

BAB II

LANDASAN TEORI

2.1 Data Mining

Secara sederhana, *data mining* merupakan ekstraksi informasi yang tersirat dalam sekumpulan data. Data mining merupakan sebuah proses untuk menggali kumpulan data dan menemukan informasi di dalamnya [12]. Data mining merupakan proses pengekstrakan informasi dari jumlah kumpulan data yang besar dengan menggunakan algoritma dan teknik gambar dari statistik, mesin pembelajaran dan sistem manajemen *database*. Penggalian data ini dilakukan pada sekumpulan data yang besar untuk menemukan pola atau hubungan yang ada dalam kumpulan data tersebut [6]. Hasil penemuan yang diperoleh setelah proses penggalian data ini, kemudian dapat digunakan untuk analisis yang lebih lanjut.

Data mining yang disebut juga dengan *Knowledge-Discovery in Database* (KDD) adalah sebuah proses secara otomatis atas pencarian data di dalam sebuah memori yang amat besar dari data untuk mengetahui pola dengan menggunakan alat seperti klasifikasi, hubungan (*association*) atau pengelompokan (*clustering*). Proses KDD ini terdiri dari langkah-langkah sebagai berikut [4]:

1. *Data Cleaning*, proses menghapus data yang tidak konsisten dan kotor.
2. *Data Integration*, penggabungan beberapa sumber data.
3. *Data Selection*, pengambilan data yang akan dipakai dari sumber data.
4. *Data Transformation*, proses dimana data ditransformasikan menjadi bentuk yang sesuai untuk diproses dalam data mining.
5. *Data Mining*, suatu proses yang penting dengan melibatkan metode untuk menghasilkan suatu pola data.
6. *Pattern Evaluation*, proses untuk menguji kebenaran dari pola data yang mewakili *knowledge* yang ada didalam data itu sendiri.

7. *Knowledge Presentation*, proses visualisasi dan teknik menyajikan *knowledge* digunakan untuk menampilkan *knowledge* hasil *mining* kepada *user*.

2.2 Metode Data Mining

Pada umumnya metode *data mining* dapat dikelompokkan kedalam dua kategori yaitu *deskriptif* dan *prediktif*. Metode *deskriptif* bertujuan untuk mencari pola yang dapat dimengeti oleh manusia yang menjelaskan karakteristik dari data. Metode *prediktif* menggunakan ciri-ciri tertentu dari data. Pada umumnya metode *data mining* dapat dikelompokkan kedalam dua untuk melakukan prediksi.

Metode-metode yang ada dalam *data mining* adalah sebagai berikut [11]:

1. *Classification*

Klasifikasi (*Classification*) merupakan proses untuk menemukan sekumpulan model yang menjelaskan dan membedakan kelas-kelas data, sehingga model tersebut dapat digunakan untuk memprediksi nilai suatu kelas yang belum diketahui pada sebuah objek. Untuk mendapatkan model, kita harus melakukan analisis terhadap data latih (*training set*). Sedangkan data uji (*test set*) digunakan untuk mengetahui tingkat akurasi dari model yang dihasilkan. Klasifikasi dapat digunakan untuk memprediksi nama atau nilai kelas dari suatu obyek data. Metode inilah yang digunakan dalam tugas akhir ini.

2. *Clustering*

Pengelompokan (*Clustering*) merupakan proses untuk melakukan segmentasi. Digunakan untuk melakukan pengelompokan secara alami terhadap atribut suatu set data, termasuk kedalam *supervised task*. Contoh *clustering* seperti mengelompokkan dokumen berdasarkan topiknya.

3. *Assosiation*

Tujuan dari metode ini untuk menghasilkan sejumlah *rule* yang menjelaskan sejumlah data yang berhubung kuat satu dengan yang lainnya. Sebagai contoh *assosiation analysis* dapat digunakan untuk menentukan produk yang datang

secara bersamaan oleh banyak pelanggan, atau bisa juga disebut dengan *basket analysis*.

4. *Regression*

Regression mirip dengan klasifikasi. Perbedaan utamanya adalah terletak pada atribut yang diprediksi berupa nilai yang kontinyu.

5. *Forecasting*

Prediksi (*Forecasting*) berfungsi untuk melakukan kejadian yang akan datang berdasarkan data sejarah yang ada.

6. *Sequence Analysis*

Tujuan dari metode ini adalah untuk mengenali pola dari data diskrit. Sebagai contoh adalah menemukan kelompok gen dengan tingkat ekspresi yang mirip.

7. *Deviation Analysis*

Tujuan dari metode ini adalah untuk menemukan penyebab perbedaan antara data yang satu dengan data yang lain dan biasa disebut sebagai *oulier detection*. Sebagai contoh adalah apakah sudah terjadi penipuan terhadap pengguna kredit dengan melihat catatan transaksi yang tersimpan dalam basis data perusahaan kartu kredit.

2.3 Klasifikasi

Klasifikasi merupakan suatu pekerjaan menilai objek data untuk memasukkannya ke dalam kelas tertentu dari sejumlah kelas yang tersedia. Dalam klasifikasi ada dua pekerjaan utama yang dilakukan, yaitu (1) pembangunan model sebagai prototipe untuk disimpan sebagai memori dan (2) penggunaan model tersebut untuk melakukan pengenalan/klasifikasi/prediksi pada suatu objek data lain agar diketahui dikelas mana objek data tersebut dalam model yang sudah disimpannya. Model dalam klasifikasi mempunyai arti yang sama dengan kotak hitam, dimana ada suatu model yang menerima masukan, kemudian mampu melakukan pemikiran

terhadap masukan tersebut dan memberikan jawaban sebagai keluaran dari hasil pemikiannya [7].

Tahapan dari klasifikasi dalam data mining terdiri dari [4] :

1. Pembangunan Model

Pada tahapan ini dibuat sebuah model untuk menyelesaikan masalah klasifikasi class atau atribut dalam data. Tahap ini merupakan fase pelatihan, dimana data latih dianalisis menggunakan algoritma klasifikasi, sehingga model pembelajaran direpresentasikan dalam bentuk aturan klasifikasi.

2. Penerapan Model

Pada tahapan ini model yang sudah dibangun sebelumnya digunakan untuk menentukan atribut/kelas dari sebuah data baru yang atribut/kelasnya belum diketahui sebelumnya. Tahap ini digunakan untuk memperkirakan keakuratan aturan klasifikasi terhadap data uji. Jika model dapat diterima, maka aturan dapat diterapkan terhadap klasifikasi data baru.

2.4 Decision Tree

2.4.1 Pengertian Decision Tree

Decision tree merupakan metode klasifikasi *data mining*. *Decision tree* dalam istilah pembelajaran merupakan sebuah struktur pohon dimana setiap *node* pohon mempresentasikan atribut yang telah diuji. Setiap cabang merupakan suatu pembagian hasil uji dan *node* daun (*leaf*) mempresentasikan kelompok kelas tertentu. [5]. Level *node* teratas dari sebuah *Decision Tree* adalah *node* akar (*root*) yang biasanya berupa atribut yang paling memiliki pengaruh terbesar pada suatu kelas tertentu. Pada umumnya *Decision Tree* melakukan strategi pencarian secara *top-down* untuk solusinya. Pada proses mengklasifikasi data yang tidak diketahui, nilai atribut akan diuji dengan cara melacak jalur dari *node* akar (*root*) sampai *node* akhir (daun) dan kemudian akan diprediksi kelas yang dimiliki oleh suatu data baru tertentu.

2.4.2 Jenis - Jenis Decision Tree

Beberapa model *decision tree* yang sudah dikembangkan antara lain C4.5 atau ID3 dan CART. Berikut ini akan dijelaskan model dari decision tree tersebut :

1. C4.5 atau ID3

Decision Tree menggunakan algoritma ID3 atau C4.5, yang diperkenalkan dan dikembangkan pertama kali oleh Quinlan yang merupakan singkatan dari *Iterative Dichotomiser 3* atau *Induction of Decision 3*. Algoritma ID3 membentuk pohon keputusan dengan metode *divide-and-conquer* data secara rekursif dari atas ke bawah. Strategi pembentukan Decision Tree dengan algoritma ID3 adalah:

- A. Pohon dimulai sebagai *node* tunggal (akar/*root*) yang merepresentasikan semua data.
- B. Sesudah *node root* dibentuk, maka data pada *node* akar akan diukur dengan *information gain* untuk dipilih atribut mana yang akan dijadikan atribut pembagiannya.
- C. Sebuah cabang dibentuk dari atribut yang dipilih menjadi pembagi dan data akan didistribusikan ke dalam cabang masing-masing.
- D. Algoritma ini akan terus menggunakan proses yang sama atau bersifat rekursif untuk dapat membentuk sebuah *Decision Tree*. Ketika sebuah atribut telah dipilih menjadi *node* pembagi atau cabang, maka atribut tersebut tidak diikuti lagi dalam penghitungan nilai *information gain*.
- E. Proses pembagian rekursif akan berhenti jika salah satu dari kondisi dibawah ini terpenuhi :
 - a. Semua data dari anak cabang telah termasuk dalam kelas yang sama.
 - b. Semua atribut telah dipakai, tetapi masih tersisa data dalam kelas yang berbeda. Dalam kasus ini, diambil data yang mewakili kelas yang terbanyak untuk menjadi label kelas pada *node* daun. Tidak terdapat data pada anak cabang yang baru. Dalam kasus ini, *node* daun akan

dipilih pada cabang sebelumnya dan diambil data yang mewakili kelas terbanyak untuk dijadikan label kelas.

Metode C4.5 dan ID3 memiliki perbedaan dalam nilai tiap atribut. Metode C4.5 menggunakan atribut yang bernilai kategorikal dan numerikal, sedangkan metode ID3 menggunakan atribut yang bernilai kategorikal. Metode *decision tree C4.5* inilah yang digunakan dalam tugas akhir ini.

2. CART

CART adalah singkatan dari *Classification And Regression Tree*. Dalam CART ada dua langkah penting yang harus diikuti untuk mendapatkan *tree* dengan performansi yang optimal. Yang pertama adalah pemecahan objek secara berulang berdasarkan atribut tertentu. Yang kedua, *prunning* (pemangkasan) dengan menggunakan data validasi.

Misalkan kita mempunyai variabel independent $x_1, x_2, x_3, \dots, x_n$ dan variabel dependent atau output y . Pemecahan secara berulang berarti kita bagi objek ke dalam kotak-kotak berdasarkan nilai variabel x_1, x_2 atau x_r . Cara ini diulang sehingga dalam suatu kotak sebisa mungkin berisi observasi dalam kelompok atau kelas yang sama.

Langkah berikutnya sesudah dilakukan pemecahan objek atau data secara berulang adalah melakukan *prunning*. Dalam *prunning* kita ingin memangkas *tree* yang mungkin terlalu besar dan terjadi fenomena *overfitting*. *Overfitting* merupakan sebuah satu buah pengelompokan yang mungkin hanya berisi satu data yang memungkinkan data tersebut merupakan *noise* yang ada di data training dan bukan pola yang mungkin terjadi dalam data testing atau data validasi. *Prunning* terdiri dari beberapa langkah pemilihan secara berulang simpul yang akan dijadikan simpul daun. Dengan mengubah simpul menjadi simpul daun artinya tidak akan dilakukan pemecahan lagi sesudah itu. Dengan demikian ukuran *tree* akan berkurang. [12]

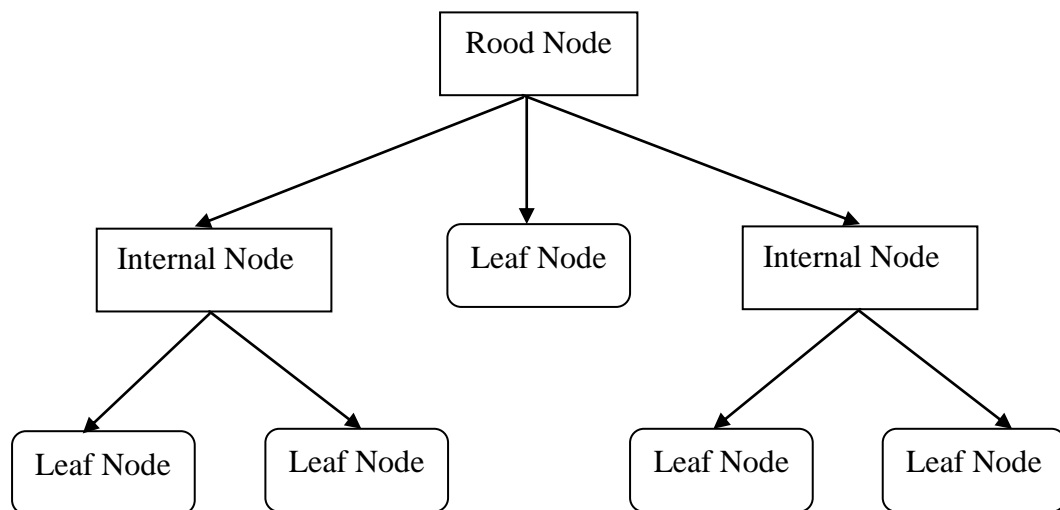
2.4.3 Model Decision Tree

Decision tree adalah *flow-chart* seperti *struktur tree*, dimana tiap *internal node* menunjukkan sebuah test pada sebuah atribut, tiap cabang menunjukkan hasil dari test, dan *leaf node* menunjukkan *class-class* atau *class distribution*.

Selain karena pembangunannya relatif cepat, hasil dari model yang dibangun mudah untuk dipahami. Pada *decision tree* terdapat 3 jenis *node*, yaitu:

- a. *Root Node*, merupakan *node* paling atas, pada *node* ini tidak ada *input* dan bisa tidak mempunyai *output* atau mempunyai *output* lebih dari satu.
- b. *Internal Node*, merupakan *node* percabangan, pada *node* ini hanya terdapat satu *input* dan mempunyai *output* minimal dua.
- c. *Leaf node* atau *terminal node*, merupakan *node* akhir, pada *node* ini hanya terdapat satu *input* dan tidak mempunyai *output*.

Contoh dari model pohon keputusan yaitu seperti pada **gambar 2.1** berikut:



Gambar 2.1 Model *Decision Tree*

2.5 Gizi

2.5.1 Antropometri Gizi

Status gizi adalah ukuran keberhasilan dalam pemenuhan nutrisi untuk anak yang diindikasikan oleh berat badan dan tinggi badan anak. Status gizi juga didefinisikan sebagai status kesehatan yang dihasilkan oleh keseimbangan antara kebutuhan dan masukan nutrien. Penelitian status gizi merupakan pengukuran yang didasarkan pada data antropometri serta biokimia dan riwayat. [10]

Antropometri adalah ilmu yang mempelajari berbagai ukuran tubuh manusia. Dalam bidang ilmu gizi digunakan untuk menilai status gizi. Ukuran yang sering digunakan adalah berat badan dan tinggi badan. Selain itu juga ukuran tubuh lainya seperti lingkar kepala, lingkar lengan atas, lingkar perut, lingkar pinggul. [9]

Adapun beberapa syarat yang mendasari penggunaan antropometri ini adalah [10] :

1. Alatnya mudah didapat dan digunakan, seperti dacin, pita, mikrotoa, dan alat pengukur panjang bayi yang dapat dibuat sendiri.
2. Pengukuran dapat dilakukan berulang- ulang dengan mudah dan objektif.
3. Pengukuran bukan hanya dilakukan dengan tenaga khusus atau profesional, juga oleh tenaga lain yang setelah dilatih untuk itu.
4. Biaya relatif murah, karena alat mudah didapat dan tidak memerlukan bahan- bahan lainnya.
5. Hasilnya mudah disimpulkan karena mempunyai ambang batas (cut off points) dan baku rujukan yang sudah pasti.
6. Secara ilmiah diakui kebenarannya. Hampir semua negara menggunakan antropometri sebagai metode untuk mengukur status gizi, khususnya untuk penampisan (*screening*) status gizi.

2.5.2 Macam-Macam Status Gizi

Status gizi terbagi menjadi dua macam, yaitu status gizi normal dan malnutrisi : [10]

1. Status Gizi Normal

Keadaan tubuh yang mencerminkan keseimbangan antara konsumsi dan penggunaan gizi oleh tubuh (*adequate*).

2. Malnutrisi

Keadaan patologis akibat kekurangan atau kelebihan secara relatif maupun absolut satu atau lebih zat gizi. Ada empat bentuk :

- a) *Under nutrition*: kekurangan konsumsi pangan secara relatif atau absolut untuk periode tertentu.
- b) *Specific deficiency*: kekurangan zat gizi tertentu, misalnya kekurangan *iodium dan Fe* (zat besi).
- c) *Over nutrition*: kelebihan konsumsi pangan untuk periode tertentu.
- d) *Imbalance*: keadaan disproporsi zat gizi, misalnya tinggi kolesterol karena tidak imbangnya kadar LDL, HDL dan VDL.

2.5.3 Jenis Parameter Gizi

Ada beberapa jenis parameter yang dilakukan untuk mengukur tubuh manusia yaitu: usia, jenis kelamin, berat badan, tinggi badan, dan lingkaran kepala.

2.5.4 Penilaian Status Gizi

Macam-macam penilaian status gizi [10].

I. Penilaian Status Gizi Secara Langsung

A. Antropometri

1. Pengertian

Secara umum antropometri artinya ukuran tubuh manusia. Ditinjau dari sudut pandang gizi, maka antropometri gizi berhubungan dengan berbagai macam pengukuran dimensi tubuh dan komposisi tubuh dari berbagai tingkat umur dan tingkat gizi.

2. Penggunaan

Antropometri secara umum digunakan untuk melihat ketidakseimbangan asupan protein dan energi. Ketidakseimbangan ini terlihat pada pola pertumbuhan fisik dan proporsi jaringan tubuh seperti lemak, otot dan jumlah air dalam tubuh.

II. Pengukuran status gizi dengan menggunakan KMS (Kartu Menuju Sehat)

1. Pengertian

KMS (Kartu Menuju Sehat) untuk balita adalah alat yang sederhana dan murah yang dapat digunakan untuk memantau kesehatan dan pertumbuhan anak. KMS berisi catatan penting tentang pertumbuhan, perkembangan anak, imunisasi, penanggulangan diare, pemberian kapsul vit A, kondisi kesehatan anak, pemberian asi eksklusif dan makanan pendamping ASI, pemberian makanan anak dan rujukan ke puskesmas. [2]

2. Manfaat KMS (Kartu Menuju Sehat)

- a) Sebagai media untuk mencatat dan memantau riwayat kesehatan balita secara lengkap meliputi: pertumbuhan, perkembangan, pelaksanaan, imunisasi, penanggulangan diare, pemberian kapsul vit A, kondisi kesehatan pemberian ASI eksklusif, dan makanan pendamping ASI.
- b) Sebagai media edukasi bagi orang tua balita tentang kesehatan anak.

- c) Sebagai sarana komunikasi yang dapat digunakan oleh petugas untuk menentukan penyuluhan dan tindakan pelayanan kesehatan gizi [2]

2.6 Algoritma Decision Tree C4.5

Algoritma C4.5 diperkenalkan oleh Quinlan pada tahun 1996 sebagai versi perbaikan dari ID3. Dalam ID3, induksi decision tree hanya bisa dilakukan pada fitur bertipe kategorikal (nominal atau ordinal), sedangkan tipe numerik (interval atau rasio) tidak dapat digunakan.

Yang menjadi hal penting dalam induksi decision tree adalah bagaimana menyatakan syarat pengujian pada node. Ada 3 kelompok penting dalam syarat pengujian node :

1. Fitur biner

Adalah Fitur yang hanya mempunyai dua nilai berbeda. Syarat pengujian ketika fitur ini menjadi node (akar maupun interval) hanya punya dua pilihan cabang.

2. Fitur kategorikal

Untuk fitur yang nilainya bertipe kategorikal (nominal atau ordinal) bisa mempunyai beberapa nilai berbeda. Secara umum ada 2 pemecahan yaitu pemecahan biner (*binary splitting*) dan (*multi splitting*).

3. Fitur numerik

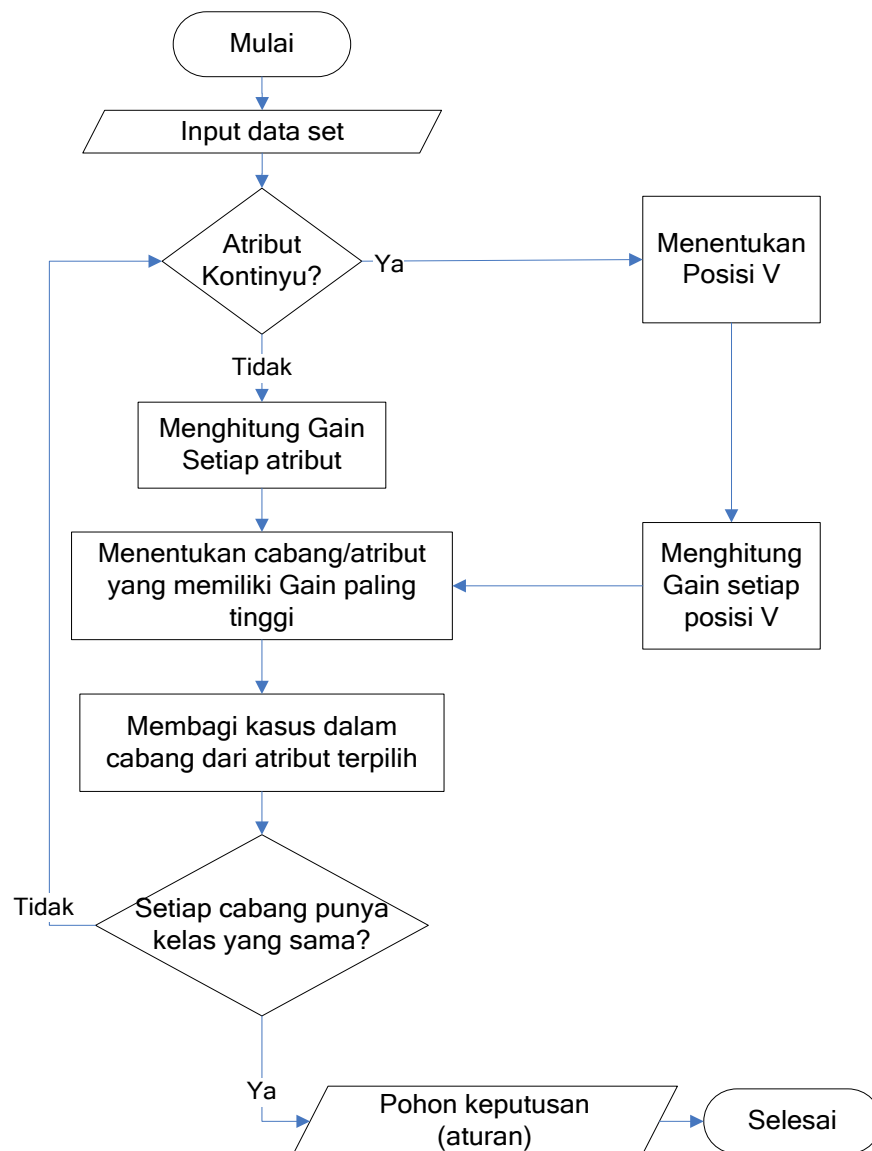
Untuk fitur bertipe numerik, Syarat pengujian dalam node (akar maupun internal) dinyatakan dengan pengujian perbandingan ($A \leq V$) atau ($A > V$) dengan hasil biner.

Secara umum algoritma C4.5 untuk membangun pohon keputusan adalah sebagai berikut:

1. Pilih atribut sebagai akar.
2. Buat cabang untuk tiap-tiap nilai.
3. Bagi kasus dalam cabang.

4. Ulangi proses untuk setiap cabang sampai semua kasus pada cabang memiliki kelas yang sama.

Berikut ini akan dijelaskan secara lebih detail algoritma C4.5 menggunakan *flowchart* yang disajikan pada **gambar 2.2**.



Gambar 2.2 Flowchart algoritma Decision Tree C4.5

Untuk memilih atribut sebagai simpul akar (*root node*) atau simpul dalam (*internal node*), didasarkan pada nilai *information gain* tertinggi dari atribut-atribut yang ada. Sebelum perhitungan *information gain*, akan dilakukan perhitungan *entropy*. *Entropy* merupakan distribusi probabilitas dalam teori informasi dan diadopsi kedalam algoritma C4.5 untuk mengukur tingkat homogenitas distribusi kelas dari sebuah himpunan data (*data set*). Semakin tinggi tingkat *entropy* dari sebuah data maka semakin homogen distribusi kelas pada data tersebut. Perhitungan *information gain* menggunakan rumus 2.1, sedangkan *entropy* menggunakan rumus 2.2.

$$Gain(S, A) = Entropy(S) - \sum_{i=1}^n \frac{|S_i|}{|S|} * Entropy(S_i) \dots \dots \dots (2.1)$$

dimana,

S : Himpunan kasus

A : Atribut

n : Jumlah partisi atribut A

|S_i|: Jumlah kasus pada partisi ke i

|S| : Jumlah kasus dalam S

$$Entropy(S) = - \sum_{i=1}^n p_i * \log_2 p_i \dots \dots \dots (2.2)$$

dimana,

S : Himpunan kasus

n : Jumlah partisi S

p_i : Proporsi dari S_i terhadap S

Selain *Information Gain* kriteria yang lain untuk memilih atribut sebagai pemecah adalah *Rasio Gain*. Perhitungan rasio gain menggunakan rumus 2.3, sedangkan split information menggunakan rumus 2.4.

$$GainRasio(S, A) = \frac{Gain(S, A)}{SplitInformation(S, A)} \dots \dots \dots (2.3)$$

$$SplitInformation(S,A) = -\sum_{i=1}^c \frac{S_i}{S} \log_2 \frac{S_i}{S} \dots\dots\dots(2.4)$$

dimana S_1 sampai S_c adalah c subset yang dihasilkan dari pemecahan S dengan menggunakan atribut A yang mempunyai sebanyak c nilai.

2.7 Contoh Perhitungan

Berikut ini akan dijelaskan ilustrasi dari alur proses perhitungan algoritma *Decision Tree C4.5*. Data set yang digunakan pada contoh ini adalah data untuk menentukan *Play* atau *Don't Play* dengan beberapa atribut yaitu atribut *outlook*, *temperature*, *humidity* dan *windy*. Dimana atribut *temperature* dan *humidity* bertipe kontinyu sedangkan *outlook* dan *windy* bertipe kategorikal, sedangkan kolom *Class* adalah kelas tujuannya atau label kelas-nya.

Tabel 2.1 Contoh data set

Outlook	Teperature	Humidity	Windy	Class
Sunny	75	70	TRUE	Play
Sunny	80	90	TRUE	Don't play
Sunny	85	85	FALSE	Don't play
Sunny	72	95	FALSE	Don't play
Sunny	69	70	FALSE	Play
overcast	72	90	TRUE	Play
overcast	83	78	FALSE	Play
overcast	64	65	TRUE	Play
overcast	81	75	FALSE	Play
Rain	71	80	TRUE	Don't play
Rain	65	70	TRUE	Don't play
Rain	75	80	FALSE	Play
Rain	68	80	FALSE	Play
Rain	70	96	FALSE	Play

Pada tabel 2.1 ini rumus yang digunakan untuk memilih atribut sebagai *node* adalah rumus *information gain*. Proses pertama adalah menghitung *entropy* untuk semua data.

Jumlah class play = 9

Jumlah class don't play = 5

Berikut adalah perhitungan *entropy* untuk semua data:

$$\begin{aligned} Entropy(S) &= -\frac{9}{14} * \log_2\left(\frac{9}{14}\right) - \frac{5}{14} * \log_2\left(\frac{5}{14}\right) \\ &= 0.940 \end{aligned}$$

Selanjutnya menghitung *gain* untuk setiap atribut. Berikut adalah contoh perhitungan *gain* untuk atribut *outlook*:

Tabel 2.2 Distribusi jumlah atribut *outlook*

Nilai Outlook	Σ Play	Σ Don't Play	Total
Sunny	2	3	5
Overcast	4	0	4
Rain	3	2	5

Berdasarkan tabel 2.2, maka nilai *information gain* untuk atribut *outlook* adalah sebagai berikut:

$$\begin{aligned} Gain(outlook) &= 0.940 - \left(\frac{5}{14} * \left(-\frac{2}{5} * \log_2\left(\frac{2}{5}\right) - \frac{3}{5} * \log_2\left(\frac{3}{5}\right) \right) \right. \\ &\quad + \frac{4}{14} * \left(-\frac{4}{4} * \log_2\left(\frac{4}{4}\right) - \frac{0}{4} * \log_2\left(\frac{0}{4}\right) \right) \\ &\quad \left. + \frac{5}{14} * \left(-\frac{3}{5} * \log_2\left(\frac{3}{5}\right) - \frac{2}{5} * \log_2\left(\frac{2}{5}\right) \right) \right) \\ &= 0.940 - 0.694 \\ &= 0.246 \end{aligned}$$

Untuk perhitungan atribut yang bertipe kontinyu, harus menentukan *posisi V* terbaik yang dinyatakan dalam perbandingan ($A \leq V$) atau ($A > V$). Berikut akan dijelaskan contoh perhitungan dari atribut *temperature*.

Misal posisi V yang akan digunakan pada atribut *temperature* adalah 65,70,75,dan 80, kemudian dihitung nilai *information gain*-nya.

Contoh perhitungan *temperature* posisi v = 65:

$$\begin{aligned} \text{Gain}(\text{temp}) &= 0.940 - \left(\frac{2}{14} * \left(-\frac{1}{2} * \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} * \log_2 \left(\frac{1}{2} \right) \right) \right. \\ &\quad \left. + \frac{12}{14} * \left(-\frac{8}{12} * \log_2 \left(\frac{8}{12} \right) - \frac{4}{12} * \log_2 \left(\frac{4}{12} \right) \right) \right) \\ &= 0.940 - 0.930 \\ &= 0.010 \end{aligned}$$

Berikut hasil perhitungan atribut numerik untuk setiap posisi yang telah ditentukan:

Tabel 2.3 Hasil perhitungan posisi V untuk atribut *temperature*

Temperature	65		70		75		80	
	≤	>	≤	>	≤	>	≤	>
Play	1	8	4	5	7	2	7	2
Don't Play	1	4	1	4	3	2	4	1
Jumlah	2	12	5	9	10	4	11	3
Entropy	1	0.918	0.722	0.991	0.881	1	0.946	0.918
Gain	0.01		0.045		0.025		0.0005	

Berdasarkan tabel 2.3, nilai gain tertinggi adalah 70, maka nilai *information gain* pada atribut *temperature* adalah 0.045. Hasil perhitungan pada setiap atribut disajikan pada tabel 2.4

Tabel 2.4 Hasil perhitungan *Information gain* untuk setiap atribut

		Jumlah	Play	Don't Play	Entropy	Gain
Total		14	9	5	0.94	
Outlook	Sunny	5	2	3	0.971	0.246
	Overcast	4	4	0	0	
	Rain	5	3	2	0.971	
Temperature	≤ 70	5	4	1	0.722	0.045

	> 70	9	5	4	0.991	
Humidity	≤ 80	9	7	2	0.764	0.102
	≤ 80	5	2	3	0.971	
Windy	TRUE	6	3	3	1	0.048
	FALSE	8	6	2	0.811	

Berdasarkan tabel 2.4 menunjukkan bahwa atribut *outlook* memiliki nilai gain tertinggi, maka atribut *outlook* akan menjadi *node*. Karena atribut *outlook* memiliki tiga nilai atribut atau lebih dari dua, maka dilakukan perhitungan rasio gain untuk memilih pilihan percabangan terbaik. Berikut adalah contoh perhitungan rasio gain untuk pilihan percabangan {sunny, overcast, rain}.

$$\begin{aligned}
 \text{Split info}(\text{Semua}, \text{overcast}) &= \left(-\frac{5}{14} * \log_2 \left(\frac{5}{14} \right) \right) + \left(-\frac{4}{14} * \log_2 \left(\frac{4}{14} \right) \right) \\
 &\quad + \left(-\frac{5}{14} * \log_2 \left(\frac{5}{14} \right) \right) \\
 &= 0.531 + 0.516 + 0.531 = 1.577
 \end{aligned}$$

$$\begin{aligned}
 \text{Rasio Gain}(\text{Semua}, \text{overcast}) &= \frac{0.246}{1.577} \\
 &= 0.156
 \end{aligned}$$

Hasil untuk perhitungan *rasio gain* lainnya ada pada tabel 2.5.

Tabel 2.5 Hasil perhitungan *Rasio gain* untuk setiap pilihan cabang

			Jumlah	Split Inf	Gain	Rasio Gain
Total			14		0.246	
Pilihan 1	sunny		5	1.577		0.156
	overcast		4			
	rain		5			
Pilihan 2	sunny		5	0.94		0.262
	overcast	Rain	9			
Pilihan 3	sunny	Overcast	9	0.94		0.262

	rain		5			
Pilihan 4	sunny	Rain	10	0.863		0.286
	overcast		4			

Dari tabel 2.5 pilihan 4 yaitu {*sunny*, *rain*} dan {*overcast*} memiliki nilai *rasio gain* tertinggi, maka atribut terpilih (*outlook*) akan dibagi menjadi dua cabang. Pembagian cabang disajikan pada tabel 2.6 dan tabel 2.7.

Tabel 2.6 Pembagian cabang (*sunny*, *rain*)

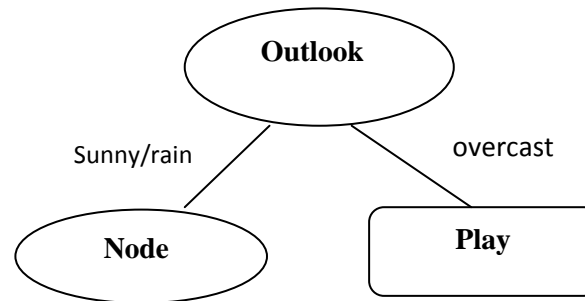
Outlook	Temperature	Humidity	Windy	Class
Sunny	75	70	TRUE	Play
Sunny	80	90	TRUE	Don't play
Sunny	85	85	FALSE	Don't play
Sunny	72	95	FALSE	Don't play
Sunny	69	70	FALSE	Play
Rain	71	80	TRUE	Don't play
Rain	65	70	TRUE	Don't play
Rain	75	80	FALSE	Play
Rain	68	80	FALSE	Play
Rain	70	96	FALSE	Play

Tabel 2.7 Pembagian cabang (*overcast*)

Outlook	Temperature	Humidity	Windy	Class
Overcast	72	90	TRUE	Play
Overcast	83	78	FALSE	Play
Overcast	64	65	TRUE	Play
Overcast	81	75	FALSE	Play

Pada cabang *overcast* memiliki kelas yang sama yaitu *Play*, maka *node* ini akan menjadi daun dengan nilai *Play*. Sedangkan cabang *sunny* dan *rain* masih ada kelas yang berbeda, maka akan memilih atribut sebagai *node* seperti ditunjukkan pada

gambar 2.3. Proses tersebut akan berulang sampai semua kasus pada cabang memiliki kelas yang sama atau menjadi daun(*leaf*).



Gambar 2.3 Hasil pembentukan cabang pada node akar

Berikut aturan IF THEN untuk *decision tree* pada penjelasan pada gambar 2.3.

IF Outlook = overcast THEN Class = Play.

2.8 Penelitian Sebelumnya

Penelitian sebelumnya yang menggunakan metode Naïve Bayes dilakukan oleh Faridatul Choiriyah (Universitas Muhammadiyah Gresik, 2015) dengan judul “Implementasi Metode Naïve Bayes Sebagai Penentu Status Gizi Balita (Study Kasus puskesmas Dukuh Kupang Surabaya)”. Algoritma yang digunakan adalah naïve bayes. Data yang dijadikan inputan dalam sistem klasifikasi status gizi balita diperoleh dari data balita yang terdapat dipegawai bidan KIA puskesmas dukuh kupang surabaya, bulan september tahun 2014 sebanyak 131 balita. Atribut yang terdapat pada tabel mewakili fitur data yang digunakan meliputi jenis kelamin, usia, berat badan, tinggi badan, dan lingkaran kepala. Jumlah data yang digunakan sebanyak 131 *record* dengan kelas Baik dan Kurang. Penelitian ini diuji sebanyak 2 kali pengujian dengan rata-rata akurasi sebesar 92%.

Penelitian selanjutnya adalah tentang metode Decision Tree C4.5 dalam penelitian berjudul “Aplikasi Klasifikasi Penentuan Penerimaan Beras Miskin (Raskin) ini Ds.Sidomulyo Kec.Deket Kab.Lamongan dengan Metode Decision Tree

C4.5”, dibuat oleh M.Basroni Rizal (Universitas Muhammadiyah Gresik, 2016). Penelitian dilakukan untuk mengklasifikasi penentuan penerimaan beras miskin. Dari 600 data kepala keluarga di Ds. Sidomulyo Kec. Deket Kab. Lamongan, data tersebut diambil 40 % yang akan dijadikan sebagai data uji dan 60 % akan menjadi data latih. Jadi jumlah pembagiannya adalah 240 data sebagai data latih dan 360 data untuk data uji. Penelitian ini diuji sebanyak 7 kali pengujian dengan rata-rata akurasi sebesar 92%.

Penelitian selanjutnya adalah tentang metode *Decision Tree C4.5* dalam penelitian yang berjudul “*System prediksi prestasi akademik mahasiswa menggunakan metode decision tree C4.5 (Studi kasus:Jurusan Teknik informatika UNMUH GRESIK)*”, dibuat oleh Aunur Rasyid (Universitas Muhammadiyah Gresik, 2014). Tujuan dari penelitian tersebut adalah untuk menghasilkan informasi perkiraan kategori prestasi mahasiswa menggunakan metode *Decision Tree C4.5* sebagai peringatan dini dan motivasi mahasiswa dalam mendapatkan prestasi yang maksimal. Atribut-atribut yang digunakan adalah instansi sekolah asal (SMK,SMA atau MA), satatus sekolah asal (Negri atau Swasta), jurusan sekolah asal (IPA,IPS,Bahasa,Teknik,Administrasi), nilai rata-rata UN, status kerja (Sudah atau Belum), dan pihak yang mempengaruhi mahasiswa dalam memilih kuliah (Sendiri,Orang tua atau Orang lain). Hasil dari penelitian tersebut, Sistim Prediksi mahasiswa yang dirancang menggunakan algoritma C4.5 dapat memprediksi prestasi mahasiswa agar mampu mempertahankan kondisinya atau melakukan perbaikan utuk mencapai prestasi yang maksimal. Hasil akurasi dari penelitian tersebut adalah 90%.