

Desain Data Warehouse dan Implementasi Data Mining Terhadap Data Nilai Mahasiswa

Mardiani
Sistem Informasi
STMIK MDP
Indonesia
mardiani@stmik-mdp.net

Abstrak— Data mahasiswa memiliki banyak informasi jika ditelusuri lebih lanjut lagi. Mulai dari data mahasiswa baru, data nilai mahasiswa sampai ke data alumni. Data yang jumlahnya banyak akan menjadi tidak berguna jika tidak di gunakan dengan tepat. Desain *data warehouse* dan implementasi *data mining* bertujuan untuk mendukung manajemen dalam mengambil keputusan yang sesuai dengan kepentingan perguruan tinggi. Data yang dapat mendukung keputusan perguruan tinggi tersebut salah satunya adalah data nilai mahasiswa, dimana didalam data nilai tersebut banyak informasi tersembunyi yang dapat digali. Data yang digunakan berupa data KHS, mata kuliah dan mahasiswa yang memiliki beberapa atribut untuk setiap tabel. Metodologi yang digunakan adalah sembilan tahapan dalam metodologi desain *database* untuk *data warehouse*. Hasil yang diperoleh dalam bentuk skema bintang sebagai hasil perancangan *data warehouse* dan beberapa hasil terkait nilai mahasiswa, menggunakan dua fungsionalitas *data mining* yaitu asosiasi dan *clustering*, dimana dari dua fungsionalitas tersebut diketahui hubungan antar atribut dan didapat kelompok-kelompok mahasiswa berdasarkan nilai.

Keywords—KHS; Mahasiswa; Asosiasi; Clustering;

I. PENDAHULUAN

Dengan semakin berkembangnya teknologi informasi, ilmu komputer cabang *data warehouse* dan *data mining* juga semakin mengambil bagian. *Data mining* adalah analisis dan pengamatan data dalam jumlah besar untuk menemukan hubungan yang tidak diketahui sebelumnya dan metode baru untuk meringkas data menjadi lebih mudah dimengerti dan berguna bagi pemilik data. Mengubah data menjadi informasi dengan menggunakan *data mining* membuat ketersediaan informasi bagi banyak orang menjadi berlimpah. Hasil penggunaan teknologi informasi di hampir semua bidang kehidupan, harus mampu dimanfaatkan dalam bentuk informasi dan pengetahuan.

Dalam dunia pendidikan, *data warehouse* dan *data mining* telah diterapkan untuk berbagai kebutuhan. Kampus memiliki berbagai data, misalnya data akademik, pemasaran dan keuangan, termasuk data mahasiswa, dosen, dan staf. Saat ini dunia pendidikan telah menggunakan *data warehouse* dan *data mining* sebagai alat untuk mendukung keputusan untuk manajemen data mahasiswa, terutama KHS, banyak informasi

yang tersembunyi yang dapat diekstraksi dari data KHS tersebut.

II. TINJAUAN PUSTAKA

Menurut [1], data persediaan buku dapat dirancang dan diolah dengan data warehouse. Kemudian menurut [2], data penjualan tentang penyewaan mobil secara lengkap juga dapat dibuat dan dirancang dengan menggunakan *data warehouse*. Sedangkan dari dunia pendidikan, ada banyak hal yang bisa di eksplorasi, terutama di tingkat perguruan tinggi. Diantaranya adalah mencari informasi tentang pemasaran data mahasiswa baru. Hal ini diantaranya telah dilakukan [3] untuk memanfaatkan data warehouse di perguruan tinggi Indonesia, dan [4] yang membuat rancang bangun data warehouse untuk menunjang evaluasi akademik di fakultas.

A. Pra Proses Data

Preprocessing data adalah proses yang harus dilakukan sebelum memasuki tahap pembuatan *data warehouse*. Data yang digunakan adalah bersifat redundan, tidak lengkap, dan tidak konsisten. Berikut adalah tahap *preprocessing* data sesuai dengan [5]:

- Integrasi data: Merupakan integrasi data dari berbagai sumber ke dalam penyimpanan data tunggal untuk data kesatuan yang koheren.
- Reduksi data: Teknik pengurangan data yang diterapkan untuk memperoleh representasi mengurangi jumlah data yang menggambarkan volume yang jauh lebih kecil.
- Pembersihan Data: Proses ini adalah langkah pembersihan data, yaitu untuk mengisi data yang hilang, memperbaiki data kotor dan rusak, mengidentifikasi atau menghapus data *outlier*, memperbaiki data yang tidak konsisten.
- Transformasi Data: Transformasi data adalah proses mengubah data menjadi bentuk yang sesuai. Proses ini dilakukan agar data tetap konsisten dan dapat digunakan untuk proses selanjutnya.

Data warehouse dibangun untuk mengatasi masalah teknis dan bisnis yang terkait dengan penggunaan data dan informasi [5]. *Data warehouse* adalah kumpulan data dari berbagai

sumber yang ditempatkan bersama-sama di tempat penyimpanan yang besar dan kemudian diolah menjadi bentuk penyimpanan multi dimensi dan dirancang untuk *query* dan pelaporan [6].

Secara umum, *data warehouse* menyimpan data historis selama beberapa tahun atau periode dan akan digunakan *query* untuk keperluan bisnis intelijen atau kegiatan analisis lainnya. Data untuk *data warehouse* didapat dari berbagai sumber data, baik internal maupun eksternal, dibersihkan dan diorganisir secara konsisten dan disesuaikan dengan kebutuhan perusahaan atau organisasi. Setelah data ditempatkan di *data warehouse*, data dapat ditransfer ke bagian-bagian atau departemen-departemen. [7].

B. Karakteristik Data warehouse

Bill Inmon, mendefinisikan: Sebuah *data warehouse* memiliki karakteristik-karakteristik, yaitu berorientasi subjek, terintegrasi, memiliki variasi waktu dan *nonvolatile*:

- Berorientasi subjek: *Data warehouse* dirancang untuk membantu dalam menganalisis data dalam jumlah besar secara keseluruhan.
- Data yang terintegrasi: Data yang terdapat dalam *data warehouse* dapat berasal dari beberapa sumber dimana semua data akan disimpan ke bagian yang sama dengan format khusus.
- Time Variant: Karena analisis - analisis yang dilakukan cenderung mengarah ke arah "analisis tren", sehingga data yang tersedia dalam jumlah besar dapat dikatakan akurat atau valid pada saat itu.
- Data Nonvolatile: Mengingat tujuan perusahaan adalah untuk membangun *data warehouse* untuk menganalisis sebuah kasus, maka data yang terkandung di dalamnya sudah tidak dapat diubah lagi.

C. Data Staging

Data Staging merupakan langkah yang harus dilakukan sebelum data yang ada disimpan dalam *data warehouse*. Seluruh data yang dikumpulkan dari sumber yang berbeda dan kemudian dimodifikasi, dan diubah menjadi format yang sama dan dirancang untuk memenuhi permintaan, penyimpanan dan analisis [8].

- Ekstraksi: Pada tahap ini proses seleksi juga dilakukan, data yang akan dimasukkan ke dalam *data warehouse* dipilih, karena tidak semua data dapat digunakan, yaitu hanya data yang dibutuhkan. Juga dilakukan proses pembersihan data pada tahap ini, untuk mencegah kesalahan penulisan dalam sumber data.
- Transformasi: Merupakan proses perubahan data mentah menjadi data siap pakai untuk dimasukkan ke dalam *data warehouse*, ada beberapa tugas yang dilakukan pada tahap ini, seperti pemilihan, penggabungan dan pemisahan, mengubah, dan meringkas data.
- Load: Adalah tahapan di mana semua data yang telah disiapkan ke dalam *data warehouse*. Loading awal

dilakukan ketika pengembang telah menyelesaikan konstruksi dan desain *data warehouse*.

D. Multidimensional Data Modelling

Tabel dimensi lebih kecil dari tabel fakta. Dalam *data warehouse*, data kubus adalah kubus dengan n-dimensi. Fakta adalah hubungan antara dimensi. Tabel fakta berisi nama-nama fakta dan kunci dari tabel dimensi yang berhubungan dengan tabel fakta. Data fakta diekstrak dari berbagai sumber. Data fakta cenderung stabil dan tidak berubah dari waktu ke waktu. Tabel fakta yang besar, memiliki sejumlah baris sesuai dengan jumlah kemungkinan kombinasi nilai dimensi dan jumlah kolom sesuai dengan jumlah dimensi yang diwakili.

Data *Cube* berasal dari berbagai dimensi. Berbentuk kubus yang dapat dibuat dengan mengambil beberapa dimensi. Jenis skema dengan multidimensi model data [5].

- Skema Bintang (*Star Schema*): Skema bintang adalah skema *data warehouse* yang paling sederhana. Skema ini disebut skema bintang karena hubungan antara tabel dimensi dan tabel fakta menyerupai bintang, di mana satu tabel fakta terhubung ke beberapa tabel dimensi. Titik tengah dari skema bintang adalah tabel fakta besar dan terhubung dengan beberapa tabel dimensi.
- Skema Kepingan Salju (*Snowflake Schema*): Skema kepingan salju adalah variasi dari skema bintang di mana beberapa tabel dimensi dinormalisasi lagi, struktur yang dihasilkan menyerupai salju dengan tabel fakta di tengah, untuk menghasilkan beberapa tabel tambahan. Keuntungan yang didapat dengan menggunakan skema ini adalah menghemat memori, tapi waktu yang dibutuhkan untuk memproses *query* menjadi lebih lama.
- Skema Galaksi (*Constellation*): Dalam skema galaksi tabel fakta dapat berjumlah lebih dari satu dan saling terhubung dengan beberapa tabel dimensi. Keuntungan menggunakan skema ini adalah untuk menghemat memori dan mengurangi kesalahan yang mungkin terjadi.

E. SQL Server 2008 Integration Services

SQL Server Integration Services atau biasanya disingkat SSIS adalah alat yang digunakan untuk melakukan *Extract, Transform, dan Load* (ETL) yang diklasifikasikan sebagai fitur *Business Intelligence* (BI). ETL adalah proses pengumpulan data dari berbagai sumber (*Extract*), membersihkannya (*Transform*), dan kemudian menyimpannya di sistem lain (*Load*). Dalam kaitannya dengan BI, SSIS adalah fitur yang digunakan untuk menarik data dari ERP (*Enterprise Resource Planning*), *database* relasional, atau file, untuk kemudian hasilnya disimpan ke dalam *data warehouse* [9].

Dari istilah itu sendiri ETL menggambarkan tiga proses utama, yaitu proses ekstraksi data dari sumber data, transformasi data dan pemuatan data ke dalam *data warehouse*. Salah satu alat yang dapat digunakan untuk melaksanakan proses ini adalah Microsoft SQL Server 2005/2008 *Integration Services* (SSIS). SSIS desainer adalah alat grafis yang digunakan untuk membuat paket dan ada *Business Intelligence*

Development Studio sebagai bagian dari proyek *Integration Services*.

SQL Server 2008 adalah teknologi yang mendukung pengembangan dan administrasi BI Application. *SQL Server 2008 Analysis Services (SSAS)* adalah elemen dari BI. *Analysis Services* adalah teknologi untuk OLAP (*Online Analytical Processing*) dan *Data mining*. Administrasi proses OLAP dilakukan dalam bentuk *SQL Server Management Studio*.

III. METODOLOGI PENELITIAN

Menurut Kimball ada sembilan tahap dalam metodologi desain database untuk *data warehouse* [10], yaitu:

- **Memilih Proses:** Sebuah proses atau fungsi, mengacu pada masalah subjek tertentu dalam *data mart*. Ketika pertama kali membangun *data mart*, *data mart* harus tepat waktu, sesuai anggaran, dan bisa dipakai untuk menjawab pertanyaan penting dalam bisnis.
- **Menyatakan Grain:** Memilih *grain* memiliki makna memutuskan apa yang dijelaskan oleh *record* dalam tabel fakta.
- **Mengidentifikasi dan menyesuaikan dimensi:** Dimensi berisi hal-hal yang akan menjadi referensi dari tabel fakta.
- **Mengidentifikasi Fakta:** *Grain* pada tabel fakta untuk menentukan fakta-fakta dapat digunakan dalam *data mart*. Semua fakta harus menunjukkan tingkat *grain* yang sama.
- **Menyimpan pra-perhitungan dalam tabel fakta:** Ketika sebuah tabel fakta telah ditentukan, maka setiap fakta harus dikaji ulang untuk menentukan apakah dapat digunakan untuk melakukan *precalculations*.
- **Rounding Out the Dimension Tables:** Pada tahap ini ditambahkan informasi selengkap mungkin dalam tabel dimensi. Pernyataan harus intuitif dan mudah dipahami oleh pengguna (*user*).
- **Memilih Panjang Database:** Durasi mengukur sejauh mana tabel fakta bisa melihat beberapa tahun ke belakang. Banyak perusahaan ingin melihat apa yang terjadi dalam periode waktu yang sama dalam satu tahun atau dua tahun sebelumnya.
- **Tracking Slowly Changing Dimension:** Perubahan atribut Dimensi karena penambahan atribut alternatif baru sehingga pengguna dari *record* lama dan baru dapat digunakan secara bersamaan pada *record* dimensi yang sama.
- **Deciding The Query Priorities and The Query Models:** Pada tahap ini, akan membahas masalah yang berada pada perancangan fisik dari sebuah penyimpanan data yang mempengaruhi cara dan sudut pandang dari pengguna, kapasitas penyimpanan data pada *disk*, cara administrasi dari data, cara pembuatan cadangan data (*back up*) dan keamanan akses dari data tersebut.

IV. HASIL DAN PEMBAHASAN

Pemilihan proses didasarkan kepada 6 ruang lingkup yaitu proses sistem kendali mutu, kurikulum, penerimaan mahasiswa baru, perkuliahan, ujian dan bisnis. Proses yang dipilih adalah proses perkuliahan dan ujian dengan fokus pada hasil nilai mahasiswa.

Pemilihan *grain* berdasarkan proses perkuliahan dan ujian adalah analisis pada hasil nilai dalam bentuk angka dan analisis pada hasil nilai dalam bentuk huruf.

Setelah proses pemilihan *grain* untuk menentukan tabel fakta dan tabel dimensi, selanjutnya identifikasi dan penyesuaian isi dari tabel-tabel tersebut. Tabel dimensi yang diambil adalah tabel mahasiswa yang berisi atribut utama NPM dan nama mahasiswa, serta tabel mata kuliah dengan atribut utama kode mata kuliah dan nama mata kuliah.

Data mentah dalam bentuk Microsoft Acces diekstrak, diubah bentuknya ke dalam format SQL. Setelah data mengalami proses transformasi selanjutnya dilanjutkan dengan pembersihan data yang memakan waktu lama karena jumlah *record* yang sangat banyak. Pembersihan data meliputi penghilangan redundansi data, membersihkan data yang tidak konsisten, misalnya menyamakan nama atribut, dan kemudian membuat hubungan antara semua tabel. Data yang diambil adalah data mahasiswa untuk beberapa angkatan dan beberapa mata kuliah penting.

Berikut ini adalah tiga tabel utama dan desain yang telah bersih, yang telah dapat digunakan untuk merancang skema data warehouse setelah melewati proses *extract, transformation dan loading*. Terdapat dua tabel dimensi yang terhubung dengan satu tabel fakta, sebagai berikut:

Column Name	Data Type	Allow Nulls
KODE_MK	nvarchar(50)	<input type="checkbox"/>
KETERANGAN	nvarchar(50)	<input checked="" type="checkbox"/>

Fig. 1. Desain Tabel Dimensi Mata Kuliah

Column Name	Data Type	Allow Nulls
npm	nvarchar(50)	<input type="checkbox"/>
nama	nvarchar(50)	<input checked="" type="checkbox"/>

Fig. 2. Desain Tabel Dimensi Mahasiswa

Column Name	Data Type	Allow Nulls
npm	nvarchar(50)	<input type="checkbox"/>
KODE_MK	nvarchar(50)	<input type="checkbox"/>
NILAI	nvarchar(50)	<input type="checkbox"/>
NILAKHIR	real	<input type="checkbox"/>

Fig. 3. Desain Tabel Fakta KHS

Pra perhitungan disimpan dalam tabel fakta. Hasil nilai kuis, tugas, uts, praktikum dan uas dirangkum menjadi satu, sehingga menghasilkan keluaran dalam bentuk nilai angka dan nilai huruf.

Pada tahapan *Rounding Out the Dimension Tables* ditambahkan informasi selengkap mungkin dalam tabel dimensi.

KODE_MK	KETERANGAN
SI207	Dasar Akuntansi
SI310	Sistem Penunjang Keputusan
SI327	Akuntansi Menengah
SI340	SIM
SI341	Pemrog. Berorientasi Objek I
SI348	Teori Keputusan
SI451	E-Business
SP345	Basis Data Terapan
SP346	Sistem Basis Data
SP429	Manajemen Hubungan Pelanggan

Fig. 4. Isi Record Tabel Dimensi Mata Kuliah

Tabel dimensi mata kuliah memiliki 10 records dengan atribut kode mata kuliah dan nama mata kuliah.

npm	nama
2008240001	Margaretha
2008240002	Rizka
2008240003	Muhammad Edwin Utama Putra
2008240004	Revita Cassandra
2008240005	Kusuma Atmaja
2008240006	Melvilia Anggi
2008240007	Putri Marcelina
2008240008	M. Eko Baranata
2008240009	Hendry Santosa
2008240010	Ananda Wanajaya
2008240011	Shintya Chandra

Fig. 5. Isi Record Tabel Dimensi Mahasiswa

Tabel dimensi mahasiswa memiliki 366 records dengan atribut NPM dan nama mahasiswa.

npm	KODE_MK	NILAI	NILAI_HURUF
2008240003	SI207	C	10
2008240004	SI207	B	66
2008240004	SI340	A	89
2008240004	SI341	C	56,5
2008240005	SP429	C	59
2008240005	SP346	B	68,3
2008240005	SI341	D	71
2008240005	SI207	B	74
2008240005	SP345	C	56
2008240005	SI348	B	69,9
2008240005	SI340	A	80
2008240006	SP346	C	56,1
2008240006	SI207	D	60
2008240007	SI207	D	51
2008240008	SI207	B	65

Fig. 6. Isi Record Tabel Fakta KHS

Tabel fakta KHS memiliki 1548 records dengan atribut NPM, kode mata kuliah, nilai huruf dan nilai angka.

Panjang durasi *database* yang dipakai adalah kurang lebih tiga tahun untuk satu angkatan mahasiswa untuk satu program studi untuk beberapa mata kuliah wajib dan mata kuliah pilihan.

Atribut pada tabel dimensi ada yang tidak berubah, misalnya atribut untuk mahasiswa, namun untuk atribut pada tabel dimensi mata kuliah dapat memiliki perubahan tergantung kepada perubahan kurikulum dan juga silabus yang dibuat oleh program studi yang biasanya dilakukan secara berkala.

Pada tahap penentuan prioritas dan model *query* menentukan kapasitas penyimpanan yang nantinya akan mengakibatkan semakin bertambahnya ruang memori yang dibutuhkan sebagai media penyimpanan. Untuk penyimpanan *database* yang terkait nilai disediakan memori khusus dan dilakukan *back up* secara berkala.

Setelah proses 9 tahap desain *database* untuk *data warehouse* diatas, selanjutnya akan menghasilkan skema dalam bentuk *star schema* yang terdiri atas satu tabel fakta dan beberapa tabel dimensi yang sudah baik dan bersih dalam tampilannya dan dapat ditelusuri selanjutnya informasi apa saja yang dapat digali dengan implementasi *data mining* menggunakan *SQL server business intelligence development studio*. Sebuah project dalam bentuk *analysis services project* dibuat menjadi kubus, dengan terlebih dahulu membuat *solution explore*, menentukan *server name*, mengambil *database* untuk *data sources*, *data source views*, memilih tabel yang sesuai, membuat *cubes* dan *dimensions*.

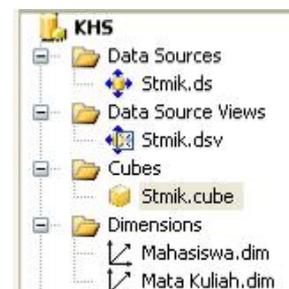


Fig. 7. Solution Explorer kubus

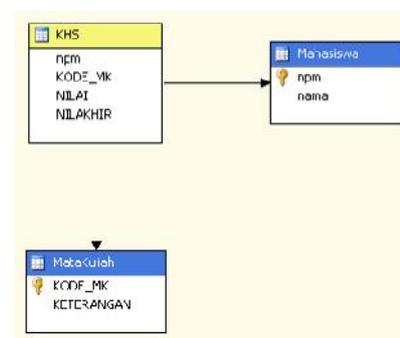


Fig. 8. Kubus Skema Bintang

Dari hasil kubus tersebut, selanjutnya ditarik *mining structures*. Fungsionalitas yang diambil adalah asosiasi dan pengelompokan. Berikut ini adalah *dependency network*, *rules* dan *itemsets* untuk fungsionalitas asosiasi.

tentang nilai-nilai mahasiswa yang tertinggi sampai terendah.

- Asosiasi *data mining* memberikan 77 *itemssets* dan juga 24 *rules*. *Itemssets* yang memiliki nilai *support* paling tinggi adalah untuk kode mata kuliah SI207 yaitu mata kuliah Dasar Akuntansi.
- Nilai A berasosiasi dengan kode mata kuliah SI 310 dan nilai akhir diatas 80,02052. Nilai B berasosiasi dengan kode mata kuliah SI 341 dan nilai akhir diantara 65,77625 dan 80,02052. Nilai C tidak berasosiasi mata kuliah tertentu dan nilai akhir diantara 53,19176 dan 65,77625. Nilai D berasosiasi dengan kode mata kuliah SI 348 dan nilai akhir diantara 29,2016 dan 53,19176. Nilai E berasosiasi dengan kode mata kuliah SP 346 dan nilai akhir dibawah 29,2016.
- Pengelompokan *data mining* sebanyak 7 *cluster*, dengan pengelompokan berdasarkan anggota kelompok yang memiliki kemiripan.
- *Cluster* dengan nilai A berjumlah 1, *cluster* dengan nilai B berjumlah 2, *cluster* dengan nilai C berjumlah 2, *cluster* dengan nilai D berjumlah 1, dan sisanya berjumlah 1 untuk nilai E.

DAFTAR PUSTAKA

- [1] Erick A., Lisangan N, dan Suswanto. S, "Perancangan Data Warehouse Pengolahan Persediaan Buku PT.Gramedia Asri Media Makassar", Munas Aptikom, Politeknik Telkom, Bandung, 2010.
- [2] Diajeng S.N., dan Indah P., "Perancangan *Data warehouse* pada PT. Olympindo Multi Finance Palembang di Area Regional Sumatera II, " STMIK MDP, 2012.
- [3] Iik W., "Pemanfaatan Data Warehouse di Perguruan Tinggi Indonesia," Jurnal Sistem Informasi MTI vol 4 no.1 ISBN 1412-896, Universitas Indonesia, 2008
- [4] Mukhlis F., dan Bayu A.T., "Rancang Bangun Data Warehouse untuk Menunjang Evaluasi Akademik di Fakultas," Konferensi Teknologi Informasi dan Aplikasinya (KNTIA), Universitas Sriwijaya, 2011.
- [5] Han J., Kamber M., dan Tung A.K.H., "Spatial Clustering Methods in Data Mining : A Survey, School of Computing Science Simon Fraser University Burnaby," Canada, 2001.
- [6] Sulianta F, and Juju D, "Data Mining Meramalkan Perusahaan ," Elex Media Komputindo, 2010.
- [7] Turban E., Aronson J. E., dan Liang T.P., "Decision Support Systems and Intelligent Systems," New Jersey Pearson Education, Inc, 2005.
- [8] Ponniah P., "Data Warehousing Fundamentals: A Comprehensive Guide for IT Professionals," John Wiley & Sons, USA, 2001.
- [9] Wirama K., Sudianto H., dan Hermawan Y., "The Essential Business Intelligence in Microsoft SQL Server 2008," Indonesia.Net Developer Community, 2009.
- [10] Connolly T.C., dan Begg C.F, "Database System : A Practical Approach to Design, Implementation, and Management," Pearson Education Inc, 2010.