

K¹ 05/02/2018

enroll: damin1718

• Kuis: P1 & P10

• Ujian: Opensheet

• UTS 30%

UAS 25%

Praktikum 25%

Tugas Akhir 20%

↳ 5% laporan progress

5% laporan akhir

10% presentasi final

- ↳ a. Clustering & deteksi anomali
- b. Teknik klasifikasi
- c. Asosiasi rule mining
- d. Teknik lanjutan: tekt
kreb
spatio-temporal data
sequential pattern

Data: min 1000 sampel
10 attribut

Pengertian Data Mining

Data Mining (knowledge discovery in database):

ekstraksi informasi/pola yg menaik (non-trivial, implicit, previously unknown dan potentially useful) dlm basis data berukuran besar.

↳ yg bukan termasuk data mining bkt:

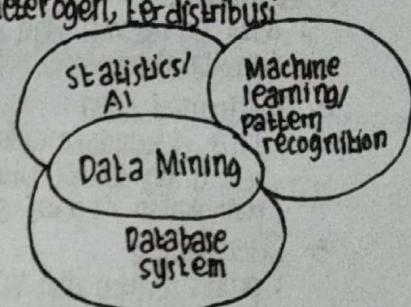
1. Pemrosesan (deductive) query.
2. Sistem pakar
3. Sistem informasi

Hubungan data mining dengan bidang lain

↳ berkaitan erat dengan bidang machine learning/AI, pattern recognition, statistika, dan sistem database.

↳ Teknik tradisional mjd bdk sesuai karena:

- data besar
- tingginya dimensi data
- data heterogen, terdistribusi



Pengantar

DATA MINING

Masalah eksplorasi data

→ menumpuknya data dalam database, data warehouse, dan penyimpanan lainnya.

→ Solusi: Data Warehousing & Data Mining

(& Online Analytical Processing (OLAP))

ekstraksi pengetahuan caturan, pola / kendala dari baris data berukuran besar.

sistem penyimpanan terpusat dari 1/lebih sumber yang berbeda untuk reporting dan analisis data

Mengapa Membanding Data?

1. Sudut pandang komersil

- telah blks data yg dikumpulkan
- teknologi komputer mjd lbh murah dan powerful
- tekanan kompetisi semakin kuat
 - : menyediakan layanan yg lbh baik, eg: customer relationship management

- web data, e-commerce, e-banking
- grocery stores
- bank/credit card transactions

2. Sudut pandang Keilmuan

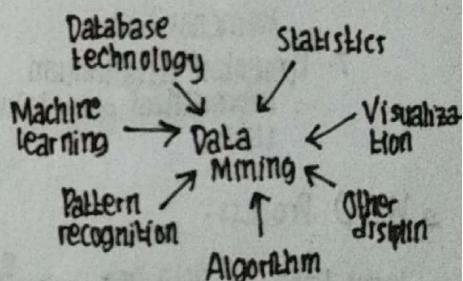
- data dikumpulkan dan disimpan dg kecepatan tinggi (GB/hour)
- teknik tradisional tdk cukup utk analisis data
- Data Mining membantu ilmuwan dlm: klarifikasi & segmentasi data, pemodelan, clustering
- remote sensor pada satelite
- microarray pembangkit gene expression data
- data spasial dalam GIS
- simulasi keilmuan

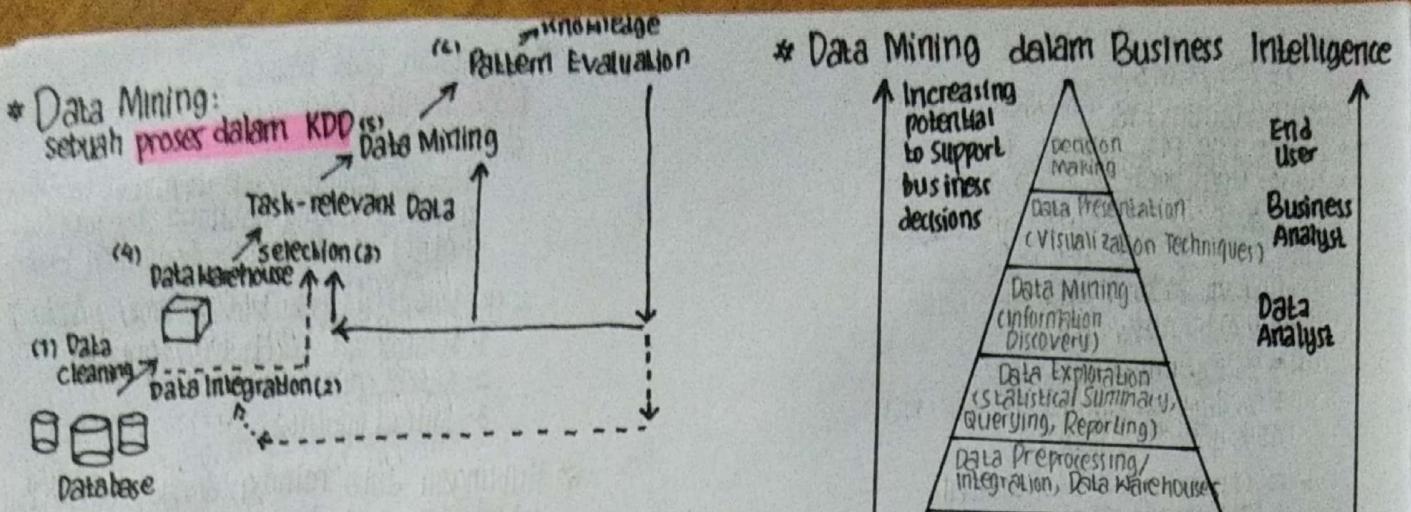
Evolution of Sciences: New Data Science Era

1990 - now: Data Science

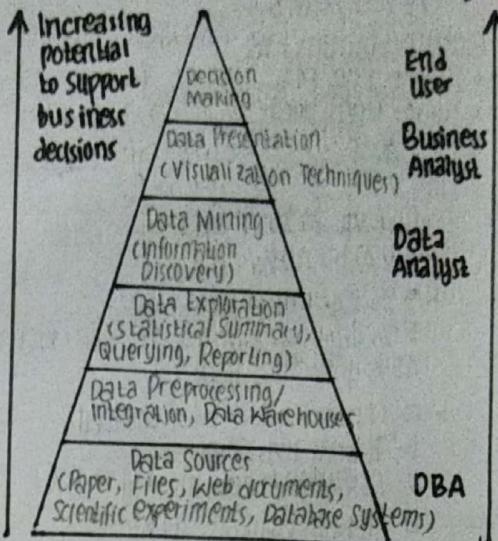
- data melimpah
- secara ekonomis: dapat disimpan dan di manage
- diakses banyak orang mel-internet
- meningkat secara linear

Data Mining:
Multi disiplin



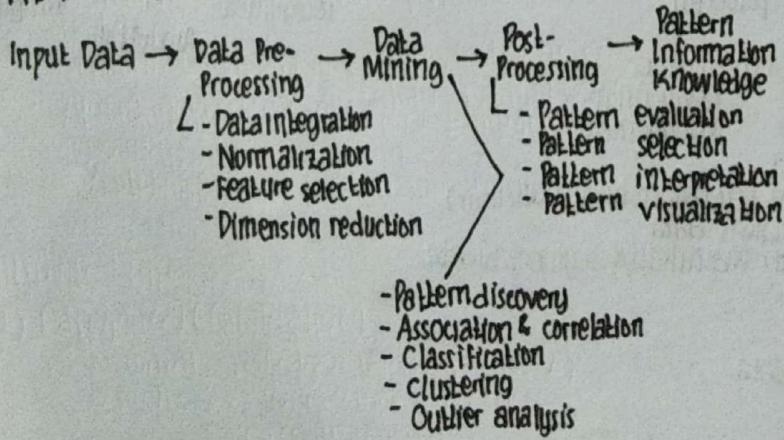


* Data Mining dalam Business Intelligence

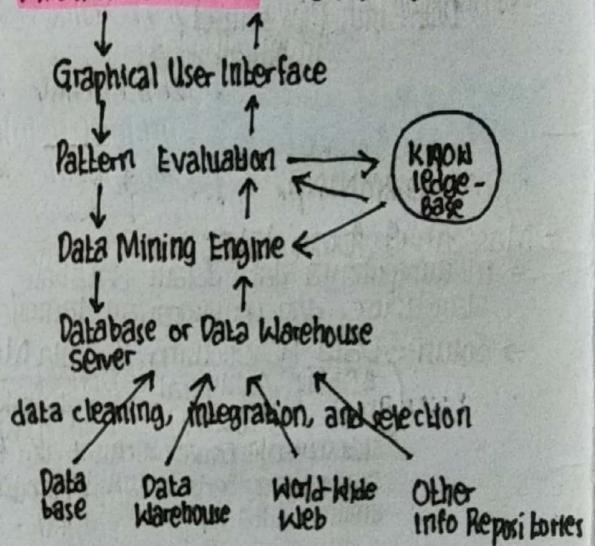


1. Pembersihan data
: menghilangkan noise dan data yg tdk konsisten.
2. Pengintegrasian data
: data digabungkan dari berbagai sumber
3. Seleksi data
: data yg relevan dg proses analisis diambil dari database.
4. Transformasi data
: data d transformasikan / digabungkan ke dm bentuk yg sesuai utk di-mine dg cara dilakukan peringkasan / operasi agregasi.
5. Data mining
: proses penting dm KDD dimana metode²x cerdas d aplikasikan utk mengekstrak pola²x data.
6. Evaluasi pola
: mengidentifikasi pola²x yg menarik yg merepresentasikan pengetahuan berdasarkan suatu ukuran kemenarikan.
7. Presentasi pengetahuan
: representasi pengetahuan yg telah digali kepada user.

* KDD Process:



* Arsitektur Sistem Data Mining



1. Basis data, data warehouse, / Tempat penyimpanan informasi lainnya
2. Basis data dan data warehouse server
: berfungsi jawab dr pengambilan data yg relevan berdasarkan permintaan pengguna.
3. Basis pengetahuan
: domain knowledge utk memandu pencarian/ evaluasi pola²x yg diharikan.
4. Data mining engine
: berdiri dari modul²x fungsional data mining, seperti karakterisasi, asosiasi, klasifikasi, dan analisis cluster.
5. Modul evaluasi pola
: menggunakan ukuran²x kemenarikan dan berinteraksi dg modul data mining dm pencarian pola²x menarik.
6. Antarmuka pengguna grafis
: media komunikasi dg pengguna dan sistem data mining.

* Data yg bisa ditambah:

- basis data relational
 - Data warehouse
 - Basis data transaksional
(eg: kasir merekam transaksi)
 - Tempat penyimpanan data lainnya:
 1. Basis data object-oriented dan basis data object-relational
 2. Basis data spasial
 3. Data Time-series dan data temporal
 4. Sequence data (incl. bio-sequences)
 5. Data teks dan basis data multimedia
 6. WWW

Tugas dalam Data Mining

1. Metode Prediksi:
 - : mengg. beberapa variabel (atribut) utk memprediksi nilai yg bdk diketahui/nilai yg akan datang dr variabel (atribut) lain.
 2. Metode Deskripsi:
 - : menemukan pola^{2x} (korelasi, trend, cluster, trayektori, dan anomali) yg meringkas hubungan dalam data.

Contoh:

1. Association (Correlation & Causality)
 - multi-dimensional vs. single-dimensional association
 - bagian sebab:
 - ada beberapa operasi
 - / prasyarat utk melakukan aksi
 - / gabungan x dan y
 - 1 kondisi yg jd prasyarat utk melakukan seratus (buy susu → buy roti)
 - snack → soft drink [0.5%, 80%] (support, confidence)
 - peluang bersyarat (peluang ini drukti ini)
 2. Klasifikasi dan Prediksi
 - : menemukan model (fungsi) yg menjelaskan dan membedakan kelas / konsep utk prediksi mendatang.
contoh: klasifikasi negara berdasarkan iklim
 - Presentasi: decision tree, classification rule, neural network
 - Prediksi: prediksi nilai numerik yg tlk diketahui / yg hilang

Pola yang menaik

Ukuran kemenangan:

- l mudah dimengerti oleh pengguna
 - l valid pd data baru / data tes dg
derajat kepastian (certainty)
 - l berguna
 - l novel (basisnya riset), atau
 - l memvalidasi hipotesis yg dicari
oleh pengguna

Pendekatan: Human-centered, query-based, focused mining

Penerapan

- Analisis halaman web
 - membangun sistem rekomendasi dan analisis kolaborasi
 - analisis basket data utk target marketing
 - analisis brologi dan medical data
 - data mining dan software engineering
 - from major dedicated data mining systems / tools to invisible data mining

Quiz

- Dividing the employee of a company according to their gender. (Bukan)
- Predicting future stock price using historical records. (Data Mining)
- Monitoring the heart rate of abnormalities. (Data Mining)
- Computing the total sales of a company. (Bukan)
- Characterize loyal buyers for a product. (Data Mining)

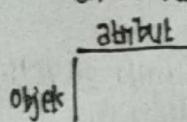
data & Eksplorasi data

Data dan Tipe Data

Data: kumpulan dari objek dan atribut

Atribut: hal yg melekat pd objek yg diamati/karakteristik dari objek (properti)

Contoh: warna mata



Atribut Values: bilangan/simbol yg diberikan pada suatu atribut.

Contoh: huruf mulu (atribut)
A, AB, B, ... (atribut values)

Perbedaan: Atribut & Atribut Value

Atribut yg sama dipetakan pada atribut value yg berbeda.

Contoh: tinggi dpt diukur dalam feet / meter.

Atribut yg berbeda dipetakan pd atribut value yg sama.

Contoh: atribut value ID dan umur = integer
- ID: tdk ada batasan
- Umur: ada batasan min dan max

Atribut: Type

1) Categorical (Qualitative)

a. Nominal: perbedaan nama saja (=, ≠)

Contoh: sex f male, female, eye color, kode pos, ID pegawai

Operasi: modus, entropy, contingency, correlation, χ^2 test

b. Ordinal: informasi dapat diurutkan. (<, >)

Contoh: street number, grades, tingkat kekerasan material (good, better, best)

Operasi: median, persentil, rank correlation, run tests, sign tests.

2.) Numerical (Quantitative)

a. Interval: perbedaan value memiliki makna. (+, -)

Contoh: calendar dates, temperature in Celsius / Fahrenheit

b. Ratio: perbedaan value dan rasio dari attribute value memiliki makna. (*, /)

Contoh: temperature in Kelvin, monetary quantities, counts, age, mass, panjang, arus listrik.

Operasi: geometric mean, harmonic mean, percent variation.

Attribute: Discrete & Continuous

Data

1. Atribut Diskret

- himpunan nilai berbatas
contoh: kode pos, counts, sekumpulan kata pd koleksi dokumen

- Variabel Integer

2. Atribut Kontinu

- atribut value: real number
contoh: temperatur, tinggi, besar
- variabel float (floating-point data)

Tipe Dataset

1. Record: kumpulan record, setiap record terdiri

himpunan atribut

- Data Matrix

- Data Dokumen

- Data Transaksi

: setiap record/transaksi membatalkan sekumpulan items.

Contoh: TRANSAKSI

TID	/item
1	Bread, Tea, Milk
2	Coffee, Bread
3	Coffee, Tea, Sugar, Milk

2. Graph: ada edges, ada node

- World Wide Web, Twitter, dll

- Struktur Molekul, interaksi \Rightarrow Protein

3. Ordered (bisa diurutkan)

- Spasial Data: koordinat longitude, latitude

- Temporal Data: pengukuran yg dilakukan dari waktu ke waktu

Time Series: sequen/tren pengukuran dari beberapa atribut / data dikumpulkan dalam periode waktu tertentu.

- Sequential Data: set transaksi dan item, nomor atribut waktu dan ID pembeli diasosiasi kan thd masing \Rightarrow transaksi

- Genetic Sequence Data

Kualitas Data

Masalah kualitas data:

- noise dan outliers
- missing values
- duplicate data

1. Noise

: berubahnya nilai atribut yg asli

Contoh: distorsi suara manusia ketika berbicara ditengah keramaian / telepon yg jelek / suara berasik dari televisi.

2. Outliers (Penculan)

: objek data yg memiliki karakteristik yg sangat berbeda dgn objek data pd umumnya.

3. Missing Values

Sebab: - informasi tsb tdk terkumpul, sng NA.

Contoh: orang menjalankan (nilai kosong) memberikan info umur dan berat badan.

- Atribut tdk dpt diterapkan pd semua kasus.
eg: atribut pendapatan tahunan tdk bisa dijawab / tdk cocok dgn anak^{2x}

Solusi: - eliminasi objek data

- estimasi missing values, melakukan interpolasi Mengganti nilai^{2x} yg ada pd atribut tsb.

- mengabaikan missing values selama tahap analisis.

4. Duplicate Data

Data cleaning: menghilangkan noise dan data yg tdk konsisten.

Statistika Ringkasan

: summarize property data, meliputi frequency, location, spread.

contoh: lokasi - mean, median

spread - standard deviation

1. Frequency and Mode

- Frequency: persentase kemunculan nilai tsb didalam dataset pada suatu waktu.

- Modus: nilai atribut yg memiliki frekuensi paling tinggi.

2. Ukuran Lokasi: Mean dan Median

sgt sensitif thd outlier.

$$\text{mean } (\bar{x}) = \bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$$

$$\text{median } (\bar{x}) = \begin{cases} x_{(r+1)} & ; m = \text{odd} : 2r+1 \\ \frac{1}{2}(x_{(r)} + x_{(r+1)}) & ; m = \text{even}: 2r \end{cases}$$

3. Ukuran Penyebaran: Range and Variance

- Range: perbedaan antara nilai max dan min dr data.

- Variance atau Standard Deviation: ukuran paling umum utk penyebaran dr suatu himpunan tsb.

$$\text{variance } (S_x^2) = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})^2$$

sensitif thd

outliers, ada pengukuran lain:

$$\text{absolute average deviation } AAD(x) = \frac{1}{m} \sum_{i=1}^m |x_i - \bar{x}|$$

$$\text{MAD}(x) = \frac{\text{median } \{|x_1 - \bar{x}|, \dots, |x_m - \bar{x}|\}}{\text{median absolute deviation}}$$

Interquartile range(x) = $x_{75\%} - x_{25\%}$

II Visualisasi

: konversi data ke dalam format visual / tabular sng karakteristik data dan keterhubungan antar item data/atribut dapat dianalisis.

: teknik yg powerful dan menarik utk eksplorasi data

- mendekripsi pola umum dan tren
- mendekripsi outliers dan pola yg tidak umum.

Representasi

: pemetaan informasi ke dalam format visual.

: objek data, atribut, keterhubungan antar objek data diterjemahkan ke dlm bentuk elemen grafik, seperti points, lines, shapes, and colors.

Arrangement

: penempatan elemen^{2x} visual dlm suatu tampilan.

: dapat membuat perbedaan signifikan dlm membaca / memahami data.

Teknik Visualisasi

1. Histogram: dibagi^{2x} kedlm bin

- distribusi nilai dari suatu variabel (single variable)

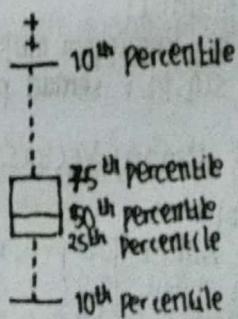
- tinggi bar: jumlah objek

- bentuk histogram tergantung jumlah bin.

2. 2-Dimensional Histogram

: joint distribution 2 atribut

3. Box Plots: membandingkan atribut^{2x} + outlier



4. Scatter Plots

- atribut value: posisi

Kemiripan dan Ketidakmiripan

• Similarity

semakin tinggi nilainya, makin tinggi kemiripannya [0, 1]

• Ketidakmiripan

semakin besar jaraknya, semakin beda.

- nilai minimum = 0

Attribute Type	Dissimilarity	Similarity
Nominal	$d = \begin{cases} 0 & p=q \\ 1 & p \neq q \end{cases}$	$s = \begin{cases} 1 & p=q \\ 0 & p \neq q \end{cases}$
Ordinal	$d = \frac{ p-q }{n-1}$ (value mapped to integers 0 to n-1, where n is the number of values)	$s = 1 - \frac{ p-q }{n-1}$
Interval/ Ratio	$d = p-q $	$s = -d, s = \frac{1}{1+d} \text{ or } s = 1 - \frac{d - \min.d}{\max.d - \min.d}$

• Euclidean Distance

$$dist = \sqrt{\sum_{k=1}^n (P_k - Q_k)^2}$$

• Minkowski Distance

$$dist = \left(\sum_{k=1}^n |P_k - Q_k|^r \right)^{\frac{1}{r}}$$

• Mahalanobis Distance

$$\text{mahalanobis}(p, q) = (p-q) \Sigma^{-1} (p-q)^T ; \Sigma : \text{matriks kovarian}$$

$$\Sigma_{j,k} = \frac{1}{n-1} \sum_{i=1}^n (X_{ij} - \bar{X}_j)(X_{ik} - \bar{X}_k)$$

data masukan X

Karakteristik Umum

a. Jarak (Distance)

1. $d(p, q) \geq 0$; $d(p, q) = 0$ hanya jika $p = q$
(Positive definiteness)
2. $d(p, q) = d(q, p)$; semua p dan $q \rightarrow$ Symmetric
3. $d(p, r) \leq d(p, q) + d(q, r)$; semua $p, q, r \rightarrow$
Triangle Inequality
; $d(p, q) = \text{jarak ketidakmimpinan antara objek } p \text{ dan } q$

b. Kemiripan (Similarity)

1. $s(p, q) = 1$; kemiripan maks hanya jika $p = q$
2. $s(p, q) = s(q, p)$; semua p dan $q \rightarrow$ Simetri

Similarity antara Binary Vectors

→ p dan q = Objek

→ M_{01} = jumlah atribut jika $p=0, q=1$

M_{10} = jumlah atribut jika $p=1, q=0$

M_{00} = jumlah atribut jika $p=0, q=0$

M_{11} = jumlah atribut jika $p=1, q=1$

• Simple Matching dari Koefisien Jaccard

$$SMC = \frac{\text{jumlah yg matches}}{\text{jumlah atribut}} = \frac{(M_{11} + M_{00})}{(M_{01} + M_{10} + M_{11} + M_{00})}$$

$$\text{Jaccard Coeffidens (J)} = \frac{M_{11}}{(M_{01} + M_{10} + M_{11})}$$

Contoh: $p = 10000000000$ $SMC = \frac{(0+7)}{2+1+7+0} = 0,7$

$$\begin{aligned} q &= 0000001001 \\ M_{01} &= 2 \quad M_{00} = 7 \\ M_{10} &= 1 \quad M_{11} = 6 \end{aligned}$$

$$J = \frac{0}{2+1+0} = 0,0$$

pra-~~process~~ data

o Why Data processing?

- Data masih kurang baik (data dirty)

1. Incomplete

- : tidak ada nilai atribut / kelengkapan data tidak utuh.
- contoh, occupation = ""

2. Noisy

- : error dan outliers.
- contoh, Salay = "-10"

3. Inconsistent

- sumber data yg berbeda
- redundansi data
- kesalahan instrumen saat mengumpulkan data
- human/computer error saat entry data
- kesalahan saat transmisi data

o Why is Data Preprocessing important?

"No quality data, no quality mining result!"

- data berkualitas → keputusan berkualitas
- Data Warehouse: perlu integrasi yg konsisten

ekstraksi data
pembersihan data
transformasi data

Major Task in Data Processing

1. Data Cleaning

- : mengisi nilai yg kosong, smooth noisy data, identify / remove outliers, resolve inconsistencies.

2. Data Integration

- : integrasi multiple DB, data cubes, / files.

3. Data Transformation

- : normalization & aggregation

4. Data Reduction

- : menghsl. jumlah data yg lbh sedikit, namun hasil analisis sama dgn ukuran aslinya.

5. Data Discretization

- : bagian dari Data Reduction, tetapi dengan kepentingan tertentu, khususnya data numerik.

1. Data Cleaning

Task:

- mengisi nilai yg kosong

- identifikasi outliers & smooth out noisy data
- memperbaiki inkonsistensi data
- mengatasi redundansi yg disebabkan oleh integrasi data.

o Missing data

- data tidak selalu available

eg: no record value for several attribute

- disebabkan oleh:

1. kerusakan peralatan
2. nilai atribut dihapus karena inkonsistensi dg data lain yg telah bersimpan.
3. data bdk diinput karena mis-understanding

4. data tertentu dianggap tdk penting saat dientry

- missing data harus tetap diisi

- Solusi:

1. mengabaikan record / sample tsb

2. isi nilai kosong manual

3. isi secara otomatis:

- global constant

- atribut mean

- atribut mean dari seluruh sample pd kelas yg sama (lebih baik)

- nilai yg plg mungkin: dugaan menng. formula Bayes / decision tree.

o Noisy data

- : keragaman / error acak pd variabel yg diukur.

- disebabkan oleh:

1. kesalahan instrumen saat pengumpulan data

2. masalah saat entry data

3. masalah saat transmisi data

4. keterbatasan teknologi

5. inkonsistensi penamaan

- Solusi:

1. binning

- : data diurutkan, dibagi ke dalam bins (sesuai inteval), smoothing (smooth by bin means, median, boundaries, etc)

2. Regresi
: smooth by fitting data into regression funct.
3. Clustering
: remove and detect outliers
4. Combined Computer & Human inspection
: deteksi nilai yg mencurigakan dan validasi oleh manusia.

Simple Discretization Methods: Binning

1. Equal-Width (distance) partitioning
: membagi nilai ke dalam range N interval dgn ukuran yg sama.
Jika A dan B: nilai atribut terkecil & terbesar, lebar interval $\Rightarrow W = \frac{B-A}{N}$
2. Equal-depth (frequency) partitioning
: membagi jumlah sampel yg sama kedalam N interval.

contoh:

data: 4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34

- Partisi: equal-frequency (equal-depth) bins

Bin 1: 4, 8, 9, 15

Bin 2: 21, 21, 24, 25

Bin 3: 26, 28, 29, 34

- Smoothing by bin means

Bin 1: 9, 9, 9, 9

Bin 2: 23, 23, 23, 23

Bin 3: 29, 29, 29, 29

- Smoothing by bin boundaries

Bin 1: 4, 4, 4, 15

Bin 2: 21, 21, 25, 25

Bin 3: 26, 26, 26, 34

2. Data Integration

: gabung data dari berbagai sumber

- Redundansi data: didetksi dgn analisis korelasi
 - Object identification: atribut/objek yg sama memiliki nama yg berbeda di DB yg berbeda.
 - Derivable data: atribut merupakan atribut "turunan" dari tabel lainnya.

Analisis Korelasi

a. Data Numerik

$$\rho_{A,B} = \frac{\sum (A - \bar{A})(B - \bar{B})}{(n-1)\sigma_A\sigma_B} = \frac{\sum (AB) - n\bar{A}\bar{B}}{(n-1)\sigma_A\sigma_B}$$

n = jumlah record

\bar{A}, \bar{B} = rataan A dan B

σ_A, σ_B = standar deviasi dari A dan B

$\sum(AB)$ = jumlah dr cross-product AB

- Jika $\rho_{A,B} > 0$ (A dan B korelasi \oplus) semakin tinggi nilai korelasi, semakin kuat korelasi A dan B. $A \uparrow, B \uparrow$

$\rho_{A,B} = 0$ (independent)

$\rho_{A,B} < 0$ (korelasi \ominus)

b. Categorical Data

χ^2 (chi-square) test

$$\chi^2 = \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}}$$

contoh:

	play chess	not play	sum (row)
like science	250(90)	200(360)	450
not like	50(210)	1000(840)	1050
sum (col)	300	1200	1500

$$\chi^2 = \frac{(250-90)^2}{90} + \frac{(50-210)^2}{210} + \frac{(200-360)^2}{360} + \frac{(1000-840)^2}{840} = 507.93$$

\therefore like science fiction dan play chess correlated in the group.

3. Data Transformation

- Smoothing: membuang noise dari data
- Aggregation: summarization, data cube construction
- Generalization: concept hierarchy climbing
- Normalization:

1. Min-max normalization: [newMinA, newMaxA]

$$V' = \frac{V - \text{min}_A}{\text{max}_A - \text{min}_A} \cdot (\text{newMax}_A - \text{newMin}_A) + \text{newMin}_A$$

eg: range pendekatan

\$12000 - \$98000 dinormalisasi ke [0, 1]. maka \$73000 akan dipetakan ke:

$$V' = \frac{(73000 - 12000)}{(98000 - 12000)} (1-0) + 0 = 0,716$$

2. Z-score normalization (μ : mean, σ : std deviation)

$$V' = \frac{V - \mu_A}{\sigma_A}$$

$$\text{eg: } \mu = 54000, \sigma = 16000 \Rightarrow V' = \frac{(73000 - 54000)}{16000} = 1.225$$

3. Normalization by decimal scaling

$$V' = \frac{V}{10^j}; j = \text{integer terkecil} \rightarrow \text{Max}(1|V'|) < 1$$

4. Data Reduction

Strategi:

- Data cube aggregation
- Dimensionality reduction: membuang atribut yg tdk penting
- Data Compression
- Numerosity reduction: eg. fil data into model
- Discretization and concept hierarchy generation

5. Diskretisasi

- 3 tipe atribut:
 1. Nominal
 - : unordered set. eg: color, profession
 2. Ordinal
 - : ordered set. eg: military / academic rank
 3. Continuous
 - : real number. eg: integer or real numbers.
- Discretization: data kontinu
 - membagi range atribut ke dalam interval
 - supervised vs. unsupervised
 - Split (top-down) vs. merge (bottom-up)
 - diterapkan dim atribut secara rekursif.
- Concept Hierarchy Generation for Categorical Data
 - Specification of a partial / total ordering of attributes explicitly at the schema level by users or experts
street < city < state < country
 - Specification of a hierarchy for a set of values by explicit data grouping
{Urbana, Champaign, Chicago} < Illinois
 - Specification of only a partial set of attributes only street < city, not others
 - Automatic generation of hierarchies (or attribute levels) by the analysis of the number of distinct values.
 - for a set of attributes: {street, city, state, country}
 - attribute with the most distinct value is placed at the lowest of the hierarchy.
eg: weekday, month, quarter, year

```

country (15 distinct values)
|
Province or state (365 ...)
|
city (3567 ...)
|
street (674339 ...)

```

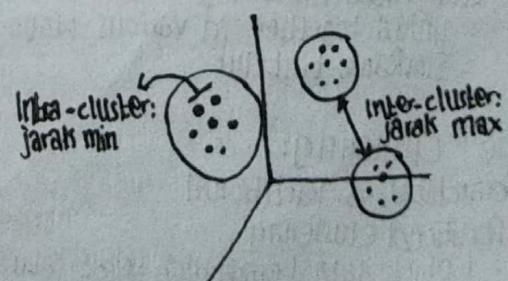
K⁴ 26/02/2018

Techniques	Data cleaning	Data integration	Data Transfo..	Data Reduct..	Discretizat
1. Binning	✓			✓	✓
2. Regression	✓			✓	
3. Clustering	✓			✓	✓
4. Correlation analysis		✓			
5. Decision tree	✓				✓
6. Normalization				✓	
7. Feature selection				✓	
8. Histogram				✓	✓
9. Sampling				✓	
10. Concept hierarchy generation			✓		✓

Analisis Cluster

Clustering

- teknik unsupervised learning
- tidak ada proper learning, hanya berdasar informasi dari data yg akan diolah.
- Tujuan: objek^{xx} yg mirip dalam cluster yg sama objek^{xx} yg berbeda dlm cluster yg sama
- Ideal: similarity / homogenitas semakin besar di dalam 1 grup & dissimilarity semakin besar di antara grup → clustering semakin baik.



Aplikasi Analisis Cluster

1. Understanding

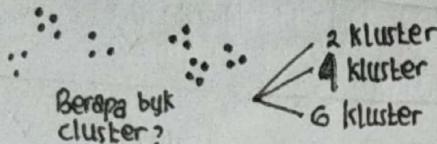
- : memahami pola / relasi antara objek data
- eg: - groups related documents for browsing.
- groups genes & proteins with similar functionality.
- groups stocks with similar price fluctuations.

2. Summarization

- : reduce the size of large data sets.

Penentuan cluster berkualitas bergantung dari kondisi data serta hasil yg diinginkan.

eg:



Major Clustering Approaches

1. Partitioning approach

- : membuat beberapa cluster berdasarkan kriteria tertentu (sudah diketahui berapa cluster)
- : K-Means, K-Medoids, CLARANS.

2. Hierarchical approach

- : dipotong berdasarkan kriteria tertentu (dekomposisi hierarki)
- : Diana, Agnes, BIRCH, CAMELEON

3. Density-based approach

- : berdasarkan pada koneksi dan density function.
- : DBSCAN, OPTICS, DenClue

4. Grid-based approach

- : berdasarkan pada multiple level granularity structure.
- : STING, WaveCluster, CLIQUE

5. Model-based

- : model hipotesis utk cluster dan dicari model terbaik (data baru langsung masuk ke cluster nantinya)
- : EM, SOM, COBWEB

6. Frequent pattern-based

- : analisis frequent patterns
- : P-Cluster

7. User-guided or constraint-based

- : by considering user specified/app specified constraint.
- : COD (obstacles), constrained clustering

8. Link-based clustering

- : linked together in various ways.
- : SimRank, LinkClus

Tipe^{**} Clustering:

Hierarchical vs. Partitional

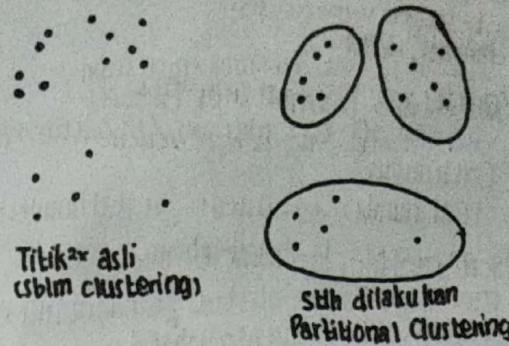
> Partitional Clustering

- : 1 objek data berada dlm tepat satu cluster
- : membagi himpunan objek data kedlm sub-himpunan (cluster) yg tdk overlap.

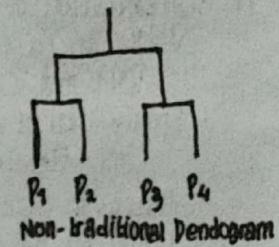
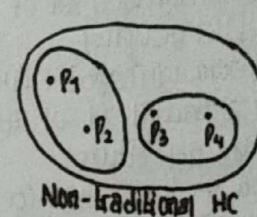
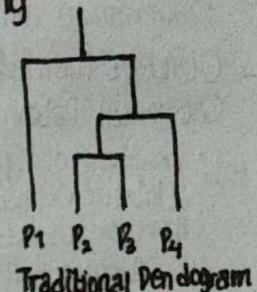
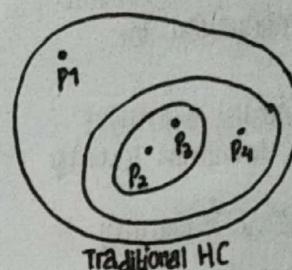
> Hierarchical clustering

- : cluster memiliki subcluster

Partitional Clustering



Hierarchical Clustering



Exclusive vs. Non-Exclusive

^l 1 objek,
^l 1 cluster

^l 1 objek bisa byk cluster/
kelas

Fuzzy vs. Non-Fuzzy

- : setiap objek mjd milik setiap cluster dg nilai keanggotaan dr antara 0 (mutlak tdk anggota cluster) dan 1 (mutlak anggota cluster)
- : penjumlahan bobot haruslah 1.

Partial vs. Complete

- ^l setiap objek tdk msk kedlm cluster
- ^l setiap objek msk kedlm cluster, karena beberapa objek dlm dataset mungkin bukan anggota kelompok yg tdk didefinisikan dg jarak.

Data Structure

1. Data Matrix : object by variable
 - : two modes \rightarrow row: objek, column: variabel
 - : N objects, p variables

$x_{11} \dots x_{1p}$
$\vdots \dots \vdots$
$x_{i1} \dots x_{ip}$
$\vdots \dots \vdots$
$x_{n1} \dots x_{np}$

2. Dissimilarity Matrix : object by object structure
 - : one mode.
 - (Δ atas dan btt nilaiya sama, pilih 1)

dim bentuk tree.

Type of data in Clustering Analysis

1. Interval-scaled variables

: standarisasi data

Similarity and Dissimilarity Between Objects

- ukuran kemiripan
- ukuran ketidakmiripan 2 objek.
- ukuran jarak
(distance)

Metode: Minkowski Distance

$$d(i,j) = \sqrt[q]{(x_{i1}-x_{j1})^q + (x_{i2}-x_{j2})^q + \dots + (x_{ip}-x_{jp})^q}$$

if $q=1 \rightarrow$ Manhattan Distance.

$q=2 \rightarrow$ Euclidean Distance.

Sifat: $d(i,j) \geq 0$

$$d(i,i) = 0$$

$$d(i,j) = d(j,i)$$

$$d(i,j) \leq d(i,k) + d(k,j)$$

2. Binary Variables

: Tabel Kontingenzi

		Object j		
		1	0	sum
Object i	1	a	b	a+b
	0	c	d	c+d
sum	a+c	b+d	p	

> Ukuran jarak utk distance variabel simetrik biner: ((Kemiripan), bukan jarak)

$$d(i,j) = \frac{b+c}{a+b+c+d}$$

> Ukuran jarak utk distance variabel asimetrik biner: ((ketidakmiripan))

$$d(i,j) = \frac{b+c}{a+b+c}$$

> Jaccard coefficient (ukuran similitudo utk variabel asimetrik biner):

$$\text{Sim}_{\text{Jaccard}}(i,j) = \frac{a}{a+b+c}$$

CONTOH:

Name	Gender	Fever	Cough	Test 1	Test 2	Test 3	Test 4
Jack	M	Y	N	P	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N	N	N	N

* Gender: atribut simetrik, lainnya asimetrik.

* Ubah Y dan P=1, N=0

$$d(\text{jack}, \text{mary}) = \frac{0+1}{2+0+1} = 0,33$$

$$d(\text{jack}, \text{jim}) = \frac{1+1}{1+1+1} = 0,67$$

$$d(\text{jam}, \text{mary}) = \frac{1+2}{1+1+2} = 0,75$$

3. Nominal Variables

: generalisasi dari variabel biner yaitu dapat mengambil lebih dari 2 nilai, contoh: red, yellow, blue, green.

Metode 1:

Simple Matching (pencocokan sederhana)

$$d(i,j) = \frac{p-m}{p}; m: \# \text{matches}$$

P: total variabel

Metode 2:

Gunakan sejumlah besar variabel biner

4. Ordinal Variables

: berupa nilai diskret / kontinu
: urutan \rightarrow penting

5. Ratio-Scaled Variables

: ukuran positif pd skala non-linear, sekitar skala eksponensial.

"Ukuran Jarak Antar Cluster"

1. Complete Linkage Clustering

: jarak maks antar elemen dim cluster.

$$d(A, B) = \max \{s(x, y) \mid x \in A, y \in B\}$$

2. Single Linkage Clustering

: jarak minimum antar elemen setiap cluster

$$d(A, B) = \min \{s(x, y) \mid x \in A, y \in B\}$$

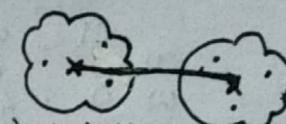
3. Average Linkage Clustering

: rata-rata jarak antar elemen setiap cluster.

$$d(A, B) = \frac{1}{n_A n_B} \sum_{x \in A} \sum_{y \in B} s(x, y)$$

4. Centroid Linkage

$$d(A, B) = s(\bar{x}, \bar{y})$$



5. Ward Linkage

$$d(A, B) = \frac{n_A n_B S_{AB}^2}{n_A + n_B}$$

jarak cluster A dan B dengan formula centroid linkage.

sensitif thd outliers.

K-Means Clustering

Complexity $O(n * k * I * d)$

I: # titik sampel

k: # cluster

I: # iterasi

d: # atribut

Partitional clustering approach
Setiap cluster terasosiasi dgn sebuah centroid (titik tengah)

Setiap titik dikumpulkan pd cluster dgn centroid terdekat

Jumlah cluster "k" hrs ditentukan
Centroid awal ditentukan: random, dari nilai rata-rata (mean) dr titik x dlm suatu cluster.

Konvergen: mengg. ukuran similitudo.
closeness: mengg. Euclidean dist, cosine simila

K-Medoids Clustering

Medoids: cari objek yg representatif didlm cluster.

PAM (Partitioning Around Medoids, 1987)

1. pilih k objek yg representatif secara acak.
2. Utk tiap pasangan objek yg tidak dipilih i h dan objek yg dipilih j, hitung total cost utk melakukan pertukaran (swapping) $T_{Gi,h}$
3. Utk tiap pasangan dari i dan h,
 - if $T_{Gi,h} < 0$, i drgantikan h
 - selanjutnya, tetapkan swap^{xx} objek yg tdk dipilih ke objek representatif yg paling minip.
4. Ulangi langkah 2-3 sampai tidak ada perubahan.

Contoh:

	A1	A2
(1)	2	6
(2)	3	4
(3)	3	8
(4)	4	7
(5)	6	2
(6)	6	4
(7)	7	3
(8)	7	4
(9)	8	5
(10)	7	6

$$K=2$$

- pilih acak 2 medoid
- eg: $(2) = (3, 4)$
- $(8) = (7, 4)$

→ Tetapkan setiap objek pd objek representatif berdekat. ((Manhattan distance))

→ Medoid $(2) = (3, 4)$

$$d(\text{data}_1, (2)) = |3-2| + |4-6| = 3$$

$$d(\text{data}_2, (2)) = |3-3| + |4-4| = 0, \text{ dst.}$$

	C1	C2	
(1)	3	7	→ Tentukan cluster 1
(2)	0	4	dan cluster 2;
(3)	4	8	lihat jarak paling
(4)	4	6	kecil masing ^{xx} objek.
(5)	5	3	
(6)	3	1	∴ Cluster 1 = { (1), (2), (3), (4) }
(7)	5	1	Cluster 2 = { (5), (6), (7), (8), (9), (10) }
(8)	4		
(9)	6	2	
(10)	6	2	

→ Hitung kriteria absolute error

(utk himpunan Medoids $((2), (8))$) → nilai dari jarak yg sdh dihitung.

$$E = \sum_{i=1}^K \sum_{p=c_i} |p - O_i| = |O_1 - O_2| + |O_3 - O_2| + |O_4 - O_2| + \\ |O_5 - O_8| + |O_6 - O_8| + |O_7 - O_8| + \\ |O_9 - O_8| + |O_{10} - O_8| \\ = (3+4+4) + (3+1+1+2+2) \\ = 20$$

→ Pilih objek acak: (7)

→ Tukar O_8 dan O_7

→ medoid $(7) = (7, 3)$

$$d(\text{data}_1, (7)) = |7-2| + |3-6| = 8 \\ d(\text{data}_2, (7)) = |7-3| + |3-4| = 5, \text{ dst.}$$

Update tabel jarak C₂:

	C1	C2
(1)	3	8
(2)	0	5
(3)	4	9
(4)	4	7
(5)	5	2
(6)	3	2
(7)	5	0
(8)	4	1
(9)	6	3
(10)	6	3

∴ clusteranya masih sama:
Cluster 1 = {1, 2, 3, 4}
Cluster 2 = {5, 6, 7, 8, 9, 10}

→ Hitung kriteria absolute error
(utk himpunan Medoids (O_2, O_7))

$$E = (3+4+4) + (2+2+1+3+3) \\ = 22$$

→ Hitung nilai cost dari fs. Swap O_8 dan O_7

$$S = \text{absolute error } [O_2, O_7] - \text{absolute error } [O_2, O_8]$$

$$= 22 - 20$$

= 2 → $S > 0 \rightarrow$ bukan ide yg baik menggantikan O_8 dg O_7 .

Kelebihan:

- PAM -

l lebih kokoh (robust) daripada K-Means (dgn adanya noise dan outlier) karena medoid kurang dipengaruhi oleh outlier atau nilai ekstrem lainnya.

l efisien utk himpunan data kecil, tdk sesuai dg himpunan data besar.

l Kompleksitas: $O(k(n-k)^2)$ utk tiap iterasi

l metode berbasis penarikan sampel (sampling): CLARA (Clustering LARge Applications).

Hierarchical Clustering

- : himpunan cluster bersarang dalam bentuk tree.
 - : Visualisasi: Dendogram
- (+) 1. Tidak perlu mengasumsikan jmh cluster.
dpt memotong dendrogram pd tingkat yg lepas
utk mendapatkan cluster yg diinginkan.
2. Menghit. klasifikasi yg bermakna

2 tipe utama:

a. Agglomerative

- 1 objek, 1 cluster (mulai dr titik sbg kluster individu): jumlah cluster awal sebanyak jmh titiknya.
- setiap langkah, gabungkan pasangan yg berdekat sampai hm ya tersisa 1 kluster (/ k cluster)

Strategi bottom-up:

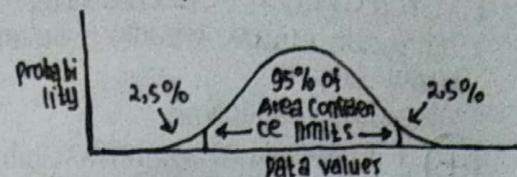
setiap objek dlm cluster kecil kemudian bergabung (merge) jd 1 cluster yg lbh besar sehingga objek berada pd cluster tunggal/sampai kondisi tertentu yg diinginkan.

b. Divisive

- semua titik disatukan dlm 1 kluster (mulai dr 1 cluster, all-inclusive cluster).
- utk setiap tahapan, pisahkan suatu cluster sampai setiap cluster mengandung 1 titik (/ berdapat sebanyak k cluster).

Strategi top-down:

menempatkan semua objek dlm 1 cluster utama, kemudian diempatkan jd sub-divisi cluster yg lbh kecil sesuai kondisi tertentu.



Detectasi Anomali

- Anomaly / Outlier
 - : sekumpulan titik data yg sangat berbeda (Tan) / sgt tidak mirip (dissimilar / Tan) dgn data lainnya.

- Anomaly Detection Schemes

- * General steps

1. Buat profile dr kondisi "normal"
 - : dpt berupa pola/ringkasan secara statistik dr seluruh populasi.
2. Menggunakan profil "normal" utk deteksi anomali.
 - : hasilnya adalah objek yg sgt berbeda dr profil "normal"

- > Tipe :

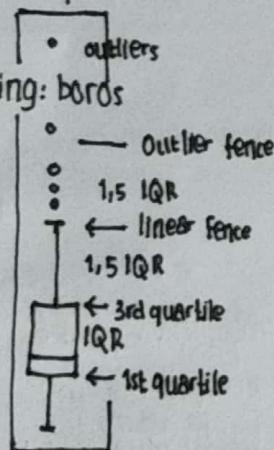
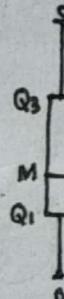
1. Graphical & Statistical-based
2. Distance-based
3. Model-based

1. Graphical Approaches

- Boxplot (1D), Scatter-plot (2D), Spin plot (3D)

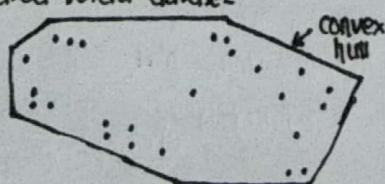
- Limitation:

- Time consuming: borders
- Subjective



Convex Hull Method

- : titik ekstrim diasumsikan sbg outlier.
- : untuk mendekati nilai ekstrim
- : pd area border dataset



Statistical Approaches

- : mengasumsikan model parametrik yg mendeskripsikan distribusi data (eg: distribusi normal).
- : Uji statistik sesuai dgn: distribusi data; parameter dr distribusi (eg: mean, variance); jml outlier yg diharapkan.

2. Distance-based Approaches

- l mengatasi kelemahan/keterbatasan metode statistik.
- l data direpresentasikan sbg vektor.
- l mengukur jarak masingx objek.

- l 3 Major Approaches:

- Nearest-neighbor based
- Density based: kerapatan sebaran data.
- Clustering based

1. Nearest-Neighbor Based Approach

- : jarak antara data ke data lain secara mengeluruh (Matrix Jarak antar pasangan Objek)

- : Mendefinisikan Pencikan:

- titik data hanya memiliki

\leftarrow jumlah tetangga yg lbh sedikit dari nilai yg dibentaskan dalam jarak D.

\leftarrow jarak > jarak \leftarrow tertinggi ke tetangga minimum berdekat ke-k.

\leftarrow jarak rata-rata \leftarrow rata-rata tertinggi ke k objek \leftarrow tetangga terdekat. threshold

2. Density-Based Local Outlier Detection

- l didasarkan pd global distance distribution
- l mengamati kepadatan sekelompok populasi.
- l Metode distance blm cukup, maka perlu: Local outlier.

Local Outlier Factor (LOF)

- l setiap objek: hitung densitas dari ketetanggaannya
- l setiap titik punya LOF
- l setiap titik hitung LOF - density objek p dan density tetangganya.
- Jika LOF terbesar: outlier

Contoh:

Diketahui 4 titik data:

a(0,0), b(0,1), c(1,1), d(3,0)

Calculate the LOF for each points & show the top 1 outlier, set k=2 and use Manhattan Distance.

Jawab:

> Step 1: Calculate distance

	a	b	c	d
a	0	1	2	3
b		0	1	4
c			0	3
d				0

> Step 2: Calculate $dist_2(o)$

: jarak maksimum utk belantara masing-masing objek. $k=2 \rightarrow dist_2$

$$dist_2(a) = dist(a, c) = 2 \quad (c \text{ is the 2nd nearest neighbor})$$

$$dist_2(b) = dist(b, a) = 1 (a/c \dots)$$

$$dist_2(c) = dist(c, a) = 2 (a \& \dots)$$

$$dist_2(d) = dist(d, a) = 3 (a/c \dots)$$

> Step 3: Calculate $N_k(o)$

$$N_k(o) = \{o' | o' \text{ in } D, dist(o, o') \leq dist_k(o)\}$$

$$N_2(a) = \{b, c\}$$

$$N_2(b) = \{a, c\}$$

$$N_2(c) = \{b, a\}$$

$$N_2(d) = \{a, c\}$$

> Step 4: Calculate $lrd_k(o)$: Local Reachability Density of o

$$||N_k(o)||$$

$$lrd_k(o) = \frac{||N_k(o)||}{\sum_{o' \in N_k(o)} reachdist_k(o' \leftarrow o)}$$

$$(o \leftarrow o') = \max \{ dist_k(o), dist(o, o') \}$$

$$\cdot lrd_2(a) = \frac{||N_2(a)||}{reachdist_2(b \leftarrow a) + reachdist_2(c \leftarrow a)}$$

$$\rightarrow reachdist_2(b \leftarrow a) = \max \{ dist_2(b), dist(b, a) \} \\ = \max \{ 1, 1 \} = 1$$

$$\rightarrow reachdist_2(c \leftarrow a) = \max \{ dist_2(c), dist(c, a) \} \\ = \max \{ 2, 2 \} = 2$$

$$lrd_2(a) = \frac{2}{1+2} = 0,667$$

$$\cdot lrd_2(b) = \frac{2}{2+2} = 0,5$$

$$\cdot lrd_2(c) = \frac{2}{1+2} = 0,667$$

$$\cdot lrd_2(d) = \frac{2}{3+3} = 0,33$$

> Step 5: Calculate $LOF_k(o)$

$$LOF_k(o) = \sum_{o' \in N_k(o)} \frac{lrd_k(o')}{||N_k(o)||}$$

$$||N_k(o)||$$

$$= \sum_{o' \in N_k(o)} lrd_k(o') \cdot \sum_{o'' \in N_k(o)} reachdist_k(o' \leftarrow o'')$$

$$\cdot LOF_2(a) = (lrd_2(b) + lrd_2(c)) * (reachdist_2(b \leftarrow a) + reachdist_2(c \leftarrow a)) \\ = (0,5 + 0,667) * (1 + 2) = 3,501$$

$$\cdot LOF_2(b) = 5,336$$

$$\cdot LOF_2(c) = 3,501$$

$$\cdot LOF_2(d) = 8,004$$

> Step 6: Sort all the $LOF_k(o)$

$$LOF_2(d) = 8,004$$

$$LOF_2(b) = 5,336$$

$$LOF_2(a) = 3,501$$

$$LOF_2(c) = 3,501$$

∴ Top 1 Outlier is point d.

- Outlier Discovery:

Deviation-Based Approach —

: identifikasi outlier dgn mengamati karakteristik dr tdk objek dlm cluster.

Sequential Exception Technique

↳ cara agar manusia dpt membedakan objek tdk umum dari serangkaian objek^{2x} dgn karakteristik tertentu.

OLAP data cube technique

↳ mengg. data cubes utk identifikasi daerah anomali dgn multi dimensi yg besar.

- 3. Clustering-Based —

: mengelompokkan data ke dalam kepadatan (density) yg berbeda.

Kandidat Outlier

↳ pilih titik^{2x} dlm cluster kecil

: hitung jarak antara titik kandidat outlier dgn titik^{2x} kluster^{2x} yg blk jadi kandidat outlier.

↳ jika titik kandidat jauh dari semua titik yg bukan kandidat, maka tsb adalah outlier.

DBSCAN

(Density-Based Spatial Clustering of Application with Noise)

: mengelompokkan titik dgn titik^{2x} yg memiliki btngga dekat.

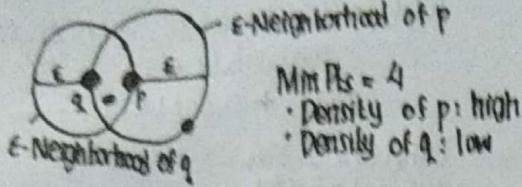
: Outlier: titik yg terletak sendirian/jauh dari titik^{2x} lainnya & density kecil.

CLASSIFICATION

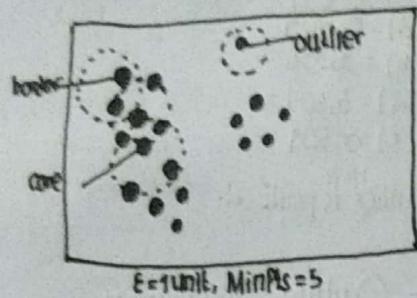
DECISION TREE

Definisi Kepadatan (Density)

- > ϵ -Neighborhood: objek dan radius ϵ dari suatu objek.
- $N_\epsilon(p) = \{q | d(p, q) \leq \epsilon\}$
- > High density: ϵ -Neighborhood of an object contains at least MinPts of object.



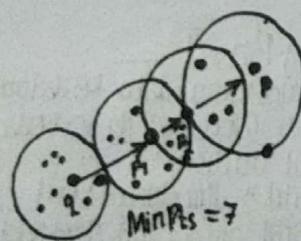
Core, Border, & Outlier



- > Core Point: titik memungkai lebih dari $Min Pts$ dlm ϵ (p)
: titik ini berada pd bagian dlm (interior) dr cluster.
- > Border Point: titik memiliki lbh sedikit dr $Min Pts$, berada dlm lingkungan dr Core Point.
- > Noise Point: bukan Core point & Border point.

Density-Reachability

: directly dan indirectly



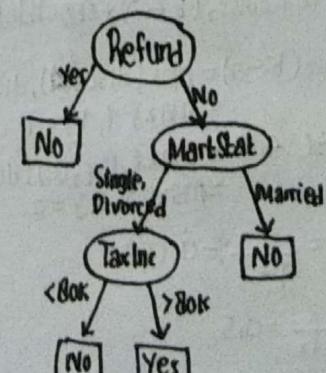
- > titik directly density-reachable dari p_2
- > p_2 directly density-reachable dari p_1
- > p_1 directly density-reachable dari q
(rancangan $p \leftarrow p_2 \leftarrow p_1 \leftarrow q$)
 - ↳ p (indirectly) density-reachable dari q
 - ↳ q is not density-reachable from p .

- (+) 1. Tahan thd noise
2. dpt menangani kluster dari bentuk dan ukuran yg berbeda
- (-) 1. Tidak mampu menangani berbagai variasi kepadatan.
2. Sensitif thd parameter
: sulit menentukan himpunan parameter yg benar Parameter.

- > Definisi: proses training + testing
 - diberikan kumpulan data (training set)
 - setiap record: ada informasi kelas
 - dilakukan proses training, menghasilkan model
(bisa untuk memprediksi data baru)
 - model diterapkan pada test set, dilakukan proses testing.
- > Classification Techniques:
 1. Decision Tree Based Methods
 2. Rule-based Methods
 - : dari hasil (1) didapatkan aturan** utk pengelompokan kelas
 3. Memory based reasoning
 4. Neural Networks
 5. Naive Bayes & Bayesian Belief Networks
 6. Support Vector Machines

— Decision Tree —

eg:	categorical		continuous	class
Training set	Refund	Marital status	Taxable Income	Cheat
1.	Yes	Single	125K	No
2.	No	Married	100K	No
3.	No	Single	70K	No
4.	Yes	Married	120K	No
5.	No	Divorced	95K	Yes
6.	No	Married	60K	No
7.	Yes	Divorced	220K	No
8.	No	Single	85K	Yes
9.	No	Married	75K	No
10.	No	Single	90K	Yes



Test Data	Refund	Marital status	Taxable Income	Cheat
1.	No	Married	80K	? ((No))

> Algoritme:

- Hunt's Algorithm (one of the earliest)
- Classification and Regression Trees (CART)
- Iterative Dichotomiser 3 (ID3), C4.5
- SLIQ, SPRINT

1. Hunt's Algorithm

- diketahui: D_t = training records
- If D_t memiliki kelas sama dgn y_t , maka t is a leaf node labeled as y_t .
 - If D_t is empty set, then t is a leaf node labeled by default class y_d .
 - If D_t memiliki >1 class, use an attribute test to split the data into smaller subset.

eg:



Tree Induction

1. Greedy Strategy

- memerlukan sampel data berdasarkan atribut yg dipilih \rightarrow diperoleh subset data training.

2. Issues:

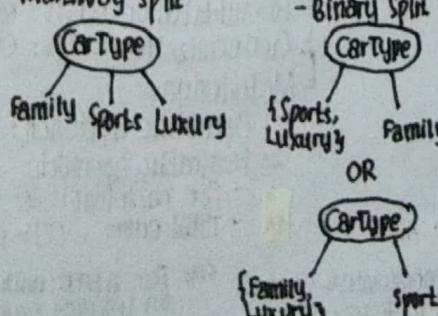
- Determine how to split the records.
 - How to specify the attribute test condition?
 - How to determine the Best split?
- Determine when to stop splitting.

1. Specify Test Condition

- bergantung tipe atribut: nominal, ordinal, continuous
- bergantung jumlah split: 2-way split (dikenal optimal partitioning)
multiway split.

eg: • Splitting Nominal Attribute

- Multiway Split

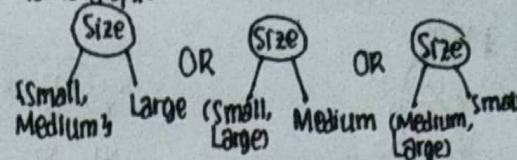


• Splitting Ordinal Attribute

- Multiway split



- Binary Split



• Splitting Continuous Attribute

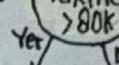
1. Proses Diskretisasi

- Statics: dikelompokkan diawal
(eg: dibagi jd 5)

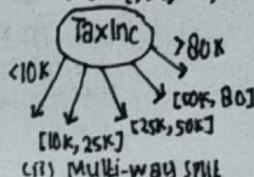
- Dynamic: range by equal interval, equal frequency (percentile), or clustering

2. Binary Decision: $(A < V)$ or $(A \geq V)$

eg: TaxInc



(i) Binary split



(ii) Multi-way split

2. How to determine the Best Split?

Ideal: semua leaf berasal dari kelas yg sama.

> Greedy Approach: node with homogeneous class distribution lebih diutamakan.

Need a measure of node impurity:

C0: 5

C1: 5

Non-Homo

geneous

(High-degree of

Impurity)

C0: 9

C1: 1

Homogeneous

(Low-degree of

Impurity)

Measure of Node Impurity

- Gini Index
- Entropy
- Misclassification error

1. GINI

- pilih nilai GINI terkecil.

node t.

$$GINI(t) = 1 - \sum_j [P(j|t)]^2$$

relative frequency of class j at node t

> Min (0.0): if all records belong to 1 Class.

> Maximum ($1 - 1/n_c$): records are equally distributed

split \leftrightarrow among all classes least interesting information

eg: $C_1: 0 \rightarrow P(C_1) = 0\% = 0$

$$P(C_2) = \frac{5}{6} = 1 \rightarrow GINI = 1 - P(C_1)^2 - P(C_2)^2$$

$$C_1: 1 \rightarrow P(C_1) = \frac{1}{6} \rightarrow GINI = 1 - (\frac{1}{6})^2 - (\frac{5}{6})^2 = 0.278$$

Splitting Based On GINI

- digunakan dalam CART, SLIQ, SPRINT

- node p is split into k partitions (children), quality of split: $\frac{n_i}{n_p}$ \rightarrow data pada child i

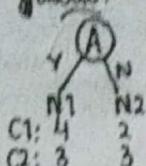
$$GINI_{split} = \sum_{i=1}^k \frac{n_i}{n_p} GINI(i)$$

eg: Parent:
 $C_1: G > \text{GINI} = 0.5$

Split: binary partition
 (2 parts)

Mana yg lebih baik?

Jawab:



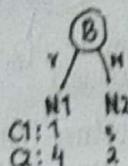
$$\text{GINI}(NT) = 1 - (4/7)^2 - (3/7)^2 = 0,4898$$

$$\text{GINI}(N_2) = 1 - (3/3)^2 - (3/3)^2 = 0,48$$

$$\text{GINI}(\text{children}) = \frac{3}{12} * 0,4898 + \frac{7}{12} * 0,48 = 0,486$$

$$\text{GINI}(B) < \text{GINI}(A)$$

∴ choose attribute B,
 karena lebih murni (pure)



$$\text{GINI}(N_1) = 0,32$$

$$\text{GINI}(N_2) = 0,488$$

$$\text{GINI}(\text{children}) = \frac{5}{12} * 0,32 + \frac{7}{12} * 0,488 = 0,375$$

3 Classification Error

: node t.

$$\text{Error}(t) = 1 - \max_i P(i|t)$$

→ Maximum ($1 - 1/n_c$)
 → Minimum ($0,0$) } = GINI

$$\text{eg: } \begin{cases} C_1: 0 \\ C_2: 6 \end{cases} \rightarrow P(C_1) = 0/6 = 0 \\ P(C_2) = 6/6 = 1 \Rightarrow \text{Error} = 1 - \max(0, 1) = 1 - 1 = 0$$

- Determine When to stop splitting -

Criteria Stop:

1. Hentikan mempartisi suatu node ketika semua record jadi anggota kelas yg sama.
2. Hentikan mempartisi suatu node ketika semua record memiliki nilai atribut yg minip.
3. Penghentian sebelum salah satu kedua kondisi diatas terpenuhi (Early termination).

Decision Tree Based

• Advantages:

1. Inexpensive
2. Fast: klasifikasi record baru yg blm diketahui kelasnya
3. Easy to interpret: if pohon yg diperlukan uk. kecil.
4. Accuracy

• Issue of Classification

1. Underfitting and Overfitting
2. Missing Values
3. Costs of Classification

1. Underfitting

- model berlalu umum → tidak dapat membedakan kelas
- kedua error hasil training dan testing nilainya besar.

Overfitting

- dalam proses training, model yg dihasilkan berlalu ideal, karena misclassified diperbaiki → training akurasi 100% → tidak mampu menggeneralisir (modelnya), shg ketika ada data test baru, error akan semakin naik.
- error ketika training set kecil, tapi error test set besar.
- disebabkan: sampel data training yg digunakan kurang

Estimating Generalization Errors

↳ Re-substitution errors: error on training ($\sum e(t)$)

↳ Generalization errors: error on testing ($\sum e'(t)$)

[Metodenya:]

1. Optimistic approach: $e'(t) = e(t)$

2. Pessimistic approach:

- For each leaf node: $e'(t) = e(t) + 0,5$

- Total error: $e'(T) = e(T) + N \times 0,5$

↳ number of leaf nodes

• Generalization error

$$= \frac{(10 + 30 \times 0,5)}{1000} = 2,5\%$$

↳ Training error = $10/1000 = 1\%$

eg: ((lihat slide hal. 18))

2. Entropy

: node t.

$$\text{Entropy}(t) = - \sum P(j|t) \log_2 P(j|t)$$

↳ Maximum ($\log_2 n_c$)
 ↳ Minimum (0-0)

$$\text{eg: } \begin{cases} C_1: 0 \\ C_2: 6 \end{cases} \rightarrow P(C_1) = 0/6 = 0 \\ P(C_2) = 6/6 = 1 \\ = -0 * \log_2 0 - \\ 1 * \log_2 1 \\ = -0 - 0 = 0$$

Splitting Based On Entropy

• Information Gain

$$(1) \text{ menghitung jumlah partisi yg besar, } GAIN_{\text{split}} = \text{Entropy}(P) - \left(\sum_{i=1}^k \frac{n_i}{n} \text{Entropy}(i) \right)$$

; P = parent node

k = # split

Measure Reduction in Entropy

: choose split that achieves most reduction (maximizes GAIN)

• Gain Ratio

$$\text{GainRatio}_{\text{split}} = \frac{\text{GAIN}_{\text{split}}}{\text{SPLITINFO}}$$

$$\text{SPLITINFO} = - \sum_{i=1}^k \frac{n_i}{n} \log_2 \frac{n_i}{n}$$

(2) mengalasi heterogenitas

homogen / pure

- > Reduced error pruning (REP)
 - : mengg. data testing (validasi) utk estimasi generalization error.

Mengatasi Overfitting

1. Pre-Pruning (Early Stopping Rule)

- l Hentikan algoritme sebelum jd "a fully-grown tree"
- l Kondisi berhenti (stopping) suatu node:
 - l Hentikan jika semua instans (records) memiliki kelas label yg sama.
 - l Hentikan ketika semua nilai atribut adalah sama.
- l Kondisi yg lebih ketat:
 - l Hentikan ketika jmlh record di suatu node lbh kecil dr threshold yg didefinisikan oleh user.
 - l Hentikan ketika distribusi kelas dari instans / record^{2x} nyg a. independent dr fitur^{2x} yg ada. (eg: gunakan χ^2 test)
 - l Hentikan ketika perluasan suatu node tdk memperbaiki derajat homogenitas (impurity). (gunakan Gini/Information Gain).

2. Post-pruning

- l Buat decision tree sampai jd pohon lengkap.
- l Potong simpul^{2x}/node^{2x} pd decision tree dimulai dr bawah (bottom-up fashion).
- l Hitung generalization error.
 - l Jika generalization error jd lbh baik slh proses pemotongan (trimming), ganti sub-tree dg n leaf node.
- l Label kelas dr leaf node ditentukan dr label kelas mayoritas dr instans^{2x}/records pd sub-tree tsb
- l Dapat mengg. Minimum Description Length (MDL) utk post-pruning.

eg: Class: Yes 20

Class: No 10

$$\text{Error} = 10/30$$

> Training Error
(Before splitting) = $10/30$

> Pessimistic Error = $\frac{(10 + 0.5)}{30} = 10.5/30$

> Training Error
(After splitting) = $9/30$

> Pessimistic Error
(After splitting) = $\frac{(9 + 4 * 0.5)}{30} = 11/30$

Relaasi Pemisah

dengan SVM, Seleksi, dan Evaluasi Model

Pendekatan dengan Machine Learning Terminologi Dasar

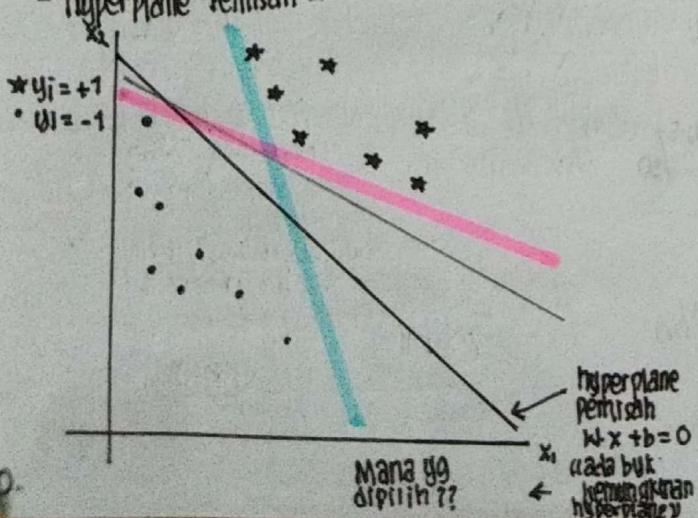
1. Training set: N data $\{x_1, \dots, x_n\}$
 - digunakan utk tuning parameter dari model
 - kategori dari data training set telah dikelahui sebelumnya
2. Target Vector: vektor unik f utk tiap target
 - merepresentasikan identitas dari data yg bersetujuan.
3. Learned function: $y(x)$
 - Training phase (fase pelatihan): proses untuk menentukan $y(x)$ berdasarkan data latih (training data)
4. Data uji (Test set): data yg tdk berdapat pd data latih (training set)
 - setelah model selesai dilatih, model dpt menentukan kategori dari data baru.
 - kemampuan utk mengkategorikan dg benar data baru yg berbeda dg n data yg digunakan pd training set
 → GENERALISASI

Definisi Masalah

- diberikan sekumpulan n titik (vectors): x_1, x_2, \dots, x_n dimana x_i adl vektor dg pjg m dan masing x_i adalah anggota salah satu dari dua kelas yg memiliki label "+1" dan "-1".
- maka training set:

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

$$\forall x_i \in \mathbb{R}^m, y_i \in \{+1, -1\}$$
- diinginkan untuk menemukan hyperplane $w^T x + b = 0$ yg memisahkan titik x tsb ke dalam 2 kelas: Positif (kelas "+1") dan Negatif (kelas "-1").
 (Asumsi: titik x dpt dipisahkan secara linier).
- Hyperplane Pemisah -



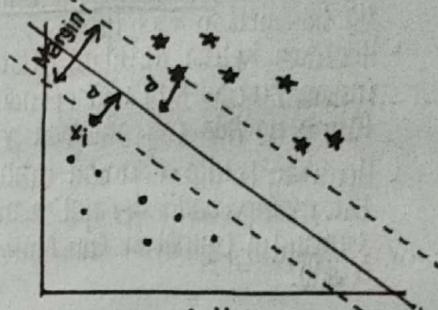
Memilih Hyperplane Pemisah

Pendekatan SVM: Kasus Linear Separable

- Hyperplane: sejauh mungkin dari titik sampel (ideal)
- sehingga data baru yg dolarat dg sampel data akan dikelasifikasikan dg benar.
 → Generalisasi Baik.
- Ide dasar dari SVM: memaksimalkan distance (jarak) antara hyperplane dan titik sampel terdekat.

Pada optimal hyperplane: jarak ke titik negatif terdekat = jarak ke titik positif terdekat.

- Tujuan SVM:
 memaksimalkan Margin yg besarnya $2 \times$ jarak "d" antara hyperplane pemisah dg sampel terdekat.



Untuk mencari Hyperplane pemisah yg optimal, daerah yg dibentuk oleh hyperplane dan margin dpt: diputar ke kanan / ke kiri
 dipersempit / diperlebar

- Tujuan: - mendapatkan hyperplane dg margin yg maksimal.
- data pd kelas positif ($y=+1$) dan kelas negatif ($y=-1$) tetap terpisah.

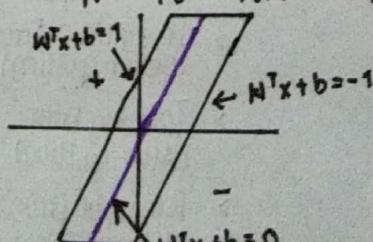
* Hyperplane dpt direpresentasikan oleh vektor normal w dan suatu skalar b :

$$w^T x + b = 0$$

- Vektor normal w : menentukan orientasi dg hyperplane.
 (ingat: gradien (m) pd $y = mx + c$)
- Bias b : menentukan pergeseran dari titik asal/pusat (origin)

* Margin dpt direpresentasikan dg 2 buah hyperplane:

$$w^T x + b = 1 \text{ (kelas positif)}$$

$$w^T x + b = -1 \text{ (kelas negatif)}$$


- Dengan mengubah sudut dg vektor normal w , margin dpt dirotasikan.
- Menggeser margin: dg memperbesar/mempertekan bias b .

- Lebar Margin (cm) = $\frac{2}{\|w\|}$

$$M = \frac{2}{\|w\|}$$

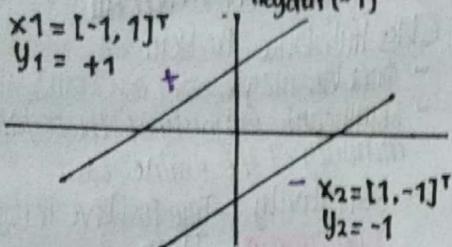
Bentuk: lebar margin M berbanding terbalik dgn panjang vektor normal ($\|w\|$)

Makin besar panjang vektor normal, makin sempit margin (cm), dan sebaliknya.

Memecahkan Masalah Optimasi:

Maksimalkan Lebar Margin (m)

misal: x_1 dan x_2 , label: y_1 dan y_2 , kelas: positif (+1) negatif (-1)



$$\text{Lebar Margin (m)} = \frac{2}{\|w\|}$$

MAX $m \rightarrow \max\left(\frac{2}{\|w\|}\right)$, artinya: panjang vektor normal w harus dikemalkan, maka:

$$\max\left(\frac{2}{\|w\|}\right) \rightarrow \min\left(\frac{\|w\|^2}{2}\right)$$

Untuk menghilangkan akar kuadrat, vektor normal harus dikuadratkan.

$$\max\left(\frac{2}{\|w\|}\right) \rightarrow \min\left(\frac{\|w\|^2}{2}\right) \rightarrow \min\left(\frac{\|w\|^2}{2}\right)$$

Konstrai:

- kelas positif: $w^T x_1 + b \geq 1$

- kelas negatif: $w^T x_2 + b \leq -1$

Selanjutnya, kalikan tiap konstrai dgn label yg sesuai.

$$y_1(w^T x_1 + b) \geq 1 \quad \text{dan} \quad y_2(w^T x_2 + b) \leq -1$$

Sehingga diperoleh:

$$\min\left(\frac{\|w\|^2}{2}\right) \text{ dgn constraint: } y_i(w^T x_i + b) - 1 \geq 0;$$

dikalikan masalah OPTIMASI

dgn Lagrange Multipliers utk membentuk The Primal-Dual Form:

$$\text{minimize } L_p(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i (y_i(w^T x_i + b) - 1)$$

$$\text{St. } \alpha_i \geq 0$$

$$\text{Selanjutnya: } \frac{\partial L_p}{\partial w} = 0 \rightarrow w = \sum_{i=1}^n \alpha_i y_i x_i$$

$$\frac{\partial L_p}{\partial b} = 0 \rightarrow \sum_{i=1}^n \alpha_i y_i = 0$$

Mendapatkan Lagrangian Dual Problem

Lagrangian Dual Problem:

$$\text{maximize } L_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$$

$$\text{St. } \alpha_i \geq 0 \text{ dan } \sum_{i=1}^n \alpha_i y_i = 0$$

// $\alpha_i = \{ \alpha_1, \alpha_2, \dots, \alpha_n \}$: variabel utk tiap sampelpoint x_i

- latihan:-

$$\text{data: } x_1 = [-2, 1]^T, x_2 = [1, -2]^T \\ y_1 = 1, y_2 = -1$$

((nanti dibelajang))

Rakteristik Solusi

- Banyak α_i yg bernilai not.

w_i = kombinasi linear dari sejumlah kecil data representasi sparse

- α_i non-zero α_i disebut support vectors (SV)

• decision boundary ditentukan oleh SV

• jika α_j ($j = 1, \dots, s$): indeks dan support, maka:

$$w = \sum_{j=1}^s \alpha_j y_j x_j$$

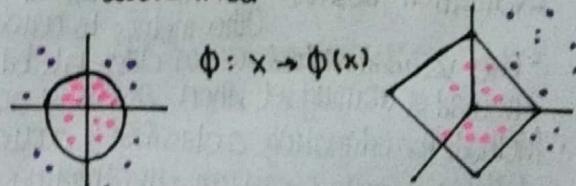
- Testing data baru z

$$\text{Hitung } w^T z + b = \sum_{j=1}^s \alpha_j y_j (x_j^T z) + b$$

dan klasifikasikan z sebagai kelas 1 jika hasil penjumlahannya positive, kelas 2 jika negatif.

Non-linear SVMs: Feature Spaces

Ide umum: original feature space dpt seluruh dipetakan ke beberapa higher-dimensional feature space dimana training set dpt dipisahkan secara linear \rightarrow VC Dimension



Perluasan ke Non-linear Decision Boundary

- Ide: transformasikan x_i ke sebuah ruang dimensi yg lebih tinggi

- Input space: ruang x_i berada

- Feature space: ruang dr $\phi(x_i)$ hasil transformasi

- Transformasi: perlu

- Linear operation pada feature space ekivalen dg non-linear operation pd input space

- Klasifikasi dpt dilakukan dg lbh mudah dg suatu transformasi yg sesuai, eg: XOR

- Masalah yg muncul: komputasi tinggi dan sulit mendapatkan estimasi yg baik

- SVM memecahkan 2 masalah tsb secara simultan.
- Kernel tricks: komputasi yg efisien
- Meminimalisasi $\|w\|^2$ utk mendapatkan a "good" classifier.

The Kernel Tricks

- Problem:
$$\max L_p(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j [\Phi(x_i^\top) \Phi(x_j)]$$

$$\text{s.t. } \alpha_i \geq 0, \quad \sum_{i=1}^n \alpha_i y_i = 0$$

fungsi keputusan:

$$f(\Phi(x)) = \text{sign}\left(\sum_{i=1}^n \alpha_i y_i [\Phi(x_i^\top) \Phi(x)] + b\right)$$

- Solusi:

- Tidak perlu memerlukan secara eksplisit mengetahui dimensi dr space baru, karena hanya perlu menghitung dot product dari vector cirinya, baik pada learning maupun test.
- Hitung: $K(x_i, x_j)$ dan tidak perlu melakukan perhitungan dalam ruang berdimensi tinggi secara eksplisit. ((Kernel Trick))

- Kernel Functions:

- Polynomial Kernel with degree d
 $K(x, y) = (x^\top y + 1)^d$
- Radial Basic Function Kernel with degree σ
 $K(x, y) = \exp(-\|x - y\|^2 / (2\sigma^2))$
 : closely related to radial basis function neural networks

- Sigmoid with Parameter k dan θ
 $K(x, y) = \tanh(kx^\top y + \theta)$
 : it doesn't satisfy the Mercer condition on all k and θ
- Research on different kernel functions in different applications is very active.

- Model Evaluation and Selection -

- Evaluation Metrics: How can we measure accuracy?
 Other metrics to consider?
 - Use validation test set of class-labeled tuples instead of training set when assessing accuracy.
 - Methods for estimating a classifier's accuracy:
 1. Holdout method, random subsampling
 2. Cross-validation
 3. Bootstrap
 - Comparing classifiers:
 - Confidence intervals

Classifier Evaluation Metrics:

Confusion Matrix

		C1	$\neg C_1$
Actual \ Predicted class	C1	True Positives (TP)	False Negatives (FN)
	$\neg C_1$	False Positives (FP)	True Negatives (TN)

Accuracy, Error Rate, Sensitivity and Specificity

Actual \ Predicted		C	$\neg C$	
C	C	TP	FN	P
	$\neg C$	FP	TN	N
		P'	N'	All

- Classifier Accuracy / Recognition rate : percentage of test set tuples that are correctly classified.

$$\text{Accuracy} = (TP + TN) / \text{All}$$

- Error Rate: 1 - accuracy, or

$$\text{Error Rate} = (FP + FN) / \text{All}$$

- Class Imbalance Problem

- One class may be rare, eg: fraud, HIV-positive
- Significant majority of the negative class and minority of the positive class
- Sensitivity: True Positive recognition rate
 $\text{Sensitivity} = TP / P$
- Specificity: True Negative recognition rate
 $\text{Specificity} = TN / N$

Precision and Recall, F-measures

- Precision: exactness - what % of tuples that the classifier labeled as positive are actually positive

$$\text{Precision} = \frac{TP}{TP + FP}$$

- Recall: completeness - what % of positive tuples did the classifier label as positive

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Perfect score} = 1.0$$

- inverse relationship between precision & recall

F-measure (F_1 or F-score): harmonic mean of precision & recall

$$F = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

F_B : weighted measure of precision and recall
 assigns β times as much weight to recall as to precision

$$F_B = \frac{(1 + \beta^2) \times \text{precision} \times \text{recall}}{\beta^2 \times \text{precision} + \text{recall}}$$

~~Latinium~~:

Actual \ Predicted class	cancer=Yes	cancer=No	Total	Recognition (%)
Actual \ Predicted class	cancer=Yes	cancer=No	Total	Sensitivity
cancer=Yes	90	210	300	
cancer=No	140	9560	9700	Specificity
Total	230	9770	10000	Accuracy

$$\cdot \text{Sensitivity} = TP / P = 90 / 300 = 30\%$$

$$\cdot \text{Specificity} = TN / N = 9560 / 9700 = 98.56\%$$

$$\cdot \text{Accuracy} = (TP + TN) / \text{All} = (90 + 9560) / 10000 = 96.5\%$$

- Precision = $\frac{TP}{TP+FP} = \frac{90}{(90+140)} = 39.13\%$

- Recall = $\frac{TP}{TP+FN} = \frac{90}{90+210} = 30\%$

1. Holdout Method

- Given data is randomly partitioned into 2 independent sets
 - Training set (eg: 2/3) for model construction
 - Test set (eg: 1/3) for accuracy estimation
- Random Sampling: a variation of holdout
 - Repeat holdout k times,
 - accuracy = avg. of the accuracies obtained

2. Cross-Validation (k -fold, where $k=10$ is most popular)

- Randomly partition the data into k mutually exclusive subsets, each approximately equal size.
- At i -th iteration, use D_i as test set and others as training set.
- Leave-one-out: k folds where $k = \#$ of tuples, for small sized data.
- Stratified Cross-Validation: folds are stratified so that class dist. in each fold is approx. the same as that in the initial data.

3. Bootstrap

- Works well with small data sets
- Samples the given training tuples uniformly with replacement.
e.g.: each time a tuple is selected, it is equally likely to be selected again and re-added to the training set.

Estimating Confidence Intervals:

• Classifier Models M_1 vs M_2

- Suppose we have 2 classifiers, M_1 and M_2 , which one is better?
- Use 10-fold cross-validation to obtain $\bar{err}(M_1)$ & $\bar{err}(M_2)$
- These mean error rates are just estimates of error on the true population of future data cases.
- What if the difference between the 2 error rates is just attributed to chance?
 - Use a test of statistical significance
 - Obtain confidence limits for our error estimates

• Null Hypothesis

- Perform 10 fold cross-validation
- Assume samples follow a t distribution with $k-1$ degrees of freedom (here, $k=10$)
- Use t -test (or Student's t -test)
- Null Hypothesis: M_1 & M_2 are the same
- If we can reject null hypothesis, then:
 - We conclude that the difference between M_1 & M_2 is statistically significant.
 - Choose model with lower error rate.

Issues Affecting Model Selection

1. Accuracy
 - classifier accuracy: predicting class label
2. Speed
 - time to construct the model (training time)
 - time to use the model (classification/prediction time)
3. Robustness
 - handling noise and missing values
4. Scalability
 - efficiency in disk-resident databases
5. Interpretability
 - understanding and insight provided by the model
6. Other measures, e.g.: goodness of rules, such as decision tree size or compactness of classification rules.

Latihan

- 1. Binning
- 2. K-Means
- 3. LOF
- 4. GINI & Decision Tree
- 5. Hyperplane

6. Entropy.

sheet: hyperplane

ASOSIASI:

Konsep Dasar dan Algoritme

- Clustering: SSE kecil
- Klasifikasi: akurasi bagus / sensitivity
↳ tinggi ↳ for imbalance class

ARM (Association Rule Mining)

: diberikan himpunan transaksi, cari akurasi yang memprediksi kemunculan sebuah item berdasarkan kemunculan item lain dalam transaksi.

eg: Transaksi Market-Basket

TID	Items
1	Bread, Milk
2	Bread, Diaper, Juice, Eggs
3	Milk, Diaper, Juice, Coke
4	Bread, Milk, Diaper, Juice
5	Bread, Milk, Diaper, Coke

$$\{ \text{Diaper} \} \rightarrow \{ \text{Juice} \}$$

$$\{ \text{Milk, Bread} \} \rightarrow \{ \text{Eggs, Coke} \}$$

$$\{ \text{Juice, Bread} \} \rightarrow \{ \text{Milk} \}$$

Implikasi: kemunculan samaan, bukan hub. sebab akibat

Definition:

1. Itemset: sekumpulan 1 atau lebih items
eg: $\{ \text{Milk, Bread, Diaper} \}$
↳ k-itemset: itemset yg mengandung k items.
2. Support Count (σ): frekuensi kemunculan sebuah itemset. \rightarrow support absolute.
eg: $\sigma(\{ \text{Milk, Bread, Diaper} \}) = 2$.
3. Support: fraksi transaksi yg mengandung sebuah itemset. \rightarrow support relative.
eg: $s(\{ \text{Milk, Bread, Diaper} \}) = 2/5$.
4. Frequent Itemset: sebuah itemset dikatakan support jika \geq minsup threshold.
5. Association Rule: ekspresi implikasi dari bentuk $X \rightarrow Y$, dimana X dan Y = itemsets.
eg: $\{ \text{Milk, Diaper} \} \rightarrow \{ \text{Juice} \}$

6. Rule Evaluation Metrics

- Support (S): fraction of transactions that contain both X and Y .
- Confidence (C): measures how often items in Y appear in transactions that contain X .

$$\text{eg: } \{ \text{Milk, Diaper} \} \rightarrow \text{Juice}$$

$$S = \frac{\sigma(\text{Milk, Diaper, Juice})}{\sigma(\text{Milk, Diaper})} = \frac{2}{5} = 0.4$$

$$C = \frac{\sigma(\text{Milk, Diaper, Juice})}{\sigma(\text{Milk, Diaper})} = \frac{2}{3} = 0.67$$

Association Rule Mining Task

- given a set of transactions T
- goal: find all rules having
 - support \geq minsup threshold
 - confidence \geq minconf threshold

Brute-force approach:

1. List all possible association rule
2. Compute the support and confidence for each rule
3. Prune rules that fall the minsup and minconf thresholds.

Mining Association Rules

Two-step approach:

1. Frequent Itemset Generation

- generate all itemsets whose support \geq minsup

2. Rule Generation

- generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset.

1. Frequent Itemset Generation

Given d unique items:

$$\text{Total number of itemsets} = 2^d$$

$$\text{Total number of possible association rules} =$$

$$R = \sum_{k=1}^{d-1} \left[\binom{d}{k} \times \sum_{j=1}^{d-k} \binom{d-k}{j} \right] = 3^d - 2^{d+1} + 1$$

$$\text{eg: } d=6 \rightarrow R = 3^6 - 2^{6+1} + 1 = 602 \text{ rules}$$

Complexity $\sim O(NM^d) \rightarrow$ expensive since $M = 2^d$

Solusi:

1. Mengurangi jumlah kandidat (M)
 - Complete search: $M=2^d$
 - Use pruning techniques to reduce M
2. Mengurangi jumlah transaksi (N)
 - Reduce size N as the size of itemset increases
 - Use vertical-partitioning of the data to apply the mining algorithms
3. Mengurangi jumlah comparisons (NM)
 - Use efficient data structures to store the candidates or transactions
 - No need to match every candidate against every transaction.

Prinsip Apriori: jika itemset: frequent, maka semua subsets harus frequent.

$$\forall X, Y: (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

Method:

- Let $k=1$
- Generate frequent itemsets of length 1
- Repeat until no new frequent itemsets are identified
 - ↳ Generate length $(k+1)$ candidate itemsets from length k frequent itemsets.
 - ↳ Prune candidate itemsets containing subsets of length k that are infrequent.
 - ↳ Count the support of each candidate by scanning the DB
 - ↳ Eliminate candidates that are infrequent, leaving only those that are frequent.

eg: $\text{minsup} = 2$

TID	Items	C1	itemset	sup
10	A, C, D		{A}	2
20	B, C, E	1st scan	{B}	3
30	A, B, C, E		{C}	3
40	B, E		{D}	1
			{E}	3

L1	itemset	sup
	{A}	2
	{B}	3
	{C}	3
	{D}	1
	{E}	3

L2	itemset	sup
	{A, C}	2
	{B, C}	2
	{B, E}	3
	{C, E}	2

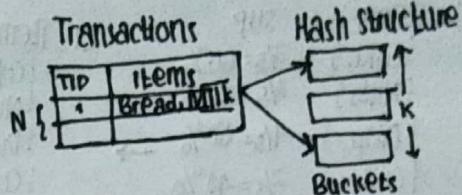
C2	itemset	sup
	{A, B}	1
	{A, C}	2
	{A, E}	1
	{B, C}	2
	{B, E}	3
	{C, E}	2

C3	itemset	sup
	{B, C, E}	2
L3	↓	
	{B, C, E}	2

Reducing Number of Comparisons

Candidate counting:

- Scan the database of transactions to determine the support of each candidate itemset.
- To reduce the number of comparisons, store the candidates in a hash structure.
 - Instead of matching each transaction against every candidate, match it against candidates contained in the hashed buckets.



Factors Affecting Complexity

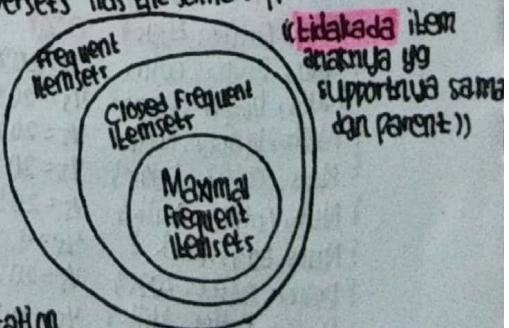
- Choice of minimum support threshold
 - semakin kecil support threshold, makin banyak frequent itemsets
 - dpt meningkatnya jumlah kandidat dan membutuhkan length of frequent itemset
- Dimensionality (number of items) of the dataset
 - butuh banyak space untuk menyimpan support count setiap item
 - jika jumlah frequent items naik, maka biaya komputasi dan I/O juga meningkat
- Size of database
 - jumlah transaksi meningkat, run time algoritme meningkat
- Average transaction width
 - transaction width increases with denser data sets.
 - menyebabkan peningkatan max length of frequent itemsets dan traversal of hash tree (number of subsets in a transaction increases with its width)

Maximal Frequent Itemset

: itemset dikatakan maximal frequent jika none of its immediate supersets is frequent.
((tidak ada item yg frequent dgn parentnya))

Closed Itemset

: itemset dikatakan closed jika none of its immediate supersets has the same support as the itemset.



Horizontal Data Layout

TID	Items
1	A, B, E
2	B, C, D
3	C, E
4	A, C, D
5	A, B, C, D
6	A, E
7	A, B
8	A, B, C
9	A, C, D
10	B

Representation of Database

A	B	C	D	E
1	1	2	1	
4	2	3	4	3
5	5	4	5	6
6	7	8	9	
7	8	9	9	
8	10			
9				

If $\text{minsup} = 20\%$, how = ...?

Apriori PR:

minsup = 20%

TID	Items
10	Coke, Nuts, Diaper
20	Coke, Coffee, Diaper
30	Coke, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk

C1	Itemsets	SUP
{Coke}	$\frac{3}{5} = 60\%$	
{Nuts}	$\frac{3}{5} = 60\%$	
{Diaper}	$\frac{4}{5} = 80\%$	\rightarrow
{Coffee}	$\frac{3}{5} = 60\%$	
{Eggs}	$\frac{3}{5} = 60\%$	
{Milk}	$\frac{2}{5} = 40\%$	

C2	Itemset	SUP
{Coke, Nuts}	$\frac{1}{5} = 20\%$	\rightarrow
{Coke, Diaper}	$\frac{3}{5} = 60\%$	
{Coke, Coffee}	$\frac{1}{5} = 20\%$	
{Coke, Eggs}	$\frac{1}{5} = 20\%$	
{Coke, Milk}	0%	
{Nuts, Diaper}	$\frac{2}{5} = 40\%$	
{Nuts, Coffee}	$\frac{1}{5} = 20\%$	
{Nuts, Eggs}	$\frac{2}{5} = 40\%$	
{Nuts, Milk}	$\frac{1}{5} = 20\%$	
{Diaper, Coffee}	$\frac{2}{5} = 40\%$	
{Diaper, Eggs}	$\frac{2}{5} = 40\%$	
{Diaper, Milk}	$\frac{1}{5} = 20\%$	
{Coffee, Eggs}	$\frac{1}{5} = 20\%$	
{Coffee, Milk}	$\frac{1}{5} = 20\%$	
{Eggs, Milk}	$\frac{2}{5} = 40\%$	

C3	Itemset	SUP
{Coke, Nuts, Diaper}	$\frac{1}{5} = 20\%$	\rightarrow
{Coke, Nuts, Coffee}	0%	
{Coke, Nuts, Eggs}	0%	
{Coke, Diaper, Coffee}	$\frac{1}{5} = 20\%$	
{Coke, Diaper, Eggs}	$\frac{1}{5} = 20\%$	
{Coke, Coffee, Eggs}	0%	
{Nuts, Diaper, Coffee}	$\frac{1}{5} = 20\%$	
{Nuts, Diaper, Eggs}	$\frac{1}{5} = 20\%$	
{Nuts, Diaper, Milk}	$\frac{1}{5} = 20\%$	
{Nuts, Coffee, Eggs}	$\frac{1}{5} = 20\%$	
{Nuts, Coffee, Milk}	$\frac{1}{5} = 20\%$	
{Nuts, Eggs, Milk}	$\frac{2}{5} = 40\%$	
{Diaper, Coffee, Eggs}	$\frac{1}{5} = 20\%$	
{Diaper, Coffee, Milk}	$\frac{1}{5} = 20\%$	
{Diaper, Eggs, Milk}	$\frac{1}{5} = 20\%$	
{Coffee, Eggs, Milk}	$\frac{1}{5} = 20\%$	

C4	Itemset	SUP
{Nuts, Diaper, Coffee, Eggs}	$\frac{1}{5} = 20\%$	
{Nuts, Diaper, Coffee, Milk}	$\frac{1}{5} = 20\%$	
{Nuts, Diaper, Coffee, Eggs, Milk}	$\frac{1}{5} = 20\%$	
{Nuts, Diaper, Eggs, Milk}	$\frac{1}{5} = 20\%$	
{Nuts, Coffee, Eggs, Milk}	$\frac{1}{5} = 20\%$	
{Diaper, Coffee, Eggs, Milk}	$\frac{1}{5} = 20\%$	

L4	Itemset	SUP
{Nuts, Diaper, Coffee, Eggs}	20%	
{Nuts, Diaper, Coffee, Milk}	20%	
{Nuts, Diaper, Coffee, Eggs, Milk}	20%	
{Nuts, Diaper, Eggs, Milk}	20%	
{Nuts, Coffee, Eggs, Milk}	20%	
{Diaper, Coffee, Eggs, Milk}	20%	

C5	Itemset	SUP
{Nuts, Diaper, Coffee, Eggs, Milk}	$\frac{1}{5} = 20\%$	

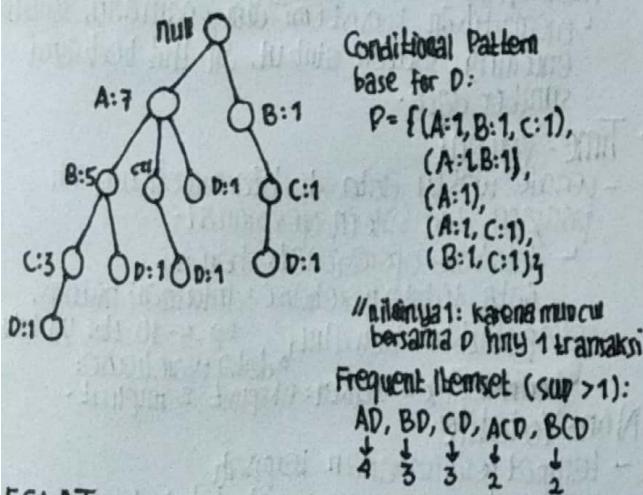
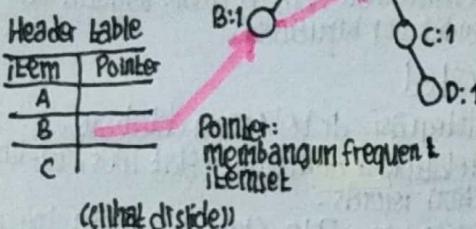
L5	Itemset	SUP
{Nuts, Diaper, Coffee, Eggs, Milk}	20%	

FP-Growth

- menggunakan FP-tree stg representasi terkompresi dr database.
- buat tree → pake pendekatan divide and conquer utk dapetin frequent & infrequent itemset.
- Caranya:
 1. Baca yg transaksi
 2. Mulai dg node null, lalu buat node item baru. Jika item sudah ada, pilih yg sama.
 3. Node item akan bertambah kemunculannya.
 4. Bikin pointer.

eg: TID Items
 1 A,B
 2 B,C,D
 3 A,C,D,E
 4 A,D,E
 5 A,B,C
 6 A,B,C,D
 7 B,C
 8 A,B,C
 9 A,B,D
 10 B,C,E

read: TID=1
 null
 A:1
 B:1
 read: TID=2
 null
 A:1
 B:1
 C:1
 D:1



ECLAT
 (lihat di slide)

Pembangunan Rule

- itemset frequent: L
- find seluruh subset yg tidak kosong fcl?

eg: {A, B, C, D} adalah frequent itemset,
 kandidat rules:

$ABC \rightarrow D, ABD \rightarrow C, ACD \rightarrow B, BCD \rightarrow A$
 $A \rightarrow BCD, B \rightarrow ACD, C \rightarrow ABD, D \rightarrow ABC$
 $AB \rightarrow CD, AC \rightarrow BD, AD \rightarrow BC, BC \rightarrow AD,$
 $BD \rightarrow AC, CD \rightarrow AB$

- if $|L| = k$, maka ada $2^k - 2$ kandidat association rules (dn mengabaikan $L \rightarrow \emptyset$ dan $\emptyset \rightarrow L$)

Pruning Berbasis Confidence

if $X \rightarrow Y$ tidak memenuhi confidence threshold, maka sembarang rule $X' \rightarrow Y - X'$, dimana X' = subset dari X , juga hrs bdk dpt memenuhi confidence threshold.

$$c(R_1) = \frac{|AB \cup CD|}{|AB|}$$

$$c(R_2) = \frac{|ABC \cup D|}{|ABC|}$$

$c(R_1) < c(R_2)$
 R_2 banyak transaksi tidak memenuhi, apalagi R_1 (transaksi).

- Confidence: tidak memiliki sifat anti-monotone

eg: $c(ABC \rightarrow D)$ bisa $>$ atau $<$ dari $c(AB \rightarrow D)$

Tapi, Confidence rules dr itemset sama punya sifat anti-monotone.

eg: $L = \{A, B, C, D\}$

$$c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD)$$

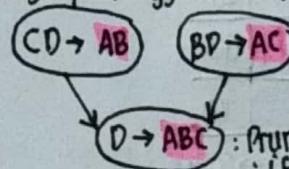
$$\begin{array}{|c|c|} \hline \text{TID} & \text{Items} \\ \hline 10 & ABCD \\ 20 & BCE \\ 30 & ABCE \\ 40 & ABE \\ \hline \end{array} \quad c(ABC \rightarrow D) = \frac{\sigma(ABC \cup D)}{\sigma(ABC)} = \frac{1}{2} = 0.5$$

$$c(AB \rightarrow CD) = \frac{\sigma(AB \cup CD)}{\sigma(AB)} = \frac{1}{3} = 0.33$$

$$c(A \rightarrow BCD) = \frac{\sigma(A \cup BCD)}{\sigma(A)} = \frac{1}{3} = 0.33$$

Pembangunan Rule: Apriori

- kandidat rule dibangun dg cara menggabungkan 2 rules yg punya prefix yg sama dibagian rule consequent.



: Prune rule $D \rightarrow ABC$
 : if subset $AD \rightarrow BC$ tidak punya high confidence.

Effect of Support Distribution

- if minsup = tinggi \rightarrow interesting item kadang tidak muncul, eg. produk harga mahal
 - If minsup = kecil \rightarrow komputasi tinggi, jmlh itemset besar
- Solusi: multiple minimum supports
 beda \times tiap barang.

eg: MS(i): minimum support for item i

$$\begin{aligned} MS(\text{Milk}) &= 5\%, MS(\text{Coke}) = 3\%, \\ MS(\text{Broccoli}) &= 0.1\%, MS(\text{Salmon}) = 0.5\% \end{aligned}$$

$$\begin{aligned} MS(\{\text{Milk}, \text{Broccoli}\}) &= \min(MS(\text{Milk}), \\ &\quad MS(\text{Broccoli})) \\ &= \min(5\%, 0.1\%) \\ &= 0.1\% \end{aligned}$$

Menghitung Interestingness Measure

Tabel Kontingensi $X \rightarrow Y$

$X \bar{Y}$	f_{11}	f_{10}	f_{1+}	f_{11} : support of X dan Y	
$\bar{X} Y$	f_{01}	f_{00}	f_{0+}	f_{01} : support of X dan \bar{Y}	
$\bar{X} \bar{Y}$	f_{+1}	f_{+0}	f_{+1} : support of \bar{X} dan Y	ukt menentukan various measures:	
				support, confidence, lift, gini, J-measure, etc.	

Statistical Independence

- Population of 1000 students
- 600 students : swim (S)
- 700 students : bike (B)
- 420 students : swim and bike (S, B)

$$\rightarrow P(S \cap B) = \frac{420}{1000} = 0.42,$$

$$P(S) \times P(B) = \frac{600}{1000} \times \frac{700}{1000} = 0.42,$$

Sehingga d:

$P(S \cap B) = P(S) \times P(B)$: statistical independence

$P(S \cap B) > P(S) \times P(B)$: korelasi positif

$P(S \cap B) < P(S) \times P(B)$: korelasi negatif

$$\cdot \text{Lift} = \frac{P(Y|X)}{P(Y)} \quad \cdot \text{Interest} = \frac{P(X, Y)}{P(X)P(Y)}$$

$$\cdot PS = P(X, Y) - P(X)P(Y)$$

$$\cdot \phi\text{-coefficient} = \frac{P(X, Y) - P(X)P(Y)}{\sqrt{P(X)[1-P(X)]P(Y)[1-P(Y)]}}$$

eg:

	Coffee	Coffee	
Tea	15	5	20
Tea	75	5	80
	90	10	100

Association Rule: Tea \rightarrow Coffee

$$\cdot \text{Confidence} = P(\text{Coffee} | \text{Tea}) \\ = P(\text{Coffee} \wedge \text{Tea})$$

$$= \frac{15}{20} = 0.75$$

$$\cdot P(\text{coffee}) = \frac{90}{100} = 0.9$$

$$\text{maka, lift} = \frac{0.75}{0.9} = 0.8333 (< 1, \text{asosiasi negatif})$$

and On-line Analytical Processing

Data Warehouse

- dikelola terpisah dari database operasional
- memrosesari informasi dgn menyediakan data historis untuk analisis.
- merupakan koleksi data yg punya karakteristik: subject-oriented, integrated, time-variant, non volatile; yg digunakan utk pengambilan keputusan.
- Data Warehousing: proses membangun dan mengg. data warehouse.

1. Subject-Oriented

- diorganisasi menurut subjek penting, eg: customer, product, sales.
- fokus: pemodelan dan analisis data utk pembuatan keputusan, bukan utk operasional harian/pemrosesan transaksi.
- tampilan simpel dan padat thd subjek dgn mengabarkan data yg tidak berguna utk proses pendukung keputusan.

2. Integrated

- berintegrasi dr beberapa database.
 - relational database, flat files, on-line transaction records.
- menerapkan: Data Cleaning & Data Integration
 - memastikan konsistensi dlm penamaan, struktur encoding, satuan atribut, dll. thd berbagai sumber data.

3. Time-Variant

- periode waktu data di data warehouse lebih panjang dari sistem operasional.
 - database operasional: hari ini
 - data di data warehouse: informasi historis, setiap struktur keg dlm, eg: 5-10 thn yg lalu.
 - berisi elemen waktu: eksplisit & implisit.

4. Non-Volatile

- tempat penyimpanan terpisah
- perbaruan data operasional tidak dilakukan update data, recovery data; yg ada hanya 2 operasi akses data: initial loading of data, access of data.

Data Warehouse terpisah:

1. kinerja tinggi utk kedua sistem:

DBMS - droptimasi utk OLTP
Warehouse - droptimasi utk OLAP

2. fungsi yg berbeda dari data yang kerjanya

- Data yg tidak dtko: perlu data historis yg tidak diketahui dlm DB operasional.
- konsolidasi data: dr sumber data yg berbeda.
- kualitas data: format yg tidak sama / mungkin berbeda, sisa partisi resorialisasi.

A Multi-dimensional data model

From Tables & Spreadsheets → Data Cubes

↳ multidimensional data model

- Dimension tables: item (item.name, brand, type), time (day, week, month, quarter, year)

- Fact Table: berisi measures dan key ke setiap dimensi tabel yg berhub.

Modeling data warehouses:
dimensions & measures

1. Star schema: tabel fakta di tengah yg terhubung sejumlah tabel dimensi.
2. Snowflake schema: perbaikan star schema, tabel fakta 1, terhubung dgn tabel dimensi yg punya anak.
3. Fact constellations: beberapa tabel fakta yg (kolaksi bintang) berbagi tabel dimensi. (galaxy)

Measures of Data Cube:

three categories

1. Distributive: if hasil turunan dr n aggregate values = turunan dr fungsi tanpa partitioning.

eg: count(), sum(), min(), max()

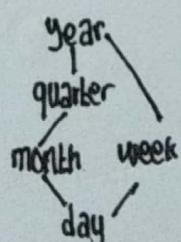
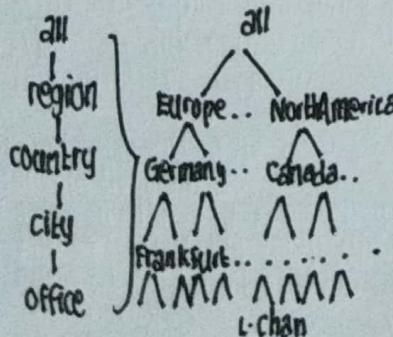
2. Algebraic: if dpt dihitung mengg. algebraic funct dgn M argumen (M = bil. int yg sgt besar), setiap perhitungan mengg. distributive aggregate funct.

eg: avg(), min(), std-deviation()

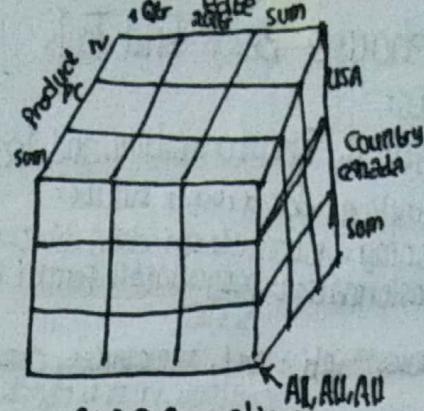
3. Holistic: if there is no constant bound on the storage size needed to describe a subaggregate.

eg: median(), mode(), rank()

A Concept Hierarchy:
dimension (location)



A Sample Data Cube
(lihat slide))

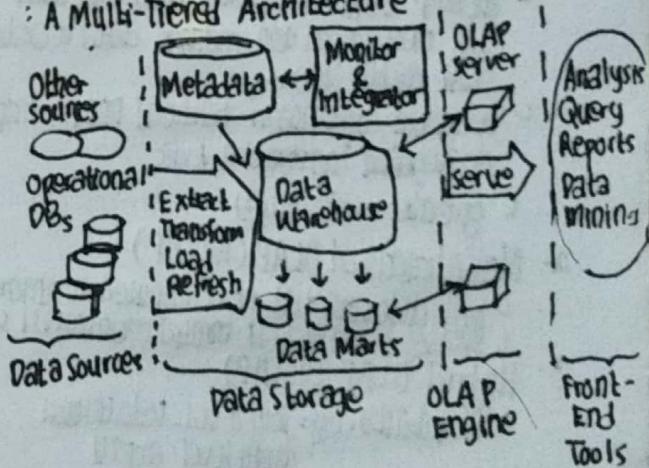


Typical OLAP Operations

1. Roll up (drill-up): summarize data by climbing up hierarchy or by dimension reduction
2. Drill down (roll-down): reverse of roll-up from higher level summary to lower level summary or detailed data, or introducing new dimensions.
3. Slice and dice: project and select.
4. Pivot (rotate): reorient the cube, visualization, 3D to series of 2D planes.
5. Other operations:
 - drill across: involving (across) more than one fact table.
 - drill through: through the bottom level of the cube to its back-end relational tables (using SQL)

Data Warehouse Architecture

A Multi-Tiered Architecture



Three Data Warehouse Models

1. Enterprise warehouse: mencakup seluruh organisasi.
2. Data Mart: kelompok pengguna berkemiru / terbatas.
eg: marketing data mart.
3. Virtual warehouse: sekumpulan view dari hasil query ult database operasional.



Data Warehouse Back-End Tools and Utilities:

ETL (Extraction, Transformation, and Learning)

1. Data extraction: dari berbagai sumber
2. Data cleaning: inkonsistensi data diberikan.
3. Data transformation: penyesuaian format agar serupa.
4. Load: proses awal: sort, summarize, consolidate, compute views, check integrity, dan build indices and partitions.
5. Refresh: update data sources ke warehouse.

2. Join Indices

From Data Warehousing to Data Mining

Data Warehouse Usage

1. Information processing
2. Analytical processing
3. Data mining

Metadata Repository

: data yg mendefinisikan warehouse objek. Berisi:

1. deskripsi struktur data warehouse: schema, view, dimensions, hierarchies, data definition, turunan, lokasi, isi data mart.
2. Operational meta-data: siapa saja yang make, kapan saja, pake cara apa
3. Algoritme utk summarization
4. Mapping

OLAP Server Architectures

1. Relational OLAP (ROLAP)

- ↳ mengg. DBMS relational atau extended-relational utk menyimpan dan manage data warehouse dan OLAP middle ware.
- ↳ termasuk optimisasi backend DBMS, implementasi aggregation navigation logic
- ↳ greater scalability.

2. Multidimensional OLAP (MOLAP)

- ↳ fast indexing utk pre-computed summarized data.
- ↳ sparse array-based multidimensional storage engine.

3. Hybrid OLAP (HOLAP)

- ↳ flexibility. e.g. low level: relational
high level: array
- ↳ e.g. Microsoft SQL Server

4. Specialized SQL Servers

- ↳ specialized support for SQL queries over star/
Snowflake schemas.

↳ e.g. Redbricks.

Data Warehouse Implementation

Indexing OLAP Data:

1. Bitmap Index

e.g.: Base Table

Cust-Region	Type	Index on Region			Index on Type			
		RecID	Asia	Europe	America	RecID	Retail	Dealer
C1	Asia	Retail	1	1	0	0	1	1
C2	Europe	Dealer	2	0	1	0	2	0
C3	Asia	Dealer	3	1	0	0	3	0
C4	America	Retail	4	0	0	1	4	1

GINI

C1: Play tennis
C2: no play

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	cool	Normal	Weak	Yes
6	Rain	cool	Normal	Strong	No
7	Overcast	cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Outlook

Sunny	Overcast	Rain
C1: 2 C2: 3	C1: 4 C2: 0	C1: 3 C2: 2

$$\text{GINI}(\text{Sunny}) = 1 - P(C1)^2 - P(C2)^2 = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 0.48$$

$$\text{GINI}(\text{Overcast}) = 1 - (4/4)^2 - (0/4)^2 = 1 - 1^2 - 0^2 = 0$$

$$\text{GINI}(\text{Rain}) = 1 - (3/5)^2 - (2/5)^2 = 0.48$$

$$\text{GINI}(\text{outlook}) = \frac{5}{14} \cdot 0.48 + \frac{4}{14} \cdot 0 + \frac{5}{14} \cdot 0.48 = 0.34 \quad (\text{Root Node})$$

Temperature

Hot	Mild	cool
C1: 2 C2: 2	C1: 4 C2: 2	C1: 3 C2: 1

$$\text{GINI}(\text{Hot}) = 1 - (2/4)^2 - (2/4)^2 = 0.5$$

$$\text{GINI}(\text{Mild}) = 1 - (4/6)^2 - (2/6)^2 = 0.44$$

$$\text{GINI}(\text{cool}) = 1 - (3/4)^2 - (1/4)^2 = 0.375$$

$$\text{GINI}(\text{Temperature}) = \frac{4}{14} \cdot 0.5 + \frac{6}{14} \cdot 0.44 + \frac{4}{14} \cdot 0.375 = 0.44$$

Humidity

High	Normal
C1: 3 C2: 4	C1: 6 C2: 1

$$\text{GINI}(\text{High}) = 1 - (3/7)^2 - (4/7)^2 = 0.49$$

$$\text{GINI}(\text{Normal}) = 1 - (6/7)^2 - (1/7)^2 = 0.25$$

$$\text{GINI}(\text{Humidity}) = \frac{7}{14} \cdot 0.49 + \frac{7}{14} \cdot 0.25 = 0.37$$

Wind

Weak	Strong
C1: 6 C2: 2	C1: 3 C2: 3

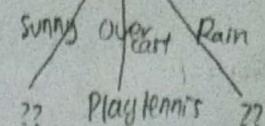
$$\text{GINI}(\text{Weak}) = 1 - (5/8)^2 - (3/8)^2 = 0.375$$

$$\text{GINI}(\text{Strong}) = 1 - (3/6)^2 - (3/6)^2 = 0.5$$

$$\text{GINI}(\text{Wind}) = \frac{8}{14} \cdot 0.375 + \frac{6}{14} \cdot 0.5 = 0.43$$

C1: Play tennis
C2: no play

Outlook



?? Play tennis ??

→ ((sunny, temperature))

(sunny, hot) (funny, mild)

C1: 0
C2: 2

C1: 1
C2: 1

C1: 1
C2: 0

$$\text{Gain}(\text{sunny}, \text{hot}) = 1 - (0/2)^2 - (2/2)^2 = 0$$

$$\text{Gain}(\text{funny}, \text{mild}) = 1 - (1/1)^2 - (1/1)^2 = 0.5$$

$$\text{Gain}(\text{funny}, \text{cool}) = 1 - (1/1)^2 - (0/1)^2 = 0$$

$$\text{GAIN}(\text{sunny}, \text{temperature}) = \frac{2}{5}(0) + \frac{3}{5}(0.5) + \frac{1}{5}(0) = 0.2$$

→ ((funny, humidity))

(funny, high) (funny, normal)

C1: 0
C2: 3

C1: 2
C2: 0

$$\text{Gain}(\text{funny}, \text{high}) = 1 - (0/3)^2 - (3/3)^2 = 0$$

$$\text{Gain}(\text{funny}, \text{normal}) = 1 - (2/2)^2 - (0/2)^2 = 0$$

$$\text{GAIN}(\text{funny}, \text{humidity}) = \frac{3}{5}(0) + \frac{2}{5}(0) = 0 \quad (\text{leaf node for funny})$$

→ ((funny, wind))

(funny, weak) (funny, strong)

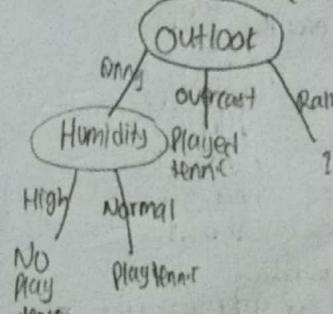
C1: 1
C2: 2

C1: 1
C2: 1

$$\text{Gain}(\text{funny}, \text{weak}) = 1 - (1/3)^2 - (2/3)^2 = 0.44$$

$$\text{Gain}(\text{funny}, \text{strong}) = 1 - (1/2)^2 - (1/2)^2 = 0.5$$

$$\text{GAIN}(\text{funny}, \text{wind}) = \frac{3}{5}(0.44) + \frac{2}{5}(0.5) = 0.464$$



→ ((rain, temperature))

(rain, hot) (rain, mild)

C1: 0
C2: 0

C1: 2
C2: 1

C1: 1
C2: 1

$$\text{GINI}(\text{rain}, \text{hot}) = 1 - (0/0)^2 - (0/0)^2 = 1$$

$$\text{GINI}(\text{rain}, \text{mild}) = 1 - (2/3)^2 - (1/3)^2 = 0.44$$

$$\text{GINI}(\text{rain}, \text{cool}) = 1 - (1/1)^2 - (1/1)^2 = 0.5$$

$$\text{GAIN}(\text{rain}, \text{temperature}) = \frac{1}{5}(1) + \frac{3}{5}(0.44) + \frac{1}{5}(0.5) = 0.464$$

→ ((rain, wind))

(rain, weak) (rain, strong)

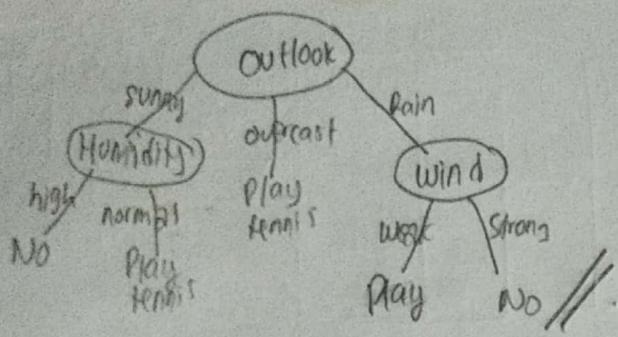
C1: 3
C2: 0

C1: 0
C2: 2

$$\text{GINI}(\text{rain}, \text{weak}) = 1 - (3/3)^2 - (0/3)^2 = 0$$

$$\text{GINI}(\text{rain}, \text{strong}) = 1 - (0/2)^2 - (2/2)^2 = 0$$

$$\text{GAIN}(\text{rain}, \text{wind}) = \frac{3}{5}(0) + \frac{2}{5}(0) = 0 \quad (\text{leaf node Rain})$$



Data Mining

#8 Analisis Asosiasi

- ARM: diberikan himpunan transaksi, cari aliran yg memprediksi kemunculan sebuah item berdasarkan kemunculan item lain dalam transaksi.
- e.g.: Transaksi Market Basket

TID	Items
1	Bread, Milk

• Definition:

1. Itemset: sekumpulan 1/ lebih items. e.g.: {Milk, Bread, Diaper's}
2. k-itemset: itemset yg mengandung k items.
3. Support Count (S): frekuensi kemunculan sebuah itemset. (support absolute).
e.g.: $S(\{\text{Milk, Bread, Diaper}\}) = 2$
4. Support (s): fraksi transaksi yg mengandung sebuah itemset. (support relative).
e.g.: $s(\{\text{Milk, Bread, ...}\}) = \frac{2}{5} = 0.4$
5. Frequent Itemset: sebuah itemset dikatakan support jika $\geq m_{\text{minsup}}$ threshold.
6. Association Rule: ekspresi implikasi dr bentuk $X \rightarrow Y$, dimana X dan Y = itemsets.
e.g.: {Milk, Diaper} \rightarrow {Juice}.

6. Rule Evaluation Metrics

- Support (S): fraction of transactions that contain both X and Y.
- Confidence (C): measures how often Items in Y appear in transactions that contain X.
e.g.: $\{ \text{Milk, Diaper} \} \Rightarrow \{ \text{Juice} \}$

$$S = \frac{S(\{\text{Milk, Diaper, Juice}\})}{S(\{\text{Milk, Diaper}\})} = \frac{2}{5} = 0.4 \quad C = \frac{S(\{\text{Milk, Diaper, Juice}\})}{S(\{\text{Milk, Diaper}\})} = \frac{2}{3} = 0.67$$

• Association Rule Mining Task

- * find all rules having: support $\geq m_{\text{minsup}}$, confidence $\geq m_{\text{minconf}}$
- * Brute-force approach:
 1. List all possible association rule.
 2. Compute the support and confidence for each rule.
 3. Prune rules that fall the minsup & minconf thresholds.

1. Frequent Itemset Generation

Given d unique items:

- Total number of itemsets = 2^d
- Total number of possible association rules:
$$R = \sum_{k=1}^{d-1} \left[\binom{d}{k} \times \sum_{j=1}^{d-k} \binom{d-k}{j} \right] = 3^d - 2^{d+1} + 1$$

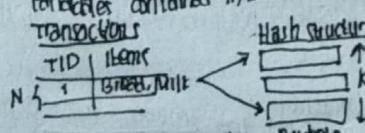
e.g.: $d=6 \rightarrow R = 3^6 - 2^{6+1} + 1 = 602$ rules.

Solusi:

1. Mengurangi jumlah kandidat (M)
 - complete search: $M = 2^d$
 - use mining technique to reduce M.
2. Mengurangi jumlah transaksi (N)
 - Reduce size N as the size of itemset increase.
 - Use vertical partitioning of the data to apply the mining algorithms.
3. Mengurangi jumlah comparisons (NM).
 - Use efficient data structures to store the candidate or transactions.
 - No need to match every candidate against every transaction.

• Reducing Number of Comparisons

1. scan the database of transactions to determine the support of each candidate itemset.
2. To reduce the number of comparisons, store the candidates in a hash structure.
 - instead of matching each transaction against every candidate, match it against candidates contained in the hashed buckets.



• Factors Affecting Complexity

1. choice of minimum support threshold
 - pilih nilai support threshold, makin kecil support threshold, makin banyak frequent itemset.
 - apt meningkatnya jumlah kandidat dan memakan banyak space.
2. Dimensionality (number of items) of the dataset.
 - butuh banyak space utk menghitung support count setiap item.
 - Jika banyak frequent items maka banyak komputasi dan I/O meningkat.
3. Size of database
 - jumlah transaksi meningkat, runtime algoritme meningkat.
4. Average transaction width
 - transaction width increases with denser datasets.
 - menyebabkan peningkatan max length of frequent itemsets dan traversals of hash tree (jumlah subsets in a transaction increases with its width).

• Mining Association Rules

two step approach:

1. Frequent Itemset Generation
 - generate all itemsets whose support $\geq m_{\text{minsup}}$
2. Rule Generation
 - generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset.

Prinsip Apriori: jika itemset frequent, maka semua subsets harus frequent.

$\forall X, Y: (X \subseteq Y) \Rightarrow S(X) \geq S(Y)$

Method:

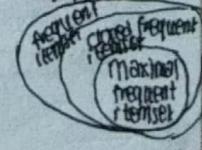
- Let $k=1$
- Generate frequent itemsets of length 1
- Repeat until no new frequent itemsets are identified:
 - Generate length $(k+1)$ candidate itemsets from length k frequent itemsets.
 - Prune candidate itemsets containing subsets of length k that are infrequent.
 - Count the support of each candidate by scanning the DB.
 - Eliminate candidates that are infrequent, leaving only those that are frequent.

• Maximal Frequent Itemset

- jika none of its immediate supersets is frequent.
- cl. tidak ada item frequent dih parentnya

• Closed Itemset

- jika none of its immediate supersets has the same support as the itemset.
- cl. tidak ada item yang supportnya sama dengan parentnya



Representasi of Database

1. Horizontal Data Layout

TID	Itemset
1	A, B, E
2	B, C, D
3	C, E
4	A, C, D
5	A, B, C
6	B, E
7	A, B, C, E
8	B, C, D, E
9	A, B, C, D, E

eg: minsup=2

TID	Items
1	A, B, E
2	B, C, D
3	C, E
4	A, C, D
5	A, B, C
6	B, E
7	A, B, C, E
8	B, C, D, E
9	A, B, C, D, E

minsup = 5

2. Vertikal Data Layout

	A	B	C	D	E
1	1	1	2	2	1
2	4	2	3	4	3
3	5	5	4	5	6
4	6	7	8	9	
5	8	9			
6	9				

minsup=2

itemset sup

15

itemset sup

5

itemset sup

2

itemset sup

2

itemset sup

3

itemset sup

3

itemset sup

2

itemset sup

2

itemset sup

3

itemset sup

3

#9. FP-Growth

- menggunakan FP-tree stg representasi berkompleks dr database.
- buat tree → proses pendekatan divide and conquer utk dptkan frequent & infrequent itemset.

- Casanya:

1. Baca per transaksi
2. Mulai dg node null, lalu buat node item baru. Jika item sudah ada, pakai.
3. Node item akan bertambah kemunculan itemnya.
4. Buat pointer.

eg: TID Items

read: TID=1

null

read TID=2

null

A:1

B:0

A:1

B:1

C:1

B:1

A:7

B:1

C:1

D:1

Footer: membandingkan frequent itemset.

HeaderTable

ITEM

Pointer

A

B

C

D

E

F

G

H

I

J

K

L

M

N

O

P

Q

R

S

T

U

V

W

X

Y

Z

ECLAT

- vertikal data format: disusun berdasarkan item.

eg: Itemset TID

1 {T100, T400, ..., 3} atau Item TID

2 {T100, T200, ..., 3}

3 {T100, T200, ..., 3}

4 {T100, T200, ..., 3}

5 {T100, T200, ..., 3}

6 {T100, T200, ..., 3}

7 {T100, T200, ..., 3}

8 {T100, T200, ..., 3}

9 {T100, T200, ..., 3}

10 {T100, T200, ..., 3}

11 {T100, T200, ..., 3}

12 {T100, T200, ..., 3}

13 {T100, T200, ..., 3}

14 {T100, T200, ..., 3}

15 {T100, T200, ..., 3}

16 {T100, T200, ..., 3}

17 {T100, T200, ..., 3}

18 {T100, T200, ..., 3}

19 {T100, T200, ..., 3}

20 {T100, T200, ..., 3}

21 {T100, T200, ..., 3}

22 {T100, T200, ..., 3}

23 {T100, T200, ..., 3}

24 {T100, T200, ..., 3}

25 {T100, T200, ..., 3}

26 {T100, T200, ..., 3}

27 {T100, T200, ..., 3}

28 {T100, T200, ..., 3}

29 {T100, T200, ..., 3}

30 {T100, T200, ..., 3}

31 {T100, T200, ..., 3}

32 {T100, T200, ..., 3}

33 {T100, T200, ..., 3}

34 {T100, T200, ..., 3}

35 {T100, T200, ..., 3}

36 {T100, T200, ..., 3}

37 {T100, T200, ..., 3}

38 {T100, T200, ..., 3}

39 {T100, T200, ..., 3}

40 {T100, T200, ..., 3}

41 {T100, T200, ..., 3}

42 {T100, T200, ..., 3}

43 {T100, T200, ..., 3}

44 {T100, T200, ..., 3}

45 {T100, T200, ..., 3}

46 {T100, T200, ..., 3}

47 {T100, T200, ..., 3}

48 {T100, T200, ..., 3}

49 {T100, T200, ..., 3}

50 {T100, T200, ..., 3}

51 {T100, T200, ..., 3}

52 {T100, T200, ..., 3}

53 {T100, T200, ..., 3}

54 {T100, T200, ..., 3}

55 {T100, T200, ..., 3}

56 {T100, T200, ..., 3}

57 {T100, T200, ..., 3}

58 {T100, T200, ..., 3}

59 {T100, T200, ..., 3}

60 {T100, T200, ..., 3}

61 {T100, T200, ..., 3}

62 {T100, T200, ..., 3}

63 {T100, T200, ..., 3}

64 {T100, T200, ..., 3}

65 {T100, T200, ..., 3}

66 {T100, T200, ..., 3}

67 {T100, T200, ..., 3}

68 {T100, T200, ..., 3}

69 {T100, T200, ..., 3}

70 {T100, T200, ..., 3}

71 {T100, T200, ..., 3}

72 {T100, T200, ..., 3}

73 {T100, T200, ..., 3}

74 {T100, T200, ..., 3}

75 {T100, T200, ..., 3}

76 {T100, T200, ..., 3}

77 {T100, T200, ..., 3}

78 {T100, T200, ..., 3}

79 {T100, T200, ..., 3}

80 {T100, T200, ..., 3}

81 {T100, T200, ..., 3}

82 {T100, T200, ..., 3}

83 {T100, T200, ..., 3}

84 {T100, T200, ..., 3}

85 {T100, T200, ..., 3}

86 {T100, T200, ..., 3}

87 {T100, T200, ..., 3}

88 {T100, T200, ..., 3}

89 {T100, T200, ..., 3}

90 {T100, T200, ..., 3}

91 {T100, T200, ..., 3}

92 {T100, T200, ..., 3}

93 {T100, T200, ..., 3}

94 {T100, T200, ..., 3}

95 {T100, T200, ..., 3}

96 {T100, T200, ..., 3}

97 {T100, T200, ..., 3}

98 {T100, T200, ..., 3}

99 {T100, T200, ..., 3}

100 {T100, T200, ..., 3}

101 {T100, T200, ..., 3}

102 {T100, T200, ..., 3}

103 {T100, T200, ..., 3}

104 {T100, T200, ..., 3}

105 {T100, T200, ..., 3}

106 {T100, T200, ..., 3}

107 {T100, T200, ..., 3}

108 {T100, T200, ..., 3}

109 {T100, T200, ..., 3}

110 {T100, T200, ..., 3}

111 {T100, T200, ..., 3}

112 {T100, T200, ..., 3}

113 {T100, T200, ..., 3}

114 {T100, T200, ..., 3}

115 {T100, T200, ..., 3}

116 {T100, T200, ..., 3}

117 {T100, T200, ..., 3}

118 {T100, T200, ..., 3}

119 {T100, T200, ..., 3}

120 {T100, T200, ..., 3}

121 {T100, T200, ..., 3}

122 {T100, T200, ..., 3}

123 {T100, T200, ..., 3}

124 {T100, T200, ..., 3}

125 {T100, T200, ..., 3}

126 {T100, T200, ..., 3}

127 {T100, T200, ..., 3}

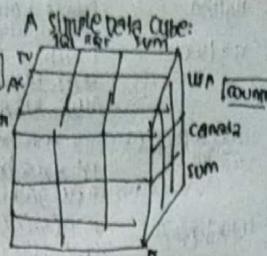
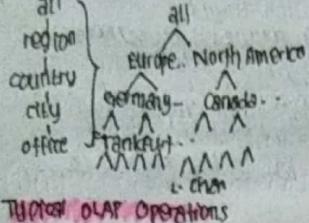
128 {T100

Measures of Data Cube: 3 categories

1. **Distributive:** if have function dr n aggregate values = function dr fungsi tanpa perhitungan.
eg: count(), sum(), min(), max()
2. **Algebraic:** if the distribution means algebraic funcn dan M arguments ($M = b/n$ int yg sdh brs).
sehingga perhitungan mengggunakan fungsi distributive aggregate fungn.
3. **Holistic:** if there no constant bound on the storage size needed to describe a subcube.
eg: median(), mode(), rank()

A Concept Hierarchy:

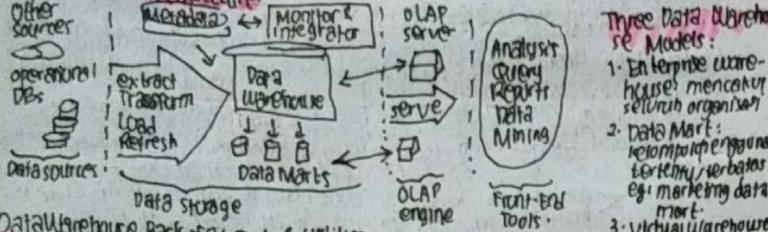
dimension (location)



Typical OLAP Operations

1. Roll up (drill-up): summarize data; dgn naik ke h/rank lbt tinggi / pengurangan dimensi.
2. Drill down (roll-down): reverse of roll up; from higher level summary to lower level summary or detailed data, or introducing new dimensions.
3. Slice and dice: project and select.
4. Pivot (rotate): mengambil kolom.
5. Other operations: - drill across: involving (crosses) more than 1 fact table.
- drill through: through the bottom level of the cube to its data warehouse Architecture. back end relational tables (using SQL).

Multi-Tiered Architecture



Data Warehouse Back-end Tools & Utilities:

1. Data extraction: dr berbagai sumber
2. Data cleaning: konsistensi data di berbagai sumber
3. Data transformation: pola/struktur format berbagai sumber
4. Load: proses ini: sort, summarize, consolidate, compute views, check integrity, dan build indexes and partitions.
5. Refresh: update dr drg sumber ke warehouse.

Metadatabase Repository

1. deskripsi struktur data warehouse: schema, view, dimensions, hierarchies, data
2. operational meta-data: definisi, turunan, lokasi, dsb data mart.
3. algoritme utk summarization
4. Mapping

OLAP Server Architectures

1. Relational OLAP (ROLAP)
 2. Multidimensional OLAP (MOLAP)
 3. Hybrid OLAP (HOLAP)
- menggabungkan DBMS relational atau extended-relational utk menyimpan & manage data warehouse dan OLAP middle ware.
 - memiliki optimasi backend DBMS, implementasi aggregation navigation logic.
 - fast indexing utk pre-computed summarized data.
 - sparse array-based multidimensional storage engine.
 - flexibility: eg: low level: relational
high level: array
 - eg: Microsoft SQL Server
4. Specialized SQL Servers: specialized support for SQL query over star / snowflake schemas.

Data Warehouse Implementation

Indexing OLAP Data:

1. Bitmap Index

eg: Base Table

Cust	Region	Type	Index on Region			Index on Type		
			RegID	Asia	Europe	America	Retail	Deale
C1	Asia	Retailer	1	0	0	0	1	0
C2	Europe	Dealer	2	0	1	0	2	0
C3	Asia	Dealer	3	1	0	0	3	0
C4	America	Retailer	4	0	0	1	4	1

2. Join Indices

From Data Warehousing to Data Mining

Data warehouse Usage:

1. Information processing
2. Analytical processing
3. Data Mining

#11. Spatiotemporal Data Mining

1. **Data spatial:** generalisasi titik geografis yg dptan mjl area terkait ada pengukuran sejumlah area geografis mengg. operasi spacial.
2. **Spatial Database Systems:** utk mengebara data spatial.
- Image db system: handling digital raster image (eg: satellite sensing, computer-photography), teknik analisis objek & ekstraksi dr image dan some spatial db functionality.
3. **Spatial (geometric, geographic) db system:** handling objek yg punya identity and well-defined extents, locations, and relationship (cid, ukuran, & hub antar data).

GIS (Geographic Information System): analysis & visualization dalam geo metric.

fungsionalitas: - kordinat, aliran.

1. **Percarian (Search):** per area.
2. **Analisis lokasi (buffer, overlay):** aliran air dr kota dr rampai air.
3. **Analisis topografi (slope, aspect, drainage network):** curva dr kota dr kota.
4. **Flow analysis (connectivity, shortest path):** jarak terpendek dr sumber ke destinasi.
5. **Distribution (nearest neighbor, proximity, change detection):** jarak terdekat.
6. **Analisis/statistics spasial (pattern, centrality, similarity, topology):** jarak terdekat.
7. **Pengukuran (distance, perimeter, shape, adjacency, direction):** jarak terdekat.

Model dari Objek spasial:

1. **Single object:** tiap entitas dptan dr area terbatas (Point, Line, Region)
2. **Materi kota, hutan, sungai**
3. **Kumpulan objek yg pdg hub. spasial:** deskrpsi langsung ciri dr point dr line.
4. **Model peng. lokasi, pembagian negara nth provinsi/distril negara** ... dr area.
1. **Point:** hanya lokasi, tanpa ukuran/area (pusat kota, stasiun kereta, pos/potongan).
2. **Line** (atau polyline/sekuens dr line segment): jalur, rute, kabel, selatan.
3. **Region:** ruang 2D (negara, daerah, teman); terdiri dr bagian yg tidak terhubung.

Representasi fitur spasial

1. **Discrete features:** fitur individual yg dapat dibedakan
- objek 1 dgn objek lain tdk ada objek ditengahnya.
2. **Continous features:** fitur yg muncul secara kontinu diantara observasi. eg: elevasi curva hujan dptan tempat curah hujan nilai beda dr

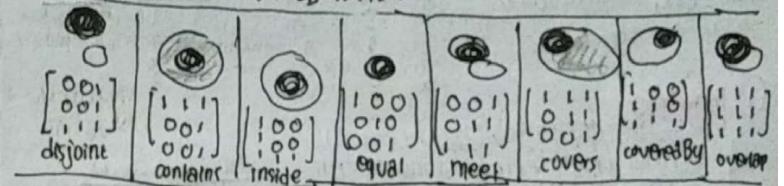
Spatial Relationship

1. **Topological relations:** toucher, overlap, contains, inside, equal, disjoint, intersect.
2. **Metric relations:** distance < d
3. **Direction relations:** north, south, west, east

Spatial Relationship drm matriks

$$J_{ij}(A, B) = \begin{bmatrix} A \cap B^o & A \cap \partial B & A \cap B^i \\ \partial A \cap B^o & \partial A \cap \partial B & \partial A \cap B^i \\ A^o \cap B^o & A^o \cap \partial B & A^o \cap B^i \end{bmatrix}$$

d = boundary
o = interior
- = exterior

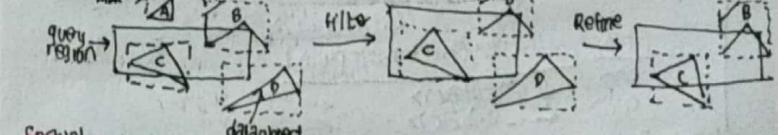


Query Processing

- spatial query language: SQL3, OGIS

- common strategy: filter & refine.

1. **Filter:** daerah query overlap drn MBR (minimum bounding rectangle/s) of B, C, D
2. **Refine:** daerah query overlap drn B, C



Spatial co-location

eg: ketertarikan hutan & pohon, kurang bersih & pemukiman, atau stasiun didekatnya ada rumah sakit, atau ada petrom, jadi objek online disekitar stasiun forensik akan ada banyak orang yg mau ke RS.

Spatial Classification

1. **Classification:** mengenal kemunculan bersama drn objek. berdekatan.
2. **Method spatial classification popular:**
 1. **Statistical Auto-regression (SAR):** 2. **Märkov Random Field (MRF):**
 3. **Algoritme Decision Tree pd spatial data:** perlu diproses, merubah algoritme/

Spatial Trend Analysis: mendekati perubahan & trend sepanjang satu dimensi spasial. tingkat kriminalitas/pengangguran.

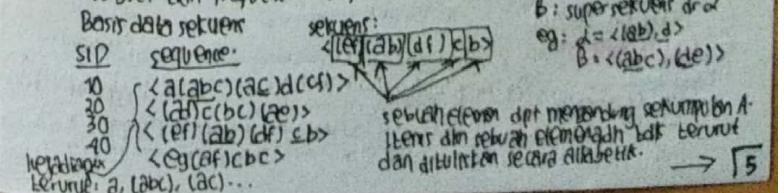
Spatial Cluster Analysis

pergelembahan daerah berdasarkan ciri yg sama, dari daerah yg dikelompokkan jd kluster cluster drn yg bersifat homogen. didefinisikan fitur unik dr daerah utk pemisahan, dibedakan berdasarkan warna.

Spatio-temporal data: data yg memiliki informasi spacial dan berubah seiring waktu. eg: ketebalan hutan, objek bergerak, badai, gelombang.

#12. Sequential Pattern Mining

1. **diberikan himpunan sekuen dan support threshold.** temukan himpunan lengkap dr frequent subsequences.



$\langle a(b)c \rangle$ adh subsequence dr $\langle b(abc) \rangle$ (ac) d (cf)

$\langle ab(bcd) \rangle$ adh supersequence dr $\langle abd \rangle$ d

dibutuhkan support threshold min-sup = 2, $\langle \langle ab \rangle c \rangle$ adh sequential pattern

- Aplikasi:
 - Untuk pengetahuan dr konsumen: pertama membeli komputer, kemudian CD-ROM dan kemudian digital camera, dlm 3 bulan.
 - pengetahuan medis: bencana alam (eg: gempa bumi), prones seismis dan letaknya, staf dian perbaikan, dll.
 - pola panggilan telepon, weblog click streams
 - sekuensi DNA dan struktur gen.

- Metode:
 - Dekatan berdasarkan Apriori \hookrightarrow GSP
 - Dekatan berbasis Pattern-Growth:
 - start dengan setiap sekuen S adh tidak frequent, maka tidak ada dari super-sekuen dr S yg frequent.
 - eg: $\langle ab \rangle$: infrequent $\rightarrow \langle hab \rangle$ dan $\langle (ab) \rangle b$ ~~infrequent~~ infrequent.

① GSP (Generalized Sequential Pattern) Mining (sama proses dan Apriori di ARM)

min-sup = 2		seq. ID	sequence	1st	kandidat	sup	
10	$\langle (b)cb(ac) \rangle$			$\langle a \rangle$			\rightarrow kemunculan biaptransisi
20	...			$\langle b \rangle$			

- Kemungkinan dr Apriori pruning: mengurangi ruang penelitian
- Bottlenecks:
 - scan basar data berulang
 - membangkitkan himpunan besar dr sekuen kandidat
 - (perlu metode mining yg lbh efisien)

② SPADE (Sequential Pattern Discovery Using Equivalence Class)

- format vertikal
- data sekuen dipotokan ke himpunan berukuran besar dr item:
 $\langle SID, EID \rangle$

SID sequence		SID	EID	Items			
1	$\langle a(b)c(b)ac \rangle$	1	1	a			
2	$\langle (ab)c(bc) \rangle$	1	2	abc	\rightarrow	SID EID	SID EID
		1	3	ac		1 1	1 2
		2	1	ad		1 2	2 3
		2	2	c		1 3	
		2	3	bc		2 1	

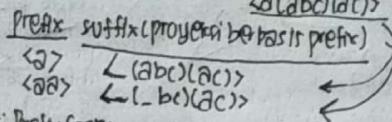
③ Prefix Span

(Prefix-Projected Sequential Pattern Growth)

- berbasis proyeksi, tetapi bukan proyeksi berbasis prefiks: proyeksi yg sedikit & jml sekuen menguras dgn cepat.

- Cara kerja: cari suffix dg polongan prefiks bertambah

eg: $\langle a \rangle$, $\langle aa \rangle$, $\langle acab \rangle$ = prefiks dr $\langle a(b)c(b)ac \rangle$



efisien: PrefixSpan:

- tidak ada kandidat sekuen yg perlu dibangkitkan
- projected databases, tetapi berukuran kecil
- blj dg prefiks span: membuat projected databases, dpt dpt perbaiki dg bi-level projection.

Kendala pd SPM

- Kendala item: cari pola weblog yg hny terkait online-bookstore.
- Kendala panjang (length): cari pola yg memiliki sedikitnya 20 barang.
- Kendala super pattern: cari super pattern dr "PC \rightarrow digital camera"
- Kendala agregat: cari pola yg tdk yg hanya barangnya melebihi \$100
- Kendala ekspresi reguler: cari pola "starting from Yahoo homepage, search for hotel in Washington DC area".
- Kendala durasi: cari pola seputar +24 jam dr pengambilan
- Kendala gap: cari pola pembelian gap antara setiap pembelian yg berurutan kurang dr 1 bulan.

#13 Text & Web Mining

- data mining yg mengolah dokumen teks sbg data.
- menugsi metode Information Retrieval (IR) utk praproses dr kumen
- proses yg secara otomatis menganalisa informasi komprehensif yg bermakna, berguna dan blm diketahui sebelumnya dr repotitoru dokumen terjauh.

- Text Mining Task
 - Diberikan: sumber dokumen teksual - kueri berbatar (berbasiskteks) yg didefinisikan
 - Temukan: kumpulan informasi relevan - ekstrak informasi relevan & abstrak informasi yg tdk relevan - hubungkan informasi & keterkaitan yg slg berhub. dlm format yg sdh difrapkan sebelumnya.

Task:	Search & retrieval	- kategorisasi
analis semantic	- feature extraction	
clustering	- pemisahan Ontologu	
	- dynamic focusing	
	Data Mining	Text Mining
objek yg dikenali	data numerik & kategorikal	text
struktur objek	basis relational	free form text
lujuhan	prediksi outcome dr situasi mendatang	benar kembali informasi yg relevan, distill the meaning, menkategorikan & penyampaikan target.
Metode	Machine learning:	Indexing, perceptron neural network khusus, linguistik, ontologies
ukuran market	SKAT, DT, NN, GA, MBR, MBA	100.000 analist pd penilaian besar & terengah
satu ini	100.000 analist pd penilaian besar & terengah	100.000.000 pekerja korporat dan pengguna individual
IMPLEMENTASI	implementasi yg ksr sejak 1994	implementasi yg luar nabi 2000

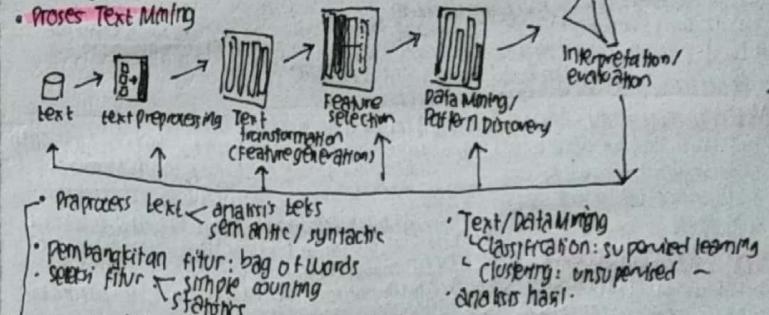
Aplikasi:

- Pemasaran: menemukan kelompok pembeli yg potensial berdasarkan profil lets pengguna. Eg: amazon
- Industri: mengidentifikasi sistem web kebanyak perang. Eg: produk posling
- Penerapan kerja: mengidentifikasi parameter dr penilaian penerapan. Eg: flipdog.com
- Search engines
- Enterprise portals
- Knowledge management system
- Virtual apps: kategori & routing email
- e-Business systems: kategori catatan call center sistem CRM.

Karakteristik Text:

- Basis dr teks, berukuran besar (pertumbuhan eksponensial $< 2.000.000.000 hal. web$)
- Dimensi tinggi (sparse input): pertumbuhan sebagi word/phrase sbg sebuah dimensi.
- Beberapa mode input:
 - web mining: informasi mengenai pengguna dr konten oleh semantik, mencari pola & outside knowledge base.
- Ketergantungan (dependensi): informasi yg relevan dthd hub. conjunction yg kompleks dr word/phrase. Eg: kategorisasi, dr ambyar pronoun.
- Ambiguitas \leftarrow ambiguitas kata: pronouns (he, she, ...), "buy" "purchase" ambiguitas semantik: the king saw the rabbit with his glasses (8 makna).
- Noisy data: eg: kesalahan ejaan
- Teks yg fdkt struktur yg baik: chat rooms, speech.

Proses Text Mining



- Text / Data Mining
 - Klasifikasi:
 - dibentuk dr record dg 1 label (training set). setiap record mngandung
 - konten: model utk kelas sbg fungsi dr nilai dg fitur.
 - lujuhan: record yg blm ada sebelumnya diberi label sejakurat mungkin.
 - test set digunakan vlc menghitung akurasi model.
 - Clustering:
 - dibentuk: setumpuk dokumen & similarity measure dr antara dokumen.
 - bentuk: cluster \rightarrow dokumen dr dalam 1 cluster dg dokumen lainnya.
 - lujuhan: kumpulan dokumen yg blm sebelumnya.
 - similarity measures: euclidean distance jika diberi dg klasifikasi

- Text / Data Mining
 - Klasifikasi:
 - dibentuk dr record dg 1 label (training set). setiap record mngandung
 - konten: model utk kelas sbg fungsi dr nilai dg fitur.
 - lujuhan: record yg blm ada sebelumnya diberi label sejakurat mungkin.
 - test set digunakan vlc menghitung akurasi model.
 - Clustering:
 - dibentuk: setumpuk dokumen & similarity measure dr antara dokumen.
 - bentuk: cluster \rightarrow dokumen dr dalam 1 cluster dg dokumen lainnya.
 - lujuhan: kumpulan dokumen yg blm sebelumnya.
 - similarity measures: euclidean distance jika diberi dg klasifikasi

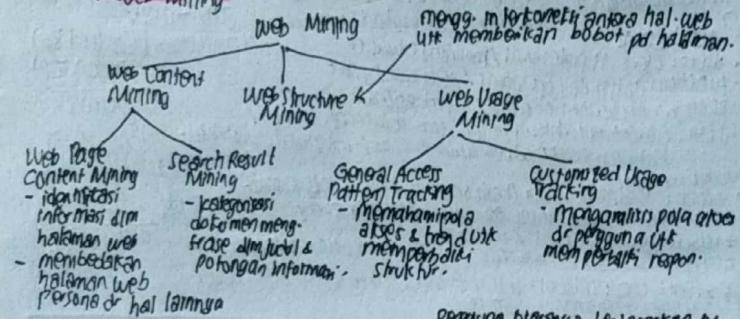
Web Mining

- mencari informasi menarik & berguna dr konten dan penggunaan web.
- 1. Web search: Google, Yahoo, MSN, Ask, ...
- 2. Specialized search: eg. google.com shopping (comparison shopping), job ads (cari lowongan)
- 3. E-commerce: rekomendasi: eg. Netflix, Amazon. memperbaiki tingkat konver: produk terbaik berikutnya yg ditawarkan.
- 4. Iklan, eg: Google AdSense
- 5. Fraud detection: click fraud detection
- 6. meningkatkan ranking & kategori website. (kelanjutan/lanjutannya web)
- 7. teknik* damping utk mencari & mengetahui secara otomatis informasi dr

- Why ?

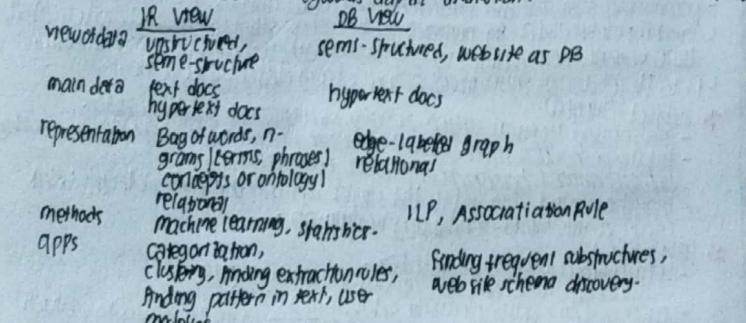
1. Web tglakan informasi teknisal
 - Book / CD / Video store eg. Amazon
 - Information retrieval eg. Zazzle
 - Harga mobil eg. Carpoint
 2. Banyaknya data pd pola akter pengguna
 3. Possible to retrieve "previously unknown" information si akter pengguna.
- Filter untuk pd web
- 1. Koleksi dokumen berukuran besar yg mengandung: informasi hyperlink & informasi
 - 2. Web sgt dinamis: hal web scr dapat berubah-ubah (dibuang)

	<u>Web Mining</u>	<u>Data Mining</u>
Struktur	web bln miring relasi	informasi teknisal dgn struktur linage
skala	dapat penggunaan berukuran besar & berukuran besar & bertumbuhscr cpt.	dala yg dibangkitkan berharap dpt dibandingkan dgn data warehouse konvensional yg terbesar.
kecepatan	segera perlu bereaksi utk mengembangkan pola penggunaan scr real-time yg merambang	manusia tidak terlibat dlm proses.
Taksonomi	Web Mining	Mengg. interkaitan antara hal-hal web utk memberikan bobot pd halaman.



// Pendekatan Web Content Mining

1. Information Retrieval: membantu & meningkatkan pencarian & filtering informasi utk diminta.
2. Basis data: utk memudahkan data pola & mengintegrasikannya dr sumber yg kompatibel selain berbasis keywords dapat dilakukan.



// Web Structure Mining

view of data: hierachical
 main idea: server log, browser log
 representation: graph
 methods: proprietary algorithms
 apps: categorization, clustering

// Web Usage Mining = Web Log Mining

view of data: interactional
 main idea: server log, browser log
 representation: relational table graph
 methods: machine learning, statistics, association rule
 apps: site construction, adaptation & management, marketing, user modeling

1. menargetkan customer yg potensial utk produk elektronik aktif.
2. memperbaik kualitas dr pengalaman Internet Information Services kepada pengguna
3. memperbaiki performa sistem web server
4. mengidentifikasi lokasi iklan yg potensial
5. memfasilitasi personalisasi & fitur adaptif
6. memperbaiki & optimis silis 7. deteksi fraud/intuition. 8. memprediksi ala pengguna