

**TODO**

Gestionnaire de bases de données

Analyse des données

Définitions:

- **Données agrégées** : données regroupées en classes (clusters), éventuellement organisées de façon hiérarchique. Possibilité d'appartenances à de multiples hiérarchies (cubes de données).
- **Analyse des données** : le but est de dégager des indicateurs à partir d'un grand ensemble de données, afin de faciliter la prise de décision.

cas d'utilisation :

- quels sont les magasins les plus rentables? doit-on ouvrir / fermer des magasins?
- où doit-on implanter un nouveau magasin?
- y a-t-il une corrélation entre le lancement d'une campagne publicitaire et les chiffres de vente? quels sont les supports les plus efficaces?
- quelle est la liste des clients à fidéliser?
- de quelle quantité doit-on approvisionner les magasins en fonction de la période de l'année?

analyse:

- quels sont les catégories de films/livres les plus fréquemment empruntés?
- réussite / taux d'embauche / salaire en fonction de la prépa d'origine / sexe / profession des parents

Les opérateurs d'agrégation permettent de réaliser des statistiques sur les données, sous forme d'histogrammes (ou camemberts) organisés selon des catégories définies par les valeurs de certains attributs:

principe :

- opérateur d'agrégation : comptage, somme, moyenne, écart-type (count, sum, mean, avg, ...)
- les données agrégées sont de type quantitatif
- les attributs définissant les classes sont de type qualitatif.

Exemples de requêtes faisant appel aux fonctions d'agrégation :

Nombre d'élèves par groupe de TD / par prepa d'origine etc..:

```
SELECT groupe_TD , COUNT(num_eleve)
FROM Eleve
GROUP BY groupe_TD
```

Donner les chiffres des ventes du magasin pour chaque mois de l'année

```
SELECT mois, SUM(montant)
FROM Vente
GROUP BY mois
```

Donner le nombre de ventes d'un montant > à 1000 euros pour chaque mois de l'année

```
SELECT mois, COUNT(num_vente)
FROM Vente
GROUP BY mois
HAVING montant >= 1000
```

Tester les disparités salariales entre hommes et femmes

```
SELECT sexe, avg( salaire )
FROM Employé
GROUP BY sexe
```

Tester les disparités salariales selon le niveau d'éducation

```
SELECT niveau_educatif, avg( salaire )
FROM Employé
GROUP BY niveau_educatif
```

4. Analyse des données et découverte d'informations

TODO : Manu



- fichiers csv
- matrice données (cube)
- excel

La découverte d'informations consiste à développer des outils de mise en forme des données facilitant leur analyse. Elle repose sur deux aspects :

- projection de données qualitatives sur des espaces vectoriels ("quantification" des données)
- production d'histogrammes dans le but d'analyser la distribution des données dans l'espace de reconstruction

Le but est de dégager :

- des **tendances** (covariables)
- des **modes** de la distribution (présence de plusieurs maxima)

à partir d'un **grand** ensemble de données (chiffre d'affaires, nb de ventes, masse salariale, ...) évoluant dans le **temps** et dans l'**espace**, afin de

- définir des **indicateurs** pertinents

- faciliter la prise de décision.

Vocabulaire anglophone généralement utilisé :

- Business Intelligence (BI)
- Data Warehouses (Entrepôts de données)
- OLAP (Online Analytical Processing) :



"Unlike Online Transaction Processing (OLTP), where typical operations read and modify individual and small numbers of records, OLAP deals with data in bulk, and operations are generally read-only."

Entrepôts de données (Data warehouses) / Magasins de données (Data Mart)

Exemples de grandes masses de données :

- Masses de données (pullulantes) : tickets de caisse, clics web, appels tel, opérations bancaires, remboursements no URSSAF, trajets SNCF...
- Données importantes : fichiers de clients, données biométriques, campagnes de mesures, sondages,...
- Données géographiquement localisées (gestion d'un "territoire") : appels tel, centres de production, consommation eau-électricité-gaz, infractions pénales, arrêts maladie, prêts bancaires, allocations chômage, jugements des TGI, accidents du travail...

Remarque : Les transactions marchandes sont un cas classique (acte d'achat bien répertorié et enregistrés, livres de comptes, ...)

Problèmes

- Quels sont les catégories de films/livres les plus fréquemment empruntés?
- Réussite / taux d'embauche / salaire en fonction de la prépa d'origine / sexe / profession des parents
- Tester les disparités salariales hommes/femmes en fonction du niveau d'éducation.
- Donner le taux de réussite par groupe de matière en fonction de la filière d'origine (MP, PSI, PC, PT, ...)
- A quelles heures de la journée la messagerie est-elle la plus sollicitée?
- Comment se répartissent géographiquement les utilisateurs de la messagerie?

Agrégation

L'agrégation consiste:

- à positionner les mesures sur des axes (temporels, géographiques,...)
- ou à les organiser en classes (et sous-classes) selon la valeur d'un ou plusieurs attributs.
- à utiliser des **opérateur d'agrégation** :

- comptage, somme, moyenne, écart-type, max, min, ...
- (count, sum, mean, avg, max, min...)
- afin de dégager :
 - des tendances (selon des dimensions)
 - des modes (analyse par histogramme)
 - des corrélations (analyse croisée)

Dimensions

- (qui?) Quels sont les magasins les plus rentables? doit-on ouvrir / fermer des magasins?
- Où doit-on implanter un nouveau magasin?
- Y a-t-il une corrélation entre le lancement d'une campagne publicitaire et les chiffres de vente? quels sont les supports les plus efficaces?
- (qui?) Quelle est la liste des clients à contacter?
- (quand?) De quelle quantité doit-on approvisionner les magasins en fonction de la période de l'année?

Problèmes :

- définir des intervalles temporels?? créneaux horaires? ⇒ distributions, histogrammes
- définir des secteurs géographiques?
 - ⇒ hiérarchies pays > département > région
- Taille du message, has_attachement?
 - ⇒ mesures sur des faits élémentaires

4.1 Tableaux de données

Organisation des données sous forme de tableaux bidimensionnels :

Schémas de données

- Une mesure est un jeu de valeurs organisé sous forme de **tuple**
- A un tuple on associe en général un **schéma de données**.

SCHEMA :

Nom

Prénom

Adresse

Âge

DONNEES :

Dubois	Martine	29, rue du Verger, Orléans	22
--------	---------	----------------------------	----

- Définir un **schéma** consiste à définir :
 - une liste d'attributs (labels) associées à chacune des valeurs du tuples.
- A chaque **attribut** correspond :
 - un *intitulé*
 - un *domaine* de valeurs (type/format des données)

Tableau de données

Un tableau de données est une liste (finie et ordonnée) de tuples, chaque tuple obéissant à un même schéma \$R\$.

Tableau de données

Nom	Prénom	Adresse	Âge
Dubois	Martine	29, rue du Verger, Orléans	22
Gilbert	Jonas	8, rue des Fleurs, Blois	23
Dalban	Pierre	13, av. du Général, Privas	22
...
Manoukian	Marianne	55, place des Bleuets, Aubagne	24

schéma

tuple

Formats d'échange

Les principaux formats d'échange de données sont :

- csv
- json
- xml



TODO

Exemples :

- [Clients](#)
- [Cours de l'Euro](#)
- [codes postaux](#)
- [codes postaux](#)

Sites de données :

- www.data.gouv.fr
- [scientific data](#)
- [Liste de ressources open data](#)
- data.gov
- [open food facts](#)

4.2 Faits élémentaires

- Notion de fait élémentaire (*fact*): transaction ou opération localisée dans le temps et dans l'espace

Exemples de "fait":

- Achat/Vente
- Opération bancaire (débit/crédit)
- Consultation (site web)
- Souscription à un contrat d'assurance
- Appel téléphonique
- Inscription

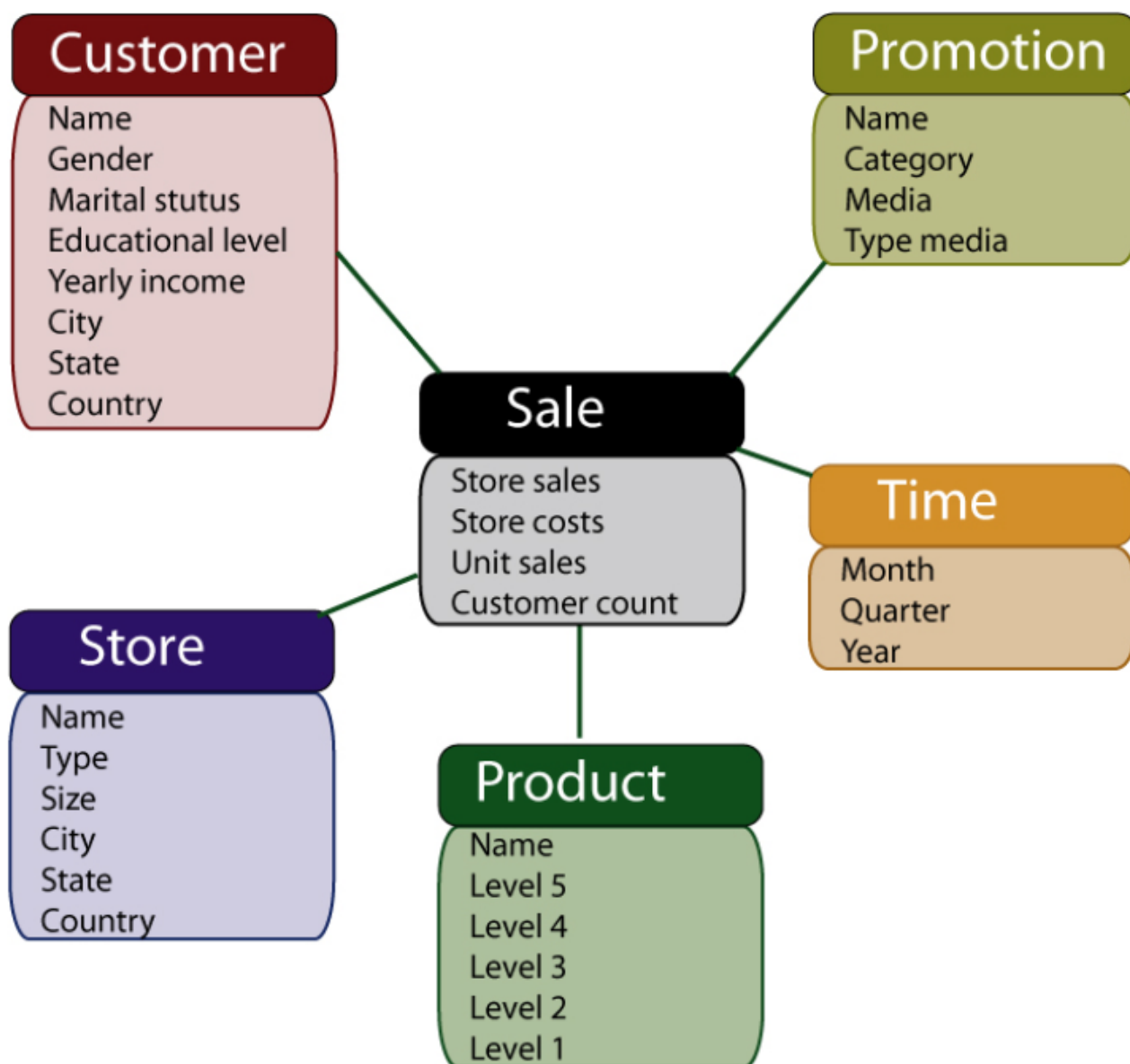
Tous ces faits peuvent être localisés. Des mesures peuvent être effectuées sur ces faits (montant d'une vente, durée d'un appel, montant d'une opération bancaire, ...)

Points clés :

- distinction entre **Dimension** et **Mesure**.
 - Notion de dimension : qui? quoi? où? quand? Comment? : associe des **coordonnées** à l'événement (géographiques, temporelles) et par extension une catégorie.
 - Notion de **mesure(s)** associées à l'événement (exemple : montant de la vente)
- distributions, répartitions, etc... cf analogie proba/stats : événement aléatoire, vecteur aléatoire, ...
 - les événements sont associés par paquets sur des intervalles réguliers ou selon des catégories discrètes.
 - Fonctions d'agrégation : réalise la mesure sur les groupe : somme, comptage, moyenne, min, max, etc...
 - histogramme : nb d'événements observés par secteur sur un maillage régulier de l'espace des coordonnées. Par extension mesure sur ce maillage par une fonction d'agrégation.

Modèle en étoile

- un fait est une association située au centre du schéma. Les attributs de l'association sont les mesures effectuées
- une dimension est une relation participant au fait. Les dimensions sont donc décrites par des attributs (ex : attributs année, trimestre, mois, jour, heure, minute, seconde,...pour une dimension temporelle)
- pour chaque dimension, on décrit une hiérarchie sur les différents attributs de la dimension en définissant un ordre, du particulier au général.



Exemples :

- sur la dimension temporelle : mois \subset trimestre \subset année
- sur la dimension promotion : nom \subset catégorie \subset média \subset type de média

etc...

Exemples

- [Ventes](#)
- [Ventes](#)
- [Pistes](#)

Tables pivot

Les tables pivot permettent d'analyser des faits selon deux dimensions organisées sur les deux axes d'un tableau

Cubes de données

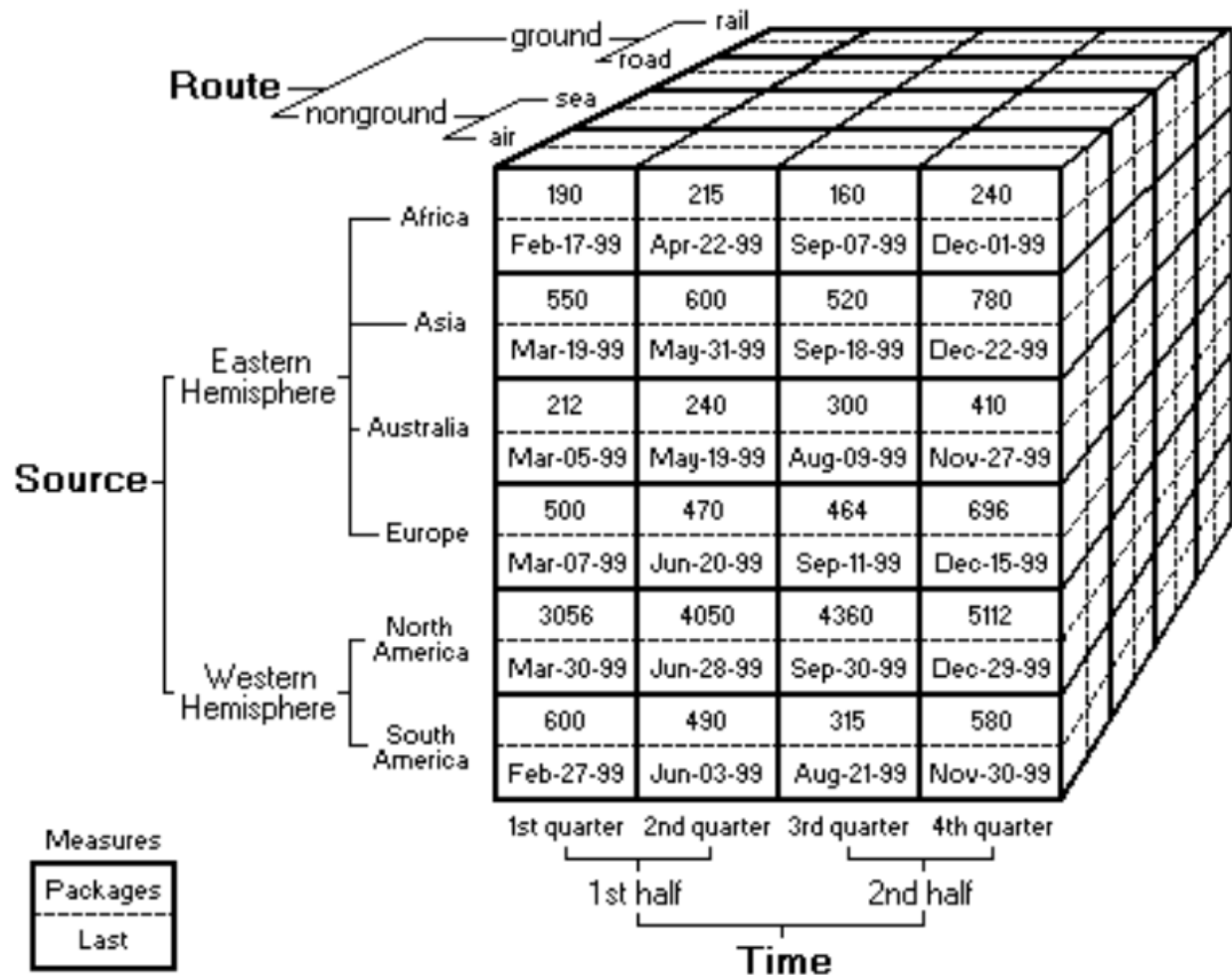
Un cube de données est une structure de données organisée sur le principe des espaces vectoriels. Différents axes sont définis, chaque axe étant associé à une dimension particulière.

- Les dimensions peuvent correspondre à des valeurs discrètes (catégories : type de produit, catégorie de client,...) ou continues (valeurs temporelles ou géographiques, ...).
- Chaque fait est décrit comme un point de l'espace vectoriel. Il est positionné dans une cellule du cube. A ce point sont associées une ou plusieurs mesures.
- Le cube est un ensemble de cellules (voir figure), chaque cellule correspondant à un intervalle (sur les axes continus) ou une valeur (sur les axes discrets).

Un élément essentiel du modèle de données est la définition de **hiérarchies** sur les dimensions du cube. Chaque dimension se divise en intervalles et sous-intervalles (pour le continu/ quantitatif) ou en catégories et sous-catégories (pour le discret/qualitatif)

Les hiérarchies sur les différentes dimensions permettent de définir le "niveau de résolution" sur les différentes dimensions.

- On peut ainsi s'intéresser à l'évolution d'une certaine grandeur au cours du temps année par année, trimestre par trimestre ou mois par mois selon le niveau de résolution choisi.
- → Hiérarchie : description arborescente d'intervalles et de sous-intervalles sur une dimension. Implémente différentes granularités sur la dimension considérée.



La structure de cube de données est adaptée pour la réalisation d'histogramme multidimensionnels, selon les axes choisis et le niveau de résolution choisi, à l'aide de fonctions d'agrégation.

- Histogramme et agrégation
 - (vue quantitative) comptage/répartition d'événements sur un intervalle (discrétisation d'une distribution d'événements)
 - (vue qualitative) comptage d'événements par catégorie
 - (vue intermédiaire) comptage d'événements par catégories hiérarchisées

4.3. Mise en oeuvre

Pandas

<http://pandas.pydata.org/pandas-docs/stable/10min.html>

XMLA / MDX

From:

<https://wiki.centrale-med.fr/informatique/> - **WiKi informatique**

Permanent link:

https://wiki.centrale-med.fr/informatique/tc_info:cm7

Last update: **2018/12/13 08:22**

