validation* serta Evaluasi Model menggunakan *Confussion Matrix*.

2.1. Pengumpulan Data

Dataset yang digunakan dalam penelitian ini adalah *Rain in Australia* dari kaggle.com *Repository*. *Dataset* diambil dari situs <https://www.kaggle.com/jsphyg/weather-dataset-rattle-package> yang merupakan data yang berisikan pengamatan cuaca harian dari berbagai stasiun cuaca Australia. *Dataset Rain in Australia* terdiri dari satu *file rar* yang berisi satu *file dataset* dengan nama *weatherAUS.csv*. Bentuk format *file* dalam *dataset* adalah *csv* atau *comma separated values* yang berarti data didalamnya dipisahkan oleh koma. *Dataset* tersebut terdiri dari 24 atribut dan 142.193 *record*. Penjelasan dari setiap atribut dapat dilihat pada **Tabel 1**.

Tabel 1. Penjelasan Atribut dalam *Dataset Rain in Australia*

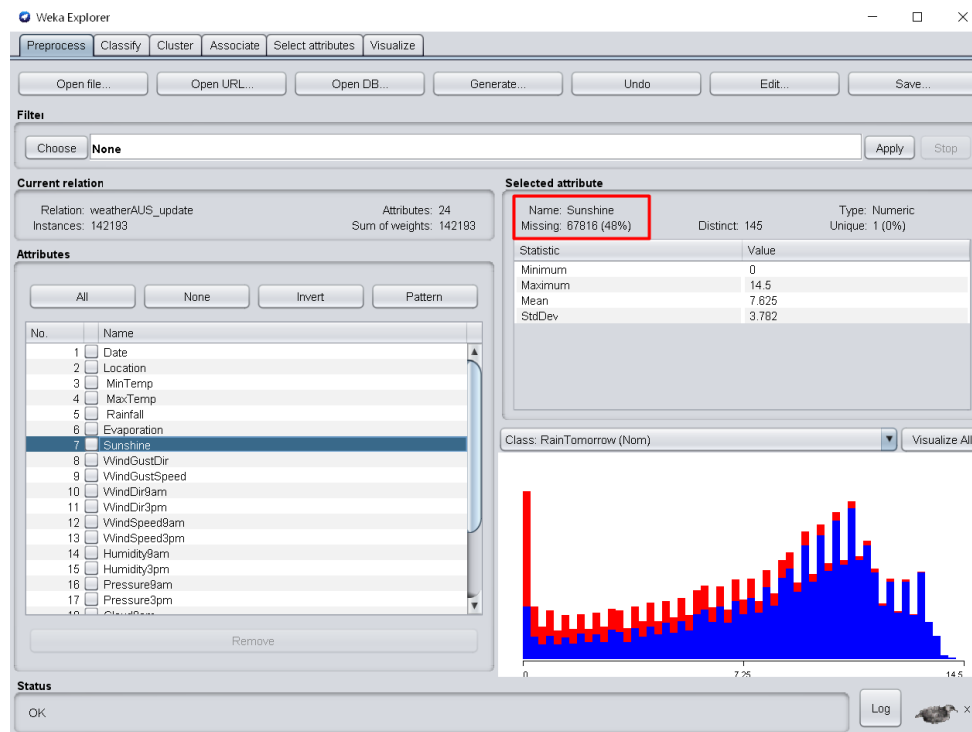
No.	Feature	Type Data	Deskripsi
1.	<i>Date</i>	<i>date</i>	Tanggal observasi prediksi pengamatan cuaca harian
2.	<i>Location</i>	<i>chr</i>	Nama umum lokasi stasiun cuaca
3.	<i>MinTemp</i>	<i>numeric</i>	Suhu minimum dalam derajat celsius
4.	<i>MaxTemp</i>	<i>numeric</i>	Suhu maksimum dalam derajat celsius
5.	<i>Rainfall</i>	<i>numeric</i>	Jumlah curah hujan yang tercatat untuk hari itu dalam mm
6.	<i>Evaporation</i>	<i>numeric</i>	Disebut juga penguapan kelas A (mm) dalam 24 jam hingga 9 pagi
7.	<i>Sunshine</i>	<i>numeric</i>	Jumlah jam sinar matahari cerah dalam satu hari
8.	<i>WindGustDir</i>	<i>chr</i>	Arah hembusan angin terkuat dalam 24 jam hingga tengah malam
9.	<i>WindGustSpeed</i>	<i>numeric</i>	Kecepatan (km/jam) hembusan angin terkuat dalam 24 jam hingga tengah malam
10.	<i>WindDir9am</i>	<i>chr</i>	Arah angin pada jam 9 pagi
11.	<i>WindDir3pm</i>	<i>chr</i>	Arah angin pada jam 3 sore
12.	<i>WindSpeed9am</i>	<i>numeric</i>	Kecepatan angin (km / jam) rata-rata lebih dari 10 menit sebelum jam 9 pagi
13.	<i>WindSpeed3pm</i>	<i>numeric</i>	Kecepatan angin (km / jam) rata-rata lebih dari 10 menit sebelum jam 3 sore
14.	<i>Humidity9am</i>	<i>numeric</i>	Kelembapan (persen) pada jam 9 pagi
15.	<i>Humidity3pm</i>	<i>numeric</i>	Kelembapan (persen) pada jam 3 sore
16.	<i>Pressure9am</i>	<i>numeric</i>	Tekanan atmosfer (hpa) berkurang hingga rata-rata permukaan laut pada pukul 9 pagi
17.	<i>Pressure3pm</i>	<i>numeric</i>	Tekanan atmosfer (hpa) berkurang hingga rata-rata permukaan laut pada pukul 3 sore
18.	<i>Cloud9am</i>	<i>numeric</i>	Bagian langit yang tertutup awan pada jam 9 pagi. Ini diukur dalam "oktas", yang merupakan satuan delapan. Ini mencatat berapa banyak awan
19.	<i>Cloud3pm</i>	<i>numeric</i>	Bagian langit yang tertutup awan pada jam 3 sore. Ini diukur dalam "oktas", yang merupakan satuan delapan. Ini mencatat berapa banyak awan
20.	<i>Temp9am</i>	<i>numeric</i>	Suhu (derajat C) pada jam 9 pagi
21.	<i>Temp3pm</i>	<i>numeric</i>	Suhu (derajat C) pada jam 3 sore
22.	<i>RainToday</i>	<i>chr</i>	Prediksi hujan hari ini
23.	<i>RISK MM</i>	<i>numeric</i>	Jumlah hujan hari berikutnya dalam mm. Digunakan untuk membuat variabel respons <i>RainTomorrow</i> . Semacam ukuran "risiko".
24.	<i>RainTomorrow</i>	<i>chr</i>	Variabel Target. Prediksi hujan besok

Selain itu, *dataset* tersebut memiliki *missing value* berkisar dari 1% hingga 48% yang terbesar pada beberapa atribut. Cara yang dapat dilakukan untuk menangani *missing value* adalah menghapus *tuple missing value* atau mengisinya dengan nilai konstanta berdasarkan rata-rata dari nilai dalam data tersebut [23]. *Missing value* dalam *dataset Rain in Australia* dapat dilihat pada **Tabel 2**.

Tabel 2. Total *Missing Value* pada Tiap Atribut

No.	Feature	Jumlah <i>Missing Value</i>	Persentase Terhadap Keseluruhan Data
1.	<i>Date</i>	0 <i>Record</i>	0%
2.	<i>Location</i>	0 <i>Record</i>	0%
3.	<i>MinTemp</i>	637 <i>Records</i>	>0%
4.	<i>MaxTemp</i>	322 <i>Records</i>	>0%
5.	<i>Rainfall</i>	1406 <i>Records</i>	1%
6.	<i>Evaporation</i>	60843 <i>Records</i>	43%
7.	<i>Sunshine</i>	67816 <i>Records</i>	48%
8.	<i>WindGustDir</i>	9330 <i>Records</i>	7%
9.	<i>WindGustSpeed</i>	9270 <i>Records</i>	7%
10.	<i>WindDir9am</i>	10013 <i>Records</i>	7%
11.	<i>WindDir3pm</i>	3778 <i>Records</i>	3%
12.	<i>WindSpeed9am</i>	1348 <i>Records</i>	1%
13.	<i>WindSpeed3pm</i>	2630 <i>Records</i>	2%
14.	<i>Humidity9am</i>	1774 <i>Records</i>	1%
15.	<i>Humidity3pm</i>	3610 <i>Records</i>	3%
16.	<i>Pressure9am</i>	14014 <i>Records</i>	10%
17.	<i>Pressure3pm</i>	13981 <i>Records</i>	10%
18.	<i>Cloud9am</i>	53657 <i>Records</i>	38%
19.	<i>Cloud3pm</i>	57094 <i>Records</i>	40%
20.	<i>Temp9am</i>	904 <i>Records</i>	1%
21.	<i>Temp3pm</i>	2726 <i>Records</i>	2%
22.	<i>RainToday</i>	1406 <i>Records</i>	1%
23.	<i>RISK MM</i>	0 <i>Record</i>	0%
24.	<i>RainTomorrow</i>	0 <i>Record</i>	0%

Saat menggunakan *WEKA*, terlihat atribut yang memiliki *missing value* di dalam *recordnya* seperti pada salah satu contoh atribut “*Sunshine*” seperti diperlihatkan **Gambar 1**.



Gambar 1. Missing Value pada Atribut Sunshine Menggunakan WEKA

2.2. Data Pre Processing

Dataset berekstensi *csv* yang terdapat pada arsip akan dilakukan tahap *pre processing* data agar lebih mudah dianalisis dan divisualisasikan. Dataset yang akan diproses pada tahapan ini adalah Dataset *Rain in Australia* yang terdiri dari 24 atribut [24].

Proses pembersihan data dari *noise* dapat dilakukan sesuai kebutuhan karena *noise* data banyak ragamnya. Contohnya, pada dataset *Rain in Australia* terdapat *record* yang memiliki *missing value* di beberapa atribut sehingga tidak akan optimal saat diproses untuk diekstraksi pada tahapan selanjutnya. Penanganan *missing value* akan dilakukan menggunakan aplikasi *RStudio*.

2.2.1. Penanganan Missing Value Menggunakan Rstudio

Penanganan *missing value* menggunakan *RStudio* melalui kode manual, dituliskan di GUI yang sudah disediakan. Tabel 2 menunjukkan tipe data pada setiap atribut dalam dataset. Terlihat atribut yang memiliki *missing value* didalamnya adalah atribut-atribut yang bertipe data numerik dan nominal, sehingga butuh perbedaan metode untuk menangani *missing value* pada atribut tersebut. Untuk menangani *missing value* pada tipe data numerik adalah dengan cara mengganti nilai yang hilang dengan nilai rata-rata dari seluruh data yang ada di dalam atributnya. Sedangkan pada tipe data nominal, *missing value* akan dihapuskan karena *record* dalam atribut tersebut terlalu beragam. Hasil penanganan *missing value* menggunakan *RStudio* dapat dilihat pada Gambar 2.

Date	Location	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindGustSpeed	WindDir9am	WindDir3pm	WindSpeed
1	2008-12-01	Albany	13.4	22.9	0.6	5.469824	7.624853	W	44	W	WNW
2	2008-12-02	Albany	7.4	25.1	0.0	5.469824	7.624853	WNW	44	NNW	WSW
3	2008-12-03	Albany	12.9	25.7	0.0	5.469824	7.624853	WSW	46	W	WSW
4	2008-12-04	Albany	9.2	28.0	0.0	5.469824	7.624853	NE	24	SE	E
5	2008-12-05	Albany	17.5	32.3	1.0	5.469824	7.624853	W	41	ENE	NW
6	2008-12-06	Albany	14.6	29.7	0.2	5.469824	7.624853	WNW	56	W	W
7	2008-12-07	Albany	14.3	25.0	0.0	5.469824	7.624853	W	50	SW	W
8	2008-12-08	Albany	7.7	26.7	0.0	5.469824	7.624853	W	35	SSE	W
9	2008-12-09	Albany	9.7	31.9	0.0	5.469824	7.624853	NNW	80	SE	NW
10	2008-12-10	Albany	13.1	30.1	1.4	5.469824	7.624853	W	28	S	SSE
11	2008-12-11	Albany	13.4	30.4	0.0	5.469824	7.624853	N	30	SSE	ESE
12	2008-12-12	Albany	15.9	21.7	2.2	5.469824	7.624853	NNE	31	NE	ENE
13	2008-12-13	Albany	15.9	18.6	15.6	5.469824	7.624853	W	61	NNW	NNW
14	2008-12-14	Albany	12.6	21.0	3.6	5.469824	7.624853	SW	44	W	SSW
15	2008-12-17	Albany	14.1	20.9	0.0	5.469824	7.624853	ENE	22	SSW	E
16	2008-12-18	Albany	13.5	22.9	16.8	5.469824	7.624853	W	63	N	WNW
17	2008-12-19	Albany	11.2	22.5	10.6	5.469824	7.624853	SSE	43	WSW	SW
18	2008-12-20	Albany	9.8	25.6	0.0	5.469824	7.624853	SSE	26	SE	NNW
19	2008-12-21	Albany	11.5	29.3	0.0	5.469824	7.624853	S	24	SE	SE
20	2008-12-22	Albany	17.1	33.0	0.0	5.469824	7.624853	NE	43	NE	N
21	2008-12-23	Albany	20.5	31.8	0.0	5.469824	7.624853	WNW	41	W	W
22	2008-12-24	Albany	15.3	30.9	0.0	5.469824	7.624853	N	33	ESE	NW
23	2008-12-25	Albany	12.6	32.4	0.0	5.469824	7.624853	W	43	E	W
24	2008-12-26	Albany	16.2	33.9	0.0	5.469824	7.624853	WSW	35	SE	WSW
25	2008-12-28	Albany	20.1	32.7	0.0	5.469824	7.624853	WNW	48	N	WNW
26	2008-12-29	Albany	19.7	27.2	0.0	5.469824	7.624853	WNW	46	NW	WSW
27	2008-12-30	Albany	12.5	24.2	1.2	5.469824	7.624853	WNW	50	WSW	SW
28	2008-12-31	Albany	12.0	24.4	0.8	5.469824	7.624853	W	39	WNW	WNW
29	2009-01-01	Albany	11.3	26.5	0.0	5.469824	7.624853	WNW	56	W	WNW
30	2009-01-02	Albany	9.6	23.9	0.0	5.469824	7.624853	W	41	WSW	SSW
31	2009-01-03	Albany	10.5	28.8	0.0	5.469824	7.624853	SSE	26	SSE	E
32	2009-01-04	Albany	12.3	34.6	0.0	5.469824	7.624853	WNW	37	SSE	NW
33	2009-01-05	Albany	12.9	35.8	0.0	5.469824	7.624853	WNW	41	ENE	NW

Gambar 2. Penanganan Missing Value pada Atribut Menggunakan Rstudio

Setelah dilakukan proses penanganan *missing value* pada seluruh atribut tersebut, *dataset* berkurang sebanyak 13% dari *dataset* awal. Sebelum *dataset* diolah, tipe data *chr* pada beberapa atribut akan diganti dengan atribut *factor* agar dapat bekerja dengan baik pada aplikasi *Rstudio*.

3. Hasil dan Pembahasan

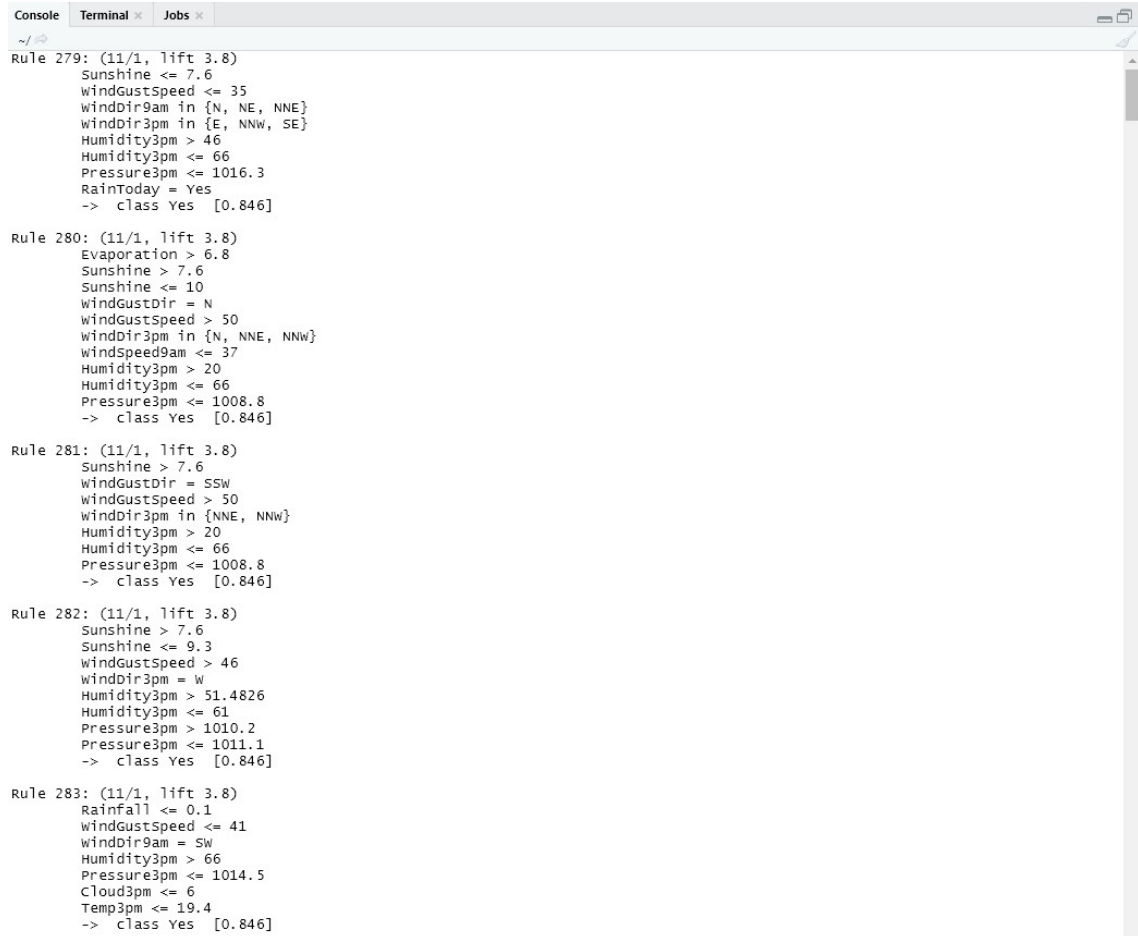
3.1. Processing and Modelling Dataset

Dataset yang sudah melalui *praprocessing* kemudian diproses menggunakan aplikasi *RStudio*. Algoritme *Decision Tree* akan diterapkan pada percobaan yang pertama. Validasi pertama adalah dengan mempartisi *dataset* ke dalam dua bagian, yaitu 80% untuk data latih dan 20% untuk data uji. Validasi kedua adalah dengan menggunakan *10-Cross Fold Validation*.

3.1.1. Algoritme Decision Tree dengan Partisi Dataset 80:20

Library Party yang terdapat pada *RStudio* dapat digunakan untuk menyelesaikan pemodelan algoritme *Decision Tree*. Sebanyak 80% partisi data latih dimodelkan dan diujikan terhadap 20% data yang sudah disiapkan dari *dataset*. *Decision Tree* yang dihasilkan dari pemodelan dapat dilihat pada Gambar 3.

Selain menghasilkan pohon keputusan, algoritme C5.0 juga menghasilkan aturan sebagai hasil dari pemodelan yang sudah dilakukan. Aturan tersebut berjumlah 373 aturan. Tampilan *Console* pada *RStudio* tidak memungkinkan untuk menampilkan hasil aturan tersebut, sehingga dapat dilihat sebagian contoh aturan tersebut pada **Gambar 5**.



```
Console Terminal Jobs
~/
Rule 279: (11/1, lift 3.8)
Sunshine <= 7.6
windGustSpeed <= 35
windDir9am in {N, NE, NNE}
windDir3pm in {E, NNW, SE}
Humidity3pm > 46
Humidity3pm <= 66
Pressure3pm <= 1016.3
RainToday = Yes
-> class Yes [0.846]

Rule 280: (11/1, lift 3.8)
Evaporation > 6.8
Sunshine > 7.6
Sunshine <= 10
windGustDir = N
windGustSpeed > 50
windDir3pm in {N, NNE, NNW}
windSpeed9am <= 37
Humidity3pm > 20
Humidity3pm <= 66
Pressure3pm <= 1008.8
-> class Yes [0.846]

Rule 281: (11/1, lift 3.8)
Sunshine > 7.6
windGustDir = SSW
windGustSpeed > 50
windDir3pm in {NNE, NNW}
Humidity3pm > 20
Humidity3pm <= 66
Pressure3pm <= 1008.8
-> class Yes [0.846]

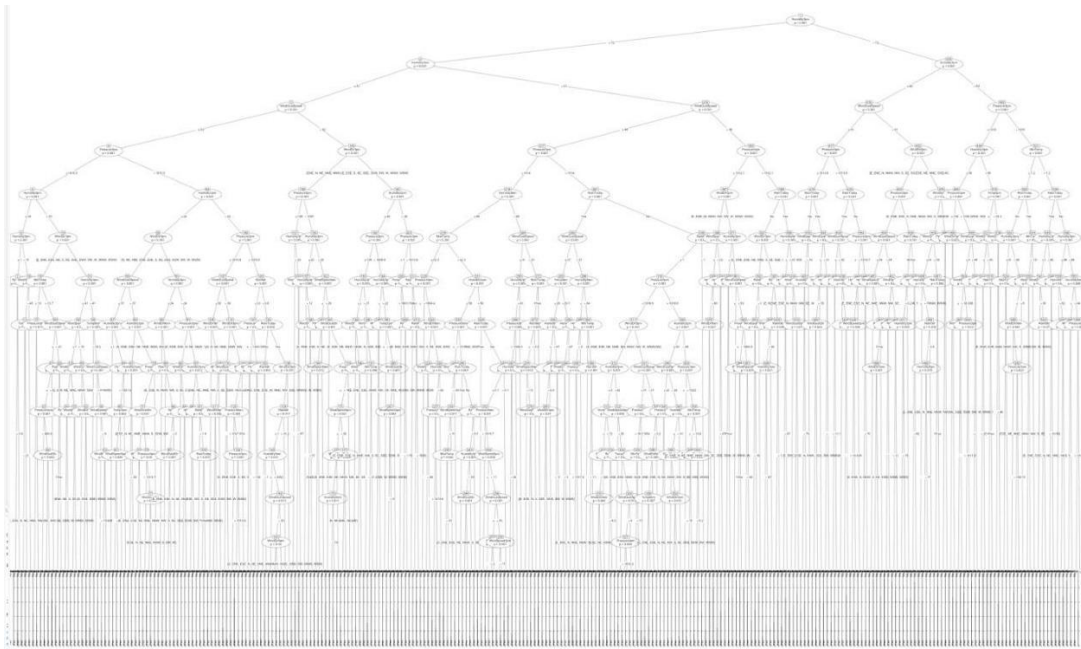
Rule 282: (11/1, lift 3.8)
Sunshine > 7.6
Sunshine <= 9.3
windGustSpeed > 46
windDir3pm = W
Humidity3pm > 51.4826
Humidity3pm <= 61
Pressure3pm > 1010.2
Pressure3pm <= 1011.1
-> class Yes [0.846]

Rule 283: (11/1, lift 3.8)
rainfall <= 0.1
windGustSpeed <= 41
windDir9am = SW
Humidity3pm > 66
Pressure3pm <= 1014.5
Cloud3pm <= 6
Temp3pm <= 19.4
-> class Yes [0.846]
```

Gambar 5. Aturan yang Dihasilkan Algoritme C5.0 dengan Partisi *Dataset* 80:20

3.1.4. Algoritme C5.0 dengan 10-Cross Fold Validation

Pemodelan berikutnya menggunakan algoritme C5.0 dengan teknik 10-Cross Fold Validation. Dengan *dataset* yang sama, *Decision Tree* yang dihasilkan dari pemodelan algoritme C5.0 dapat dilihat pada **Gambar 6**.



Gambar 6. Algoritme *Decision Tree* dan Algoritme C5.0 dengan Teknik *10-Cross Fold Validation*

Selain menghasilkan *Decision Tree*, algoritme C5.0 juga menghasilkan aturan sebagai hasil dari pemodelan yang sudah dilakukan. Aturan tersebut berjumlah 417 aturan. Tampilan *Console* pada *RStudio* tidak memungkinkan untuk menampilkan hasil aturan tersebut, sehingga dapat dilihat sebagian contoh aturan tersebut pada **Gambar 7**.

```
Console Terminal Jobs
~/
Pressure3pm <= 1015.258
Cloud3pm > 3
RainToday = Yes
-> class Yes [0.839]

Rule 325: (29/4, lift 3.7)
RainFall > 0.2
Sunshine > 1.6
windgustDir = SW
windgustSpeed > 48
windgustSpeed <= 57
windDir3pm in {E, S, W, WNW, WSW}
windspeed9am <= 31
windspeed3pm <= 31
Humidity3pm > 66
Humidity3pm <= 81
Pressure3pm > 1010.4
Temp9am > 4.6
-> class Yes [0.839]

Rule 326: (71/11, lift 3.7)
Sunshine <= 7
windgustDir = E
windgustSpeed > 54
Humidity3pm > 52
Pressure3pm <= 1014.3
-> class Yes [0.836]

Rule 327: (59/9, lift 3.7)
Sunshine <= 10
windgustDir = NNW
windgustSpeed > 52
windDir3pm in {N, NNE, NNW, NW}
windspeed3pm > 31
windspeed3pm <= 41
Humidity3pm > 24
Pressure3pm <= 1005.8
-> class Yes [0.836]

Rule 328: (125/20, lift 3.7)
Sunshine <= 7.5
windgustDir in {ENE, NNE, NW, S, SW, WSW}
windgustSpeed > 54
windDir9am = N
Humidity3pm > 14
Pressure3pm <= 1006.4
-> class Yes [0.835]

Rule 329: (119/19, lift 3.7)
Sunshine <= 7.3
windgustSpeed > 56
windDir9am = NNE
Humidity3pm > 14
Pressure3pm <= 1006.4
-> class Yes [0.835]

Rule 330: (10/1, lift 3.7)
```

Gambar 7. Aturan yang Dihasilkan Algoritme C5.0 dengan Teknik 10-Cross Fold Validation

3.2. Evaluation

Setelah pemodelan menggunakan metode algoritme *Decision Tree* dan Algoritme C5.0 selesai, selanjutnya akan dijelaskan mengenai *Confusion Matrix* dari masing-masing algoritme. Dalam sistem klasifikasi, *Confusion matrix* berisi informasi mengenai klasifikasi aktual dan prediksi. Performa dari sistem tersebut pada umumnya dilakukan evaluasi menggunakan data dalam matriks. Terdapat dua kelas pada *Confusion matrix* yang mengklasifikasikan kelas positif dan kelas negatif [25] ditunjukkan oleh **Tabel 3**.

Tabel 3. *Confusion Matrix*

		Aktual	
		Negatif	Positif
Prediksi	Negatif	a	b
	Positif	c	d

Keterangan:

a adalah jumlah prediksi yang benar bahwa sebuah *instance* negatif,

b adalah jumlah prediksi yang salah bahwa sebuah *instance* positif,

c adalah jumlah salah prediksi yang sebuah contoh negatif, dan

d adalah jumlah prediksi yang benar bahwa sebuah *instance* bernilai positif.

3.2.1. Akurasi Algoritme *Decision Tree* dan Algoritme C5.0

Confusion Matrix dari setiap pemodelan yang sudah dilakukan, yaitu dari algoritme *Decision Tree* dan algoritme C5.0 dapat dilihat pada **Tabel 4-9**.

Tabel 4. *Confusion Matrix* Algoritme *Decision Tree* Split 80:20 Terhadap Data Uji

Accuracy: 84,55%

	No	Yes
No	73206	10535
Yes	3962	11397

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} = \frac{11397+73206}{73206+10535+3962+11397} = 0.853713421$$

Tabel 5. *Confusion Matrix* Algoritme *Decision Tree* 10-Cross Fold Validation

Accuracy: 87,35%

	No	Yes
No	9690	1111
Yes	454	1116

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} = \frac{1116 + 9690}{9690 + 1111 + 454 + 1116} = 0.873494463$$

Tabel 6. *Confusion Matrix* Algoritme C5.0 Split 80:20 Terhadap Data Uji (*Tree-Based Model*)

Accuracy: 83,99%

	No	Yes
No	17864	2654
Yes	1286	2806

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} = \frac{2806 + 17864}{2806 + 17864 + 1286 + 2654} = 0.839902479$$

Tabel 7. Confusion Matrix Algoritme C5.0 Split 80:20 Terhadap Data Uji (Rule-Based Model)

Accuracy: 84,58%

	No	Yes
No	18018	2664
Yes	1132	2796

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} = \frac{2796 + 18018}{18018 + 2664 + 1132 + 2796} = 0.845753759$$

Tabel 8. Confusion Matrix Algoritme C5.0 10-Cross Fold Validation (Tree-Based Model)

Accuracy: 85,95%

	No	Yes
No	9475	1069
Yes	669	1158

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} = \frac{1158 + 9475}{9475 + 1069 + 669 + 1158} = 0.859510145$$

Tabel 9. Confusion Matrix Algoritme C5.0 10-Cross Fold Validation (Rule-Based Model)

Accuracy: 86,85%

	No	Yes
No	9608	1091
Yes	536	1136

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} = \frac{1136 + 9608}{9608 + 1091 + 536 + 1136} = 0.868482742$$

Semua hasil tersebut dapat digabungkan dalam satu tabel, yaitu **Tabel 10**.

Tabel 10. *Confusion Matrix* Algoritme *Decision Tree* dan Algoritme C5.0

No.	Algoritme	Akurasi
1.	<i>Decision Tree (Split 80:20)</i>	84.55%
2.	<i>Decision Tree (10-Cross Fold Validation)</i>	87.35%
3.	<i>C5.0 Tree-Based Model (Split 80:20)</i>	83.99%
4.	<i>C5.0 Rule-Based Model (Split 80:20)</i>	84.58%
5.	<i>C5.0 Tree-Based Model (10-Cross Fold Validation)</i>	85.95%
6.	<i>C5.0 Rule-Based Model (10-Cross Fold Validation)</i>	86.85%

3.2.2. Evaluasi Model Split 80:20 dengan 10-Cross Fold Validation

Pada Tabel 10 dapat disimpulkan bahwa dalam kasus klasifikasi prediksi hujan di Australia, skema model menggunakan *10-Cross Fold Validation* lebih unggul dalam semua kasus dibandingkan dengan model split 80:20. Keunggulan *10-Cross Fold Validation* dialami oleh kedua algoritme yaitu algoritme *Decision Tree* dan algoritme C5.0.

3.2.3. Hasil Aturan dari Model yang Memiliki Akurasi Paling Tinggi yaitu C5.0 10-Cross Fold Validation

Hasil aturan dari model algoritme C5.0 dengan Teknik *10-Cross Fold Validation (Tree- Based Model)* dapat dilihat pada **Gambar 8**.

```

Console Terminal Jobs
~/
SubTree [s1]
windGustSpeed <= 37: No (5/1)
windGustSpeed > 37: Yes (25/3)
SubTree [s2]
windGustSpeed <= 37: No (11)
windGustSpeed > 37: Yes (5/1)
SubTree [s3]
windGustSpeed <= 41: Yes (3)
windGustSpeed > 41: No (2)
SubTree [s4]
windGustSpeed <= 33: No (3)
windGustSpeed > 33: Yes (8)
SubTree [s5]
windDir9am in {E,ENE,ESE,N,NE,NNE,S,SE,SSE,SSW,W,WNW,WSW}: No (8)
windDir9am in {NNW,NW,SW}: Yes (3)
SubTree [s6]
Humidity9am <= 89: No (7)
Humidity9am > 89: Yes (5/1)
SubTree [s7]
windDir9am in {ESE,SSE,SSW,SW}: No (4)
windDir9am in {E,ENE,N,NE,NNE,NNW,NW,SE,S,SE,SSE,SSW,W,WNW,WSW}: Yes (4)
windDir9am = S:
...Pressure9am <= 1025.1: Yes (3)
...Pressure9am > 1025.1: No (2)
SubTree [s8]
Sunshine <= 4.7: Yes (4/1)
Sunshine > 4.7: No (8)
SubTree [s9]
windGustDir in {E,ENE,ESE,N,NE,NNE,NNW,NW,S,SE,SSE,SW,WNW}: Yes (10)
windGustDir in {SSW,W,WSW}:
...Humidity3pm > 85: No (8)
Humidity3pm <= 85:
...Wintemp <= 13.5: Yes (3)
...Wintemp > 13.5: No (2)
SubTree [s10]
windGustDir in {E,ENE,ESE,NE,NNE,NNW,NW,S,SSW,SW,W,WNW}: Yes (12/2)
windGustDir in {N,SE,SSE,WSW}: No (9/1)
    
```

Gambar 8. Contoh Aturan yang Dihasilkan Algoritme C5.0 dengan Teknik *10-Cross FoldValidation (Tree-Based Model)*

Sedangkan hasil aturan dari model algoritme C5.0 dengan Teknik 10-Cross Fold Validation (Rule-Based Model) dapat dilihat pada **Gambar 9**.

```

Console Terminal Jobs
~/
Pressure3pm <= 1015.258
Cloud3pm > 3
RainToday = Yes
-> c class Yes [0.839]

Rule 325: (29/4, 11ft 3.7)
Rainfall > 0.2
Sunshine > 1.6
windGustDir = SW
windGustSpeed > 48
windGustSpeed <= 57
windDir3pm in {E, S, W, WNW, WSW}
windSpeed9am <= 31
windSpeed3pm <= 31
Humidity3pm > 66
Humidity3pm <= 81
Pressure3pm > 1010.4
Temp9am > 4.6
-> c class Yes [0.839]

Rule 326: (71/11, 11ft 3.7)
Sunshine <= 7
windGustDir = E
windGustSpeed > 54
Humidity3pm > 52
Pressure3pm <= 1014.3
-> c class Yes [0.836]

Rule 327: (59/9, 11ft 3.7)
sunshine <= 10
windGustDir = NNW
windGustSpeed > 52
windDir3pm in {N, NNE, NNW, NW}
windSpeed3pm > 31
windSpeed3pm <= 41
Humidity3pm > 24
Pressure3pm <= 1005.8
-> c class Yes [0.836]

Rule 328: (125/20, 11ft 3.7)
Sunshine <= 7.5
windGustDir in {ENE, NNE, NW, S, SW, WSW}
windGustSpeed > 54
windDir9am = N
Humidity3pm > 14
Pressure3pm <= 1006.4
-> c class Yes [0.835]

Rule 329: (119/19, 11ft 3.7)
Sunshine <= 7.3
windGustSpeed > 56
windDir9am = NNE
Humidity3pm > 14
Pressure3pm <= 1006.4
-> c class Yes [0.835]

Rule 330: (10/1, 11ft 3.7)

```

Gambar 9. Contoh Aturan yang Dihasilkan Algoritme C5.0 dengan Teknik 10-Cross FoldValidation (Rule-Based Model)

3.2.4. Aturan-Aturan Berpengaruh dari Model yang Memiliki Akurasi Paling Tinggi yaitu C5.0 10-Cross Fold Validation

Aturan-aturan yang berpengaruh dalam menentukan label kelas berdasarkan model berbasis aturan dan model berbasis *tree* dapat dilihat pada **Gambar 10**.

	id LHS	RHS	support	confidence
	<int> <chr>	<chr> <int>	<dbl>	<
1	1 Rainfall <= 0.2 & sunshine > 10.3 & windDir3pm %in% c('E'~	No	12771	0.975
2	2 MinTemp <= 10.6 & Rainfall <= 0.5 & Sunshine > 8.2 & Pres~	No	7231	0.974
3	3 Rainfall <= 2.9 & windGustSpeed <= 46 & windDir9am == 'S'~	No	71	0.973
4	4 windGustSpeed <= 41 & windDir9am == 'SE' & Cloud3pm <= 3	No	1438	0.956
5	5 Rainfall <= 2.9 & windGustSpeed <= 46 & windDir3pm == 'NW~	No	428	0.942
6	6 MinTemp <= 19.9 & windGustDir == 'S' & windDir3pm %in% c(~	No	75	0.935
7	7 Rainfall <= 2.9 & windGustDir %in% c('SSW', 'W', 'WSW') &~	No	1625	0.935
8	8 windDir9am == 'SW' & Humidity3pm <= 95 & Pressure3pm > 10~	No	146	0.932
9	9 MinTemp <= 4.4 & windDir9am == 'E' & Pressure3pm > 1014.9~	No	723	0.932
10	10 Rainfall <= 2.2 & WindGustSpeed <= 26 & Cloud3pm <= 5	No	9569	0.926

Gambar 10. Aturan-Aturan Berpengaruh dari Model yang Memiliki Akurasi Paling Tinggi yaitu C5.0 10-Cross Fold Validation

4. Kesimpulan dan Saran

Berdasarkan hasil yang diperoleh pada penelitian, evaluasi menggunakan 10-Cross Fold Validation lebih unggul yaitu memiliki akurasi paling tinggi sebesar 87.35% untuk algoritme

Decision Tree dan akurasi sebesar 86.85% untuk algoritme *C5.0 Rule-Based Model*, dibandingkan dengan metode Split 80:20 pada kasus prediksi hujan di Australia. Meskipun hasil akurasi yang diperoleh untuk kedua metode tersebut rata-rata di atas 80.00%.

Untuk penelitian selanjutnya, disarankan menggunakan algoritme klasifikasi lainnya untuk ikut dibandingkan dalam prediksi hujan di Australia. Sehingga dapat diperoleh hasil akurasi prediksi yang lebih baik dari penelitian sebelumnya.

Daftar Pustaka

- [1] Anonim, "Cuaca di Australia," *Tourism Australia*, 2020. <http://australiaxy.com/id/id/facts-and-planning/weather-in-australia.html>.
- [2] P. Yasmin, "Karakteristik, Iklim dan Daftar Negara di Benua Australia," *Detiktravel*. Detik.com, 2020, [Online]. Available: <https://travel.detik.com/travel-news/d-5158330/karakteristik-iklim-dan-daftar-negara-di-benua-australia>.
- [3] C. Thirumalai, K. S. Harsha, M. L. Deepak, and K. C. Krishna, "Heuristic prediction of rainfall using machine learning techniques," in *Proceedings - International Conference on Trends in Electronics and Informatics, ICEI 2017*, 2018, vol. 2018-January, pp. 1114–1117, doi: 10.1109/ICOEI.2017.8300884.
- [4] A. M. Bagirov, A. Mahmood, and A. Barton, "Prediction of monthly rainfall in Victoria, Australia: Clusterwise linear regression approach," *Atmos. Res.*, vol. 188, pp. 20–29, 2017, doi: 10.1016/j.atmosres.2017.01.003.
- [5] A. M. Bagirov and A. Mahmood, "A Comparative Assessment of Models to Predict Monthly Rainfall in Australia," *Water Resour. Manag.*, vol. 32, no. 5, pp. 1777–1794, 2018, doi: 10.1007/s11269-018-1903-y.
- [6] S. Aftab, M. Ahmad, N. Hameed, M. S. Bashir, I. Ali, and Z. Nawaz, "Rainfall prediction in Lahore City using data mining techniques," in *International Journal of Advanced Computer Science and Applications*, 2018, vol. 9, no. 4, pp. 254–260, doi: 10.14569/IJACSA.2018.090439.
- [7] M. P. Darji, V. K. Dabhi, and H. B. Prajapati, "Rainfall forecasting using neural network: A survey," in *Conference Proceeding - 2015 International Conference on Advances in Computer Engineering and Applications, ICACEA 2015*, 2015, no. December, pp. 706–713, doi: 10.1109/ICACEA.2015.7164782.
- [8] G. Sethupathi M, Y. S. Ganesh, and M. M. Ali, "Efficient Rainfall Prediction and Analysis using Machine Learning Techniques," *Turkish J. Comput. Math. Educ.*, vol. 12, no. 6, pp. 3467–3474, 2021.
- [9] D. A. Kristiyanti, E. Purwaningsih, E. Nurelasari, A. Al Kaafi, and A. H. Umam, "Implementation of Neural Network Method for Air Quality Forecasting in Jakarta Region," *J. Phys. Conf. Ser.*, vol. 1641, no. 1, 2020, doi: 10.1088/1742-6596/1641/1/012037.
- [10] A. Dikshit, B. Pradhan, and A. M. Alamri, "Short-term spatio-temporal drought

- forecasting using random forests model at New South Wales, Australia,” *Appl. Sci.*, vol. 10, no. 12, 2020, doi: 10.3390/app10124254.
- [11] V. A. Vuyyuru, G. Apparao, and S. Anuradha, “Prediction Of Rainfall With A Machine Learning Approach,” *Turkish J. Comput. Math. Educ.*, vol. 12, no. 7, pp. 1762–1776, 2021.
- [12] S. Zainudin, D. S. Jasim, and A. A. Bakar, “Comparative analysis of data mining techniques for malaysian rainfall prediction,” *Int. J. Adv. Sci. Eng. Inf. Technol.*, vol. 6, no. 6, pp. 1148–1153, 2016, doi: 10.18517/ijaseit.6.6.1487.
- [13] M. Marjanović, M. Krautblatter, B. Abolmasov, U. Đurić, C. Sandić, and V. Nikolić, “The rainfall-induced landsliding in Western Serbia: A temporal prediction approach using Decision Tree technique,” *Eng. Geol.*, vol. 232, no. February 2017, pp. 147–159, 2018, doi: 10.1016/j.enggeo.2017.11.021.
- [14] A. Geetha and G. M. Nasira, “Data mining for meteorological applications: Decision trees for modeling rainfall prediction,” in *2014 IEEE International Conference on Computational Intelligence and Computing Research, IEEE ICCIC 2014*, 2015, pp. 0–3, doi: 10.1109/ICCIC.2014.7238481.
- [15] R. N, S. S, and K. S, “Comparison of Decision Tree Based Rainfall Prediction Model with Data Driven Model Considering Climatic Variables,” *Irrig. Drain. Syst. Eng.*, vol. 05, no. 03, 2016, doi: 10.4172/2168-9768.1000175.
- [16] J. Dou *et al.*, “Assessment of advanced random forest and decision tree algorithms for modeling rainfall-induced landslide susceptibility in the Izu-Oshima Volcanic Island, Japan,” *Sci. Total Environ.*, vol. 662, pp. 332–346, 2019, doi: 10.1016/j.scitotenv.2019.01.221.
- [17] W. Wei, Z. Yan, and P. D. Jones, “A decision-tree approach to seasonal prediction of extreme precipitation in eastern China,” *Int. J. Climatol.*, vol. 40, no. 1, pp. 255–272, 2020, doi: 10.1002/joc.6207.
- [18] N. Prasad, P. Kumar, and M. M. Naidu, “An approach to prediction of precipitation using Gini Index in SLIQ decision tree,” *Proc. - Int. Conf. Intell. Syst. Model. Simulation, ISMS*, pp. 56–60, 2013, doi: 10.1109/ISMS.2013.27.
- [19] E. Kurniawan, F. Nhita, A. Aditsania, and D. Saepudin, “C5.0 algorithm and synthetic minority oversampling technique (SMOTE) for rainfall forecasting in bandung regency,” *2019 7th Int. Conf. Inf. Commun. Technol. ICoICT 2019*, vol. 4, pp. 1–5, 2019, doi: 10.1109/ICoICT.2019.8835324.
- [20] N. Patil, R. Lathi, and V. Chitre, “Customer Card Classification Based on C5 . 0 & CART Algorithms,” *Int. J. Eng. Res. Appl.*, vol. 66, no. 3, pp. 37–39, 2012.
- [21] J. Young, “Rain in Australia,” *Kaggle.com*, 2018. <https://www.kaggle.com/jsphyg/weather-dataset-rattle-package>.
- [22] J. M. Frederic Lardinois, Matthew Lynley, “Google is acquiring data science community Kaggle,” 2017. .

- [23] I. Aprillani, A. A. Suryani, F. T. Informatika, and U. Telkom, “Analisis Penanganan Missing Value Dengan Metode Predictive Mean Matching (Pmm),” 2012.
- [24] State of New York, “New York State Index Crimes,” *Kagle.com*, 2019. <https://www.kaggle.com/new-york-state/new-york-state-index-crimes/metadata>.
- [25] F. Provost and T. Fawcett, “Analysis and Visualization of Classifier Performance: Comparison under Imprecise Class and Cost Distributions,” in *THE THIRD INTERNATIONAL CONFERENCE ON KNOWLEDGE DISCOVERY AND DATA MINING*, 1997, pp. 43–48.