

## BAB II

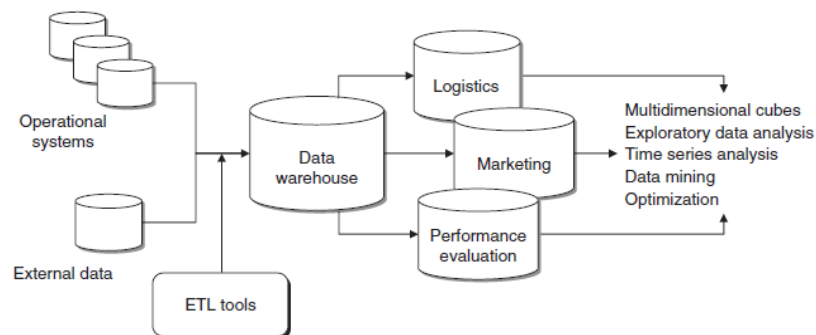
### LANDASAN TEORI

#### 2.1 *Data Mining*

##### 2.1.1 Pengertian *Data Mining*

Menurut Han dan Kamber(2006, p.9) *data mining* merupakan proses menggali informasi yang tersembunyi dari data berjumlah besar untuk pengambilan keputusan, prediksi ataupun pemecahan masalah. Dimana melibatkan integrasi teknik dari berbagai disiplin ilmu seperti database dan teknologi data warehouse, statistik, *machine learning*, performa tinggi komputasi, pengenalan pola, jaringan saraf, visualisasi data, penemuan informasi, gambar dan pemrosesan sinyal, dan analisis data spasial atau temporal.

Sedangkan Vercellis (2009, p.9-11), menggambarkan *data mining* sebagai bagian dari arsitektur sistem *business intelligence* (BI).



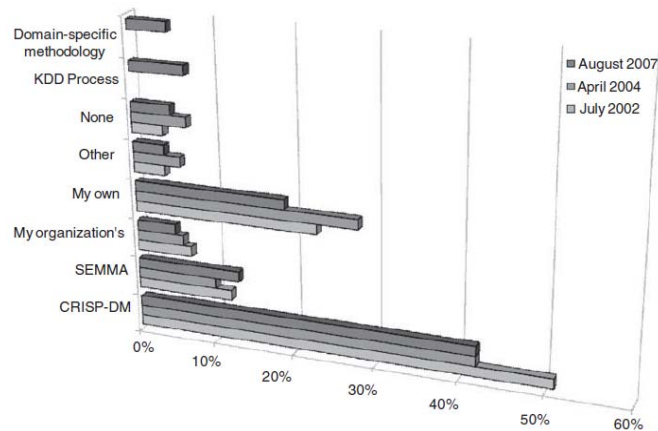
Gambar 2-1 Arsitektur sistem BI (Vercellis, 2009)

Dari gambar di atas, dapat dilihat karakteristik *data mining* yang tergantung pada prosedur untuk mengumpulkan data *history* yang berasal dari

berbagai sumber (*Operational systems* dan *External data*) kemudian memasukkannya ke dalam *database* (*data warehouse* atau *data marts*) melalui *ETL tools*. *Data mining* pada sistem BI merupakan metodologi aktif dengan tujuan untuk ekstraksi informasi dan pengetahuan dari data yang ada di *data warehouse/ data marts*. *Data mining* termasuk model matematika untuk pengenalan pola, *machine learning* dan teknik *data mining*. *Data mining* tidak memerlukan hipotesis sebelumnya untuk diverifikasi tapi sebaliknya untuk memperluas pengetahuan pembuat keputusan.

### **2.1.2 Cross-Industry Standard Process for Data Mining (CRISP-DM)**

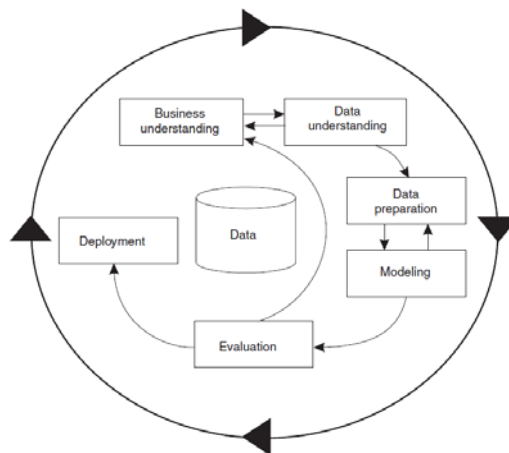
Dalam menanggapi isu-isu umum dan kebutuhan dalam proyek *data mining* dipertengahan 90-an, kelompok organisasi yang terlibat dalam *data mining* (Teradata, SPSS-ISL, Daimler-Chrysler dan Ohra) mengusulkan panduan referensi untuk mengembangkan proyek-proyek *data mining*, bernama CRISP-DM (*Cross Industry Standard Process for Data Mining*). Mariscal, Marba dan Fernandez (2010) menyatakan CRISP-DM sebagai *de facto* standar untuk pengembangan proyek *data mining* dan *knowledge discovery* karena paling banyak digunakan dalam pengembangan *data mining*. Hal tersebut dapat terlihat dari survei yang dilakukan terhadap penggunaan metodologi dalam proyek *data mining*:



Gambar 2-2 Survei Penggunaan Metodologi Data Mining (Mariscal, Marban, & Fernandez, 2010)

Hasil survei “Penggunaan Metodologi dalam Proyek *Data Mining*”, memperlihatkan pengguna CRISP-DM di tahun 2002 mencapai 51%, kemudian menurun menuju 41% di tahun 2004. Meskipun persentasi penggunaan CRISP-DM menurun 10%, jumlah pengguna metodologi ini masih terbilang lebih banyak daripada pengguna metodologi lain.

Model proses CRISP-DM memberikan gambaran tentang siklus hidup proyek *data mining*. Siklus hidup proyek *data mining* dalam CRISP-DM terdiri dari enam fase (Gambar 2-3).



Gambar 2-3 Fase CRISP-DM (Chapman, dkk, 2000)

Pada CRISP-DM, urutan fase tidak kaku, dapat bergerak bolak-balik antar fase yang berbeda. Tanda panah pada gambar 2-3, menunjukkan frekuensi ketergantungan antara fase. Berikut gambaran dari garis besar setiap fase disertai dengan tugas-tugas (*bold*) dan output (*italic*).

<b>Business Understanding</b>	<b>Data Understanding</b>	<b>Data Preparation</b>	<b>Modeling</b>	<b>Evaluation</b>	<b>Deployment</b>
<b>Determine Business Objectives</b> <i>Background Business Objectives Business Success Criteria</i>	<b>Collect Initial Data</b> <i>Initial Data Collection Report</i>	<b>Select Data</b> <i>Rationale for Inclusion/Exclusion</i>	<b>Select Modeling Techniques</b> <i>Modeling Technique Modeling Assumptions</i>	<b>Evaluate Results</b> <i>Assessment of Data Mining Results w.r.t. Business Success Criteria Approved Models</i>	<b>Plan Deployment</b> <i>Deployment Plan</i>
<b>Assess Situation</b> <i>Inventory of Resources Requirements, Assumptions, and Constraints Risks and Contingencies Terminology Costs and Benefits</i>	<b>Describe Data</b> <i>Data Description Report</i>	<b>Clean Data</b> <i>Data Cleaning Report</i>	<b>Generate Test Design</b> <i>Test Design</i>	<b>Review Process</b> <i>Review of Process</i>	<b>Plan Monitoring and Maintenance</b> <i>Monitoring and Maintenance Plan</i>
<b>Determine Data Mining Goals</b> <i>Data Mining Goals Data Mining Success Criteria</i>	<b>Explore Data</b> <i>Data Exploration Report</i>	<b>Construct Data</b> <i>Derived Attributes Generated Records</i>	<b>Build Model</b> <i>Parameter Settings Models Model Descriptions</i>	<b>Determine Next Steps</b> <i>List of Possible Actions Decision</i>	<b>Produce Final Report</b> <i>Final Report Final Presentation</i>
<b>Produce Project Plan</b> <i>Project Plan Initial Assessment of Tools and Techniques</i>	<b>Verify Data Quality</b> <i>Data Quality Report</i>	<b>Integrate Data</b> <i>Merged Data</i>	<b>Assess Model</b> <i>Model Assessment Revised Parameter Settings</i>		<b>Review Project</b> <i>Experience Documentation</i>
		<b>Format Data</b> <i>Reformatted Data Dataset Dataset Description</i>			

Gambar 2-4 Tugas dan *output* dari model CRISP-DM (Chapman, dkk,2000)

### 2.1.2.1 Business Understanding

#### 2.1.2.1.1 Determine Business Objectives

Pada tahap awal, tim proyek harus benar-benar memahami dari perspektif bisnis, apa yang ingin dicapai. Adapun tahap ini bertujuan untuk memahami bidang masalah, menghasilkan solusi yang tepat, dan mengungkapkan faktor penting yang berpengaruh pada hasil proyek.

*Output* dari tahap ini merupakan analisis dari hasil observasi, wawancara dan dokumen-dokumen yang dimiliki perusahaan. Berikut penjelasan dari *output* yang dihasilkan:

- **Background**

*Output* ini berisi catatan informasi yang diketahui tentang situasi bisnis, organisasi pada awal proyek dan gambaran dasar dari konteks proyek. Konteks proyek yang dimaksud seperti proyek ini terjadi di bidang apa, identifikasi masalah, dan mengapa data mining muncul untuk memberikan solusi. Aktivitas pada point ini:

	Aktivitas
Organisasi	Memahami bagan organisasi, identifikasi divisi, departemen, dan kelompok proyek
	Identifikasi nama dan tanggung jawab manajer
	Mengidentifikasi orang-orang penting dalam bisnis dan peran mereka.
	Mengidentifikasi yang menjadi sponsor internal (sponsor keuangan dan user ahli/ domain utama). Sponsor proyek merupakan orang yang memberikan arahan atas proyek atau sebagai <i>client</i> dari proyek.
	Mengidentifikasi unit bisnis yang terkena proyek data mining (misalnya, pemasaran, penjualan, keuangan)
Identifikasi masalah	Mengidentifikasi bidang masalah (misalnya, pemasaran, layanan pelanggan, pengembangan bisnis, dan lain-lain)
	Periksa status proyek (misalnya, apakah proyek data mining akan dilakukan dalam unit bisnis? Atau apakah data mining perlu dipromosikan sebagai kunci teknologi dalam bisnis?)

	Menjelaskan prasyarat proyek (misalnya, Apa motivasi proyek? Apakah bisnis sudah menggunakan data mining?)
	Mengidentifikasi sasaran untuk hasil proyek (misalnya, Apakah kita diharapkan memberikan laporan untuk manajemen puncak atau sistem operasional yang digunakan oleh <i>end user</i> ?)
	Mengidentifikasi kebutuhan dan harapan pengguna
Solusi sekarang	Jelaskan solusi saat ini yang digunakan untuk mengatasi masalah
	Jelaskan keuntungan dan kerugian dari solusi saat ini dan tingkat yang diterima oleh pengguna

- ***Business objectives***

*Output* ini menjelaskan tujuan utama konsumen, dari perspektif bisnis. Di samping tujuan utama bisnis, biasanya ada pertanyaan lain yang berhubungan dengan bisnis. Misalnya, tujuan bisnis mungkin untuk mempertahankan konsumen saat ini dengan memprediksi kapan mereka pindah ke pesaing. Contoh pertanyaan-pertanyaan bisnis yang berhubungan, "Apakah channel utama yang digunakan (misalnya, ATM dan Internet) mempengaruhi perginya konsumen?" atau "Apakah menurunkan biaya ATM secara signifikan mengurangi jumlah konsumen yang pergi?". Aktivitas pada point ini terdiri dari:

- Gambarkan masalah yang harus dipecahkan.
- Spesifikasikan semua pertanyaan bisnis secara cermat mungkin.

- Spesifikasikan kebutuhan bisnis lainnya (misalnya, tidak ingin kehilangan pelanggan)
- Spesifikasikan manfaat yang diharapkan dalam istilah bisnis
- ***Business success criteria***

*Output* ini menjelaskan kriteria sukses atau kegunaan hasil proyek dari sudut pandang bisnis. Penjelasan dilakukan dengan spesifik dan dapat diukur secara obyektif, misalnya, pengurangan pelanggan yang pergi ke tingkat tertentu, atau secara subyektif, seperti "memberikan wawasan yang berguna". Namun harus ditunjukkan penilaian subjektif dibuat oleh siapa. Aktivitas dari point ini:

  - Spesifikasikan kriteria keberhasilan usaha (misalnya, Meningkatkan respond dalam promosi mailing sebesar 10% dan tingkat pendaftaran sebesar 20%).
  - Mengidentifikasi siapa yang menilai kriteria keberhasilan.

#### **2.1.2.1.2 Assess situation**

Tugas ini mencakup penemuan fakta yang lebih rinci tentang semua sumber daya, kendala, asumsi, dan faktor lain yang harus dipertimbangkan dalam menentukan tujuan analisis data dan rencana proyek. Adapun *output* dari tugas ini adalah:

- ***Inventory of resources***

*Output* ini berisi daftar sumber daya yang tersedia untuk proyek, termasuk personil (*business experts, data experts, technical support, data mining experts*), data (*fixed extracts, access to live, warehoused, or operational data*), sumber daya komputasi (*hardware platforms*), dan

perangkat lunak (*data mining tools*, perangkat lunak lainnya yang relevan).

Berikut aktivitas dari point ini:

	Aktivitas
Sumber daya hardware	Mengidentifikasi perangkat keras
	Menetapkan ketersediaan perangkat keras dasar untuk proyek
	Periksa apakah jadwal pemeliharaan hardware konflik dengan ketersediaan perangkat keras untuk proyek data mining
	Mengidentifikasi perangkat keras yang tersedia untuk <i>tool</i> data mining yang akan digunakan (jika <i>tool</i> yang digunakan sudah diketahui pada tahap ini)
Sumber daya data dan pengetahuan	Mengidentifikasi sumber data
	Identifikasi jenis sumber data (sumber online, ahli, dokumentasi tertulis, dll)
	Mengidentifikasi sumber-sumber pengetahuan
	Identifikasi jenis sumber pengetahuan (sumber online, tenaga ahli, dokumentasi tertulis, dll)
	Periksa alat dan teknik yang tersedia
	Jelaskan latar belakang pengetahuan yang relevan (atau informal)
Sumber daya manusia	Identifikasi sponsor proyek (jika berbeda dari sponsor internal dalam bagian 2.1.2.1.1.1)
	Mengidentifikasi sistem administrator, database



	administrator, dan stafdukungan teknis untuk pertanyaan lebih lanjut
	Mengidentifikasi analis pasar, ahli data mining, dan ahli statistik, dan memeriksa ketersediaan mereka
	Periksa ketersediaan pakar domain untuk tahap selanjutnya

- ***Requirements, assumptions, and constraints***

*Output* yang dimaksud berupa:

Daftar semua *requirements* proyek, termasuk jadwal penyelesaian, komprehensibilitas dan kualitas hasil, keamanan, serta masalah hukum seperti memastikan ijin untuk penggunaan data.

Daftar asumsi yang dibuat oleh proyek. Ini mungkin asumsi tentang data yang dapat diverifikasi selama data mining, atau asumsi non-verifikasi tentang bisnis yang terkait dengan proyek.

Daftar batasan pada proyek. Daftar ini berisi batasan pada ketersediaan sumber daya, termasuk teknologi seperti ukuran data set yang digunakan untuk pemodelan.

Berikut aktivitas dari point ini:

	Aktivitas
<i>Requirements</i>	Spesifikasi target
	Menangkap semua keperluan pada penjadwalan
	Menangkap keperluan komprehensibilitas, akurasi, menyebarkan kemampuan, pemeliharaan, dan

	pengulangandari proyekdata miningdan model yang dihasilkan
	Menangkapkeperluankeamanan,batasan hukum, privasi,pelaporan, danjadwal proyek
Asumsi	Menjelaskansemua asumsi(termasuk yang implisit) dan menjadikannya eksplisit
	Daftarkanasumsi pada kualitasdata (misalnya, akurasi, ketersediaan)
	Daftarkanasumsi pada faktor eksternal(misalnya, isu-isu ekonomi, produk yang kompetitif, kemajuan teknis)
	Menjelaskanasumsi yang mengarahke salah satu perkiraan(misalnya, harga alat tertentu diasumsikan lebih rendah dari \$ 1.000)
	Daftar semuaasumsi tentang apakah perlu untuk memahami dan menggambarkan atau menjelaskan model(misalnya, bagaimana seharusnya menyampaikan model dan hasil kepada manajemen senior/sponsor)
Batasan	Periksa batasan umum (misalnya, masalah hukum, anggaran, rentang waktu, dan sumber daya)
	Periksa hak akses ke sumber data(misalnya, pembatasan akses, password)
	Periksa aksesibilitas teknis data (sistem operasi, sistem

	manajemendata,file ataudatabase format)
	Periksa apakahpengetahuan yang relevandapat diakses
	Periksa kendalaanggaran(biaya tetap, biaya pelaksanaan, dll)

- ***Risks and contingencies***

*Output* ini berisi daftarresikoyang mungkin menundaataumenyebabkan proyekgagal.Termasuk rencana kontingensi, yaitu tindakan apa yang akan diambil jikaresikoterjadi. Aktivitas pada point ini:

	Aktivitas
Identifikasi resiko	Identifikasiresiko bisnis(misalnya, pesaing memiliki hasil yang lebih baik)
	Mengidentifikasi resikoorganisasi (misalnya, departemen yang meminta proyektidak memilikidana untukproyek)
	Identifikasiresiko keuangan(misalnya, pendanaan lebih lanjutbergantung padahasilawal data mining)
	Identifikasiresiko teknis
	Mengidentifikasi resiko yang bergantung pada data dan sumber data (misalnya, kualitas buruk dan cakupan)
Mengembangkan rencana kontigensi	Menentukankondisi di manaresikodapat terjadi
	Mengembangkan rencanakontigensi

- ***Terminology***

*Output* ini berisidaftar terminologiyang relevan denganproyek.Dua komponen yang mungkin ada pada point ini :

1. Glossarybisnis, yangmerupakan bagian daripemahaman bisnisuntukproyek
2. Glossarydata mining, diilustrasikandengan contohyang relevan dari pertanyaan bisnis

- ***Costs and benefits***

Point ini terdiri dari analisis biaya dan laba untuk proyek, yang membandingkanbiaya proyekdengan potensibagibisnisjika proyek berhasil.Analisis biaya dan laba menghasilkan ROI (*Return of Investment*) yang didapatkan dari proyek. Berikut adalah rumus dari perhitungan ROI:

$$(total\ discounted\ benefits - total\ discounted\ costs) / discounted\ costs$$

$$Discounted\ benefits = benefit \times discount\ factors$$

$$Discounted\ costs = cost \times discount\ factors$$

$$Discount\ factors = 1 / (1 + r)^t$$

*Discount factors* adalah pengali untuk setiap tahun berdasarkan *discount rate* dan tahun, yang digunakan untuk perhitungan biaya dan laba setiap tahun.

$r = discount\ rate$  (nilai yang digunakan untuk memotong *cash flow* di masa depan)

$t =$  tahun ke-

### **2.1.2.1.3 Determine data mining goals**

Tujuan *data mining* bagian dari tujuan proyek dalam istilah-istilah teknis. Sebagai contoh, tujuan bisnis "Meningkatkan penjualan katalog

untuk pelanggan yang sudah ada". Maka tujuan *data mining* "Memprediksi berapa banyak *widget* pelanggan akan membeli, melihat pembelian mereka tiga tahun terakhir, informasi demografis (umur, gaji, kota, dll), dan harga item". *Output* yang dihasilkan dalam tugas ini:

- ***Data mining goals***  
Menjelaskan *output* proyek yang diharapkan untuk mencapai tujuan bisnis.
- ***Data mining success criteria***  
Menentukan kriteria sukses dalam hal teknis misalnya, akurasi prediksi di tingkat tertentu. Seperti halnya kriteria keberhasilan bisnis, mungkin kriteria sukses digambarkan dalam hal subyektif, sehingga orang atau pihak yang melakukan penilaian harus diidentifikasi.

#### **2.1.2.1.4 Produce project plan**

Tugas ini menjelaskan rencana yang ditetapkan untuk mencapai tujuan *data mining*. Rencana yang dibuat berisikan langkah-langkah yang akan dilakukan sampai akhir proyek, termasuk pemilihan alat dan teknik.

- ***Project plan***  
*Output* ini berisi daftar tahap yang akan dieksekusi dalam proyek, durasi, sumber daya yang dibutuhkan, *input, output*, dan dependensi. Bagian dari *project plan* adalah untuk menganalisis ketergantungan antar waktu dan resiko. *Project plan* adalah dokumen yang dinamis dalam arti bahwa pada akhir setiap tahap, diperlukan *review* kemajuan, prestasi dan *updates* sesuai jadwal.
- ***Initial assessment of tools and techniques***  
Penilaian alat dan teknik dilakukan karena dapat mempengaruhi seluruh proyek. Aktivitas pada point ini:

- Membuat daftar kriteria untuk seleksi alat dan teknik
- Pilih alat dan teknik yang potensial
- Evaluasi kesesuaian teknik
- *Review* dan *memprioritaskan* teknik yang diterapkan menurut evaluasi dari alternatif

## 2.1.2.2 Data Understanding

### 2.1.2.2.1 Collect initial data

Pada tugas ini dilakukan akses ke data yang tercantum dalam sumber daya proyek, mencakup *load data*. *Output* yang dihasilkan:

- *Initial Data Collection Report*

Menjelaskan semua data yang digunakan dalam proyek, termasuk kebutuhan *selection* untuk data yang lebih detail. Selain itu, ditentukan atribut yang lebih penting daripada atribut lain. Untuk itu, dilakukan observasi ke perusahaan untuk mengetahui prosedur pengambilan data dan ketersediaan data. Aktivitas pada point ini:

	Aktivitas
<i>Data requirements planning</i>	Mendaftarkan informasi yang diperlukan (misalnya, memberikan atribut, atau informasi tambahan tertentu)
	Periksa apakah semua informasi yang dibutuhkan tersedia
<i>Selection criteria</i>	Tentukan kriteria <i>select</i> (misalnya, atribut yang diperlukan untuk tujuan data mining? apakah

	relevan? berapa banyak atribut yang dapat ditanganidengantekniky yang dipilih?)
	Pilih tabel/file berkepentingan
	Pilih data dalam tabel/file
	Tentukan menggunakan data <i>history</i> berapa lama (misalnya, data 18 bulan yang tersedia, hanya 12 bulan mungkin diperlukan untuk latihan)

### 2.1.2.2.2 Describe Data

Tugas ini memeriksa properti data yang diperoleh dan melaporkan hasilnya. *Output* yang dihasilkan:

- *Data Description Report*

Menjelaskan data yang telah diperoleh, format data, jumlah data (misalnya, jumlah *record* dan *field* dalam setiap tabel), identitas *field*, dan fitur lainnya yang ditemukan. Aktivitas pada point ini:

	Aktivitas
Analisis volume data	Identifikasi data dan metode <i>capture</i>
	Akses sumber data
	Gunakan analisis statistik jika diperlukan
	Laporan tabel dan hubungan mereka
	Periksa volume data, kompleksitas
	Catat data yang berisi teks bebas
Jenis dan nilai atribut	Periksa aksesibilitas dan ketersediaan atribut
	Periksa jenis atribut (numerik, simbolik, dll)
	Periksa rentang nilai atribut
	Menganalisis korelasi atribut
	Memahami arti dari setiap atribut dan nilai atribut dalam

	istilah bisnis
	Untuk setiap atribut, menghitung statistik dasar (misalnya, menghitung distribusi, rata-rata, max, min, standar deviasi, varians, modus, kemiringan, dll)
	Menganalisis statistik dasar dan makna hasilnya yang berhubungan dengan bisnis
	Putuskan apakah atribut relevan untuk tujuan data mining
	Menentukan apakah makna atribut yang digunakan konsisten
	Wawancara dengan pakar domain untuk mendapatkan pendapat mereka tentang relevansi atribut
	Menentukan apakah perlu untuk menyeimbangkan data (berdasarkan teknik pemodelan yang akan digunakan)
Keys	Menganalisis hubungan antar kunci
	Periksa jumlah nilai atribut kunci yang tumpang tindih di tabel lain
Tinjauan asumsi/tujuan	Bila perlu, perbaharui daftar asumsi

### 2.1.2.2.3 Explore Data

Tugas ini menanganikan pertanyaan-pertanyaan *data mining* tentang *query*, visualisasi, dan teknik pelaporan. Analisis ini dapat secara langsung menjawab tujuan *data mining*, kontribusi perbaikan deskripsi data dan kualitas laporan, dan menjadi masukan bagi transformasi dan langkah-langkah persiapan data lainnya yang diperlukan sebelum analisis lebih lanjut dapat terjadi.

*Output* yang dihasilkan:

- *Data Exploration Report*



Menjelaskan hasil dari tugas ini, termasuk hipotesis awal dan dampaknya terhadap sistem proyek. Laporan ini juga dapat mencakup grafik dan plot yang menunjukkan karakteristik data atau data subset yang menarik untuk pemeriksaan lebih lanjut. Aktivitas dari tahap ini:

	Aktivitas
<i>Data exploration</i>	Menganalisis properti dari atribut secara rinci (misalnya, statistik dasar, sub-populasi yang menarik)
	Mengidentifikasi karakteristik sub-populasi
Anggapan untuk analisis masa depan	Mempertimbangkan dan mengevaluasi informasi dan temuan dalam laporan deskripsi data
	Bentuk hipotesis dan mengidentifikasi tindakan
	Bila memungkinkan, ubah hipotesis menjadi tujuan data mining
	Memperjelas tujuan data mining atau membuatnya lebih tepat. Misalnya, pencarian “blind” tidak selalu berguna, tapi pencarian lebih diarahkan ke tujuan bisnis adalah lebih baik.
	Lakukan analisis dasar untuk memverifikasi hipotesis

#### **2.1.2.2.4 Verify data Quality**

Pada tugas ini dilakukan pemeriksaan kualitas data, menangani pertanyaan-pertanyaan seperti: Apakah data lengkap? (tidak untuk semua kasus), Apakah ada kesalahan? Jika ada kesalahan, berapa banyak kesalahan

tersebut? Apakah ada nilai yang hilang dalam data? Jika demikian, bagaimana hal tersebut direpresentasikan?. *Output* yang dihasilkan:

- *Data Quality Report*

Berisi daftar hasil verifikasi kualitas data, jika ada masalah kualitas, daftarkan juga solusi yang dapat dilakukan.

	Aktivitas
<i>Review key dan atribut</i>	Periksa luasan cakupan yang digunakan (misalnya, apakah mungkin semua nilai diwakili)
	Periksa <i>Key</i>
	Pastikan bahwa makna dari atribut dan nilai-nilai yang terkandung cocok satu sama lain
	Mengidentifikasi atribut hilang dan <i>field</i> kosong
	Menetapkan makna dari data yang hilang
	Periksa apakah atribut dengan nilai yang berbeda namun memiliki arti yang sama (misalnya, rendah lemak, diet)
	Periksa ejaan dan format nilai (misalnya, nilai yang sama tapi kadang-kadang dimulai dengan huruf kecil, kadang-kadang dengan huruf kapital)
	Periksa penyimpangan, apakah penyimpangan adalah " <i>noise</i> " atau mungkin menunjukkan fenomena menarik
	Periksa apakah nilai-nilai yang ada masuk akal, (misalnya, semua <i>field</i> harus sama atau mendekati nilai yang sama)

## 2.1.2.3 Data Preparation

### 2.1.2.3.1 Select data

Pada tugas ini ditentukan data yang akan digunakan. Kriteria yang dipakai meliputi relevansi dengan tujuan *data mining*, kualitas, dan kendala teknis seperti batas pada volume data atau jenis data. *Output* yang dihasilkan :

- *Rationale for inclusion/exclusion*

Berupa daftar data yang akan digunakan/dikeluarkan dan alasan untuk keputusan ini. Aktivitas pada point ini:

- Mengumpulkan data tambahan yang sesuai (dari berbagai sumber-internal maupun eksternal)
- Melakukan signifikansi dan korelasitas untuk memutuskan apakah suatu field harus disertakan
- Mempertimbangkan kembali kriteria *select data* (bagian 2.1.2.1.3.1) dengan melihat kualitas data dan data eksplorasi (misalnya, mungkin ingin menyertakan data set yang lain)
- Mempertimbangkan kembali kriteria *select data* (bagian 2.1.2.1.3.1) dengan melihat pemodelan (mungkin model menunjukkan dibutuhkannya data set lain)
- Pilih *data set* yang berbeda (misalnya, atribut yang berbeda, hanya data yang memenuhi kondisi tertentu)
- Pertimbangkan penggunaan teknik sampling (misalnya, sebuah solusi cepat mungkin memerlukan *splitting test* dan data set pelatihan atau mengurangi ukuran *dataset*, jika alat tidak dapat menangani ukuran *dataset*)
- Dasar pemikiran untuk inklusi / eksklusi
- Periksa teknik yang tersedia untuk data sampel

### 2.1.2.3.2 Clean data

Tugas ini bertujuan meningkatkan kualitas data ke tingkat yang dibutuhkan oleh teknik analisis. Melibatkan pemilihan *subset* data yang bersih, penyisipan *default* yang cocok. *Output* yang dihasilkan:

- *Data Cleaning Report*

Menjelaskan keputusan dan tindakan yang diambil untuk mengatasi masalah kualitas data yang dilaporkan pada tahap *verify data quality*. Aktivitas pada point ini:

- Mempertimbangkan kembali bagaimana menangani *noise*.
- Memperbaiki, menghapus, atau mengabaikan *noise*.
- Putuskan bagaimana menangani nilai-nilai khusus dan maknanya. Contoh nilai-nilai khusus yang bisa timbul melalui pengambilan hasil survei di mana beberapa pertanyaan tidak ditanyakan atau tidak menjawab. Hal ini dapat mengakibatkan nilai 99 untuk data yang tidak diketahui. Sebagai contoh, 99 untuk status perkawinan atau afiliasi politik. nilai-nilai khusus juga bisa muncul ketika data dipotong - misalnya, 00 orang untuk 100 tahun atau semua mobil dengan 100.000 km di odometer.
- Mempertimbangkan kembali kriteria seleksi data (Lihat Tugas 2.1.2.1.3.1) dalam pembersihan data (misalnya, memasukkan/mengecualikan *dataset* lain).

### 2.1.2.3.3 Construct data

Tugas ini meliputi operasi persiapan konstruktif data seperti mendapatkan atribut, menyelesaikan rekor baru, atau perubahan nilai-nilai untuk atribut yang ada. *Output* yang dihasilkan:

- *Derived attributes*

*Derived attributes* adalah atribut baru yang dibangun dari satu atau lebih atribut yang ada di record yang sama. Contoh: luas = panjang \* lebar. *Derived attributes* dibangun karena:

- Latar belakang meyakinkan bahwa beberapa fakta penting dan harus diwakili meskipun saat ini tidak memiliki atribut yang mewakili.
- Hasil dari fase pemodelan menunjukkan bahwa fakta-fakta tertentu yang tidak tercakup

Aktivitas pada point ini terdiri dari:

- Putuskan apakah setiap atribut harus dinormalisasi (misalnya, bila menggunakan algoritma *clustering* dengan usia dan pendapatan, dalam mata uang tertentu, pendapatan akan mendominasi)
- Pertimbangkan untuk menambahkan informasi baru tentang pentingnya relevan atribut dengan menambahkan atribut baru (misalnya, bobot atribut, normalisasi bobot)
- Bagaimana atribut hilang dibangun atau diperhitungkan?
- Tambahkan atribut baru untuk data yang diakses.

- *Generated records*

*Generated records* adalah record baru yang menambah pengetahuan baru.

### **2.1.2.3.4 Integrate data**

Pada tugas ini dilakukan penggabungan informasi dari beberapa tabel atau sumber informasi lain untuk membuat *record* atau nilai-nilai baru. *Output* yang dihasilkan:

- *Merged data:*

*Merged Table* mengacu pada bergabungnya dua atau lebih tabel yang memiliki informasi yang berbeda dalam objek yang sama. Contoh: pada rangkaian ritel memilikisatu tabel dengan informasi tentang karakteristik umum masing-masing toko (misalnya, ruang lantai, jenis mall), tabel lain berisi ringkas data penjualan (misalnya, laba, perubahan persen penjualan dari tahun sebelumnya), dan informasi lain tentang demografi daerah sekitarnya. Masing-masing tabel tersebut berisi satu *record* untuk setiap toko. Tabel ini dapat digabungkan ke tabel baru dengan *record* untuk setiap toko merupakan penggabungan *fields* dari tabel sumber.

- *Aggregation*  
Merupakan operasi di mana nilai-nilai baru dihitung dengan meringkas informasi dari beberapa *record* atau tabel. Aktivitas yang dilakukan:
  - Periksa apakah fasilitas integrasi mampu mengintegrasikan sumber masukan yang diperlukan.
  - Mengintegrasikan sumber dan hasil penyimpanan
  - Mempertimbangkan kembali kriteria seleksi data (bagian 2.1.2.1.3.1) mengingat pengalaman integrasi data (misalnya, Anda mungkin ingin untuk memasukkan / mengeluarkan *dataset* yang lain)

### **2.1.2.3.5 Format Data**

Transformasi format mengacu pada sintak modifikasi yang dibuat untuk tidak mengubah makna data, tapi mungkin diperlukan oleh alat pemodelan. *Output* yang dihasilkan:

- *Reformatted data.*

Beberapa alat memiliki persyaratan pada urutan atribut, seperti *field* pertama menjadi unik pengidentifikasi untuk setiap record atau *field* terakhir menjadi bidang hasil model adalah untuk memprediksi. Aktivitas pada point ini:

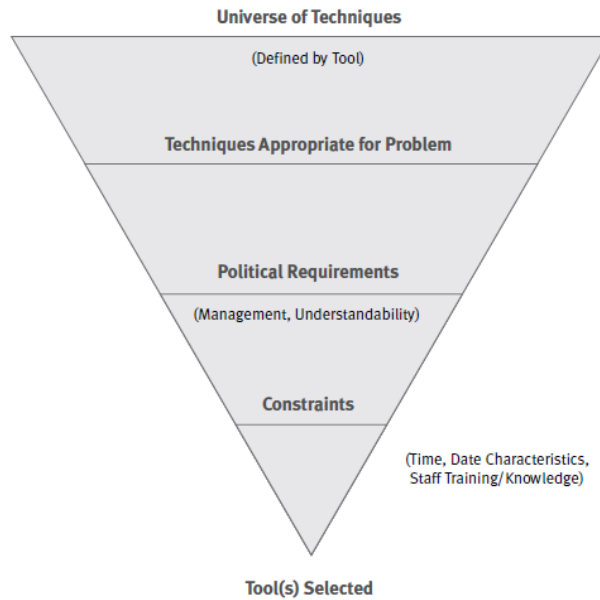
	Aktivitas
<i>Rearranging attributes</i>	Beberapa alat memiliki persyaratan pada urutan atribut, seperti <i>field</i> pertama menjadi unik pengidentifikasi untuk setiap record atau <i>field</i> terakhir menjadi bidang hasil model adalah untuk memprediksi.
<i>Reordering records</i>	Ini mungkin penting untuk mengubah urutan catatan dalam <i>dataset</i> . Mungkin alat pemodelan membutuhkan bahwa catatan diurutkan sesuai dengan nilai hasil atribut.
<i>Reformatted within-value</i>	Perubahan sintaksis yang dibuat untuk memenuhi persyaratan alat pemodelan tertentu
	Mempertimbangkan kembali kriteria seleksi data (bagian 2.1.2.1.3.1) berdasarkan pengalaman pembersihan data.

#### 2.1.2.4 Modelling

##### 2.1.2.4.1 Select modeling technique

Pada tahap *business understanding* telah dipilih alat dan teknik yang akan digunakan. Untuk tahap ini dipilih teknik pemodelan yang lebih spesifik, misalnya, *decision-tree* dengan 5.0, atau *neural network generation* dengan *back*

*propagation*. Jika beberapa teknik diterapkan, lakukan tugas ini secara terpisah untuk masing-masing teknik.



Gambar 2-5 *Universe of Techniques*

Untuk masalah-masalah tertentu, hanya beberapa teknik yang sesuai. "*Political Requirements*" dan *Constraints* dapat membatasi pilihan teknik data mining. Mungkin hanya salah satu alat atau teknik tersedia untuk memecahkan masalah yang dihadapi dan mungkin alat yang dipilih tidak benar-benar yang terbaik, dari segi teknis.

*Output* yang dihasilkan pada tugas ini:

- *Modeling Technique*

Catatan tentang teknik pemodelan yang digunakan. Tentukan teknik yang tepat, mengingat alat yang dipilih.

- *Modeling Assumptions*



Banyak teknik pemodelan membuat asumsi tertentu tentang data, contohnya semua atribut memiliki distribusi seragam, tidak ada nilai-nilai yang hilang diperbolehkan, atribut kelas harus simbolik, dll. Catat asumsi yang dibuat.

#### **2.1.2.4.2 Generate test design**

Sebelum membangun sebuah model, diperlukan prosedur atau mekanisme untuk menguji kualitas dan validitas model. Misalnya, dalam tugas *supervised data mining* seperti klasifikasi, digunakan tingkat kesalahan sebagai ukuran kualitas untuk model data mining. Oleh karena itu, kita biasanya memisahkan *dataset* ke dalam pelatihan dan tes, membangun model di set pelatihan, dan memperkirakan kualitas di set tes terpisah. *Output* yang dihasilkan:

- *Test desain*

Jelaskan rencana yang ditujukan untuk pelatihan, pengujian, dan evaluasi model. Sebuah komponen utama dari rencana adalah menentukan bagaimana membagi data yang tersedia dalam dataset pelatihan, pengujian, dan validasi.

#### **2.1.2.4.3 Build model**

Jalankan alat pemodelan pada *dataset* untuk membuat satu atau lebih model. *Output* yang dihasilkan:

- *Parameter Settings*

Dengan alat pemodelan, sering ada sejumlah besar parameter yang dapat disesuaikan. Buatlah daftar parameter dan nilai-nilai yang mereka pilih, bersama dengan alasan untuk pilihan pengaturan parameter.

- Model

Ini adalah model yang sebenarnya yang dihasilkan oleh alat pemodelan, bukan laporan.

- Model descriptions

Jelaskan model yang dihasilkan.

Laporkan interpretasi model dan dokumen kesulitan yang ditemui dengan maknanya.

#### **2.1.2.4.4 Assess model**

*Engineer data mining* menginterpretasikan model menurut pengetahuan domain-nya, kriteria keberhasilan data mining, dan rancangan yang diinginkan. *Engineer data mining* menilai keberhasilan penerapan pemodelan dan secara teknis teknik penemuan, menghubungkan analisis bisnis dan ahli data mining dalam rangka untuk membahas hasil data mining dalam konteks bisnis. Harap dicatat bahwa tugas ini hanya mempertimbangkan model, sedangkan tahap evaluasi memperhitungkan semua hasil lain yang diproduksi dalam perjalanan proyek. *Engineer data mining* mencoba untuk menentukan peringkat model. Dia menilai model sesuai dengan kriteria evaluasi. Sebaiknya mungkin, ia juga memperhitungkan tujuan bisnis dan kriteria kesuksesan bisnis. Dalam kebanyakan proyek data mining, *Engineer data mining* menerapkan teknik tunggal lebih dari sekali, atau menghasilkan hasil *data mining* dengan beberapa teknik yang berbeda. Dalam tugas ini, ia juga membandingkan semua hasil sesuai dengan kriteria evaluasi. *Output* yang dihasilkan:

- Penilaian Model

Merangkum hasil dari tugas ini, daftar kualitas model yang dihasilkan (misalnya, dalam hal akurasi), dan peringkat kualitas mereka dalam hubungan satu sama lain.

- Revisi pengaturan parameter

Berdasarkan penilaian model, merevisi pengaturan parameter untuk jangka berikutnya ada pada tugas *Build Model*. Iterasi pembangunan dan penilaian model sampai Anda percaya bahwa telah menemukan model yang terbaik. Dokumentasikan semua revisi dan penilaian tersebut.

### **2.1.2.5 Evaluation**

#### **2.1.2.5.1 Evaluation Result**

Langkah-langkah evaluasi ditangani dengan faktor-faktor seperti akurasi dan generalitas model. Langkah ini menilai sejauh mana model memenuhi tujuan bisnis dan berusaha untuk menentukan apakah ada beberapa alasan mengapa model ini kurang. Pilihan lain adalah untuk menguji model dengan testing dalam aplikasi nyata.

Hasil *data mining* melibatkan model yang selalu berhubungan dengan tujuan bisnis dan semua temuan lain yang tidak selalu berhubungan dengan tujuan bisnis asli, tetapi mungkin juga mengungkap tantangan tambahan, informasi, atau petunjuk untuk arah masa depan. *Output* yang dihasilkan:

- Penilaian hasil *data mining* berdasarkan kriteria keberhasilan bisnis.

Meringkas hasil penilaian dalam hal kriteria keberhasilan usaha, termasuk pernyataan akhir mengenai apakah proyek tersebut sudah memenuhi tujuan bisnis awal.

- Persetujuan model.

Setelah menilai model berhubungan dengan kriteria keberhasilan bisnis, model yang dihasilkan yang memenuhi kriteria yang dipilih menjadi model yang disetujui.

### **2.1.2.5.2 Review process**

Pada point ini, model yang dihasilkan tampak memuaskan dan untuk memenuhi kebutuhan bisnis. Tepat untuk melakukan tinjauan lebih menyeluruh tentang keterlibatan *data mining* untuk menentukan apakah ada faktor penting atau tugas yang telah diabaikan. Ulasan ini juga mencakup kualitas sistem jaminan, contoh: Apakah kita benar membangun model? Apakah kita hanya menggunakan atribut yang diijinkan untuk digunakan dan yang tersedia untuk analisis masa depan?

Output yang dihasilkan adalah tinjauan ke bagian proses, ringkasan proses tinjauan dan kegiatan yang telah dilewatkan dan harus diulang.

### **2.1.2.5.3 Determine next steps**

Tergantung pada hasil penilaian dan tinjauan proses, tim proyek memutuskan bagaimana melanjutkannya. Tim memutuskan apakah akan menyelesaikan proyek ini dan beralih ke *deployment*, memulai iterasi lanjut, atau mengatur proyek *data mining* baru. Tugas ini meliputi analisis sumber daya yang tersisa dan anggaran, yang dapat mempengaruhi keputusan. *Output* yang dihasilkan:

- Daftar dari tindakan yang mungkin dilakukan

Daftar potensi tindakan lebih lanjut, bersama dengan alasan setiap pilihan.

- Keputusan  
Jelaskan keputusan tentang bagaimana untuk melanjutkan, bersama dengan dasar pemikiran.

### **2.1.2.6 Deployment**

#### **2.1.2.6.1 Plan deployment**

Tugas ini membutuhkan hasil evaluasi dan menentukan strategi untuk *deploy*. Jika prosedur telah diidentifikasi untuk menciptakan model yang relevan, prosedur ini didokumentasikan di sini untuk penyebaran nanti. *Output* yang dihasilkan:

- *Deployment Plan*  
Merangkum strategi penyebaran, termasuk langkah-langkah yang diperlukan dan cara untuk melakukan itu.

#### **2.1.2.6.2 Plan monitoring and maintenance**

Pengawasan dan pemeliharaan merupakan isu-isu penting jika hasil *data mining* menjadi bagian dari bisnis dan lingkungannya. Persiapan yang cermat dari strategi pemeliharaan membantu untuk menghindari yang tidak perlu dalam jangka waktu yang lama untuk penggunaan yang salah dari hasil *data mining*. Untuk mengawasi penyebaran hasil *data mining*, dibutuhkan rencana proses pengawasan yang rinci. Rencana ini memperhitungkan jenis spesifik dari penyebaran. *Output* yang dihasilkan:

- *Monitoring and maintenance plan*  
Merangkum pengawasan dan strategi pemeliharaan, termasuk langkah-langkah yang diperlukan dan cara untuk melakukan itu.

### **2.1.2.6.3 Produce final report**

Pada akhir proyek, tim proyek menulis sebuah laporan akhir. Tergantung pada rencana penyebaran, laporan mungkin hanya ringkasan proyek dan pengalaman (jika belum didokumentasikan) atau mungkin presentasi akhir dan komprehensif dari hasil *data mining*. *Output* yang dihasilkan:

- *Final Report*

Ini adalah laporan tertulis akhir dari keterlibatan *data mining*. Ini mencakup semua kiriman sebelumnya, meringkas dan mengatur hasil.

- *Final Presentation*

Menyampaikan hasilnya kepada pelanggan.

### **2.1.2.6.4 Review project**

Menilai apa yang benar dan apa yang salah, apa yang telah dilakukan dengan baik dan apa yang perlu ditingkatkan. *Output* yang dihasilkan:

- *Experience Documentation*

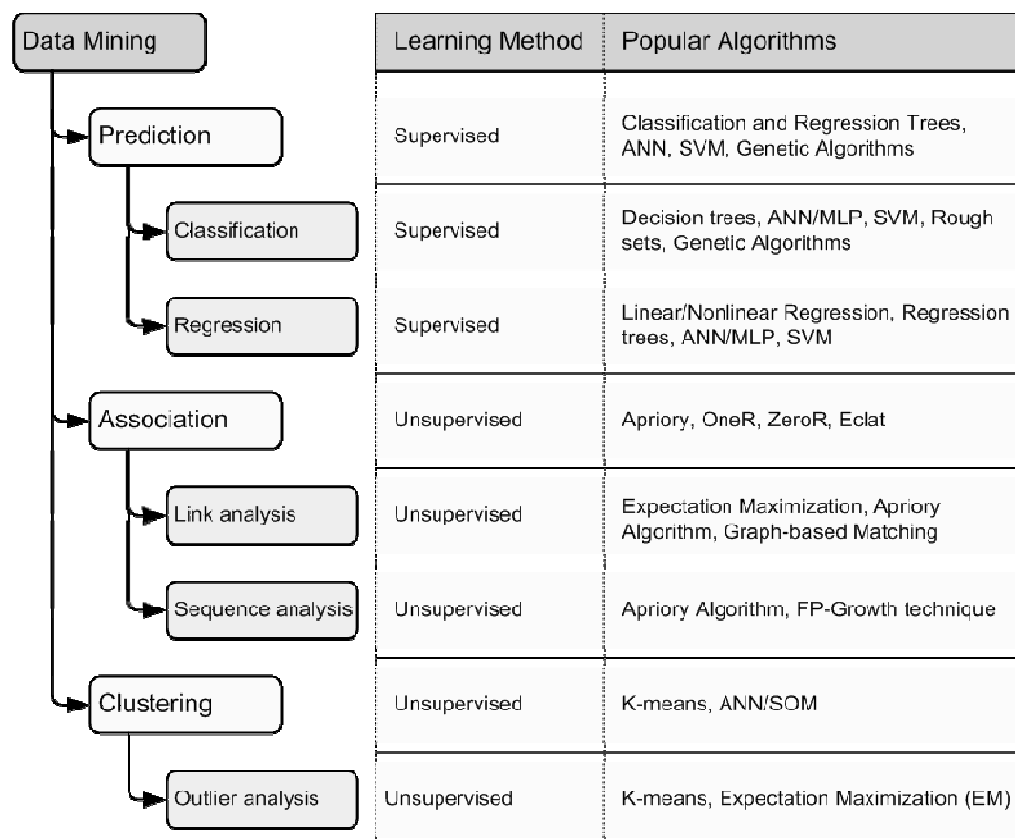
Meringkas pengalaman penting yang diperoleh selama proyek. Sebagai contoh, pendekatan yang menesatkan, atau petunjuk untuk memilih yang terbaik sesuai teknik *data mining* dalam situasi yang sama bisa menjadi bagian dari dokumentasi. Dalam proyek-proyek yang ideal, dokumentasi pengalaman juga mencakup setiap laporan yang telah ditulis oleh anggota proyek individu selama fase proyek.

## **2.1.3 Tugas Data Mining**

Turban, Sharda, Delen, dan King (2011, p.162) menjelaskan bahwa dengan menggunakan data yang relevan, *data mining* membangun model untuk

mengidentifikasi pola-pola di antara atribut yang disajikan dalam *dataset*. Model adalah representasi matematis (hubungan linear sederhana dan/atau hubungan non-linear yang kompleks) yang mengidentifikasi pola-pola di antara atribut dari objek.

Berdasarkan cara untuk mendapatkan pola dari data historis, algoritma pembelajaran dari metode *data mining* terbagi menjadi *supervised* dan *unsupervised*. Dengan algoritma *supervised learning*, data pelatihan mencakup atribut deskriptif (variabel independen) serta atribut kelas (variabel *output* atau variabel hasil). Sebaliknya, dengan *unsupervised learning*, data pelatihan hanya menyertakan atribut deskriptif



Gambar 2-6 Taxonomy berdasarkan tugas *data mining*

(Turban, Sharda, Delen, dan King, 2011)

Secara umum, data mining berusaha untuk mengidentifikasi empat jenis tugas:

- a. Asosiasi menemukan pengelompokan umum dimana terjadi pada saat bersamaan, seperti birdan popok akan bersama-sama dalam analisis *market-basket*.
- b. Prediksi memberitahu yang terjadi di masa depan pada peristiwa tertentu, berdasarkan apa yang telah terjadi masa lalu, seperti meramalkan pemenang *super bowl* atau peramal temperatur hari tertentu.
- c. Cluster mengidentifikasi pengelompokan berdasarkan karakteristik yang diketahui, seperti menempatkan pelanggan dalam segmen yang berbeda berdasarkan demografi dan perilaku pembelian masa lalu.
- d. Hubungan sekuensial menemukan waktu-mengurutkan kejadian, seperti meramalkan bahwa nasabah perbank yang ada, yang sudah memiliki rekening giro akan membuka rekening tabung diikuti oleh akun investasi dalam setahun.

### **2.1.3.1 Klasifikasi**

Turban, Sharda, Delen, dan King (2011, p.178) memaparkan klasifikasi sebagai *machine learning* yaitu pembelajaran untuk menemukan pola dari data *history* (satu set informasi-ciri, variabel, fitur pada karakteristik yang sebelumnya diberi label) dalam rangka untuk menempatkan data baru (tidak diketahui labelnya) ke dalam kelas. Misalnya, menggunakan klasifikasi awan untuk memprediksi apakah cuaca pada hari tertentu akan "cerah", "hujan" atau "berawan". Selain itu tugas klasifikasi digunakan juga untuk persetujuan



keuntungan (misalnya baik atau buruk), lokasi toko (misalnya baik, moderat, buruk), target pemasaran (misalnya kemungkinan menjadi pelanggan atau tidak).

Klasifikasi mempelajari fungsi antara karakteristik (variabel *independen/ input*) dan keanggotaan mereka (variabel *output*) melalui proses *supervised learning*, di mana kedua jenis variabel (*input dan output*) digunakan oleh algoritma. Berikut adalah contoh data yang digunakan pada klasifikasi:

	Variabel Independent				Variabel Output
	Sepal Length (cm)	Sepal Width (cm)	Petal Length (cm)	Petal Width (cm)	Type
1	5.1	3.5	1.4	0.2	Iris setosa
2	4.9	3.0	1.4	0.2	Iris setosa
3	4.7	3.2	1.3	0.2	Iris setosa
4	4.6	3.1	1.5	0.2	Iris setosa
5	5.0	3.6	1.4	0.2	Iris setosa

Gambar 2-7 Contoh data klasifikasi

Dua langkah dalam metodologi klasifikasi tipe prediksi:

- Model pembangunan/pelatihan

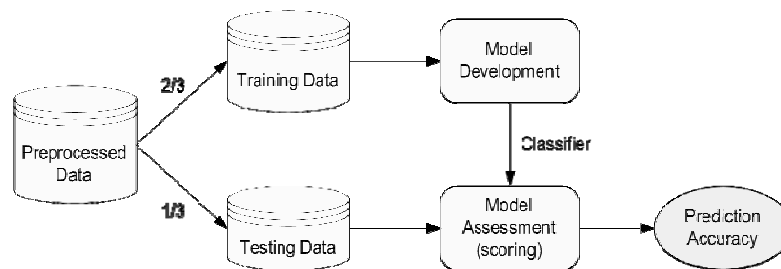
Dalam model pembangunan/pelatihan dikumpulkan data yang akan digunakan terdiri dari variabel input dan label kelas yang sebenarnya. Kemudian dilakukan pelatihan (*training*) menggunakan algoritma yang tersedia.

- Model pengujian/penyebaran.

Setelah dilatih, model diuji terhadap sampel terakhir untuk penilaian akurasi dan digunakan untuk memprediksi kelas sebuah data baru (dimana label kelas tidak diketahui) dengan melihat nilai dari variabel *input* yang dimiliki.

### 2.1.3.1.1 Simple Split

*Simple Split* (atau *holdout* atau estimasi sampel uji) melakukan partisipasi data menjadi dua subset eksklusif disebut set pelatihan dan set uji. Metode ini membagi  $2/3$  dari data sebagai set pelatihan dan sisanya  $1/3$  sebagai set uji. Set pelatihan digunakan oleh *model builder* dan membangun klasifikasi kemudian diuji dengan sampel uji. Aturan ini digunakan dengan asumsi bahwa data yang ada dalam dua subset memiliki jenis yang sama (misalnya, memiliki sifat yang benar-benar sama).



Gambar 2-8 Metode *Simple Split*  
(Turban, Sharda, Delen, dan King, 2011)

### 2.1.3.1.2 Principal Component Analysis (PCA)

Principal Component Analysis (PCA) melakukan pemetaan/transformatasi set data dari dimensi lama ke dimensi baru (yang relative berdimensi lebih rendah) dengan memanfaatkan teknik dalam aljabar linier, tanpa memerlukan masukan parameter tertentu dalam memberikan keluaran hasil pemetaannya.

PCA memerlukan masukan data yang mempunyai sifat zero-mean pada setiap fitur nya. Sifat zero mean pada setiap fitur data bisa didapatkan dengan mengurangi semua nilai dengan rata-ratanya. Set data  $X$  dengan dimensi  $M \times N$ , dimana  $M$  adalah jumlah data dan  $N$  adalah jumlah fitur, akan tampak seperti berikut:

$$X = \begin{bmatrix} x_{11} & x_{12} & x_{1j} & x_{1N} \\ x_{21} & \cdot & \cdot & x_{2N} \\ x_{i1} & \cdot & \cdot & x_{iN} \\ x_{M1} & x_{M2} & x_{Mj} & x_{MN} \end{bmatrix}$$

Untuk fitur ke-j, semua nilai pada kolom tersebut dikurang dengan rata-ratanya, diformulasikan dengan:

$$x_{ij} = x_{ij} - \bar{x}_j$$

$i = 1, 2, \dots, M$  dan  $j$  adalah kolom ke-j

selanjutnya dilakukan perhitungan matriks kovarian dari matriks X, yaitu  $C_x$ . Formula yang digunakan adalah dot-product pada setiap fitur.

$$C_x = \frac{1}{M} X^T \cdot X \dots \dots \dots (3.1)$$

N adalah jumlah fitur, sedangkan  $X^T$  adalah matriks transpos dari X.

$$\begin{aligned} C_x &= \frac{1}{M} x \begin{bmatrix} x_{11} & x_{21} & x_{i1} & x_{M1} \\ x_{12} & \cdot & \cdot & x_{M2} \\ x_{1j} & \cdot & \cdot & x_{Mi} \\ x_{1N} & x_{2N} & x_{iN} & x_{NM} \end{bmatrix} x \begin{bmatrix} x_{11} & x_{12} & x_{1j} & x_{1N} \\ x_{21} & \cdot & \cdot & x_{2N} \\ x_{i1} & \cdot & \cdot & x_{iN} \\ x_{M1} & x_{M2} & x_{Mj} & x_{MN} \end{bmatrix} \\ &= \frac{1}{M} x \begin{bmatrix} x_{11} & x_{21} & x_{1j} & x_{1N} \\ x_{21} & \cdot & \cdot & x_{2N} \\ x_{i1} & \cdot & \cdot & x_{iN} \\ x_{N1} & x_{N2} & x_{Nj} & x_{NN} \end{bmatrix} \end{aligned}$$

Pada matriks  $C_x$ , elemen ke-ij adalah inner-product antara baris matriks  $X^T$  dengan kolom matriks X. Sifat-sifat yang dimiliki oleh matriks  $C_x$  adalah sebagai berikut:

- $C_x$  adalah matriks simetris bujur sangkar berukuran  $N \times N$
- Bagian diagonal utama (dari kiri atas ke kanan bawah) adalah nilai varian masing-masing fitur sesuai dengan indeks kolomnya.
- Bagian selain diagonal utama adalah kovarian di antara pasangan dua fitur yang bersesuaian.

Jadi, matriks  $C_x$  menangkap kovarian diantara semua pasangan yang mungkin dari fitur data set data matriks X. Nilai kovarian merefleksikan noise dan redundansi pada fitur:

- Dalam diagonal utama, asumsinya adalah nilai yang tinggi berkorelasi dengan struktur yang penting.
- Dalam elemen selain diagonal utama, nilai jarak yang besar menandakan redundansi yang tinggi.

Jika Y adalah matriks set data hasil pemetaan dan  $C_y$  adalah matriks kovarian dari Y, yang diharapkan dalam PCA adalah:

- Semua elemen selain diagonal utama dalam  $C_y$  harus nol. Maka,  $C_y$  harus matriks diagonal. Dengan kata lain, Y adalah matriks terdekorelasi.
- Peletakan dimensi dalam Y dari kiri ke kanan diurutkan secara menurun.

Cara yang umum digunakan untuk mendapatkan  $C_y$  adalah dengan eigenvalue dan eigenvector. Eigenvalue dan eigenvektor dari matriks X berturut-turut adalah nilai skala  $\lambda$  dan vector  $u$  yang memenuhi persamaan berikut:

$$Xu = \lambda u$$

Dengan mencari matriks ortonormal P dimana  $Y = PX$  dan  $C_y = \frac{1}{M} YY^T$  adalah matriks diagonal, dan kolom dari P adalah komponen utama (principal components) dari X, persamaan  $C_y$  bisa dijabarkan sebagai berikut:

$$\begin{aligned} C_y &= \frac{1}{M} YY^T \\ &= \frac{1}{M} (PX)(PX)^T \\ &= \frac{1}{M} PXX^T P^T \end{aligned}$$

$$= P \left( \frac{1}{M} XX^T \right) P^T$$

Dengan mensubstitusikan persamaan 3.1, kita mendapatkan matriks  $C_y$  berdimensi  $N \times N$ :

$$C_y = P C_x P^T$$

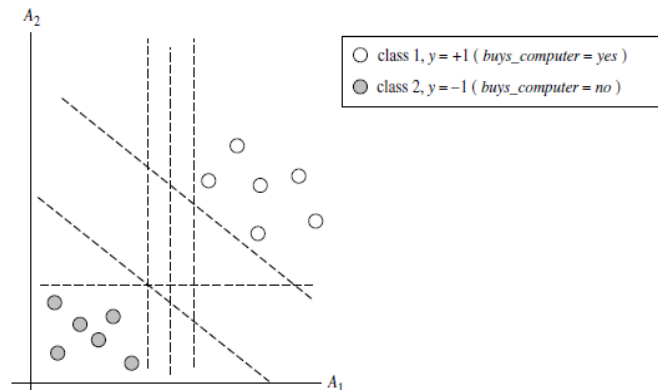
Setiap baris matriks  $P$  adalah eigen vector  $C_x$

### 2.1.3.1.3 Algoritma Klasifikasi

#### 2.1.3.1.3.1 Support Vector Machine(SVM)

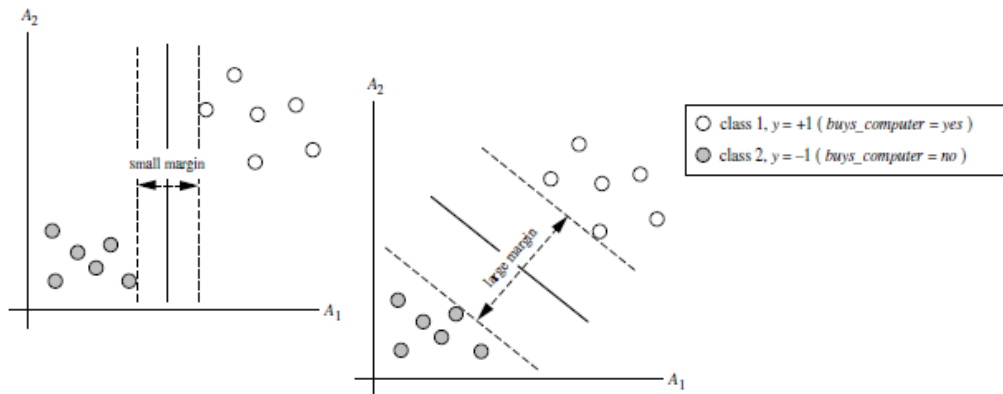
SVM adalah algoritma klasifikasi baik linier maupun nonlinear, pola data yang kompleks dan menghindari *overfitting*, situasi di mana model menghafal polanya yang relevan dengan kasus-kasus analisis tertentu. Han dan Kamber (2006, p.337-343) memaparkan cara kerja algoritma SVM dengan pemetaan nonlinier yaitu mengubah data pelatihan asli ke dimensi yang lebih tinggi. Dalam hal ini, dimensi baru akan mencari linier optimal yang memisahkan data ("decision boundary" memisahkan tupel dari satu kelas dengan yang lain). SVM menemukan *hyperplane* menggunakan *support vektor* ("esensial" pelatihan tuple) dan margin.

Berikut adalah penjelasan SVM dalam dua kelas dimana kelas yang dipisahkan secara linear. Kumpulan data  $D$  terdiri dari  $(x_1, y_1), (x_2, y_2), \dots, (x_{|d|}, y_{|d|})$ , dimana  $X_i$  adalah himpunan tupel pelatihan dengan label kelas terkait,  $Y_i$ . Setiap  $Y_i$  dapat mengambil salah satu dari dua nilai, baik  $+1$  atau  $-1$  (yaitu,  $Y_i \in \{+1, -1\}$ ). Untuk membantu dalam visualisasi, digunakan dua atribut masukan,  $A_1$  dan  $A_2$ , seperti gambar berikut:



Gambar 2-9 Data *training* dua dimensi (Han & Kamber, 2006)

Dari grafik, kita melihat bahwa 2-D data linear, garis lurus dapat ditarik untuk memisahkan semua tuple kelas +1 dari semua tuple kelas -1. Sebuah pendekatan untuk masalah menemukan *hyperplane* terbaik adalah dengan mencari MMH.



Gambar 2-10 Kemungkinan dua pemisah *hyperplane* dan margin yang terkait (Han & Kamber, 2006)

Kedua *hyperplane* dapat mengklasifikasikan semua data tuple yang diberikan. *Hyperplane* dengan margin yang lebih besar lebih akurat mengklasifikasikan masa depan data tuple daripada *hyperplane* dengan margin yang lebih kecil.

Sebuah *hyperplane* pemisah dapat ditulis sebagai:

$$WX + b = 0$$

W adalah vektor bobot, yaitu,  $W = \{w_1, w_2, \dots, w_n\}$ ,  $n$  adalah jumlah atribut, dan  $b$  adalah skalar, sering disebut sebagai bias. Pada

tupel pelatihan 2-D, misalnya,  $X = (x_1, x_2)$ , dimana  $x_1$  dan  $x_2$  adalah nilai-nilai atribut  $A_1$  dan  $A_2$ , masing-masing. Jika kita berpikir tentang bobot tambahan,  $w_0$ , kita dapat menulis ulang *hyperplane* menjadi:

$$w_0 + w_1x_1 + w_2x_2 = 0$$

Dengan demikian, setiap titik yang terletak di atas *hyperplane* memenuhi:

$$w_0 + w_1x_1 + w_2x_2 > 0$$

Demikian pula, setiap titik yang terletak di bawah *hyperplane*:

$$w_0 + w_1x_1 + w_2x_2 < 0$$

Bobot dapat disesuaikan sehingga *hyperplane* mendefinisikan "sisi" dari margin ditulis sebagai:

$$H_1 : w_0 + w_1x_1 + w_2x_2 \geq 1 \text{ untuk } y_i = +1$$

dan

$$H_2 : w_0 + w_1x_1 + w_2x_2 \leq -1 \text{ untuk } y_i = -1$$

Kedua persamaan di atas digabungkan, maka menghasilkan:

$$y_i (w_0 + w_1x_1 + w_2x_2) \geq 1, \forall i.$$

Setiap tuple pelatihan yang jatuh pada *hyperplane*  $H_1$  atau  $H_2$  (yaitu, "sisi" yang mendefinisikan margin) yang memenuhi persamaan disebut *vektor Support*. Artinya, mereka sama-sama dekat dengan MMH. Untuk ukuran margin maksimal, jarak dari *hyperplane* untuk setiap titik pada  $H_1$  adalah:

$$\frac{1}{\|w\|}$$

Dimana  $\|W\|$  adalah *Euclidean norm* dari  $W$ , yaitu  $\sqrt{W \cdot W}$  (jika

$W = \{w_1, w_2, \dots, w_n\}$  maka  $\sqrt{W \cdot W} = \sqrt{w_1^2 + w_2^2 + \dots + w_n^2}$ ).

Dengan definisi, ini sama dengan jarak dari setiap titik pada  $H_2$  ke *hyperplane*.

Oleh karena itu, margin maksimal adalah:

$$\frac{2}{\|W\|}$$

Setelah menemukan *support vektor* dan MMH, kami memiliki *support vektor* machine terlatih. MMH adalah batas kelas linear, sehingga SVM yang sesuai dapat digunakan untuk mengklasifikasikan data linear.

Pendekatan yang dijelaskan untuk SVM linear dapat diperluas untuk membuat SVM nonlinear. Ada dua langkah utama. Pada langkah pertama, mengubah input data asli ke dalam ruang dimensi yang lebih tinggi menggunakan pemetaan nonlinear. Beberapa pemetaan nonlinear umum dapat digunakan dalam langkah ini. Setelah data telah berubah menjadi ruang baru yang lebih tinggi, pencarian langkah kedua untuk *hyperplane* pemisahan linear dalam ruang baru.

Optimasi *quadratic* dapat diselesaikan dengan menggunakan perumusan linear SVM. MMH ditemukan di ruang yang baru sesuai dengan *hypersurface* nonlinear dalam ruang aslinya.

Sebuah masukan 3D vector  $X = (x_1, x_2, x_3)$  dipetakan ke dalam ruang 6D,  $Z$ , dengan menggunakan pemetaan  $\varphi_1(X) = x_1$ ,  $\varphi_2(X) = x_2$ ,  $\varphi_3(X) = x_3$ ,  $\varphi_4(X) = (x_1)^2$ ,  $\varphi_5(X) = x_1 x_2$ , dan  $\varphi_6(X) = x_1 x_3$ . Sebuah *hyperplane* dalam ruang baru adalah:

$$d(Z) = WZ + b$$

$W$  dan  $Z$  adalah vektor. Kami memecahkan  $W$  dan  $b$  dan kemudian mengganti kembali sehingga *hyperplane* keputusan linear dalam ruang ( $Z$ ) yang baru sesuai dengan nonlinier- orde kedua polinomial dalam ruang 3-D masukan asli:

$$\begin{aligned} d(Z) &= w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 (x_1)^2 + w_5 x_1 x_2 + w_6 x_1 x_3 + b \\ &= w_1 z_1 + w_2 z_2 + w_3 z_3 + w_4 z_4 + w_5 z_5 + w_6 z_6 + b \end{aligned}$$



Tupel pelatihan hanya muncul dalam bentuk titik,  $\Phi(X_i) \cdot \Phi(X_j)$ , di mana  $\Phi(X)$  hanya fungsi pemetaan non-linear diterapkan untuk mengubah tupel pelatihan. Daripada menghitung perkalian titik pada tupel data yang berubah, ternyata secara matematis setara dengan menerapkan fungsi kernel,  $K(X_i, X_j)$ , dengan data masukan yang asli. Artinya,

$$K(X_i, X_j) = \Phi(X_i) \cdot \Phi(X_j)$$

Tiga fungsi kernel yang dapat diterima meliputi:

*Polynomial kernel of degree h:*  $K(X_i, X_j) = (X_i \cdot X_j + 1)^h$

*Gaussian radial basis function kernel:*  $K(X_i, X_j) = e^{-\|X_i - X_j\|^2 / 2\sigma^2}$

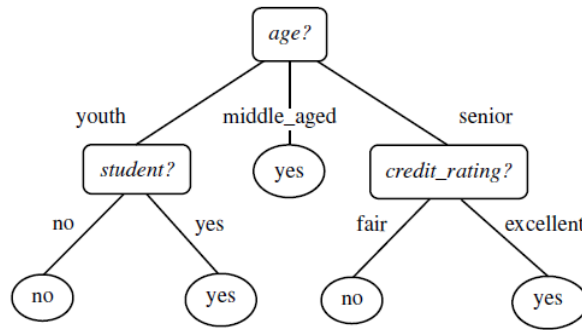
*Sigmoid kernel:*  $K(X_i, X_j) = \tanh(K(X_i, X_j) \cdot \delta)$ .

Kelebihan pada SVM adalah solusi yang diberikan adalah global optima, tidak seperti Artificial Neural Network yang solusinya sering masuk pada wilayah local optima. Hal ini berarti SVM selalu memberikan model yang sama dan solusi dengan margin maksimal.

### **2.1.3.1.3.2 Decision Tree**

*Decision tree* merupakan bagian dari "machine learning", bagian penting dari kecerdasan buatan (Lin, Shiue, Chen, & Cheng, 2009). Metode ini mendekati angka atau simbol-fungsi nilai sasaran yang kuat untuk data yang *noise* dan

mampu belajar ekspresi disjungtif. Sedangkan Han dan Kamber (2006, p.337-343) menjelaskan *Decision tree* adalah flowchart- seperti struktur pohon, di mana setiap simpul internal (simpul tan pada daun) menunjukkan pengujian pada atribut, setiap cabang merupakan hasil tes, dan setiap simpul daun (atau simpul terminal) memiliki label kelas. Simpul yang paling atas adalah simpul akar.



Gambar 2-11 Struktur Decision Tree

Satu set sampel dalam partisi  $S$ , atribut  $X$  dipilih untuk partisi set ke  $S_1, S_2, \dots, S_L$  dan ini ditambahkan ke pohon keputusan sebagai anak-anak dari simpul untuk  $S$ . Simpul untuk  $S$  diberi label dengan uji  $X$ , dan partisi  $S_1, S_2, \dots, S_L$  kemudian rekursif dipartisi. Berikut langkah dan perhitungan dalam algoritma *decision tree*:

- Hitung  $Info(S)$  untuk mengidentifikasi kelas pada dataset pelatihan  $S$ .

$$Info(S) = - \sum_{i=1}^k \left\{ \left[ \frac{freq(C_i, S)}{|S|} \right] \log_2 \left[ \frac{freq(C_i, S)}{|S|} \right] \right\}$$

Dimana  $|S|$  adalah jumlah kasus pada pelatihan set.  $C_i$  adalah kelas,  $i=1, 2, \dots, k$ .  $k$  adalah jumlah kelas dan  $freq(C_i, S)$  adalah nomor dari kasus yang termasuk dalam  $C_i$ .

- Hitung nilai informasi yang diharapkan,  $Info_x(S)$  untuk fitur  $X$  ke partisi  $S$ .

$$Info_x(S) = - \sum_{i=1}^L \left[ \left( \frac{|S_i|}{|S|} \right) info(S_i) \right],$$

dimana  $L$  adalah jumlah output untuk fitur  $X$ ,  $S_i$  adalah bagian dari  $S$  yang sesuai dengan output ke  $i$  dan  $|S_i|$  adalah jumlah kasus yang subset  $S_i$ .

- Hitung information gain setelah mempartisi berdasarkan fitur  $X$ .

- Hitung nilai informasi partisi,  $\text{SplitInfo}(X)$  diperoleh dari  $S$  yang dipartisi menjadi subset  $L$ .

- Hitung rasio gain dari  $\text{Gain}(X)$  terhadap  $\text{SplitInfo}(X)$

$$\text{GainRatio}(X) = \text{Gain}(X) / \text{SplitInfo}(X)$$

The  $\text{GainRatio}(X)$  mengimbangititik lemah dari  $\text{Gain}(X)$ , yang merupakan jumlah dari informasi yang diberikan oleh  $X$  dalam pelatihan set. Oleh karena itu, fitur dengan  $\text{GainRatio}$  tertinggi ( $X$ ) diambil sebagai akar pohon keputusan.

#### 2.1.3.1.4 Akurasi Metode Klasifikasi

Dalam klasifikasi, perhitungan akurasi menggunakan *confusion matrix*, disebut juga *classification matrix* atau *contingency table* (Gambar 2-12).

		True Class	
		Positive	Negative
Predicted Class	Positive	True Positive Count (TP)	False Positive Count (FP)
	Negative	False Negative Count (FN)	True Negative Count (TN)

Gambar 2-12 *Confusion Matrix*  
(Turban, Sharda, Delen, dan King, 2011)

Angka sepanjang diagonal dari kiri atas ke kanan bawah mewakili keputusan yang benar, dan jumlah sisi luardiagonal mewakili kesalahan. Dari data *confusion matrix*, dapat dihitung akurasi umum untuk model klasifikasi.

Tabel 2-1 Perhitungan akurasi untuk metode klasifikasi

Metric	Description
True Positive Rate = $\frac{TP}{TP+FN}$	The ratio of correctly classified positives divided by total positive count (i.e., hit rate or recall)
True Negative Rate = $\frac{TN}{TN+FP}$	The ratio of correctly classified negative divided by the total negative count (i.e., false alarm rate)
Accuracy = $\frac{TP+TN}{TP+TN+FP+FN}$	The ratio of correctly classified instances (positive and negative) divided by total number of instances

Perhitungan akurasi model klasifikasi (atau *classifier*) dengan algoritma *supervised learning* penting, karena dua alasan. Pertama, dapat digunakan untuk memperkirakan akurasi prediksi masa depan, yang berarti tingkat kepercayaan terhadap satu *output classifier* dalam sistem prediksi. Kedua, dapat digunakan untuk memilih *classifier* (mengidentifikasi model klasifikasi terbaik).

## 2.2 Data Warehouse

Vercellis (2009, p.45) menjelaskan bahwa *data warehouse* merupakan repositori utama untuk menyediakan data dalam mengembangkan arsitektur *business intelligencedan decision support systems*. Dimana *data warehouse* memiliki karakteristik: data statis, data saat ini dan lampau, data agregasi dan konsolidasi, dan berfungsi untuk analisis.

Sedangkan *data mart* berisi semua data yang dibutuhkan oleh suatu departemen tertentu dalam perusahaan untuk melakukan analisis *business intelligencedan* mengeksekusi aplikasi *decision support* sesuai fungsinya sendiri.

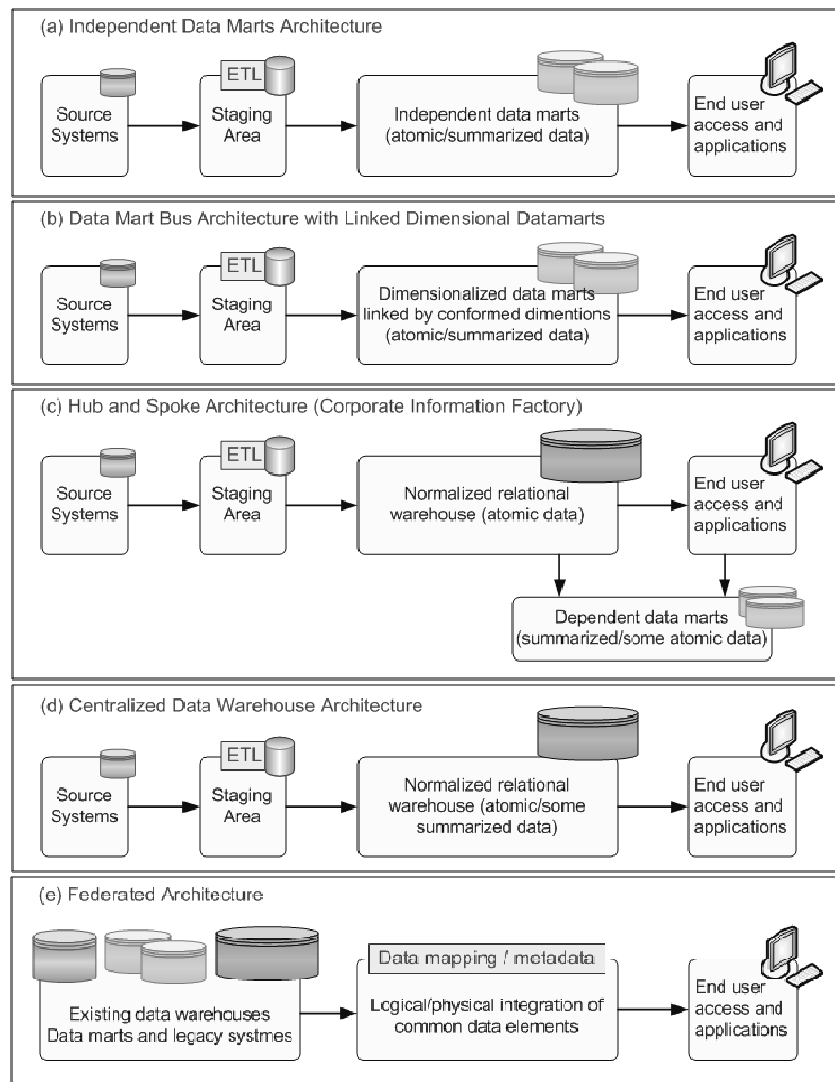
Dari pengertian di atas, *data mart* dapat dianggap sebagai fungsionalitas *data warehouse* departemen, dengan ukuran yang lebih kecil dan jenis yang lebih spesifik daripada *data warehouse*.

Ada beberapa alternatif arsitektur *data warehouse* (Turban, Sharda, Delen, dan King, 2011) yang diusulkan dalam merancang dan mengembangkan *data warehouse*, yaitu:

- a. *Independent data mart*, arsitektur ini bisa dibayangkan paling sederhana dan alternatif arsitektur paling mahal. *Data mart* yang dikembangkan memiliki tujuan beroperasi secara *independent* dari satu sama lain untuk melayani kebutuhan unit organisasi masing-masing. Karena kemandirian, memungkinkan untuk memiliki definisi yang tidak konsisten, dimensi dan ukuran data yang berbeda, sehingga sulit untuk menganalisis data seluruh *data mart*.
- b. *Data mart bus architecture*, arsitektur ini merupakan alternatif untuk *data mart independent* dimana individu dihubungkan satu sama lain melalui beberapa jenis *middleware*. Karena data terkait antara individu *data mart*, ada kesempatan yang lebih baik menjaga konsistensi data seluruh perusahaan. Meskipun memungkinkan untuk permintaan data yang kompleks di seluruh *data mart*, kinerja untuk jenis analisis mungkin tidak pada tingkat yang memuaskan.
- c. *Hub-and-spoke architecture*, pada arsitektur ini diperhatikan dan difokuskan pada pembangunan skalabilitas dan kemampuan *maintain* infrastruktur yang mencakup *data warehouse* terpusat dan beberapa *dependent data mart* (masing-masing untuk unit organisasi). Arsitektur ini memungkinkan

untuk kemudahan dan kustomisasi antarmuka penggunaan laporan. Di sisi negatif, arsitektur ini memiliki pandangan menyeluruh perusahaan, dan dapat menyebabkan redundansi data dan latency data.

- d. *Centralized data warehouse*, arsitektur *data warehouse* terpusat mirip dengan *hub-and-spoke* architecture kecuali tidak ada *data mart dependent*, melainkan ada *data warehouse* raksasa yang melayani untuk kebutuhan semua unit organisasi. Pendekatan terpusat menyediakan pengguna dengan akses ke semua data dalam *data warehouse*. Selain itu, mengurangi jumlah data tim teknis yang harus mentransfer atau mengubah, sehingga menyederhanakan pengelolaan data dan administrasi. Jika dirancang dan dilaksanakan dengan baik, arsitektur ini memberikan pandangan tepat waktu dan menyeluruh dari perusahaan untuk siapa, kapan, dan di mana pun mereka berada dalam organisasi.
- e. *Federated data warehouse*, menggunakan segala cara yang mungkin untuk mengintegrasikan sumber daya analitis dari berbagai sumber untuk memenuhi perubahan kebutuhan atau kondisi bisnis. Pendekatan federasi didukung oleh *vendor middleware* yang mengusulkan permintaan distribusi dan kemampuan untuk bergabung. Karena masalah kualitas kinerja dan data, kebanyakan ahli sepakat bahwa pendekatan federasi bekerja dengan baik untuk melengkapi *data warehouse*, tidak menggantikannya.



Gambar 2-13 Arsitektur Data Warehouse  
(Turban, Sharda, Delen, dan King, 2011)

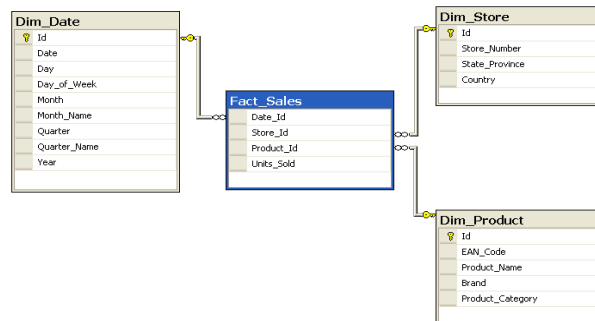
Beberapa perusahaan lebih memilih untuk merancang dan mengembangkannya dengan cara bertahap serangkaian *datamart* daripada sebuah *centralized data warehouse*, dalam rangka untuk mengurangi waktu pelaksanaan.

Step dalam pembuatan *data warehouse* yang paling utama adalah proses untuk ETL (*Extract Transform Load*). *Extract* yang dimaksud adalah proses untuk pengambilan data dari database sumber. *Transform* merupakan proses untuk

perubahan data yang telah diambil menyesuaikan dengan kebutuhan analisa. Dan *Load* adalah proses untuk menaikkan data hasil olahan ketarget.

Dalam proses pembuatan *data warehouse*, diperlukan untuk mendesain schema dari database. Schema ini biasanya terdiri dari *fact* dan *measurements*, dimana *Fact* adalah sekumpulan data real yang sudah ada di sumber database, sedangkan *Measurement* adalah data olahan yang terdiri atas satu atau sekumpulan *fact* yang telah diolah menjadi suatu point tertentu.

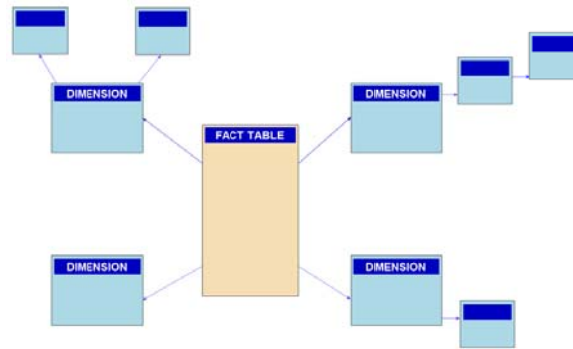
Design *data warehouse* terdiri dari *Star Schema* dan *Snowflake Schema*. *Star schema* adalah skema yang paling banyak digunakan untuk desain *Data Warehouse*, dimana sekumpulan *Fact* dan *Measurements* dikelilingi oleh dimensi.



Gambar 2-14 Star Schema

Snowflake Schema adalah skema lanjutan dari *star scheme*, dimana dimensi yang ada *data warehouse* bisa berlapis dikarenakan data dimensi tersebut sudah dinormalisasi.





Gambar 2-15 Snowflake Schema

## 2.3 Multi Finance dan Kredit

Perusahaan *Multifinance* merupakan perusahaan yang berperan serta dalam kegiatan yang berhubungan dengan produk-produk yang ada dalam pasar *multi finance*. Produk-produk yang ada yaitu sewa guna usaha atau *leasing* (salah satu cara perusahaan memperoleh asset atau kepemilikan tanpa harus melalui proses yang berkepanjangan), anjak piutang, dan *consumer financing*.

Anjak piutang merupakan salah satu instrument yang sering disebut dengan *factoring*, didefinisikan sebagai transaksi pembelian, penagihan serta pengurusan piutang atau tagihan jangka pendek klien (penjual) kepada perusahaan *factoring*. Kemudian akan ditagih oleh perusahaan anjak piutang kepada pembeli karena adanya pembayaran kepada klien oleh perusahaan *factoring*.

*Consumer financing* atau pembiayaan konsumen dimana dalam hal ini ada perusahaan yang bertindak sebagai penjamin dalam pemberian kredit kepada konsumen yang tidak menjadi satu kesatuan dalam perbankan melainkan perusahaan ini berdiri sendiri. Pembiayaan konsumen ini adalah suatu pinjaman atau kredit yang diberikan oleh suatu perusahaan kepada debitur untuk pembelian barang atau jasa yang akan langsung digunakan oleh konsumen.

Kredit berasal dari bahasa Italia yaitu kata *credere* yang memiliki arti percaya, dimana pihak yang mendapat kredit berarti orang tersebut dipercaya untuk mendapatkan pinjaman. Kasmir (2008 pp 103-105) menjelaskan unsur-unsur dalam kredit adalah:

- a. Kepercayaan, yang berarti bahwa pemberi kredit yakin prestasi yang diberikannya baik dalam bentuk uang, barang, atau jasa, akan benar-benar diterimanya kembali dalam jangka waktu tertentu di masa yang akan datang.
- b. Kesepakatan, di mana dituangkan dalam suatu perjanjian dan masing-masing pihak menandatangani hak dan kewajiban masing-masing.
- c. Jangka Waktu, dimana mencakup masa pengembalian kredit yang telah disepakati.
- d. Resiko, faktor resiko kerugian dapat diakibatkan dua hal yaitu resiko kerugian yang diakibatkan nasabah sengaja tidak mau membayar kreditnya padahal mampu dan resiko kerugian yang diakibatkan karena nasabah tidak sengaja yaitu akibat terjadinya musibah seperti bencana alam. Penyebab tidak tertagih sebenarnya dikarenakan adanya suatu tenggang waktu pengembalian (jangka waktu). Semakin panjang jangka waktu suatu kredit semakin besar risikonya tidak tertagih, demikian pula sebaliknya.
- e. Balas Jasa, dimana dalam bentuk bunga, biaya provisi, dan komisi serta biaya administrasi kredit ini merupakan keuntungan Lembaga Keuangan.

## **2.4 Analisis Kredit, Resiko Kredit dan Credit Scoring**

Analisis kredit merupakan proses analisis atau menilai permohonan kredit yang diajukan dengan menggunakan pendekatan-pendekatan dan rasio-rasio keuangan sehingga dapat memberikan keyakinan kepada pihak penyedia kredit

bahwa proyek kredit yang diajukan cukup layak untuk disetujui. Adapun tujuan dari analisis kredit adalah mengurangi resiko kredit yang ada.

Purwanto (2011) menjelaskan resiko kredit sebagai resiko terjadinya kerugian akibat kegagalan pembayaran oleh debitur atau terjadinya kemerosotan kualitas kemampuan membayar pihak debitur. Resiko kredit menjadi kritis apabila nasabah dengan pemberian dana yang besar tidak dapat membayar pinjamannya sehingga menimbulkan kerugian dan kesulitan dalam likuiditas.

Sabato (2010) menulis bahwa sejak tahun 1960 organisasi besar telah memanfaatkan *credit scoring* dengan cepat dan menilai tingkat risiko dari prospek pelamar dan konsumen yang sudah ada terutama dalam bisnis kredit konsumen. *Credit scoring* memprediksi kemungkinan bahwa calon konsumen akan sesuai dengan standar atau menunggak selama jangka waktu yang tetap dengan tujuan untuk menentukan pinjaman kredit. Sabato (2010) mendeskripsikan bahwa model ini memberikan peringkat risiko calon konsumen berdasarkan asumsi jika konsumen yang ada memiliki perilaku tertentu, ada kemungkinan calon konsumen dengan karakteristik serupa menunjukkan perilaku yang sama.