

# Qualité des données

## Etude de cas

M2 DataScale

Zoubida Kedad



## Utilisation d'un ETL pour évaluer et améliorer la qualité des données

- Mise en œuvre d'un ETL
  - Talend
- Définition d'un scénario comportant
  - Des processus d'évaluation de la qualité des données
  - Des processus d'amélioration de la qualité des données

## Scénario : intégration de données de pollution

- Plusieurs sources de données distinctes
- Des objets cibles qui seront calculés à partir des sources de données

Zoubida Kedad

3

## Le schéma cible – Pollution de l'air

Paris (ID Polluant, Taux\_moyen\_jour, NBS, NBC, Statut)

Hauts\_de\_Seine (ID Polluant, Taux\_moyen\_jour, NBS, NBC, Statut)

Yvelines (ID Polluant, Taux\_moyen\_jour, NBS, NBC, Statut)

Stations (ID Station, Adresse, Tel, Contact\_mail)

- NBS : nombre de stations utilisées pour le calcul
- NBC : nombre de capteurs mobiles utilisés pour le calcul
- Statut : « Alerte pollution » ou « Normal »

Zoubida Kedad

4

## Les sources

### ■ Source S1 – Mesures de stations fixes

Mesures (ID\_Polluant, Date, ID\_Station, Taux\_relevé)

Station (ID\_Station, Num, Rue, Ville, Code\_postal, Tel, Contact\_Mail)

### ■ Source S2 – Mesures de stations fixes

Mesures (ID\_Polluant, Date, ID\_Station, Taux\_relevé)

Station (ID\_Station, Num, Rue, Ville, Code\_postal, Tel, Contact\_Mail)

Zoubida Kedad

5

## Les sources (suite)

### ■ Source S3 – Polluants

Polluants (ID\_Polluant, Description, Seuil\_toléré)

### ■ Source S4 – Mesures de capteurs mobiles

Mesures (ID\_Polluant, Date, ID\_Capteur, Localisation, Taux\_relevé)

### ■ Source S5 – Mesures de capteurs mobiles

Mesures (ID\_Polluant, Date, ID\_Capteur, Localisation, Taux\_relevé)

Zoubida Kedad

6

## Problèmes de qualité

- Conformité à un format, une codification
- Hétérogénéité des échelles, de la granularité
- Complétude des données
- Détection et élimination de doublons

Zoubida Kedad

7

## Résultats attendus

- Bilan sur l'audit de la qualité des sources et des données calculées
  - Pour les facteurs de qualité considérés
- Workflow de traitement permettant d'améliorer la qualité des données en utilisant l'ETL choisi
- Date de fin : 12/11/2021
- Livrables
  - Compte rendu oral par groupe pendant les séances de cours
  - Démonstration
  - Compte rendu écrit (**pas de rapport !**)

Zoubida Kedad

8