



# DATA MINING

ABU SALAM, M.KOM



# PROFIL

## ■ Pendidikan

- SD N Kedungtukang 1 Brebes
- SMP N 5 Brebes
- SMA N 1 Brebes
- S1 dan S2 (Universitas Dian Nuswantoro)

## ■ Research Interest

- Software Engineering (Web App)
- Data Mining

## ■ Activity

- Dosen Fasilkom UDINUS (2009 - sekarang)
- Kepala DIV Software PT DINUSTECH (2008 - sekarang)
- CEO CV Desa Media (2012 - sekarang)

# CONTACT

ALAMAT :

PERUM PERMATA TEMBALANG KAVLING DAHLIA NO 11, KRAMAS TEMBALANG

EMAIL :

■ [masaboe@gmail.com](mailto:masaboe@gmail.com)

■ [abu.salam@dsn.dinus.ac.id](mailto:abu.salam@dsn.dinus.ac.id)

YM : mas\_aboe@yahoo.com

FB : [masaboe@yahoo.com](mailto:masaboe@yahoo.com)

HP : 0817244958

# OUTLINE

1. **Pengenalan Data Mining**
2. Proses Data Mining
3. Evaluasi dan Validasi pada Data Mining
4. Metode dan Algoritma Data Mining
5. Penelitian Data Mining

# KOMPONEN PENILAIAN

■ Kehadiran : 75 %

■ Tugas : 30%

■ UTS : 35%

■ UAS : 35%

## Range Nilai

**A** : 85 - 99

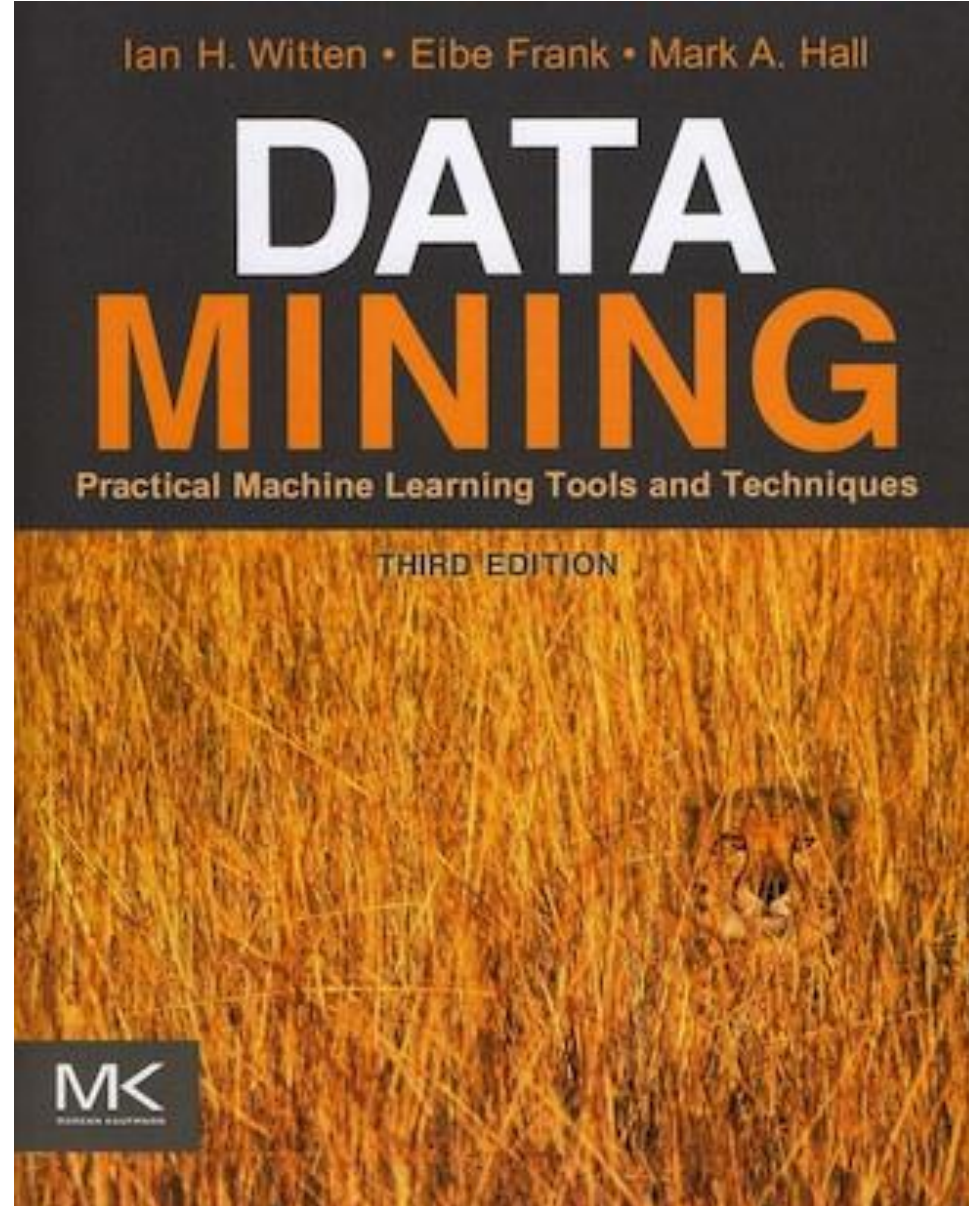
**B** : 70 - 84

**C** : 56 – 69

**D** : 40 – 55

**E** : 0 - 39

# TEXTBOOKS



# PRETEST

1. Jelaskan apa yang dimaksud dengan **data mining**?
2. Sebutkan **peran data mining** dan **algoritma apa saja** yang mendukung peran data mining tersebut?
3. Berikan contoh **penerapan ataupun penelitian data mining** ?

# PENGENALAN DATA MINING





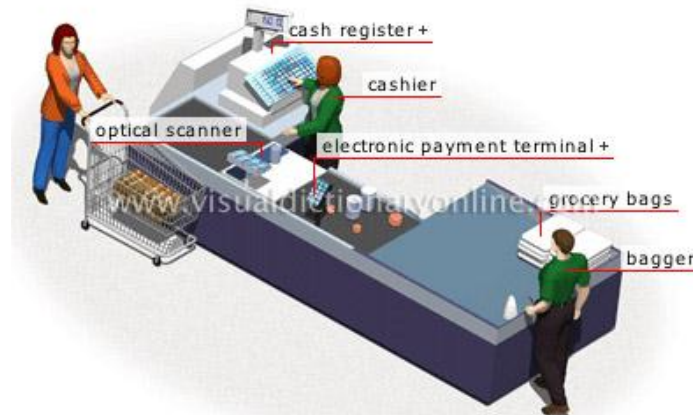
# PENGENALAN DATA MINING

1. Apa itu Data Mining?
2. Peran Utama Data Mining?
3. Algoritma Data Mining?

# MINING? WAREHOUSING?



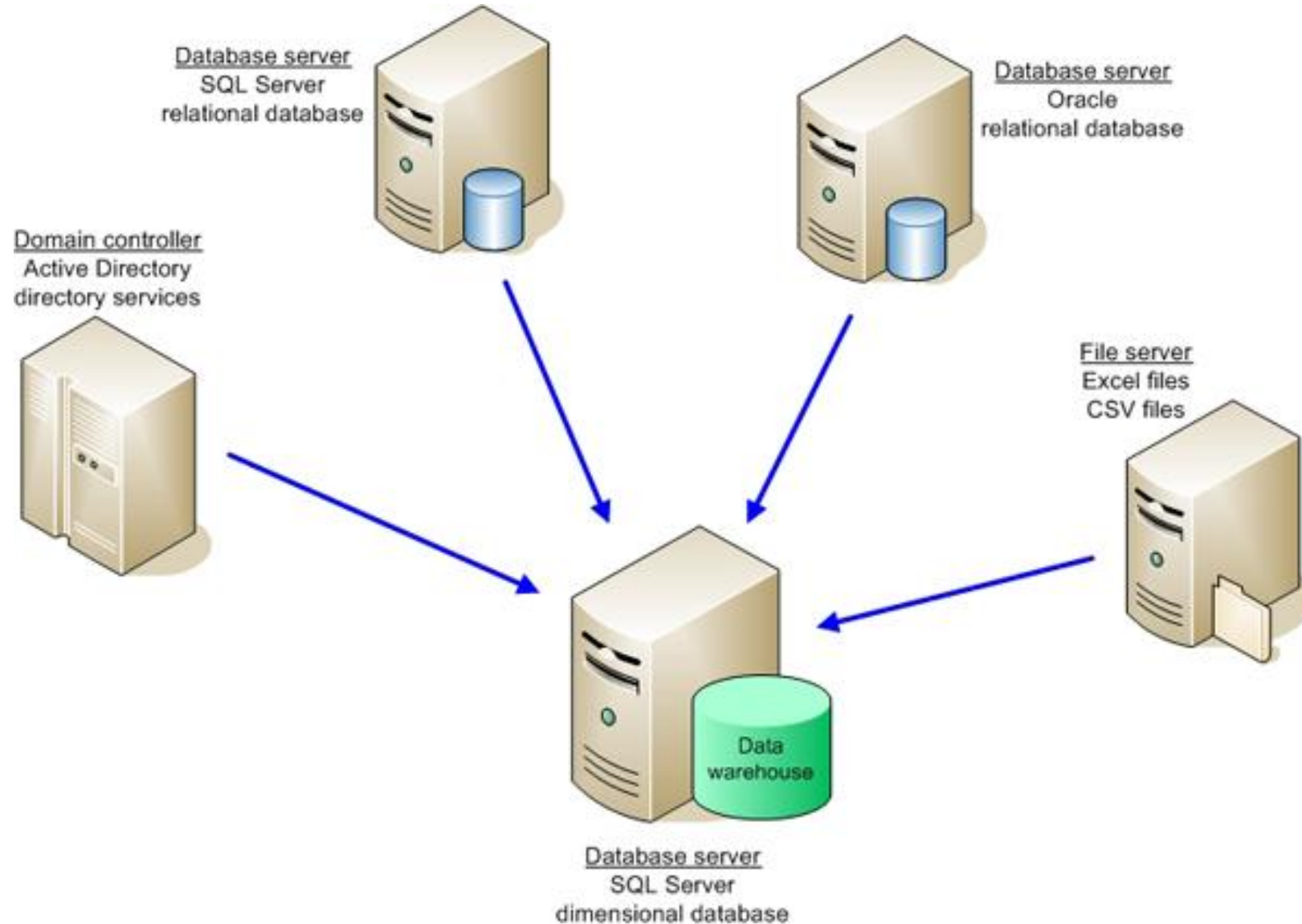
# THE WORLD OF DATA







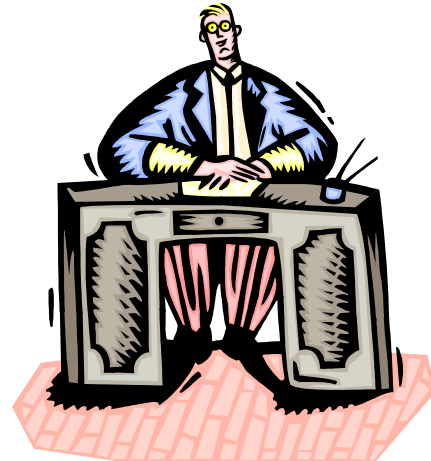
# HETEROGENEOUS DATA



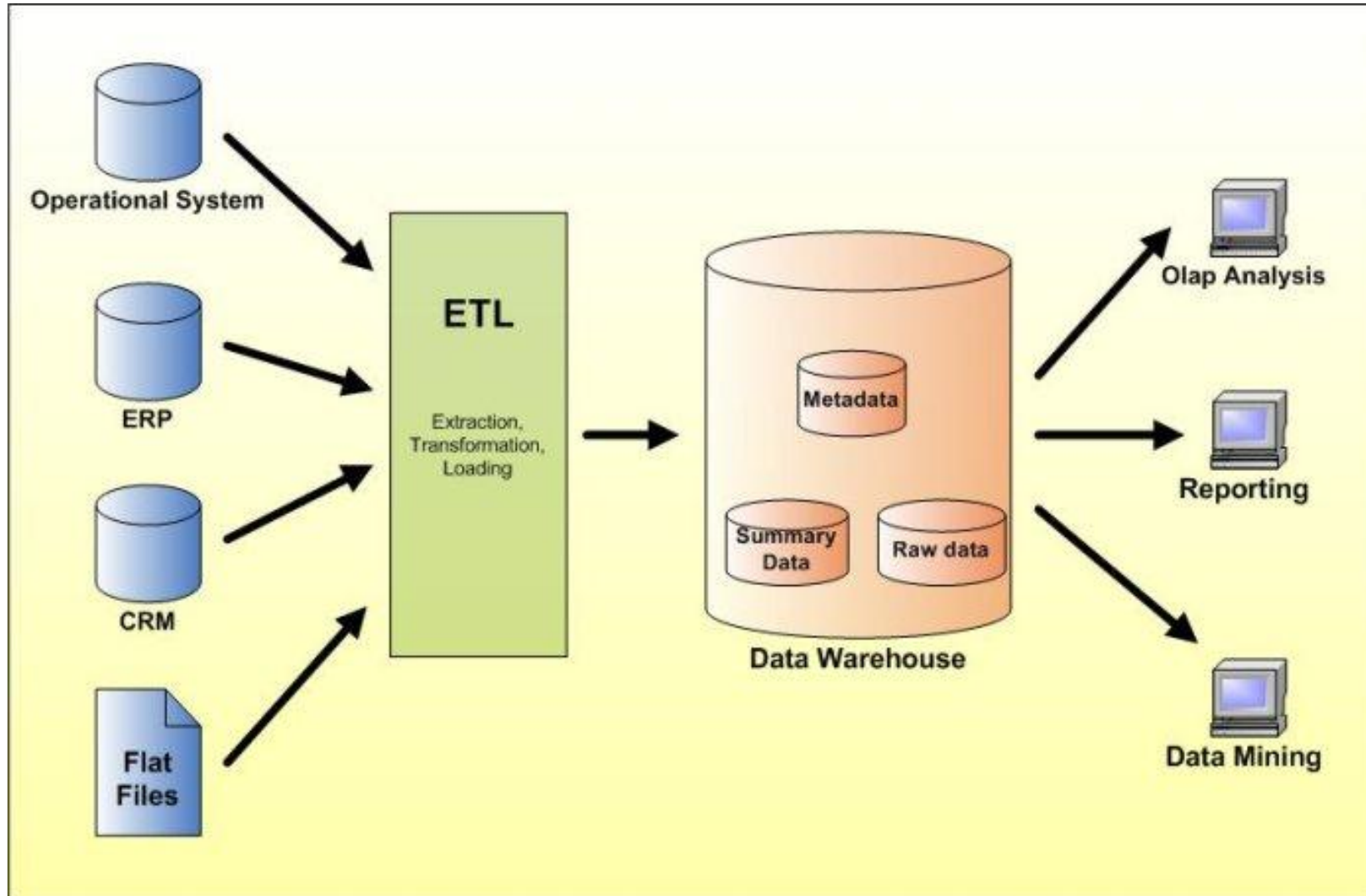
# DATA RICH, INFORMATION POOR



# BUSINESS INTELLIGENCE

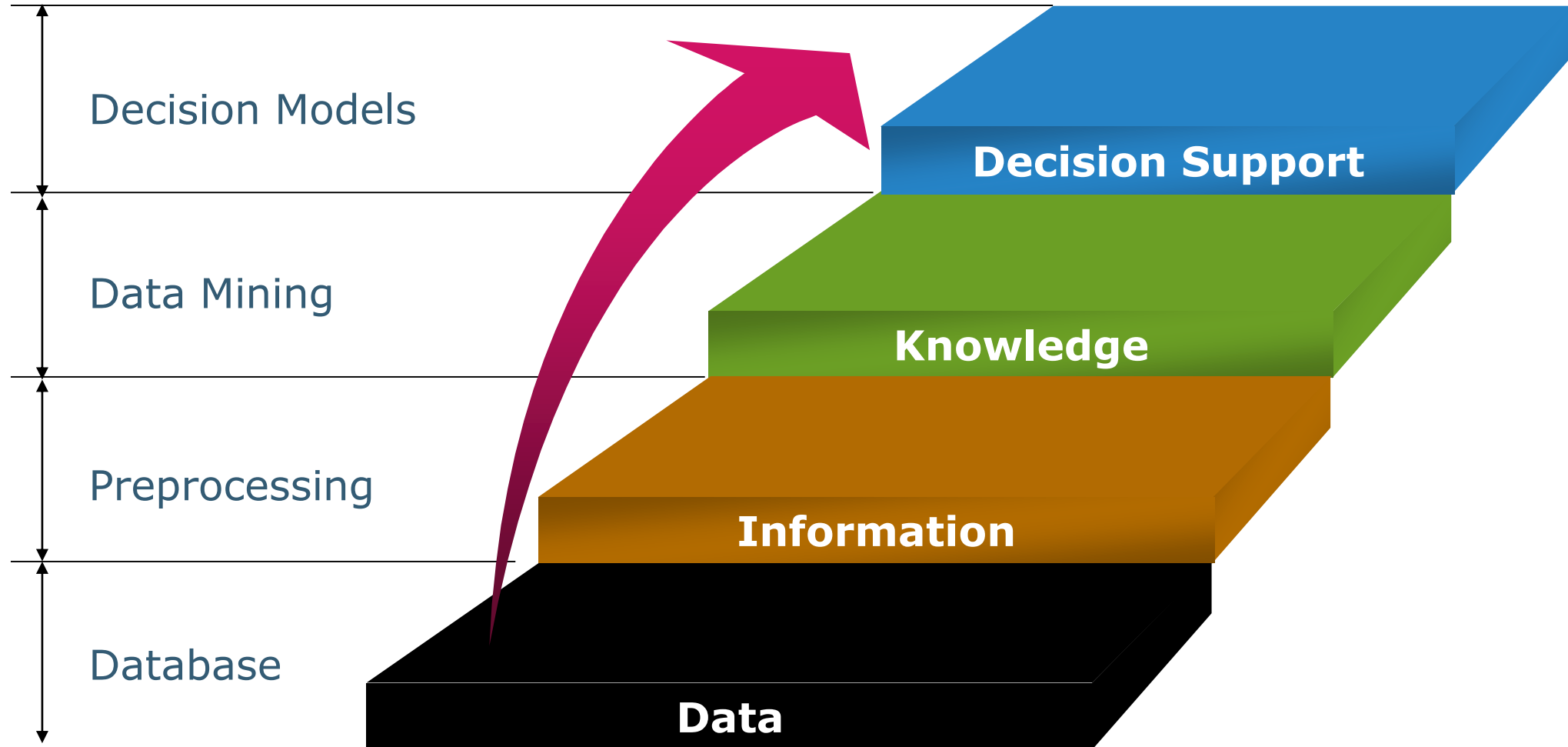


# DATA INTEGRATION & ANALYSIS





# FROM DATA TO INTELLIGENCE



# IT IS ALL ABOUT DATA ...

Retail

Financial Institutions

WWW

Healthcare

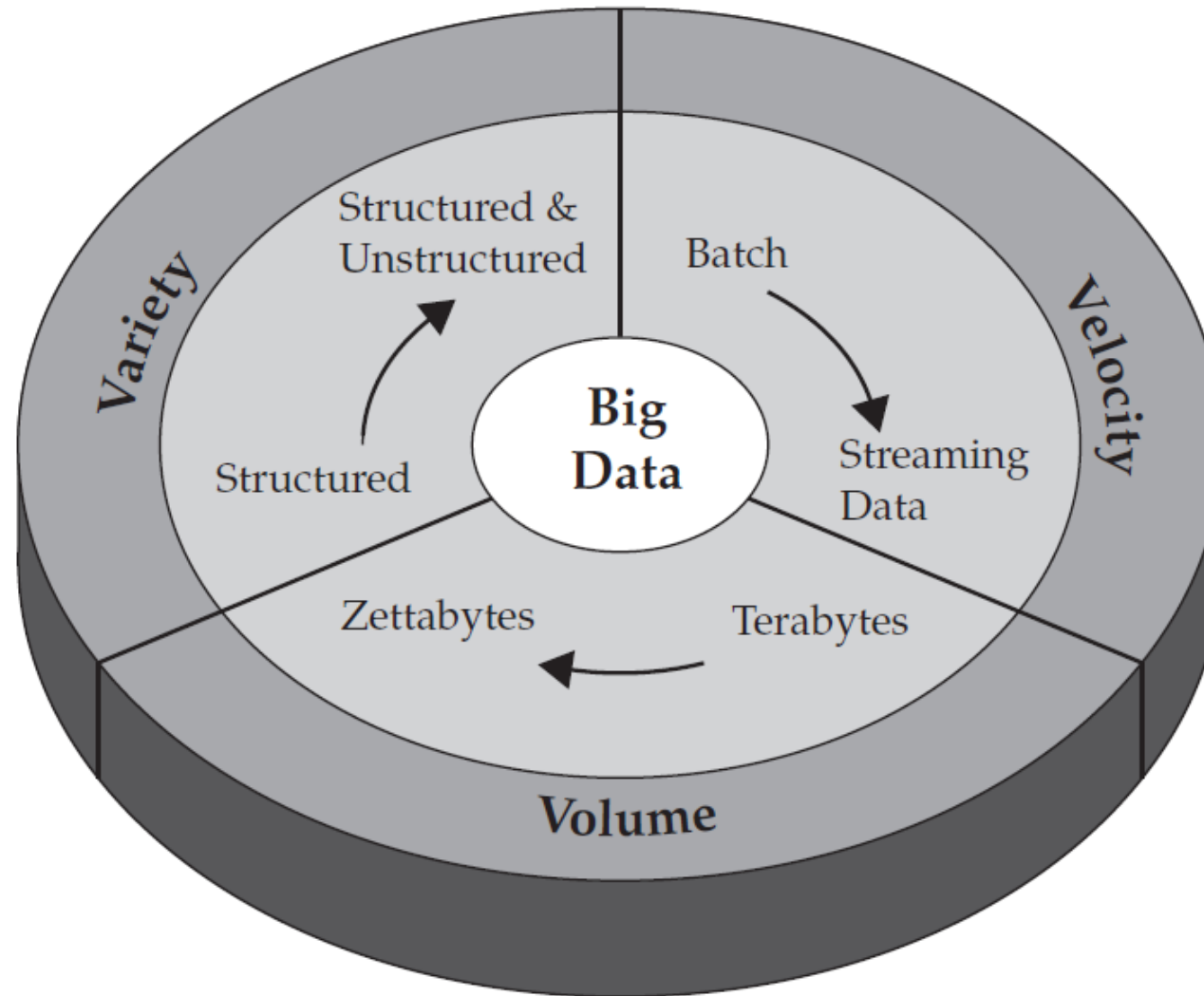
Consulting Companies

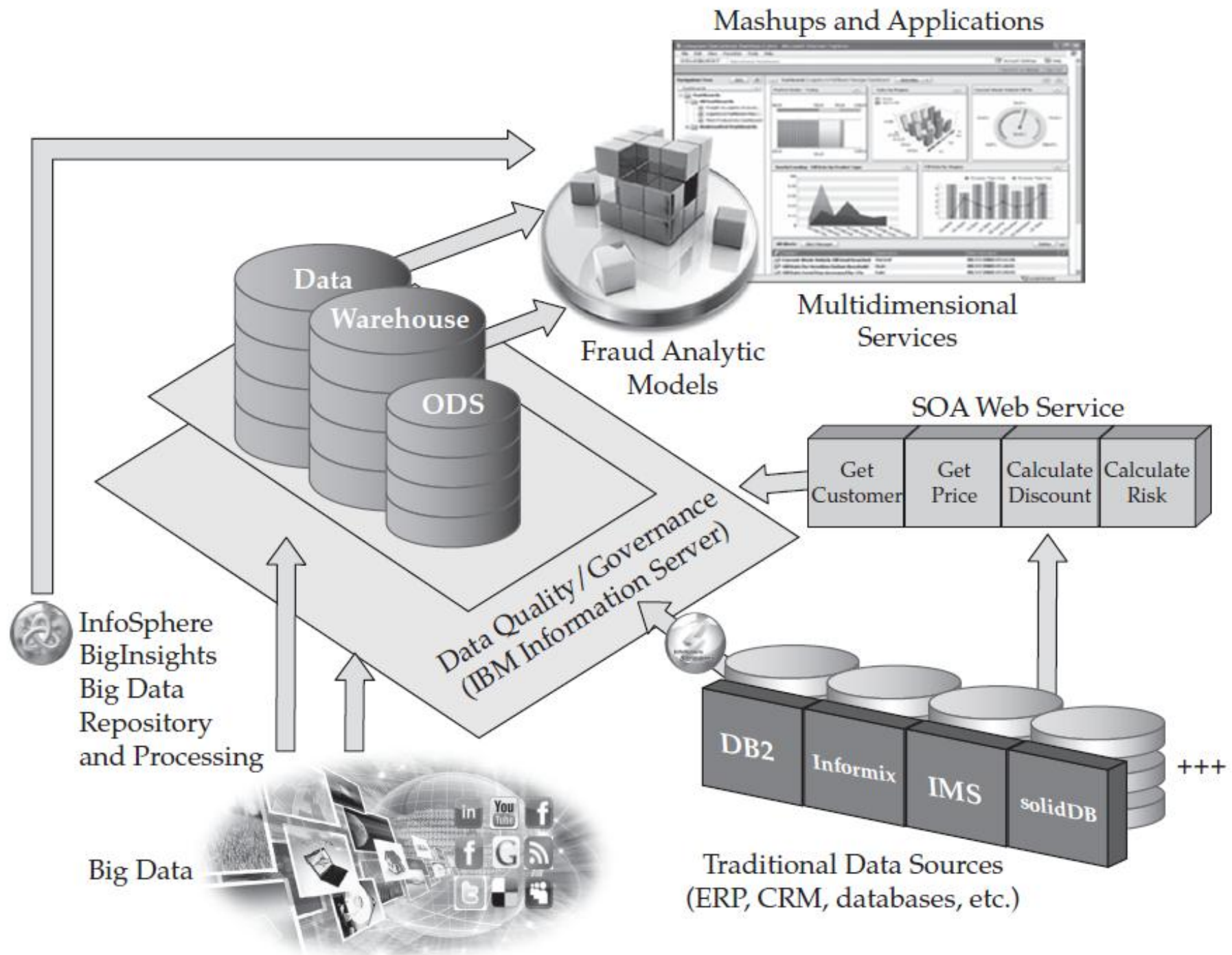
Government

Bioinformatics

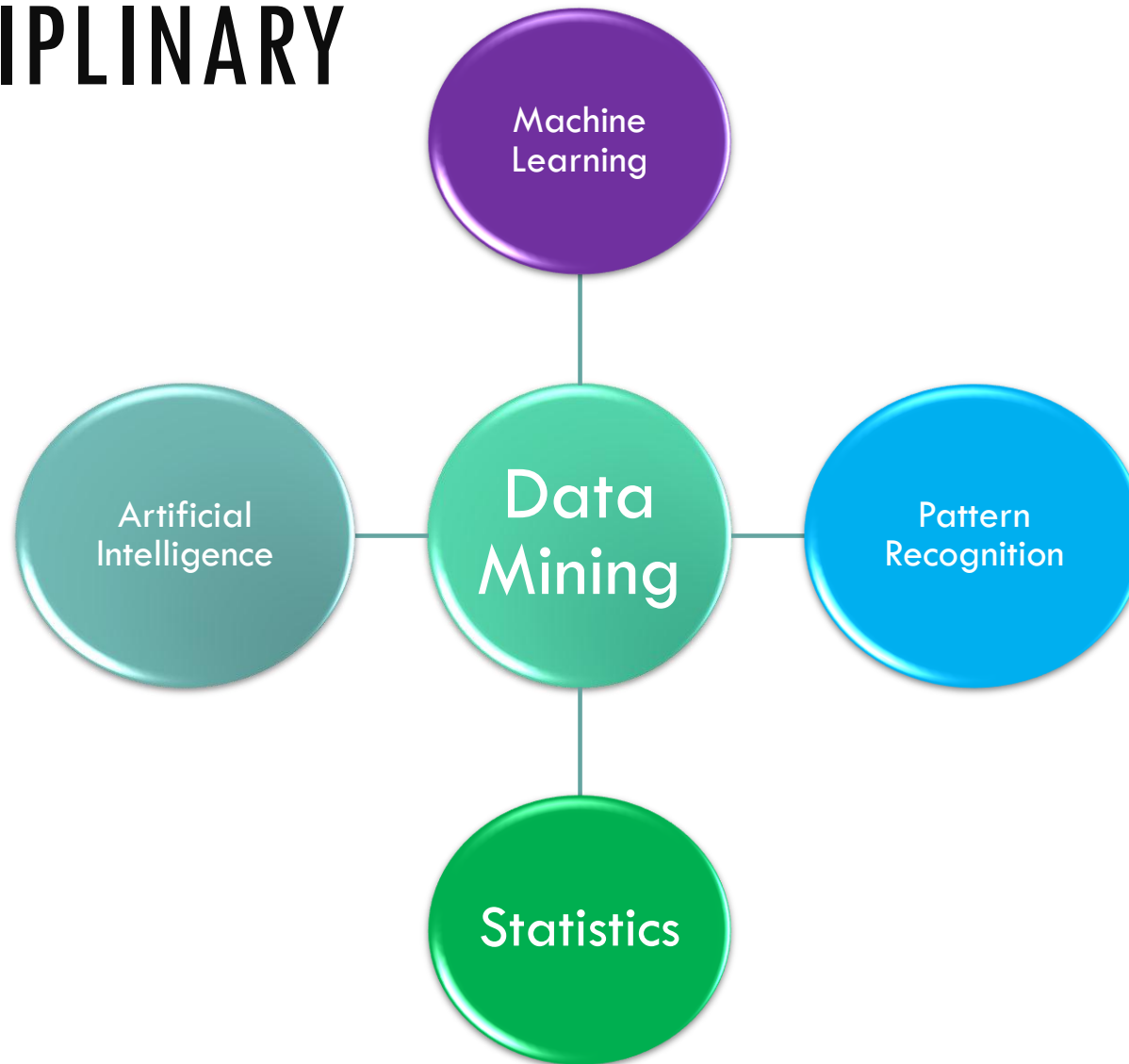
Telecommunication

# BIG DATA

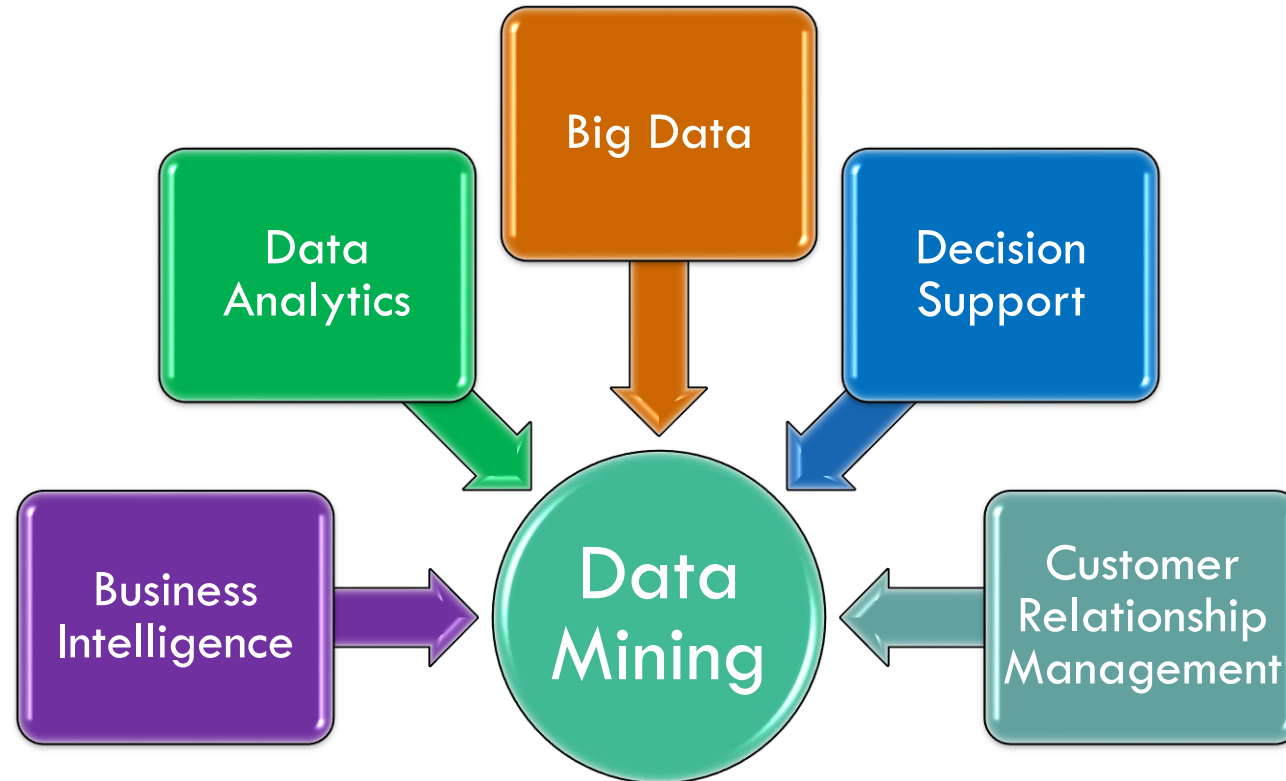




# INTERDISCIPLINARY



# UBIQUITOUS



**APA ITU DATA MINING ?**



# MENGAPA DATA MINING?

Manusia dalam suatu organisasi, sadar atau tidak sadar telah memproduksi berbagai data yang jumlahnya sangat besar

- Contoh data: bisnis, kedokteran, ekonomi, geografi, olahraga, ...

Pada dasarnya, data adalah entitas yang tidak memiliki arti, meskipun kemungkinan memiliki nilai di dalamnya



# APA ITU DATA MINING?

Disiplin ilmu yang mempelajari **metode untuk mengekstrak** pengetahuan atau **menemukan pola** dari suatu data

1. **Data**: fakta yang terekam dan tidak membawa arti
2. **Pengetahuan**: pola, aturan atau model yang muncul dari data

Sehingga Data mining sering disebut **Knowledge Discovery in Database (KDD)**

Konsep Transformasi

**Data** → **Informasi** → **Pengetahuan**



# DATA

- **Tidak membawa arti**, merupakan kumpulan dari fakta-fakta tentang suatu kejadian
- Suatu catatan terstruktur dari suatu transaksi
- Merupakan materi penting dalam **membentuk informasi**

# PENGETAHUAN

Gabungan dari suatu **pengalaman, nilai, informasi kontekstual dan juga pandangan pakar** yang memberikan suatu framework untuk mengevaluasi dan menciptakan pengalaman baru dan informasi (*Thomas H. Davenport, Laurence Prusak*)

Bisa berupa **solusi pemecahan suatu masalah, petunjuk suatu pekerjaan** dan ini bisa ditingkatkan nilainya, dipelajari dan juga bisa diajarkan kepada yang lain

## ***Data - Informasi – Pengetahuan***

### **Data** Kehadiran Pegawai

NIP	TGL	DATANG	PULANG
1103	02/12/2004	07:20	15:40
1142	02/12/2004	07:45	15:33
1156	02/12/2004	07:51	16:00
1173	02/12/2004	08:00	15:15
1180	02/12/2004	07:01	16:31
1183	02/12/2004	07:49	17:00

## ***Data - Informasi – Pengetahuan***

**Informasi** Akumulasi Bulanan Kehadiran Pegawai

NIP	Masuk	Alpa	Cuti	Sakit	Telat
1103	22				
1142	18	2		2	
1156	10	1	11		
1173	12	5			5
1180	10			12	

## ***Data - Informasi – Pengetahuan***

**Informasi** Kondisi Kehadiran Mingguan Pegawai

	Senin	Selasa	Rabu	Kamis	Jumat
Terlambat	7	0	1	0	5
Pulang Cepat	0	1	1	1	8
Izin	3	0	0	1	4
Alpa	1	0	2	0	2

## ***Data - Informasi – Pengetahuan***

- Pengetahuan tentang kebiasaan pegawai dalam jam datang/pulang kerja
- Pengetahuan tentang bagaimana teknik meningkatkan kehadiran pegawai → **kebijakan**



## DATA - INFORMASI - PENGETAHUAN - KEBIJAKAN

**Kebijakan penataan jam kerja** karyawan khusus untuk hari senin dan jumat

Peraturan jam kerja:

- Hari Senin dimulai jam 10:00
- Hari Jumat diakhiri jam 14:00
- Sisa jam kerja dikompensasi ke hari lain:
  1. Senin pulang setelah maghrib, toh jalanan jakarta macet total di sore hari (**bayar hutang 2 jam**)
  2. Rabu dan kamis bayar hutang setengah jam di pagi hari dan setengah jam di sore hari (**bayar hutang 2 jam**)



# DEFINISI DATA MINING

Melakukan **ekstraksi** untuk mendapatkan **informasi penting** yang sifatnya **implisit** dan sebelumnya tidak diketahui, dari suatu data (*Witten et al., 2011*)

Kegiatan yang meliputi pengumpulan, pemakaian data historis untuk **menemukan keteraturan, pola dan hubungan** dalam set data berukuran besar (*Santosa, 2007*)

# DEFINISI DATA MINING

The analysis of (often large) observational data sets **to find unsuspected relationships** and to **summarize the data** in novel ways that are both understandable and useful to the data owner (*Han & Kamber, 2001*)

The process of **discovering meaningful new correlations, patterns and trends** by sifting through large amounts of data stored in repositories, using pattern recognition technologies as well as statistical and mathematical techniques (*Gartner Group*)

# IRISAN BIDANG ILMU DATA MINING

## 1. Statistik:

- Lebih bersifat teori
- Fokus ke pengujian hipotesis

## 2. Machine Learning:

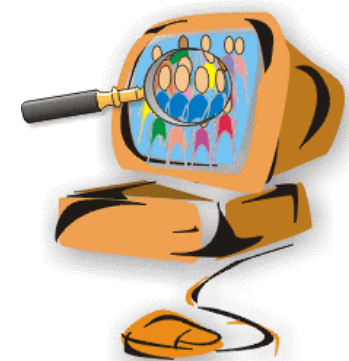
- Lebih bersifat heuristik
- Fokus pada perbaikan performansi dari suatu teknik learning

## 3. Data Mining:

- Gabungan teori dan heuristik
- Fokus pada seluruh proses penemuan knowledge dan pola
- Termasuk data cleaning, learning dan visualisasi hasilnya

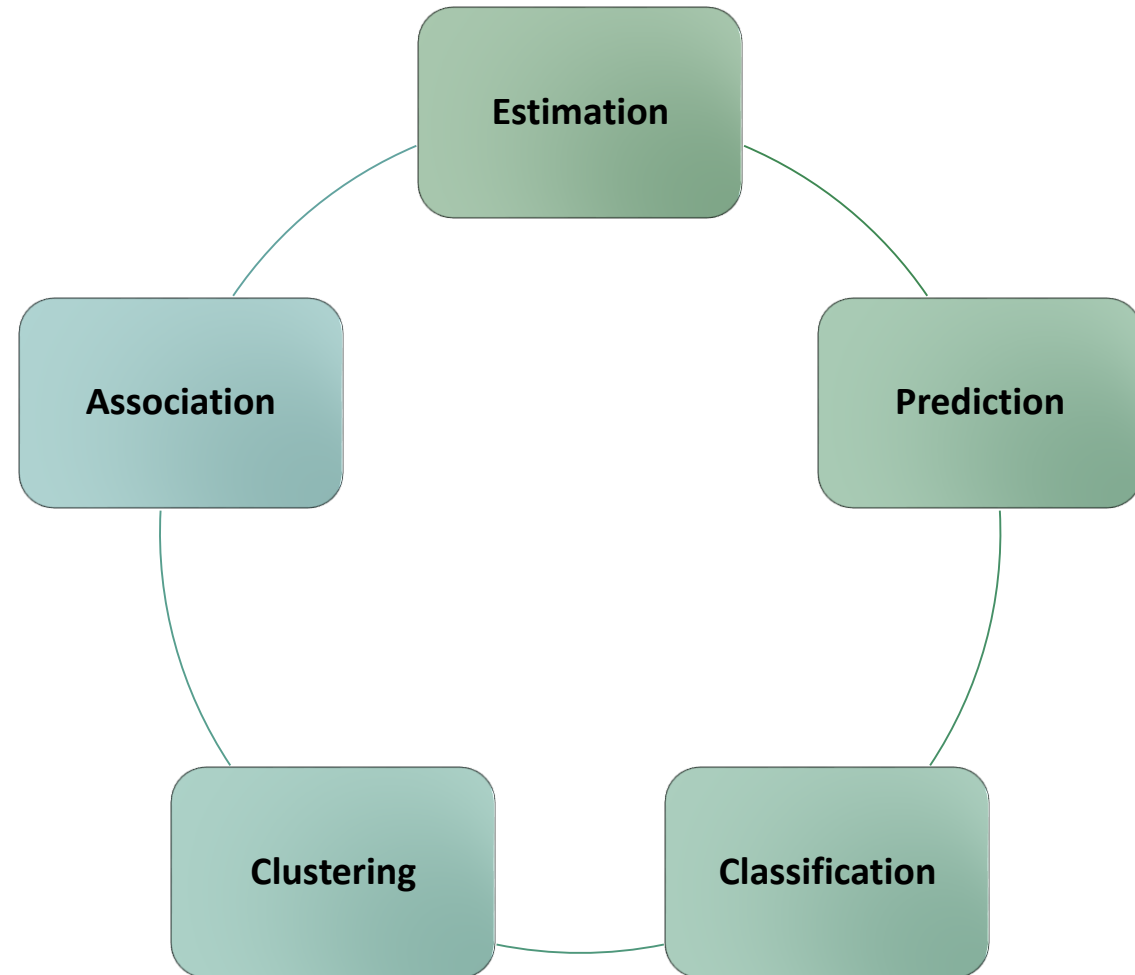


# PERAN UTAMA DATA MINING




# PERAN UTAMA DATA MINING

1. Estimation
2. Prediction
3. Classification
4. Clustering
5. Association



# DATASET WITH ATTRIBUTE AND CLASS



Attribute

Class/Label

	Sepal Length (cm)	Sepal Width (cm)	Petal Length (cm)	Petal Width (cm)	Type
1	5.1	3.5	1.4	0.2	<i>Iris setosa</i>
2	4.9	3.0	1.4	0.2	<i>Iris setosa</i>
3	4.7	3.2	1.3	0.2	<i>Iris setosa</i>
4	4.6	3.1	1.5	0.2	<i>Iris setosa</i>
5	5.0	3.6	1.4	0.2	<i>Iris setosa</i>
...					
51	7.0	3.2	4.7	1.4	<i>Iris versicolor</i>
52	6.4	3.2	4.5	1.5	<i>Iris versicolor</i>
53	6.9	3.1	4.9	1.5	<i>Iris versicolor</i>
54	5.5	2.3	4.0	1.3	<i>Iris versicolor</i>
55	6.5	2.8	4.6	1.5	<i>Iris versicolor</i>
...					
101	6.3	3.3	6.0	2.5	<i>Iris virginica</i>
102	5.8	2.7	5.1	1.9	<i>Iris virginica</i>
103	7.1	3.0	5.9	2.1	<i>Iris virginica</i>
104	6.3	2.9	5.6	1.8	<i>Iris virginica</i>
105	6.5	3.0	5.8	2.2	<i>Iris virginica</i>
...					

## ESTIMASI WAKTU PENGIRIMAN PIZZA

Customer	Jumlah Pesanan (P)	Jumlah Bangjo (B)	Jarak (J)	Waktu Tempuh (T)
1	3	3	3	16
2	1	7	4	20
3	2	4	6	18
4	4	6	8	36
...				
1000	2	4	2	12

$$\text{Waktu Tempuh (T)} = 0.48P + 0.23B + 0.5J$$

## PENENTUAN KELULUSAN MAHASISWA

NIM	Gender	Nilai UN	Asal Sekolah	IPS1	IPS2	IPS3	IPS 4	...	Lulus Tepat Waktu
10001	L	28	SMAN 2	3.3	3.6	2.89	2.9		Ya
10002	P	27	SMA DK	4.0	3.2	3.8	3.7		Tidak
10003	P	24	SMAN 1	2.7	3.4	4.0	3.5		Tidak
10004	L	26.4	SMAN 3	3.2	2.7	3.6	3.4		Ya
...									
...									
11000	L	23.4	SMAN 5	3.3	2.8	3.1	3.2		Ya

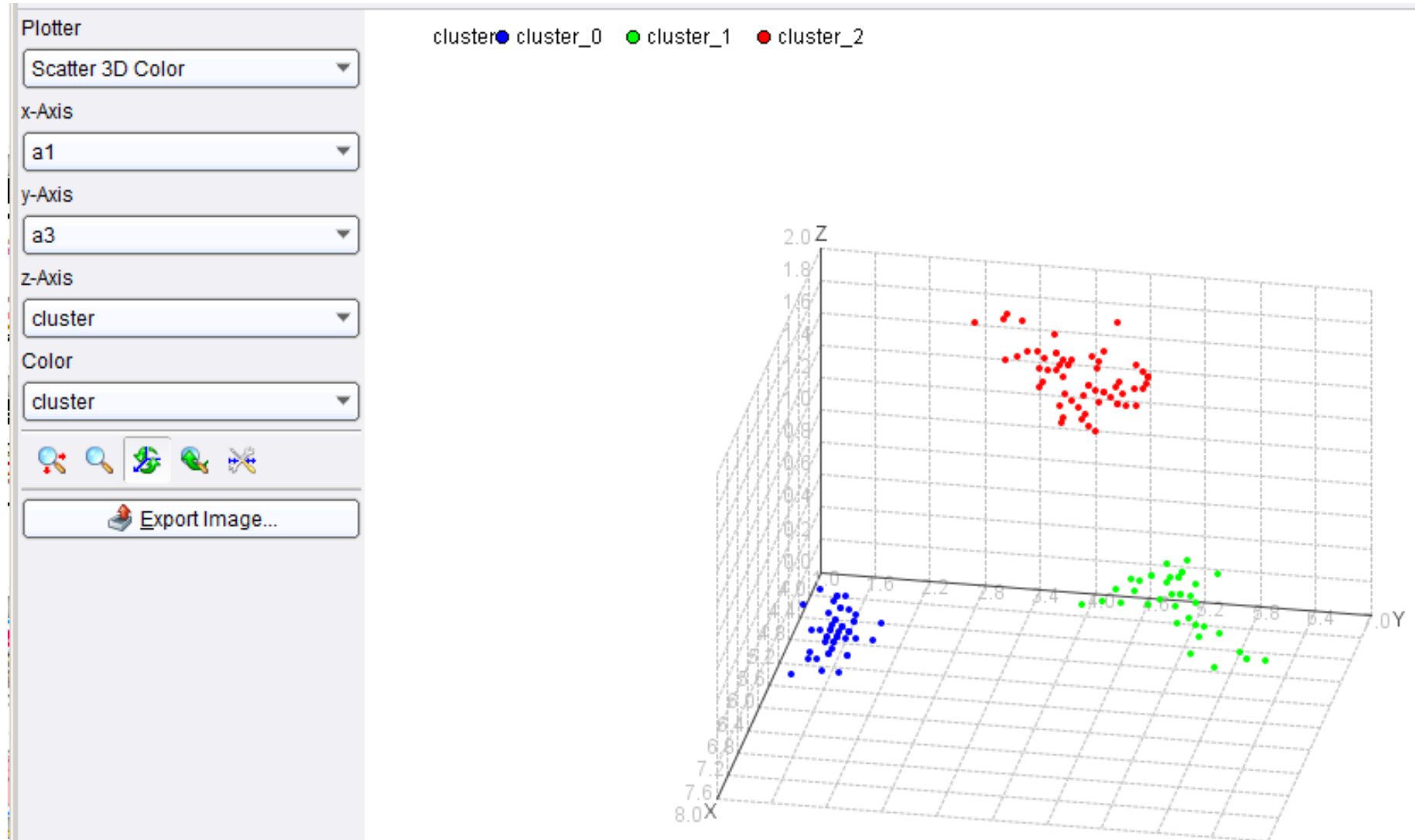


# KLASTERING BUNGA IRIS

ExampleSet (150 examples, 2 special attributes, 4 regular attributes)

Row No.	id	label	a1	a2	a3	a4
1	id_1	Iris-setosa	5.100	3.500	1.400	0.200
2	id_2	Iris-setosa	4.900	3	1.400	0.200
3	id_3	Iris-setosa	4.700	3.200	1.300	0.200
4	id_4	Iris-setosa	4.600	3.100	1.500	0.200
5	id_5	Iris-setosa	5	3.600	1.400	0.200
6	id_6	Iris-setosa	5.400	3.900	1.700	0.400
7	id_7	Iris-setosa	4.600	3.400	1.400	0.300
8	id_8	Iris-setosa	5	3.400	1.500	0.200
9	id_9	Iris-setosa	4.400	2.900	1.400	0.200
10	id_10	Iris-setosa	4.900	3.100	1.500	0.100
11	id_11	Iris-setosa	5.400	3.700	1.500	0.200
12	id_12	Iris-setosa	4.800	3.400	1.600	0.200
13	id_13	Iris-setosa	4.800	3	1.400	0.100
14	id_14	Iris-setosa	4.300	3	1.100	0.100
15	id_15	Iris-setosa	5.800	4	1.200	0.200
16	id_16	Iris-setosa	5.700	4.400	1.500	0.400
17	id_17	Iris-setosa	5.400	3.900	1.300	0.400
18	id_18	Iris-setosa	5.100	3.500	1.400	0.300
19	id_19	Iris-setosa	5.700	3.800	1.700	0.300
20	id_20	Iris-setosa	5.100	3.800	1.500	0.300
21	id_21	Iris-setosa	5.400	3.400	1.700	0.200
22	id_22	Iris-setosa	5.100	3.700	1.500	0.400
23	id_23	Iris-setosa	4.600	3.600	1	0.200
24	id_24	Iris-setosa	5.100	3.300	1.700	0.500

# KLASTERING BUNGA IRIS



# ALGORITMA DATA MINING (DM)

## 1. **Estimation** (Estimasi):

Linear Regression, [Neural Network](#), Support Vector Machine, etc

## 2. **Prediction/Forecasting** (Prediksi/Peramalan):

Linear Regression, [Neural Network](#), Support Vector Machine, etc

## 3. **Classification** (Klasifikasi):

Naive Bayes, K-Nearest Neighbor, [C4.5](#), ID3, CART, Linear Discriminant Analysis, etc

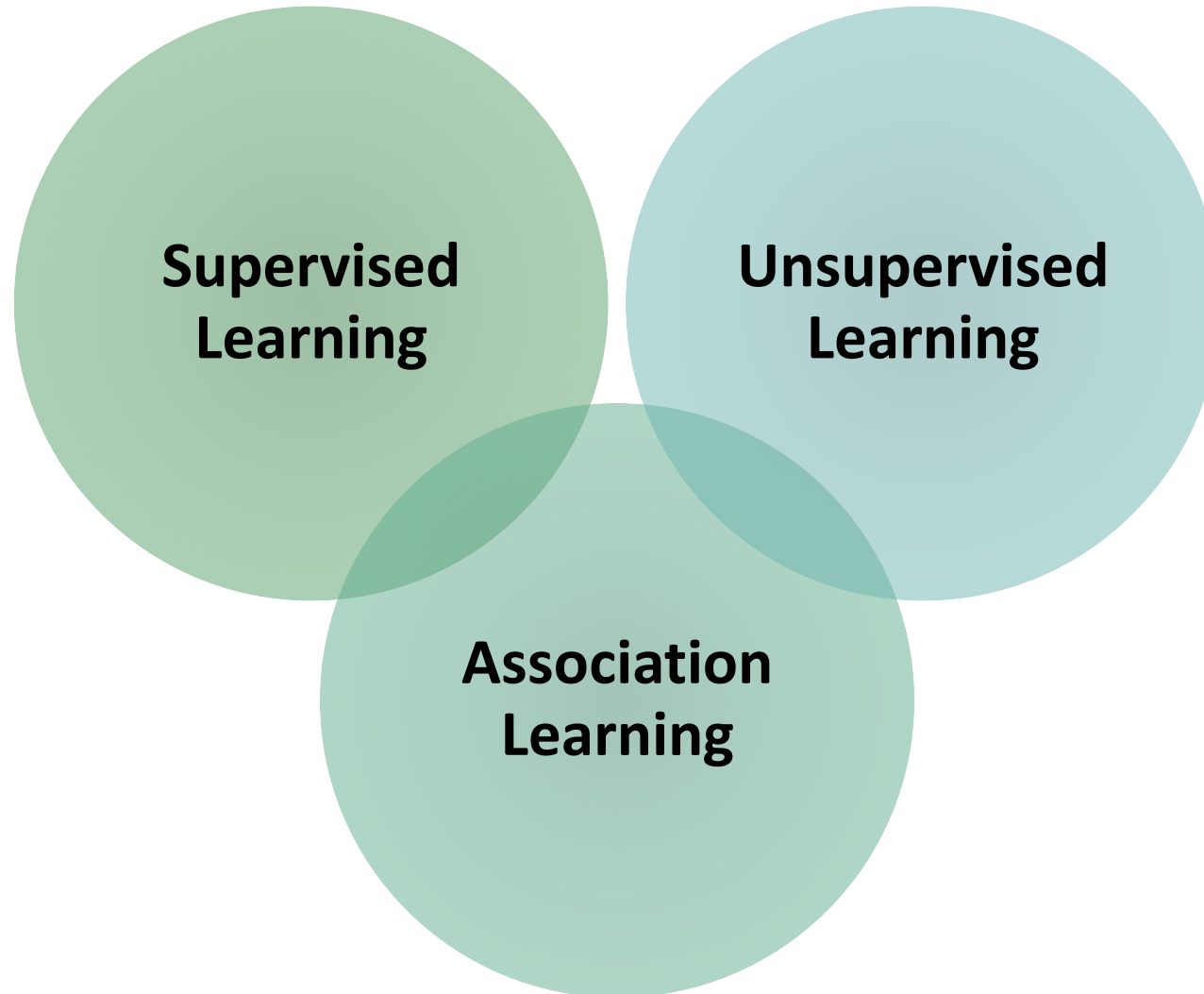
## 4. **Clustering** (Klastering):

[K-Means](#), K-Medoids, Self-Organizing Map (SOM), Fuzzy C-Means, etc

## 5. **Association** (Asosiasi):

FP-Growth, [A Priori](#), etc

# METODE LEARNING PADA ALGORITMA DM



# METODE LEARNING PADA ALGORITMA DM

## 1. **Supervised** Learning (Pembelajaran dengan Guru):

- Sebagian besar algoritma data mining (estimation, prediction/forecasting, classification) adalah supervised learning
- Variabel yang menjadi **target/label/class** ditentukan
- Algoritma melakukan proses belajar berdasarkan **nilai dari variabel target** yang terasosiasi dengan nilai dari variable prediktor


# METODE LEARNING PADA ALGORITMA DM

## 2. **Unsupervised** Learning (Pembelajaran tanpa Guru):

- Algoritma data mining mencari pola dari **semua variable (atribut)**
- Variable (atribut) yang menjadi **target/label/class tidak ditentukan (tidak ada)**
- Algoritma **clustering** adalah algoritma unsupervised learning

# DATASET WITH ATTRIBUTE (NO CLASS)

Attribute



	Sepal Length (cm)	Sepal Width (cm)	Petal Length (cm)	Petal Width (cm)
1	5.1	3.5	1.4	0.2
2	4.9	3.0	1.4	0.2
3	4.7	3.2	1.3	0.2
4	4.6	3.1	1.5	0.2
5	5.0	3.6	1.4	0.2
...				
51	7.0	3.2	4.7	1.4
52	6.4	3.2	4.5	1.5
53	6.9	3.1	4.9	1.5
54	5.5	2.3	4.0	1.3
55	6.5	2.8	4.6	1.5
...				
101	6.3	3.3	6.0	2.5
102	5.8	2.7	5.1	1.9
103	7.1	3.0	5.9	2.1
104	6.3	2.9	5.6	1.8
105	6.5	3.0	5.8	2.2
...				

# METODE LEARNING PADA ALGORITMA DM

## 3. **Association** Learning (Pembelajaran untuk Asosiasi Atribut)

- Proses learning pada algoritma asosiasi (*association rule*) agak berbeda karena tujuannya adalah untuk mencari **atribut yang muncul bersamaan dalam satu transaksi**
- Algoritma asosiasi biasanya untuk analisa transaksi belanja, dengan konsep utama adalah mencari “**produk/item mana yang dibeli bersamaan**”
- Pada pusat perbelanjaan **banyak produk yang dijual**, sehingga pencarian seluruh asosiasi produk memakan **cost tinggi**, karena sifatnya yang **kombinatorial**
- Algoritma *association rule* seperti **apriori algorithm**, dapat memecahkan masalah ini dengan efisien



# DATASET TRANSACTION

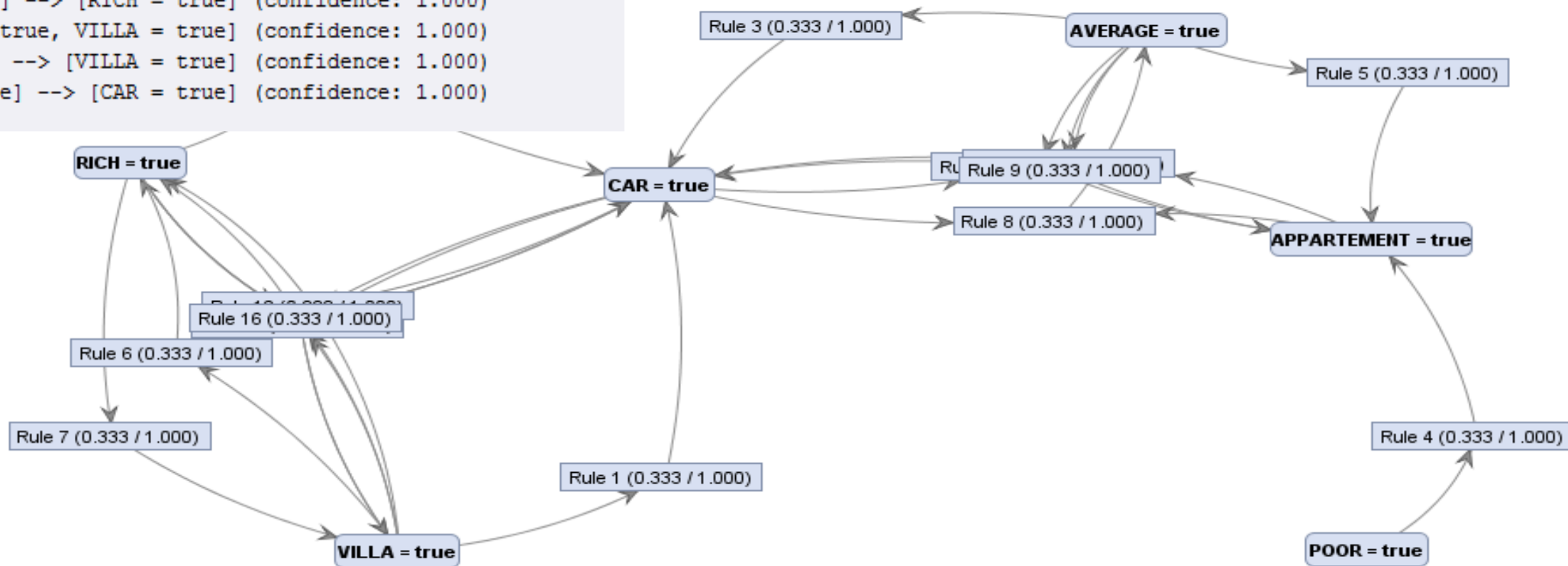
ExampleSet (3 examples, 0 special attributes, 6 regular attributes)						
Row No.	CAR = true	APPARTEMENT = true	VILLA = true	POOR = true	AVERAGE = true	RICH = true
1	false	true	false	true	false	false
2	true	true	false	false	true	false
3	true	false	true	false	false	true

# Association Rules

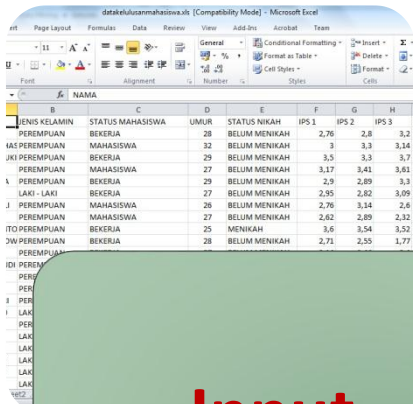
## Association Rules

```
[VILLA = true] --> [CAR = true] (confidence: 1.000)
[RICH = true] --> [CAR = true] (confidence: 1.000)
[AVERAGE = true] --> [CAR = true] (confidence: 1.000)
[POOR = true] --> [APPARTEMENT = true] (confidence: 1.000)
[AVERAGE = true] --> [APPARTEMENT = true] (confidence: 1.000)
[VILLA = true] --> [RICH = true] (confidence: 1.000)
[RICH = true] --> [VILLA = true] (confidence: 1.000)
[CAR = true, APPARTEMENT = true] --> [AVERAGE = true] (confidence: 1.000)
[AVERAGE = true] --> [CAR = true, APPARTEMENT = true] (confidence: 1.000)
[CAR = true, AVERAGE = true] --> [APPARTEMENT = true] (confidence: 1.000)
[APPARTEMENT = true, AVERAGE = true] --> [CAR = true] (confidence: 1.000)
[VILLA = true] --> [CAR = true, RICH = true] (confidence: 1.000)
[CAR = true, VILLA = true] --> [RICH = true] (confidence: 1.000)
[RICH = true] --> [CAR = true, VILLA = true] (confidence: 1.000)
[CAR = true, RICH = true] --> [VILLA = true] (confidence: 1.000)
[VILLA = true, RICH = true] --> [CAR = true] (confidence: 1.000)
```

# ASSOCIATION RULES



# PROSES UTAMA PADA DATA MINING



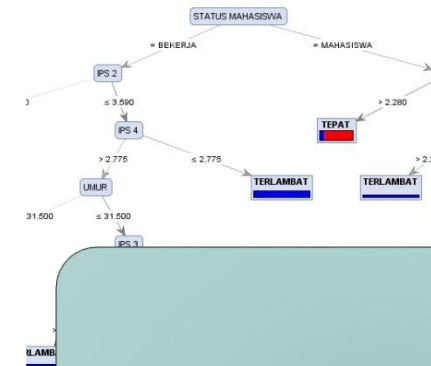
The screenshot shows a Microsoft Excel spreadsheet with the following data:

	B	C	D	E	F	G	H
	SEX	NAMA	STATUS MAHASISWA	UMUR	STATUS NIKAH	IPS 1	IPS 2
1	PEREMPUAN	BEKERJA	28	BELUM MENIKAH	2,76	2,8	3,2
2	PEREMPUAN	MAHASISWA	32	BELUM MENIKAH	3	3,3	3,14
3	PEREMPUAN	BEKERJA	29	BELUM MENIKAH	3,5	3,3	3,7
4	PEREMPUAN	MAHASISWA	27	BELUM MENIKAH	3,17	3,41	3,61
5	PEREMPUAN	BEKERJA	29	BELUM MENIKAH	2,9	2,89	3,3
6	LAKI - LAKI	BEKERJA	27	BELUM MENIKAH	2,95	2,82	3,09
7	PEREMPUAN	MAHASISWA	26	BELUM MENIKAH	2,76	3,14	2,6
8	PEREMPUAN	MAHASISWA	27	BELUM MENIKAH	2,62	2,89	2,32
9	PEREMPUAN	BEKERJA	25	MENIKAH	3,6	3,54	3,52
10	PEREMPUAN	BEKERJA	28	BELUM MENIKAH	2,71	2,55	1,77

**Input**  
(Data)

$$\int_a^b f(x) dx = \lim_{n \rightarrow \infty} \frac{b-a}{n} \sum_{k=1}^n f\left(a + \frac{b-a}{n} \cdot k\right)$$
$$= \left( -m_0 \int \tan(\phi) \right) \left[ l - \frac{r^2}{4l} + r \left( \cos(\omega t) + \frac{r}{4l} \cos(2\omega t) \right) \right]$$
$$= R_1 e^{\left( -\zeta + \sqrt{\zeta^2 - 1} \right) \omega t} \left( -\zeta - \sqrt{\zeta^2 - 1} \right) \omega t$$

**Metode**  
(Algoritma  
Data Mining)



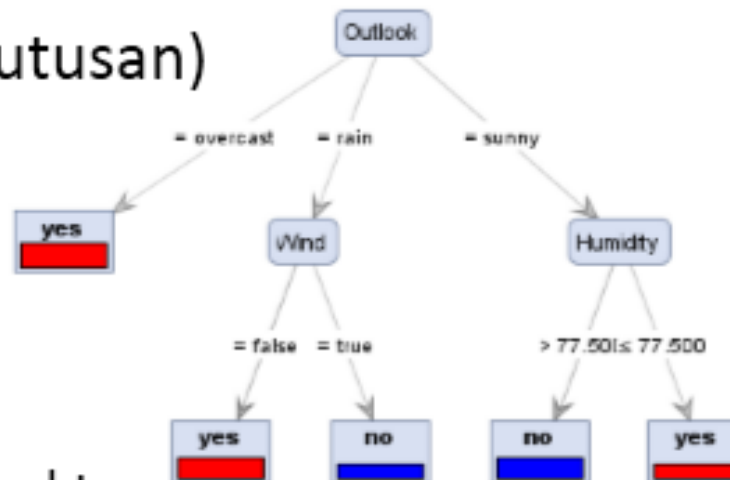
**Output**  
(Pola/Model)

# OUTPUT/POLA/MODEL/KNOWLEDGE

## 1. Formula/Function (Rumus atau Fungsi Regresi)

- $\text{WAKTU TEMPUH} = 0.48 + 0.6 \text{ JARAK} + 0.34 \text{ LAMPU} + 0.2 \text{ PESANAN}$

## 2. Decision Tree (Pohon Keputusan)

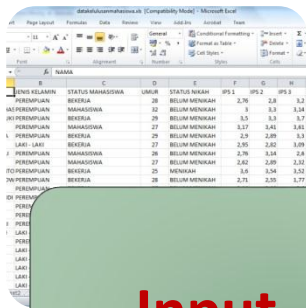


## 3. Rule (Aturan)

- IF  $\text{ips3}=2.8$  THEN lulustepatwaktu

## 4. Cluster (Klaster)

# INPUT – METODE – OUTPUT – EVALUATION

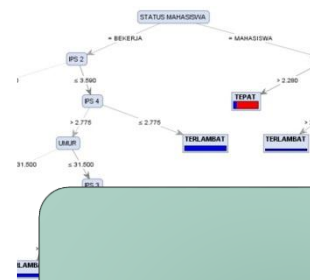


No	NAMA	C	D	E	F	G	H
1	BENI KILAMIN	STATUS MAHASISWA	UMUR	STATUS MAHASISWA	IPS 1	IPS 2	IPS 3
2	PEREMPUAN	BEKERJA	28	BEKUM MENENGAH	2.76	2.8	3.2
3	MAHASISWA	BEKUM MENENGAH	32	BEKUM MENENGAH	3	3.3	3.4
4	PEREMPUAN	BEKERJA	29	BEKUM MENENGAH	3.5	3.5	3.7
5	MAHASISWA	BEKERJA	27	BEKUM MENENGAH	3.17	3.41	3.81
6	PEREMPUAN	BEKERJA	29	BEKUM MENENGAH	3.3	3.49	3.5
7	LARI LARI	BEKERJA	27	BEKUM MENENGAH	3.35	3.62	3.89
8	PEREMPUAN	MAHASISWA	26	BEKUM MENENGAH	3.76	3.84	3.8
9	PEREMPUAN	MAHASISWA	27	BEKUM MENENGAH	3.82	3.89	3.82
10	PEREMPUAN	BEKERJA	25	BEKUM	3.6	3.54	3.52
11	PEREMPUAN	BEKERJA	28	BEKUM MENENGAH	3.75	3.75	3.77

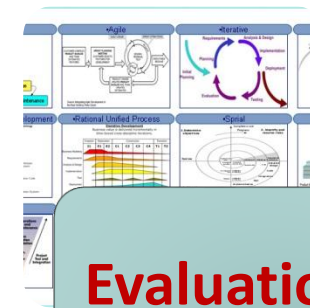
**Input**  
(Data)

$$f(x) dx = \lim_{n \rightarrow \infty} \frac{b-a}{n} \sum_{k=1}^n f\left(a + \frac{b-a}{n} \cdot k\right)$$
$$= (-m_p \sin(\phi)) \left[ t - \frac{r^2}{4t} + r \left( \cos(\omega t) + \frac{r}{4t} \cos(2\omega t) \right) \right]$$
$$+ R_1 e^{\left( -\zeta + \sqrt{\zeta^2 - 1} \right) \log t} \quad \left( -\zeta - \sqrt{\zeta^2 - 1} \right) \log t$$

**Metode**  
(Algoritma  
Data Mining)



**Output**  
(Pola/Model)



**Evaluation**  
(Akurasi, AUC,  
RMSE, etc)



# ALGORITMA DATA MINING

# ALGORITMA ESTIMASI

- Algoritma estimasi mirip dengan algoritma klasifikasi, tapi **variabel target adalah berupa bilangan numerik (kontinyu)** dan bukan kategorikal (nominal atau diskrit)
- Estimasi nilai dari variable target ditentukan **berdasarkan nilai dari variabel prediktor** (atribut)
- Algoritma estimasi yang biasa digunakan adalah: **Linear Regression, Neural Network, Support Vector Machine**

# CONTOH: ESTIMASI PERFORMANSI CPU

**Example:** 209 different computer configurations

	Cycle time (ns)	Main memory (Kb)		Cache (Kb)	Channels		Performance
	MYCT	MMIN	MMAX	CACH	CHMIN	CHMAX	PRP
1	125	256	6000	256	16	128	198
2	29	8000	32000	32	8	32	269
...							
208	480	512	8000	32	0	0	67
209	480	1000	4000	0	0	0	45

**Linear regression** function

$$\begin{aligned} \text{PRP} = & -55.9 + 0.0489 \text{ MYCT} + 0.0153 \text{ MMIN} + 0.0056 \text{ MMAX} \\ & + 0.6410 \text{ CACH} - 0.2700 \text{ CHMIN} + 1.480 \text{ CHMAX} \end{aligned}$$



# ALGORITMA PREDIKSI

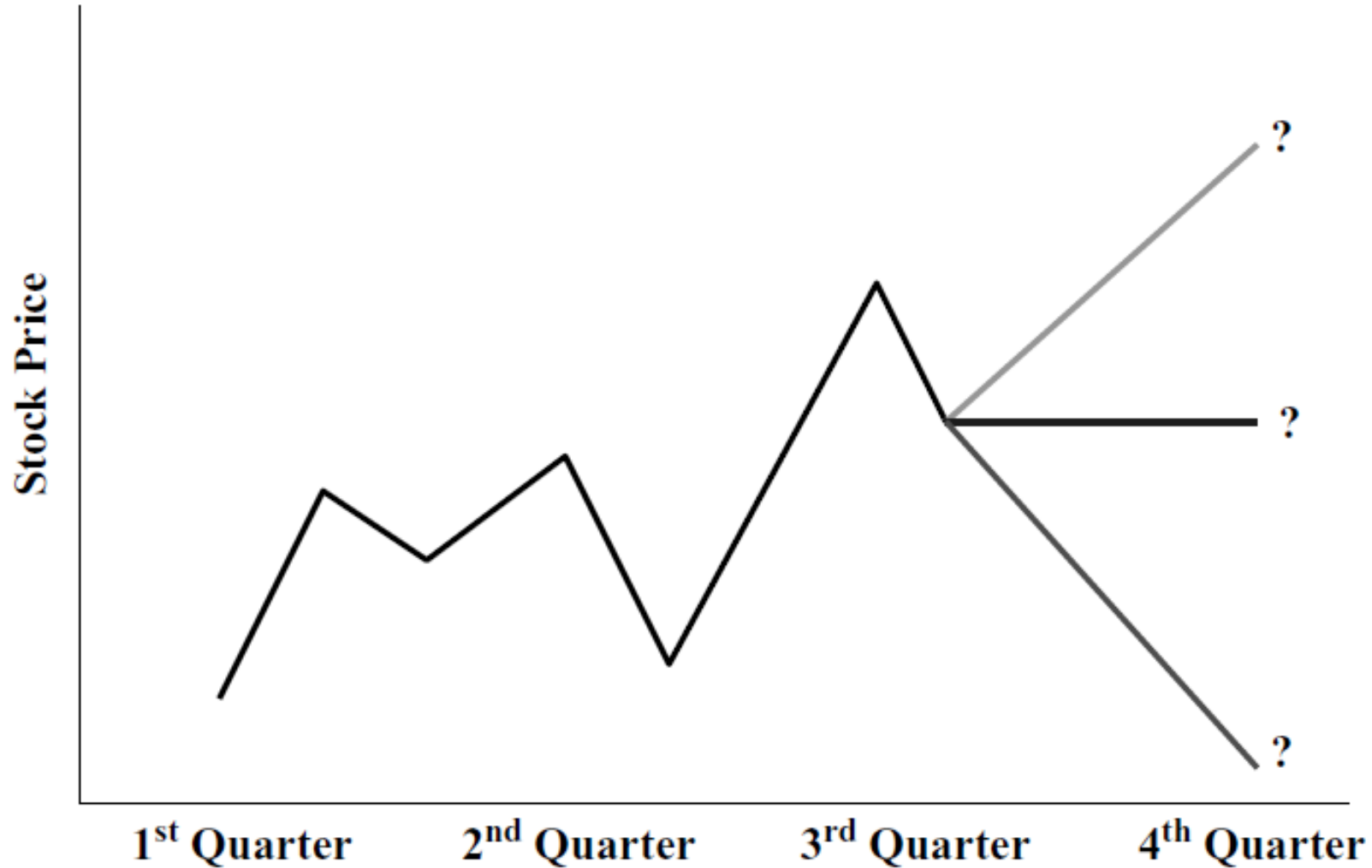
- Algoritma prediksi/forecasting sama dengan algoritma estimasi di mana label/target/class bertipe numerik, bedanya adalah data yang digunakan merupakan data rentet waktu (data time series)
- Istilah prediksi kadang digunakan juga untuk klasifikasi, tidak hanya untuk prediksi time series, karena sifatnya yang bisa menghasilkan class berdasarkan berbagai atribut yang kita sediakan
- Semua algoritma estimasi dapat digunakan untuk prediksi/forecasting

# CONTOH: PREDIKSI HARGA SAHAM

Row No.	Close	Date	Open	High	Low	Volume
1	1286.570	Apr 11, 2006	1296.600	1300.710	1282.960	2232880000
2	1288.120	Apr 12, 2006	1286.570	1290.930	1286.450	1938100000
3	1289.120	Apr 13, 2006	1288.120	1292.090	1283.370	1891940000
4	1285.330	Apr 17, 2006	1289.120	1292.450	1280.740	1794650000
5	1307.280	Apr 18, 2006	1285.330	1309.020	1285.330	2595440000
6	1309.930	Apr 19, 2006	1307.650	1310.390	1302.790	2447310000
7	1311.460	Apr 20, 2006	1309.930	1318.160	1306.380	2512920000
8	1311.280	Apr 21, 2006	1311.460	1317.670	1306.590	2392630000
9	1308.110	Apr 24, 2006	1311.280	1311.280	1303.790	2117330000
10	1301.740	Apr 25, 2006	1308.110	1310.790	1299.170	2366380000
11	1305.410	Apr 26, 2006	1301.740	1310.970	1301.740	2502690000
12	1309.720	Apr 27, 2006	1305.410	1315	1295.570	2772010000
13	1310.610	Apr 28, 2006	1309.720	1316.040	1306.160	2419920000
14	1305.190	May 1, 2006	1310.610	1317.210	1303.460	2437040000
15	1313.210	May 2, 2006	1305.190	1313.660	1305.190	2403470000
16	1308.120	May 3, 2006	1313.210	1313.470	1303.920	2395230000
17	1312.250	May 4, 2006	1307.850	1315.140	1307.850	2431450000
18	1325.760	May 5, 2006	1312.250	1326.530	1312.250	2294760000
19	1324.660	May 8, 2006	1325.760	1326.700	1322.870	2151300000
20	1325.140	May 9, 2006	1324.660	1326.600	1322.480	2157290000
21	1322.850	May 10, 2006	1324.570	1325.510	1317.440	2268550000
22	1305.920	May 11, 2006	1322.630	1322.630	1303.450	2531520000
23	1291.240	May 12, 2006	1305.880	1305.880	1290.380	2567970000
24	1294.500	May 15, 2006	1291.190	1294.810	1284.510	2505660000

Dataset harga saham dalam bentuk time series (rentet waktu) harian

# CONTOH: PREDIKSI HARGA SAHAM (PLOT)



# CONTOH: PREDIKSI HARGA SAHAM (PLOT)

Series ▼

Lower Bound None ▼

Upper Bound None ▼

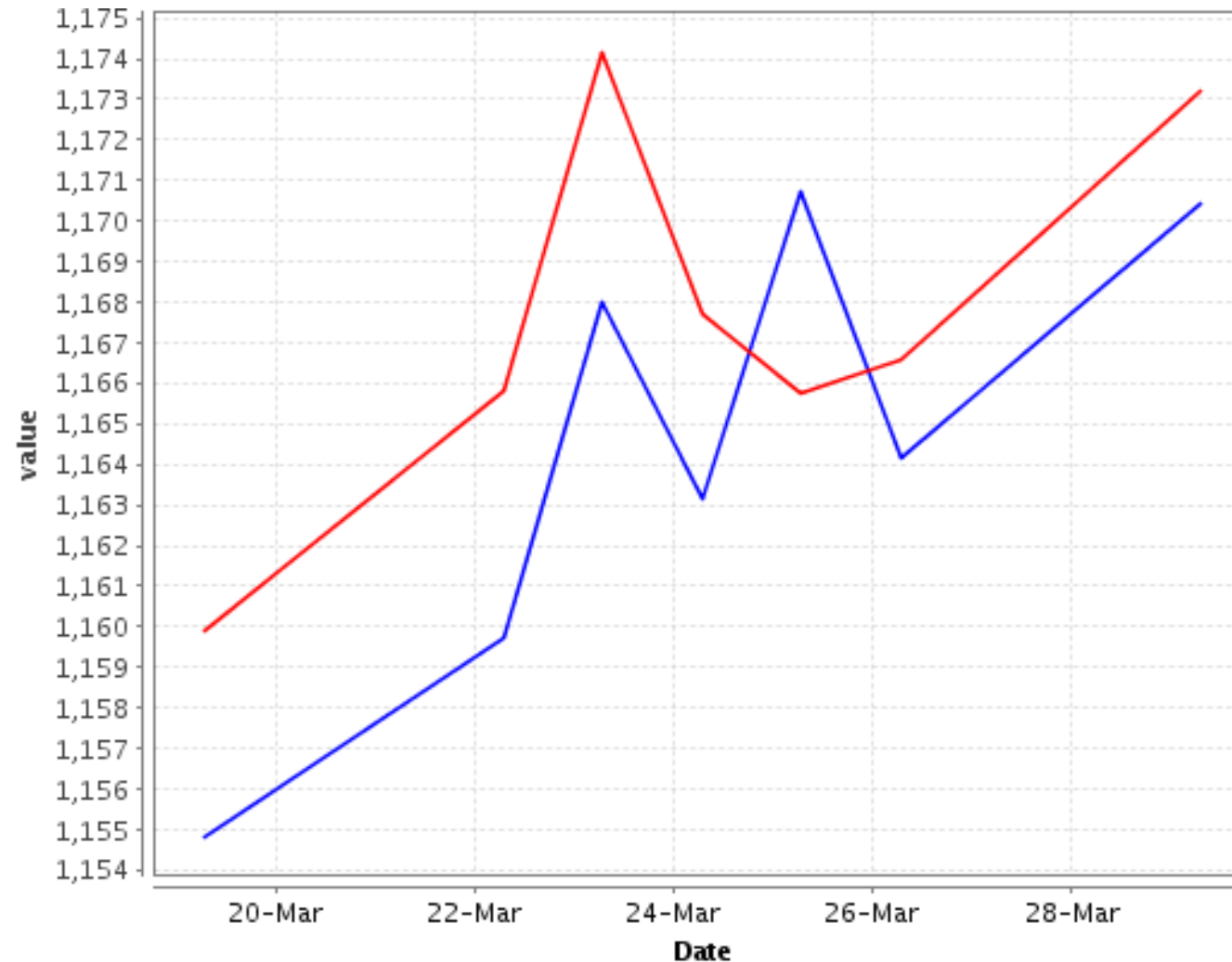
Index Dimension Date ▼

Plot Series

- Date
- Open
- High
- Low
- Volume
- Close
- prediction(Close)

☐ Rotate Labels

 Export Image...



# ALGORITMA KLASIFIKASI

- Klasifikasi adalah algoritma yang menggunakan data dengan **target/class/label berupa nilai kategorikal (nominal)**
- Contoh, apabila **target/class/label** adalah pendapatan, maka bisa digunakan nilai nominal (kategorikal) sbb: pendapatan besar, menengah, kecil
- Contoh lain adalah rekomendasi contact lens, apakah menggunakan yang jenis **soft, hard** atau **none**
- Algoritma klasifikasi yang biasa digunakan adalah: Naive Bayes, K-Nearest Neighbor, C4.5, ID3, CART, Linear Discriminant Analysis, etc

# CONTOH: REKOMENDASI MAIN GOLF

## Input:

Outlook	Temperature	Humidity	Windy	Play
Sunny	hot	high	false	no
Sunny	hot	high	true	no
Overcast	hot	high	false	yes
Rainy	mild	high	false	yes
Rainy	cool	normal	false	yes
Rainy	cool	normal	true	no
Overcast	cool	normal	true	yes
Sunny	mild	high	false	no
Sunny	cool	normal	false	yes
Rainy	mild	normal	false	yes
Sunny	mild	normal	true	yes
Overcast	mild	high	true	yes
Overcast	hot	normal	false	yes
Rainy	mild	high	true	no

## Output (Rules):

- If outlook = sunny and humidity = high then play = no
- If outlook = rainy and windy = true then play = no
- If outlook = overcast then play = yes
- If humidity = normal then play = yes
- If none of the above then play = yes

# CONTOH: REKOMENDASI MAIN GOLF

## Input (Atribut Nominal dan Numerik):

Outlook	Temperature	Humidity	Windy	Play
Sunny	85	85	false	no
Sunny	80	90	true	no
Overcast	83	86	false	yes
Rainy	70	96	false	yes
Rainy	68	80	false	yes
Rainy	65	70	true	no
Overcast	64	65	true	yes
Sunny	72	95	false	no
Sunny	69	70	false	yes
Rainy	75	80	false	yes
Sunny	75	70	true	yes
Overcast	72	90	true	yes
Overcast	81	75	false	yes
Rainy	71	91	true	no

## Output (Rules):

If outlook = sunny and humidity = high then play = no

If outlook = sunny and humidity > 83 then play = no

If outlook = rainy and windy = true then play = no

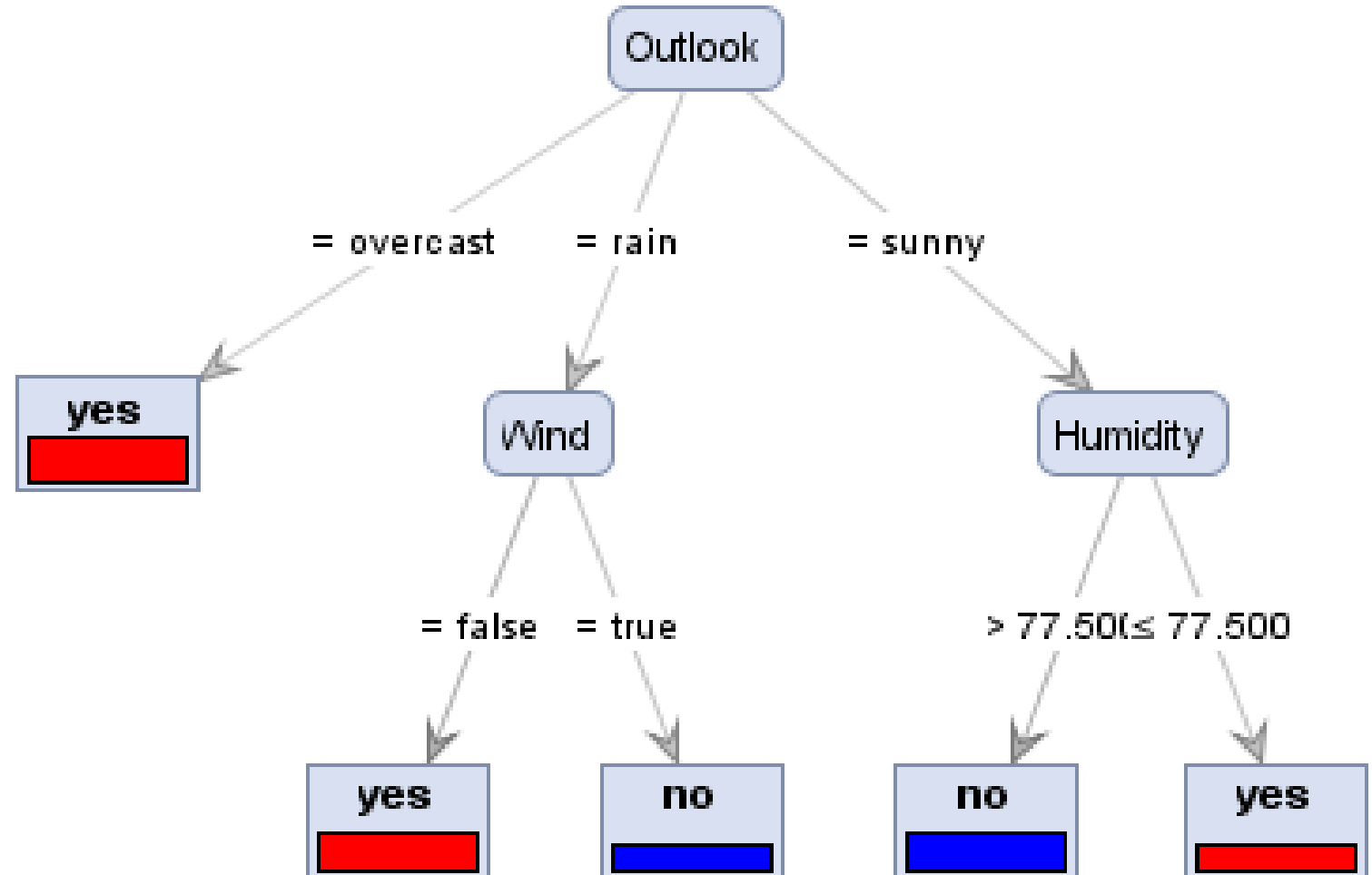
If outlook = overcast then play = yes

If humidity < 85 then play = yes

If none of the above then play = yes

# CONTOH: REKOMENDASI MAIN GOLF

**Output (Tree):**





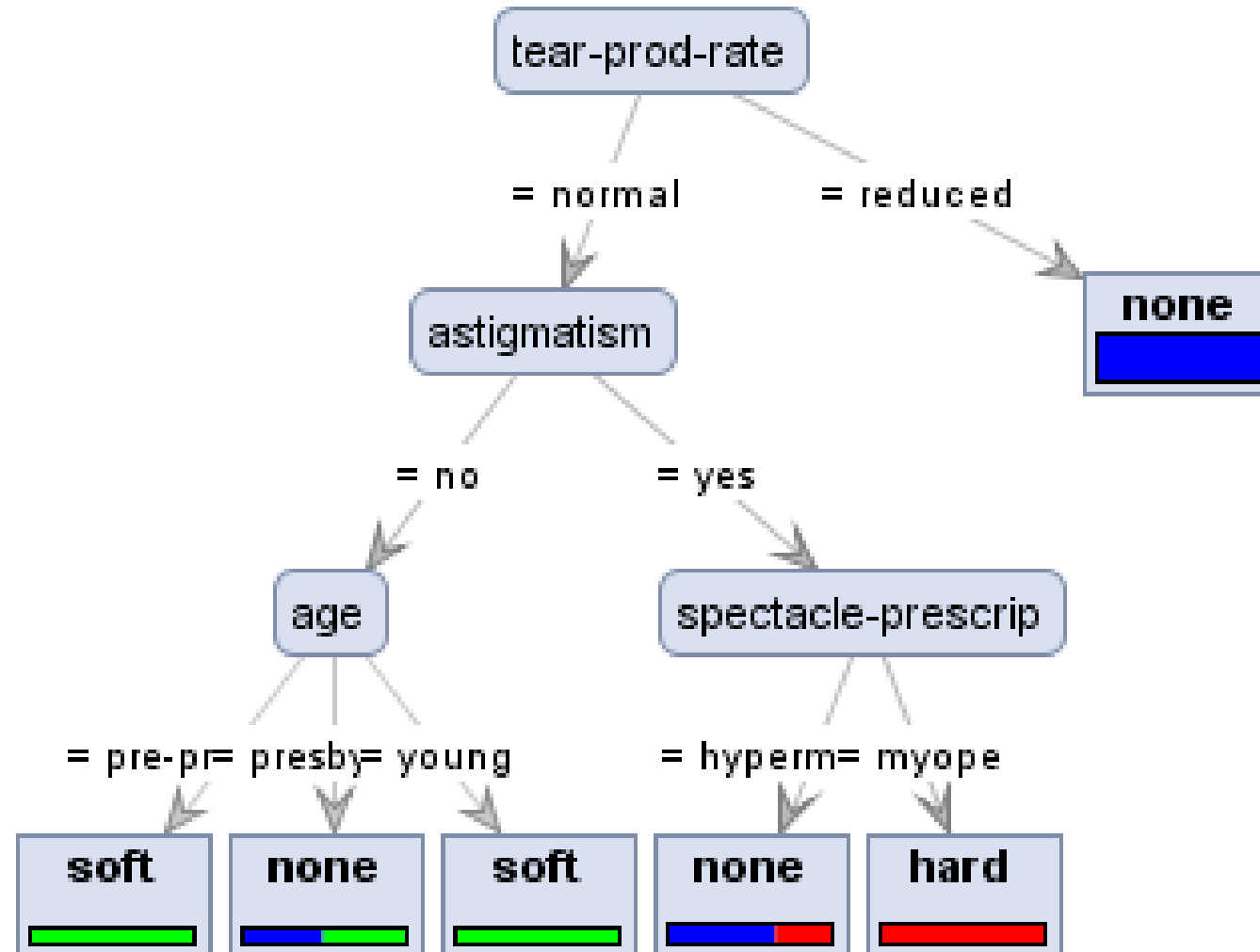
# CONTOH: REKOMENDASI CONTACT LENS

Input:

Age	Spectacle Prescription	Astigmatism	Tear Production Rate	Recommended Lenses
young	myope	no	reduced	none
young	myope	no	normal	soft
young	myope	yes	reduced	none
young	myope	yes	normal	hard
young	hypermetrope	no	reduced	none
young	hypermetrope	no	normal	soft
young	hypermetrope	yes	reduced	none
young	hypermetrope	yes	normal	hard
pre-presbyopic	myope	no	reduced	none
pre-presbyopic	myope	no	normal	soft
pre-presbyopic	myope	yes	reduced	none
pre-presbyopic	myope	yes	normal	hard
pre-presbyopic	hypermetrope	no	reduced	none
pre-presbyopic	hypermetrope	no	normal	soft

# CONTOH: REKOMENDASI CONTACT LENS

Output/Model (Tree):



# CONTOH: PENENTUAN JENIS BUNGA IRIS

**Input:**

	Sepal Length (cm)	Sepal Width (cm)	Petal Length (cm)	Petal Width (cm)	Type
1	5.1	3.5	1.4	0.2	<i>Iris setosa</i>
2	4.9	3.0	1.4	0.2	<i>Iris setosa</i>
3	4.7	3.2	1.3	0.2	<i>Iris setosa</i>
4	4.6	3.1	1.5	0.2	<i>Iris setosa</i>
5	5.0	3.6	1.4	0.2	<i>Iris setosa</i>
...					
51	7.0	3.2	4.7	1.4	<i>Iris versicolor</i>
52	6.4	3.2	4.5	1.5	<i>Iris versicolor</i>
53	6.9	3.1	4.9	1.5	<i>Iris versicolor</i>
54	5.5	2.3	4.0	1.3	<i>Iris versicolor</i>
55	6.5	2.8	4.6	1.5	<i>Iris versicolor</i>
...					
101	6.3	3.3	6.0	2.5	<i>Iris virginica</i>
102	5.8	2.7	5.1	1.9	<i>Iris virginica</i>
103	7.1	3.0	5.9	2.1	<i>Iris virginica</i>
104	6.3	2.9	5.6	1.8	<i>Iris virginica</i>
105	6.5	3.0	5.8	2.2	<i>Iris virginica</i>
...					

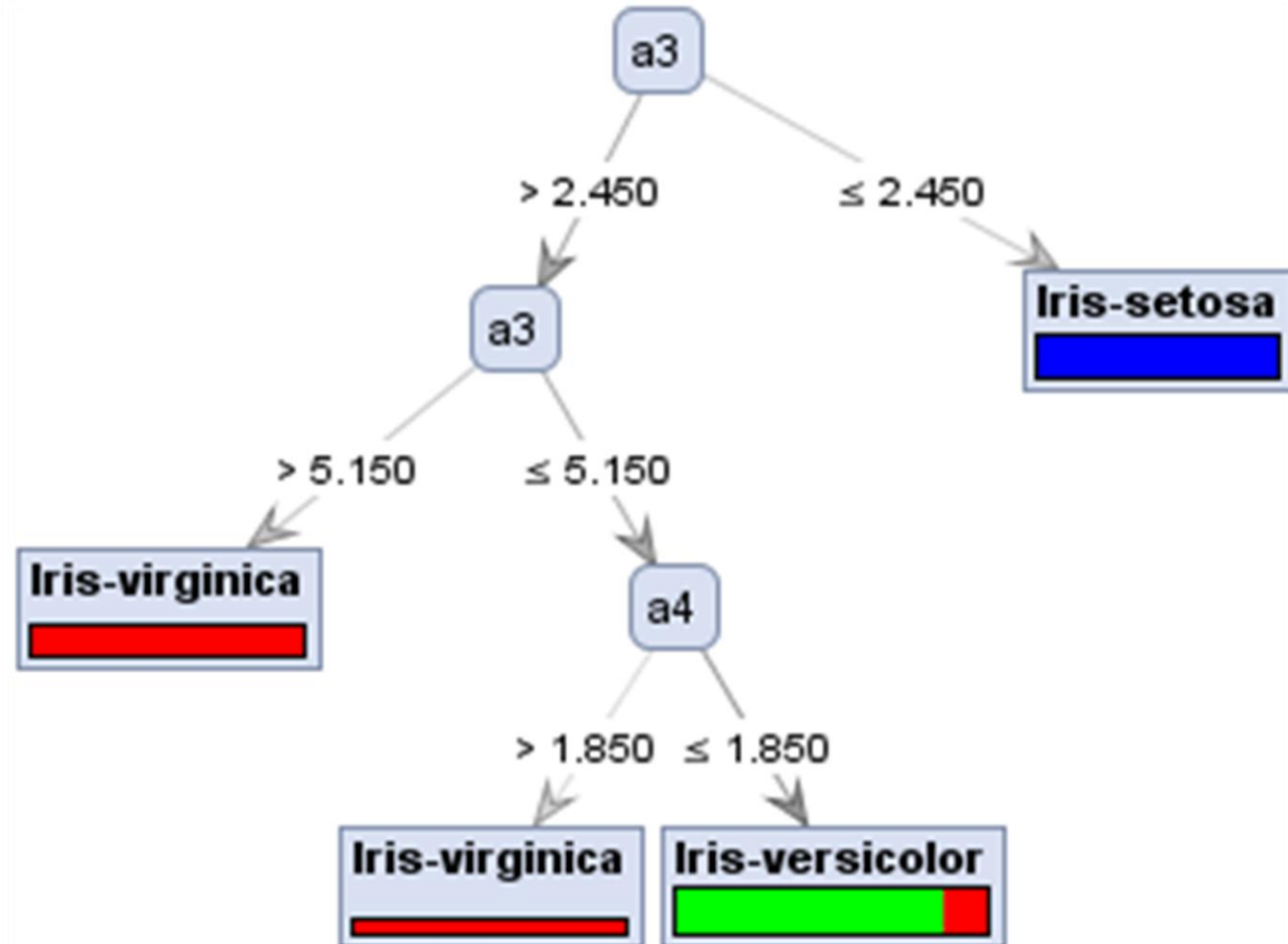
# CONTOH: PENENTUAN JENIS BUNGA IRIS

## Output (Rules):

```
If petal-length < 2.45 then Iris-setosa
If sepal-width < 2.10 then Iris-versicolor
If sepal-width < 2.45 and petal-length < 4.55 then Iris-versicolor
If sepal-width < 2.95 and petal-width < 1.35 then Iris-versicolor
If petal-length ≥ 2.45 and petal-length < 4.45 then Iris-versicolor
If sepal-length ≥ 5.85 and petal-length < 4.75 then Iris-versicolor
If sepal-width < 2.55 and petal-length < 4.95 and
    petal-width < 1.55 then Iris-versicolor
If petal-length ≥ 2.45 and petal-length < 4.95 and
    petal-width < 1.55 then Iris-versicolor
If sepal-length ≥ 6.55 and petal-length < 5.05 then Iris-versicolor
If sepal-width < 2.75 and petal-width < 1.65 and
    sepal-length < 6.05 then Iris-versicolor
If sepal-length ≥ 5.85 and sepal-length < 5.95 and
    petal-length < 4.85 then Iris-versicolor
If petal-length ≥ 5.15 then Iris-virginica
If petal-width ≥ 1.85 then Iris-virginica
If petal-width ≥ 1.75 and sepal-width < 3.05 then Iris-virginica
If petal-length ≥ 4.95 and petal-width < 1.55 then Iris-virginica
```

# CONTOH: PENENTUAN JENIS BUNGA IRIS

**Output** (Tree):



# ALGORITMA KLASSTERING

- Klastering adalah **pengelompokkan data**, hasil observasi dan kasus ke dalam **class yang mirip**
- Suatu klaster (cluster) adalah **koleksi data yang mirip** antara satu dengan yang lain, dan **memiliki perbedaan** bila dibandingkan dengan data dari klaster lain
- Perbedaan utama algoritma klastering dengan klasifikasi adalah **klastering tidak memiliki target/class/label**, jadi termasuk *unsupervised learning*
- Klastering sering digunakan sebagai **tahap awal dalam proses data mining**, dengan hasil klaster yang terbentuk akan menjadi input dari algoritma berikutnya yang digunakan

# CONTOH: KLAUSTERING JENIS GAYA HIDUP

- Claritas, Inc. provide a **demographic profile of each of the geographic areas in the country**, as defined by zip code. One of the clustering mechanisms they use is the PRIZM segmentation system, which **describes every U.S. zip code area in terms of distinct lifestyle types (66 segments)**. Just go to the company's Web site, enter a particular zip code, and you are shown the most common PRIZM clusters for that zip code.
- What do these clusters mean? For illustration, let's look up the clusters for zip code 90210, Beverly Hills, California. The **resulting clusters for zip code 90210** are:
  - Cluster 01: *Blue Blood Estates*
  - Cluster 10: *Bohemian Mix*
  - Cluster 02: *Winner's Circle*
  - Cluster 07: *Money and Brains*
  - Cluster 08: *Young Literati*

What is Nielsen  
PRIZM?

Features and Benefits

Lifestyle Segmentation

Urbanization Classes

Social Groups

Lifestage Classes

Lifestage Groups

Summary

## Features and Benefits

and market to them with tailored messages and products designed just for them. Captured by catchy names, images and behavior snapshots that bring the segments to life for marketers, PRIZM segments are memorable and summarize complex consumer profiles in a way that is intuitive and easy to communicate.

For example, PRIZM Segment number 16 is known as *Bohemian Mix*. We can describe both the demographic traits, as well as the lifestyle characteristics of the households in this segment. You can review these segment descriptors in the image at right.

## Bohemian Mix

16



### Y2 Young Achievers

Upper-Mid Middle Age Family Mix

<55

Renters

White-Collar, Mix

College Graduate

White, Black, Asian, Hispanic

Eat at Au Bon Pain

Buy Spanish/Latin music

Read *The Economist*

Watch soccer

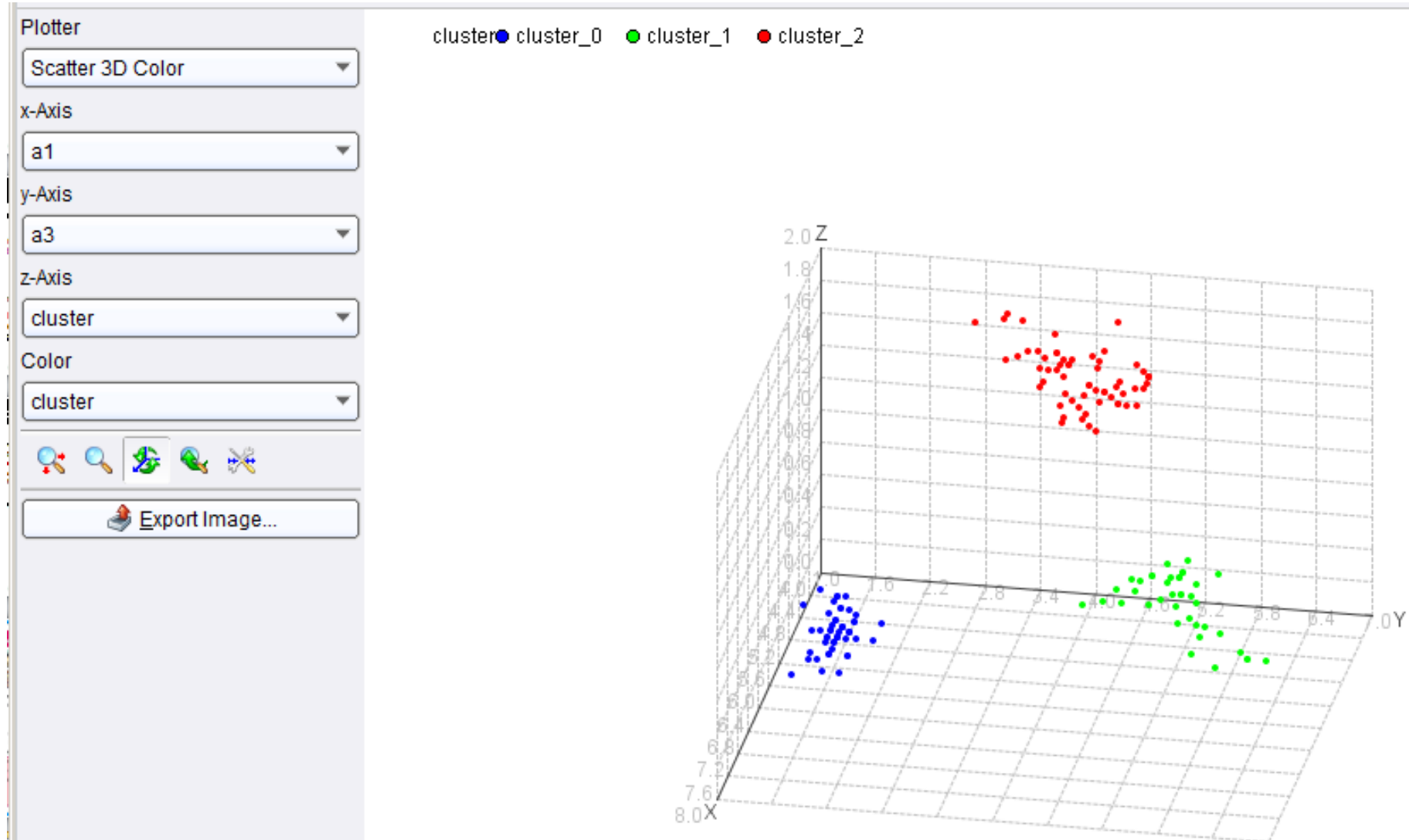
Audi A4



# CONTOH: KLASSTERING BUNGA IRIS

ExampleSet (150 examples, 2 special attributes, 4 regular attributes)						
Row No.	id	label	a1	a2	a3	a4
1	id_1	Iris-setosa	5.100	3.500	1.400	0.200
2	id_2	Iris-setosa	4.900	3	1.400	0.200
3	id_3	Iris-setosa	4.700	3.200	1.300	0.200
4	id_4	Iris-setosa	4.600	3.100	1.500	0.200
5	id_5	Iris-setosa	5	3.600	1.400	0.200
6	id_6	Iris-setosa	5.400	3.900	1.700	0.400
7	id_7	Iris-setosa	4.600	3.400	1.400	0.300
8	id_8	Iris-setosa	5	3.400	1.500	0.200
9	id_9	Iris-setosa	4.400	2.900	1.400	0.200
10	id_10	Iris-setosa	4.900	3.100	1.500	0.100
11	id_11	Iris-setosa	5.400	3.700	1.500	0.200
12	id_12	Iris-setosa	4.800	3.400	1.600	0.200
13	id_13	Iris-setosa	4.800	3	1.400	0.100
14	id_14	Iris-setosa	4.300	3	1.100	0.100
15	id_15	Iris-setosa	5.800	4	1.200	0.200
16	id_16	Iris-setosa	5.700	4.400	1.500	0.400
17	id_17	Iris-setosa	5.400	3.900	1.300	0.400
18	id_18	Iris-setosa	5.100	3.500	1.400	0.300
19	id_19	Iris-setosa	5.700	3.800	1.700	0.300
20	id_20	Iris-setosa	5.100	3.800	1.500	0.300
21	id_21	Iris-setosa	5.400	3.400	1.700	0.200
22	id_22	Iris-setosa	5.100	3.700	1.500	0.400
23	id_23	Iris-setosa	4.600	3.600	1	0.200
24	id_24	Iris-setosa	5.100	3.300	1.700	0.500

# CONTOH: KLASTERING BUNGA IRIS (PLOT)



# CONTOH: KLASSTERING BUNGA IRIS (TABLE)

ExampleSet (150 examples, 3 special attributes, 4 regular attributes)							View
Row No.	id	label	cluster	a1	a2	a3	a4
1	id_1	Iris-setosa	cluster_0	5.100	3.500	1.400	0.200
2	id_2	Iris-setosa	cluster_0	4.900	3	1.400	0.200
3	id_3	Iris-setosa	cluster_0	4.700	3.200	1.300	0.200
4	id_4	Iris-setosa	cluster_0	4.600	3.100	1.500	0.200
5	id_5	Iris-setosa	cluster_0	5	3.600	1.400	0.200
6	id_6	Iris-setosa	cluster_0	5.400	3.900	1.700	0.400
7	id_7	Iris-setosa	cluster_0	4.600	3.400	1.400	0.300
8	id_8	Iris-setosa	cluster_0	5	3.400	1.500	0.200
9	id_9	Iris-setosa	cluster_0	4.400	2.900	1.400	0.200
10	id_10	Iris-setosa	cluster_0	4.900	3.100	1.500	0.100
11	id_11	Iris-setosa	cluster_0	5.400	3.700	1.500	0.200
12	id_12	Iris-setosa	cluster_0	4.800	3.400	1.600	0.200
13	id_13	Iris-setosa	cluster_0	4.800	3	1.400	0.100
14	id_14	Iris-setosa	cluster_0	4.300	3	1.100	0.100
15	id_15	Iris-setosa	cluster_0	5.800	4	1.200	0.200
16	id_16	Iris-setosa	cluster_0	5.700	4.400	1.500	0.400
17	id_17	Iris-setosa	cluster_0	5.400	3.900	1.300	0.400
18	id_18	Iris-setosa	cluster_0	5.100	3.500	1.400	0.300
19	id_19	Iris-setosa	cluster_0	5.700	3.800	1.700	0.300
20	id_20	Iris-setosa	cluster_0	5.100	3.800	1.500	0.300
21	id_21	Iris-setosa	cluster_0	5.400	3.400	1.700	0.200
22	id_22	Iris-setosa	cluster_0	5.100	3.700	1.500	0.400
23	id_23	Iris-setosa	cluster_0	4.600	3.600	1	0.200
24	id_24	Iris-setosa	cluster_0	5.100	3.300	1.700	0.500

## Cluster Model

Cluster 0: 50 items  
Cluster 1: 39 items  
Cluster 2: 61 items  
Total number of items: 150

# ALGORITMA ASOSIASI

- Algoritma *association rule* (aturan asosiasi) adalah algoritma yang menemukan atribut yang “**muncul bersamaan**”
- Dalam dunia bisnis, sering disebut dengan *affinity analysis* atau *market basket analysis*
- Algoritma asosiasi akan mencari aturan yang **menghitung hubungan diantara dua atau lebih atribut**
- Algoritma association rules berangkat dari pola “**If antecedent, then consequent,**” bersamaan dengan pengukuran **support** (coverage) dan **confidence** (accuracy) yang terasosiasi dalam aturan

# ALGORITMA ASOSIASI

- Contoh, pada hari Kamis malam, 1 000 pelanggan telah melakukan belanja di supermarket ABC, dimana:
  - 200 orang membeli Sabun Mandi
  - dari 200 orang yang membeli sabun mandi, 50 orangnya membeli Fanta
- Jadi, association rule menjadi, “Jika membeli sabun mandi, maka membeli Fanta”, dengan nilai support =  $200/1000 = 20\%$  dan nilai confidence =  $50/200 = 25\%$
- Algoritma association rule diantaranya adalah: A priori algorithm, FP-Growth algorithm, GRI algorithm

# CONTOH PENERAPAN DATA MINING

- Penentuan kelayakan aplikasi peminjaman uang di bank
- Penentuan pasokan listrik PLN untuk wilayah Jakarta
- Diagnosis pola kesalahan mesin
- Perkiraan harga saham dan tingkat inflasi
- Analisis pola belanja pelanggan
- Memisahkan minyak mentah dan gas alam
- Pemilihan program TV otomatis
- Penentuan pola pelanggan yang loyal pada perusahaan operator telepon
- Deteksi pencucian uang dari transaksi perbankan
- Deteksi serangan (intrusion) pada suatu jaringan

# COGNITIVE-PERFORMANCE TEST

1. Sebutkan **5 peran utama** data mining!
2. **algoritma apa** saja yang dapat digunakan untuk 5 peran utama data mining di atas?
3. Jelaskan perbedaan **estimasi** dan **prediksi**!
4. Jelaskan perbedaan **estimasi** dan **klasifikasi**!
5. Jelaskan perbedaan **klasifikasi** dan **klustering**!
6. Jelaskan perbedaan **klustering** dan **prediksi**!
7. Jelaskan perbedaan **supervised** dan **unsupervised** learning!
8. Sebutkan **tahapan utama proses** data mining!

# REFERENSI

1. Ian H. Witten, Frank Eibe, Mark A. Hall, *Data mining: Practical Machine Learning Tools and Techniques 3rd Edition*, Elsevier, 2011
2. Daniel T. Larose, *Discovering Knowledge in Data: an Introduction to Data Mining*, John Wiley & Sons, 2005
3. Florin Gorunescu, *Data Mining: Concepts, Models and Techniques*, Springer, 2011
4. Jiawei Han and Micheline Kamber, *Data Mining: Concepts and Techniques Second Edition*, Elsevier, 2006
5. Oded Maimon and Lior Rokach, *Data Mining and Knowledge Discovery Handbook Second Edition*, Springer, 2010
6. Warren Liao and Evangelos Triantaphyllou (eds.), *Recent Advances in Data Mining of Enterprise Data: Algorithms and Applications*, World Scientific, 2007
7. Santosa Budi, *Teknik Pemanfaatan Data Untuk Keperluan Bisnis*, Graha Ilmu, 2007