

PENERAPAN ALGORITMA KLASIFIKASI DATA MINING C4.5 PADA DATASET CUACA WILAYAH BEKASI

Adhika Novandya

AMIK BSI Bekasi
e-mail: adhika.avn@bsi.ac.id

Abstrak

Cuaca adalah keadaan udara pada saat tertentu dan pada wilayah tertentu yang relatif sempit dan pada jangka waktu yang singkat. Prakiraan cuaca pada umumnya sering disebut peramalan cuaca yang merupakan penggunaan ilmu dan teknologi untuk memperkirakan atmosfer bumi pada masa akan datang untuk suatu tempat tertentu. Data yang digunakan pada penelitian didapat dari *World Weather Online*, merupakan sebuah situs yang memberikan data dan informasi mengenai kondisi cuaca sehari-hari. Data yang dipakai memiliki interval waktu setiap 3 jam terhitung mulai tanggal 12 Agustus 2016 pukul 01.00 hingga tanggal 20 Agustus 2016 pukul 22.00. Penelitian bertujuan untuk mendapatkan pola klasifikasi cuaca dengan menggunakan algoritma klasifikasi data mining yaitu algoritma C4.5. Hasil pengujian algoritma C4.5 menggunakan *10-fold cross validation* dan dibuktikan dengan pembuatan aplikasi web untuk pengujian sehingga menghasilkan nilai akurasi sebesar 88.89%.

Keywords: Cuaca, Prakiraan Cuaca, Data Mining, Algoritma Klasifikasi

1. Pendahuluan

Cuaca dan iklim merupakan dua kondisi yang hampir sama tetapi berbeda pengertian, khususnya terhadap kurun waktu. Cuaca merupakan bentuk awal yang dihubungkan dengan penafsiran dan pengertian akan kondisi fisik udara sesaat pada suatu lokasi dan suatu waktu, sedangkan iklim merupakan kondisi lanjutan dan merupakan kumpulan dari kondisi cuaca yang kemudian disusun dan dihitung dalam bentuk rata-rata kondisi cuaca dalam kurun waktu tertentu (Winarso, 2003). Cuaca adalah keadaan udara pada saat tertentu dan pada wilayah tertentu yang relatif sempit dan pada jangka waktu yang singkat. Menurut *World Climate Conference*, cuaca adalah keadaan atmosfer secara menyeluruh termasuk perubahan, perkembangan, dan menghilangnya suatu fenomena.

Banyaknya parameter dalam menentukan suatu cuaca menyebabkan ketepatan dan kecepatan dalam memprediksikan cuaca kurang terpenuhi (A. Joshi, 2015). Metode klasifikasi data mining merupakan sebuah teknik yang dilakukan untuk memprediksi *class* atau properti dari data itu sendiri (Larose, 2006). Adapun metode klasifikasi data mining memiliki beberapa algoritma salah satunya yaitu algoritma C4.5.

Penelitian ini bertujuan untuk mendapatkan pola klasifikasi cuaca yang terjadi dengan jumlah interval selama 3 jam sekali, sehingga nantinya dapat digunakan untuk memprediksi cuaca pada keesokan harinya atau dalam periode waktu tertentu.

2. Metode Penelitian Penelitian Sebelumnya

Proses prakiraan cuaca memerlukan banyak komponen data cuaca, jumlah data yang besar serta kemampuan prakirawan. Hal tersebut menyebabkan ketepatan dan kecepatan prakiraan kurang terpenuhi. Untuk memecahkan masalah tersebut, dilakukan penelitian model prediksi menggunakan beberapa teknik data mining yaitu *Association rule*, C4.5, *Classification* dan *Random Forest*. Penelitian menghasilkan bahwa model prediksi C4.5 memiliki tingkat akurasi 68.5% (Mujiasih, 2011).

Peramalan cuaca merupakan suatu proses memprediksikan bagaimana kondisi atmosfer berubah. Untuk memprediksi suatu cuaca digunakan algoritma *decision tree* untuk mengklasifikasikan parameter cuaca seperti temperatur maksimum, temperatur minimum, curah hujan, penguapan, dan kecepatan angin dengan menggunakan data dari situs cuaca *wonderground* mulai dari

tahun 2001 sampai 2013. Hasilnya didapat bahwa parameter tersebut mempunyai pengaruh yang berarti (A. Joshi, 2015). Peramalan cuaca adalah aplikasi yang paling penting dalam meteorologi dan telah menjadi salah satu yang paling ilmiah dan menjadi permasalahan teknologi yang menantang. Algoritma klasifikasi pohon keputusan C5 digunakan untuk menghasilkan pohon keputusan dan aturan klasifikasi parameter cuaca pada data yang didapat dari stasiun meteorologi Ibadan dari tahun 2000 sampai 2009, Artificial Neural Networks (ANN) dapat mendeteksi hubungan antara variabel input dan menghasilkan output berdasarkan pola observasi data (Olaiya, 2012).

Algoritma *FP Growth* digunakan untuk menghasilkan pohon keputusan. Data yang digunakan didapat dari departemen cuaca Nagpur periode 2010 sampai dengan 2014. Algoritma *FP growth* dengan evaluasi MAE, MSE, dan SD menampilkan hasil yang lebih akurat dibandingkan dengan algoritma Neural Net (NN), dimana *FP Growth* menghasilkan prediksi curah hujan yang benar setiap bulannya (Mohod, 2015).

Data Mining

Data mining adalah proses penting dimana metode kecerdasan diaplikasikan untuk mengekstrak pola data (J. Han, 2012).

Data mining adalah analisis observasional sekumpulan data untuk menemukan hubungan tidak terduga dan untuk meringkas data dengan cara baru yang dapat dipahami dan berguna bagi pemilik data (Larose, 2006).

Klasifikasi

Klasifikasi merupakan suatu proses menemukan kumpulan pola atau fungsi yang mendeskripsikan serta memisahkan kelas data yang satu dengan yang lainnya untuk menyatakan objek tersebut masuk pada kategori tertentu yang sudah ditentukan.

Klasifikasi adalah bentuk analisis data yang mengekstrak model yang menggambarkan kelas data (J. Han, 2012).

Algoritma C4.5

Algoritma C4.5 adalah ekstensi Quinlan untuk algoritma ID3 untuk menghasilkan pohon keputusan, algoritma C4.5 rekursif mengunjungi setiap node keputusan, memilih split optimal sampai tidak ada perpecahan lanjut yang memungkinkan (Larose, Discovering knowledge in data, 2005).

Pada dasarnya konsep dari algoritma C4.5 adalah mengubah data menjadi pohon keputusan dan aturan-aturan keputusan (*rule*). C4.5 adalah algoritma yang cocok untuk masalah klasifikasi dan data mining. C4.5 memetakan nilai atribut menjadi kelas yang dapat diterapkan untuk klasifikasi baru (Xindong, 2009).

Ada beberapa tahapan dalam membangun sebuah pohon keputusan dengan algoritma C4.5 yaitu (Kusrini, 2009).

1. Menyiapkan data *training*. Data *training* biasanya diambil dari data histori yang pernah terjadi sebelumnya dan sudah dikelompokkan ke dalam kelas-kelas tertentu.
2. Menentukan akar dari pohon. Akar akan diambil dari atribut yang terpilih, dengan cara menghitung nilai *gain* dari masing-masing atribut, nilai *gain* yang paling tinggi yang akan menjadi akar pertama.
3. Sebelum menghitung nilai *gain* dari atribut, hitung dahulu nilai entropi. Untuk menghitung nilai entropi digunakan rumus:

$$Entropy(S) = - \sum_{i=1}^N p_i \log_2 p_i \quad (1)$$

Dimana:

S = Himpunan Kasus
n = Jumlah partisi S
Pi = Proporsi Si terhadap S

4. Kemudian hitung nilai *gain* yang menggunakan rumus:

$$Gain(S, A) = entropy(S) - \sum_{i=1}^n \frac{|S_i|}{S} * entropy(S_i) \quad (2)$$

Dimana:

S = Himpunan Kasus
A = Fitur
n = Jumlah Partisi Atribut A
|Si| = Proporsi Si terhadap S
|S| = Jumlah Kasus dalam S

5. Ulangi langkah ke-2 hingga semua *record* terpartisi.
6. Proses partisi pohon keputusan akan berhenti saat:
 - a. Semua *record* dalam simpul N mendapat kelas yang sama.
 - b. Tidak ada atribut di dalam *record* yang dipartisi lagi.
 - c. Tidak ada *record* di dalam cabang

yang kosong.

Model Validasi

Penelitian menggunakan *stratified 10-fold cross-validation* untuk melakukan pengujian terhadap *dataset*. Peneliti melakukan 10 kali pengujian terhadap data untuk melihat performa dari masing-masing algoritma klasifikasi yang digunakan. Adapun bentuk model *stratified 10 fold cross validation* dapat dilihat pada gambar dibawah ini:

Tabel 1. Stratified 10 Fold Cross Validation (R. S. Wahono, 2014).

n-validation	Dataset's Partiton
1	
2	
3	
4	
5	
6	
7	
8	
9	
10	

Model Evaluasi

Model evaluasi yang digunakan oleh peneliti yaitu *Confusion Matrix* yang menghasilkan nilai akurasi dari validasi algoritma terhadap *dataset* yang ada. *Confusion Matrix* adalah alat visualisasi yang biasa digunakan pada *supervised learning*. Tiap kolom pada matriks adalah contoh kelas prediksi, sedangkan tiap baris mewakili kejadian di kelas yang sebenarnya (Goronescu, 2011). Hasil dari proses perhitungan *confusion matrix* yaitu 4 keluaran diantaranya *recall*, *precision*, *accuracy*, dan *error rate*. Dapat dilihat pada tabel dibawah ini:

Tabel 2. Confusion Matrix

	Prediksi			
	Negatif	A	C	
Aktual	Positif	B	D	

Keterangan:

1. A = jumlah prediksi yang tepat bersifat negatif.
2. B = jumlah prediksi yang salah bersifat positif.
3. C= jumlah prediksi yang salah bersifat negatif.
4. D = Jumlah prediksi yang tepat bersifat positif.

Beberapa persyaratan yang telah didefinisikan untuk matrik klasifikasi diantaranya sebagai berikut:

1. *Accuracy* merupakan proporsi jumlah prediksi benar. Rumus akurasi adalah:

$$AC = (A + D) / A + B + C + D$$

2. *Recall* atau tingkat positif benar (TP) adalah proporsi kasus positif yang diidentifikasi dengan benar, yang dapat dihitung dengan persamaan:

$$TP = D / C + D$$

3. Tingkat positif salah (FP) adalah proporsi kasus negatif yang salah diklasifikasikan sebagai positif, yang dapat dihitung dengan menggunakan persamaan:

$$FP = B / A + B$$

4. Tingkat negatif sejati (TN) didefinisikan sebagai proporsi kasus negatif yang diklasifikasikan dengan benar, dapat dihitung dengan menggunakan persamaan:

$$TN = A / A + B$$

5. Tingkat negatif palsu (FN) adalah proporsi kasus positif yang salah diklasifikasikan sebagai negatif, yang dihitung dengan menggunakan persamaan:

$$FN = C / C + D$$

6. *Precision* (P) adalah proporsi prediksi kasus positif yang benar, yang dihitung dengan menggunakan persamaan:

$$P = D / B + D$$

3. Pembahasan

3.1. Dataset

Dataset yang digunakan pada penelitian ini dapat dilihat pada gambar 1.

ID	Date	Time	Desc	Weather	Temp (Celcius)	Rain (mm)	Wind (mph)	Dir	Cloud (%)	Humidity (%)	Pressure (hPa)
1	22 August 2016	1	night	clear	26	0.0	5.4	0	0	80	1014
2	22 August 2016	4	morning	clear	25	0.0	4.5	0	0	82	1012
3	22 August 2016	7	morning	clear	25	0.0	4.5	0	0	75	1011
4	22 August 2016	10	daylight	clear	26	0.0	5.4	0	0	69	1012
5	22 August 2016	13	daylight	clear	26	0.0	6.3	0	0	61	1012
6	22 August 2016	16	afternoon	partly cloudy	25	0.0	12.3	0	10	57	1010
7	22 August 2016	19	night	partly rain-morally	24	0.4	7.5	0	14	52	1011
8	22 August 2016	22	night	partly rain-morally	27	0.4	5.4	0	16	50	1010
9	22 August 2016	1	night	clear	27	0.0	4.5	0	0	49	1012
10	22 August 2016	4	morning	clear	28	0.0	4.5	0	0	40	1020
11	22 August 2016	7	morning	clear	28	0.0	5.4	0	0	36	1018
12	22 August 2016	10	daylight	clear	28	0.0	5.4	0	0	32	1018
13	22 August 2016	13	daylight	partly rain-morally	25	0.0	6.3	0	0	40	1016
14	22 August 2016	16	afternoon	light rain shower	24	0.7	6.3	0	0	40	1014
15	22 August 2016	19	night	light rain shower	25	0.0	6.3	0	0	36	1015
16	22 August 2016	22	night	partly rain-morally	27	0.0	5.4	0	17	31	1012
17	22 August 2016	1	night	partly cloudy	28	0.0	5.4	0	0	24	1011
18	22 August 2016	4	morning	partly rain-morally	28	0.4	6.3	0	0	20	1010
19	22 August 2016	7	morning	partly rain-morally	26	0.4	5.4	0	0	19	1011
20	22 August 2016	10	daylight	clear	26	0.0	5.4	0	0	16	1012
21	22 August 2016	13	daylight	clear	26	0.0	5.4	0	0	12	1010

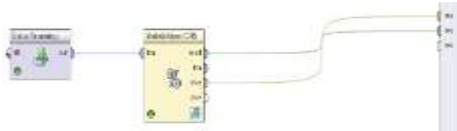
Gambar 1. Dataset

Dataset yang diolah untuk penelitian berdasarkan informasi yang ditampilkan pada situs *World Weather Online*. *Dataset* memiliki beberapa atribut diantaranya yaitu *Date*, *Time*, *Desc*, *Weather*, *Temp (Celcius)*, *Rain (mm)*, *Wind (mph)*, *Dir*, *Cloud (%)*,

Humidity (%), dan *Pressure (mdb)*. Atribut yang menjadi *class* atau label pada *dataset* yaitu atribut *weather*.

3.2. Pemodelan Proses

Bentuk pemodelan proses yang digunakan pada penelitian menggunakan software Rapid Miner Studio dimana model proses dapat dilihat pada gambar 2.

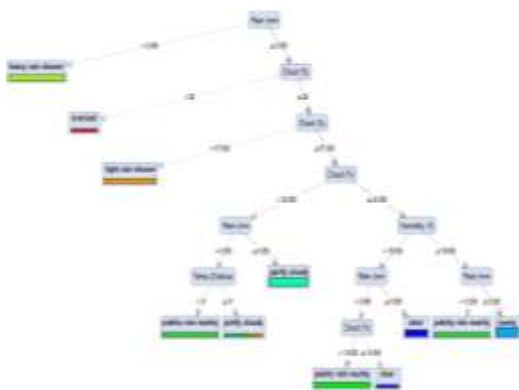


Gambar 2. Pemodelan Proses

Berdasarkan gambar 2 dapat dijelaskan dengan menggunakan *software Rapid Miner*, *dataset* yang digunakan diuji dengan menggunakan operator *X-Validation* dimana tipe sampling yaitu *stratified* dengan jumlah validasi yang dilakukan terhadap *dataset* sebanyak 10 kali.

3.3. Pola Pohon Keputusan

Setelah melewati proses pengujian, maka dihasilkan sebuah pola yang berbentuk pohon keputusan (*decision tree*) dikarenakan algoritma C4.5 termasuk ke dalam algoritma *decision tree*. Pohon keputusan yang dihasilkan dapat dilihat pada gambar 3.



Gambar 3. Pohon Keputusan

Berdasarkan gambar 3 dapat diketahui bahwa atribut-atribut apa saja yang dapat memberikan pengaruh dalam melakukan proses pengklasifikasian *dataset* sehingga menyebabkan terbentuknya pola pengetahuan yang dapat digunakan sebagai acuan dalam memprediksikan cuaca yang akan datang. Atribut tersebut adalah *Rain*, *Cloud*, *Humidity*, dan *Temp*.

3.4. Pengujian dengan Aplikasi Web

Setelah mendapatkan pola pengetahuan terhadap proses klasifikasi cuaca, penelitian dibuktikan dalam bentuk aplikasi berbasis web menggunakan bahasa pemrograman PHP. Aplikasi web tersebut terdiri dari 2 konten halaman yaitu halaman *single record test* dan halaman *multi record test*. Halaman *single record test* digunakan untuk membandingkan klasifikasi cuaca *user* dengan klasifikasi cuaca hasil dari algoritma C4.5 dengan hanya menggunakan satu *record* data yang dapat dilihat pada gambar 4.

[illegible]

Gambar 4. *Single Record Test*

Berdasarkan gambar 4, terdapat 4 input data yang dimasukkan oleh *user*. Sebagai contoh, *user* menginput nilai atribut *rain* sebesar 3.0, dan atribut lainnya dengan nilai *null*. Sesuai dengan pola pengetahuan yang didapatkan maka algoritma menghasilkan cuaca yaitu *heavy rain shower*.

Berikutnya yaitu pengujian untuk *multi record test* yang dapat dilihat pada gambar 5.

Machine Classification using DCS Algorithm

Search for the question

Question No. | Age (years) | Personality (No.) | Salary (K) | Education | Religion (1 & 0) | Compensation

1	2	30	25	none	none	150.0
2	3	35	20	none	none	150.0
3	3	30	25	none	none	150.0
4	3	35	20	none	none	150.0
5	3	30	25	none	none	150.0
6	3	40	30	none	none	150.0
7	3	37	32	partly cloudy	partly cloudy	150.0
8	3	34	30	partly sun mostly	partly sun mostly	150.0
9	3	33	27	partly sun mostly	partly cloudy	150.0
10	3	30	27	cloud	cloud	150.0

Showing 10 entries (1 - 10)

Gambar 5. *Multi Record Test*

Berdasarkan gambar 5, *user* dapat memasukkan sebuah file bertipe CSV yang didalamnya terdiri dari banyak *record* data, dimana data tersebut nantinya akan diuji dengan algoritma yang ada dan didapatkan nilai akurasinya.

4. Simpulan

Penelitian menggunakan *dataset* yang dibentuk dari informasi yang dihasilkan pada situs peramalan cuaca yaitu *World Weather Online* terhitung sejak tanggal 12 Agustus 2016 pukul 01:00 sampai dengan 20 Agustus 2016 pukul 22:00. Akurasi dari algoritma klasifikasi C4.5 menghasilkan nilai sebesar 88.89% yang telah dibuktikan melalui program yang dibuat.

Pengembangan pekerjaan yang akan datang dapat mempertimbangkan tidak hanya nilai *accuracy* dan *kappa* dari algoritma tersebut, tetapi memperhatikan nilai AUC yang dihasilkan. Untuk meningkatkan *accuracy* maka dapat digunakan metode optimasi salah satunya dengan menggunakan metode PSO (*Particle Swarn Optimization*) agar hasil yang didapat lebih akurat.

Referensi

- A. Joshi, B. K. (2015). Weather Forecasting and Climate Changing Using Data Mining Application Rain Effects on Speed. *Int. J. Adv. Res. Comput. Commun. Eng.* , 19–21.
- Goronescu, F. (2011). *Data Mining: Concepts, Models and Techniques*. Verlag Berlin Heidel: Springer.
- J. Han, M. K. (2012). *Data Mining Concepts and Techniques 3rd Edition*. USA: Morgan Kauffman.
- Kusrini, E. T. (2009). *Algoritma Data Mining*. Yogyakarta: Penerbit Andi.
- Larose, D. T. (2006). *Data Mining Methods and Models*. New Jersey: John Wiley & Sons, inc.
- Larose, D. T. (2005). *Discovering knowledge in data*. New Jersey: John Wiley & Sons, inc.
- Mohod, A. A. (2015). Applications of Data Mining in Weather Forecasting Using Frequent Pattern Growth Algorithm. *Int. J. Sci. Res.* , 3048–3051.
- Mujiasih, S. (2011). Utilization of Data Mining for Weather Forecastin. *Journal of Meteorol dan Geofis* , 189-195.
- Olaiya, F. (2012). Application of Data Mining Techniques in Weather Prediction and Climate Change Studies. *I.J. Inf. Eng. Electron. Bus.* , 51-59.
- R. S. Wahono, N. S. (2014). A comparison framework of classification models for software defect prediction. *Adv. Sci. Lett.* , 1945–1950.
- Winarso, P. A. (2003). *Pengelolaan Bencana Cuaca dan Iklim untuk masa mendatang*. Indonesia: KLIH.
- Xindong, W. K. (2009). *The Top Ten Algorithms in Data Mining*. USA: Taylor & Francis Group, LLC.