

ห้องสมุดคณะเทคโนโลยีสารสนเทศ สจล.

การพัฒนาระบบคลังข้อมูลการประกันวินาศภัย
Data Warehouse for Insurance

โดย

นายชินพัฒน์ เคลือวัลย์

รหัส 45066068



H002169

อาจารย์ที่ปรึกษา

ผศ.ดร.วรพงษ์ กริสุระเดช

วัน เดือน ปี.....	0 6 ก.พ. 2550
เลขทะเบียน.....	02169
เลขเรียกหนังสือ.....	อภ. ๒๕๔๖/๒๕๔๖
"ห้องสมุดคณะเทคโนโลยีสารสนเทศ สจล."	

รายงานนี้เป็นส่วนหนึ่งของวิชาโครงการพัฒนาระบบงาน
หลักสูตรวิทยาศาสตรมหาบัณฑิต สาขาวิชาเทคโนโลยีสารสนเทศ
ภาคเรียนที่ 2 ปีการศึกษา 2546
คณะเทคโนโลยีสารสนเทศ

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

สารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านกา
ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไป

ชื่อหัวข้อ	ระบบคลังข้อมูลการประกันวินาศภัย
นักศึกษา	นายชินพัฒน์ เคลือวัลย์
อาจารย์ที่ปรึกษา	ผศ.ดร. วรพจน์ กรีสุระเดช
ระดับการศึกษา	วิทยาศาสตร์มหาบัณฑิต สาขาวิชาเทคโนโลยีสารสนเทศ
แขนงวิชา	วิทยาการสารสนเทศ
ปีการศึกษา	2546

บทคัดย่อ

การวิเคราะห์ข้อมูลให้อยู่ในรูปแบบสารสนเทศ เป็นสิ่งที่มีบทบาทมากในธุรกิจปัจจุบัน ซึ่งข้อมูลที่มีอยู่ในฐานข้อมูลนั้นส่วนใหญ่จะอยู่ในรูปของรายการที่เกิดจากกิจกรรมทางธุรกิจ และการวิเคราะห์ข้อมูลในมุมมองต่างๆ ไม่สามารถทำได้อย่างมีประสิทธิภาพ สำหรับโครงการนี้เป็น การพัฒนาระบบวิเคราะห์ข้อมูลการประกันวินาศภัย โดยจะทำการเปลี่ยนรูปแบบข้อมูลให้เป็น ข้อมูลสรุปซึ่งจัดเก็บในดาต้าแวร์เฮาส์ และจะทำการวิเคราะห์ข้อมูลเพื่อนำเสนอในรูปแบบรายงาน

Title	Data Warehouse for Insurance
Student	Mr. Chinnapat Kluewan
Advisor	Asst. Prof. Dr. Worapoj Kreesuradej
Level of Study	Master of Science in Information Technology
Major	Information Science
Academic Year	2003

ABSTRACT

Information analysis is most important for business today. While data in database is business transaction and information analysis for decision support is non-efficiency. This project is providing data warehouse for insurance data and information analysis for present into exclusive report.

สารบัญ

	หน้า
บทคัดย่อภาษาไทย.....	I
บทคัดย่อภาษาอังกฤษ.....	II
สารบัญ.....	III
สารบัญภาพ.....	V
สารบัญตาราง.....	VII
บทที่	
1. บทนำ.....	1
1.1 ความเป็นมา.....	1
1.2 วัตถุประสงค์.....	2
1.3 ขอบเขตการดำเนินงาน.....	2
1.4 ขั้นตอนและวิธีการดำเนินงาน.....	3
1.5 ประโยชน์ที่คาดว่าจะได้รับ.....	3
2. เทคโนโลยีดาต้าแวร์เฮาส์.....	4
2.1 เบื้องต้นของดาต้าแวร์เฮาส์.....	4
2.2 สถาปัตยกรรมและการประมวลผลแบบ End-to-End.....	6
2.3 เครื่องมือ Back End และยูทิลิตี้.....	7
2.3.1 การปรับแก้ข้อมูล (Data Cleaning).....	7
2.3.2 การโหลดข้อมูล.....	8
2.3.3 รีเฟรชแวร์เฮาส์.....	9
2.3.4 Conceptual Model และเครื่องมือ Front End.....	9
2.4 กระบวนการออกแบบฐานข้อมูล.....	11
2.5 เซิร์ฟเวอร์สำหรับแวร์เฮาส์.....	13
2.5.1 โครงสร้างอินเด็กซ์และการใช้งาน.....	13
2.5.2 Materialized View และการใช้งาน.....	15
2.5.3 การประมวลผลแบบขนาน.....	16
2.5.4 สถาปัตยกรรมของเซิร์ฟเวอร์สำหรับการประมวลผลคิวรี.....	16
2.5.5 คำสั่ง SQL เพิ่มเติม.....	17

สารบัญ (ต่อ)

2.5.6 เมตะดาต้าและการจัดการแวร์เฮาส์	18
3. โมเดลสำหรับฐานข้อมูลแบบมุมมองหลายมิติ	20
3.1 โมเดลข้อมูล (Data Model)	20
3.2 โอเปอเรเตอร์ (Operators)	21
3.2.1 Push	22
3.2.2 Pull	23
3.2.3 Destroy Dimension	23
3.2.4 Restriction	24
3.2.5 Join	25
3.2.6 Merge	28
3.3 Multidimensional Query Language	29
4. การพัฒนาระบบคลังข้อมูลการประกันวินาศภัย	33
4.1 วัตถุประสงค์ในการจัดทำดาต้าแวร์เฮาส์	33
4.2 การเตรียมข้อมูลและการทำดาต้าคลีนนิ่ง	33
4.2.1 การกำหนดโครงสร้างฐานข้อมูลสำหรับดาต้าแวร์เฮาส์	33
4.2.2 การทำดาต้าคลีนนิ่ง (Data Cleaning)	35
4.3 การออกแบบและทำการประมวลผลข้อมูล OLAP	37
4.3.1 การระบุแหล่งที่มาของข้อมูล (Data Source)	37
4.3.2 การสร้างไคเมนชัน (Dimensions)	38
4.3.3 การสร้างคิวบ์ (Cube)	43
4.3.4 การออกแบบที่จัดเก็บข้อมูล (Storage Design)	49
4.3.5 การประมวลผลข้อมูล (Process or Reprocess)	50
4.4 การพัฒนาระบบส่วนการสอบถามและนำเสนอสารสนเทศ	53
5. การทดสอบระบบและผลสรุป	57
5.1 การทดสอบระบบ	57
5.2 การสรุปผล	58
บรรณานุกรม	59
ประวัติผู้เขียน	60

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญภาพ

ภาพที่	หน้า
2-1 แสดงสถาปัตยกรรมทั่วไปของการทำดาต้าแวร์เฮาส์.....	6
2-2 แสดงข้อมูลแบบมุมมองหลายมิติ.....	10
2-3 แสดง Star Schema.....	12
2-4 แสดง Snowflake Schema.....	13
3-1 แสดง Logical cube ที่มียอดขายเป็นมิติ.....	21
3-2 แสดงโอเปอเรชัน push ในมิติของสินค้า.....	22
3-3 การดึง (Pull) ค่าสมาชิกตัวแรกของแต่ละค่าให้เป็นมิติของยอดขาย.....	23
3-4 แสดงผลของโอเปอเรชัน Restrict.....	25
3-5 แสดงการjoinค่าสองคู่.....	26
3-6 แสดงการทำ associate ของสองคู่.....	27
3-7 แสดงการเมิร์จมิติ Date และ Product โดยใช้ $f_{elem} = SUM$	28
3-8 แสดงภาพเบื้องต้นของการทำงาน MDX Query.....	29
3-9 แสดงภาพโครงสร้างลำดับชั้น ที่กำหนดเป็นรูปแบบของ Time-dimension.....	30
3-10 แสดงภาพโครงสร้างข้อมูลลำดับชั้นบน Time-dimension.....	30
3-11 แสดงภาพโครงสร้างของการจัดเก็บเมตาดาต้าในCUBE Schema.....	31
4-1 แสดงภาพ Data warehouse schema ของระบบคลังข้อมูลการประกันวินาศภัย.....	34
4-2 แสดงการสร้าง Data Source (โดยเลือก Provider Connection).....	37
4-3 แสดงการสร้าง Data Source (ในการระบุ Parameter สำหรับ Connection).....	38
4-4 แสดงการสร้างไดเมนชัน (ขั้นตอนที่ 1).....	39
4-5 แสดงการสร้างไดเมนชัน (ขั้นตอนที่ 2).....	39
4-6 แสดงการสร้างไดเมนชัน (ขั้นตอนที่ 3).....	40
4-7 แสดงการสร้างไดเมนชัน (ขั้นตอนที่ 4).....	40
4-8 แสดงการสร้างไดเมนชัน (ขั้นตอนที่ 5).....	41
4-9 แสดงการสร้างไดเมนชัน (ขั้นตอนที่ 6).....	41
4-10 แสดงการสร้างคู่ (ขั้นตอนที่ 1).....	44
4-11 แสดงการสร้างคู่ (ขั้นตอนที่ 2).....	44
4-12 แสดงการสร้างคู่ (ขั้นตอนที่ 3).....	45

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญภาพ (ต่อ)

4-13	แสดงการสร้างคูป (ขั้นตอนที่ 4)	45
4-14	แสดงการสร้างคูป (ขั้นตอนที่ 5)	46
4-15	แสดงการสร้าง Virtual Cube (ขั้นตอนที่ 1)	46
4-16	แสดงการสร้าง Virtual Cube (ขั้นตอนที่ 2)	47
4-17	แสดงการสร้าง Virtual Cube (ขั้นตอนที่ 3)	47
4-18	แสดงการสร้าง Virtual Cube (ขั้นตอนที่ 4)	48
4-19	แสดงการสร้าง Virtual Cube (ขั้นตอนที่ 5)	48
4-20	แสดงขั้นตอนการออกแบบพื้นที่จัดเก็บข้อมูล (ขั้นตอนที่ 1)	49
4-21	แสดงขั้นตอนการออกแบบพื้นที่จัดเก็บข้อมูล (ขั้นตอนที่ 2)	50
4-22	แสดงหน้าจอการประมวลผลข้อมูลของคูปทั้งหมด	51
4-23	แสดงหน้าจอการล็อกอินเข้าใช้งานระบบ	53
4-24	แสดงหน้าจอเมนูหลักรายงาน	54
4-25	แสดงรายงาน “สรุปผลการรับประกันภัย” ที่ออกแบบ	54
4-26	แสดงการบันทึกรายงานที่ออกแบบไว้	55
4-27	แสดงหน้าจอข้อมูลที่ใช้อธิบายที่มาของข้อมูล	55
4-28	แสดงหน้าจอคำศัพท์ธุรกิจสำหรับระบบประกันวินาศภัย	56

สารบัญตาราง

ตารางที่	หน้า
4-1 แสดงชื่อของมุมมองข้อมูลของกลุ่มข้อมูลจากคาด้าแวร์เฮาส์ และ โดเมนข้อมูลที่จะทำคลีนนิ่ง.....	36
4-2 แสดงรายการ Dimension และ Level ที่ทำสำหรับการประมวลผลข้อมูล.....	42
4-3 แสดงตัวอย่างของการออกแบบ “Claim_Profile” Cube Schema.....	52



บทที่ 1

บทนำ

1.1 ความเป็นมา

การประกันวินาศภัยรถยนต์ภาคสมัครใจ จัดอยู่ในหมวดหนึ่งของประเภทการประกันวินาศภัย ซึ่งเป็นธุรกิจที่มีการเติบโตค่อนข้างสูง ทำให้มีการแข่งขันในกลุ่มของบริษัทในธุรกิจนี้สูงด้วยเช่นกัน ดังนั้นจึงมีความต้องการที่จะทำวิธีที่มีประสิทธิภาพเข้ามาใช้ในการวิเคราะห์ข้อมูลต่างๆ ที่เป็นรายการที่เกิดขึ้นจากอดีต ถึงปัจจุบัน เพื่อช่วยในการตัดสินใจ และเป็นแนวทางในการดำเนินธุรกิจในอนาคต ดังนั้นกระบวนการดังกล่าว จึงเป็นเรื่องที่น่าสนใจในการศึกษา เพื่อเป็นการสนับสนุนสารสนเทศให้แก่กลุ่มบริษัทประกันภัย

การจัดเตรียมข้อมูลให้อยู่ในรูปแบบที่เหมาะสมสำหรับประมวลผล เป็นเรื่องที่มีความสำคัญสำหรับการทำให้การประมวลผลมีประสิทธิภาพ เนื่องจากข้อมูลที่อยู่ในระบบงานแหล่งต่างๆ นั้นเป็นข้อมูลที่มีรูปแบบที่ดีหรือและไม่ดีรูปแบบที่ไม่เหมาะสมในการประมวลผล ดังนั้นจึงมีความจำเป็นที่จะต้องทำการเลือกและจัดข้อมูลให้อยู่ในรูปแบบที่ง่ายและมีประสิทธิภาพในการประมวลผล สำหรับข้อมูลที่จะนำมาใช้ในการประมวลผลได้แก่ ข้อมูลการทำประกันวินาศภัย ข้อมูลการจ่ายค่าสินไหมทดแทน และข้อมูลเชิงสถิติต่างๆ ซึ่งข้อมูลเหล่านี้จะต้องเปลี่ยนรูปแบบรายการที่เป็นทรานแซกชัน (Transaction) ให้อยู่ในรูปแบบของข้อมูลอนุกรมเวลาในรูปแบบต่างๆ ให้เหมาะสมกับความต้องการในการวิเคราะห์ ดังเช่น การวิเคราะห์อัตราการเติบโตของธุรกิจ สถิติตามช่วงเวลา แนวโน้มสำหรับโอกาสทางธุรกิจ และการวิเคราะห์ความเสี่ยง เป็นต้น

คาต้าแวร์เฮาส์ (Data Warehouse) เป็นหมวดหมู่ของเทคโนโลยีที่สนับสนุนการตัดสินใจซึ่งมีวัตถุประสงค์ให้ผู้ทำงานสามารถใช้ในการบริหาร จัดการ และวิเคราะห์ได้ดียิ่งขึ้นและรวดเร็วกว่าเดิม ในปัจจุบันคาต้าแวร์เฮาส์เป็นจุดที่สร้างความสนใจในธุรกิจประเภทต่างๆ เนื่องจากข้อมูลที่มีจำนวนมาก เช่นรายการขายจากระบบประกันวินาศภัย มีจำนวนมากในระดับกิกะไบต์ ซึ่งในการประมวลผลของคาต้าเบสที่มีอยู่ในปัจจุบัน ส่วนใหญ่เป็น RDBMS นั้นได้รับความนิยมมาก แต่ความสามารถในการประมวลผลรายการเพื่อสนับสนุนความต้องการสารสนเทศแบบปรับเปลี่ยน (Ad Hoc Report or Query) นั้นไม่สามารถทำงานได้อย่างมีประสิทธิภาพ เนื่องจากข้อมูลนั้นมีข้อจำกัดในเรื่องของปริมาณข้อมูล ความถูกต้องของข้อมูล และความซับซ้อนของประเภทรายการเหล่านี้เป็นปัจจัยทำให้ความต้องการรายงานหรือการค้นหานั้นไม่มีประสิทธิภาพ จึงจำเป็นที่หาวิธีที่เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จะช่วยแก้ปัญหาดังกล่าว ซึ่งเป็นการลดรูปของข้อมูลและจัดรูปแบบให้เหมาะสมอย่างมีความสัมพันธ์กัน ให้อยู่ในดาต้าแวร์เฮาส์ ซึ่งดาต้าแวร์เฮาส์นี้เองจะมีการแยกออกจากระบบจัดการฐานข้อมูลทั่วไป เนื่องจากเหตุผลหลาย อย่างเช่น การสนับสนุนการทำงานประมวลผลเชิงวิเคราะห์แบบออนไลน์ หน้าที่การทำงาน และความต้องการความสามารถสูงในการประมวลผล

สำหรับทฤษฎีในการสร้างดาต้าแวร์เฮาส์นั้นมีหลักการหลายอย่างที่จะต้องคำนึงถึง ทั้งในเรื่องของประเภทข้อมูล ช่วงความเหมาะสม ต่างๆ เพื่อให้ได้ข้อมูลที่ถูกต้องนั้น ยังพบว่าปัญหาของการสร้างดาต้าแวร์เฮาส์ยังมีอีกหลายอย่างซึ่งในโครงการนี้จะกล่าวถึงปัญหาที่อาจจะเกิดและแนวทางแก้ไข เพื่อเป็นแนวทางในการพัฒนาระบบ และในทฤษฎีของการสร้างฐานข้อมูลสำหรับดาต้าแวร์เฮาส์จะใช้โมเดลของฐานข้อมูลแบบหลายมิติ (Modeling Multidimensional Database) ซึ่งจะสามารถนำมาพัฒนาระบบประกันวินาศภัยรถยนต์ภาคสมัครใจ โดยใช้ดาต้าแวร์เฮาส์ ได้

1.2 วัตถุประสงค์

- 1.2.1 เพื่อศึกษาหลักการในการจัดทำดาต้าแวร์เฮาส์ ปัญหาและการแก้ไขที่เหมาะสม
- 1.2.2 เพื่อศึกษาทฤษฎีของโมเดลของฐานข้อมูลแบบหลายมิติ ถึงหลักการในการสร้างฐานข้อมูลอย่างมีประสิทธิภาพ
- 1.2.3 เพื่อนำทฤษฎีมาทำการพัฒนาระบบประกันวินาศภัยรถยนต์ภาคสมัครใจ โดยใช้ดาต้าแวร์เฮาส์

1.3 ขอบเขตการดำเนินงาน

โครงการนี้เป็นการศึกษาเพื่อพัฒนาระบบประกันวินาศภัยรถยนต์ภาคสมัครใจโดยใช้ดาต้าแวร์เฮาส์ โดยมีขอบเขตของการศึกษาดังต่อไปนี้

- 1.3.1 เป็นการพัฒนาบบโดยใช้โมเดลของฐานข้อมูลแบบหลายมิติ ในการสร้างดาต้าแวร์เฮาส์
- 1.3.2 ข้อมูลที่ใช้ในการวิเคราะห์เป็นข้อมูลการทำประกันวินาศภัยรถยนต์ภาคสมัครใจที่ได้จากกรมการประกันภัย สำหรับข้อมูลระหว่างปีพ.ศ. 2544 ถึง พ.ศ. 2546
- 1.3.3 ระบบวิเคราะห์ข้อมูลเพื่อทำการสร้างดาต้าแวร์เฮาส์ พัฒนาด้วย VBScript ซึ่งจะมีขั้นตอนการเลือกและการตรวจสอบการเปลี่ยนแปลงของข้อมูล แล้วทำการปรับปรุงดาต้าแวร์เฮาส์

1.4 ขั้นตอนและวิธีการดำเนินงาน

ในการสร้างระบบประกันวินาศภัยรถยนต์ภาคสมัครใจโดยใช้ดาต้าแวร์เฮาส์ ของโครงการพัฒนาระบบงานที่จัดทำขึ้นนี้มีขั้นตอนการดำเนินงานดังต่อไปนี้

- 1.4.1 ศึกษาเทคโนโลยีของดาต้าแวร์เฮาส์ เพื่อเป็นแนวทางในการดำเนินงาน
 - 1.4.2 ศึกษาหลักการการสร้างดาต้าแวร์เฮาส์ ปัญหาและวิธีการแก้ไข
 - 1.4.3 ศึกษาโมเดลของฐานข้อมูลแบบหลายมิติ เพื่อดำเนินการสร้าง Schema ของดาต้าแวร์เฮาส์
 - 1.4.4 ดำเนินการรวบรวมข้อมูลการทำประกันวินาศภัยรถยนต์ภาคสมัครใจจากระบบฐานข้อมูลที่เกี่ยวข้อง
 - 1.4.5 ดำเนินการพัฒนาระบบงาน โดยแบ่งส่วนการทำงานเป็น 4 ขั้นตอนดังนี้คือ
 - 1.4.5.1 การวิเคราะห์รูปแบบความต้องการของข้อมูลประเภทต่างๆ เพื่อเป็นขอบเขตในการเลือกข้อมูล
 - 1.4.5.2 การล้างข้อมูลและการนำข้อมูลเข้าดาต้าแวร์เฮาส์
 - 1.4.5.3 การสร้างแบบสอบถามหรือรายงาน เพื่อเรียกใช้ข้อมูลจากดาต้าแวร์เฮาส์
 - 1.4.5.4 การแสดงผลหรือการนำเสนอข้อมูล
- #### 1.5 ประโยชน์ที่คาดว่าจะได้รับ
- 1.5.1 สามารถเรียนรู้และเข้าใจทฤษฎีของดาต้าแวร์เฮาส์ ในเรื่องของความแตกต่างจากระบบฐานข้อมูลทั่วไป
 - 1.5.2 สามารถสร้างดาต้าแวร์เฮาส์ได้ตามโมเดลของฐานข้อมูลแบบหลายมิติ และการสร้างข้อมูลในดาต้าแวร์เฮาส์
 - 1.5.3 สามารถพัฒนาระบบที่สนับสนุนระบบประกันวินาศภัยรถยนต์ภาคสมัครใจโดยใช้ดาต้าแวร์เฮาส์

บทที่ 2

เทคโนโลยีดาต้าแวร์เฮาส์

ดาต้าแวร์เฮาส์เป็นระบบความสำคัญอย่างมากที่ช่วยในเรื่องของการสนับสนุนการตัดสินใจ (Decision Support) ซึ่งในปัจจุบันอุตสาหกรรมดาต้าเบสได้พยายามมุ่งเน้นมากขึ้น ดังจะพบว่าหลายผลิตภัณฑ์ต่างมุ่งเน้นสำหรับสินค้าและบริการที่เป็นไปได้ที่เป็นเชิงการค้าให้มีการทำดาต้าแวร์เฮาส์เพื่อสนับสนุนเรื่องของการประมวลผลเพื่อช่วยในระบบการตัดสินใจ เช่นเทคโนโลยีการประมวลผลเชิงวิเคราะห์แบบออนไลน์ (On-Line Analytical Processing : OLAP) ซึ่งมีการเสนอเครื่องมือต่างๆ ในการช่วยสร้าง แต่ทั้งนี้ ในเรื่องของหลักการสร้างที่เหมาะสมยังเป็นศิลป์ที่ผู้ทำการสร้างจะต้องเข้าใจและออกแบบให้เหมาะสม

2.1 เบื้องต้นของดาต้าแวร์เฮาส์

ดาต้าแวร์เฮาส์ เป็นเทคโนโลยีที่ช่วยสนับสนุนการตัดสินใจ โดยมีเป้าหมายเพื่อให้คนทำงาน (ผู้บริหาร, ผู้จัดการ, นักวิเคราะห์) สามารถทราบว่าจะทำอย่างไรให้ดี และรวดเร็ว จากอดีตที่ผ่านมาจะพบว่าสินค้าและบริการมีการเพิ่มจำนวนยอดขายอย่างรวดเร็วจึงทำให้ดาต้าแวร์เฮาส์มีการพัฒนาอย่างมากเช่นกัน และเหตุผลหลักที่สำคัญคือ ดาต้าแวร์เฮาส์สนับสนุนการประมวลผลเชิงวิเคราะห์แบบออนไลน์ ฟังก์ชัน และประสิทธิภาพสูง ซึ่งแตกต่างกับระบบการประมวลผลรายการ (On-Line Transaction Processing : OLTP) ที่จะประมวลผลกับระบบฐานข้อมูลสำหรับงานทั่วไป

ในความแตกต่างของดาต้าแวร์เฮาส์ คือ เป้าหมายของระบบช่วยสนับสนุนการตัดสินใจ ด้านประวัติ ข้อมูลสรุป และข้อมูลในการตัดสินใจจะมีความสำคัญมากกว่ารายละเอียด ซึ่งจะมีความสามารถมากในเรื่องของช่วงเวลา ซึ่งทำงานได้ดีกว่าระบบฐานข้อมูลทั่วไป โดยเฉพาะโครงการที่มีขนาดข้อมูลเป็นร้อยกิกะไบต์ การทำงานจะมีประสิทธิภาพมากถ้าเป็นคิวรีแบบ Ad Hoc เพราะการที่ทำงานของคิวรีที่มีความซับซ้อน จะสามารถเข้าถึงข้อมูลในจำนวนล้านเรคคอร์ดด้วยฟังก์ชันการทำงานของ Join และ Aggregates แล้วให้ผลลัพธ์ของคิวรีตามเวลาที่ไม่นาน ซึ่งเป็นเรื่องที่มีความสำคัญมากกว่า ผลของ รายการทรานแซคชัน

หลักการทั่วไปของโมเดลในดาต้าแวร์เฮาส์ ก็คือ Multidimensional หรือ มุมมองหลายมิติ จะช่วยให้สามารถทำการวิเคราะห์ข้อมูลที่ซับซ้อน และเห็นภาพได้ง่าย ตัวอย่างเช่น ยอดขาย

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่นอนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

กรรมธรรม์ในดาต้าแวร์เฮาส์ จะมีรายละเอียดเกี่ยวกับประเภทกรรมธรรม์ ประเภทรถยนต์ ระยะเวลา ความคุ้มครอง อัตราเบี้ยประกันภัย และค่าสินไหมทดแทน การที่จะสรุปข้อมูลนั้นขึ้นอยู่กับมุมมองที่เราสนใจว่าต้องการเลือกข้อมูลอะไร ช่วงเวลาแบบใด เช่น วัน เดือน ไตรมาส ปี ในส่วนประเภทกรรมธรรม์ เช่น ประกันภัยประเภท 1, 2, 3 เป็นต้น

ในการทำงานของฐานข้อมูลที่มีการออกแบบอย่างดีเพื่อให้สามารถทำงานแบบ OLTP (On-Line Transaction Processing) จะมีการประมวลผลที่ซับซ้อนกว่า OLAP (On-Line Analytical Processing) และประสิทธิภาพในการแสดงผลที่นั้นยอมรับไม่ได้สำหรับใช้ในระบบสนับสนุนการตัดสินใจ ซึ่งอาจจะเกิดจากการทำงานผิดพลาดของฐานข้อมูลเอง หรือการทำคิวรีซ้ำของข้อมูลเก่าที่ไม่ถูกต้อง เพราะวาระบบฐานข้อมูลจะทำการบันทึกเพียงข้อมูลปัจจุบันเท่านั้น ระบบสนับสนุนการตัดสินใจโดยปกติจะใช้ข้อมูลจากแหล่งที่มีความแตกต่างกัน ซึ่งอาจจะรวมถึงฐานข้อมูลภายนอกในระบบอื่นๆ ซึ่งแหล่งข้อมูลที่แตกต่างกันอาจจะมีคุณภาพของข้อมูลที่แตกต่างกัน หรือไม่สามารถนำเสนอได้ รหัส หรือรูปแบบ ซึ่งจะมีความขัดแย้งกันได้ สรุปได้ว่าการสนับสนุนของโมเดลข้อมูลแบบหลายมิติ และการทำงาน OLAP มีความต้องการข้อมูลที่มีการจัดโครงสร้างโดยเฉพาะ วิธีการเข้าถึงข้อมูล และวิธีการนำมาใช้งาน ซึ่งไม่ได้เตรียมไว้ในระบบฐานข้อมูลในเชิงธุรกิจทั่วไป จึงเป็นเหตุผลที่ต้องแยกดาต้าแวร์เฮาส์ออกจากระบบการทำงานของฐานข้อมูล

ดาต้าแวร์เฮาส์อาจจะมีนำมาใช้บนมาตรฐาน หรือเพิ่มเติมใน RDBMSs ที่เรียกว่า รีเลชันนอล OLAP (ROLAP) ซึ่งเครื่องที่ทำการติดตั้งจะทำการบันทึกข้อมูลเสมือนบันทึกในฐานข้อมูลแบบรีเลชันนอล และจะสนับสนุนคำสั่งเพิ่มเติมของ SQL และการเข้าถึงข้อมูลและวิธีการที่มีประสิทธิภาพที่ใช้กับโมเดลข้อมูลแบบหลายมิติ ซึ่งจะแตกต่างจากการใช้ มัลติไดเมนชัน OLAP ที่เครื่องจะทำการบันทึกข้อมูลแบบมุมมองหลายมิติ และเป็น โครงสร้างพิเศษ (เช่น อาร์เรย์) และการใช้งาน OLAP จะต้องอยู่ภายใต้โครงสร้างข้อมูลเฉพาะนี้

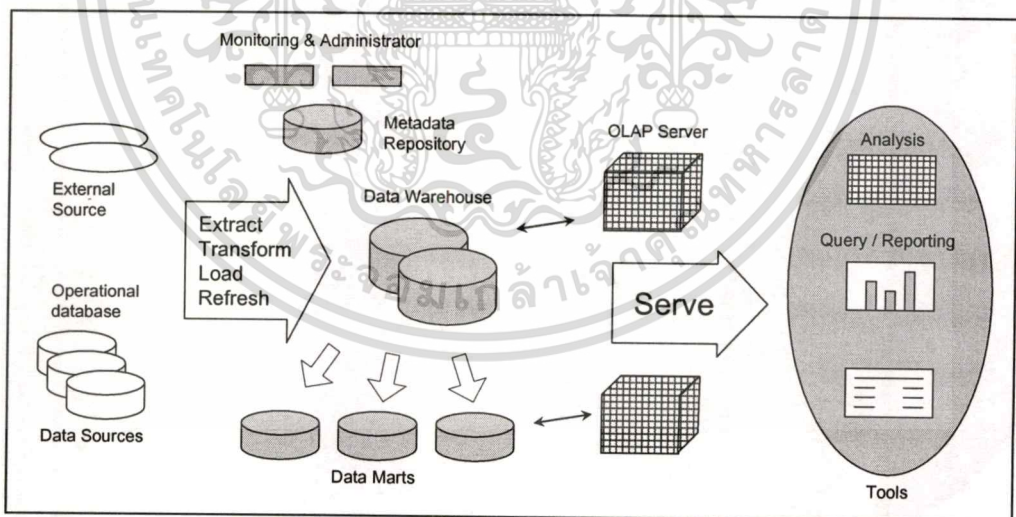
มีวิธีการสร้างและปรับปรุงดาต้าแวร์เฮาส์ที่มากกว่าติดตั้ง OLAP การออกแบบ schema และการทำคิวรีที่ซับซ้อน เนื่องจากความแตกต่างของปัจจุบันสถาปัตยกรรมที่มีอยู่สามารถเลือกได้ และหลายองค์กรต้องการที่จะรวมระบบให้เป็นเอนเตอร์ไพรส์แวร์เฮาส์ (Enterprise Warehouse) ที่จะรวมทั้งหมดเข้าด้วยกัน เช่น ลูกค้า สินค้า ทرفฟ์สิน เป็นต้น ซึ่งอยู่ในระบบต่างๆ อย่างไรก็ตาม การสร้างเอนเตอร์ไพรส์แวร์เฮาส์จะต้องใช้เวลานานและมีความซับซ้อนในการประมวลผล และใช้ค่าใช้จ่ายสูง ซึ่งขึ้นกับนโยบายทางธุรกิจของระดับบริหาร และใช้เวลาหลายปีกว่าจะประสบความสำเร็จ บางองค์กรอาจจะแยกข้อมูลออกเป็นดาต้ามาร์ท (Data Mart) ซึ่งจะทำแยกแผนกโดยเลือกเฉพาะส่วนที่สนใจในการสร้าง ดาต้ามาร์ทเหล่านี้สามารถทำงานได้เร็วและไม่ต้องการความ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สอดคล้องกันในระหว่างคาต้ามาร์ท ขององค์กรที่มีขนาดใหญ่ แต่อาจจะเกิดปัญหาในการรวมเข้าด้วยกัน และใช้เวลานานในการทำงาน ถ้าธุรกิจที่ต้องการความสมบูรณ์ของระบบ โมเดลนี้ไม่เหมาะสม

2.2 สถาปัตยกรรมและการประมวลผลแบบ End-to-End

จากรูป 2-1 เป็นการรวมเครื่องมือสำหรับดึงข้อมูลจากฐานข้อมูลหลายๆ แหล่ง และแหล่งข้อมูลภายนอก สำหรับการปรับแก้ข้อมูลให้อยู่ถูกต้อง (Clearing) การเปลี่ยนรูปข้อมูล และการรวมข้อมูลเข้าไว้ในคาต้าแวร์เฮาส์ โดยจะมีการรีเฟรชเป็นระยะซึ่งเป็นผลที่เกิดจากการปรับปรุงของแหล่งข้อมูล และจากการลบข้อมูลออกจากคาต้าแวร์เฮาส์ ในส่วนที่เพิ่มเติมเข้าคาต้าแวร์เฮาส์ หลักอาจจะมีหลายคาต้ามาร์ท ข้อมูลที่อยู่ในแวร์เฮาส์และคาต้ามาร์ทจะถูกบันทึกและจัดการโดยหนึ่งหรือหลายเซิร์ฟเวอร์แวร์เฮาส์ ซึ่งจะอยู่ในมุมมองหลายมิติของข้อมูลที่มีความแตกต่างของเครื่องมือที่สร้าง (front end tools) ซึ่งได้แก่ query tools, report writers, analysis tools และ data mining tools ซึ่งระบบเหล่านี้จะมีที่จัดเก็บและจัดการเมตะคาต้า และเครื่องมือสำหรับจัดการแวร์เฮาส์ในการบริหารและติดตาม



รูปที่ 2-1 แสดงสถาปัตยกรรมทั่วไปของการทำคาต้าแวร์เฮาส์

แวร์เฮาส์อาจจะเป็นรูปแบบกระจายเพื่อให้เกิดความสมดุลในการทำงาน (Load balancing) มากที่สุดเท่าที่จะทำได้ ดังตัวอย่างเช่นสถาปัตยกรรมแบบกระจาย การสร้างฐานข้อมูลเพื่อเก็บเมตะคาต้าของแต่ละส่วนที่กระจายอยู่ (Fragment) ของแวร์เฮาส์ เพื่อที่จะสามารถนำมารวมกันได้ ซึ่งการเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เลือกสถาปัตยกรรมแบบนี้อาจจะต้องมีค่าใช้จ่ายที่สูงขึ้นเช่นกัน สำหรับมุมมองที่กลับกัน การรวมหรือการแยกแต่ละแวร์เฮาส์เป็นดาต้ามาร์ท ก็ต้องมีแหล่งเก็บแยกจากกัน และต้องมีระบบจัดการแบบกระจายด้วย

การออกแบบและการปรับให้ดาต้าแวร์เฮาส์ลดความซับซ้อนของการประมวลผลจะต้องคำนึงกิจกรรมดังนี้

- กำหนดรูปแบบสถาปัตยกรรม วางแผนความจุของข้อมูล และเลือกเซิร์ฟเวอร์ที่จะใช้บันทึกเป็นฐานข้อมูลกับ OLAP เซิร์ฟเวอร์ และเลือก tools ที่จะใช้งาน
- การปรับให้เซิร์ฟเวอร์ แหล่งเก็บข้อมูลและเครื่องใช้งานให้เข้ากัน
- ออกแบบ schema และวิวของแวร์เฮาส์
- กำหนดโครงสร้างของแวร์เฮาส์ที่ต้องการ การปรับแก้ข้อมูล การแบ่งกลุ่มข้อมูล และวิธีการเข้าถึงข้อมูล
- เชื่อมต่อแหล่งข้อมูล โดยใช้เกตเวย์ ผ่าน ODBC หรือใช้ตัวกลางอื่นๆ (Wrappers)
- ออกแบบและประยุกต์สคริปต์สำหรับการแยกข้อมูล การปรับแก้ข้อมูล การเปลี่ยนรูปข้อมูล การโหลดและการรีเฟรชแวร์เฮาส์
- ทำการข้อมูลจาก schema และวิวที่กำหนดไว้, สคริปต์
- ออกแบบและพัฒนา end-user แอปพลิเคชัน
- ตรวจสอบแวร์เฮาส์และแอปพลิเคชัน

2.3 เครื่องมือ Back End และยูทิลิตี้

ระบบการดาต้าแวร์เฮาส์จะใช้ความหลากหลายของการแยกข้อมูล และเครื่องมือในการปรับแก้ข้อมูล และโปรแกรมช่วยการโหลดและรีเฟรช สำหรับทำการย้ายแวร์เฮาส์ การแยกข้อมูลจากแหล่งข้อมูลที่ต่างกัน โดยส่วนมากจะทำโดยผ่านเกตเวย์และอินเตอร์เฟสมาตรฐาน เช่น Information Builders EDA/SQL, ODBC, Oracle Open Connect, Sybase Enterprise Connect, Informix Enterprise Gateway เป็นต้น

2.3.1 การปรับแก้ข้อมูล (Data Cleaning)

เนื่องจากดาต้าแวร์เฮาส์เป็นการใช้งานสำหรับการสร้างการตัดสินใจ ดังนั้นจึงมีความสำคัญที่ว่าข้อมูลในแวร์เฮาส์จะต้องถูกต้อง แต่เนื่องจากขนาดของข้อมูลจากแหล่งต่างๆ ที่เกี่ยวข้องมีโอกาสที่จะเกิดความผิดพลาดและผิดพลาดในตัวของข้อมูล ซึ่งเครื่องมือที่ช่วยในการค้นหาข้อมูลที่ผิดพลาดไม่สามารถทำการตรวจหาได้ จึงจำเป็นจะต้องล้างเพิ่มเริ่มทำงานใหม่

ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตัวอย่างเช่น ขนาดของฟิลด์ที่ไม่คงที่ คำอธิบายที่ไม่ถูกต้อง ค่าของสัญลักษณ์ไม่ถูกต้อง มีการป้อนค่าผิด และมีการฝ่าฝืนความถูกต้องของข้อมูล (Integrity constraints) สำหรับสิ่งที่เกิดโดยไม่คาดคิดมาก่อน เช่นการเลือกฟิลด์จากแหล่งข้อมูลที่ไม่เหมาะสม ข้อมูลมี 3 ความสัมพันธ์ แต่มีอะไรที่แตกต่างกัน

การแบ่งประเภทของเครื่องมือที่ทำการปรับแก้ข้อมูล ได้แก่ data migration tools เป็นการย้ายข้อมูลตามกฎที่กำหนด เช่นการแทนที่คำว่า “gender” ด้วย “sex” ซึ่งระบบจัดการแวร์เฮาส์ของ Prism เป็นตัวอย่างของกรณีนี้, data scrubbing tools ใช้ความรู้เฉพาะขอบเขตที่กำหนด เช่นรหัสไปรษณีย์ ซึ่งจะใช้หลักการในการแยกและใช้เทคนิคพีชชีแมทซึ่ง เครื่องมือบางชนิดอาจจะทำการหาความสัมพันธ์ของแหล่งข้อมูล สำหรับ data auditing tools อาจจะค้นหากฎของความสัมพันธ์ของข้อมูล โดยการสแกนข้อมูล ซึ่งเครื่องมือดังกล่าวอาจจะพิจารณาจากความแตกต่างของ data mining tools

2.3.2 การโหลดข้อมูล

หลังจากทำการแยกข้อมูล การแก้ไขและการเปลี่ยนรูปของข้อมูลแล้ว จะต้องทำการโหลดข้อมูลเข้าดาต้าแวร์เฮาส์ สิ่งที่จะต้องทำก่อนการประมวลผลก็คือ การตรวจสอบกฎความถูกต้อง (Integrity constraints) การเรียงลำดับ การคำนวณผลรวม การสรุป และการคำนวณอื่นที่จะต้องบันทึกในตารางที่สร้างไว้ในแวร์เฮาส์ การสร้าง index สำหรับการเข้าถึงข้อมูล การแบ่งพื้นที่ในการจัดเก็บ โดยทั่วไปจะใช้วิธีการโหลดแบบ batch สำหรับยูทิลิตี้ที่ช่วยในการโหลดข้อมูลเข้าแวร์เฮาส์ จะต้องสามารถให้ผู้ใช้ดูแลระบบสามารถทำการดูสถานะ เพื่อที่จะยกเลิก สั่งให้หยุดหรือทำงานต่อ และการเริ่มต้นทำงานใหม่ถ้าเกิดการผิดพลาดเนื่องจากข้อมูลสูญเสียความถูกต้อง

ยูทิลิตี้ที่ช่วยในการโหลดข้อมูลสำหรับดาต้าแวร์เฮาส์ ต้องใช้เวลาในการโหลดในเวลาจำกัด (ปกติจะทำช่วงกลางคืน) เมื่อแวร์เฮาส์ offline เพื่อทำการรีเฟรชข้อมูล การโหลดจะใช้เวลานานมาก ถ้าข้อมูลเป็นเทราไบต์ อาจจะใช้เวลาเป็นสัปดาห์หรือเป็นเดือน ดังนั้นการทำ pipeline และการแบ่งงานทำแบบคู่ขนาน ก็จะช่วยลดระยะเวลาการทำงานได้มาก โดยการทำเป็นลักษณะ batch transaction ขณะที่ทำการประมวลผลรายการปัจจุบันเสร็จก็จะทำการ commits และระบบจะเริ่มประมวลผลใหม่ถ้าเกิดความผิดพลาดที่จุด checkpoint ล่าสุด

อย่างไรก็ดี การทำงานแบบคู่ขนานก็ยังคงใช้เวลานาน โดยส่วนมากยูทิลิตี้ที่ใช้ในงานด้านธุรกิจจะใช้วิธีการโหลดข้อมูลเมื่อมีการเพิ่มหรือลดขนาดจำนวนข้อมูลพร้อมกัน ซึ่งเป็นเพียงการ

update ค่าเมื่อมีการ insert ถึงแม้วิธีการนี้ก็จะเป็นการเพิ่มการทำงานของคิวรี แต่การทำงานของระบบฐานข้อมูลในการเข้าถึงข้อมูลนั้นมีกระบวนการที่ดีและการใช้อินเด็กซ์

2.3.3 รีเฟรชแวร์เฮาส์

การรีเฟรชแวร์เฮาส์ หรือการปรับปรุงข้อมูลในแวร์เฮาส์ให้ทันสมัยจากแหล่งข้อมูลที่ทำางานร่วมกันของข้อมูลหลักและบันทึกข้อมูลเหล่านั้นในแวร์เฮาส์ มีสองสิ่งที่ต้องพิจารณา คือ เมื่อไรที่จะทำการรีเฟรช และจะทำการรีเฟรชอย่างไร โดยปกติ เวิร์เฮาส์จะมีการรีเฟรชเป็นระยะเวลา เช่น ทุกวัน หรือทุกสัปดาห์ ในบางระบบอาจจะมีควมจำเป็นที่ต้องมีการปรับปรุงเสมอเพื่อให้เป็นข้อมูลปัจจุบัน จึงจำเป็นที่จะต้องทำการปรับปรุงตามทุกครั้งที่มีการเปลี่ยนแปลงนโยบายการรีเฟรชจะถูกกำหนดโดยผู้บริหารระบบ โดยพิจารณาจากความต้องการของผู้ใช้งานและความคับคั่งของเครือข่าย และอาจจะแตกต่างกันระหว่างแหล่งข้อมูลที่ต่างกัน

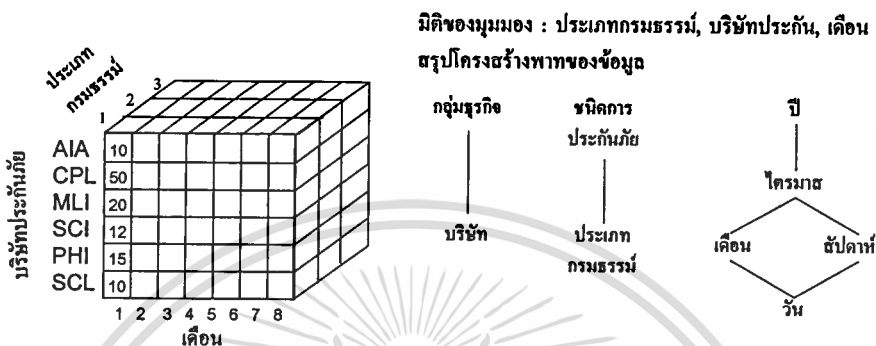
เทคนิคของการรีเฟรชอาจจะขึ้นอยู่กับลักษณะของแหล่งข้อมูลและความจุของฐานข้อมูลประเภทของแหล่งข้อมูลอาจจะเป็นระบบไฟล์ หรือ เป็นฐานข้อมูลที่มีราคาแพงก็ได้ ซึ่งในการต้องทำงานกับระบบงานเก่า เราจะต้องทราบว่าระบบงานนั้นสามารถสนับสนุนกับการทำเรพลิเคท (Replicate) หรือ ไม่เมื่อแหล่งข้อมูลมีการเปลี่ยนแปลง ตัวอย่างของเทคนิคเรพลิเคท เช่น การทำ log ไฟล์ ซึ่งอาจจะเป็น table ในฐานข้อมูลแล้วสร้าง trigger ทำเมื่อมีการเปลี่ยนแปลงรายการข้อมูลก็จะไปเพิ่มบันทึก log และเมื่อระบบเรพลิเคททราบว่าเกิดการเปลี่ยนแปลง โดยดูจาก log ก็จะทำการปรับปรุงข้อมูลทันที ซึ่งวิธีจะเป็นการเพิ่มการทำงานให้กับฐานข้อมูลมากขึ้น เนื่องจากไม่สามารถหาวิธีอื่นที่จะทำการตรวจสอบการเปลี่ยนแปลงข้ามระบบฐานข้อมูลที่เป็น DBMSs ที่ต่าง Vender เพราะยังไม่มาตรฐาน APIs สำหรับการเข้าไปดึงข้อมูลใน log ของทรานเซกชันได้ ดังนั้นการทำเรพลิเคทเซิร์ฟเวอร์จะถูกใช้ในการรีเฟรชดาต้าแวร์เฮาส์ แต่อย่างไรก็ดีในแต่ละรอบที่ทำการรีเฟรชจะขึ้นกับจำนวนข้อมูลที่ต้องทำของยูทิลิตี้ที่ทำการ โหลดข้อมูลด้วย

2.3.4 Conceptual Model และเครื่องมือ Front End

โมเดลในการย้ายข้อมูลนั้นจะหมายถึง เครื่องมือ front end การออกแบบฐานข้อมูล และคิวรีเอนจินสำหรับ OLAP ที่เป็น วิวข้อมูลในแวร์เฮาส์แบบหลายมิติ ในดาต้าโมเดลแบบหลายมิติ จะมีการตัวเลขในการวัดค่าของการวิเคราะห์ ตัวอย่างเช่น การวัดจำนวนยอดขาย รายได้ ค่าใช้จ่าย และผลตอบแทนการลงทุน ตัวเลขแต่ละค่าที่ใช้วัดจะขึ้นอยู่กับชุดของไดเมนชัน (dimension) ซึ่งเตรียมข้อมูลสำหรับการวัด เช่น การวัดค่าความเกี่ยวข้องของยอดขายในเมือง ชื่อสินค้า และวันที่เมื่อมีการขาย สมมติว่ามีไดเมนชันเดียวที่จะทำการวัด ซึ่งค่าในมิติเป็นที่ว่าง และแต่ละมิติ จะถูกอธิบายด้วยเซตของแอตทริบิวต์ ตัวอย่างเช่น ประเภทของการประกันภัย แบบของการทำประกัน, บริษัท

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่นอนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ประกันภัย, เดือนที่ทำการเริ่มคุ้มครอง โดยสมมติว่า บริษัท A ได้ทำการขายกรมธรรม์ประเภท 1 ในเดือนเมษายน ได้จำนวน 5,000,000 บาท สามารถทำการวัดได้ว่าเป็นสัดส่วนของข้อมูลในทั้งกลุ่ม โดยจะวัดได้จากความสัมพันธ์ระหว่างข้อมูลประเภทเดียวกันจากมิติที่มอง



รูปที่ 2-2 แสดงข้อมูลแบบมุมมองหลายมิติ

สำหรับความสามารถอื่นๆ ของคอนเซ็ปต์โมเดลตามตัวอย่างนี้ เป็นการสรุปโดย aggregate ในการวัดหนึ่งหรือหลายมิติ ที่ใช้เป็นตัวหลักในการทำงาน เช่น การคำนวณและแบ่งผลรวมยอดขาย ตามบริษัท ในแต่ละเดือน สำหรับการคำนวณโดยใช้สองตัววัด อย่างเช่น จำนวนกรมธรรม์ และเบี้ยประกัน ซึ่งจะมีการวัดค่าโดยใช้ลักษณะมิติได้เหมือนกัน

มิติแบบแบ่งช่วงเวลา (Time Dimension) เป็นรูปแบบที่สนับสนุนการตัดสินใจ เช่นการวิเคราะห์แนวโน้ม ซึ่งจะถูกใช้บ่อยมาก โดยจะมีการเพิ่มความสามารถลงไปในปฏิทินและฟังก์ชันต่างๆ ที่เกี่ยวกับมุมมองมิติแบบแบ่งช่วงเวลา

เครื่องมือ Front End

เครื่องมือ front end เป็นส่วนที่ใช้ในการนำเสนอข้อมูลจากโมเดลข้อมูลหลายมิติ โดยส่วนมากนิยมใช้โปรแกรม spreadsheet ในการใช้สำหรับการวิเคราะห์ทางธุรกิจ เนื่องจากกระดาษทำการ (Spreadsheet) จะสามารถมีพื้นที่ว่างในการใช้ประโยชน์เพื่อคำนวณอื่นๆ เพิ่มเติมได้ และเป็นความท้าทายที่จะทำให้ OLAP สามารถสนับสนุนการทำงานบนกระดาษทำการได้อย่างมีประสิทธิภาพบนฐานข้อมูลที่มีขนาดใหญ่หลายกิกะไบต์ได้ ซึ่ง Microsoft Excel ก็เป็น front end ถูกพัฒนาเพื่อรองรับการทำงานสำหรับมุมมองหลายมิตินี้

อย่างที่ทราบว่าคุณสมบัติที่จะนำกระดาษทำการมาใช้งานเนื่องจากเป็นโปรแกรมที่สนับสนุนมุมมองหลายมิติ อีกงานหนึ่งที่เราเห็นได้ก็คือการทำ pivoting หรือการสรุปสาระสำคัญโดยการเปลี่ยนระดับการมอง เมื่อพิจารณาจากรูปที่ 2 ซึ่งสามารถแทนลงในกระดาษทำการ เมื่อแต่ละ

แถวอธิบายถึงบริษัทที่เกี่ยวข้องกับยอดขาย โดยมีหนึ่งคอลัมภ์สำหรับแต่ละมุมมอง และคอลัมภ์พิเศษที่แทนค่ายอดขาย สำหรับตัวอย่างมุมมองแบบสรุปรายละเอียด (Pivot) จะมีการเลือกสองมิติที่ต้องการวัดค่าตำแหน่ง (x, y) ที่จะทำการสรุปค่าที่จะวัดจากมิติแรกโดยมีค่า x และมิติที่สองเป็นค่า y จากตัวอย่าง ถ้าเลือกมิติเป็นประเภทกรรมธรรม์ และ เดือน ดังนั้น แกน x อาจจะแทนค่าทั้งหมดของประเภทกรรมธรรม์ และแกน y อาจจะแทนค่าทั้งหมดของเดือน ดังนั้นจุด (x, y) จะแทนค่าสรุป (aggregate) ยอดขาย สำหรับสำหรับกรรมธรรม์ประเภท x ในเดือน y และค่าอะไรก็ตามในแถวหรือคอลัมภ์ที่อยู่ในตารางข้อมูลเดิม ก็จะไปปรากฏเป็นตัวเลือกในตารางสรุปรายละเอียด

การทำงานอื่นๆ ที่เกี่ยวข้องกับการสรุปรายละเอียดก็คือการ rollup หรือ drill-down, rollup เป็นการทำงานในลักษณะการทำข้อมูลปัจจุบันและทำการจัดกลุ่ม (group by) ด้วยมิติหนึ่งของทั้งหมด เช่นการ rollup ข้อมูลยอดขาย แล้วใช้ฟังก์ชันสรุปค่าของประเภทกรรมธรรม์ และชื่อบริษัท ส่วน drill-down เป็นการทำงานย้อนกลับของ rollup

Slice_and_dice เป็นการทำงานที่ลดข้อมติของข้อมูล เช่นการทำ projection ของข้อมูลให้เป็นซับเซตของมิติทั้งหมด ตัวอย่างเช่น การทำ slice_and_dice ข้อมูลยอดขาย สำหรับการเลือกเฉพาะบริษัทประกันภัย แล้วสร้างเป็นตารางที่เป็นมิติสำหรับประเภทกรรมธรรม์ และเดือนของยอดขาย นอกจากนี้ยังมีวิธีการอื่นๆ อีก เช่นการจัดลำดับ (ranking) และการคำนวณฟิลด์ข้อมูล

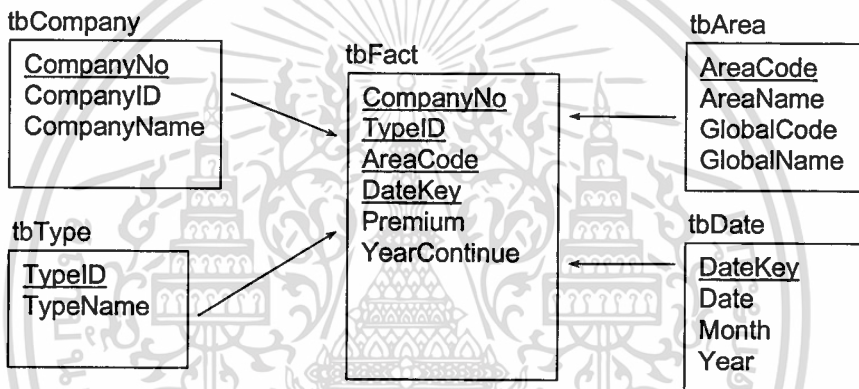
ทว่าสิ่งที่ให้มากกว่ามุมมองหลายมิติของกระดาษทำการนั้นมีความน่าสนใจอย่างยิ่ง เนื่องจากการเพิ่มความสามารถอย่างมากให้ผู้ใช้งานสามารถทำการวิเคราะห์ข้อมูลธุรกิจ แต่สิ่งนี้ไม่ได้เป็นการเข้ามาแทนที่ระบบที่ทำการวิเคราะห์ที่อยู่เดิม เช่นการค้นหา หรือการทำการวิเคราะห์ที่ซับซ้อนสำหรับใช้งานจากระบบฐานข้อมูลเดิม ดังนั้น เครื่องมือที่สามารถให้ผู้ใช้งานระบบสามารถทำการวิเคราะห์เพื่อดูเฉพาะข้อมูลที่ต้องการ ก็จะต้องเป็นโปรแกรมที่เข้าถึงข้อมูลดิบได้อย่างมีประสิทธิภาพตามการควบคุมของระบบฐานข้อมูลหลัก แต่ใช้กระบวนการของเครื่องมือ เช่น Microsoft Access จะสามารถสร้างคิวรีแบบ ad hoc SQL ที่เรียกว่า “pointing-and-clicking” ในความหลากหลายของเครื่องมือสำหรับดาต้าไมนิ่งเหล่านี้จะมีการใช้ เป็นเครื่องมือ front end สำหรับดาต้าแวร์เฮาส์

2.4 กระบวนการออกแบบฐานข้อมูล

โมเดลของข้อมูลหลายมิติมุมมอง จากรายละเอียดก่อนหน้าจะสามารถนำมาใช้งานโดยตรงกับเซิร์ฟเวอร์แบบ MOLAP ได้ แต่อย่างไรก็ดี ถ้ามีการใช้แบบ ROLAP เราจะต้องหาความสัมพันธ์ระหว่างโมเดลหลายมิติมุมมองกับการปฏิบัติงาน ซึ่งการออกแบบฐานข้อมูลแบบเชิงสัมพันธ์นี้จะสะท้อนภาพมุมมองหลายมิติของข้อมูลได้

ปกติ ER diagram (Entity Relationship) และเทคนิคการทำ Normalization เป็นวิธีการที่นิยมใช้ในการออกแบบฐานข้อมูลใน OLAP ต่างอย่างไรก็ดี การออกแบบฐานข้อมูลโดยใช้ ER diagram เป็นวิธีที่ไม่เหมาะสมสำหรับระบบสนับสนุนการตัดสินใจ เพราะมีส่วนที่ต้องคำนึงถึงประสิทธิภาพในการคิวรีและการโหลดข้อมูลซึ่งเป็นสิ่งสำคัญมาก

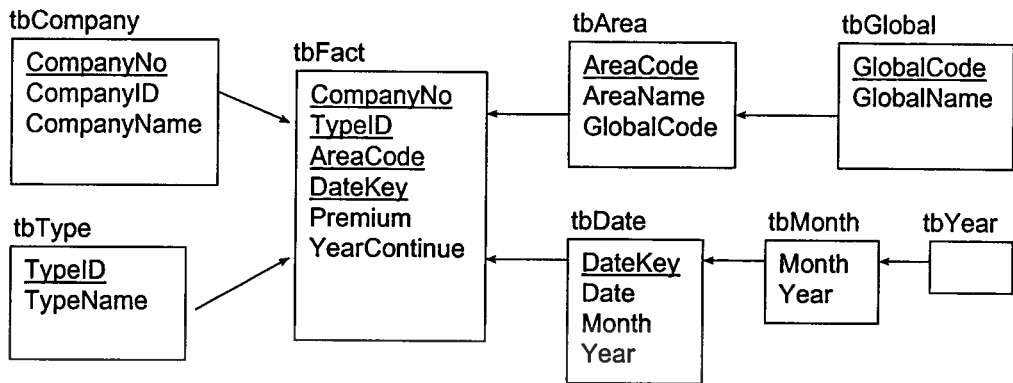
โดยส่วนมาก คาด้าแวร์เหล่านี้จะใช้ star schema ในการแทนความหมายของโมเดลข้อมูลหลายมิติ ซึ่งฐานข้อมูลที่เป็นข้อมูลจริงจะเป็น ตารางเดี่ยว สำหรับแต่ละมิติ แต่ละแถวในตารางจะมีจุดเชื่อม (foreign key) ไปยังแต่ละมิติที่ทำงานร่วมกันและสิ่งที่บันทึกเป็นตัวเลขที่ใช้ในการทำงานร่วมกัน ดังแสดงตัวอย่าง star schema ตามรูปที่ 3



รูปที่ 2-3 แสดง Star Schema

Star Schema ไม่มีความชัดเจนในการเตรียมให้สนับสนุนสำหรับเอทริบิวต์แบบหลายชั้น แต่ Snowflake schema เป็นการปรับปรุงให้มีโครงสร้างลำดับชั้นของมิติมีความชัดเจน โดยการทำ นอมอลไลซ์ ตารางมิติดังแสดงตามรูปที่ 2-4 ซึ่งวิธีการนี้เป็นประโยชน์ในการปรับปรุงตารางมิติได้ อย่างไรก็ตามถ้าโครงสร้างของตารางมิติไม่มีการทำ นอมอลไลซ์ โครงสร้างก็อาจจะเหมาะสมมากกว่า สำหรับการใช้แค่ค้นหาข้อมูล

Fact constellations เป็นตัวอย่างของ โครงสร้างตารางมิติที่แยกตารางมิติค่าจริง (fact table) ออกเป็นหลายตาราง สำหรับตัวอย่างอาจจะมีการแยกตารางออก และตารางที่แยกออกก็อาจจะเป็น ตารางจริงได้



รูปที่ 2-4 แสดง Snowflake Schema

2.5 เซฟเวอร์สำหรับแวร์เฮาส์

ดาต้าแวร์เฮาส์ อาจจะบรรจุข้อมูลเป็นจำนวนมาก เพื่อตอบคิวิรีได้อย่างมีประสิทธิภาพ ซึ่งเมื่อเราต้องการความประสิทธิภาพสูงในวิธีการเข้าถึงข้อมูลและเทคนิคการประมวลผลคิวิรี ดังนั้นจึงมีหลายเรื่องที่ต้องดำเนินการ สิ่งแรกก็คือ ดาต้าแวร์เฮาส์ต้องใช้โครงสร้างที่ซ้ำซ้อนกัน เช่นการใช้อินเด็กซ์และการใช้วิว การเลือกอินเด็กซ์จะใช้สำหรับสร้าง และวิวจะใช้ในการแก้ไขปัญหาการออกแบบ ในแนวความคิดถัดไปเพื่อให้เกิดประสิทธิภาพในการใช้อินเด็กซ์และวิวที่มีอยู่เพื่อใช้ตอบคิวิรี การลดความซับซ้อนของคิวิรีเป็นปัญหาที่มีความสำคัญต่อมา ดังนั้น ขณะที่จะคิวิรีข้อมูลการใช้อินเด็กซ์ในการสแกนหา อาจจะมีประสิทธิภาพมาก ส่วนข้อมูลที่คงที่และไม่มีมาก คิวิรีอาจจะใช้วิธีสแกนแบบลำดับ จึงเป็นข้อสรุปเพื่อประสิทธิภาพในการสแกนหาข้อมูล ความจำเป็นในการทำงานแบบขนานเพื่อเป็นการลดเวลาในการประมวลผล คงจะเป็นไปไม่ได้ เพราะจำเป็นอย่างยิ่งที่จะต้องทำการวิเคราะห์ในแต่ละส่วน

2.5.1 โครงสร้างอินเด็กซ์และการใช้งาน

จำนวนของเทคนิคในการประมวลผลคิวิรีนั้นควรจะใช้ประโยชน์จากอินเด็กซ์ สำหรับตัวอย่าง เช่น การเลือกจากหลายเงื่อนไข สามารถใช้อินเด็กซ์ทำการอินเตอร์เซค (Intersection) นอกจากนี้ยังสามารถใช้อินเด็กซ์ทำการยูเนียน (Union) ซึ่งอินเด็กซ์เหล่านี้เป็นคำที่ใช้ในการเข้าถึงข้อมูลหลัก

เซฟเวอร์สำหรับแวร์เฮาส์ สามารถใช้อินเด็กซ์แบบ bit map ซึ่งจะช่วยให้การทำยูเนียนและอินเตอร์เซค มีประสิทธิภาพเพิ่มขึ้น

เซิร์ฟเวอร์สำหรับดาต้าแวร์เฮาส์สามารถใช้บิตแมปอินเด็กซ์ ซึ่งสนับสนุนการทำงานประเภทยูเนียน หรืออินเตอร์เซ็คชันของอินเด็กซ์ได้อย่างมีประสิทธิภาพ การพิจารณาแต่ละหน้าในโครงสร้างอินเด็กซ์ที่เป็นขอบเขตของค่า d ซึ่งหน้านั้นจะบรรจุรายการของรหัสเรคคอร์ด (RIDs) ที่เป็นเรคคอร์ดที่บรรจุค่า d แต่ว่าบิตแมปอินเด็กซ์จะถูกใช้เป็นทางเลือกในการนำเสนอค่าที่สัมพันธ์กับรายการ RID ที่เป็นบิตเวกเตอร์ ซึ่งแต่ละบิตของเรคคอร์ด ซึ่งเป็นการกำหนดไว้ในขอบเขตของค่าที่เรคคอร์ดเป็น d ในแนวคิดนี้ บิตแมปอินเด็กซ์จะไม่ใช้โครงสร้างอินเด็กซ์ใหม่ แต่เป็นทางเลือกในการนำมาใช้งานของรายการ RID การที่นิยมใช้บิตแมปอินเด็กซ์เพื่อให้ได้ค่าจริงในการแทนค่าของรายการ RID สามารถเพิ่มความเร็วของการทำงาน Intersection, Union, Join และ Aggregation สำหรับตัวอย่าง ถ้าเรามีคิวรีของเงื่อนไข $column1 = d$ และ $column2 = d'$ ดังนั้นเราจะเจาะจงรายการโดยใช้ AND ของสองบิตเวกเตอร์ ขณะที่เราแทนค่านั้นจะมีประสิทธิภาพสำหรับลดปัญหาการขัดแย้งของขอบเขต เช่น “gender” ซึ่งจะมีประสิทธิภาพสูงสำหรับขอบเขตของข้อมูลที่มีการบีบอัดไว้ (เพราะเวลารันระบบจะต้องถอดรหัสเสมอ) บิตแมปอินเด็กซ์นั้นแรกๆ จะถูกใช้ในโมเดล 204 แต่ผลิตภัณฑ์เกี่ยวกับฐานข้อมูลจำนวนมากพึ่งสนับสนุนการทำงานเร็วๆ นี้ เช่น Sybase IQ คำถามที่น่าสนใจคือ จะเลือกแอตทริบิวต์ใดมาทำอินเด็กซ์ เพราะโดยทั่วไปการจะเลือกแอตทริบิวต์ใดต้องขึ้นอยู่กับการออกแบบฐานข้อมูล

ในการเพิ่มอินเด็กซ์ในตาราง ที่ทำสำหรับ star schema เพื่อทำการ JOIN อินเด็กซ์ โดยเฉพาะสำหรับระบบสนับสนุนการตัดสินใจ ขณะที่การแมปอินเด็กซ์ค่าในแถวนั้น จะต้องทำการ JOIN อินเด็กซ์ระหว่างฟอร์เรนจ์คีย์กับไพรมารีคีย์ ซึ่งใน star schema จะมีการ JOIN อินเด็กซ์ที่สามารถกำหนดความสัมพันธ์ของค่าหนึ่งหรือหลายแอตทริบิวต์ของตารางมิติที่ทำการแมทช์แถวกับตารางค่าจริง ยกตัวอย่างจากรูปที่ 3 การทำการ JOIN อินเด็กซ์ “Area” สำหรับแต่ละ area ที่อยู่ใน RIDs ของแถวในตารางค่าจริงนั้นจะสอดคล้องกับขอบข่ายของ area ดังนั้นการ JOIN อินเด็กซ์จะมีการคำนวณล่วงหน้าแบบ binary join สำหรับการ JOIN อินเด็กซ์ที่เป็นมัลติคีย์ (Multitkey) สามารถทำการคำนวณล่วงหน้าโดยการ JOIN แบบ n-way ตัวอย่างเช่น ขอบข่ายทั้งหมดในฐานข้อมูลสามารถทำการ JOIN อินเด็กซ์จาก (ประเภทกรรมธรรม์, รหัสบริษัท) ในตารางค่าจริง เช่นการกำหนดค่า ('1', 'AIA') จุดที่ RIDs ของแถวทั้งหมดในตารางขอบข่าย จะถูกรวมกันก่อน การใช้วิธีการดังกล่าว การ JOIN อินเด็กซ์ของมุมมองหลายมิตินี้สามารถทำการเตรียมไว้ในบางเวลาก่อนได้ เพราะจะเป็นการประหยัดเวลาในการทำ intersection ในการแยกอินเด็กซ์ของประเภทกรรมธรรม์ และรหัสบริษัทประกกันกันได้ ดังนั้นการ JOIN อินเด็กซ์สามารถใช้สำหรับรายการ RID เพื่อเพิ่มประสิทธิภาพในการประมวลผลของการ JOIN ได้

ฐานข้อมูลของระบบสนับสนุนการตัดสินใจควรมีสิ่งที่สามารถสรุปรายละเอียดเกี่ยวกับตัวอักษรและอินเด็กซ์ตัวเองจะสนับสนุนในการค้นหาข้อความได้อย่างดี

2.5.2 Materialized View และการใช้งาน

คิวรีจำนวนมากบนดาต้าแวร์เฮาส์ต้องการข้อมูลแบบสรุป และ การใช้ aggregate ซึ่งการทำอินเด็กซ์เพิ่มก็จะช่วยลดเวลาในการสรุปข้อมูลของคิวรีทั่วไปได้ ตัวอย่าง เช่น เรื่องของ สิ้นไหมค่าทดแทน จะมีการทำคิวรีที่ต้องใช้ทรัพยากรมากในการทำงานเพื่อจะให้ได้คำตอบสำหรับรายไตรมาสหรือรายเดือน ดังนั้นการทำสรุปข้อมูลไว้จะเป็นการช่วยเพิ่มความเร็วในการประมวลผลได้

ในความคิดเกี่ยวกับการใช้งาน materialized view นี้จะไม่เหมือนกับการใช้งานอินเด็กซ์ คือ (1) การระบุไปว่าวิวเป็น materialize (2) เป็นการใช้งาน materialized view เพื่อเป็นคำตอบของคิวรี (3) เกิด ประสิทธิภาพในระหว่างการโหลดและรีเฟรชก็จะทำการปรับปรุง materialized view สำหรับปัจจุบันการนำแนวคิดนี้มาใช้งานยังมีปัญหา เมื่อพิจารณาความสัมพันธ์ของโครงสร้าง ดังเช่น view ของการจอยน์ในตารางค่าจริง ซึ่งเป็นซับเซตของตารางมิติ (อาจจะเป็นไปได้หลังจากที่มีการเลือกมิติเหล่านี้) กับฟังก์ชันในการวัดค่าโดยการ group by ค่าของแอตทริบิวท์ จากตารางมิติ เนื่องจากโครงสร้างของวิวเหล่านี้จะมีขนาดที่เล็กที่ซับซ้อนมาก (ดูจากโครงสร้างแบบ snowflake)

แม้ว่าจะมีการรูปแบบวิวที่เป็น materialize จะเป็นทางเลือกในการใช้งาน แต่โดยลักษณะของการทำงานจะเป็นการเพิ่มการทำงานด้าน update และต้องใช้พื้นที่ในการจัดเก็บมากขึ้น จากข้อสมมติฐานอย่างง่ายด้วย greedy algorithm แสดงให้เห็นว่าประสิทธิภาพดี จากปัญหาที่เกี่ยวข้องภายใต้การปรับปรุงประสิทธิภาพก็อาจจะเลือก materialized view เพื่อมาใช้งานสำหรับการหา aggregation ตอนสร้างความสัมพันธ์

ถึงแม้การใช้งานจะง่าย แต่ประโยชน์ในการใช้ materialized view ก็มีมากในการ select ข้อมูล หรือการ rollup โดยการ group และการเพิ่มฟังก์ชัน ตัวอย่างเช่น สมมติว่า materialized view ประกอบด้วยผลรวมยอดขายสำหรับแต่ละไตรมาส ในแต่ละรหัสบริษัทประกันภัย ซึ่ง materialized view นี้สามารถใช้ตอบคิวรีที่ถามเกี่ยวกับผลรวมยอดขายของ “AIA” สำหรับปีโดย ถ้าจะแสดงเป็นปีจะต้อง rolling up จากยอดรวมไตรมาสเป็นยอดรวมปี

โดยทั่วไป อาจจะมีหลายทางวิธีในการใช้งาน materialized view เพื่อเป็นคำตอบคิวรี ตัวอย่างเช่น วิว V เป็นเซตของมุมมองมิติของ Q ถ้ามีการเลือกบางส่วนใน Q แล้วปรากฏว่าส่วนที่เลือกนั้นอยู่ใน V และถ้าทำการ group-by คอลัมภ์ใน V เพื่อเป็นซับเซตของ group-by คอลัมภ์ใน Q ก็จะเห็นได้ว่าวิว V สามารถทำงานแทน Q ที่เกิดใหม่ได้ โดยกำหนดให้เป็นเซต materialized view เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

M และคิวรี Q เราสามารถกำหนดได้ว่า M' เป็นเซตที่น้อยที่สุดสำหรับ Q ตัวอย่างการใช้งานเช่น การถามผลรวมยอดขายกรมธรรม์ประเภท 3 โดยมีวิวที่ใช้งาน 2 วิวในการสร้าง คือ (1) ผลรวมยอดขายแต่ละจังหวัดสำหรับแต่ละประเภทกรมธรรม์ และ (2) ผลรวมที่ขายแต่ละภาคสำหรับแต่ละบริษัทประกันภัย ด้วยวิธีการแบบนี้จะเป็นการทำให้การค้นหาแคบลงและ materialized view เหมาะสมในการใช้งาน ขณะที่เราทำให้เป็นคำตอบที่น้อยที่สุดของคิวรี แต่ปัญหาโดยทั่วไปคือการทำ multi materialized view ก็จะมีควมซับซ้อนมากขึ้น ในกรณีที่ทำ Select-Project-Join

2.5.3 การประมวลผลแบบขนาน

การประมวลผลแบบขนานสำหรับฐานข้อมูลที่มีขนาดใหญ่กับข้อมูล Teradata เป็นเรื่องหนึ่งของเทคโนโลยี โดยผู้จำหน่ายระบบจัดการฐานข้อมูลทั้งหมด ขณะนี้ได้เสนอเทคโนโลยีในการจัดการข้อมูลแบบแบ่งส่วน (partitioning) และการประมวลผลแบบขนาน (parallel) โดยเทคนิคที่เกี่ยวข้องตัวหนึ่งที่น่าสนใจคือ การใช้ทรัพยากรของระบบสนับสนุนการตัดสินใจแบบอ่านอย่างเดียว (read-only) ที่เป็นลักษณะของสแกนแบบ piggybacking โดยจะทำงานลักษณะนี้เมื่อถูกสั่งงานแบบมัลติเพิลคิวรี (ใช้ใน Redbrick)

การสแกนแบบ piggybacking จะลดเวลาในการทำงานทั้งหมดได้เป็นอย่างดี โดยจะทำงานสำหรับคำสั่งที่มาพร้อมๆ กันแบบ overlapping

2.5.4 สถาปัตยกรรมของเซิร์ฟเวอร์สำหรับการประมวลผลคิวรี

เซิร์ฟเวอร์แบบรีเรชั่นนอลแบบเดิมไม่ค่อยมีความฉลาดในการใช้งานอินเด็กซ์และความต้องการอื่นๆ ที่สนับสนุนวิวและข้อมูลแบบหลายมิติ โดยผู้จำหน่าย RDBMS ทั้งหมดต่างก็พยายามให้ระบบสามารถสนับสนุนได้อย่างรวดเร็ว โดยการเพิ่มตามความต้องการ ดังนั้นสำหรับผู้พัฒนาระบบต้องพิจารณาสามเรื่องที่ต้องมีในระบบสนับสนุนการตัดสินใจ

- Specialized SQL Servers : Redbrick เป็นตัวอย่างของเซิร์ฟเวอร์ที่มีวัตถุประสงค์เตรียมไว้สำหรับภาษาคิวรีขั้นสูงและสนับสนุนการประมวลผลคิวรีบน star และ snowflake schema ภายใต้การทำงานแบบ read-only
- ROLAP Servers : มีการจัดเก็บที่อยู่ของข้อมูลระหว่าง relational back end server (จะบันทึกที่อยู่ของข้อมูลที่อยู่ในแวร์เฮาส์) กับ client front end tools สำหรับ Microstrategy เป็นตัวอย่างของเซิร์ฟเวอร์แบบนี้ ซึ่งจะไปเพิ่มความสามารถของระบบจัดการฐานข้อมูลรุ่นเก่า โดยระบบจะทำงานเป็นตัวกลาง (Middleware) ที่จะสนับสนุนให้คิวรีแบบ OLAP ทำงานได้อย่างมีประสิทธิภาพ และการทำงานเหมาะสมกับ back

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

end relational servers โดยจะทำการสร้างวิวให้เป็น materialized และทำการเปลี่ยนรูปแบบของการใช้คิวรีให้เหมาะสมกับ materialize view แล้วสร้าง multi-statement SQL สำหรับ back end server ดังนั้นการเตรียมเซิร์ฟเวอร์ ที่เพิ่มนั้นจะต้องทำการจัดการตารางของคิวรีและการกำหนดให้ใช้ทรัพยากรให้เหมาะสม ซึ่งแนวโน้มจะมีการเปลี่ยนมาใช้ ROLAP มาก และด้วยจุดเด่นของเซิร์ฟเวอร์ ROLAP นี้จะสามารถขยายระบบและคุณสมบัติการทำทรานแซกชันของระบบรีเลชันนอลได้ อย่างไรก็ดี การทำงานระหว่างคิวรีในรูปแบบ OLAP และ SQL นั้นก็ยังไม่ดี เช่น ลำดับในการประมวลผลยังเป็นแบบลำดับ (sequential) หรือการใช้ aggregate ยังคงทำกับคอลัมน์จึงเป็นสาเหตุทำให้เกิด bottleneck ของเซิร์ฟเวอร์ OLAP

- MOLAP Servers : เซิร์ฟเวอร์ประเภทนี้จะสนับสนุนข้อมูลและวิวแบบมุมมองหลายมิติโดยตรง เนื่องจากมีเอ็นจินสำหรับเก็บข้อมูลแบบหลายมิติ ความเห็นไปได้ในการนำมาใช้คิวรีด้วย front-end ในมุมมองหลายมิติบนเลเยอร์ของที่เก็บข้อมูลจะสามารถทำการแมปได้เลย ตัวอย่างของเซิร์ฟเวอร์นี้คือ Essbase ดังเช่นการใช้ประโยชน์จากคุณสมบัติของอินเด็กซ์จะทำได้ดีมาก แต่การเตรียมพื้นที่จัดเก็บข้อมูลจะค่อนข้างไม่มีประสิทธิภาพ โดยเฉพาะอย่างยิ่งเมื่อชุดข้อมูลอยู่กระจาย สำหรับเซิร์ฟเวอร์ MOLAP หลายผลิตภัณฑ์สามารถนำมาใช้งานเป็น 2-Level โดยจะทำการปรับการกระจายของข้อมูลให้อยู่รวมกันอย่างต่อเนื่อง ซึ่งในพื้นที่เก็บข้อมูลของ 2-Level นี้จะถูกทำให้เป็นซับบาเรย์ หนึ่งหรือสองมิติ ซึ่งจะทำให้ข้อมูลอยู่รวมกันต่อเนื่อง แต่ว่าการใช้งานของเครื่องมือออกแบบหรือการที่ผู้ใช้งานกำหนดเอง ให้อยู่ในรูปแบบอาเรย์ ซึ่งจะทำอินเด็กซ์ให้เอง มีหลายเทคนิคที่เป็นกลไกของระบบฐานข้อมูลทั่วไปได้นำมาใช้งานในเซิร์ฟเวอร์ MOLAP

2.5.5 คำสั่ง SQL เพิ่มเติม

มีหลายคำสั่งของ SQL ที่เพิ่มเข้ามาซึ่งง่ายต่อการใช้งานและมีวัตถุประสงค์เพื่อประมวลผลคิวรีของ OLAP หรือการนำมาใช้เพิ่มเติมในเซิร์ฟเวอร์แบบรีเลชันนอล อธิบายบางส่วนของคำสั่งที่เพิ่มดังนี้

- กลุ่มฟังก์ชัน aggregate กลุ่มนี้จะรวมความสามารถในการทำการเรียงลำดับ (rank) และการแจกแจงแบบ percentile เช่น all Product in Top 10 Percentile หรือ Top 10 Products by total Sale และยังสนับสนุนฟังก์ชันที่รู้จักกันดีและใช้งานเกี่ยวกับการเงิน เช่น mean, mode, median

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- คุณสมบัติการทำรายงาน การสร้างรายงานสำหรับการวิเคราะห์ข้อมูลเชิงธุรกิจ ความต้องการคุณสมบัติ aggregate ที่นำมาใช้เป็นลักษณะ time window เช่น แสดงการเคลื่อนที่ของค่าเฉลี่ย ซึ่งคุณสมบัติที่เพิ่มนี้มีความสำคัญที่จะกำหนดจุดหยุด (breakpoint) และทำการหาผลรวมที่เกิดขึ้น โดยคุณสมบัติ SQL นี้มีใน Redbrick
- Multiple Group-By เครื่องมือที่เป็น front end เช่น กระจายทำการมุมมองหลายมิติ ต้องการที่จะทำการ group by ค่าของแอตทริบิวต์ที่แตกต่างกัน ด้วยเงื่อนไขจะสามารถทำการทดลองกำหนดให้คำสั่งของ SQL ทำการสแกนว่ามีข้อมูลใดบ้างที่สามารถอยู่ในรูปแบบของ multiple times แต่วิธีการนี้ไม่มีประสิทธิภาพ ในปัจจุบันมี 2 คำสั่งใหม่ที่จะทำ ก็คือ Rollup และ Cube ซึ่งมีวัตถุประสงค์เพื่อทำการเชื่อม SQL การทำ Rollup ของรายการแอตทริบิวต์ (รหัสบริษัท, เดือน, ประเภทกรรมกรรม) แล้วข้อมูลที่เป็นผลลัพธ์ก็จะแสดงคำตอบตามรูปแบบ Group by คือ (a) group by (รหัสบริษัท, เดือน, ประเภทกรรมกรรม), (b) group by (รหัสบริษัท, เดือน), (c) group by (รหัสบริษัท) จากรายการจำนวนทั้งหมดจะแสดงผลจำนวนเป็น k คอลัมน์ ส่วนการทำ Cube เป็นการเตรียม group-by สำหรับแต่ตัวที่เป็นตัวประกอบ (Combination) การทำงาน 2K ของคอลัมน์ ดังนั้นคำสั่ง group-by แบบหลายจำนวนสามารถประมวลผลได้มีประสิทธิภาพมากกว่า โดยการตรวจสอบการทำงานระหว่างค่าของข้อมูลด้วยกันเอง ซึ่ง Microsoft SQL Server สนับสนุนการทำงานของทั้ง Cube และ Rollup
- การเปรียบเทียบ เนื่องจากความไม่สมบูรณ์ของ SQL จึงทำให้เกิดการเปรียบเทียบ ซึ่งเป็นเรื่องปกติในทางธุรกิจ ตัวอย่างเช่นการเปรียบเทียบระหว่างผลรวมยอดขายทั้งหมดและผลรวมยอดขายแต่ละไตรมาส ซึ่งปัจจุบันยังต้องทำการวิจัยเพื่อให้สนับสนุนการเปรียบเทียบ

2.5.6 เมตาดาต้าและการจัดการแวร์เฮาส์

เนื่องจากดาต้าแวร์เฮาส์เป็นการพัฒนาให้เป็นโมเดลทางธุรกิจของกิจการ โดยปัจจัยที่มีความสำคัญสำหรับสถาปัตยกรรมการทำแวร์เฮาส์ก็คือ การจัดการเมตาดาต้า มีหลายวิธีที่แตกต่างกันในการจัดการเมตาดาต้า

- Administrative Metadata ประกอบด้วย สารสนเทศที่จำเป็นสำหรับการตั้งค่าและการใช้แวร์เฮาส์ ซึ่งประกอบด้วย รายละเอียดของที่มาของฐานข้อมูล เครื่องมือที่เป็น front-end, back-end การให้คำจำกัดความของ warehouse schema มิติและลำดับชั้นของ

ข้อมูล การกำหนดคิวรีและรายงาน ที่จัดเก็บค่าตัวมาร์ท นโยบายการรีเฟรชและการล้างค่าตัวแวร์เฮาส์ และสิทธิการใช้งานระบบต่างๆ รวมทั้ง Log ของผู้ใช้ระบบ

- Business Metadata ประกอบด้วย คำจำกัดความและความหมายเกี่ยวกับธุรกิจ ผู้เป็นเจ้าของข้อมูล และนโยบายการเปลี่ยนแปลง
- Operational Metadata ประกอบด้วยสารสนเทศที่เกิดขึ้นระหว่างการทำงานของดาต้าแวร์เฮาส์ ได้แก่ การย้ายข้อมูล และการเปลี่ยนรูปของข้อมูล ข้อมูลปัจจุบันในแวร์เฮาส์ (สถานะ คงอยู่ หรือทำการลบ) การข้อมูลสำหรับการเฝ้าติดตาม เช่นข้อมูลสถิติ รายงานความผิดพลาด และการตรวจสอบ

การสร้างและการจัดการระบบแวร์เฮาส์เป็นเรื่องที่ยาก เนื่องจากความแตกต่างของเครื่องมือที่ใช้งานมีมาก ในการพัฒนาเครื่องมือที่จะใช้ออกแบบและแก้ไขสคีมา วิว สตรีปต์ ญู คิวรี และรายงาน จะต้องมีการวางแผนและวิเคราะห์เครื่องมือที่จะใช้ว่าสำหรับเรื่องต่างๆ ว่าอะไรใช้ทำการแก้ไขสคีมา หรืออัปเดตการรีเฟรช การวางแผนความจุ เครื่องมือจัดการแวร์เฮาส์จะถูกใช้ในการเฝ้าติดตามแวร์เฮาส์ การทำรายงานสถิติ และให้คำแนะนำผู้ดูแลระบบ โดยใช้แยกและตารางสรุป เวลาในการประมวลผลคิวรี ชนิดความถี่ในการ drill downs หรือ rollup เพื่อรู้ว่าผู้ใช้งานต้องการอะไร จุดสูงสุดและค่าเฉลี่ยของเวลาการทำงาน การทำรายงานเพิ่มเติม ค้นหาคิวรีที่กำลังทำงานอยู่ และตัววัดคุณภาพของบริการ สำหรับเครื่องมือในการจัดการระบบและเครือข่าย จะถูกใช้วัดความคับคั่งระหว่างไคล์แอนท์ และเซิร์ฟเวอร์ ระหว่างเซิร์ฟเวอร์แวร์เฮาส์และระบบฐานข้อมูลเพื่อใช้ในการพิจารณาสำหรับการจัดการเรื่องของการประมวลผลเช่น โหลด รีเฟรช ในขั้นตอนของการประมวล ก็จะทำให้เกิดความเหมาะสม และสามารถทำงานในลักษณะกำหนดเวลาทำงานหรือทำงานเมื่อต้องการ ซึ่งระบบจะแสดงให้เห็นลำดับของการทำงานและสถานะของงานว่าทำสำเร็จหรือ เกิดความผิดพลาด ก็จะสามารถกู้สิ่งที่ดำเนินการด้วยการ roll back

บทที่ 3

โมเดลสำหรับฐานข้อมูลแบบมุมมองหลายมิติ

ระบบฐานข้อมูลแบบมุมมองหลายมิติ เป็นระบบที่สามารถทำการหมุนเพื่อดูมุมมองข้อมูลแบบต่างๆ ซึ่งใช้สร้างแอปพลิเคชันสำหรับระบบสนับสนุนการตัดสินใจ โดยความต้องการของข้อมูลจะเป็นลักษณะเชิงสรุป เป็นการรวบรวมข้อมูล การสร้างมุมมอง หรือการประยุกต์ใช้สูตร กับข้อมูลที่มีหลายมิติ ซึ่งประกอบด้วย คิวบ์โมเดลที่เป็น hypercube ซึ่งเป็นชุดคำสั่งพื้นฐานในการออกแบบวิธีที่ใช้กันอยู่ในปัจจุบัน และฟังก์ชันที่เพิ่มเติม โดยวัตถุประสงค์ของโมเดลที่ทำให้เกิดความสมดุลสำหรับทุกมิติของข้อมูลรวมถึงการวัดค่า ด้วยโมเดลที่ช่วยให้ยืดหยุ่นในการเตรียมให้สามารถทำงานได้ในแต่ละมิติและสนับสนุน ad hoc aggregate

3.1 โมเดลข้อมูล (Data Model)

โมเดลข้อมูลนี้เป็นรูปแบบที่สนับสนุนสำหรับข้อมูลที่มีมุมมองหลายมิติ และคำสั่งที่มีใช้ อยู่ในระบบฐานข้อมูลแบบมุมมองหลายมิติทั่วไป สำหรับหลักการของโมเดลนี้จะอธิบายวิธีการทำ มิติและการวัดค่าที่ได้แต่ละสัดส่วน และทำให้การทำงานมีความยืดหยุ่น ทั้งแบบหลายลำดับชั้น การทำ ad hoc aggregate โดย logical model ของข้อมูลจะมีโครงสร้างเป็นหนึ่งหรือหลายคิวบ์ โดยแต่ละคิวบ์จะประกอบด้วย

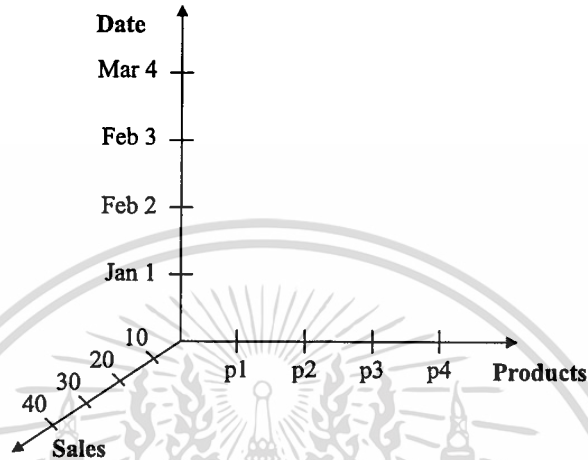
- k มิติ และแต่ละมิติจะมีชื่อ D_i อยู่ในโดเมน dom_i แล้วจะมีค่าปรากฏ
- ค่าที่กำหนดให้เป็นคู่ $E(C)$ จาก $dom_1 \times \dots \times dom_k$ ไปจนถึง n -tuple (n แถว) ถ้า $E(C)(d_1, \dots, d_k)$ จะเป็นการอ้างถึงค่า ณ ตำแหน่งที่ d_1, \dots, d_k ของ cube C
- ส่วนของเมตาดेटาที่เป็น n -tuple ของชื่อ โดยในแต่ละชื่อจะอธิบายเป็นหนึ่งในสมาชิกของค่า n -tuple ที่อยู่ใน cube

แต่ละค่าของคิวบ์สามารถเป็น 0, 1 หรือ n -tuple $\langle X_1, \dots, X_n \rangle$ ถ้าค่าตรงกันกับ $E(C)(d_1, \dots, d_k)$ จะเป็น 0 แสดงว่าค่านั้น ไม่มีอยู่ในฐานข้อมูล ถ้าเป็น 1 แสดงว่าค่านั้นตรงกัน และการแสดง n -tuple นั้นเกิดจากการเพิ่มค่าที่เป็นไปได้สำหรับค่าที่ตรงกันในแต่ละมิติ ถ้าบางค่าของคิวบ์เป็น 1 แล้วจะสามารถมีค่าที่จะเป็น n -tuple เราจะแทนค่าเหล่านั้นให้ครบมิติของคิวบ์ ซึ่งจะมีอย่างน้อยหนึ่งค่าของคิวบ์ที่ไม่ใช่ 0 แต่ถ้าค่าทั้งหมดของ cube เป็น 0 แสดงว่าเป็นคิวบ์ว่าง หรือในกรณีถ้าโดเมน dom_i ของมิติ D_i ไม่มีค่าเลขคิวบ์ก็ว่างเช่นกัน

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3.2 โอเปอเรเตอร์ (Operators)

ในการแสดงมุมมองข้อมูลจะแสดงเป็น 3 มิติ เพื่อให้เห็นมุมมองของมิติ จากรูปที่ 3-1 จะแสดง สินค้า, วันที่, ยอดขาย ซึ่ง ยอดขายไม่ได้ใช้วัดค่า แต่เป็นมิติสำหรับโมเดลนี้



รูปที่ 3-1 แสดง Logical cube ที่มียอดขายเป็นมิติ

โอเปอเรเตอร์บน logical cube มิติของยอดขายอาจจะไม่แสดง ซึ่งอยู่ในรูปดั่งนั้นจำนวนของยอดขายสามารถวัดค่าได้โดยมิติของสินค้าและวันที่ ซึ่งจะอธิบายภายหลังว่าดึงมาแสดงได้อย่างไร สำหรับตอนนี้จะใช้คู่ด้วยจำนวนยอดขายที่มีสมาชิกค่าของ cube ดั่งนั้นค่า $\langle 15 \rangle$ สำหรับ “Date = Mar 4” และ “Product = p1” ในรูปที่ 3-2 ที่แสดงนั้นใน logical cube รูปที่ 2 จะเป็นค่าที่ตรงกันกับ “Date = Mar 4”, “Product = p1” และ “Sales = 15” ซึ่งมีค่าเป็น “1” ซึ่งจะแสดงเมตาดาต้าของแต่ละค่าที่เป็นคำอธิบายในรูป ดั่งนั้น $\langle \text{Sales} \rangle$ ในรูปที่ 3-2 จะแสดงแต่ละค่าในรูปที่เป็นค่าของยอดขาย

หมายเหตุ สิ่งที่กำหนดเป็นโอเปอเรเตอร์ในการใช้ cube C ด้วย k มิติ โดยอ้างถึงมิติที่เป็น D_1, \dots, D_k ในการใช้ D_i เพื่ออ้างถึงโดเมนของมิติ D_i และจะใช้ตัวอักษรเป็นตัวเล็ก a,b,c เพื่ออ้างถึงค่าคงที่

ค่าของมิติในฟังก์ชันโมเดลข้อมูลนี้จะมีอ้างค่าของรูป เป็นผลของโปรแกรมที่เกิดจากโอเปอเรชัน อาจจะมีการแมปค่าแบบบหนึ่งต่อหนึ่งหรือหลายค่า ที่เป็นคำตอบของรูป ค่าเหล่านี้จะเป็นการรวมเข้าด้วยกันเป็นหนึ่งค่าด้วยฟังก์ชันใด ก็จะถูกเรียกว่า element combining function โดยแทนค่าเป็น f_{elem}

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

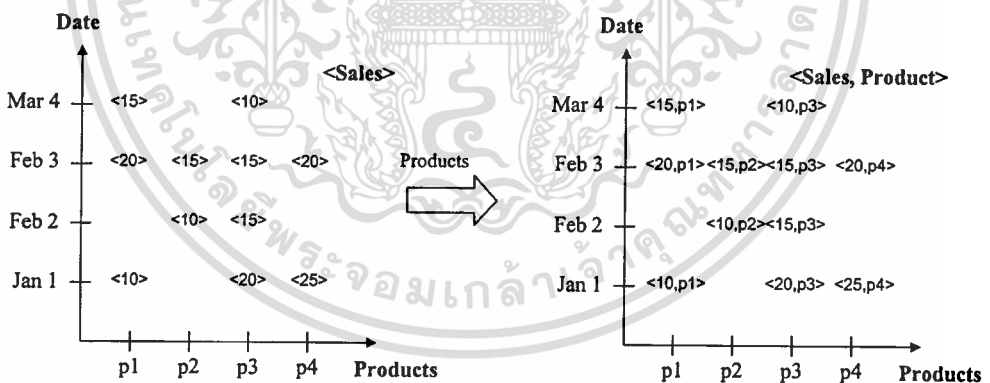
บางครั้งการเมิร์จ (merge) ของค่าตลอดชุดมิติ จะเรียกว่า dimension merging function โดยแทนด้วย f_{merge}

3.2.1 Push

Push เป็น โอเปอร์เรเตอร์ (ดูจากรูป 3-2) เป็นการคอนเวิร์ทค่าในมิติโดยการสลับตำแหน่ง โดยใช้ฟังก์ชัน f_{elem} ซึ่งโอเปอร์เรเตอร์นี้จะยอมให้มิติและการวัดค่าเป็นรูปแบบเฉพาะ

- Input : C, D_i
- Output : C ด้วยค่าที่ไม่เป็น 0 และถูกขยายโดยสมาชิกของค่าที่เพิ่ม ซึ่งเป็นค่ามิติของ D_i สำหรับแต่ละค่า ถ้าค่าไม่ใช่ n -tuple แต่เป็น "1" เมื่อทำการคอนเวิร์ทไปเป็น 1-tuple ก็จะบรรจุค่าที่ตรงกันของมิติ
- Mathematically : $push(C, D_1) = C_{ans}$

$E(C_{ans})(d_1, \dots, d_k) = g \oplus \langle d_i \rangle$ เมื่อ $g = E(C)(d_1, \dots, d_k)$ โดยเครื่องหมาย \oplus กำหนดให้เป็น 0 ถ้า $g=0$ และจะมีค่าเป็น $\langle d_i \rangle$ ถ้า $g=1$ และทุกกรณีจะทำต่อเนื่องกัน 2 แถว



รูปที่ 3-2 แสดง โอเปอร์เรชัน push ในมิติของสินค้า

3.2.2 Pull

Pull เป็นโอเปอร์เรเตอร์ที่ใช้ย้อนค่ากลับจากผลที่เกิดจากการทำโอเปอร์เรเตอร์ push โอเปอร์เรเตอร์ Pull จะทำการสร้างมิติใหม่สำหรับแยกสมาชิกในแต่ละค่า โดยการทำงานจะคอนเวิร์ตค่าที่อยู่ในมิติ ดังนั้นค่าสามารถทำการเมิร์จหรือทำการจอยน์ โอเปอร์เรเตอร์นี้จะใช้ได้ก็ต่อเมื่อมิติมีขนาดเท่ากัน (symmetric)

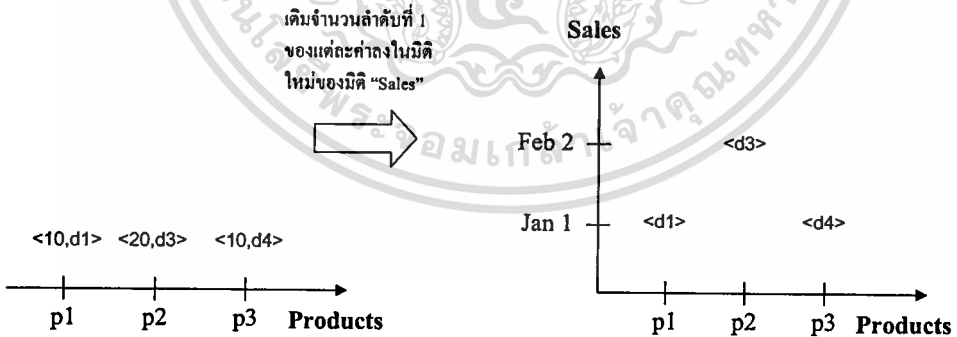
- Input : C , ชื่อใหม่ของมิติ D , Integer i
- Output : C_{ans} จะเป็นมิติ D ที่ถูกเพิ่ม โดยจะดึงค่าในตำแหน่ง i ของแต่ละของเมทริกซ์
- Constraint : ค่าทั้งหมดที่ไม่เป็น 0 ของ C จะเป็น n -tuples เพราะแต่ละค่าที่ไม่เป็น 0 ก็จะมีสมาชิกอย่างน้อยหนึ่งค่าที่สามารถสร้างมิติใหม่ได้
- Mathematically : $pull(C, D, i) = C_{ans}, 1 \leq i \leq n$

D จะกลายเป็นมิติที่ $k + 1$ ของคูปล

$$dom_{k+1}(C_{ans}) = \{ e \mid e \text{ เป็นสมาชิก } i \text{ ของบางค่าใน } E(C)(d_1, \dots, d_k) \}$$

$$E(C_{ans})(d_1, \dots, d_k) = \langle e_1, \dots, e_{i-1}, e_{i+1}, \dots, e_n \rangle \text{ ถ้า } E(C)(d_1, \dots, d_k) = \langle e_1, \dots, e_i, \dots, e_n \rangle,$$

มีเช่นนั้น $E(C_{ans})(d_1, \dots, d_k) = 0$ หมายถึง: ถ้าค่าของผลลัพธ์ไม่มีจำนวนสมาชิกจะถูกแทนค่าด้วย 1



รูปที่ 3-3 การดึง (Pull) ค่าสมาชิกตัวแรกของแต่ละค่าให้เป็นมิติของยอดขาย

3.2.3 Destroy Dimension

บ่อยครั้งที่มีความจำเป็นต้องลดขนาดจำนวนมิติในคูปล ซึ่งโอเปอร์เรเตอร์นี้จะทำการลบมิติ D นั้นออก ซึ่งอยู่ในโดเมนที่เป็นค่าเดี่ยว ความหมายของค่าเดียวนั้นให้หมายถึงค่าที่คงเหลือจาก

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

มิติ $k-1$ ซึ่งมิติของคูป $k-1$ จะไม่ซ้ำกัน ถ้ามิติ D ถูกลบออกแล้ว ผลของมิติ $k-1$ ก็จะเกิดช่องว่าง และจะถูกลดขนาด

- Input : C , ชื่อใหม่ของมิติ D_i
- Output : C_{ans} ด้วยมิติ D_i ที่ถูกทำลาย
- Constraint : D_i มีเพียงหนึ่งค่า
- Mathematically : $\text{destroy}(C, D_i)$

C_{ans} มีจำนวนมิติ $k-1$

ถ้ามิตินั้นมีหลายค่าจะไม่สามารถทำการลบโดยตรงได้เนื่องจากมีค่าที่เกี่ยวข้องกับฟังก์ชันที่ให้ค่านั้น ในกรณีที่มีหลายค่าและต้องการ destroy จะต้องทำการ merge เรียง และทำคำสั่งนี้อีกครั้ง

3.2.4 Restriction

โอเปอเรเตอร์ Restrict เป็นการทำงานบนมิติของคูปและลบค่าในมิติของคูปที่เงื่อนไขไม่ถูกต้องออก จากรูปที่ 3-4 เป็นผลที่เกิดจากการทำ restrict หมายเหตุ : คำสั่งนี้ก็คือ slicing/dicing ในการจัดการคูปของฐานข้อมูลมุมมองหลายมิติ

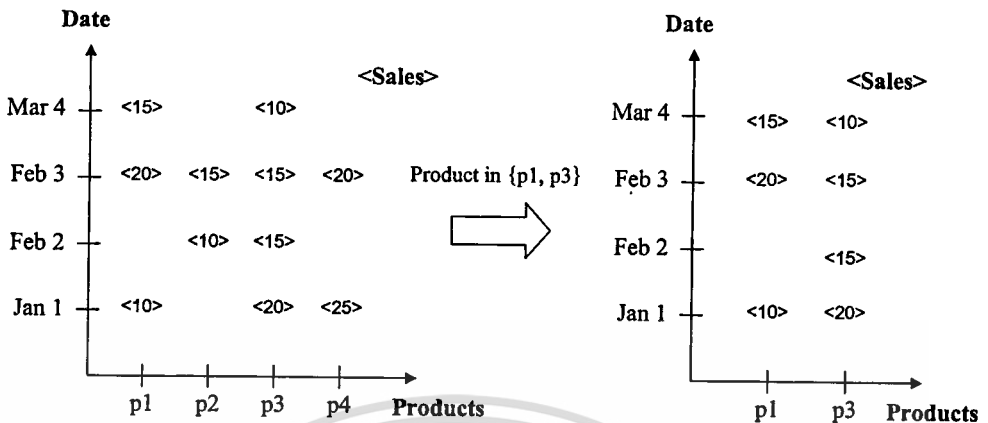
- Input : คูป C และค่า \mathcal{P} ที่อ้างอยู่บน D_i
- Output : คูปใหม่ C_{ans} ที่ทำการย้ายค่าที่มีการอ้างถึง \mathcal{P} ในมิติ D_i ออกจาก C

หมายเหตุ \mathcal{P} เป็นสิ่งที่ประกอบเป็นชุดข้อมูล และไม่ได้เป็นค่าเดี่ยว แต่ว่าถ้า \mathcal{P} เป็นสิ่งที่ถูกกำหนดโดเมนของ D_i เช่น “top 5 values” ถ้าไม่มีค่าของมิติ D_i แล้ว C_{ans} ก็จะเป็นคูปว่าง

- Mathematically : $\text{restrict}(C, D_i, \mathcal{P}) = C_{ans}$

$\text{dom}_j(C_{ans}) = \text{dom}_j(C)$ ถ้า $1 \leq j \leq k$ & $j \neq i$ ไม่เช่นนั้น $\text{dom}_j(C_{ans}) = \mathcal{P}(\text{dom}_j(C))$

$E(C_{ans})(d_1, \dots, d_k) = E(C)(d_1, \dots, d_k)$



รูปที่ 3-4 แสดงผลของโอเปอเรชั่น Restrict

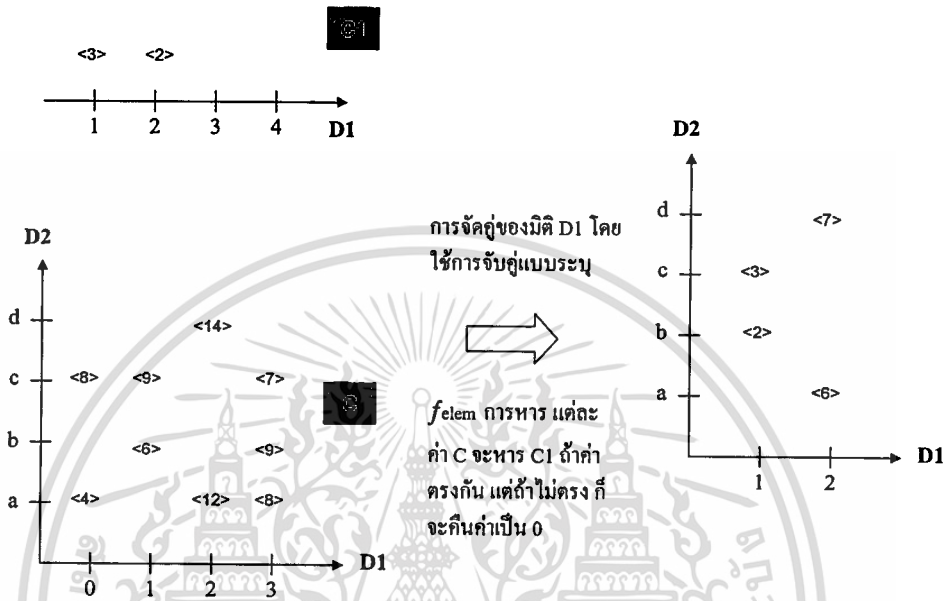
3.2.5 Join

โอเปอเรชั่นการจอยน์เป็นการใช้สร้างความสัมพันธ์ของสารสนเทศในสองคูป โดยผลของการจอยน์จะเป็น m มิติ เมื่อคูป C มีจำนวน n มิติ คูป C_1 บน k มิติ จะเรียกการจอยน์มิติว่าเป็นคูป C_{ans} ด้วย $m+n-k$ มิติ ในการจอยน์แต่ละ D_i มิติของ C จะรวมเป็นหนึ่งมิติ D_x ของ C_1 ผลลัพธ์ของมิติที่ตรงกันจะมีค่าเป็นยูเนียนของค่าใน D_i และ D_x สำหรับการทำให้ Transformation สามารถนำมาใช้เพื่อเปลี่ยนค่าของมิติ D_i และ D_x ก่อนที่จะทำการแมปเพื่อให้ผลลัพธ์ สำหรับค่าเป็นผลลัพธ์ของคูป C_{ans} จะถูกรวมเข้าด้วยกันโดยใช้ฟังก์ชัน f_{elem} ทุกค่าของ C และ C_1 นั้นจะทำการแมปไปยังค่าของคูป C_{ans}

รูปที่ 3-5 ภาพแสดงการจอยน์คูป C กับคูป C_1 บนมิติของ D_1 (ไม่มีการใช้ Transformation) มิติ D_1 ของผลลัพธ์คูปจะมีเพียงสองค่า ฟังก์ชัน f_{elem} ทำการหารทุกค่าจากคูป C ด้วยค่าจาก C_1 ถ้าแต่ละค่าเป็น 0 ก็จะทำให้ผลลัพธ์ได้ 0 ด้วย ค่าของผลลัพธ์ในมิตินั้นจะมีเพียงค่า 0 เท่านั้นก็จะไม่แสดงใน C_{ans} (เหมือนค่า 0 และ 3 ในมิติ D_1) ในการขจัดค่าจะทำตามลำดับที่อยู่ในคูปโดยตลอดแนว มิตินั้นจะต้องมีอย่างน้อยหนึ่งค่าที่ไม่ใช่ค่า 0

- Input : C กับมิติ $D_1 \dots D_m$ และ C_1 กับมิติ $D_{m-k} \dots D_n$ โดยมิติ $D_{m-k} \dots D_m$ เป็นมิติที่จอยน์ ด้วยฟังก์ชันการแมป $2k$ ฟังก์ชัน f_{m-k}, \dots, f_m นี้ถูกประกาศค่าของมิติ $D_{m-k} \dots D_m$ ของ C และ ฟังก์ชัน f'_{m-k}, \dots, f'_m เป็นการประกาศบนมิติ $D_{m-k} \dots D_m$ ของ C_1 การแมป f_i จะถูกใช้กับค่า $v \in dom_i(C)$ เพื่อสร้างค่าสำหรับมิติ D_i ใน C_{ans} ในลักษณะเดียวกัน f'_i จะถูกใช้กับค่า $v' \in dom_i(C_1)$ เพื่อสร้างค่าสำหรับ

มิติ D_i ใน C_{ans} ดังนั้นสิ่งที่ต้องมีสำหรับการใช้ฟังก์ชัน f_{elem} เป็นการรวมค่าที่มาจาก C และ $C1$ เพื่อให้ได้ผลลัพธ์เป็นค่าของ C_{ans}



รูปที่ 3-5 แสดงการjoinค่าสองคู่

- Output : C_{ans} กับ n มิติ โดย D_1, \dots, D_n ค่าต่างๆ ใน C และ $C1$ สามารถที่จะทำการแมปกับ C_{ans} ที่เป็นผลรวมของฟังก์ชัน f_{elem} ที่สร้างค่าผลลัพธ์ของ C_{ans} ถ้ามีค่า v บางค่าของมิติ D_i ค่า $E(C_{ans})(x_1, \dots, v, x_{i+1}, \dots, x_n)$ เป็น 0 สำหรับทุกค่าของมิติอื่นๆด้วย ดังนั้น v จึงไม่ควรรวมอยู่ในมิติ D_i ของ C_{ans}

- Mathematically : $join(C, C1, [f_{m-k_1}, \dots, f_m, f'_{m-k_1}, \dots, f'_m], f_{elem}) = C_{ans}$

$$dom_i(C_{ans}) = dom_i(C) \text{ ถ้า } 1 \leq i \leq m - k_1 - 1$$

$$dom_i(C_{ans}) = dom_i(C1) \text{ ถ้า } m \leq i \leq n$$

$$dom_i(C_{ans}) = \{ d^a \mid d^a \in dom_i(C) \text{ หรือ } d^a \in f'_i(d'), d' \in dom_i(C1) \}$$

$$E(C_{ans})(d_1, \dots, d_{m-k_1}, \dots, d_m, \dots, d_n) = f_{elem}(\{t1\}, \{t2\}) \text{ ดังนั้น}$$

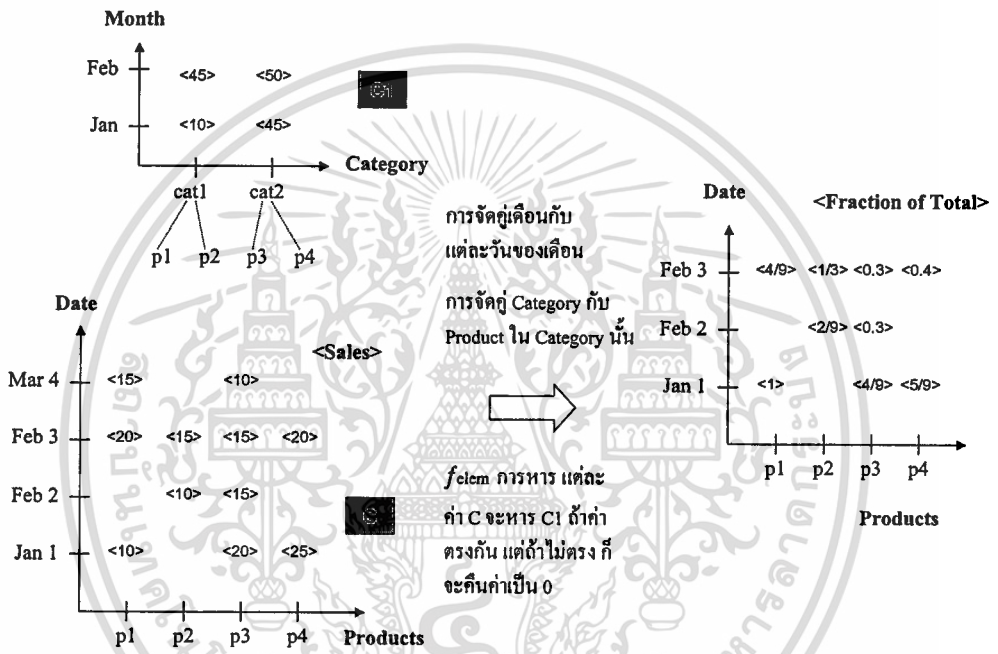
$$t1 = E(C)(d_1, \dots, d_{m-k_1}, \dots, d_m), t2 = E(C1)(d'_{m-k_1}, \dots, d'_{m-k_1}, \dots, d'_{m+1}, \dots, d'_n)$$

และ

$$d^a_i \in f_i(d_i) \text{ หรือ } d^a_i \in f'_i(d'_i) \text{ สำหรับค่า } m - k_1 \leq i \leq m$$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สำหรับโอเปอเรเตอร์จอยน์ มี 2 กรณีที่ไม่สามารถทำได้ คือ cartesian product และ associate ในกรณีของ cartesian product คือ มีสองคู่ที่ไม่มีมิติปกติที่ทำการจอยน์กันได้ ส่วน associate จะมีใช้ในระบบของ OLAP สำหรับการคำนวณที่มีลักษณะ “ทำคว้น สำหรับหายอดขายของแต่ละเดือน คิดเป็นเปอร์เซ็นต์จากยอดขายรายไตรมาส” โดย associate เป็น asymmetric และต้องการแต่ละมิติของ $C1$ มาทำการจอยน์กับบางมิติของ C รูปที่ 3-6 เป็นการแสดงการทำ associate สำหรับคู่ $C1$ กับ C



รูปที่ 3-6 แสดงการทำ associate ของสองคู่

ในแต่ละเดือนในมิติ $C1$ และแต่ละวันในมิติ C ทำการจอยน์โดยแมปที่ product ของ C_{ans} สำหรับมิติของเดือน จะมีการแมปไปยังวันที่ทั้งหมดของเดือนนั้น สำหรับมิติของ categories จะมีค่าเป็น cat1 ที่จะแมปไปยัง product p1, p2, และ cat2 จะแมปไปยัง p3, p4 สำหรับมิติวันที่และสินค้าจะมีการระบุวิธีการแมปโดยใช้ฟังก์ชัน felem การหารค่าจากคู่ C ด้วยค่าในคู่ $C1$ ถ้าแต่ละค่าเป็น 0 ผลลัพธ์ที่ได้ก็จะเป็น 0 ดังจะเห็นว่าค่าใน Mar4 ที่อยู่ใน C_{ans} ไม่มีเนื่องจากค่าทั้งหมดไม่ตรงกันเลยจึงได้ค่าเป็น 0

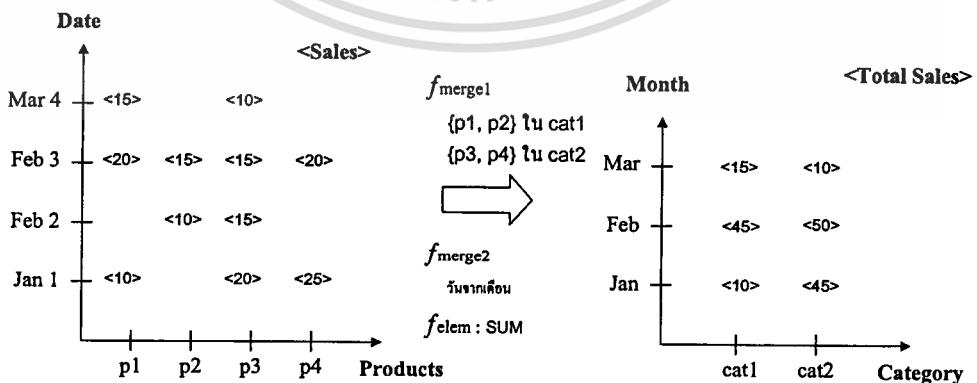
3.2.6 Merge

โอเปอเรชัน Merge เป็นการทำสรุป ดังรูปที่ 3-7 ซึ่งแสดงให้เห็นว่าลำดับชั้นในฐานข้อมูลมุมมองหลายมิติเป็นอย่างไร เพื่อที่จะทำการประยุกต์ใช้สำหรับโอเปอเรชันเมิร์จ

ฟังก์ชันในการเมิร์จเป็นการใช้การแมปชื่อ product หลายๆ ค่าที่อยู่ในหนึ่งหรือหลาย categories และฟังก์ชันอื่นๆ ที่ถูกใช้เพื่อแมปข้อมูลที่มีค่าตรงกันของเดือน ดังนั้นหลายๆค่าบนแต่และมิติที่ถูกเมิร์จจะถูกสร้างเป็นมิติใหม่ซึ่งเป็นไปได้ว่าจะมีขนาดเล็กกว่าโดเมนเดิม ผลของการเมิร์จแต่ละมิติ จะมีการรวมค่าและสามารถใส่ฟังก์ชันเฉพาะเข้าไปเพื่อหาผลสรุปได้ จากตัวอย่างของฟังก์ชัน aggregate ที่ใช้ f_{elem} คือ SUM

โดยทั่วไปการใช้ฟังก์ชันเมิร์จมิติอาจจะเป็นความสัมพันธ์แบบ หนึ่งต่อหลาย ทำให้ค่าในระบบล่างจะรวมอยู่ในค่าของระดับที่สูงกว่าของลำดับชั้น

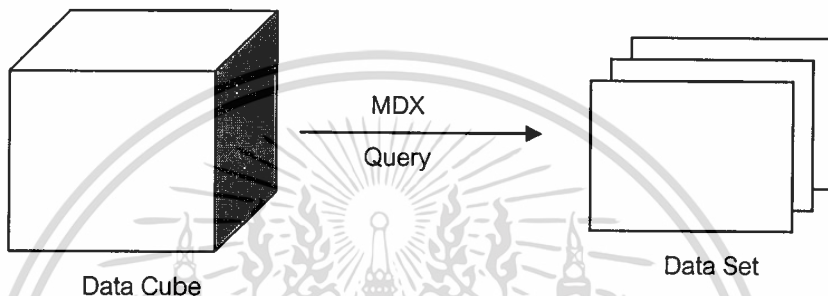
- Input : C ฟังก์ชัน f_{elem} สำหรับการเมิร์จค่าและคู่ m (มิติ, ฟังก์ชัน)
สมมติว่าคู่ m เป็น $[D_1, f_{merge_1}], \dots, [D_m, f_{merge_m}]$
- Output : คู่ C_{ans} ของบางส่วนที่มีลักษณะใกล้เคียง C จะมีการเมิร์จมิติ D_i โดยการใช้ฟังก์ชัน f_{merge_i} โดยค่าที่ตรงกันจะถูกรวมให้เป็นค่าของ f_{elem}
- Mathematically : $merge(C, \{ [D_1, f_{merge_1}], \dots, [D_m, f_{merge_m}] \}, f_{elem}) = C_{ans}$
 $dom_x(C_{ans}) = \{ f_{merge_i}(e) \mid e \in dom_x(C) \}$ ถ้า $1 \leq i \leq m$
ไม่เช่นนั้น $dom_x(C_{ans}) = dom_x(C)$
 $E(C_{ans})(d_1, \dots, d_k) = f_{elem}(\{ t \mid t = E(C)(d_1, \dots, d_k) \})$
เมื่อ $f_{merge_i}(d'_i) = d_i$ ถ้า $1 \leq i \leq m$ ไม่เช่นนั้น $d'_i = d_i$



รูปที่ 3-7 แสดงการเมิร์จมิติ Date และ Product โดยใช้ $f_{elem} = \text{SUM}$

3.3 Multidimensional Query Language

ภาษาสำหรับการค้นหาข้อมูลหลายมิติ เป็นภาษาที่คิดค้นเพิ่มเติมจาก SQL พื้นฐาน โดยรูปแบบของภาษาจะมีความแตกต่างกันไป ขึ้นกับผลิตภัณฑ์ที่ใช้งาน ในรายงานนี้ได้นำเสนอ Multidimensional Expression (MDX) ซึ่งสนับสนุน OLE DB สำหรับ OLAP ซึ่งลักษณะของการทำงานของ MDX Query จะเป็นการดึงข้อมูลจากคิวบ์เพื่อให้ได้ค่าตัดเซตสำหรับนำมาใช้งาน ดังรูปที่ 3-8



รูปที่ 3-8 แสดงภาพเบื้องต้นของการทำงานของ MDX Query

คำจำกัดความของ MDX Data Cube

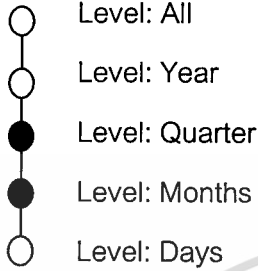
Dimension	หมายถึง เซตของกลุ่มข้อมูลที่มีชนิดเดียวกัน
Member	หมายถึง สมาชิกของ Dimension
Hierarchy	หมายถึง กลุ่มสมาชิกของ Dimension
Member-Properties	หมายถึง รายละเอียดเพิ่มเติมของแต่ละ Member

ตัวอย่าง MDX Data Cube

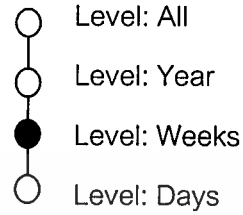
4 Dimensions	SalesPerson, Time, Measures, Product
SalesPerson-dimension	Komol, Wanchai, ...,Siththisak
Time-dimension	2001, Jan 2002, Qtr2 2003, ...
Measures-dimension	Sales, Profit-Loss
Product-dimension	Computer, Printer, ...
Member-Properties of SalesPerson-dimension	Name, Phone, ...

รูปแบบ Hierarchy-Schemata บน Time-dimension

Hierarchy: Quarters



Hierarchy: Weeks

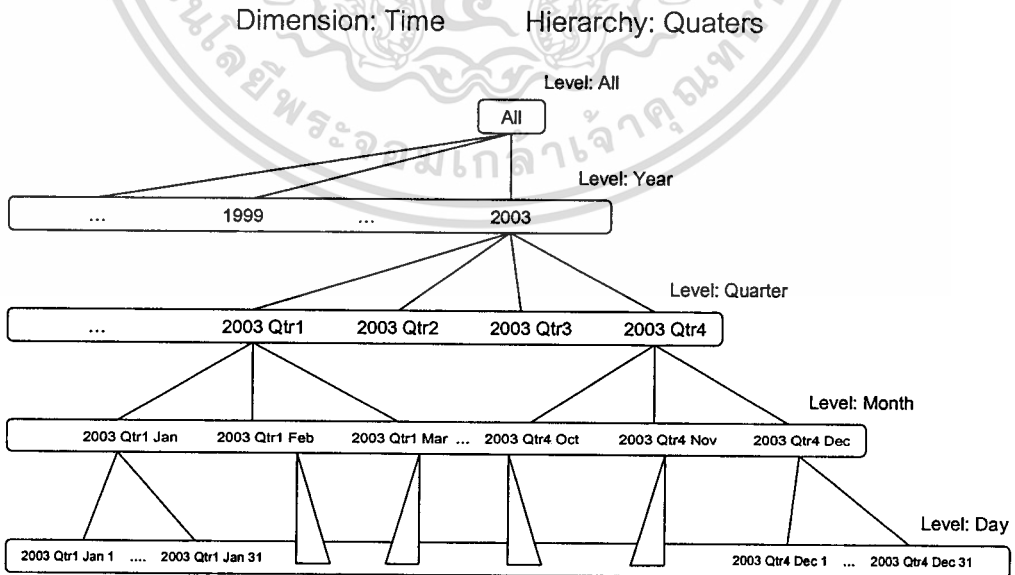


รูปที่ 3-9 แสดงภาพโครงสร้างลำดับชั้น ที่กำหนดเป็นรูปแบบของ Time-dimension

จากรูปที่ 3-9 แสดงถึงลำดับของลำดับชั้นของเวลา โดยแบ่งประเภทของแต่ละลำดับดังนี้

- หมายถึง สมาชิกของ โดเมนชั้นที่บรรจุค่าลำดับชั้นที่แตกต่างกัน
- หมายถึง สมาชิกของ โดเมนชั้นที่บรรจุค่าเพียงหนึ่งลำดับชั้น

ตัวอย่าง โครงสร้างข้อมูลลำดับชั้นบน Time-dimension

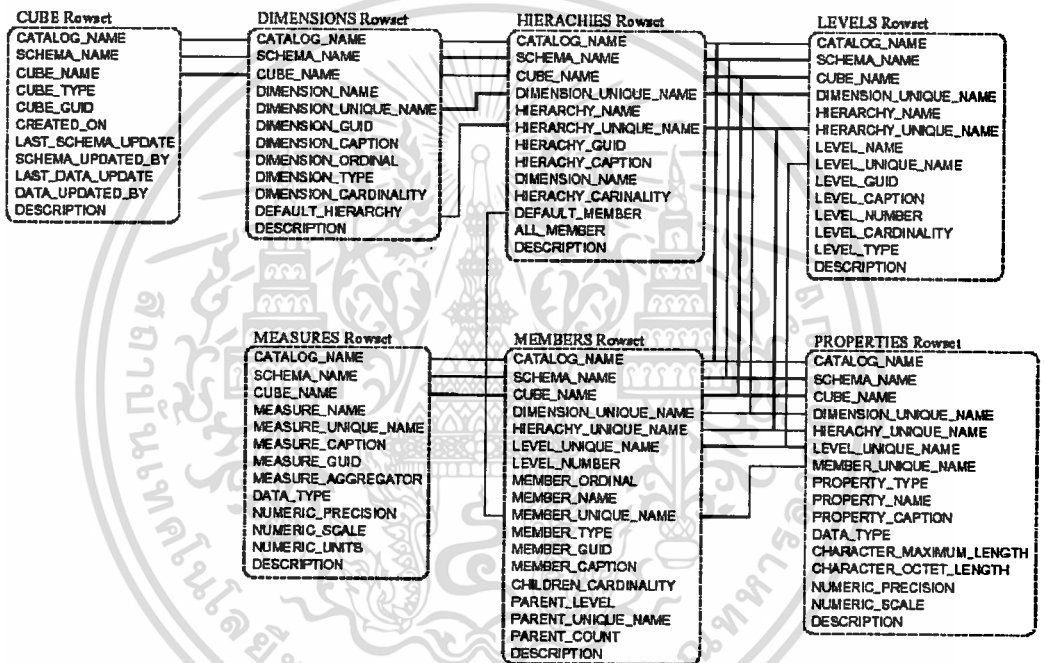


รูปที่ 3-10 แสดงภาพ โครงสร้างข้อมูลลำดับชั้นบน Time-dimension

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จากรูปที่ 3-10 เป็นภาพแสดงถึงความสัมพันธ์ในการจัดกลุ่มของลำดับชั้นของไคเมนชัน โดยแบ่งลำดับชั้นตามรูปแบบ Quarter ที่แสดงดังรูปที่ 3-9 แล้วนั้น ซึ่งลำดับของข้อมูลที่เห็นได้ แบ่งเป็นลำดับของความสัมพันธ์คือชั้นบนสุดจะแสดงความสัมพันธ์ทั้งหมดของปี และแต่ละปี จะมีลำดับชั้นความสัมพันธ์ของไตรมาสของปีนั้น และแต่ละไตรมาสจะสัมพันธ์กับเดือนและแต่ละเดือนจะสัมพันธ์กับช่วงวันของเดือนนั้นด้วย

โดยลำดับความสัมพันธ์ของ Dimension จะมีโครงสร้างของการจัดเก็บเมตาดาต้าสำหรับ Rowset ของ CUBE ดังรูปที่ 3-11



รูปที่ 3-11 แสดงภาพโครงสร้างของการจัดเก็บเมตาดาต้าในCUBE Schema

โครงสร้างภาษาสำหรับ MDX

รูปแบบของภาษาของคิวรีที่ใช้ในการสอบถามข้อมูลใน Cube มีโครงสร้างดังนี้

```
[WITH <formula_spec>]
SELECT <axis_spec> [, <axis_spec>, ...]
FROM <cube_spec> [, <cube_spec>, ...]
WHERE < slicer_spec> [cell_properties]
```

คำอธิบาย Syntax ใน MDX

- WITH** : ใช้เพื่อสร้างโดเมนชั้นใหม่ ด้วยสิ่งที่อยู่ในประโยคของ Expression เช่น
 WITH MEMBER Measures.[% Change from Last Year's PolicyCount] AS
 PolicyCount / (PolicyCount, Year.PREVMEMBER)
- <axis_spec>** : การกำหนดโดเมนชั้นให้อยู่แกนมิติที่กำหนด เช่น
 2543, 2544 ON COLUMNS
- <cube_spec>** : การระบุข้อมูลที่อยู่ในคิวบ์ (ส่วนมากจะระบุเพียงตัวเดียว)
- < slicer_spec >** : จำกัดโดเมนชั้น ไม่ให้อยู่ <axis_spec> ที่กำหนดไว้

ตัวอย่าง MDX Query

```
SELECT CROSSJOIN({PolicyType.CHILDREN},
  {VehicleType.CHILDREN,
    110, 120})
  ON COLUMNS
  {BrandModel.CHILDREN}
  ON ROWS
FROM Policy_Profile
WHERE (UYr, [2543], PolicyCount.All)
```

บทที่ 4

การพัฒนาระบบคลังข้อมูลการประกันวินาศภัย

จากทฤษฎีเกี่ยวกับเทคโนโลยีดาต้าแวร์เฮาส์ และโมเดลในการพัฒนาฐานข้อมูลหลายมิติ ดังที่กล่าว ดาต้าแวร์เฮาส์ ที่ได้กล่าวไปแล้ว สามารถนำมาใช้ในการพัฒนาระบบคลังข้อมูลการประกันวินาศภัยได้ ซึ่งในโครงการพัฒนาระบบนี้ จะประกอบด้วยกระบวนการเตรียมข้อมูล การออกแบบ การประมวลผลข้อมูล และการนำเสนอสารสนเทศมุมมองหลายมิติ เมื่อได้สารสนเทศที่ผ่านการประมวลผลแล้ว จะเป็นขั้นตอนในการเรียกใช้สารสนเทศในรูปแบบ Ac Hoc Query

4.1 วัตถุประสงค์ในการจัดทำดาต้าแวร์เฮาส์

ในการจัดทำดาต้าแวร์เฮาส์สำหรับระบบงานประกันวินาศภัยรถยนต์ภาคสมัครใจ มีวัตถุประสงค์ในการลดเวลาของการประมวลผลข้อมูล ช่วยให้ง่ายต่อการสอบถาม ซึ่งสารสนเทศจากการวิเคราะห์จะช่วยให้ผู้ใช้งานข้อมูลสามารถนำไปใช้ร่วมสำหรับการตัดสินใจ สำหรับ แนวโน้มและทิศทางทางการตลาดที่เป็นไปได้ในอนาคต โดยข้อมูลดังกล่าวจะเป็นภาพรวมของตลาด ทั้งนี้ผู้ใช้งานสามารถนำข้อมูลที่มีอยู่ขององค์กรตนเองมาใช้เปรียบเทียบได้ง่ายขึ้น

4.2 การเตรียมข้อมูลและการทำดาต้าคัตลิ่ง

เนื่องด้วยข้อมูลจากฐานข้อมูลในระบบงานประกันวินาศภัยเดิมจะเป็นข้อมูลในลักษณะทรานแซกชัน และจะต้องมีการประมวลผลเชิงสถิติเพื่อออกรายงาน ซึ่งข้อมูลดังกล่าวยังไม่มี ความเหมาะสมในเรื่องของข้อมูล เช่น ไม่มีการระบุรุ่นของรถยนต์ หรือค่าของข้อมูลบางค่าอยู่นอก โดเมนของข้อมูลที่ต้องการวัดค่าทำให้เกิดปัญหาในการประมวลผลที่นานมาก จึงมีความจำเป็นที่ จะต้องเตรียมรายการทรานแซกชัน เหล่านี้ให้อยู่ใน View หรือเป็นข้อมูลที่พร้อมในการประมวลผล ดังมีขั้นตอนดังนี้

4.2.1 การกำหนดโครงสร้างฐานข้อมูลสำหรับดาต้าแวร์เฮาส์

เป็นขั้นตอนที่สำคัญ เพราะความถูกต้องของระบบจะต้องมาจากพื้นฐานการได้มาของ ข้อมูลที่ถูกต้อง ซึ่งข้อมูลที่จะต้องสร้างได้แก่ ข้อมูลการรับประกันภัย และข้อมูลสินไหมทดแทน เป็นส่วนสำคัญในการวิเคราะห์ระบบงาน

POLICY_PROFILE	CLAIM_PROFILE
UYr	UYr
InsNo	InsNo
FleetNo	PolicyNo
Fleet_Count	AcctYr
PolicyNo	IssueDate
AcctYr	PolType
IssueDate	BeginDate
PolType	ExpireDate
BeginDate	ExposureUnit
ExpireDate	TotPremPerExpUnit
ExposureUnit	MainCoverPremPerExpUnit
TotPremPerExpUnit	::
::	VehType
VehType	VehGroup
VehGroup	VehBrand
VehBrand	VehModel
VehModel	VehSubModel
VehSubModel	VehYear
VehYear	VehAge
VehAge	VehSize
VehSize	::
AddEquip	OD_ClaimCount
ODFT_SI_Band	OD_LI_Total
ODFT_SI	OD_Paid_Total
PD_LIMIT	Fire_TotalLossCount
BI_PersonLimit	Fire_ClaimCount
BI_EventLimit	Theft_ClaimCount
BasePrem	THEFT_LI
::	THEFT_Paid
FleetDiscRate	THEFT_OS
FleetDisc_PolicyCount	::
FleetDiscAmt	Recovery
::	
PolStatus	

รูปที่ 4-1 แสดงภาพ Data warehouse schema ของระบบคลังข้อมูลการประกันวินาศภัย

- ข้อมูลการรับประกันภัย

เป็นข้อมูลเกี่ยวกับการรับประกันภัยรถยนต์ สำหรับข้อมูลที่ใช้ในการวัดค่า
ได้แก่ เบี้ยรับประกันภัยและเบี้ยส่วนเพิ่มสำหรับความคุ้มครองประเภทต่างๆ, เบี้ย
ประกันภัยรวม, วงเงินความคุ้มครอง และ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ข้อมูลที่ใช้ในการกำหนดโดเมนของข้อมูลได้แก่ ประเภทของกรรมธรรม์, บริษัทรับประกันภัย, ปีรับประกันภัย, ประเภทรถยนต์, ยี่ห้อ-รุ่นของรถยนต์, อายุรถยนต์, ระดับความเสี่ยง

- ข้อมูลสินไหมค่าทดแทน

เป็นข้อมูลเกี่ยวกับการจ่ายค่าสินไหมทดแทน เมื่อเกิดเหตุตามปีรับประกันภัยและเลขที่กรรมธรรม์ที่ตรงกับการข้อมูลการรับประกันภัย สำหรับข้อมูลที่ใช้ในการวัดค่า ได้แก่ วงเงินคุ้มครอง, ความเสียหายต่อทรัพย์สิน, ความเสียหายต่อร่างกาย, ความเสียหายที่เกิดจากการโจรกรรมหรือไฟไหม้, สินไหมรวม

ข้อมูลที่ใช้ในการกำหนดโดเมนของข้อมูลได้แก่ ประเภทของกรรมธรรม์, บริษัทรับประกันภัย, ปีรับประกันภัย, ประเภทรถยนต์, ยี่ห้อ-รุ่นของรถยนต์, อายุรถยนต์, ระดับความเสี่ยง ในการเตรียมข้อมูลสำหรับข้อมูลการรับประกันภัยและสินไหมค่าทดแทน จะทำการเตรียมด้วยกระบวนการ Data Transformation และ Store Procedure

4.2.2 การทำดาต้าคลีนนิ่ง (Data Cleaning)

กระบวนการทำดาต้าคลีนนิ่ง หมายถึง การลดความผิดพลาดในตัวข้อมูลให้สามารถประมวลผลได้ถูกต้อง หรือใกล้เคียงมากที่สุด สำหรับปัญหาของการเกิดความผิดพลาดของข้อมูลในระบบมาเนื่องจากที่มาของข้อมูล ทั้งข้อมูลการรับประกันภัยและข้อมูลสินไหม มาจากบริษัทรับประกันภัย ซึ่งมีการป้อนข้อมูลด้วยวิธีที่แตกต่างกัน เมื่อรวบรวมข้อมูลเหล่านั้นไว้ด้วยกัน ย่อมทำให้เกิดความหลากหลายของข้อมูล อย่างเช่น ยี่ห้อรถยนต์ จะไม่มีมาตรฐานของการบันทึกชื่อ เป็นต้น ดังนั้นในขั้นตอนการเตรียมทำ Data Cleaning มีดังต่อไปนี้

- พิจารณาข้อมูลที่จะใช้ในการสอบถาม หรือใช้เป็นมุมมองของรายงานที่ต้องการ ซึ่งเหล่านี้ถือเป็นโดเมนชั้น ซึ่งทั้งข้อมูลการรับประกันภัย และข้อมูลสินไหมทดแทน ซึ่งได้แก่ ปีรับประกันภัย, ประเภทกรรมธรรม์, ประเภทรถยนต์, ยี่ห้อ-รุ่นของรถยนต์ เป็นต้น
- ตรวจสอบค่าของแต่ละมุมมองที่สนใจ แล้วพิจารณาค่าของข้อมูลเหล่านั้นว่าอยู่นอกโดเมนที่กำหนดหรือไม่

ตารางที่ 4-1 แสดงชื่อของมุมมองข้อมูลของกลุ่มข้อมูลจากดาต้าแวร์เฮาส์ และ โดเมนข้อมูลที่จะทำ
คลีนนิ่ง

โดเมนชั้น	โดเมนที่กำหนด	ค่าที่อยู่นอกโดเมน
ปีปรับประกันภัย	..., 2540, 2541, ...	
ประเภทกรมธรรม์	1, 2, 3	
ประเภทรถยนต์	1, 2, 3, 4, 5, {ไม่ระบุ}	<Null> จะเปลี่ยนเป็น ไม่ระบุ
ยี่ห้อ-รุ่นรถยนต์	Toyota, Honda, Nissan, Mazda, Ford, BMW, Volvo, Benz, {Other}	ทุกค่าที่อยู่นอกจาก List ที่สนใจ จะถือว่าเป็น Other
อายุรถ	1, 2, ..., 10, {99}	ค่าอื่นๆ ที่เกิน 10 จะเปลี่ยนเป็น 99
จังหวัดทะเบียนรถ	กท, อย, ... {NA}	<Null> หรือ ไม่อยู่ใน List จะถือว่าเป็นค่า NA
หน่วยความเสี่ยง	-1, {0}, 1, ..., 10	<Null> จะตั้งค่าเป็น 0 ถ้าค่า > 10 จะให้ค่าเป็น 10 ถ้าค่า < -1 จะให้ค่าเป็น -1
กรมธรรม์ระบุชื่อ คนขับ	{0}, 1, 2	<Null> จะตั้งค่าเป็น 0 คือไม่ระบุ

จากตารางที่ 4-1 แสดงให้ทราบชื่อของมุมมองข้อมูลที่ต้องเตรียมความพร้อมสำหรับการประมวลผล ซึ่งค่าของกลุ่มโดเมนจะมีทั้งที่เป็นช่วง หรือเป็นค่าที่กำหนด และค่าที่เป็นวงเล็บ {} จะถือว่าเป็นค่า Default สำหรับใช้ในการทำ Data Cleaning ถ้าค่าข้อมูลที่ตรวจสอบอยู่นอกโดเมนดังกล่าว

จากตัวอย่างเช่น ยี่ห้อ-รุ่นรถยนต์ ที่มีอยู่ในปัจจุบันมีหลายยี่ห้อและหลายรุ่นมาก แต่ในการสรุปข้อมูล ผู้สอบถามต้องการรู้เพียงไม่กี่ยี่ห้อ ซึ่งต้องการเฉพาะยี่ห้อตลาดเท่านั้น เช่น Toyota, Honda, Nissan, Mitsubishi, Benz, BMW, Volvo เป็นต้น ส่วนยี่ห้ออื่นๆ เช่น Daihatsu มีจำนวนน้อยมากเมื่อเทียบกับตลาดรถยนต์ ยี่ห้ออื่น ดังนั้น ค่าข้อมูลที่อยู่นอกการที่ระบุ จะเปลี่ยนเป็น Other ทำให้สามารถประมวลผลข้อมูลได้ครบทุกรายการ

4.3 การออกแบบและทำการประมวลผลข้อมูล OLAP

เนื่องจากข้อมูลที่อยู่ในดาต้าแวร์เฮาส์ เราไม่สามารถใช้งานข้อมูลเหล่านั้นโดยตรงได้อย่างมีประสิทธิภาพ จึงจำเป็นที่จะต้องนำข้อมูลเหล่านั้นมาผ่านกระบวนการประมวลผลที่เหมาะสม จึงจะทำให้การเรียกใช้ข้อมูลทำงานได้มีประสิทธิภาพสูง

โดยกระบวนการในการออกแบบเป็นการกำหนดเกี่ยวกับเมตาดาต้าของโครงสร้างไดเมนชัน, คิว แล้วจึงทำการประมวลผลข้อมูล

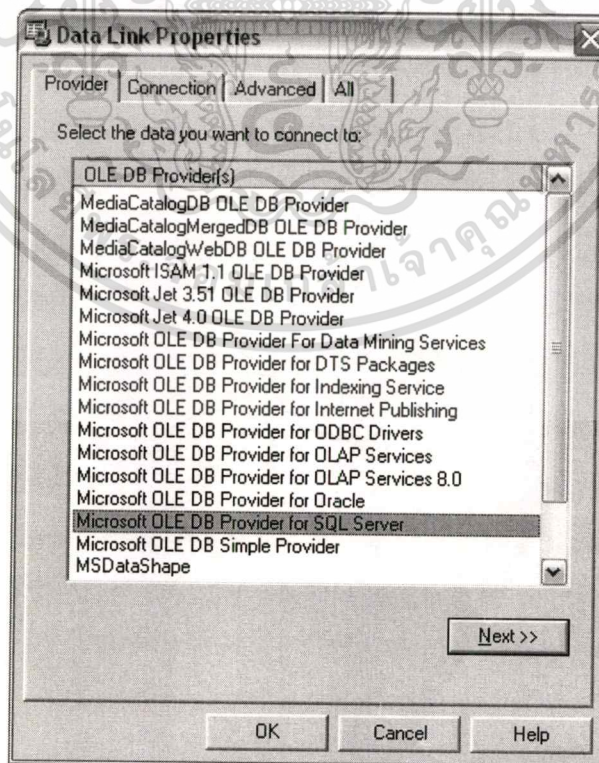
4.3.1 การระบุแหล่งที่มาของข้อมูล (Data Source)

แหล่งที่มาของข้อมูล เป็นข้อมูลที่จะใช้ในการประมวลผล ซึ่งอาจจะเป็น Text File, Excel, ระบบฐานข้อมูล ซึ่งจะสามารถสร้างตัวเชื่อมโยงแหล่งข้อมูลได้ตามที่มีอยู่ใน Provider ของระบบ เช่น การใช้ ODBC เพื่อกำหนด Data Source เชื่อมโยงฐานข้อมูล Oracle เป็นต้น สำหรับระบบนี้ จะทำการเชื่อมโยงฐานข้อมูล Microsoft SQL

ขั้นตอนที่ 1: เลือกการเชื่อมต่อโดยเลือก “Microsoft OLE DB Provider for SQL”

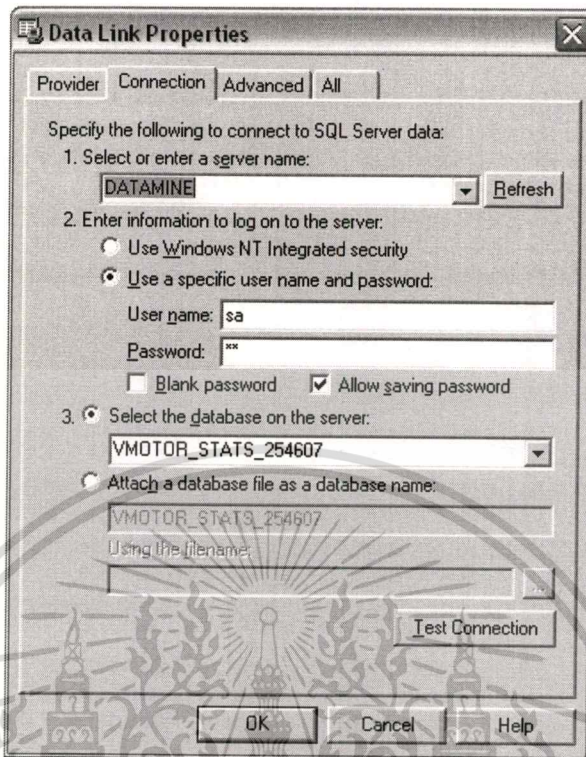
ขั้นตอนที่ 2: ระบุชื่อ Server, User Name, Password, ชื่อฐานข้อมูล

ขั้นตอนที่ 3: ตั้งชื่อ Data Source เพื่อใช้ในการอ้างอิงการเชื่อมโยงข้อมูล



รูปที่ 4-2 แสดงการสร้าง Data Source (โดยเลือก Provider Connection)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น มิอนุญาตให้เผยแพร่ไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4-3 แสดงการสร้าง Data Source (ในการระบุ Parameter สำหรับ Connection)

4.3.2 การสร้างโดเมนชั้น (Dimensions)

สำหรับระบบงานที่พัฒนา จะเลือกการสร้างโดเมนชั้นที่บันทึกค่าไว้เป็นแบบ Shared Dimension เพราะค่าข้อมูลเหล่านี้ จะมีการถูกใช้ร่วมกันในรูปต่างๆ ที่สร้าง

ขั้นตอนที่ 1: เลือกประเภทของ Dimension ซึ่งแต่ละประเภทจะมีวิธีการกำหนดค่าที่แตกต่างกัน ในที่นี้เลือกแบบ Star Schema

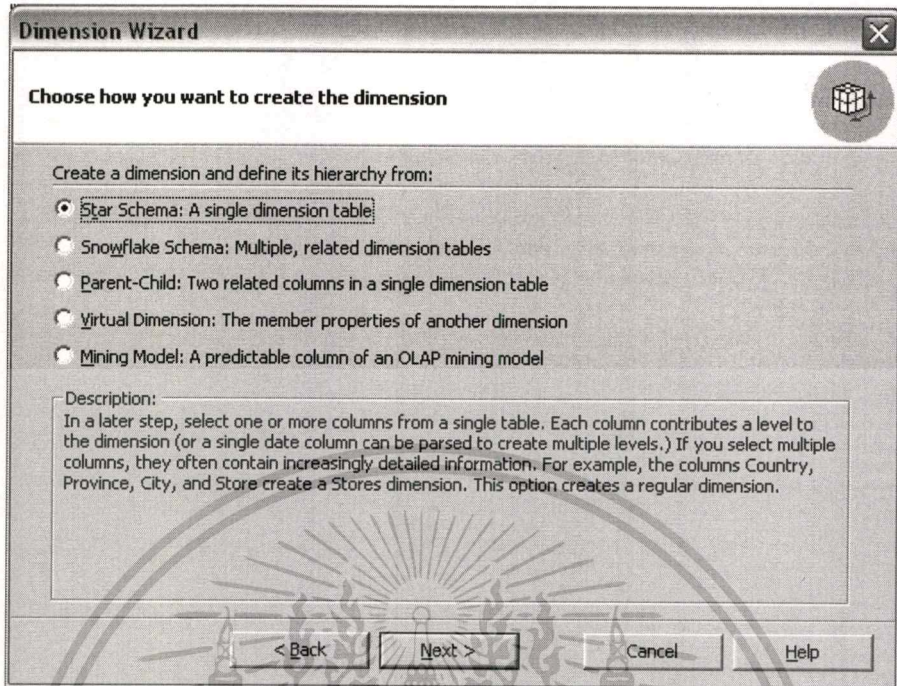
ขั้นตอนที่ 2: เลือก Table ของจาก Data Source ที่ใช้เป็น Dimension

ขั้นตอนที่ 3: กำหนด Level จากรายการ Field ที่อยู่ใน Table ที่เลือก โดย Field นี้จะใช้สำหรับ แสดงค่าข้อมูลในมิตินี้

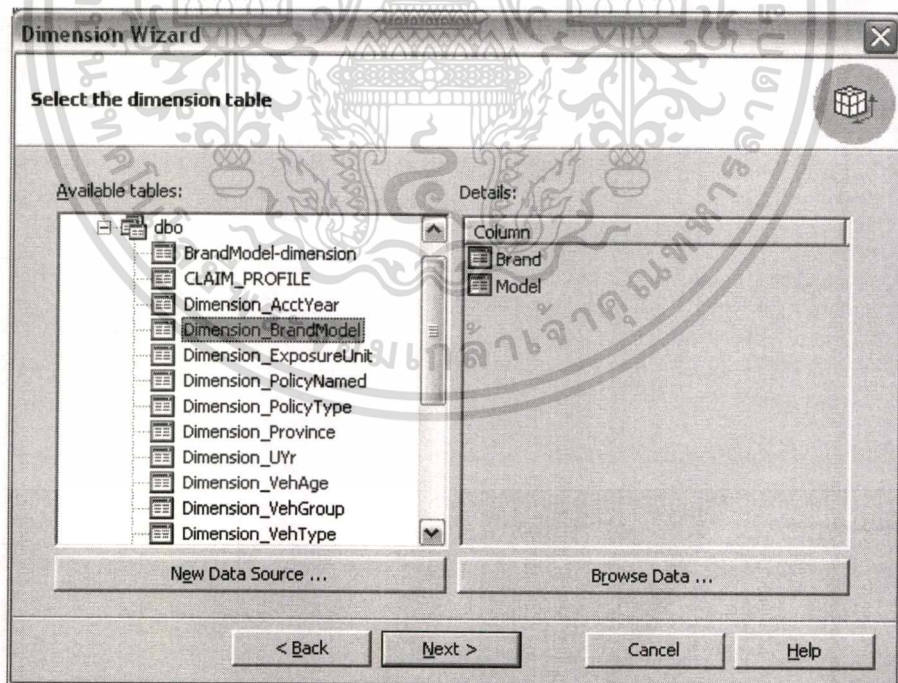
ขั้นตอนที่ 4: กำหนดฟิลด์ข้อมูลที่ใช้เป็นคีย์ ซึ่งคีย์อาจจะไม่ใช่ชื่อเดียวกับที่เลือกเป็น Level

ขั้นตอนที่ 5: บันทึกชื่อโดเมนชั้น

ขั้นตอนที่ 6: ปรับแก้ค่าของ Dimension มุมมองการออกแบบ

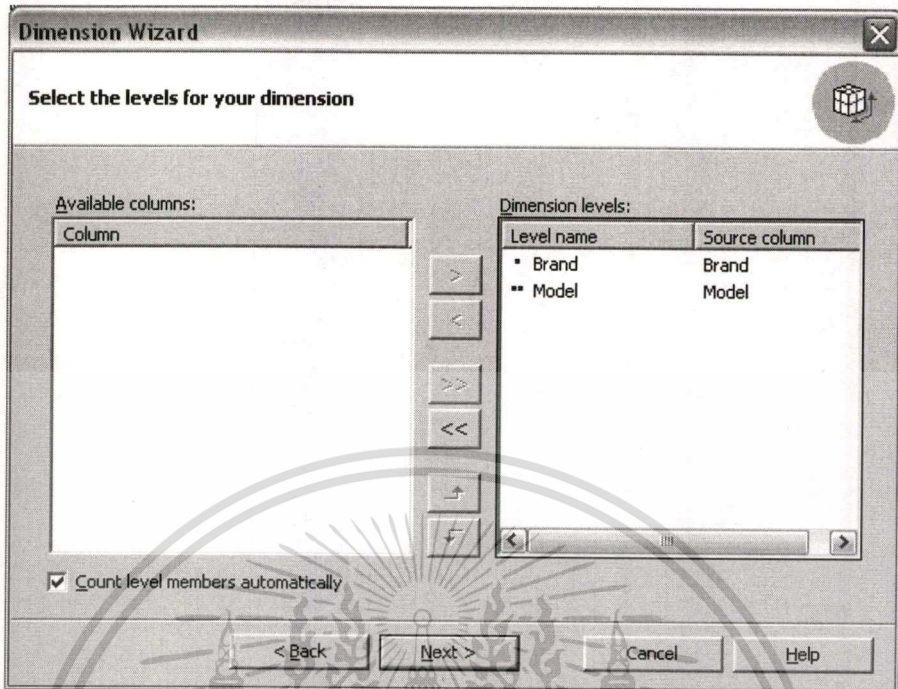


รูปที่ 4-4 แสดงการสร้างไดเมนชัน (ขั้นตอนที่ 1)

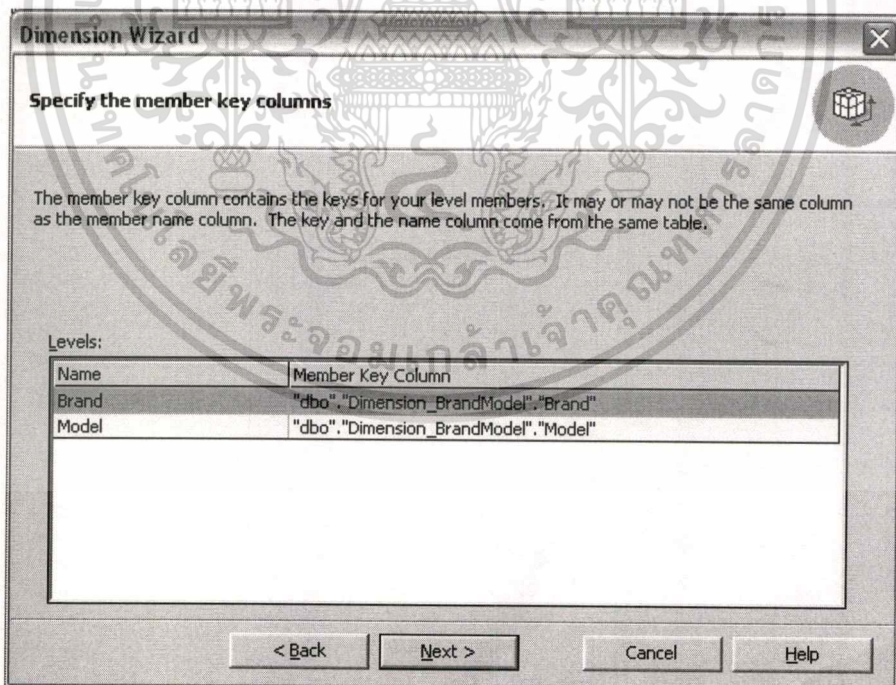


รูปที่ 4-5 แสดงการสร้างไดเมนชัน (ขั้นตอนที่ 2)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

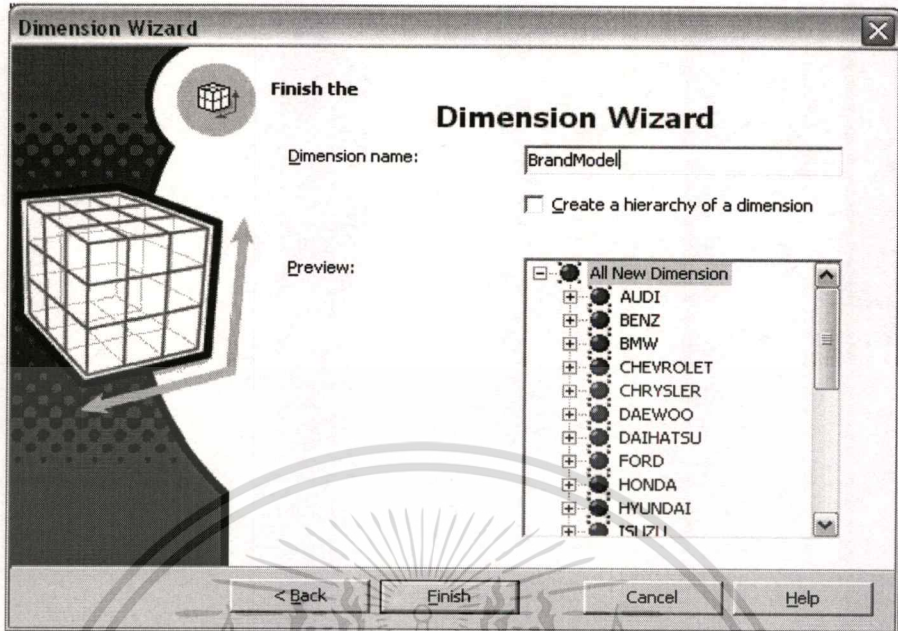


รูปที่ 4-6 แสดงการสร้างไคเมนชัน (ขั้นตอนที่ 3)

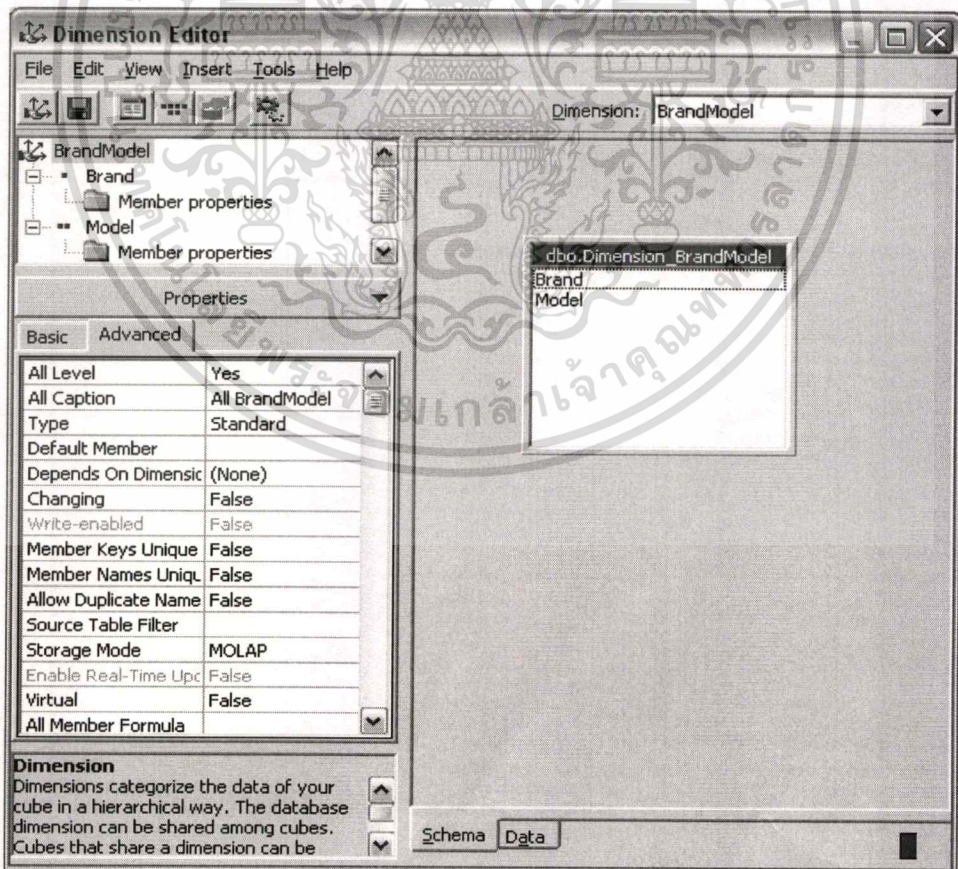


รูปที่ 4-7 แสดงการสร้างไคเมนชัน (ขั้นตอนที่ 4)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4-8 แสดงการสร้างโดเมนชั้น (ขั้นตอนที่ 5)



รูปที่ 4-9 แสดงการสร้างโดเมนชั้น (ขั้นตอนที่ 6)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4-2 แสดงรายการ Dimension และ Level ที่ทำสำหรับการประมวลผลข้อมูล

Dimension Name	Level Name	Description
UnderwritingYear	UYR	ปีรับประกันภัย
BrandModel	Brand Model	ยี่ห้อ - รุ่น รถยนต์
DriverNamed	Policy Named	การระบุชื่อผู้ขับขี่
ExposureUnit	Exposure Unit	หน่วยความเสี่ยง
PolicyType	Policy Type	ประเภทกรมธรรม์
Province	Province	จังหวัดที่จดทะเบียนรถยนต์
VehicleType	Veh Type	ประเภทรถยนต์
VehicleGroup	Veh Group	กลุ่มรถยนต์ (ตามลักษณะรถ และการใช้งาน)
VehicleAge	Veh Age	อายุรถยนต์

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

4.3.3 การสร้างคิวบ์ (Cube)

ตามทฤษฎีของการสร้างคิวบ์ซึ่งจะทำการเก็บข้อมูลได้หลายมิติ ซึ่งสามารถสร้างได้ทั้งแบบบันทึกข้อมูลที่ประมวลผลจริง และการสร้างคิวบ์ในลักษณะวิว (มุมมองย่อยในคิวบ์) สำหรับการสร้างคิวบ์ทั่วไปมีขั้นตอนดังนี้

ขั้นตอนที่ 1: เลือก Table ที่ต้องการสร้างคิวบ์จาก Data Source

ขั้นตอนที่ 2: เลือก Fields สำหรับที่ใช้ทำ Measures ซึ่งได้แก่ Fields เป็นข้อมูลที่ต้องการนำเสนอ และสามารถคำนวณวัดค่าได้

ขั้นตอนที่ 3: กำหนด Dimension ที่ใช้ในการสอบถาม / ประมวลผล

ขั้นตอนที่ 4: บันทึกชื่อ Cube

ขั้นตอนที่ 5: ปรับแต่งค่าของการ Join สำหรับ Dimensions, กำหนด Aggregate สำหรับแต่ละ Measure ที่เลือก ได้แก่ Max, Min, Sum, Average เป็นต้น

ขั้นตอนที่ 6: ทำการประมวลผลข้อมูล

ในการสร้างคิวบ์ สามารถทำแบบ Virtual Cube (วิวของคิวบ์) ได้หลังจากที่ทำการประมวลผลข้อมูลเสร็จแล้ว โดยขั้นตอนสำหรับการสร้าง Virtual Cube มีดังนี้

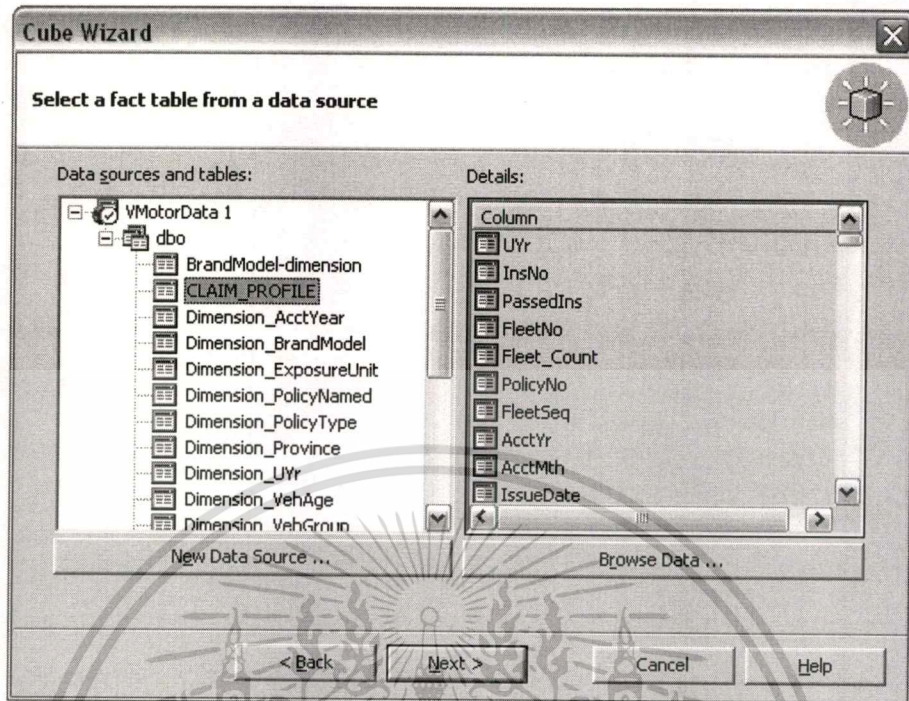
ขั้นตอนที่ 1: เลือก Cube ที่ต้องการสร้าง View ขึ้นใน Cube นั้น

ขั้นตอนที่ 2: เลือก Measure ที่ต้องการแสดง

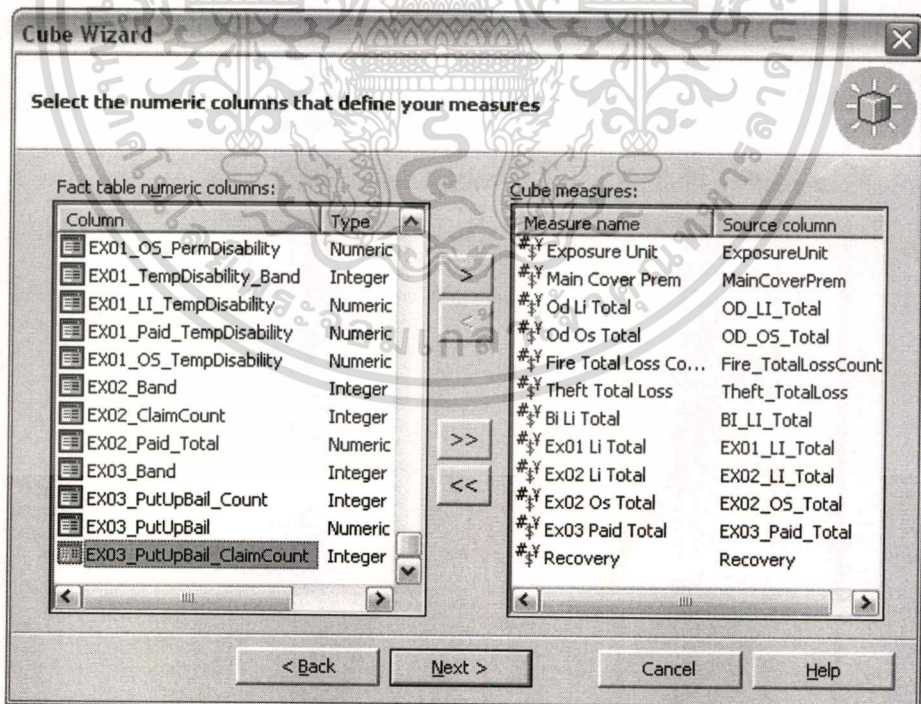
ขั้นตอนที่ 3: เลือก Dimension ที่ต้องการให้เป็นเงื่อนไขการแสดงผลข้อมูล

ขั้นตอนที่ 4: บันทึกชื่อ Virtual Cube

ขั้นตอนที่ 5: ปรับแก้ไขมุมมองการออกแบบ

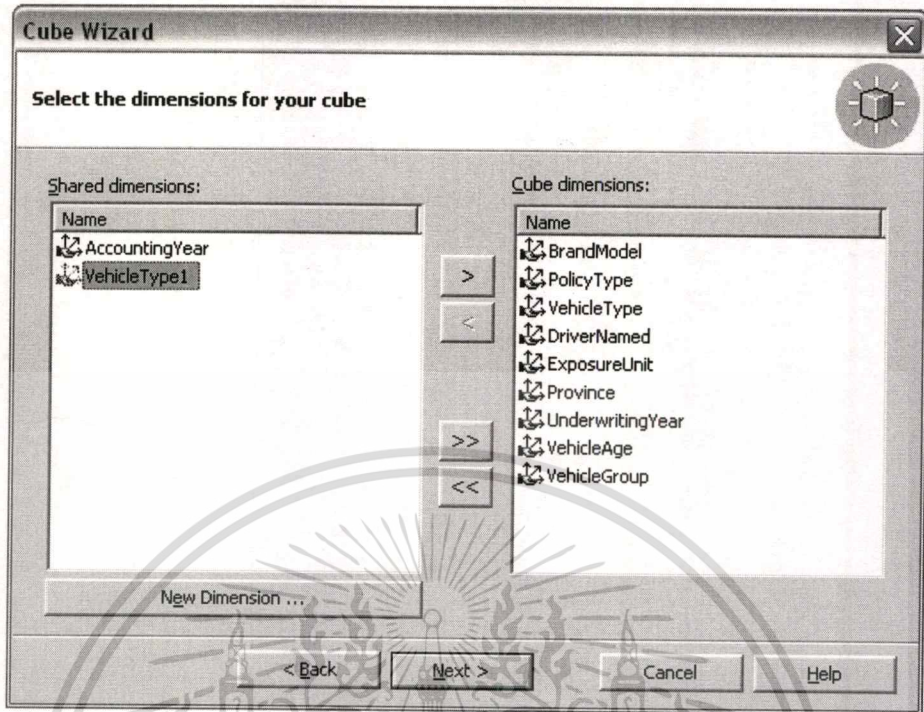


รูปที่ 4-10 แสดงการสร้างคิวบ์ (ขั้นตอนที่ 1)

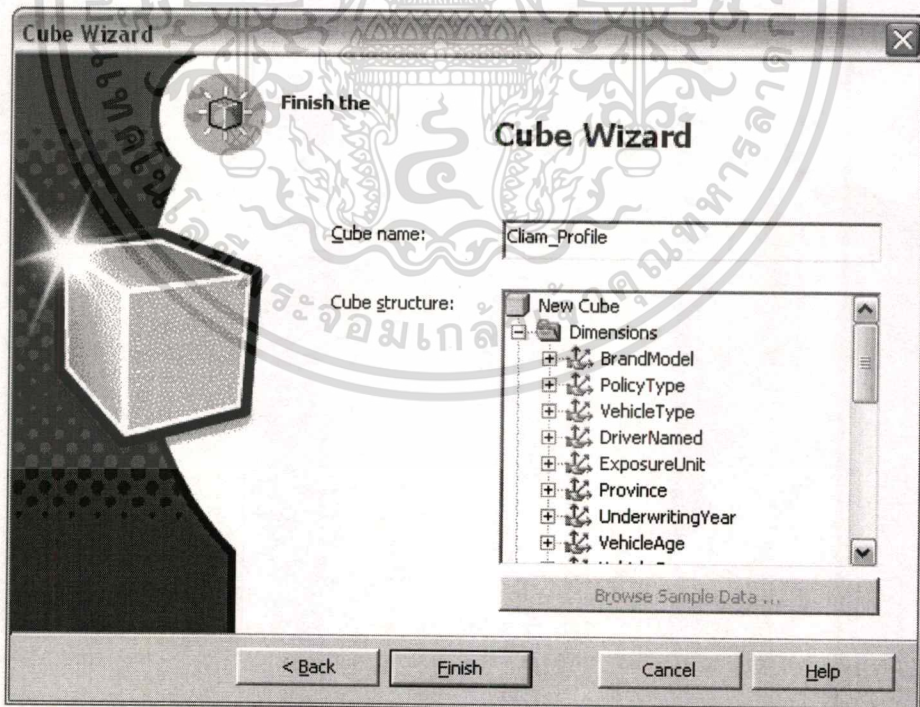


รูปที่ 4-11 แสดงการสร้างคิวบ์ (ขั้นตอนที่ 2)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

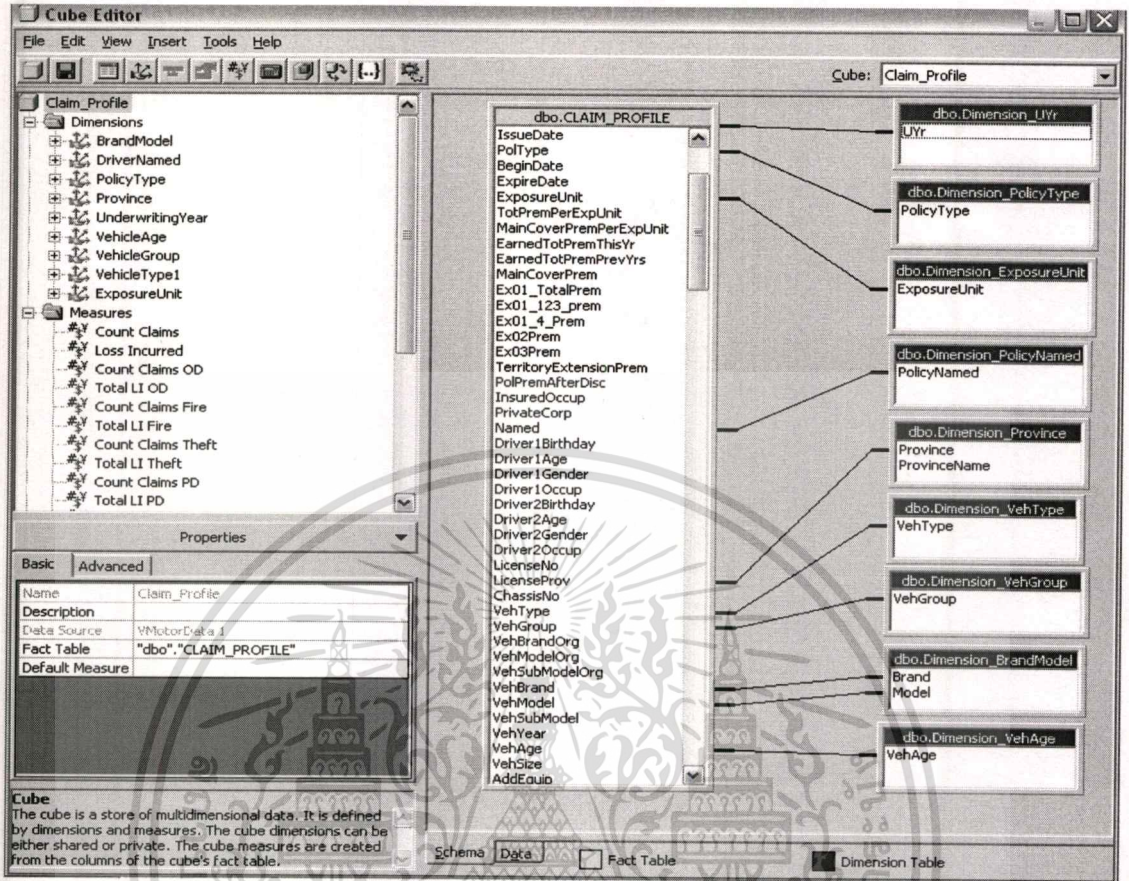


รูปที่ 4-12 แสดงการสร้างคิวบ์ (ขั้นตอนที่ 3)

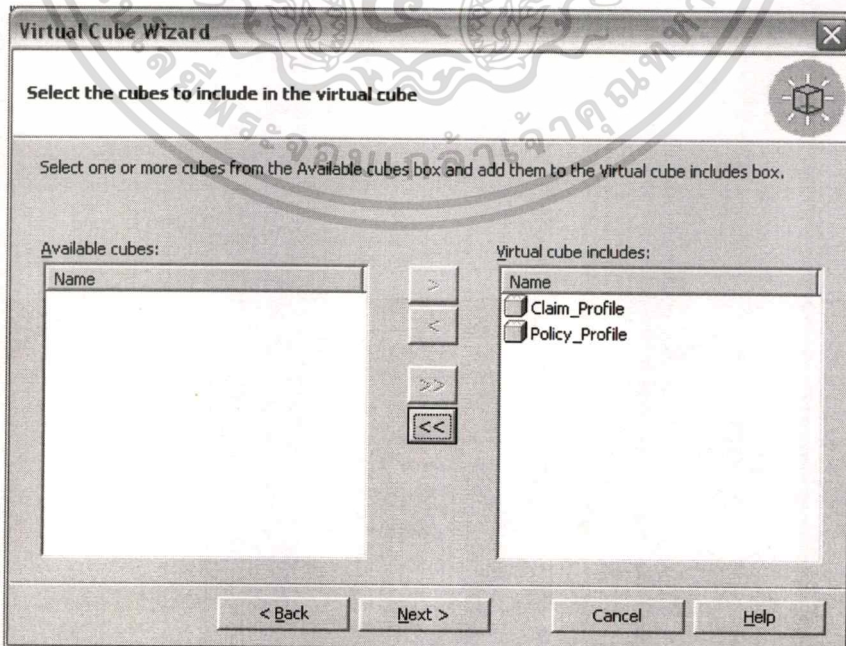


รูปที่ 4-13 แสดงการสร้างคิวบ์ (ขั้นตอนที่ 4)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

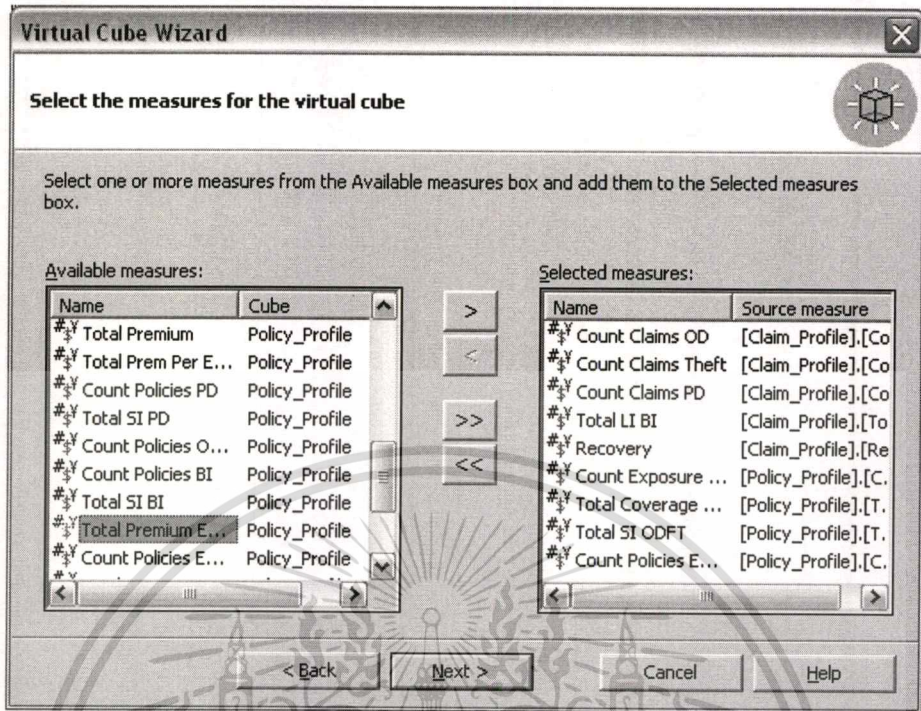


รูปที่ 4-14 แสดงการสร้างคิวบ์ (ขั้นตอนที่ 5)

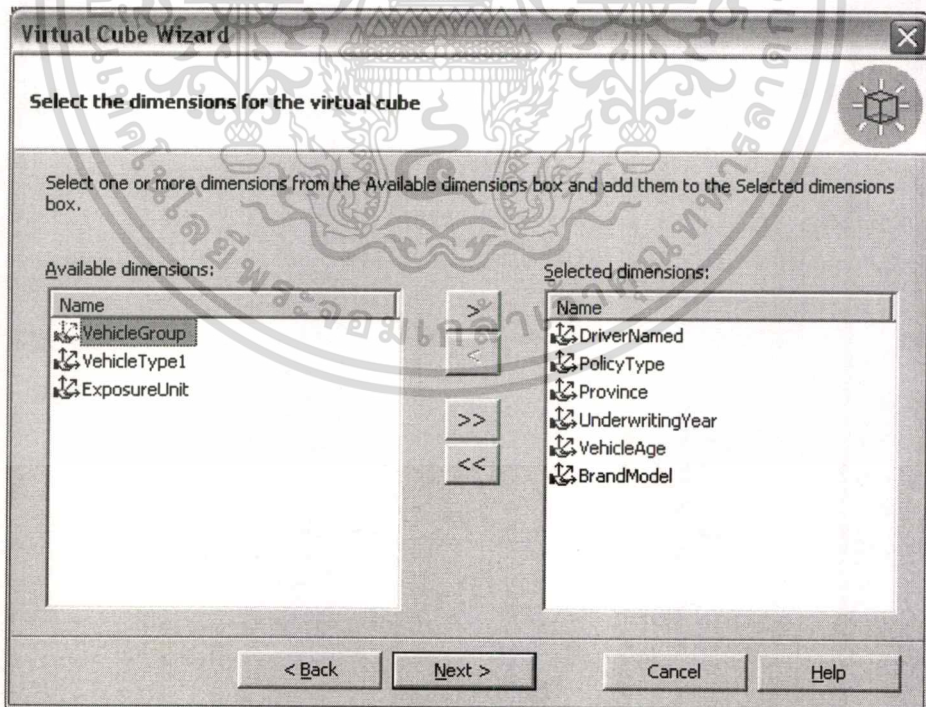


รูปที่ 4-15 แสดงการสร้าง Virtual Cube (ขั้นตอนที่ 1)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับใช้ภายในเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้เผยแพร่ไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

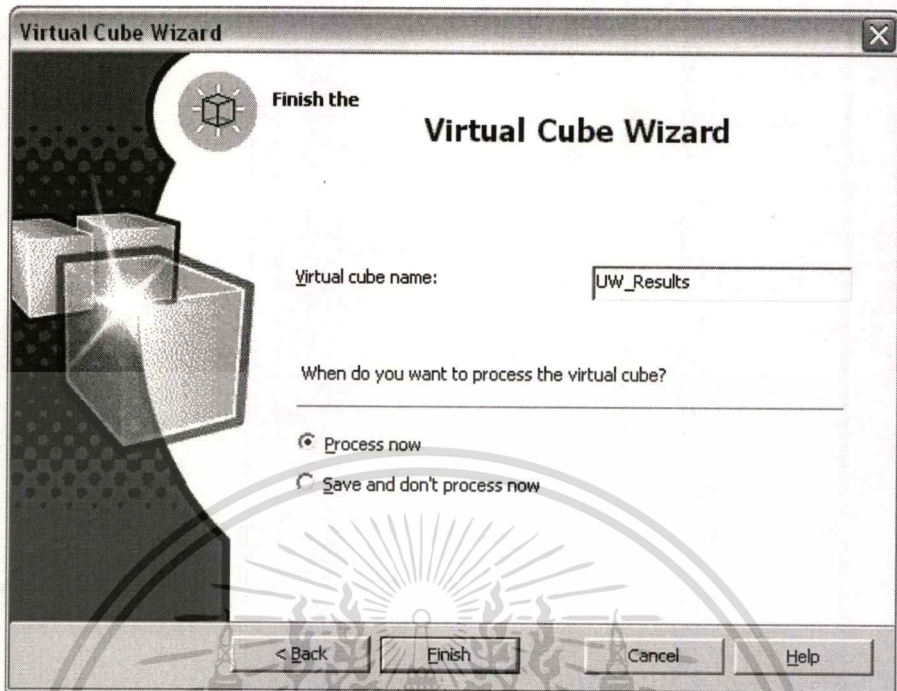


รูปที่ 4-16 แสดงการสร้าง Virtual Cube (ขั้นตอนที่ 2)



รูปที่ 4-17 แสดงการสร้าง Virtual Cube (ขั้นตอนที่ 3)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4-18 แสดงการสร้าง Virtual Cube (ขั้นตอนที่ 4)

Brand	Count Policies	Count Exposure Unit	Total Premium
All BrandModel	164,042	164,042	2,909,406,043.00
+ AUDI	856	856	26,473,010.00
+ BENZ	7,264	7,264	222,794,957.00
+ BMW	3,999	3,999	117,900,219.00
+ CHEVROLET	448	448	19,227,808.00
+ CHRYSLER	417	417	10,638,514.00
+ DAEWOO	292	292	3,354,085.00
+ DAIHATSU	323	323	2,999,251.00
+ FORD	1,676	1,676	27,687,796.00
+ HONDA	25,515	25,515	471,640,477.00
+ HYUNDAI	1,566	1,566	21,258,534.00
+ ISUZU	18,558	18,558	310,397,517.00
+ JEEP CHEROKEE	668	668	22,656,356.00
+ KIA	406	406	9,182,354.00

รูปที่ 4-19 แสดงการสร้าง Virtual Cube (ขั้นตอนที่ 5)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

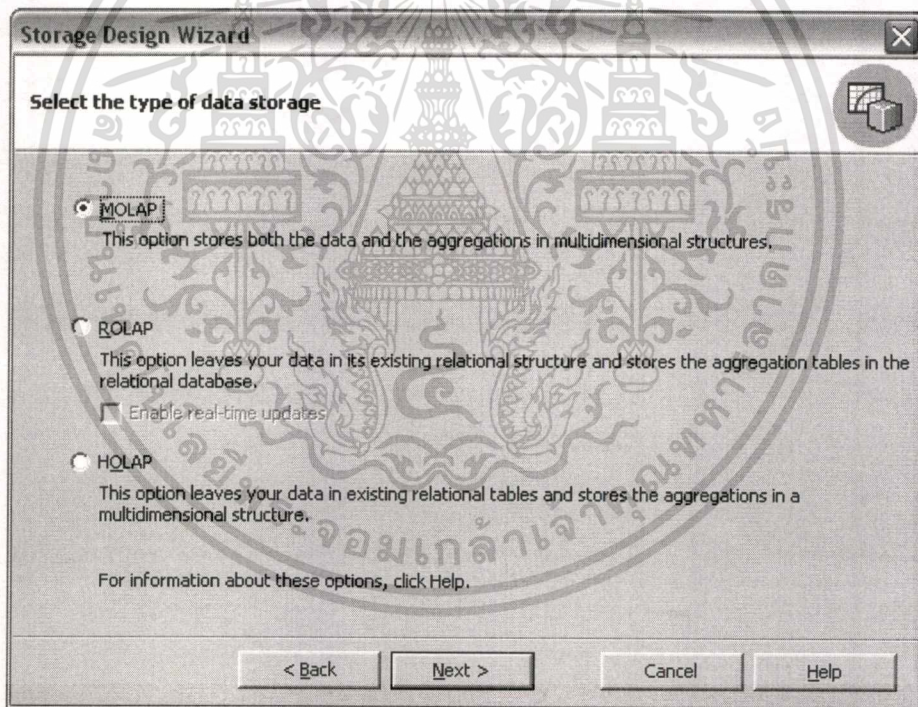
4.3.4 การออกแบบที่จัดเก็บข้อมูล (Storage Design)

โดยการออกแบบวิธีการจัดเก็บข้อมูล จะเป็นขั้นตอนที่ทำให้ผู้ออกแบบระบบสามารถกำหนดปัจจัยที่ช่วยในการเพิ่มประสิทธิภาพ โดยจะทำการเปรียบเทียบระหว่างประสิทธิภาพความเร็วในการเข้าถึงข้อมูลและพื้นที่จัดเก็บข้อมูลที่มีขั้นตอนดังนี้

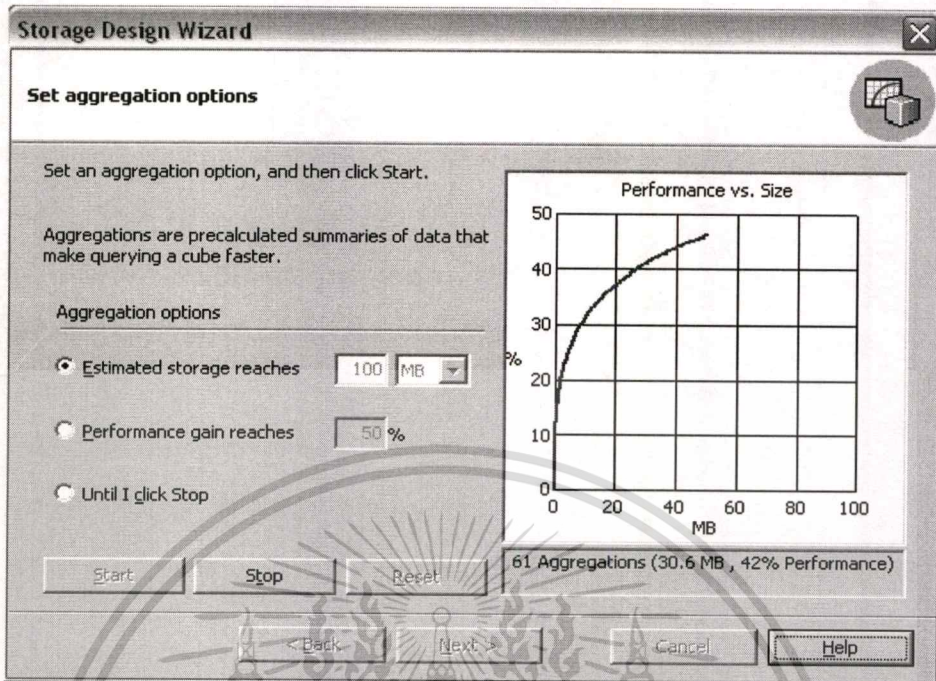
ขั้นตอนที่ 1: เลือกประเภทของวิธีการบันทึกข้อมูลได้แก่ MOLAP, ROLAP, HOLAP ซึ่งแต่ละวิธีการจัดเก็บจะแตกต่างกันในเรื่องของแหล่งที่จัดเก็บข้อมูล และการรวมข้อมูลกับ Aggregates

ขั้นตอนที่ 2: กำหนดความพื้นที่สำหรับบันทึกข้อมูล หรือกำหนดระดับความเร็วใจการเข้าถึงข้อมูล เพื่อเปรียบเทียบความเร็วและพื้นที่จัดเก็บข้อมูลที่เหมาะสม

ขั้นตอนที่ 3: ทำการ Reprocess เพิ่มเริ่มดำเนินการประมวลผลข้อมูลใหม่ทั้งหมด



รูปที่ 4-20 แสดงขั้นตอนการออกแบบพื้นที่จัดเก็บข้อมูล (ขั้นตอนที่ 1)



รูปที่ 4-21 แสดงขั้นตอนการออกแบบพื้นที่จัดเก็บข้อมูล (ขั้นตอนที่ 2)

4.3.5 การประมวลผลข้อมูล (Process or Reprocess)

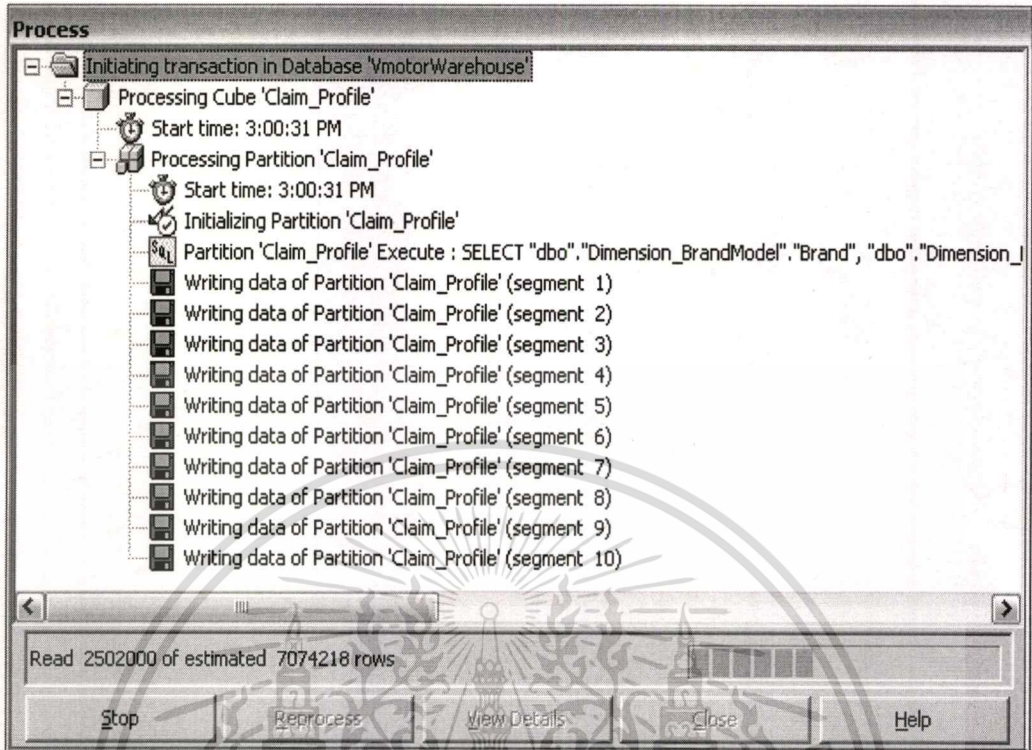
สำหรับคิวบ์, โดเมนชั้นที่ได้ออกแบบไปนั้น จะยังไม่สามารถแสดงข้อมูลได้จนกว่าจะทำการประมวลผลข้อมูล(Process / Reprocess) เพื่อทำการบันทึกข้อมูล โดยสามารถเลือกประมวลผลเฉพาะ คิวบ์ หรือเลือกประมวลผลทั้งหมดได้ตามรูปที่ 4-22

ในระบบประกันวินาศภัยรถยนต์ภาคสมัครใจได้สร้างไว้ 3 กลุ่ม คือ

1) Policy Profile (Policy_Profile) เป็นข้อมูลของการรับประกันภัย ซึ่งจะประกอบด้วยข้อมูลเกี่ยวกับประเภทกรมธรรม์, ประเภทรถยนต์ ซึ่งมีตัววัดค่า เช่น เบี้ยประกันภัย, ความคุ้มครอง เป็นต้น

2) Claim Profile (Claim_Profile) เป็นข้อมูลของการจ่ายสินไหมทดแทน ซึ่งจะประกอบด้วยข้อมูลเกี่ยวกับประเภทกรมธรรม์, ประเภทรถยนต์ ซึ่งมีตัววัดค่า เช่น วงเงินคุ้มครอง, จำนวนสินไหมต่อตัวทรัพย์สิน, ความเสียหายต่อร่างกาย, ความเสียหายจากการลักทรัพย์และไฟไหม้, จำนวนสินไหมค้างจ่าย เป็นต้น

3) Under Writing เป็นข้อมูลผลการรับประกันภัย ซึ่งจะเป็นการสรุปข้อมูลที่เกิดของ Policy_Profile และข้อมูลใน Claim_Profile ซึ่งจะเป็นการสร้างคิวบ์ประเภท Virtual Cube



รูปที่ 4-22 แสดงหน้าจอกการประมวลผลข้อมูลของคิวบ์ทั้งหมด

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4-3 แสดงตัวอย่างของการออกแบบ “Claim_Profile” Cube Schema

Measure	Measure	Dimension Name	Description
UYr		UnderwritingYear	ปีรับประกันภัย
InsNo			รหัสบริษัทประกันภัย
PassedIns			สถานะของบริษัทประกันภัยผ่านเกณฑ์
FleetNo			เลขที่กลุ่ม
Fleet_Count			จำนวนในกลุ่ม
PolicyNo	Count		เลขที่กรมธรรม์
FleetSeq			ลำดับกลุ่ม
AcctYr			ปีบัญชี
AcctMth			เดือนบัญชี
IssueDate			วันที่เริ่มมีผลบังคับ
PolType		PolicyType	ประเภทกรมธรรม์
BeginDate			วันที่เริ่มต้นคุ้มครอง
ExpireDate			วันที่สิ้นสุดคุ้มครอง
ExposureUnit	Sum		หน่วยความเสี่ยง
TotPremPerExpUnit	Sum		เบี้ยรวมต่อหน่วยความเสี่ยง
MainCoverPremPerExpUnit	Sum		เบี้ยความคุ้มครองหลักต่อหน่วยความเสี่ยง
EarnedTotPremThisYr			รายได้รวมของปีนี้
MainCoverPrem			เบี้ยความคุ้มครองหลัก
::			
VehBrand		BrandModel.Brand	ยี่ห้อรถยนต์
VehModel		BrandModel.Model	รุ่นรถยนต์
VehAge		VehicleYear	อายุรถยนต์
Recovery	Sum		จำนวนเงินคืน

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

4.4 การพัฒนาระบบส่วนการสอบถามและนำเสนอสารสนเทศ

ระบบสอบถามและนำเสนอสารสนเทศนี้เป็นการพัฒนาด้วย Active Server Page (ASP) และให้บริการผ่านเครือข่ายอินเทอร์เน็ต ในการนำเสนอข้อมูลแบบสำหรับโครงสร้างฐานข้อมูลแบบหลายมิติ สามารถที่จะแสดงข้อมูลให้อยู่ได้ทั้งในรูปแบบตาราง (2 มิติ) และการแสดงรายงานลักษณะ Pivot Table ซึ่งมีความสามารถในลักษณะ Dynamic คือการกรองข้อมูล การเปลี่ยนแกนข้อมูล ซึ่งจะนำมาใช้ในโครงการนี้

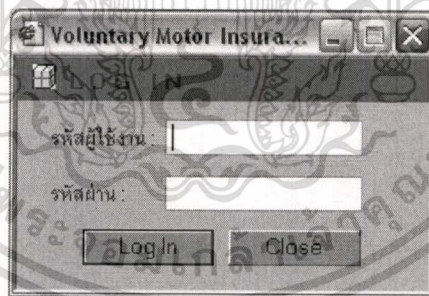
การพัฒนาการสอบถามและการแสดงรายงานสารสนเทศ

ขั้นตอนที่ 1: การล็อกอินเข้าระบบเพื่อระบุสิทธิ์ของผู้ใช้งาน (รูปที่ 2-23)

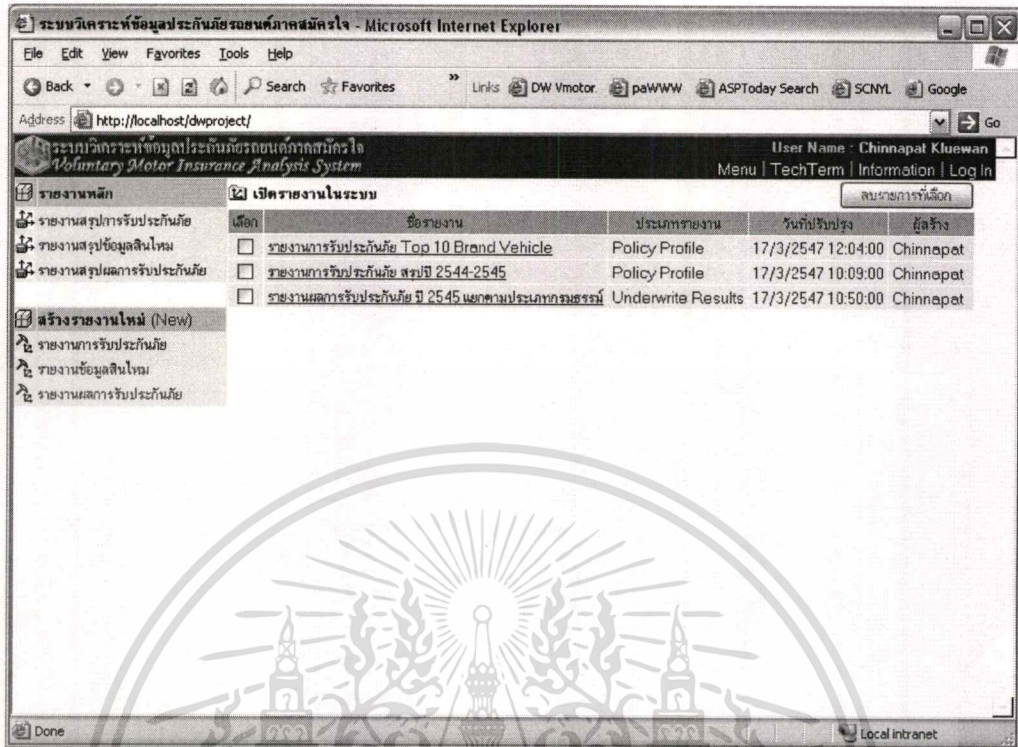
ขั้นตอนที่ 2: หน้าจอแสดงเมนูรายงานที่เป็นลักษณะรายงานมาตรฐาน และรายงานของผู้ใช้งานระบบ (รูปที่ 2-24)

ขั้นตอนที่ 3: เลือกรายงานที่ต้องการ เพื่อทำการออกแบบและทำการบันทึกการ

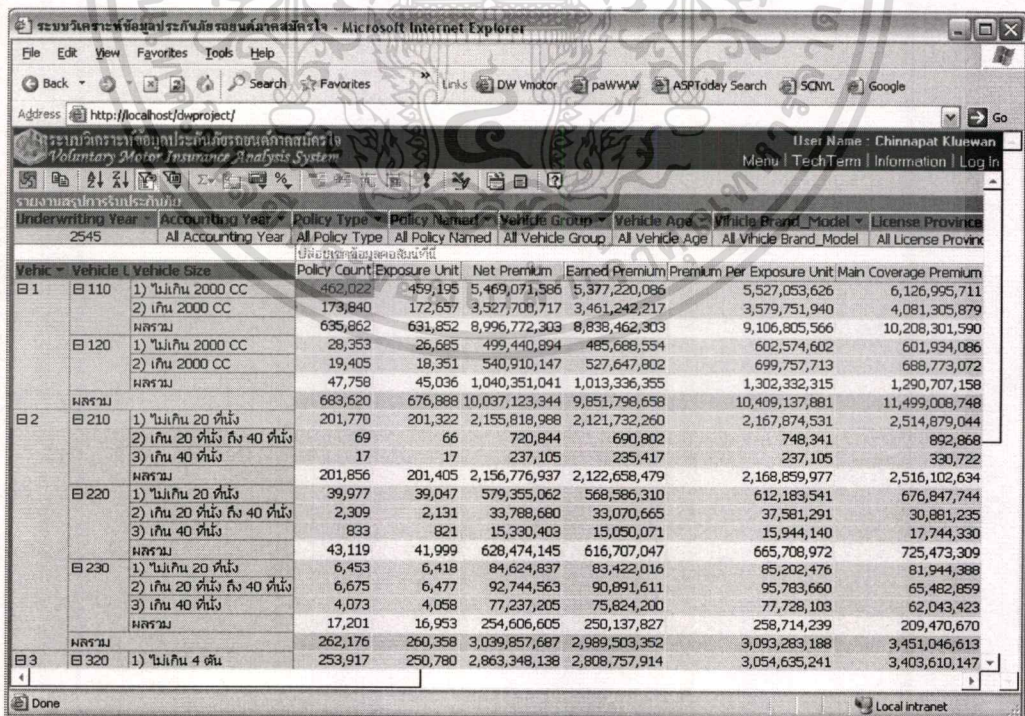
- เลือกรายงานใหม่ หรือเลือกรายงานมาตรฐาน
- ออกแบบหรือปรับแก้ไขรายงานให้อยู่ในรูปแบบที่ต้องการ (รูปที่ 2-25)
- ทำการบันทึกรายงาน โดยตั้งชื่อรายงานใหม่ และกำหนดสิทธิ์ในการเข้าถึงข้อมูล (รูปที่ 2-26)



รูปที่ 4-23 แสดงหน้าจอการล็อกอินเข้าใช้งานระบบ



รูปที่ 4-24 แสดงหน้าจอเมนูหลักรายงาน



รูปที่ 4-25 แสดงรายงาน “สรุปผลการรับประกันภัย” ที่ออกแบบ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Underwriting Year	Accounting Year	Policy Type	Policy Named	Vehicle Group	Vehicle Age	Vehicle Brand Model	License Province	Policy Count	Exposure Unit	Net Premium	Earned Premium	Premium Per Exposure Unit	Main Coverage Premium	OD Policy Count	F&T Policy
2545	All Accounting Year	All Policy Type	All Policy Named	All Vehicle Group	All Vehicle Age	All Vehicle Brand Model	All License Province								
110	120	รวม						635,862	631,852	8,996,772,303	8,838,462,303	9,106,805,566	10,208,301,590	449,285	41,261
2	3	4	5	6	7	8	รวมทั้งหมด	262,176	260,358	3,039,857,687	2,989,503,352	3,093,283,188	3,451,046,613	161,899	1,013,336,355
								388,898	383,535	4,731,714,463	4,639,220,171	5,014,215,184	5,505,996,294	244,083	2,989,503,352
								13,112	12,771	217,122,249	212,829,050	228,499,992	230,223,860	5,878	212,829,050
								24,469	23,978	146,399,338	143,517,330	150,224,674	187,533,933	8,977	143,517,330
								38,700	37,026	70,980,049	69,326,671	80,279,356	64,471,014	9,917	69,326,671
								43,896	43,875	252,553,523	249,272,835	252,717,260	250,525,053	873	249,272,835
								2,662	2,612	23,754,016	23,073,308	24,432,032	29,550,319	713	23,073,308
								1,457,533	1,441,043	18,519,504,669	18,178,541,375	19,252,789,567	21,218,355,634	922,886	1,013,336,355

รูปที่ 4-26 แสดงการบันทึกรายงานที่ออกแบบไว้

นอกจากส่วนของระบบที่เป็นรายงานแล้ว จำเป็นจะต้องมีหมายเหตุและคำอธิบายของข้อมูล เพื่อให้ผู้ใช้งานสามารถทราบที่มาของข้อมูล และคำอธิบายสำหรับตัวข้อมูล ดังแสดงรูปที่ 4-27 อธิบายที่มาของข้อมูลและ 4-28 อธิบายคำศัพท์ธุรกิจ

คำอธิบายในการใช้ข้อมูล
ระบบวิเคราะห์ข้อมูลประกันภัยรถยนต์ภาคสมัครใจ

1. เบี้ยที่ถือเป็นรายได้คำนวณถึง 31 ธันวาคม 2546
2. ข้อมูลที่นำมาแสดงไม่ใช่ผลรวมของทุกบริษัทประกันภัยที่ได้รับประกันภัยรถยนต์ภาคสมัครใจ แต่เป็นข้อมูลเฉพาะบริษัทที่ผ่านเกณฑ์ (กรมธรรม์ผ่านการตรวจสอบ >= 80%, สินไหมผ่านการตรวจสอบ >= 70% และรวมเบี้ยประกันภัยซึ่งต่ำกว่าเกณฑ์ 70% ของตลาด)
ปรากฏว่าการทำรายงานนี้ใช้ข้อมูลจาก 36 ใน 62 บริษัทที่ทำประกันภัยรถยนต์ภาคสมัครใจ และ 36 บริษัทเหล่านี้มีเบี้ยประกันภัยรวมกันเป็น 69.26 % ของธุรกิจประกันภัยรถยนต์ภาคสมัครใจทั้งหมด และกรมธรรม์ผ่านการตรวจสอบโดยเฉลี่ย 99.13 % สินไหมผ่านการตรวจสอบโดยเฉลี่ย 95 %
3. Adjusted Loss Ratio ได้จากการปรับฐานของทั้งกรมธรรม์และสินไหมให้เป็น 100% โดยใช้ข้อมูลอัตราผ่านการตรวจสอบกรมธรรม์ และสินไหมในของ 2.
4. ผลการรับประกันปี 2545 กรมธรรม์สิ้นสุดแล้ว
5. การพิจารณาผลการรับประกันภัยที่แท้จริง ควรจะต้องพิจารณา Combined Ratio และรายงานผลการรับประกันภัยตามปีบัญชีประกอบด้วย

แสดงทุกครั้งเมื่อเข้าระบบ ปิดหน้าต่าง

รูปที่ 4-27 แสดงหน้าจอข้อมูลที่ใช้อธิบายที่มาของข้อมูล

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Microsoft Internet Explorer

Glossary: ABCDEFGHIJKLMNOPQRSTUVWXYZ ปิดหน้าต่าง

B	
Body Injury	ความคุ้มครองต่อร่างกาย (BI)
C	
Coverage Premium	เบี้ยความคุ้มครอง
E	
Exposure Unit	หน่วยความเสี่ยงภัย ได้แก่ ความเสี่ยงต่อความเสียหาย, ระดับความเสี่ยงที่วัดได้, ความเสี่ยงต่อทรัพย์สินที่เอาประกันภัยต่อทรัพย์สินอันมีสาเหตุจากสิ่งแวดล้อม
Extended 01	ความคุ้มครองเพิ่มเติม ร.ย. 01
Extended 02	ความคุ้มครองเพิ่มเติม ร.ย. 02
Extended 03	ความคุ้มครองเพิ่มเติม ร.ย. 03
Extended 04	ความคุ้มครองเพิ่มเติม ร.ย. 04
Extended Territory	ความคุ้มครองเพิ่มเติมตาม ร.ย. 04
F	
Fire & Theft	ความคุ้มครองรถยนต์จากการเกิดภัยไฟไหม้ และขโมย
Fleet	การทำประกันภัยกลุ่ม
Fleet Bonus	โบนัสจากการทำประกันภัยกลุ่ม

รูปที่ 4-28 แสดงหน้าจอคำศัพท์ธุรกิจสำหรับระบบประกันวินาศภัย

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 5

การทดสอบระบบและสรุปผล

5.1 การทดสอบระบบ

หลังจากการพัฒนาระบบคลังข้อมูลการประกันวินาศภัย ตามขั้นตอนการทำงานในบทที่ 4 แล้ว สามารถทำการศึกษาและทดสอบระบบ เพื่อวัดประสิทธิภาพและความถูกต้อง โดยมีรายละเอียดในการทดสอบระบบดังนี้

ที่มาของข้อมูล : ฐานข้อมูลสำหรับดาต้าแวร์เฮาส์ ซึ่งเกิดจากการเตรียมข้อมูลมาจาก ฐานข้อมูลที่เป็นทรานแซคชัน

CUBE : POLICY_PROFILE

จำนวนข้อมูลการรับประกันภัย :	6,099,815	รายการ
การกำหนดค่าของการออกแบบพื้นที่จัดเก็บข้อมูล		
ประสิทธิภาพในการทำงาน	60	%
ความจุของดิสก์ที่ใช้	198	MB
ผลการประมวลผลข้อมูล		
จำนวน Segment ในการแบ่งประมวลผล	14	Segments
เวลาที่ใช้ในการประมวลผล	16 นาที	12 วินาที

CUBE : CLAIM_PROFILE

จำนวนข้อมูลสินไหมค่าทดแทน :	7,074,218	รายการ
การกำหนดค่าของการออกแบบพื้นที่จัดเก็บข้อมูล		
ประสิทธิภาพในการทำงาน	60	%
ความจุของดิสก์ที่ใช้	220	MB
ผลการประมวลผลข้อมูล		
จำนวน Segment ในการแบ่งประมวลผล	25	Segments
เวลาที่ใช้ในการประมวลผล	19 นาที	18 วินาที

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เปรียบเทียบความเร็วในการการประมวลผลจาก RDBMS กับ OLAP

คำสั่งในการทดสอบเพื่อเปรียบเทียบ

```
-- // หาผลรวมของเบี้ยประกันภัยรวม //--
SELECT SUM(MainCoverPrem) AS TotalAllPremium
FROM POLICY_PROFILE
```

การประมวลผลจาก RDBMS

เวลาที่ใช้ในการประมวลผล	1 นาที 25 วินาที
จำนวน CPU ที่ใช้	14,912 CPU ms.

การประมวลผลจาก OLAP

อ้างถึงเวลาในการประมวลผลข้อมูลรูป “POLICY_PROFILE” ซึ่งการทำงานจะเป็นลักษณะ Preprocessing โดยจะทำการคำนวณ Aggregates ต่างๆ ซึ่งซับซ้อนมากกว่าคำสั่งในการหาค่าเฉพาะผลรวมเบี้ยประกันภัย ตามตัวอย่างด้านบน ซึ่งใช้เวลา 16 นาที 12 วินาที ปรากฏว่า เวลาในการคำนวณค่าข้อมูลสามารถแสดงได้ทันที ซึ่งอาจจะใช้เวลาเพียงเล็กน้อยในการหาผลรวมเฉพาะค่าที่อยู่ใน Dimensions ที่กำหนดเงื่อนไขที่กรองค่า (เวลาที่ใช้เพียง 1-2 วินาที)

หมายเหตุ : ในการทดสอบวัดประสิทธิภาพการประมวลผลได้ใช้ SQL Profiler เป็น Tool ในการทดสอบการประมวลผลคำสั่ง SQL

5.2 การสรุปผล

จากผลการทดสอบการประมวลผลข้อมูลประกันวินาศภัยรถยนต์ภาคสมัครใจ พบว่า ประสิทธิภาพในการประมวลผลและความเร็วในการทำงานของรูปขึ้นอยู่กับกรอบพื้นที่จัดเก็บข้อมูลลงดิสก์ ซึ่งเป็นปัจจัยสำคัญ อีกทั้งในเรื่องของความถูกต้องของข้อมูลในการนำเสนอเป็นเป้าหมายหลักของการพัฒนาระบบดาต้าแวร์เฮาส์ ดังนั้นการลดความข้อผิดพลาดในการประมวลผลจะต้องกำหนดกระบวนการคลีนนิ่งข้อมูลให้เหมาะสม เพื่อป้องกันการสูญหายของข้อมูลดังกล่าวด้วย และผลที่ได้จากการพัฒนาระบบจะเป็นแนวทางในการพัฒนาระบบดาต้าแวร์เฮาส์สำหรับระบบงานประกันวินาศภัยด้านอื่นๆ ให้มีประสิทธิภาพต่อไป

บรรณานุกรม

Rakesh Agrawal. 2002. **Modeling Multidimensional Databases** [Online]. Available:

http://infolab.usc.edu/csci599/Fall2002/paper/I3_agrawal95modeling.pdf

Robert Fenk. 1999. **MDX – Multidimensional Expressions**. [Online]. Available:

<http://www3.informatik.tu-muenchen.de/lehre/SS2001/DAWA-bayer/DWH-Ch10-1.pdf>

David Willson. 2001. **The Microsoft Data Warehousing Framework**. [Online]. Available:

<http://msdn.microsoft.com/library/default.asp?url=/library/en-us/dnolap/html/intromdx.asp>

Carl Nolan. 1999. **Introduction to Multidimensional Expressions (MDX)**. [Online]. Available:

<http://msdn.microsoft.com/library/default.asp?url=/library/en-us/dnolap/html/intromdx.asp>

ประวัติผู้เขียน

ชื่อผู้เขียน นายชินพัฒน์ เคลือวัลย์ นักศึกษาปริญญาโท สาขาวิชาเทคโนโลยีสารสนเทศ คณะเทคโนโลยีสารสนเทศ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

วุฒิการศึกษาระดับปริญญาตรีจาก มหาวิทยาลัยมหิดล คณะวิทยาศาสตร์ ภาควิชาวิทยาการคอมพิวเตอร์ ปี พ.ศ.2545 และมหาวิทยาลัยเกษมบัณฑิต คณะบริหารธุรกิจ ภาควิชาคอมพิวเตอร์ ธุรกิจ ปี พ.ศ. 2539

ปัจจุบันมีตำแหน่งเป็นนักวิเคราะห์ระบบ บริษัท ไทยอินชัวร์เรอส์ คาด้าเน็ต จำกัด ซึ่ง ออกแบบและพัฒนาโครงการประกันภัยเชื้ออาหาร และระบบประมวลผลข้อมูลประกันภัยรถยนต์ ภาคสมัครใจ



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้