

Solution to COMP5121

Assignment 3

Group Questions

11500811G QING Pei

11500478G Fang Nan

2011/12/4

Table Of Contents

| | |
|---|----|
| Table Of Contents..... | 2 |
| Question 1 | 3 |
| Understanding of the situation | 3 |
| Creating a target data set | 3 |
| Data cleaning and preprocessing | 3 |
| Data reduction and projection | 4 |
| Data mining..... | 6 |
| Revenue analysis | 7 |
| Profit analysis | 12 |
| Marketing analysis..... | 13 |
| Consolidating discovered knowledge..... | 14 |
| Question 2 | 18 |
| Data preparation | 18 |
| a) | 19 |
| b)..... | 19 |
| c) | 30 |
| d)..... | 31 |

Question 1

We follow the general KDD process to find knowledge in the given data. This report is organized according to the purpose of each step.

Understanding of the situation

The goal is to help improve its competitiveness. Competitiveness may be:

- Better service quality (call quality)
- Larger revenue
- More profit
- etc.

For any improvement, current state should be analyzed and described. Based on the knowledge found, strategies can be carried out to make a change.

Creating a target data set

Data set is given as 4 tables: PhoneCallDetail, CDemographics, PaymentRecord and MarketingData. Without investigation, all columns these 4 tables appear to be related to competitiveness in that they are either about customers who judge the service quality or payments and offerings which define the financial state of the company.

Although all the columns are retained, many of them cannot be directly used for analysis. More data processing and transformation need to be done to the raw data.

Data cleaning and preprocessing

PhoneCallDetails.csv

1. A number of call destination numbers are NULL. To preserve the type of this field (integer), these cells are filled with a 0.
2. Some regions of ending the phone call are missing. They are filled with "Unknown".
3. Two pairs of inconsistent notations for disconnection exist. "Y/N" and "Yes/No". The values are reclassified to "Y/N" only.

CDemographics.csv

1. 5 customers' sex is missing. The information cannot be guessed from any existing data. The action taken is to use the last valid value in the table to fill the blanks.
2. 4 customers' birthdays and 7's joining dates are missing. These are also filled using the last non-blank values.
3. Other than single character representations, "Sec" also exists in education level column. It is reclassified to "S".
4. Some types occur in JoinFrom column and are reclassified to the correct spellings.
5. Phone numbers of customers are clearly separated into two ranges, one starting with 98 and another starting with 99. An according binning action is applied.

MarketingData.csv

1. No explicit data error or missing value exists. However, only 234 customers have ever received a special offer. Another 22 customers, who exist in other tables, do not exist in this table. There will be some missing values after joining the table.
2. The blank cells in marketing data are filled with 0 for number of offers and "No" for whether a specific offer has ever been received by a customer.

PaymentRecord.csv

1. Missing payment dates and due dates are filled with last non-blank dates. If the generation time of payment records are related to payment date or due date, this may work better. In that way, payments with similar PID tend to have similar payment/due dates. In this table, PIDs and both dates have no correlation and the interpolation function can be randomly selecting a date in the data range.
2. Other than digital representation of payment method, textual ones also exist. "cred", "PPS", "Cash" and "onli" are reclassified to the according digital forms of representation.

Data reduction and projection

General

1. Duplicate records still exist after processing each single table. To minimize the table size, these records are deduplicated with a Distinct

operator. Meanwhile, a record count column is added to keep the frequency information.

2. For columns with very uneven distribution, balanced nodes are generated to assist feature selection.
3. If some columns are never used in later process, a filter operator is used to avoid these data from being processed by following operators or models.

PhoneCallDetails.csv

1. Time of making the call contains both date and time. This is transformed into 4 different 6-hour time slots in 24 hours.
2. Call duration is derived from call start/end times and then binned to three groups, short/normal/long calls, according to standard deviation by 1 sigma.
3. Each customer's total/average call duration is aggregated from the transactional data.

CDemographics.csv

1. Age is derived from birthday. As the date of data collection is not given, the age is calculated according to today's date.
2. Customers are then binned into two groups of the same size, young and old. The age distribution does not resemble the real-world one. It is not practical to guess the mapping from the data to customers' real ages. Then we discard the values, leave the order information alone.
3. With an evident peak at the data audit results, we decide to reduce the number of JoinFrom companies. Apple and FoolishTone are grouped as "AnF". The rest are grouped as "Other". The motivation is the difference in payment statistics between the two groups of customers.
4. Similar to the previous one, education levels are reduced to two kinds, "TD" for tertiary or degree, "Other" for all the rest.
5. Date of joining ABC Telecom is transformed into months of service. Then divided to two equal size groups, "long" and "short".

MarketingData.csv

1. The original data is transactional. Each transaction is restructured to be flags of whether special offer A/B/C is offered. This gives a Boolean value.
2. The restructured transactions are then aggregated to show the number of special offer A/B/C each customer has ever been offered. A sum of all types of offers is also aggregated. This gives a quantitative measurement of special offers.

PaymentRecord.csv

1. A new column UnpaidPreviousCharge is derived as “Balance-Charge”. This is related to financial operations of the company.
2. A new column DueTime is derived as the difference in dates between payment date and due date. Early or late payment patterns may affect the operation of ABC Telecom.
3. Charge, UnpaidPreviousCharge and DueTime are aggregated to provide total/average statistics for each customer.
4. The aggregation results of last step are all binned to three groups according to standard deviation by 1 sigma.

Merged table

1. CDemographics and PhoneCallDetails are merged by PhoneNo.
2. The merged result of last step is merged with PaymentRecord and MarketingData by CustomerID.

Profit measurement

DEFINITION: THE PROFIT INDEX OF A CUSTOMER IS DEFINED AS THE TOTAL AMOUNT OF MONEY PAID BY THAT CUSTOMER DIVIDED BY THE TOTAL AMOUNT OF RESOURCE HE OR SHE USED.

$$Profit = \frac{\Sigma(charge)}{\Sigma(call\ duration)}$$

CUSTOMERS WHO PAY MORE BUT MAKE LESS PHONE CALLS ARE MORE PROFITABLE. EFFORT IN IDENTIFYING THIS KIND OF CUSTOMERS IS REWARDING.

Data mining

Service quality analysis

The goal of this analysis is to know the existing call quality.

According to the data, 41.05% calls are with disconnection. This is a high rate for local calls. Two classification models are used to illustrate the call quality issues:

1. C5.0, input={Cell, Cell2, SameCell(whether Cell=Cell2)} output={Disconnect}
2. C5.0, input={PhoneNo, Destination} output={Disconnect}

The result of the 1st model tells that calls ending in unknown cells are 100% with disconnection.

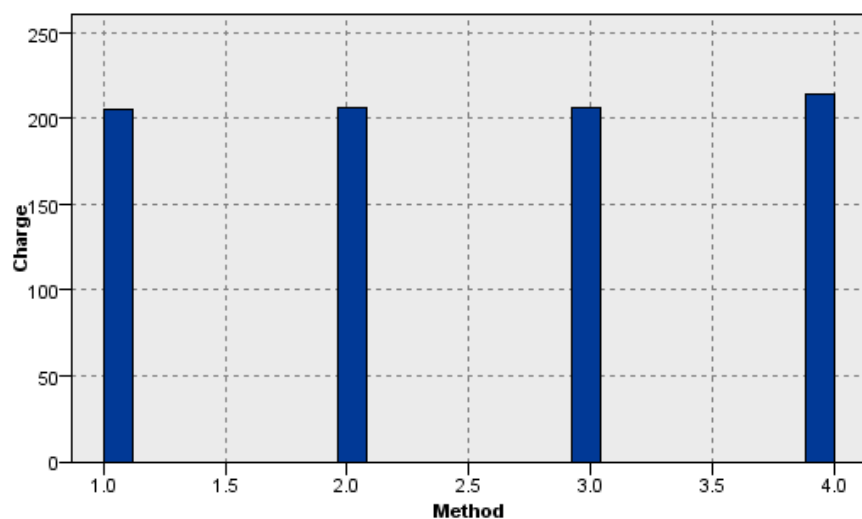
The result of the 2nd model tells that calls to unknown destination number are 100% with disconnection. Calls made from 98***** are less likely to have disconnection than from 99***** (28% versus 48%). Call quality is extremely bad when the call is from 99***** to 91*****.

The company should improve signal coverage to more regions. This will reduce the chance of customers ending the call in an unknown cell.

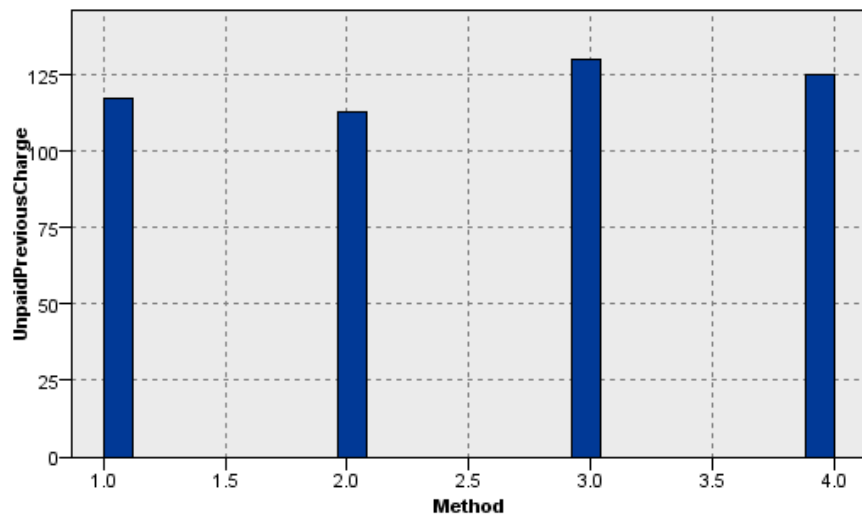
The difference in phone number ranges should also be paid attention to. The problem may lie in the switching hardware/software. Further tests or data collection is necessary to find out the cause of different call quality for different phone number ranges.

Revenue analysis

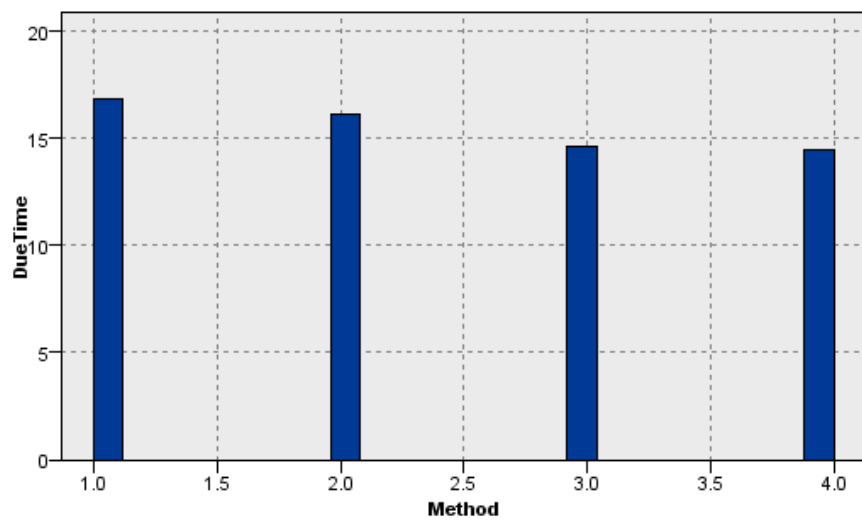
First of all, we examined the current revenue characteristics.



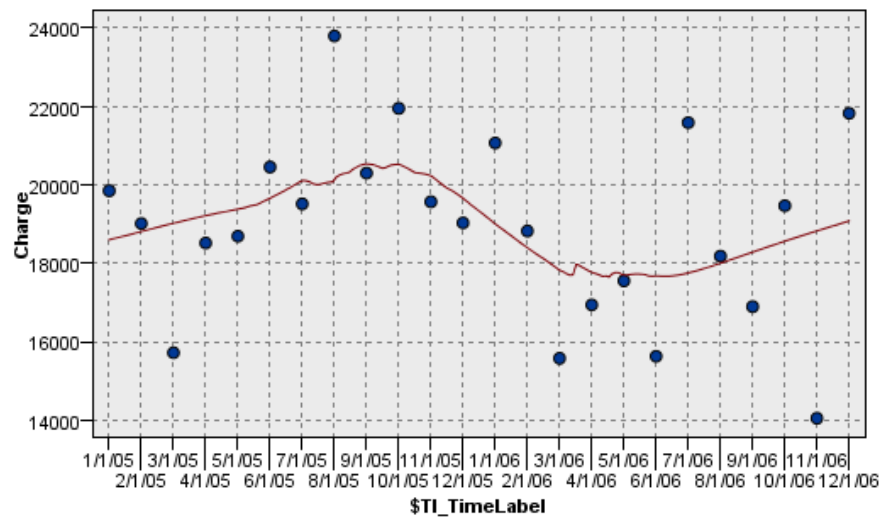
Average charge of different payment methods are almost the same. Online payments are slightly higher.



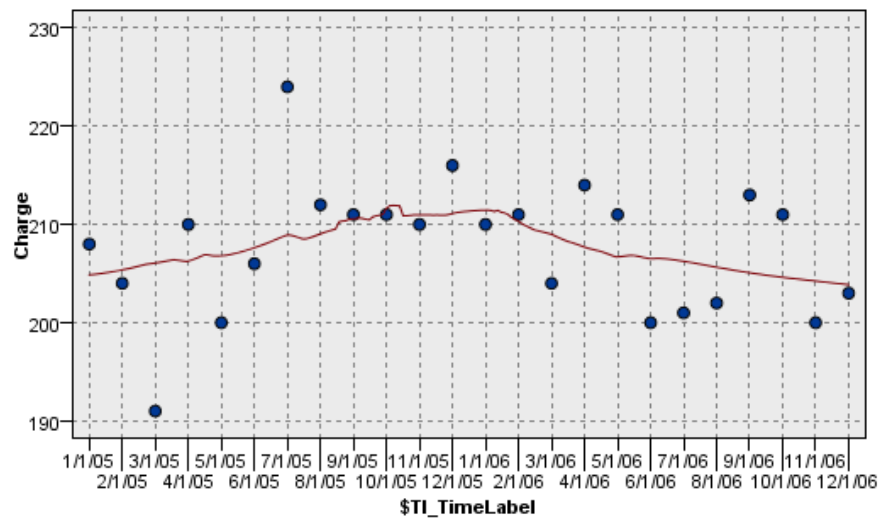
Unpaid previous charge (balance – charge) is different among payment methods. Customers who pay by cash tend to pay less frequently than PPS users.



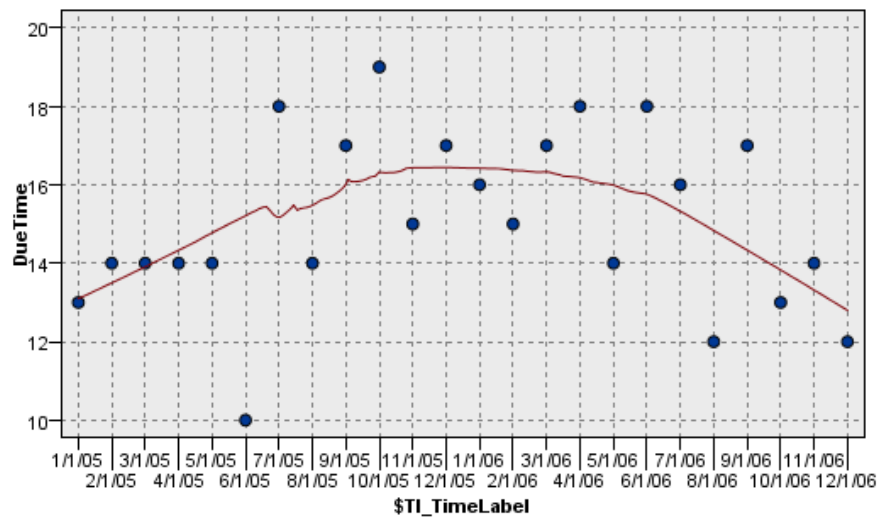
Average due time in days is also different among payment methods. Credit card or PPS users prefer late payments.



Monthly revenue increased in the first three quarters and then dropped for two quarters. In the last two quarters, the figures did not have a stable trend, ranging from as low as 14000 to as high as 22000.



Average charge of customers is decreasing slightly over time.



Customers nowadays finish their payments more timely. However, they still pay the bills a dozen days later than the due date.

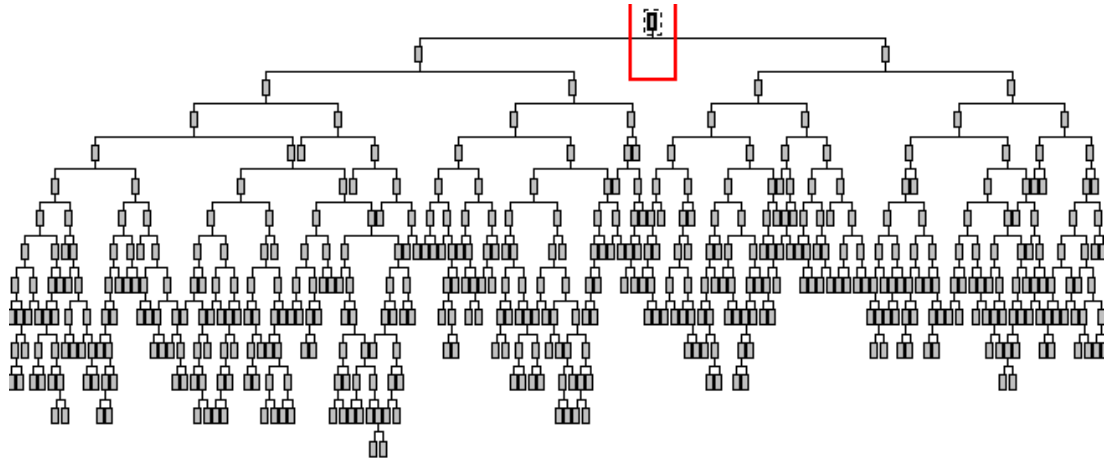
After discretizing the Charge_Sum, Charge_Mean, UnpaidPreviousCharge_Mean and DueTime_Mean, C5.0 models are used to classify customers into different types.

Candidate inputs are all the variables in control of ABC Telecom: customer demographics and special offers. The company can choose customers based on their demographic information or the agents may decide whether or not to give a particular customer a special offer.

The 4 models provide decision trees with

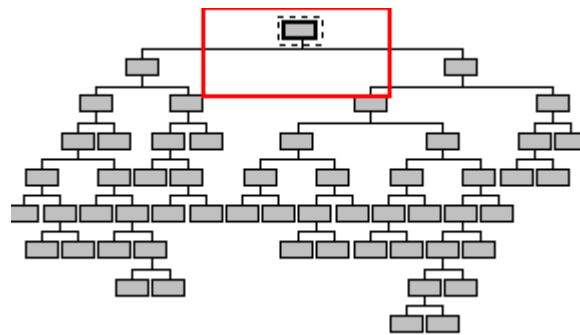
- Approximately 96% accuracy
- Tree depth of 14~17
- 13~21 seconds of building time

Take the decision tree of UnpaidPreviousCharge_Mean for example; a tree like this is given:

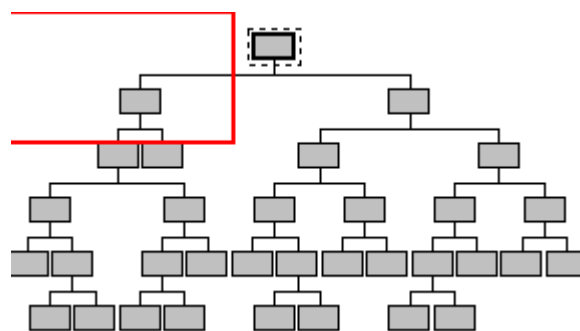


The model can be used by computers. For manual verification, however, we need to simplify the tree.

By limiting the smallest record count of a node to be 800, we get a simpler tree with a depth of 10.



By further limiting the smallest record count of a node to be 2000 and increase pruning severity from 75 to 80, we get a tree with a depth of 7.



Here is the comparison of the three trees:

| Modeling Mode | Simple | Expert | Expert |
|------------------|--------|--------|--------|
| Pruning severity | / | 75 | 80 |

| Minimum records per child branch | / | 800 | 2000 |
|---|----------|------------|-------------|
| Tree depth | 14 | 10 | 7 |
| Accuracy | 96.674% | 88.322% | 81.219% |
| Building time | 19s | 17s | 13s |

The tradeoff is between simplicity of the model and the accuracy.

For computer aided processing, the first model appears to be the best. With a little more building time, it offers much better accuracy.

For manual processing, the last model's simplicity is preferred.

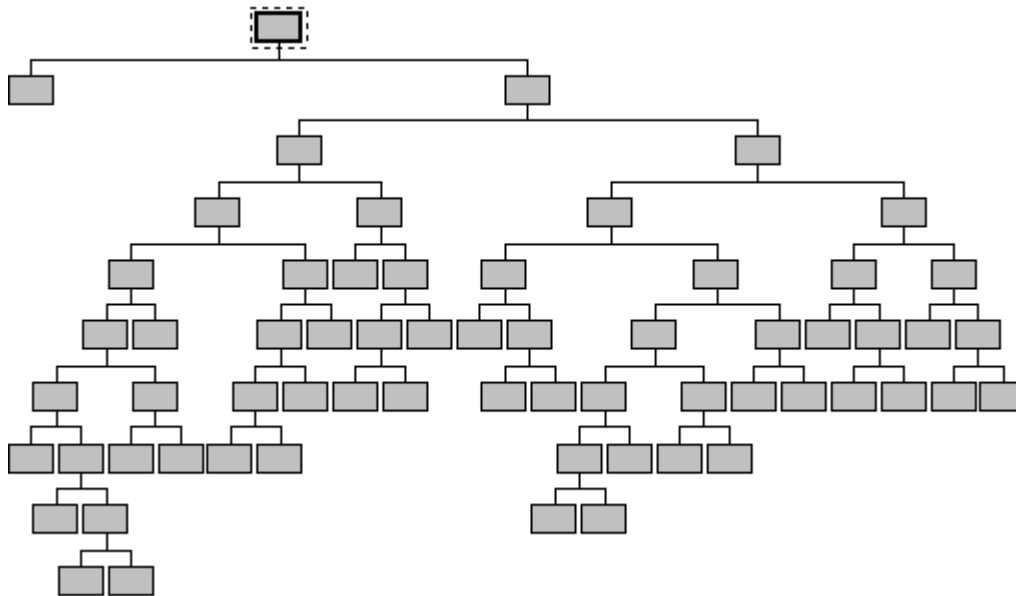
Similar tradeoffs apply to all the other decision tree building options.

Profit analysis

Revenue is not the only pursuit of the company. What's more important is the profit rate. With our definition of profit, we bin customers into two groups: more profitable and less profitable ones.

We use C5.0 again to try to classify customers into these two groups. If a model with acceptable accuracy is available, ABC Telecom can identify those customers who are potentially more profitable and use some strategies such as special offers to attract them.

By limiting the minimum records per child branch to be 600, we are able to get a decision tree of profit index as below.



The accuracy of the 10-level tree is 84%.

Marketing analysis

The following table shows the effect of special offer A:

| Field | Mean (N) | Mean (Y) | Importance |
|----------------------------|----------|----------|------------|
| CallDurationInMinutes_Sum | 2648.606 | 2340.59 | 0.997 |
| CallDurationInMinutes_Mean | 83.136 | 82.73 | 0.394 |
| Charge_Sum | 495.508 | 497.055 | 0.119 |
| Charge_Mean | 207.396 | 211.585 | 0.931 |
| UnpaidPreviousCharge_Mean | 115.406 | 130.823 | 0.999 |
| DueTime_Mean | 15.113 | 15.506 | 0.556 |

Customers with special offer A, compared to those without:

- ☹ Make 308 minutes *fewer* phone calls
- ☹ Have \$15 more unpaid previous charge
- ☺ Pay \$4 more each month

At first glance, the fact that customers make fewer phone calls is bad news. However, they pay even more for fewer phone calls, which turns this to be rather good news.

The following table shows the effect of special offer B:

| Field | Mean (N) | Mean (Y) | Importance |
|----------------------------|----------|----------|------------|
| CallDurationInMinutes_Sum | 2286.102 | 2549.131 | 0.99 |
| CallDurationInMinutes_Mean | 83.138 | 82.956 | 0.182 |
| Charge_Sum | 474.462 | 509.872 | 0.999 |
| Charge_Mean | 207.018 | 211.871 | 0.966 |
| UnpaidPreviousCharge_Mean | 130.382 | 121.924 | 0.941 |
| DueTime_Mean | 16.247 | 14.822 | 0.995 |

Customers with special offer B, compared to those without:

- ☺ Make 263 minutes more calls
- ☺ Pay \$5 more each month
- ☺ Pay \$35 more in total
- ☺ Have \$8 less unpaid previous charge
- ☺ Pay bills 1.4 days earlier

The following table shows the effect of special offer C:

| Field | Mean (N) | Mean (Y) | Importance |
|----------------------------|----------|----------|------------|
| CallDurationInMinutes_Sum | 2512.324 | 2418.658 | 0.653 |
| CallDurationInMinutes_Mean | 83.102 | 83.352 | 0.255 |
| Charge_Sum | 519.055 | 477.607 | 1 |
| Charge_Mean | 215.208 | 205.711 | 1 |
| UnpaidPreviousCharge_Mean | 129.089 | 121.805 | 0.906 |
| DueTime_Mean | 15.115 | 15.567 | 0.635 |

Customers with special offer C, compared to those without:

- ☺ Have \$7 less unpaid previous charge
- ☹ Pay \$10 *less* each month
- ☹ Pay \$42 *less* in total

Consolidating discovered knowledge

With all the data mining results, we can summarize the knowledge found as the following:

Service quality:

The company should improve signal coverage to more regions. This will reduce the chance of customers ending the call in an unknown cell.

The difference in phone number ranges should also be paid attention to. The problem may lie in the switching hardware/software. Further tests or data

collection is necessary to find out the cause of different call quality for different phone number ranges.

Revenue

From the current revenue characteristics, we suggest:

- ① As for the online payment method is slightly efficient than others and from the perspective of unpaid previous charge, we recommend the company to build a better website and construct a simple but secure channel for customer to pay. This measurement may not only improve the efficiency in paying rate ,but also helps company to improve the rate of cost recovery
- ② From the classification result of unpaid pervious charge mean by using C5.0, we can conclude that the type of special offer and the age distribution affect most. So, if the company wants to decrease the unpaid pervious charge mean, it should take more relevant measures from the two parts.
 - Suggestion 1. The company would provide discount purchase of phone to young male customers. This method will encourage this group of customers to decrease unpaid charge to some extent.
 - Suggestion 2. The company could provide the group of customers whose educational level is Tertiary or Degree and also using phone number staring with 98 more priorities in payment method especially through online and cash method
 - Suggestion 3. Don't give discount purchase of phone to the group of customer whose education level is not Tertiary or Degree and always pay the bill with cash or online method. For this social group, we recommend the company consider the other two special offers to protect them from delaying payment.
 - Suggestion 4. The company would give free MTR, Tunnel & License fee and discount purchase of phone to the old male customers with phone no. starting in 99.
- ③ Take the classification result of charge sum/mean using C5.0, we found that the attribute of join from , the type of special offer and the group phone no. are the top 3 critical factors .
 - Suggestion 1. Customers who join from the company of AnF tend to pay more. So we want to benefit this condition from two aspects: first, we need to attract more customers from the above company to join us such as provides more priorities than AnF to the customers through advertisements or other media to propaganda our company. Second, we need to learn some operation strategies from these companies. For example, for some support parts of core transaction we could consider outsourcing strategy to drop fixed cost etc.

---Suggestion 2. We should recommend the company do not provide more than one special offer to customers for receiving high charge sum.

- ④ From the perspective of the result of due time mean, we recommend that decrease the amount of special offer and female customers tend to pay in time.

Profit analysis

As for we defined the profit index of a customer as the total amount of money paid by that customer divided by the total amount of resource he or she used.

$$Profit = \frac{\Sigma(charge)}{\Sigma(call\ duration)}$$

Customers who pay more but make less phone calls are more profitable. Effort in identifying this kind of customers is rewarding.

It is evident that those customers with less charge mean must bring less profit to the company. So our target group should be focused on customers who have the higher charge mean.

---Target group 1. We could give more special offer to no due time female customers joining from AnF and encourage them to pay using credit card or PPS.

---Target group 2. Young customers joining from AnF

---Target group 3. We recommend give free MTR, Tunnel & License fee to customers joining from companies except AnF, but do not give them discount purchase of phone.

Marketing analysis

We want to analysis this field from macroscopic view and microscopic view respectively.

MACROSCOPIC VIEW

Given the advantages and disadvantages among three kind of special offers, the aim is to provide more customers a certain special offer with the least cost. How can we attain this goal? Seeing from the statistics shown before, we found that special offer B bringing the most benefits which promotes capital backflow and stimulates customers making more phone calls constructed a virtuous circle on revenue. Based on the above facts, we strongly recommend the company to provide special offer B to the majority of customer who owns the qualification of special offer. It is evident that when a customer is given the additional min each month, he or she always does not to make more phone calls than usual for the

additional minutes. At the same time, the statistics showed that they pay more each month resulting in less debt and pay earlier than usual winning time value of fund for the company.

MICROSCOPIC VIEW

Special offer A (free MTR, Tunnel & License fee)

| Field | Mean (N) | Mean (Y) | Importance |
|----------------------------|-------------|-------------|------------|
| CallDurationInMinutes_Sum | 2648.606 | 2340.59 | 0.997 |
| CallDurationInMinutes_Mean | 83.136 | 82.73 | 0.394 |
| Charge_Sum | 495.508 | 497.055 | 0.119 |
| Charge_Mean | 207.396 | 211.585 | 0.931 |
| UnpaidPreviousCharge_Mean | 115.406 | 130.823 | 0.999 |
| DueTime_Mean | 15.113 | 15.506 | 0.556 |

Special offer (Additional min)

| Field | Mean (N) | Mean (Y) | Importance |
|----------------------------|-------------|-------------|------------|
| CallDurationInMinutes_Sum | 2286.102 | 2549.131 | 0.99 |
| CallDurationInMinutes_Mean | 83.138 | 82.956 | 0.182 |
| Charge_Sum | 474.462 | 509.872 | 0.999 |
| Charge_Mean | 207.018 | 211.871 | 0.966 |
| UnpaidPreviousCharge_Mean | 130.382 | 121.924 | 0.941 |
| DueTime_Mean | 16.247 | 14.822 | 0.995 |

Special offer C (Discount purchase of phone)

| Field | Mean (N) | Mean (Y) | Importance |
|----------------------------|-------------|-------------|------------|
| CallDurationInMinutes_Sum | 2512.324 | 2418.658 | 0.653 |
| CallDurationInMinutes_Mean | 83.102 | 83.352 | 0.255 |
| Charge_Sum | 519.055 | 477.607 | 1 |
| Charge_Mean | 215.208 | 205.711 | 1 |
| UnpaidPreviousCharge_Mean | 129.089 | 121.805 | 0.906 |
| DueTime_Mean | 15.115 | 15.567 | 0.635 |

Compared the three kinds from the perspective of importance:

| Field | special offer A | special offer B | special offer C |
|----------------------------|-----------------|-----------------|-----------------|
| CallDurationInMinutes_Sum | ✓ | ✓ | |
| CallDurationInMinutes_Mean | | | |
| Charge_Sum | | ✓ | ✓ |
| Charge_Mean | | ✓ | ✓ |
| UnpaidPreviousCharge_Mean | ✓ | ✓ | ✓ |

We can conclude that UnpaidPreviousCharge Mean is one of the most important factors which affected largely among all the special offers followed by all the attributes about charge. CallDurationInMinutes_Mean and DueTime_Mean have nothing with special offer.

So, no matter which kind of special offer the company wants to provide to customers, the first factor to consider is the UnpaidPreviousCharge Mean. The index of CallDurationInMinutes_Mean affects the company to give special offer A or B to a certain group of customers. For the same reason, charge decide the provision of special offer B or C

Question 2

Data preparation

We extracted useful data from the given tables and build the following three new tables for our design.

| City_Date_Prod |
|----------------|
| ProductName |
| OrderPrice |
| OrderDate |
| City |

| Cat_Prod |
|--------------|
| CategoryName |
| ProductName |

| Countyr_City |
|--------------|
| City |
| Country |

The preparation process is done in PASW Modeler.

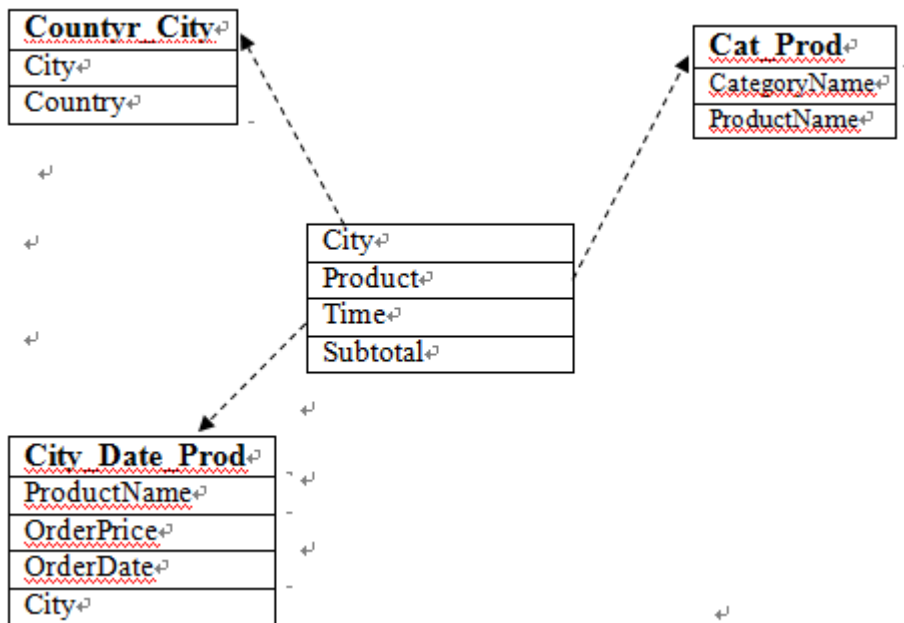
1. Join **Category** with **Products** by *CategoryID*.
2. Then join the result with **OrderDetails** by *ProductID*.
3. Join with **Orders** by *OrderID*.
4. Join with **Customers** by *CustomerID*.
5. Derive a new column called *OrderPrice* by formula:

$$\text{OrderPrice} = \text{UnitPrice} * \text{Quantity} * (1 - \text{Discount})$$

6. Use 3 filters to select the 3 groups of columns we need.
7. Use **distinct** operator to remove duplicate rows.
8. Export each result to a flat file.

a)

Star schema:



The above star schema consists of two types of tables: a fact table and three dimensional tables. We design the fact table contained three dimensional tables (city, product, time) and a measure (subtotal) . As for the requirements of a sales report, we need the following four aspects of information: city ,product ,time and sales. So we extract the information of city and country from the original table Customers to build a new table . For the same reason ,we build another three dimensional tables from Category , OrderDetails , Orders and Products.

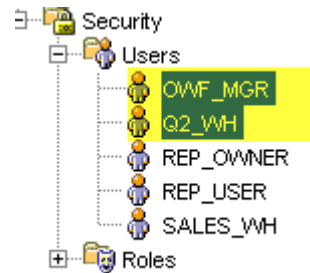
The primary benefit of this star schema is its simplicity for users to write, and databases to process: queries are written with simple inner joins between the facts and three dimensions. Star joins are simpler than possible in snowflake schema. Where conditions need only to filter on the attributes desired, and aggregations are fast. We users can see clearly from the hierarchical design and get the target information.

b)

To construct the data warehouse we designed, we follow these steps:

0. Prerequisites

- a) Install Workflow Server 2.6.4, and then add the “owf_mgr” to the current workspace.
- b) Create a new user “q2_wh”.



1. Flat files



- a) Create a source file location where all the flat files are stored.

A screenshot of a dialog box titled 'Edit File System Location: Q2_FILE_LOC'. The dialog has a blue title bar with a close button. It contains the following fields and controls:

- Name:** A text box containing 'Q2_FILE_LOC'.
- Description:** A large empty text box.
- Type:** A dropdown menu set to 'General'.
- Path:** A text box containing 'C:\Data' with a 'Browse...' button to its right.
- Buttons:** 'Help', 'OK', and 'Cancel' buttons at the bottom.

- b) Import flat files.
 - i. Specify date mask.

| Edit Flat File: CITY_DATE_PROD_SUBTOTAL_CSV | | | | | | | |
|--|-------------|-------------|--------|-----------|-------|------------|----|
| Name General Structure | | | | | | | |
| A length, precision or scale of 0 implies the use of SQL*Loader default value. | | | | | | | |
| Fields: | | | | | | | |
| | Name | Type | Length | Precision | Scale | Mask | NU |
| 1 | ProductName | CHAR | 255 | | | | |
| 2 | OrderPrice | DECIMAL ... | 0 | | | | |
| 3 | OrderDate | DATE | 0 | | | YYYY/MM/DD | |
| 4 | City | CHAR | 255 | | | | |
| | | | | | | | |

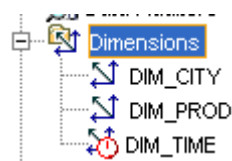
2. External tables



- New external table for each flat file.
- Specify the data file for each external table.

| Configuration Properties: DEFAULT_CONFIGURATION | |
|---|--|
| <div> <div>CAT_PROD_CSV</div> <div> <div>Data Files</div> <div>NEW_DATA_FILE_1</div> <div>External Table Columns</div> </div> </div> | |
| <div>NEW_DATA_FILE_1</div> <div> <div>Data File Name</div> <div>cat_prod.csv</div> </div> <div> <div>Data File Location</div> <div>Q2_FILE_LOC</div> </div> | |

3. Dimensions



- Create a new time dimension with wizard and use the calendar Year/Quarter/Month/Day hierarchy.

| DIM_TIME | | |
|--------------------|--|-----|
| Attributes | | |
| ID | | 78g |
| DAY | | 31 |
| CODE | | 78g |
| START_DATE | | 31 |
| END_DATE | | 31 |
| TIME_SPAN | | 78g |
| JULIAN_DATE | | 78g |
| DESCRIPTION | | abc |
| NAME | | abc |
| DAY_OF_CAL_WEEK | | 78g |
| DAY_OF_CAL_MONTH | | 78g |
| DAY_OF_CAL_QUARTER | | 78g |
| DAY_OF_CAL_YEAR | | 78g |
| CAL_MONTH_NUMBER | | 78g |
| MONTH_OF_QUARTER | | 78g |
| MONTH_OF_YEAR | | 78g |
| CAL_QUARTER_NUMBER | | 78g |
| QUARTER_OF_YEAR | | 78g |
| CAL_YEAR_NUMBER | | 78g |
| Levels | | |
| DAY | | |
| CALENDAR_MONTH | | |
| CALENDAR_QUARTER | | |
| CALENDAR_YEAR | | |
| Hierarchies | | |

b) Create a new city dimension to specify the hierarchy of country/city.

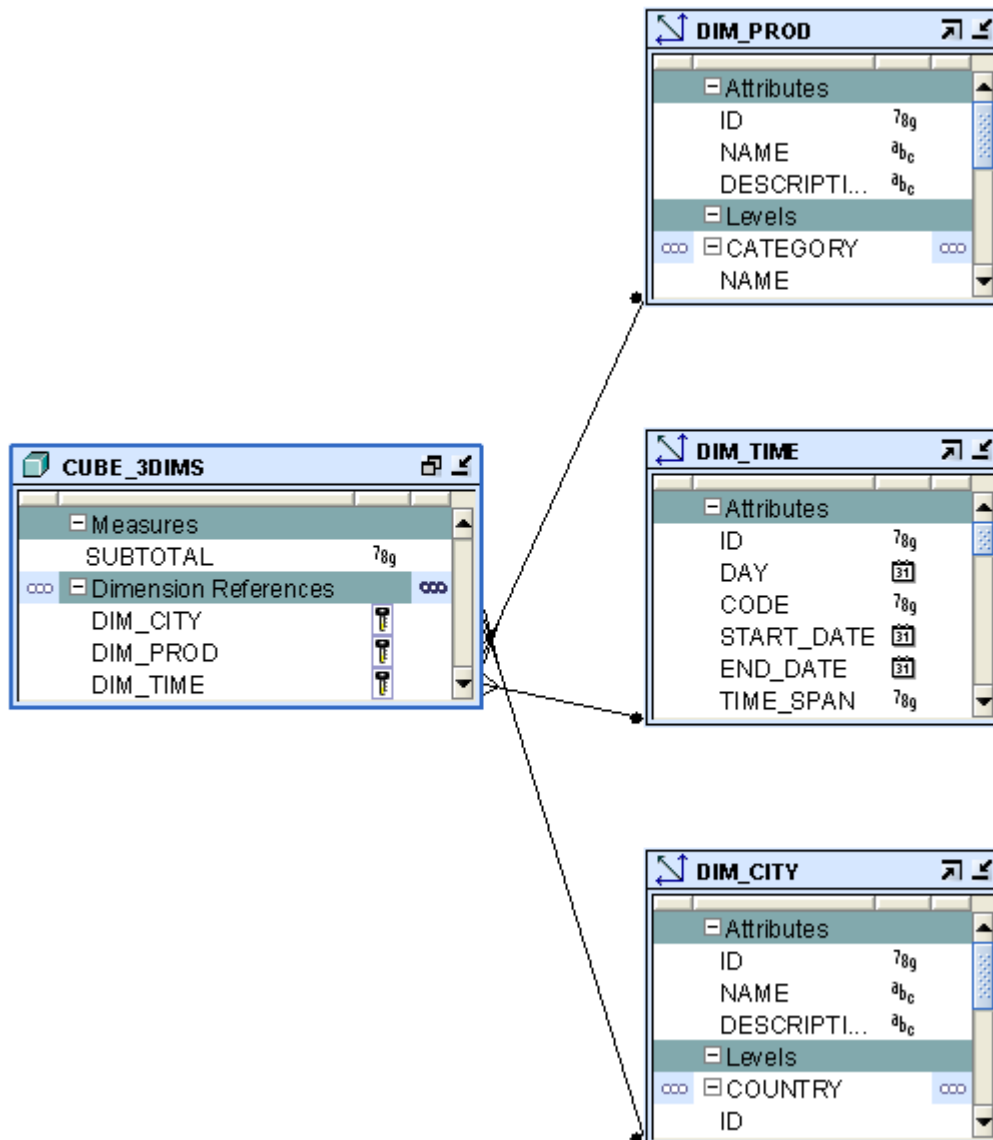
| DIM_CITY | | |
|-------------|--|-----|
| Attributes | | |
| ID | | 78g |
| NAME | | abc |
| DESCRIPTION | | abc |
| Levels | | |
| COUNTRY | | |
| ID | | |
| DESCRIPTION | | |
| NAME | | |
| CITY | | |
| NAME | | |
| DESCRIPTION | | |
| ID | | |
| Hierarchies | | |
| STANDARD | | |
| COUNTRY | | |
| CITY | | |

| DIM_PROD | | |
|-------------|--|-----|
| Attributes | | |
| ID | | 78g |
| NAME | | abc |
| DESCRIPTION | | abc |
| Levels | | |
| CATEGORY | | |
| NAME | | |
| ID | | |
| DESCRIPTION | | |
| PRODUCT | | |
| NAME | | |
| DESCRIPTION | | |
| ID | | |
| Hierarchies | | |
| STANDARD | | |
| CATEGORY | | |
| PRODUCT | | |

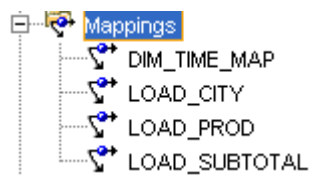
c) Create a new product dimension to specify the hierarchy of category/product.

4. Cube

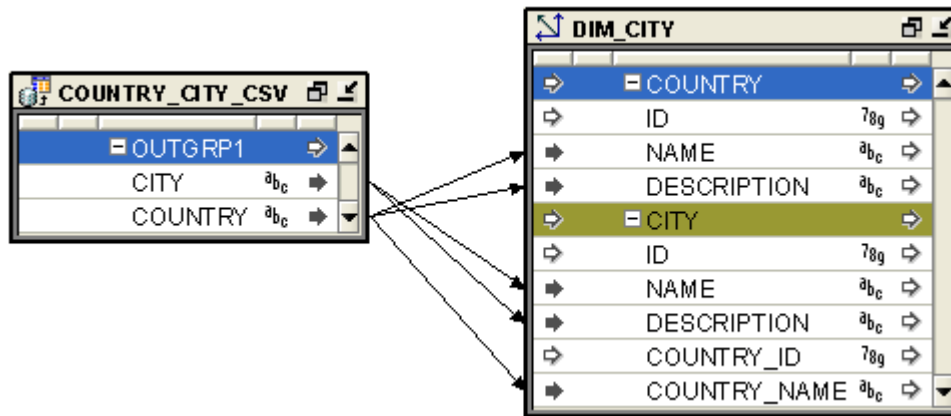
- a) Create a cube that has the above 3 dimensions.



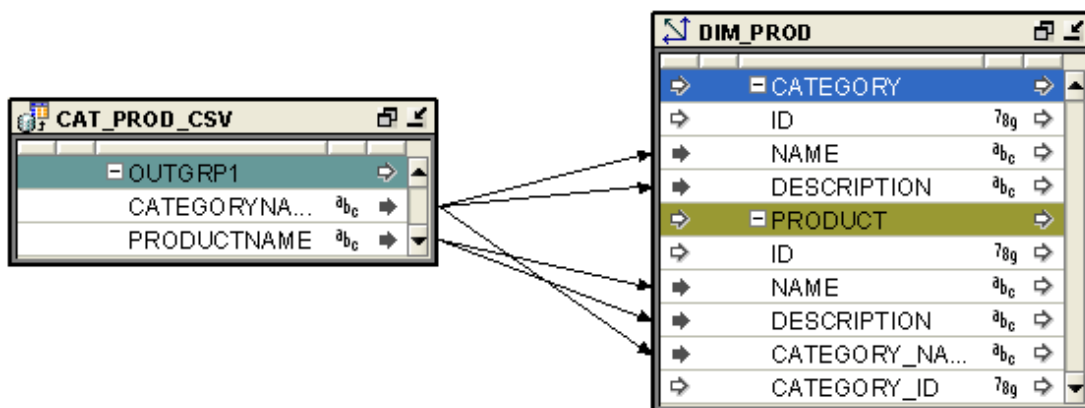
5. Mappings



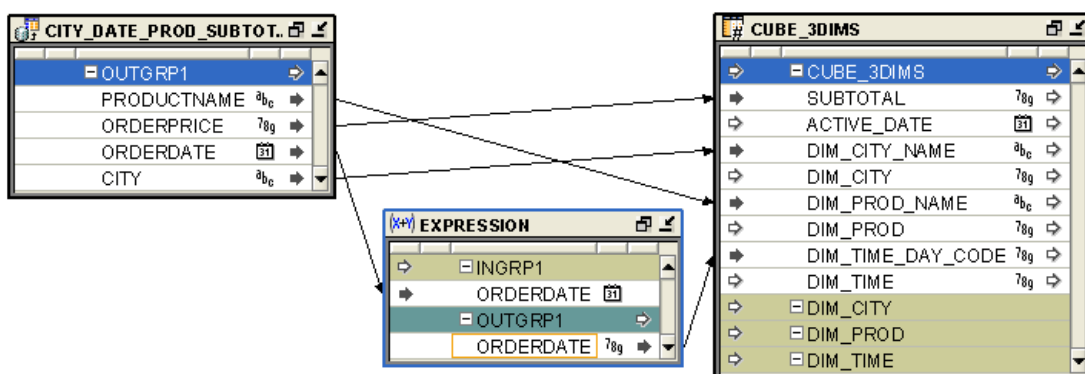
- a) The time mapping is automatically generated.
 b) Add a mapping to load city/country data from external table.



c) Add a mapping to load product/category data from external table.



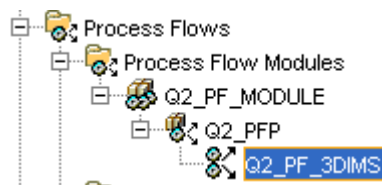
d) Add one more mapping to load the orders (city, date, product, subtotal) data to the cube. The date input should be converted to the same form as in the dimension by using an expression.



Expression for ORDERDATE 78g

```
1 TO_NUMBER(TO_CHAR(INGRP1.ORDERDATE, 'YYYYMMDD'), '99999999')
```

6. Process flow



- a) Create a new process flow. This may require a new process flow module, a new process flow package and a process flow location.

Edit Oracle Workflow Location: Q2_PF_MODULE_LOCATION

Name: Q2_PF_MODULE_LOCATION

Description:

Password: *****

Type: HOST:PORT:SERVICE

Host: localhost

Port: 1521

Service Name: orcl

Schema: OWF_MGR

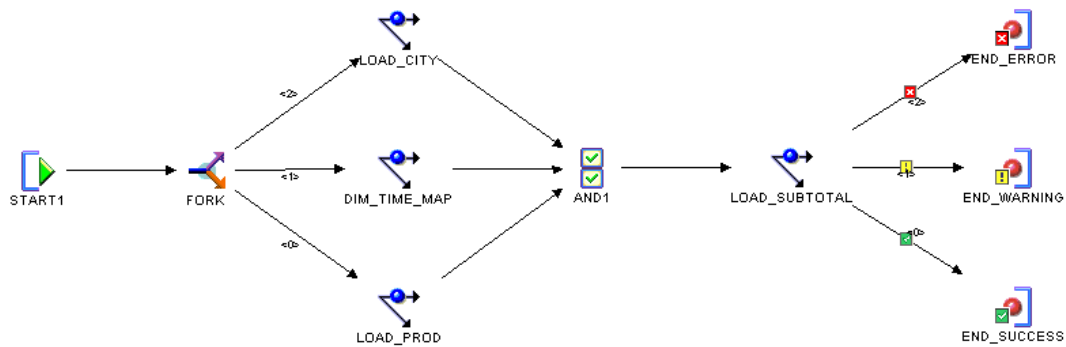
Version: 2.6.4

Test Connection

Test Results: Successful!

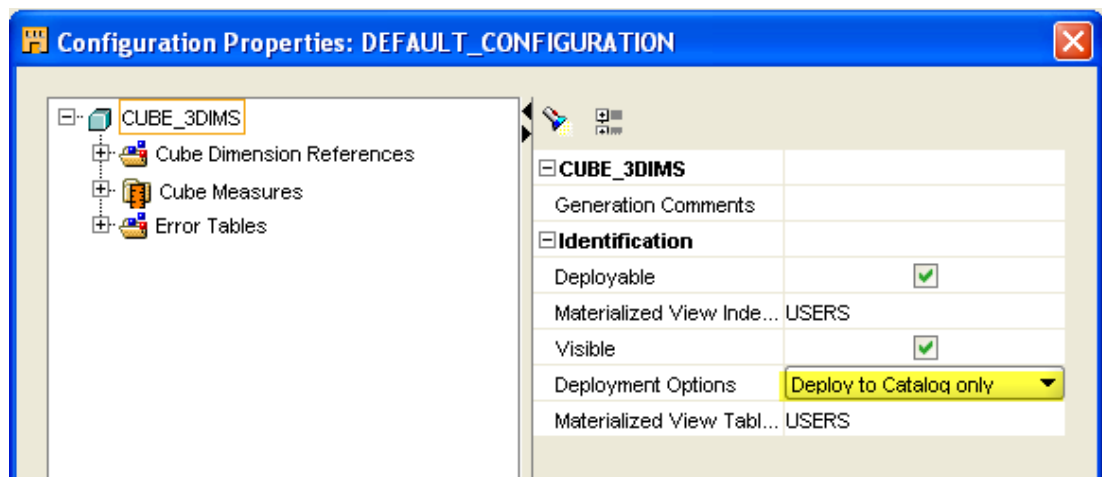
Help OK Cancel

- b) Design the flow.



7. Deployment

- a) Configure all the dimensions and cubes to be deployed to catalog only.



- b) Set actions to all the external tables, tables, sequences, dimensions, cubes, mappings and process flows to "Create". If old deployments exist, drop them before create, or use the equivalent "Replace" action.

Control Center: DEFAULT_CONTROL_CENTER

File Edit View Window Help

View: All Objects

Object Details

Details History

| Object | Design Status | Deploy Action | Deployed | Deploy Status | Location | Module |
|----------------------|---------------|---------------|------------------|---------------|--------------|--------------|
| Q2_PFP | Unchanged | Create | 12/4/11 11:54 AM | Success | Q2_PF_MOD... | Q2_PF_MOD... |
| DIM_TIME_MAP | Unchanged | Create | 12/4/11 11:54 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| LOAD_CITY | Unchanged | Create | 12/4/11 11:54 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| LOAD_PROD | Unchanged | Create | 12/4/11 11:54 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| LOAD_SUBTOTAL | Unchanged | Create | 12/4/11 11:54 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_CITY | Unchanged | Create | 12/4/11 11:54 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_PROD | Unchanged | Create | 12/4/11 11:54 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_TIME | Unchanged | Create | 12/4/11 11:54 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| CUBE_3DIMS | Unchanged | Create | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_CITY | Unchanged | Create | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_PROD | Unchanged | Create | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_TIME | Unchanged | Create | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| CAT_PROD_CSV | Unchanged | Create | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| CITY_DATE_PROD_SU... | Unchanged | Create | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| COUNTRY_CITY_CSV | Unchanged | Create | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_CITY_SEQ | Unchanged | Create | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_PROD_SEQ | Unchanged | Create | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_TIME_SEQ | Unchanged | Create | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| Q2_FILE_LOC | Unchanged | Create | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | |

Default Actions Reset Actions

c) Execute the actions and confirm all the executions are successful.

Control Center: DEFAULT_CONTROL_CENTER

File Edit View Window Help

View: All Objects

Object Details

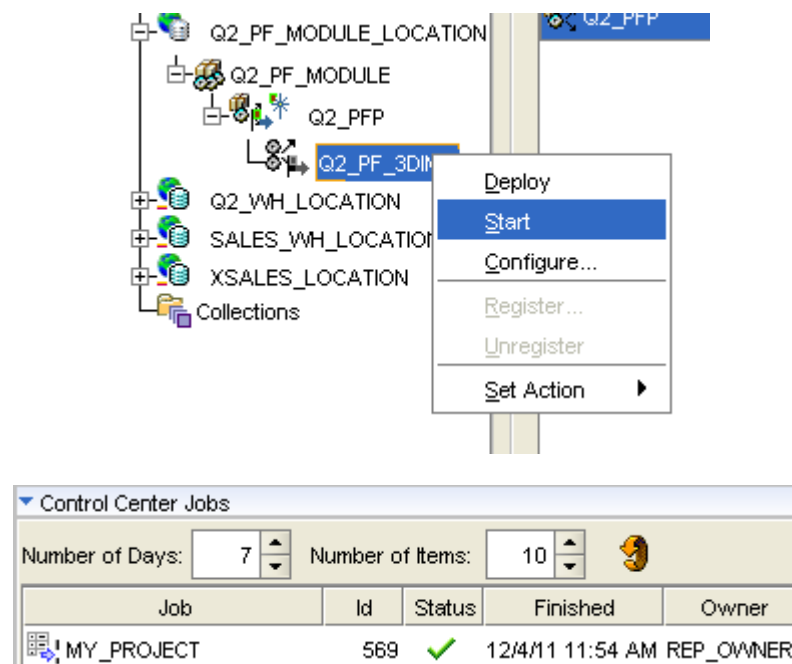
Details History

| Object | Design Status | Deploy Action | Deployed | Deploy Status | Location | Module |
|----------------------|---------------|---------------|------------------|---------------|--------------|--------------|
| Q2_PFP | Unchanged | None | 12/4/11 11:54 AM | Success | Q2_PF_MOD... | Q2_PF_MOD... |
| DIM_TIME_MAP | Unchanged | None | 12/4/11 11:54 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| LOAD_CITY | Unchanged | None | 12/4/11 11:54 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| LOAD_PROD | Unchanged | None | 12/4/11 11:54 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| LOAD_SUBTOTAL | Unchanged | None | 12/4/11 11:54 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_CITY | Unchanged | None | 12/4/11 11:54 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_PROD | Unchanged | None | 12/4/11 11:54 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_TIME | Unchanged | None | 12/4/11 11:54 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| CUBE_3DIMS | Unchanged | None | 12/4/11 11:54 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| CUBE_3DIMS | Unchanged | None | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_CITY | Unchanged | None | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_PROD | Unchanged | None | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_TIME | Unchanged | None | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| CAT_PROD_CSV | Unchanged | None | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| CITY_DATE_PROD_SU... | Unchanged | None | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| COUNTRY_CITY_CSV | Unchanged | None | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_CITY_SEQ | Unchanged | None | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_PROD_SEQ | Unchanged | None | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| DIM_TIME_SEQ | Unchanged | None | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | Q2_VWH... |
| Q2_FILE_LOC | Unchanged | None | 12/4/11 11:53 AM | Success | Q2_VWH_LO... | |

Default Actions Reset Actions

Active Configuration: DEFAULT CONFIGURATION

8. Start the process flow
 - a) Select the process flow and start it.



9. View the data
 - a) Data in dimensional tables should be visible now.

| Relational Data Viewer : DIM_PROD | | | | | | | |
|--|---------------|----------------|-------------|----------------------|--------------|---------------------|------------|
| Object Window Help | | | | | | | |
| Execute Query Get More Where Clause... | | | | | | | |
| | DIMENSION_KEY | CATEGORY_NAME | CATEGORY_ID | CATEGORY_DESCRIPTION | PRODUCT_NAME | PRODUCT_DESCRIPTION | PRODUCT_ID |
| 1 | -10 | Produce | -10 | Produce | | | |
| 2 | -11 | Seafood | -11 | Seafood | | | |
| 3 | -12 | Beverages | -12 | Beverages | | | |
| 4 | -13 | Condiments | -13 | Condiments | | | |
| 5 | -14 | Confections | -14 | Confections | | | |
| 6 | -15 | CategoryName | -15 | CategoryName | | | |
| 7 | -16 | Meat/Poultry | -16 | Meat/Poultry | | | |
| 8 | -17 | Dairy Products | -17 | Dairy Products | | | |
| 9 | -18 | Grains/Cereals | -18 | Grains/Cereals | | | |
| 10 | 97 | Beverages | -12 | Beverages | Chai | Chai | 97 |
| 11 | 98 | Produce | -10 | Produce | Tofu | Tofu | 98 |
| 12 | 99 | Beverages | -12 | Beverages | Chang | Chang | 99 |
| 13 | 100 | Seafood | -11 | Seafood | Ikura | Ikura | 100 |
| 14 | 101 | Seafood | -11 | Seafood | Konbu | Konbu | 101 |
| 15 | 102 | Dairy Products | -17 | Dairy Products | Geitost | Geitost | 102 |
| 16 | 103 | Confections | -14 | Confections | Pavlova | Pavlova | 103 |
| 17 | 104 | Grains/Cereals | -18 | Grains/Cereals | Filo Mix | Filo Mix | 104 |
| 18 | 105 | Confections | -14 | Confections | Maxilaku | Maxilaku | 105 |
| 19 | 106 | Grains/Cereals | -18 | Grains/Cereals | Tunnbr?d | Tunnbr?d | 106 |
| 20 | 107 | Confections | -14 | Confections | Chocolade | Chocolade | 107 |
| 21 | 108 | Seafood | -11 | Seafood | Spegesild | Spegesild | 108 |
| 22 | 109 | Meat/Poultry | -16 | Meat/Poultry | Tourti?re | Tourti?re | 109 |

Relational Data Viewer : CUBE_3DIMS

Object Window Help

Execute Query Get More Where Clause...

| | SUBTOTAL | DIM_CITY | DIM_PROD | DIM_TIME |
|----|----------|----------|----------|----------|
| 1 | 100.3 | 137 | 109 | 1515 |
| 2 | 648.72 | 137 | 153 | 1560 |
| 3 | 68.85 | 137 | 107 | 1763 |
| 4 | 778 | 137 | 119 | 1712 |
| 5 | 197.62 | 137 | 119 | 1654 |
| 6 | 110.4 | 137 | 114 | 1627 |
| 7 | 196 | 137 | 170 | 2027 |
| 8 | 216 | 137 | 138 | 1640 |
| 9 | 306 | 137 | 133 | 1712 |
| 10 | 1237.9 | 137 | 161 | 2027 |
| 11 | 36 | 137 | 144 | 1515 |

b) Cube data is visible.

Cube Data Viewer : CUBE_3DIMS

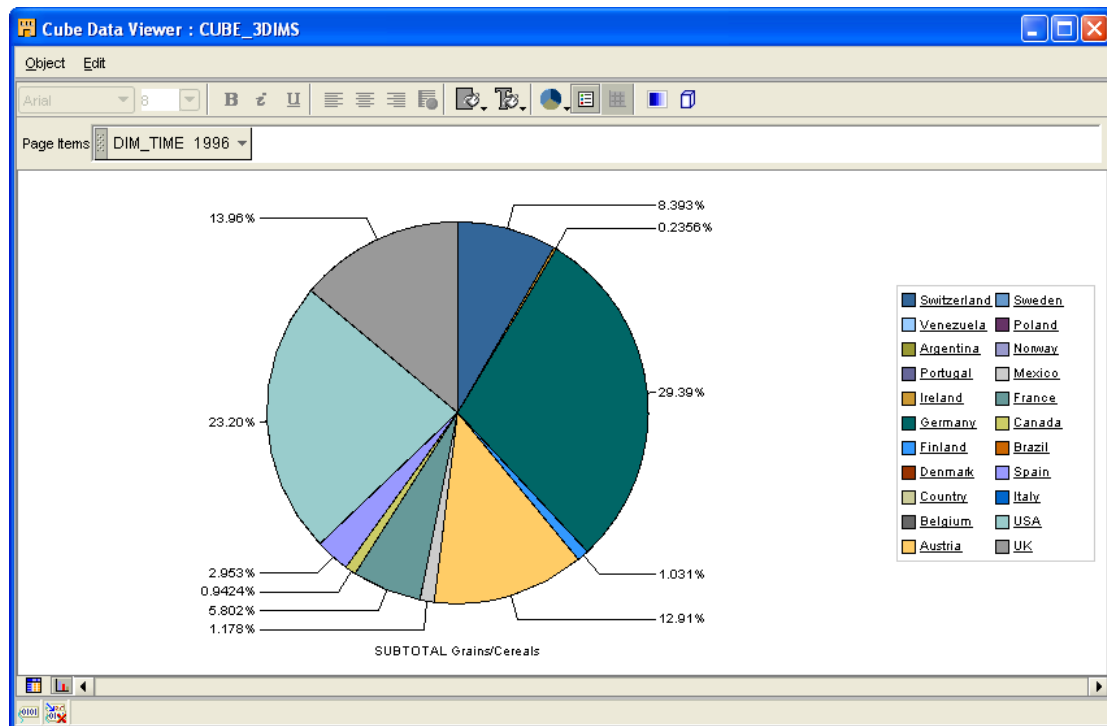
Object Edit

Arial 8 B I U

Page Items DIM_TIME 1996

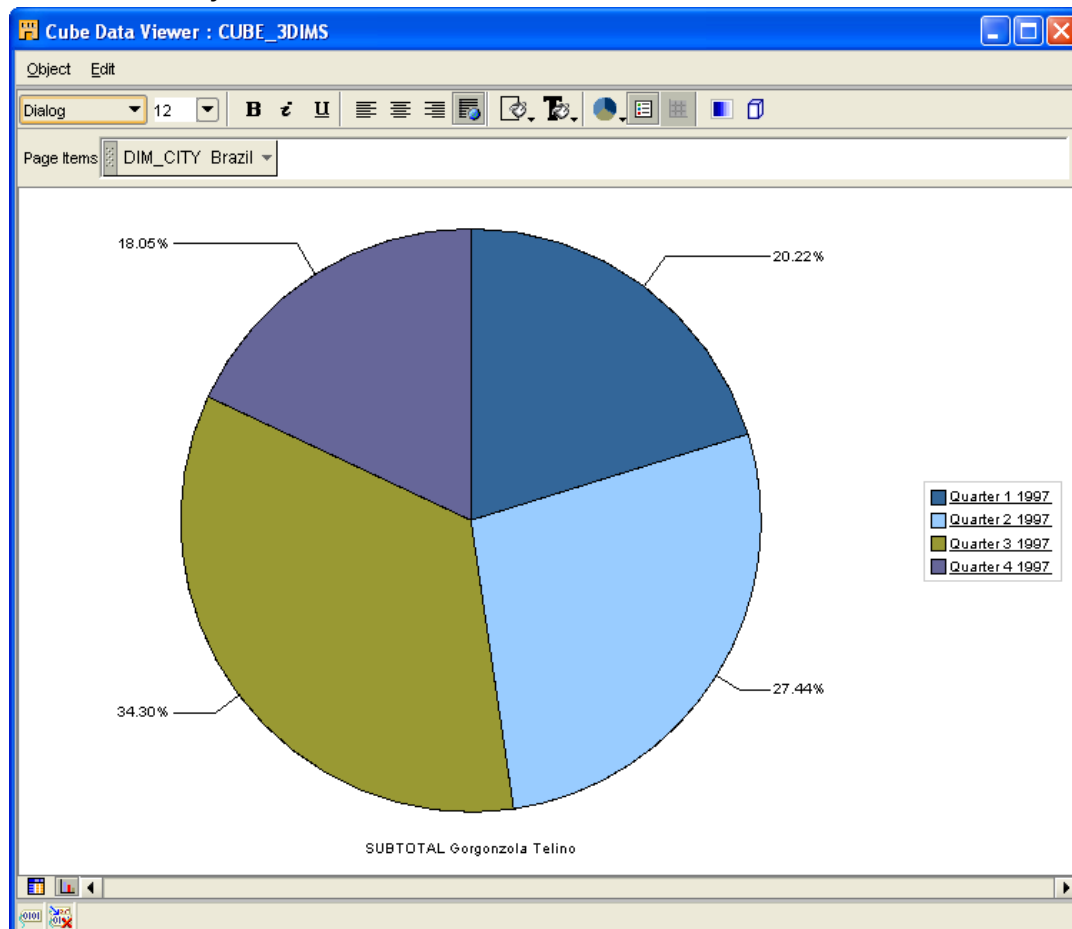
| | SUBTOTAL | | | | | | |
|-------------|----------------|----------------|--------------|--------------|-------------|------------|-----------|
| | Grains/Cereals | Dairy Products | Meat/Poultry | CategoryName | Confections | Condiments | Beverages |
| Switzerland | 798 | 1,320 | 343 | | 486 | | 534 |
| Venezuela | | 2,775 | 2,641 | | 2,883 | 404 | 86 |
| Argentina | | | | | | | |
| Portugal | | 396 | 100 | | 168 | 957 | 36 |
| Ireland | 22 | 2,199 | 472 | | 591 | 952 | 1,789 |
| Germany | 2,795 | 9,542 | 2,474 | | 3,708 | 2,748 | 5,277 |
| Finland | 98 | 1,824 | | | 72 | | 920 |
| Denmark | | 869 | 106 | | 600 | 470 | 907 |
| Country | | | | | | | |
| Belgium | | 1,136 | 1,248 | | 2,462 | | 442 |
| Austria | 1,228 | 2,985 | 1,386 | | 662 | 2,721 | 13,664 |
| Sweden | | 1,843 | 384 | | 1,814 | 529 | 1,065 |
| Poland | | 300 | | | | | |
| Norway | | 786 | | | | | 54 |
| Mexico | 112 | 906 | 646 | | 560 | 176 | 738 |

c) Switch to from table view to charts and select pie chart.



c)

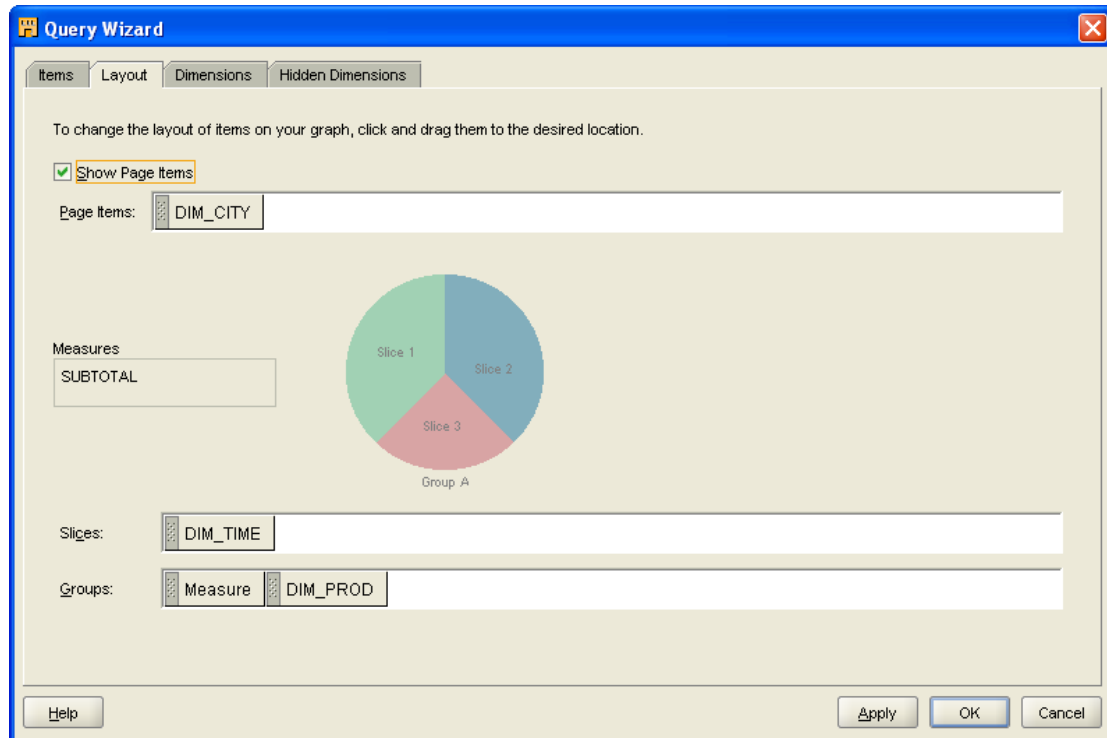
The following is a generated pie chart of “Quarterly Sales” of “Gorgonzola Telino” in Brazil for the year 1997.



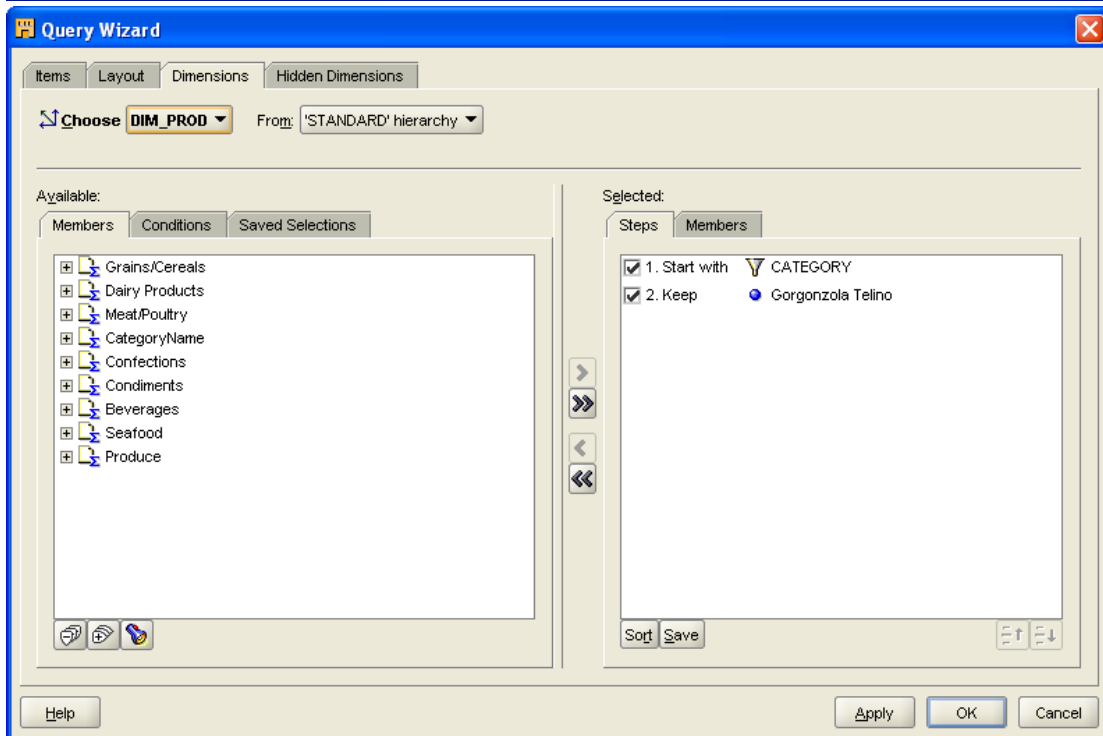
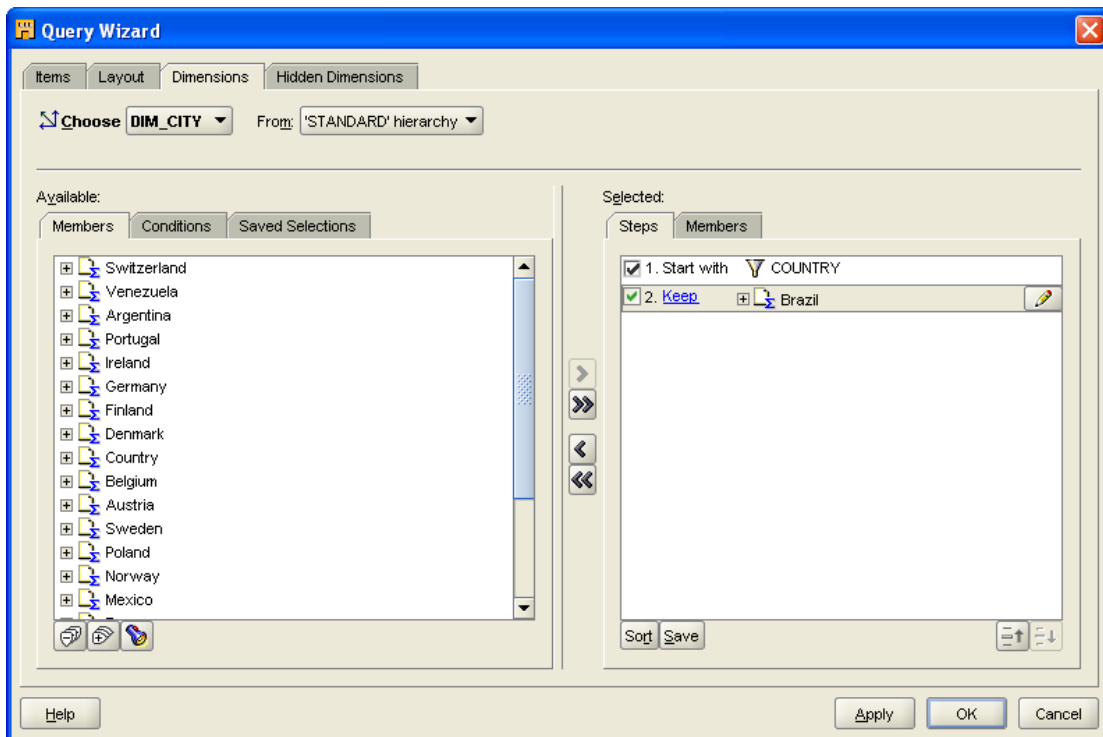
Actually, arbitrary country/city, product/category or time span can be specified to generate such a pie chart. That is shown in the next part.

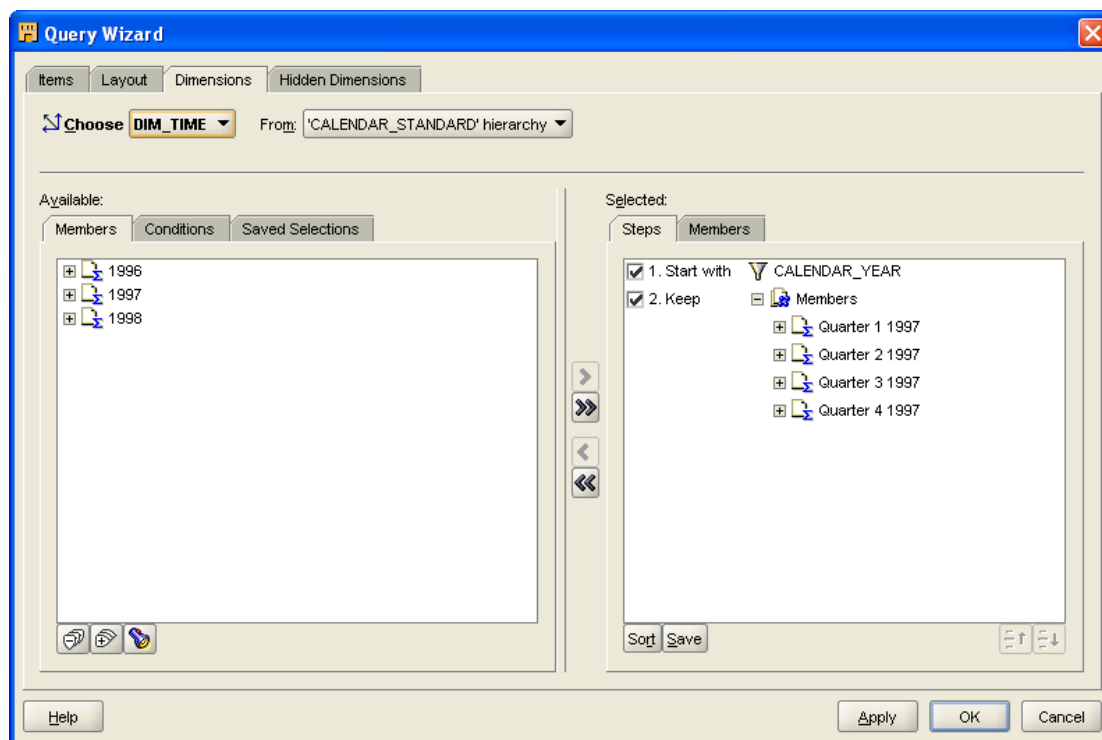
d)

When viewing data, use query builder to retrieve sales figures interested in. The layout basically defines the slicing and grouping criteria of the data view.



Specific slicing or dicing operations can be done to each dimension. For example, the user may keep only Brazil in DIM_CITY, the four quarters in 1997 in DIM_TIME and Gorgonzola Telino in DIM_PROD. Then the pie chart of the previous part will show up.





The user may specify freely what to be shown in the chart by applying filters to each dimension. In this way, a variety of sales figures can be obtained from this application.

Cube Data Viewer : CUBE_3DIMS

Object Edit

Dialog 12 B I U

Page Items DIM_PROD Grains/Cereals

| | | SUBTOTAL | | | |
|---------------|----------------|----------------|----------------|----------------|----------------|
| | | Quarter 1 1997 | Quarter 2 1997 | Quarter 3 1997 | Quarter 4 1997 |
| | Grains/Cereals | | | | |
| | Dairy Products | | | | |
| | Meat/Poultry | | | | |
| USA | CategoryName | 2,926 | 1,416 | 4,921 | 3,197 |
| Boise | Confections | | 1,216 | 2,718 | 2,771 |
| Butte | Condiments | | | | |
| Elgin | Beverages | | | | |
| Eugene | Seafood | | | 1,463 | |
| Lander | Produce | | | | |
| Seattle | | 578 | | 152 | 426 |
| Kirkland | | 200 | 200 | | |
| Portland | | | | | |
| Anchorage | | 588 | | 588 | |
| Albuquerque | | 2,926 | 2,926 | | |
| Walla Walla | | | | | |
| San Francisco | | | | | |