

Modélisation des ED : Concepts de base

Modélisation des BD : Entité et relation.

Modélisation des DW : dimension et mesure.

-**les dimensions :** les points de vues depuis lesquels les mesures peuvent être observées/ une table qui contient les axes d'analyse selon lesquels on veut étudier des données observables.

-**les mesures :** valeurs numériques que l'on compare. C'est le résultat d'une opération d'agrégation des données

-**les propriétés des mesures :**

- Fait additif : additionable suivant toutes les dimensions. ex quantité vendue - chiffre d'affaire.
- Fait semi-additif : additionable selon certains des dimensions. ex niveau de stock - nombre de transactions.
- Fait non-additif : non-additionable, recalculer. ex MxCa pour l'ensemble des magasins.

-Les faits :

-**Un fait** représente un sujet d'analyse, la valeur d'une mesure, mesurée ou calculée, selon un membre de chacune des dimensions.

-**Les mesures** sont stockées dans les tables de faits.

-**La table de fait** contient les données observables(les faits)/les mesures, les clés vers les tables de dimensions, les dimensions dégénérées(sans attribut).

-Dans un **entrepot de données** les "Faits" sont normalement numériques puisqu'elle sont d'ordre quantitatif.

-**Caractéristiques d'une table de faits :**

- contient les valeurs numériques de ce qu'on désire.
- contient les clés étrangères.

-contient un nombre de colonnes réduit.

-contient plus d'enregistrement qu'une table de dimensions

-**Caractéristiques des info de la table de faits :**

- elles sont numériques et sont utilisées pour faire des SUM, AVG
- les données sont additives/semi-additives

-les mesures de la table doivent référer et avoir un lien direct aux clés de dimension.

-**une dimension peut être définie comme** un thème ou un axe selon lequel les données seront analysées. C'est un axe d'analyse

-**une dimension contient** des membres organisées en hiérarchie. (temps: année, semestre, mois, jour).

-**structure de base d'une table de dimension :**

=>Clés de substitution(cléprimaire)/Clés d'affaires(clénaturelle), attributs de la dim , Clés spéciales (pour la gestion de l'historique de la dimension).

-**Caractéristiques d'une dimension:**

- contient le détail sur les faits.
- contient les info descriptives des valeurs numériques de la table de faits.

-les attributs sont souvent utilisés comme "lignes" et "colonnes" dans un rapport ou résultat de requête.

=>**Clé de substitution** : clé non-intelligente utilisée dans un DW pour remplacer et compléter la clé artificielle du sys opérationnel afin de rendre un élément unique dans la dimension.

=>**Clé naturelle** est en générale composée de plusieurs colonnes.

-**Dans un système opérationnel on utilise une clé artificielle** afin d'identifier d'une façon unique un élément de l'entité.

=>**La clé de substitution ne doit pas être confondue avec la clé artificielle attribuée par le sys opérationnel.**

Fonctionnalité des clés de substitution :

- Remplacer la clé naturelle : la clé de substitution remplace la clé artificielle
- Compléter l'information : la clé de substit est utilisée dans l'entrepot de données seulement

Avantages des clés de substitution :

- Performance : accélère l'accès aux données du moment où l'on va utiliser un index numérique
- Indépendance du sys source : on ne peut garantir que la clé d'affaire ne change pas dans les systèmes sources

-Historique des changements et granularité infinie : gérer la clé de substitution pour garder l'historique des changements de la dimension selon certains critères

Table d'une BD multidimensionnelle :

-Date retrait : date où l'enregistrement a été retiré.

-Indicateur effectif: "O" si l'enregistrement actif, "N" sinon.

(facultatif) :

-Type de l'évolution (SCD)

-Valeur avant le changement

Les dimensions: Clés spéciales:

-Date effective: date de création de l'enregistrement.

Evolution des dimensions :

-**Lente** : gestion de la solution

-> **Ecrasement de l'ancienne valeur:**

(correction des info erronées)

-facile à mettre en œuvre (avantage)

-perte de la trace des val antérieures

des attributs. (inconvenient)

-perte de la cause de l'évolution des les

faits mesurés. (inconvenient)

-> **Versionnement** (ajout d'un nouvel

enregistrement, utilisation clé primaire)

-permet de suivre l'évolution des

attributs (avantage)

-permet de segmenter la table de faits

en fonction de l'historique. (avantage)

-accroît le volume de la table.

(inconvenient)

-> **Valeur d'origine/courante** (ajout d'un

nouvel attribut) :

-avoir deux visions simultanées des

données. (avantage)

-voir les données comme si le

changement n'avait pas eu lieu (avantage)

-inadapté pour suivre plsr valeurs

d'attributs intermédiaires. (inconvenient)

-**Rapide** :

-> Subit des changements très fréquents

dont on veut préserver l'historique

-> **Solution** : isoler les attributs qui

changent rapidement.

Modélisation des ED : Modèle d'un DW

-Modèle en étoile : table de faits au centre du schéma

-Modèle en flocon de neige

-Modèle en constellation/galaxy

Modèle en flocon de neige :
-> les dimensions sont décomposées en hiérarchie, c'est une méthode de normalisation des tables de dimensions.
-> le seul but est de minimiser l'espace disque sinon il n'est pas recommandé.
-> si les tables sont très volumineuses -> réduction du volume + possibilité de réaliser des analyses par palier (drill down).
-> Avantages : - Conçu pour des requêtes flexibles sur des dimensions et des relations complexes
- Performance des requêtes
- Grande évolutivité
-> Inconv : - requêtes plus complexes.
- plus de temps d'exécution.
- navigation difficile.
- nombreuses jointures.
- maintenance supplémentaire.

Modèle en constellation/galaxy:
=> Fusionner plsr modèles en étoile qui utilisent des dimensions communes.
=> Structure logique d'entrepot de données.

Démarche de conception :

étape 1 : -choisir le processus à modéliser. Processus Métier

étape 2 : -choisir le grain de faits, données granulaire.

-décider de ce que représente une ligne de la table de faits Fait à observer
-niveau de détail: transaction individuelles,...

étape 3 : -identifier les dimensions. (typiquement le temps, le client...).

Dimensions

étape 4 : -identifier les mesures de faits (de préférence additives).

Mesures

Analyse Multidimensionnelles

Le Data Warehouse est le socle indispensable pour obtenir les réponses aux questions essentielles à la prise de décision et au pilotage de l'entreprise.

Quitter un modèle de base de données opérationnelle et se baser sur un modèle de base de données décisionnelle revient à : transformer les éléments du modèle de base de données opérationnelle en dimensions et fait

Copyright al Ikhwan msawba min taraf sidkom Riad

BD Multidimensionnelles : Cube

- Un cube OLAP est une structure de données multidimensionnelle stockant les faits comme des mesures indexées par plusieurs dimensions.
- Chaque cellule d'un cube représente la mesure ou valeur quantitative d'un fait sur le croisement de plusieurs dimensions.
- L'intérêt d'un cube OLAP est d'offrir à l'utilisateur la capacité de faire des analyses multidimensionnelles ou des agrégations par axe de dimension dans l'espace.

BD Multidimensionnelles = super-tableur.

Technologies OLAP

- **ROLAP** - Relationnel OLAP : OLAP sur du relationnel
 - **MOLAP** - Multidimensionnel OLAP : OLAP sur un DW dimensionnel
 - **HOLAP** - Hybride OLAP : Mélange des deux
- =>ROLAP :
- Utilisent un SGBD relationnel classique avec des adaptations spécifiques à l'OLAP
 - organisées en schémas en forme d'étoiles ou en flocon de neige.
 - Peuvent conduire à des temps de réponses élevés.
- s'appuie sur la maturité de la tech relationnelle

- L'approche HOLAP consiste à
- Utiliser les tables comme structure permanente de stockage des données
 - Manipuler les informations du DW avec un moteur ROLAP
 - Exploiter les Data Marts selon une approche multidimensionnelle avec un système MOLAP.

généraliser des requetes plus compacte que les requetes SQL

-> select axis1 ON COLUMNS, axis2 ON ROWS, axis3 ON AXIS(0).

pour définir un axe et présenter sur l'axe tous les membres d'une dimension :

=> <dimension name>.MEMBERS

pour voir apparaitre tous les membres d'une dimension a un certain niveau :

=> <dimension name><level name>.MEMBERS

Ds un slice on peut avoir plsr membres, mais ils doivent appartenir a des dimension différentes.

// ou -- => commentaire en fin de ligne.

/* */ => commentaire sur plsr lignes

BD multidimensionnelle, Opérateurs :

- Les opérateurs appliqués sur un cube sont algébriques.
- Le résultat est un autre cube.
- Les opérateurs permettent :
 - Des extractions Slicing : Prendre une tranche du cube.
 - Des extractions Dicing : Extraire un sous-cube.
- Des changements de granularité d'une dimension :
 - Roll-up (agrégation d'une dimension -> Résumé)
 - Drill-down (informations plus détaillées).

dont : les structures sont optimisées + l'accès est rapide en lecture/écriture.

- MOLAP nécessite le pre-calcul, sur tous les niveaux de hiérarchies des dimensions.
 - Très rapide et performant mais avec des limitations de taille
 - Exemples : Essbase, SAS OLAP Server, ...
- =>HOLAP (HOLAP permet d'avoir des DW de taille importante tout en ayant des temps de réponse satisfaisants)
- Les données multidimensionnelles sont stockées et traitées en se basant sur le SGBD Relationnel et le SGBD multidimensionnel, afin d'éviter les problèmes des systèmes MOLAP et ROLAP.

- HOLAP permet d'avoir des DW de taille importante tout en ayant des temps de réponse satisfaisants.
- Exemples : Oracle, IBM DB2 OLAP Server ...

Cellule = le tuple permet d'identifier les cellules ds un cube.

Mesure = ds un tuple les mesures sont traitées comme une dimension particulière.

Set = ensemble ordonné de tuples définit sur une même dimension.(le mot apres ON)

SQL	MDX
SELECT column1, column2, ..., columnn FROM table	SELECT axis1 ON COLUMNS, axis2 ON ROWS FROM cube
FROM : une ou plusieurs tables	FROM : un cube
SELECT : • une vue des données en 2 dimensions : lignes (rows) et colonnes (columns)	SELECT : • nombre quelconque de dimensions pour former les résultats de la requête
Colonne : chaîne de caractère ou valeur numérique	Niveau : Level
Plusieurs colonnes liées ou une table de dimension	Dimension

Exemple de requête MDX :	Signification
SELECT { [Time].[2003].[Q1], [Time].[2004].CHILDREN } ON COLUMNS, { [Markets].[APAC].[Australia], [Markets].[EMEA].[France] } ON ROWS FROM CubeSales WHERE { [MEASURES].[Sales], [Product].[Ships], [Customers].[All Customers] }	← Les colonnes : Q1 2003, Q1 2004, Q2 2004, Q3 2004, Q4 2004 ← Les lignes : Australia, France ← Cube en question ← Agrégation de la mesure Sales avec la fonction SUM ← Sélection de la dimension Product (Ships seulement) ← Sélection de la dimension Customers (l'ensemble des clients)

SELECT : Description des axes du cube résultat	Chaque dimension du résultat est : <ul style="list-style-type: none">• associée à un rôle correspondant à sa représentation dans le tableau retourné par la requête MDX• Exemples : ON COLUMNS, ON ROWS• sur un ou plusieurs niveaux de l'hiérarchie :<ul style="list-style-type: none">• Exemple 1 : { [Markets].[APAC].[Australia], [Markets].[EMEA].[France] } de la dimension Markets, niveau Pays• Exemple 2 : { [Time].[2003].[Q1], [Time].[2004].[Q1].CHILDREN } de la dimension Temps, niveaux trimestre et mois (tous les mois du trimestre)
FROM : Spécification du/des cube/s de départ	<ul style="list-style-type: none">• Ensemble de cubes nécessaires à la création du cube résultat• Si plusieurs cubes nécessaires, cela implique une jointure multidimensionnelle : chaque paire de cubes doit alors posséder au moins une dimension concordante
WHERE : Restriction sur le/s cube/s de départ	<ul style="list-style-type: none">• Restrictions sur le/s cube/s de départ de la clause FROM• Spécification des restrictions par une liste de noeuds de la hiérarchie d'une dimension nommée slicer-dimension

Conception des Entrepôts : Exercice 2

- Concevoir un modèle en étoile qui permet d'analyser les ventes d'une entreprise de restauration rapide.
- Le principe est de mesurer les ventes grâce aux quantités vendues et aux bénéfices, en fonction des ventes réalisées par jour, dans un restaurant donné, pour un aliment donné.
- L'objectif est de pouvoir analyser les ventes :
 - par jour,
 - par semaine,
 - par mois
 - et par année.
- Les restaurants peuvent être regroupés en fonction de leur ville et de leur pays.

Dans cette requête, on ne s'intéresse qu'aux quantités vendues des produits Automotive et Energy en 2015 et 2016.

Deux dimensions sont considérées :

- Niveau Category de la dimension Product
- et Niveau Year de la dimension Time

```
SELECT  
{ ([Measures].[Quantity], [Product].[Automotive]),  
([Measures].[Quantity], [Product].[Energy]) }  
ON COLUMNS,  
{ [Time].[2015], [Time].[2016] }  
ON ROWS  
FROM [Sales]
```

- Maintenant nous on s'intéresse aux quantités vendues en 2015 et 2016 dans les différents continents. Deux dimensions sont considérées :
 - Niveau Year de la dimension Temps
 - et Niveau Continent de la dimension Customer
- Aussi, on s'intéresse uniquement aux ventes des produits Television. On définit alors le nouvel axe :
 - { [Product].[Television] }

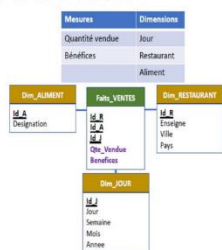
```
SELECT  
{ [Time].[2015], [Time].[2016] }  
ON COLUMNS,  
{ [Customer].[Continent].Members }  
ON ROWS  
FROM [Sales]  
WHERE { [Product].[Television] }
```

- On pourrait ajouter le 3ème axe à la requête, mais les outils OLAP ne pourraient pas le visualiser !
- Ici, on préfère utiliser une opération de Slice (filtre)
- Utilisation du Where

Conception des Entrepôts de Données :

Exercice 2, suite

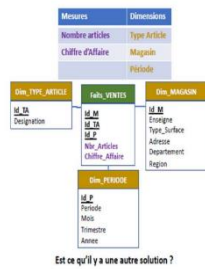
- Concevoir un modèle en étoile qui permet d'analyser les ventes d'une entreprise de restauration rapide.
- Le principe est de mesurer les ventes grâce aux quantités vendues et aux bénéfices, en fonction des ventes réalisées par jour, dans un restaurant donné, pour un aliment donné.
- L'objectif est de pouvoir analyser les ventes :
 - par jour,
 - par semaine,
 - par mois
 - et par année.
- Les restaurants peuvent être regroupés en fonction de leur ville et de leur pays.



Est ce qu'il y a une autre solution ?

Exercice 1, suite

- Une entreprise de fabrication de vaisselle jetable souhaite mettre en place un système d'information décisionnel sous la forme d'un Data Mart pour observer son activité de ventes au niveau des différents lieux de distributions de ses articles et cela dans plusieurs villes.
- Ces lieux de distributions sont renseignés par :
 - leur enseigne,
 - leur type (en fonction de leur surface),
 - leur adresse (code postal et ville),
 - leur département,
 - leur région.
- Les ventes sont renseignées selon une période qui se décline en mois, en trimestre et année.
- Les ventes sont observées par le nombre d'articles selon le type et le chiffre d'affaire



Est ce qu'il y a une autre solution ?