

# Konsep dan Teknik Data Mining



Klasifikasi dan Prediksi

1/50

## Klasifikasi dan Prediksi

Apa Itu Klasifikasi dan apa itu prediksi?

Isu yang berkenaan dengan klasifikasi dan prediksi

Klasifikasi dengan induksi pohon keputusan

Klasifikasi bayesian

Klasifikasi dan Prediksi

2/50

## Klasifikasi vs. Prediksi

- **Klasifikasi:**
  - Meramalkan kategori label kelas (diskrit atau nominal)
  - Mengklasifikasikan (membuat suatu model) berdasarkan himpunan pelatihan dan nilai-nilai (**label kelas**) dalam suatu atribut klasifikasi dan menggunakannya didalam mengklasifikasikan data baru
- **Prediksi:**
  - Memodelkan fungsi bernilai kontinu, artinya menaksir nilai-nilai yang tidak diketahui atau nilai-nilai yang hilang

Klasifikasi dan Prediksi

3/50

## Klasifikasi vs. Prediksi

- Bentuk umum aplikasi
  - Persetujuan kredit
  - Target marketing
  - Diagnosa medis
  - Analisis keefektifan tindakan

Klasifikasi dan Prediksi

4/50

## Klasifikasi—Suatu Proses 2 Step

1. **Konstruksi model**: menguraikan suatu himpunan kelas yang ditentukan sebelumnya
  - Setiap tuple/sample dimasukkan masuk kedalam suatu kelas yang didefinisikan sebelumnya, seperti yang ditetapkan melalui **label atribut kelas**
  - Himpunan dari tuple yang digunakan untuk konstruksi model dinamai **himpunan pelatihan**
  - Model disajikan sebagai kaidah klasifikasi, pohon keputusan, atau rumus matematika

Klasifikasi dan Prediksi

5/50

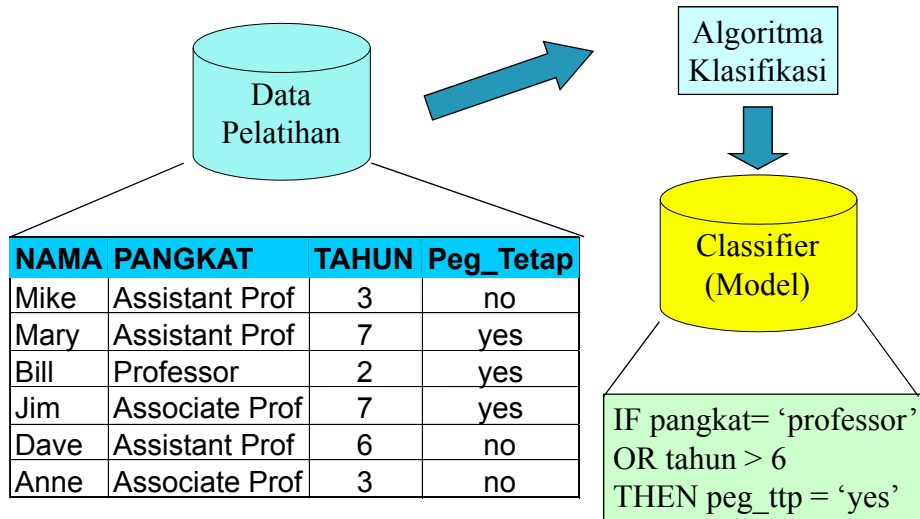
## Klasifikasi—Suatu Proses 2 Step

2. **Kegunaan model**: untuk klasifikasi kedepan atau untuk objek-objek yang tidak dikenal
  - Menaksir akurasi dari model
    - Label yang dikenali melalui uji sampel dibandingkan dengan hasil klasifikasi dari model
    - Tingkat akurasi adalah persentasi dari himpunan sampel uji yang secara benar diklasifikasikan oleh model tersebut
  - **Jika akurasi bisa diterima, gunakan model tersebut untuk mengklasifikasikan tuple data yang label kelasnya tidak diketahui**

Klasifikasi dan Prediksi

6/50

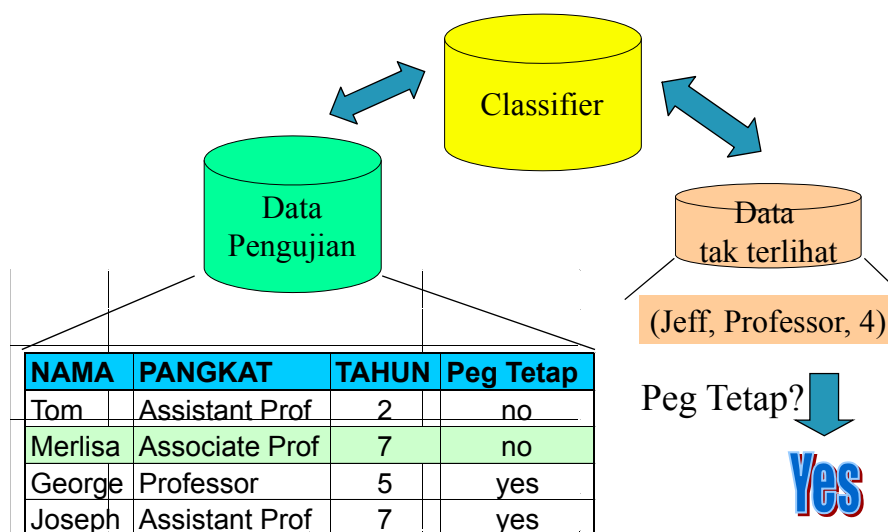
## Proses Klasifikasi (1): Konstruksi Model



Klasifikasi dan Prediksi

7/50

## Proses Klasifikasi (2): Menggunakan Model Dalam Prediksi



Klasifikasi dan Prediksi

8/50

## Pembelajaran Dengan Supervisi vs. Tanpa Supervisi

- **Pembelajaran dengan supervisi (klasifikasi)**
  - Supervisi: data pelatihan (observasi, pengukuran, dsb.) disertakan dengan label yang menunjukkan kelas dari pengamatan
  - Data baru diklasifikasikan berdasarkan himpunan pelatihan
- **Pembelajaran tanpa supervisi (clustering)**
  - Label-label kelas dari data pelatihan tidak dikenal
  - Diberikan sekumpulan pengukuran, observasi, dsb. Dengan tujuan menetapkan keberadaan kelas atau cluster dalam data

Klasifikasi dan Prediksi

9/50

## Isu Yang Berkenaan dgn Klasifikasi & Prediksi (1):Penyiapan Data

- Pembersihan data
  - Memproses awal data dalam upaya untuk mengurangi noise dan menangani nilai-nilai yang hilang
- Analisis relevansi (seleksi fitur)
  - Membuang atribut yang tak relevan atau redundan
- Transformasi data
  - Menggeneralisasi dan/atau menormalisasi data (data punya skala yang sama).

Klasifikasi dan Prediksi

10/50

## Isu Yang Berkenaan dgn Klasifikasi & Prediksi (2):Evaluasi Metoda Klasifikasi

- Akurasi prediksi
- Kecepatan dan skalabilitas
  - Waktu untuk membuat model
  - Waktu untuk menggunakan model
- Kekokohan
  - Penanganan noise dan nilai-nilai yang hilang
- Skalabilitas
  - Efisiensi dalam database yang berada dalam disk

Klasifikasi dan Prediksi

11/50

## Isu Yang Berkenaan dgn Klasifikasi & Prediksi (2):Evaluasi Metoda Klasifikasi

- Penafsiran:
  - Pemahaman dan pengertian yang disediakan oleh model
- Kebaikan dari kaidah
  - Ukuran pohon keputusan
  - Kepadatan dari kaidah klasifikasi

Klasifikasi dan Prediksi

12/50

## Pohon Keputusan

- Pohon keputusan merupakan salah satu tool paling populer untuk klasifikasi karena hasilnya yang bisa dipahami dalam bentuk kaidah keputusan
- Klasifikasi adalah proses penetapan suatu objek baru kedalam kategori atau kelas yang telah didefinisikan sebelumnya.
  - Diberikan sekumpulan record berlabel
  - Buat suatu model (pohon keputusan)
  - Taksir label-label untuk melabeli record-record tak berlabel

Klasifikasi dan Prediksi

13/50

## Pohon Keputusan

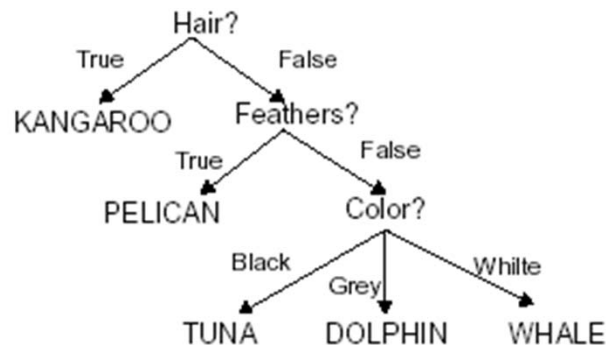
- Suatu **pohon keputusan** biasanya berupa suatu graph berarah dari simpul-simpul dan panah berarah.
- Simpul-simpul seringkali berhubungan dengan suatu pertanyaan atau suatu pengujian

Klasifikasi dan Prediksi

14/50

## Pohon Keputusan

- Sebagai contoh, suatu pohon keputusan digunakan untuk mengklasifikasikan suatu binatang baru apakah memiliki Hair, Feather (bulu), dan Color-nya



Klasifikasi dan Prediksi

15/50

## Struktur Dari Pohon Keputusan

- Suatu pohon keputusan dibangun dari 3 tipe dari simpul: simpul root, simpul perantara, dan simpul leaf.
- Simpul leaf memuat suatu keputusan akhir atau kelas target untuk suatu pohon keputusan
- Simpul root adalah titik awal dari suatu pohon keputusan
- Setiap simpul perantara berhubungan dengan suatu pertanyaan atau pengujian

Klasifikasi dan Prediksi

16/50



## Data Set Pelatihan (org beli komputer)

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Klasifikasi dan Prediksi

17/50

## Attribute Selection by Information Gain Computation

- Class P: buys\_computer = "yes"
- Class N: buys\_computer = "no"
- $I(p, n) = I(9, 5) = 0.940$
- Compute the entropy for *age*:

age	$p_i$	$n_i$	$I(p_i, n_i)$
<=30	2	3	0.971
30...40	4	0	0
>40	3	2	0.971

$$E(\text{age}) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$$

$\frac{5}{14} I(2,3)$  means "age <=30" has 5 out of 14 samples, with 2 yes'es and 3 no's. Hence

$$\text{Gain}(\text{age}) = I(p, n) - E(\text{age}) = 0.246$$

Similarly,

$$\text{Gain}(\text{income}) = 0.029$$

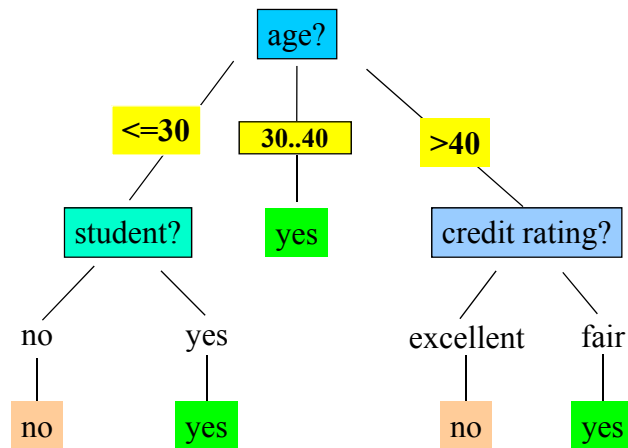
$$\text{Gain}(\text{student}) = 0.151$$

$$\text{Gain}(\text{credit\_rating}) = 0.048$$

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

16

## Output: Pohon Keputusan Untuk “Membeli Komputer”



Klasifikasi dan Prediksi

19/50

## Algoritma Untuk Induksi Pohon Keputusan

- Algoritma dasar (suatu greedy algorithm)
  - Pohon dibangun dalam suatu **metoda rekursif top-down divide-and-conquer** (conquer = menaklukkan)
  - 1. Seluruh contoh pelatihan dimulai dari simpul root, lalu dilakukan pengujian
  - 2. Mencabang ke jalur yang benar berdasarkan hasil pengujian
  - 3. Apakah simpul leaf ditemukan? Jika yes, masukkan contoh ini ke kelas target, jika tidak kembali ke langkah 1.
  - Atribut-atribut berada dalam suatu kategori (jika bernilai kontinu, nilai-nilai tersebut didiskritkan terlebih dahulu)

Klasifikasi dan Prediksi

20/50

## Algoritma Untuk Induksi Pohon Keputusan

- Contoh-contoh dipartisi secara rekursif berdasarkan atribut terpilih
- Atribut-atribut uji dipilih berdasarkan heuristik atau pengukuran statistik (misal, **information gain**)
- Kondisi-kondisi penghentian partisi
  - Seluruh sampel untuk suatu simpul yang diberikan masuk ke kelas yang sama
  - Tidak ada atribut sisa untuk partisi selanjutnya -- **majority voting** diperlakukan untuk klasifikasi leaf
  - Tidak ada sampel tertinggal

Klasifikasi dan Prediksi

21/50

## Ukuran Seleksi Atribut: Information Gain

- Memilih atribut dengan information gain terbesar
- S memuat tupel-tupel  $s_i$  dari kelas  $C_i$  untuk  $i = \{1, \dots, m\}$
- **Ukuran-ukuran informasi** info diperlukan untuk mengklasifikasikan tuple sebarang apapun

$$I(s_1, s_2, \dots, s_m) = - \sum_{i=1}^m \frac{s_i}{S} \log_2 \frac{s_i}{S}$$

Klasifikasi dan Prediksi

22/50

## Ukuran Seleksi Atribut: Information Gain

- **entropy** dari atribut A dengan nilai-nilai  $\{a_1, a_2, \dots, a_v\}$

$$E(A) = \sum_{j=1}^v \frac{s_{1j} + \dots + s_{mj}}{S} I(s_{1j}, \dots, s_{mj})$$

- **Informasi yang diperoleh** melalui pencabangan pada atribut A

$$Gain(A) = I(s_1, s_2, \dots, s_m) - E(A)$$

Klasifikasi dan Prediksi

23/50

## Kaidah Ekstraksi Klasifikasi Dari Pohon

- Sajikan pengetahuan dalam bentuk kaidah **IF-THEN**
- Satu kaidah dibuat untuk setiap path dari root ke leaf
- Setiap pasangan nilai atribut sepanjang path membentuk suatu fabungan
- Simpul leaf menyimpan prediksi kelas
- Kaidah-kaidah lebih mudah dipahami manusia

Klasifikasi dan Prediksi

24/50

## Kaidah Ekstraksi Klasifikasi Dari Pohon

- Contoh-contoh

IF *age* = " $\leq 30$ " AND *student* = "no" THEN  
*buys\_computer* = "no"

IF *age* = " $\leq 30$ " AND *student* = "yes" THEN  
*buys\_computer* = "yes"

IF *age* = "31...40" THEN  
*buys\_computer* = "yes"

IF *age* = " $> 40$ " AND *credit\_rating* = "excellent"  
 THEN *buys\_computer* = "yes"

IF *age* = " $\leq 30$ " AND *credit\_rating* = "fair" THEN  
*buys\_computer* = "no"

Klasifikasi dan Prediksi

25/50

## Contoh klasifikasi Decision Tree: Algo ID3

Database T :

Attribute1	Attribute2	Attribute3	Class
A	70	True	Class1
A	90	True	Class2
A	85	False	Class2
A	95	False	Class2
A	70	False	Class1
B	90	True	Class1
B	78	False	Class1
B	65	True	Class1
B	75	False	Class1
C	80	True	Class2
C	70	True	Class2
C	80	False	Class1
C	80	False	Class1
C	96	False	Class1

Klasifikasi dan Prediksi

26/50

# Klasifikasi Decision Tree: Algo ID3

Basis Data T yang diatas memiliki 14 sampel dengan tiga buah atribut penghitung yang akan dibagi menjadi dua kelas yaitu "Class1" dan "Class2". Class1 memiliki 9 sampel dan Class2 memiliki 5 sampel. Jadi entropy penghitungnya adalah

$$\begin{aligned} I_T(9.5) &= -\sum_{i=1}^2 \frac{S_i}{S} \log_2 \frac{S_i}{S} \\ &= -\left[ \left( \frac{9}{14} \log_2 \frac{9}{14} \right) + \left( \frac{5}{14} \log_2 \frac{5}{14} \right) \right] \\ &= 0.94027 \end{aligned}$$

## Hitungan Nilai Informasi

$$\begin{aligned} \frac{9}{14} \log_2 \frac{9}{14} &= \frac{9}{14} \log_2 (0,64285) \\ &= \frac{9}{14} \frac{\log 0,64285}{\log 2} \\ &= \frac{9}{14} \frac{-0,19189}{0,30103} = -0,40977 \\ \frac{5}{14} \log_2 \frac{5}{14} &= \frac{5}{14} \frac{\log 0,35714}{\log 2} \\ &= \frac{5}{14} \frac{-0,44716}{0,30103} = -0,53050 \end{aligned}$$

$$I_T(9,5) = - \left[ \left( \frac{9}{14} \text{Log}_2 \frac{9}{14} \right) + \left( \frac{5}{14} \text{Log}_2 \frac{5}{14} \right) \right] = 0,94027$$

## Entropy untuk attribute ke 1

Sesudah mendapatkan nilai diatas tersebut, maka penghitungan akan menghitung pada atribut masing-masing untuk mendapatkan nilai gain yang terbesar. Untuk nilai informasi Attribute1 adalah :

Attribute1	Banyaknya sampel pada "Class1"	Banyaknya sampel pada "Class2"	I(Class1,Class2)
A	2	3	0,9709
B	4	0	0
C	3	2	0,9709

## Entropy dan Information Gain untuk attribute ke 1

$$E(\text{Attribute1}) = \frac{S_{ij}}{|S_j|} I_A(2,3) + \frac{S_{ij}}{|S_j|} I_B(4,0) + \frac{S_{ij}}{|S_j|} I_C(3,2)$$

$$E(\text{Attribute1}) = \frac{5}{14}(0,9709) + \frac{4}{14}(0) + \frac{5}{14}(0,9709)$$

$$E(\text{Attribute1}) = 0,6935$$

$$\text{Gain}(\text{Attribute1}) = I(9,5) - E(\text{Attribute1})$$

$$\text{Gain}(\text{Attribute1}) = 0,9402 - 0,6935$$

$$\text{Gain}(\text{Attribute1}) = 0,2467$$

## Nilai Informasi untuk atribut ke 2 (numerik)

Untuk nilai informasi pada Attribute2, data harus diurutkan terlebih dahulu, seperti rumus yang telah kita ketahui diatas. Nilai Attribute2 sesudah diurut adalah: 65,70,75,78,80,85,90,95,96.

Attribute2	Banyaknya sampel pada "Class1"	Banyaknya sampel pada "Class2"	I(Class1,Class2)	E(Attribute2)	Gain (Attribute2)
$\leq 65$	1	0	0	-	-
$> 65$	0	0	-		
$\leq 70$	3	1	0,8112	0,9251	0,0149
$> 70$	6	4	0,9708		
$\leq 75$	4	1	0,7218	0,8947	0,0450
$> 75$	5	4	0,8947		
$\leq 78$	5	1	0,6499	0,8499	0,0901
$> 78$	4	4	1		
$\leq 80$	7	2	0,7641	0,8379	0,1020
$> 80$	2	3	0,9708		
$\leq 85$	7	3	0,8812	0,9151	0,0240
$> 85$	2	2	0,9151		
$\leq 90$	8	4	0,9182	0,9298	0,0102
$> 90$	1	1	1		
$\leq 95$	0	1	0	-	-
$> 95$	0	0	-		

Klasifikasi dan Prediksi

31/50

## Nilai Informasi untuk atribut ke 3 (logical)

Setelah dihitung dan mendapatkan nilai *gain* masing-masing, maka nilai 80 yang akan dijadikan batas pembagiannya, karena memiliki informasi *gain* yang tertinggi yaitu 0,1020.

Untuk nilai *gain* pada attribute3 adalah

Attribute3	Banyaknya sampel pada "Class1"	Banyaknya sampel pada "Class2"	I(Class1,Class2)	E(Attribute3)	Gain (Attribute3)
True	3	3	1	0,892	0,048
False	6	2	0,8112		

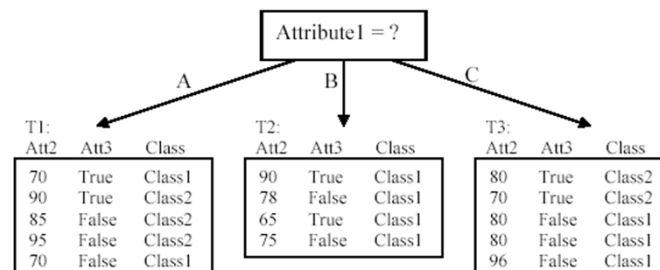
Klasifikasi dan Prediksi

32/50



## Klasifikasi Decision Tree: Algo ID3

Sesudah menghitung informasi *gain* dari semua atribut, maka attribute1 yang dipilih untuk menjadi atribut pembagi karena memiliki nilai informasi *gain* yang tertinggi yakni 0,2467. Cabang pohon yang dibentuk adalah 3 cabang, karena attribute1 memiliki 3 nilai unik(kelas). Setelah dibagi, setiap anak cabang akan memiliki sampel yang nilai attribute1 nya sama dengan cabangnya, seperti yang dapat dilihat pada gambar 2.5.

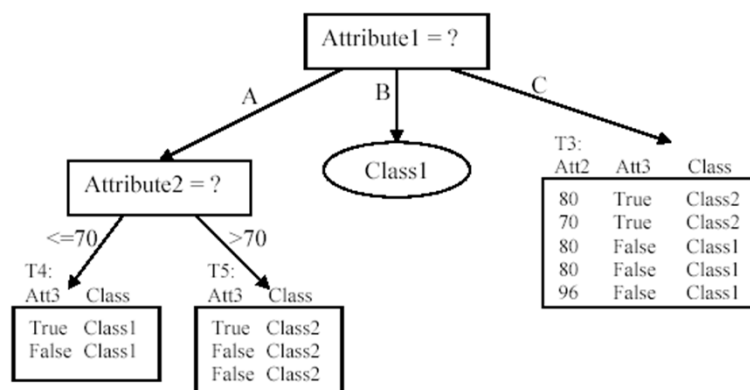


Gambar 2.5 Pohon keputusan awal yang dibentuk  
Klasifikasi dan Prediksi

33/50

## Klasifikasi Decision Tree: Algo ID3

Nilai informasi *gain* yang tertinggi adalah pada attribute2 dengan batas nilai 70 yaitu 0,9709. Pohon keputusan yang dibentuk adalah seperti pada gambar 2.6.

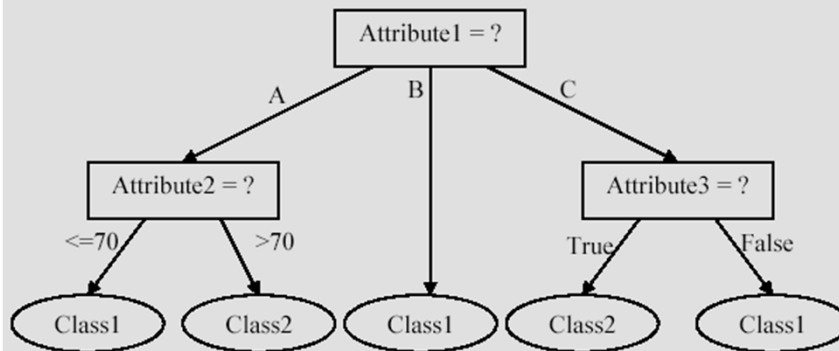


Gambar 2.6 Pohon keputusan yang dibentuk setelah T1 dihitung  
Klasifikasi dan Prediksi

34/50

## Klasifikasi Decision Tree: Algo ID3

Nilai informasi *gain* yang tertinggi adalah pada attribute3 yaitu 0,9709, sehingga didapat pohon keputusan yang baru yaitu seperti pada gambar 2.7.



Gambar 2.7 Pohon keputusan yang dibentuk dari sampel contoh

Klasifikasi dan Prediksi

35/50

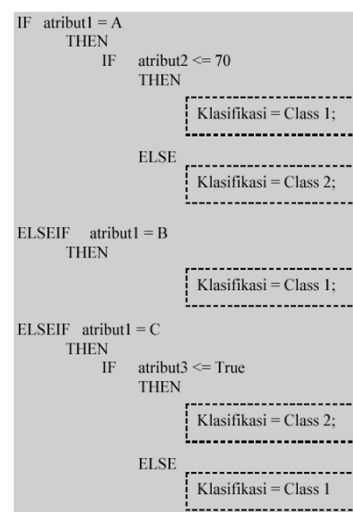
## Decision Rule dan Kegunaannya

- Bila ada pertanyaan :

- Attribute 1 = A
- Attribute 2 = 75
- Attribute 3 = F

- Maka menurut rule yang telah di training hasilnya adalah:

**masuk Class 2**



Klasifikasi dan Prediksi

36/50

## Klasifikasi Dalam Database Besar

- Mengapa induksi pohon keputusan dalam data mining?
  - Kecepatan pembelajaran yang relatif lebih cepat(dari metoda klasifikasi lainnya)
  - Bisa dikonversi ke kaidah klasifikasi yang sederhana dan mudah dipahami
  - Bisa menggunakan query SQL untuk pengaksesan database
  - Akurasi klasifikasi yang bisa dibandingkan dengan metoda lainnya

Klasifikasi dan Prediksi

37/50

## Studi Metoda Dalam Data Mining: Induksi Pohon Keputusan yg Skalabel

- **SLIQ** (EDBT'96 — Mehta dkk.)
  - Membangun suatu indeks untuk setiap atribut dan hanya daftar kelas dan daftar atribut sekarang yang ada dalam memori
- **SPRINT** (VLDB'96 — J. Shafer dkk.)
  - Membangun suatu struktur data daftar atribut
- **PUBLIC** (VLDB'98 — Rastogi & Shim)
  - Mengintegrasikan pohon terpisah dan pohon terpangkas: menghentikan pertumbuhan pohon lebih awal

Klasifikasi dan Prediksi

38/50

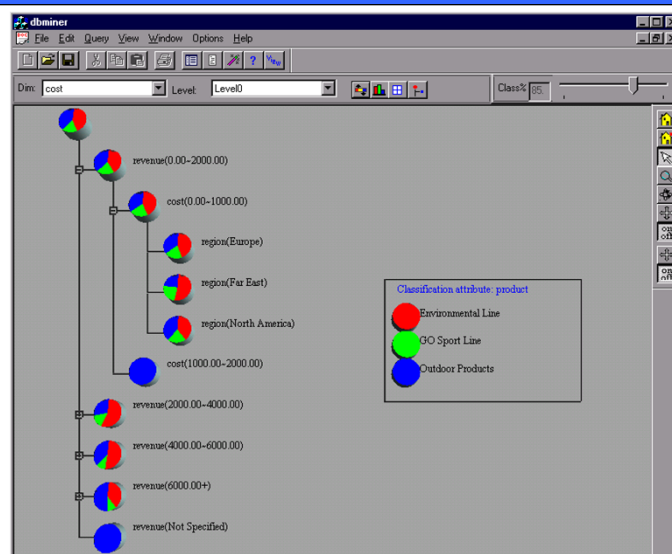
## Studi Metoda Dalam Data Mining: Induksi Pohon Keputusan yg Skalabel

- **RainForest** (VLDB'98 — Gehrke, Ramakrishnan & Ganti)
  - Memisahkan aspek-aspek skalabilitas dari kriteria yang menentukan kualitas dari pohon
  - Membangun suatu daftar AVC (atribut, value, lebel class)

Klasifikasi dan Prediksi

39/50

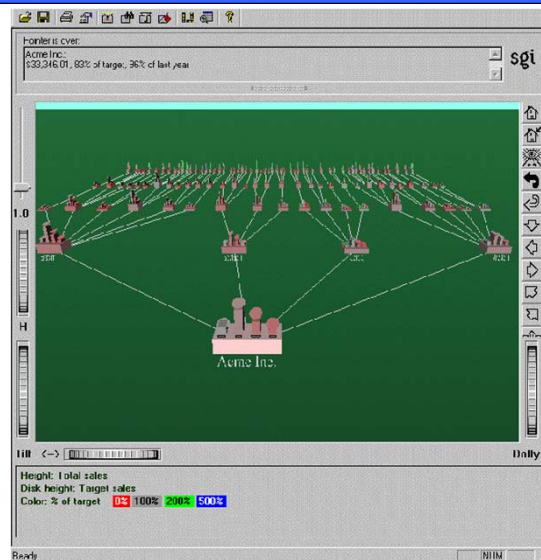
## Presentasi Hasil Klasifikasi



Klasifikasi dan Prediksi

40/50

## Visualisasi Dari Suatu Pohon Keputusan Dalam SGI/MineSet 3.0



Klasifikasi dan Prediksi

41/50

## Teorema Bayesian: Dasar

- Misal  $X$  merupakan suatu sampel data yang label kelasnya tidak dikenal
- Misal  $H$  suatu hipotesa dimana  $X$  masuk ke kelas  $C$
- Untuk problem klasifikasi, tentukan  $P(H/X)$ : peluang dimana hipotesa berlaku untuk sampel data pengamatan  $X$  diberikan

Klasifikasi dan Prediksi

42/50

## Teorema Bayesian: Dasar

- $P(H)$ : peluang hipotesa sebelumnya  $H$  (artinya peluang awal sebelum kita mengamati data apapun, menggambarkan latar belakang pengetahuan)
- $P(X)$ : peluang bahwa sampel data diamati
- $P(X|H)$  : peluang pengamatan sampel  $X$ , apabila diberikan bahwa hipotesa berlaku

Klasifikasi dan Prediksi

43/50

## Teorema Bayesian

- Diberikan data pelatihan  $X$ , *peluang yang mengikuti suatu hipotesa  $H$ ,  $P(H|X)$  memenuhi teorema Bayes*

$$P(H|X) = \frac{P(X|H)P(H)}{P(X)}$$

- Secara informal, ini bisa ditulis sebagai  
posterior = likelihood x prior / evidence
- Hipotesis MAP (maximum posteriori)  
$$h_{MAP} \equiv \arg \max_{h \in H} P(h|D) = \arg \max_{h \in H} P(D|h)P(h).$$
- Prakteknya sukar: membutuhkan pengetahuan awal dari banyak peluang, biaya komputasi cukup berarti

Klasifikasi dan Prediksi

44/50

## Classifier Bayes Naif

- Suatu penyederhanaan asumsi: atribut-atribut bebas bersyarat:

$$P(X | C_i) = \prod_{k=1}^n P(x_k | C_i)$$

- Perkalian dari kejadian, katakan 2 elemen  $x_1$  dan  $x_2$ , dimana kelas sekarang diberikan  $C$ , adalah perkalian dari peluang masing-masing elemen yang diambil terpisah, diberikan kelas yang sama  $P([y_1, y_2], C) = P(y_1, C) * P(y_2, C)$

Klasifikasi dan Prediksi

45/50

## Classifier Bayes Naif

- Tak ada relasi kebergantungan diantara atribut
- Sangat besar mengurangi biaya komputasi, hanya menghitung distribusi kelas.
- Apabila peluang  $P(X|C_i)$  sudah diketahui, berikan  $X$  ke kelas dengan  $P(X|C_i)*P(C_i)$  maksimum

Klasifikasi dan Prediksi

46/50

## Dataset Pelatihan

Class:

C1:buys\_computer=  
'yes'

C2:buys\_computer=  
'no'

Sampel Data

X =(age<=30,  
Income=medium,  
Student=yes  
Credit\_rating=  
Fair)

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
30...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Klasifikasi dan Prediksi

47/50

## Classifier Bayes Naif: Contoh

- Hitung  $P(X/C_i)$  untuk masing-masing kelas

$$P(\text{age}=\text{"<30"} \mid \text{buys\_computer}=\text{"yes"}) = 2/9=0.222$$

$$P(\text{age}=\text{"<30"} \mid \text{buys\_computer}=\text{"no"}) = 3/5 =0.6$$

$$P(\text{income}=\text{"medium"} \mid \text{buys\_computer}=\text{"yes"})= 4/9 =0.444$$

$$P(\text{income}=\text{"medium"} \mid \text{buys\_computer}=\text{"no"}) = 2/5 = 0.4$$

$$P(\text{student}=\text{"yes"} \mid \text{buys\_computer}=\text{"yes"})= 6/9 =0.667$$

$$P(\text{student}=\text{"yes"} \mid \text{buys\_computer}=\text{"no"})= 1/5=0.2$$

$$P(\text{credit\_rating}=\text{"fair"} \mid \text{buys\_computer}=\text{"yes"})=6/9=0.667$$

$$P(\text{credit\_rating}=\text{"fair"} \mid \text{buys\_computer}=\text{"no"})=2/5=0.4$$

**X=(age<=30 ,income =medium, student=yes,credit\_rating=fair)**

$$P(X|C_i) : P(X|\text{buys\_computer}=\text{"yes"})= 0.222 \times 0.444 \times 0.667 \times 0.667 =0.044$$

$$P(X|\text{buys\_computer}=\text{"no"})= 0.6 \times 0.4 \times 0.2 \times 0.4 =0.019$$

**P(X|C<sub>i</sub>)\*P(C<sub>i</sub>) :**

$$P(X|\text{buys\_computer}=\text{"yes"}) * P(\text{buys\_computer}=\text{"yes"})=0.044*(9/14)=0.028$$

$$P(X|\text{buys\_computer}=\text{"no"}) * P(\text{buys\_computer}=\text{"no"})=0.019*(5/14)=0.007$$

**Jadi : X masuk ke kelas "buys\_computer=yes"**

Klasifikasi dan Prediksi

48/50



## Classifier Bayes Naif: Komentar

- Keuntungan:
  - Mudah diimplementasikan
  - Dalam banyak kasus, hasilnya baik
- Kerugian
  - Asumsi: kelas bebas bersyarat, jadi kehilangan akurasi
  - Dalam prakteknya, kebergantungan ada diantara variabel
  - Misal rumah sakit: pasien: Profil: usia, family history, misalkan DeBrito-Cowok, Stece-Cewek  
Symptom: demam, batuk dsb., Penyakit: lung cancer, diabetes dsb
  - Kebergantungan diantara variabel ini tidak dapat dimodelkan dengan Classifier Bayes Naif

Klasifikasi dan Prediksi

49/50

## Metoda Klasifikasi Lainnya

- k-nearest neighbor classifier
- Penalaran berbasis kasus
- Algoritma genetika
- Pendekatan himpunan rough
- Pendekatan himpunan Fuzzy



Klasifikasi dan Prediksi

50/50