



Compte rendu  
04 Novembre 2022

---

Projet Qualité des données

---

**Encadré par**

M. ZOUBIDA KEDAD

**Réalisé par**

AOUCHICHE SALAH  
HADJ MAHFOUD KENZA  
BOUDJEMAI ALI  
BELHANNACHI MOUNJJI HOCINE

















- Lors de notre implémentation sur Talend nous avons adopté la métrique qui consistait à évaluer la complétude en calculant le rapport entre le nombre de tuples dont certaines valeurs d'attributs sont nulles sur le nombre total de tuple, cela nous a permis d'avoir le taux de tuples correctes. La formule est la suivante :

$$\text{Taux\_évalué} = 1 - [(\text{Nbr\_tuple\_attribut\_null}) / (\text{nbr\_tuple\_total})]$$

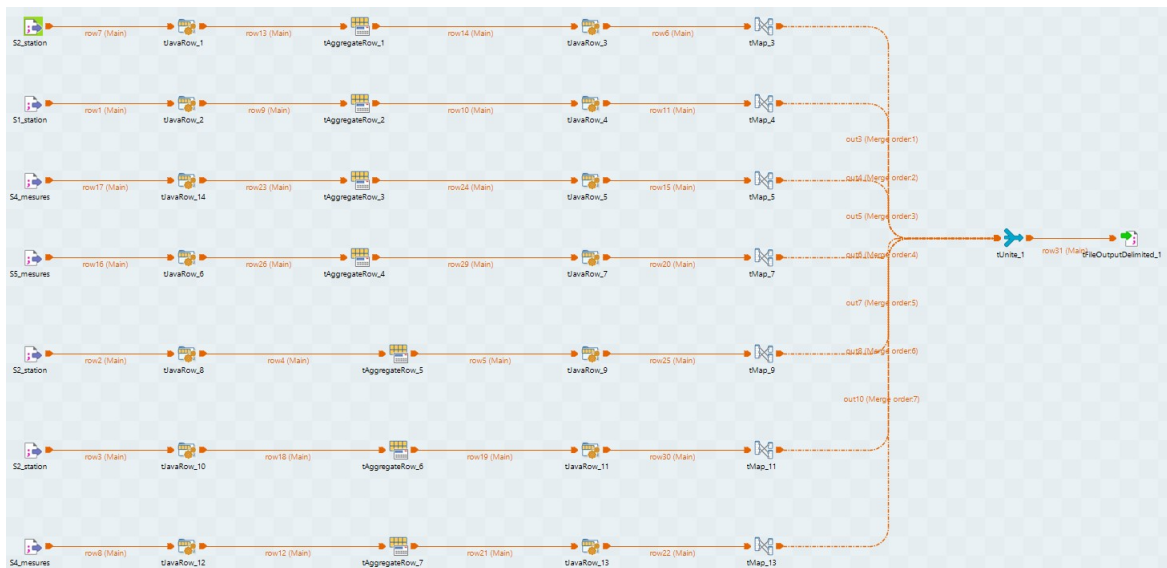


FIGURE 9 – Job d'évaluation de la complétude

#### 4.4 Détection et élimination de doublons

Au niveau des mappings, on va éliminer les doublons obtenus dans les résultats des unions et des jointures à l'intérieur des mappings.

- Dans le cas de l'union des deux station, on peut avoir des même stations mais avec des ID\_station différents vu qu'elles viennent de deux BD différentes. Donc pour détecter et supprimer les stations identiques il suffit juste de vérifier le numéro de téléphone ou bien l'adresse mail de la station).
- Nous avons suivis la métrique suivante : Tout d'abord nous avons fait l'union entre les tables S1.Station et S2.Sation puis nous avons évalué le nombre de tuple sans doublons en utilisant la formule suivante :

$$\text{Taux\_évalué} = 1 - [(\text{Nbr\_tuple\_avec\_doublons}) / (\text{nbr\_tuple\_total})]$$

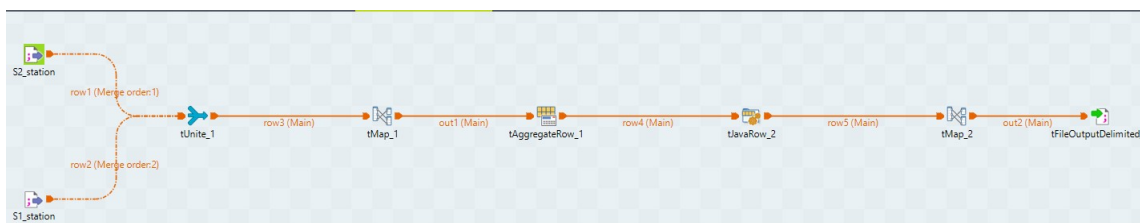


FIGURE 10 – Job de détection de doublons

- **Résumé des différentes valeurs des différentes métriques pour les facteurs de qualité :**

Nous avons stockés sous forme d'historique les pourcentages obtenus lors de l'étude des facteurs de qualité (complétude, doublons, conformité, granularité et hétérogénéité).

– La figure suivante montre l'historique de l'évaluation de la complétude.

Date	Source	Attribut	Type	Taux
27/10/2022 23:24:27	S2 Station	Numero	Complétude	50.0
27/10/2022 23:24:27	S1 Station	Numero	Complétude	71.42857
27/10/2022 23:24:27	S4	Localisation	Complétude	86.666664
27/10/2022 23:24:27	S5	Taux releve	Complétude	83.33333
27/10/2022 23:24:27	S2 Station	Ville	Complétude	83.33333
27/10/2022 23:24:27	S2 Station	Mail contact	Complétude	83.33333
27/10/2022 23:24:27	S4	Taux releve	Complétude	86.666664
27/10/2022 23:24:35	S2 Station	Numero	Complétude	50.0
27/10/2022 23:24:35	S1 Station	Numero	Complétude	71.42857
27/10/2022 23:24:35	S4	Localisation	Complétude	86.666664
27/10/2022 23:24:35	S5	Taux releve	Complétude	83.33333
27/10/2022 23:24:35	S2 Station	Ville	Complétude	83.33333
27/10/2022 23:24:35	S2 Station	Mail contact	Complétude	83.33333
27/10/2022 23:24:35	S4	Taux releve	Complétude	86.666664
28/10/2022 00:07:16	S2 Station	Numero	Complétude	50.0
28/10/2022 00:07:16	S1 Station	Numero	Complétude	71.42857
28/10/2022 00:07:17	S4	Localisation	Complétude	86.666664
28/10/2022 00:07:17	S5	Taux releve	Complétude	83.33333
28/10/2022 00:07:17	S2 Station	Ville	Complétude	83.33333

FIGURE 11 – Historique l'étude de la complétude

## 5 Amélioration de la qualité des données

### 5.1 Complétude des données :

Dans cette partie nous allons étudier la complétude de nos données sur certains attributs.

#### 5.1.1 Localisation :

Nous avons opté pour la solution qui consiste remplacer la valeur nulle par la localisation d'une mesure faite **à quelques minutes près (la différence des deux timestamp ne dépasse pas les 60min)** par le même capteur sur un autre polluant.

Si on trouve aucune localisation remplaçante, on supprime le tuple concerné de

---

notre table.

Exemple : dans la table Mesure4 on a une valeur nulle pour la localisation faite par Can45 à "02/10/2021 15 :00" sur le polluant PM10. Le but est donc de chercher dans la table une autre localisation faite par ce même capteur pour un polluant quelconque à quelques minutes près ou au même moment.

#### 5.1.2 Taux relevé :

En Essayant de chercher d'autres valeurs de taux\_relevé pour le même polluant et le même capteur avec une date à quelques minutes près d'un taux manquant sur les sources S4 et S5, on remplace la valeur nulle par un des taux trouvés, et dans le cas échéant on supprime le tuple.

#### 5.1.3 Numéro :

Au niveau de la station 1 et 2, lorsque nous détectons qu'un numéro de rue d'une station est nul, nous allons chercher la station qui a le même mail ou bien le même le numéro de téléphone que cette dernière et récupérer son Numéro de rue (Nous corrigeons par rapport aux doublons)

#### 5.1.4 Contact\_Mail :

Nous allons étudier la complétude de l'attribut Contact\_Mail au niveau de Station2 et Station1.

les mails suivent le formats suivant "xxx@airparif.fr" ou "xxx" correspond à la ville donc pour compléter les champs vides. Nous avons deux cas :

- Premier cas : Si la ville c'est **Paris**, nous allons regarder les deux derniers chiffres du code postal par exemple **75019**, ensuite nous allons suivre le format "@airparif.fr" donc la valeur nulle sera remplacé par "paris19"@airparif.fr".
- Deuxième cas : Si la ville est différente de **Paris**, nous allons récupérer la ville et faire respecter le format imposé.  
Exemple : Pour le tuple 6 nous avons Contact\_Mail = nul donc nous allons récupérer ici la ville qui est "putaux" ensuite nous allons suivre le format "@airparif.fr" donc la valeur nulle sera remplacé par "puteaux"@airparif.fr".

#### 5.1.5 Ville :

Nous allons étudier la complétude de l'attribut Ville au niveau de Station2 et Station1.

Pour cela nous allons récupérer le code postal, nous avons 3 cas, si le code postal commence par :

- 75 : On affecte la valeur **Paris** à l'attribut Ville.

- 78 : On affecte la valeur **Yvelines** à l'attribut Ville.
- 92 : On affecte la valeur **Hauts-de-seine** à l'attribut Ville.

## 5.2 Conformité à un format, une codification

Nous allons dans cette partie améliorer la qualité de nos données en respectant le critère de conformité. Ce dernier est nécessaire pour avoir des données cohérentes au niveau des sources.

### 5.2.1 Format d'un mail :

Au niveau de la station 1 et 2, nous allons mettre tous les mails sous le format suivant "xxx@airparif.fr" ou "xxx" correspond à la ville donc pour corriger le format du mail. Nous avons deux cas :

- Premier cas : Si la ville c'est **Paris**, nous allons regarder les deux derniers chiffres du code postal par exemple **75019**, ensuite nous allons suivre le format "@airparif.fr" donc la valeur nulle sera remplacé par "paris19@airparif.fr".
- Deuxième cas : Si la ville est différente de **Paris**, nous allons récupérer la ville et faire respecter le format imposé.  
Exemple : Pour le tuple 6 nous avons Contact\_Mail = nul donc nous allons récupérer ici la ville qui est "putaux" ensuite nous allons suivre le format "@airparif.fr" donc la valeur nulle sera remplacé par "puteaux@airparif.fr".

### 5.2.2 Format du numéro de téléphone :

Au niveau de la station 1 et 2, lorsque nous détectons qu'un numéro d'une station n'est pas dans le bon format, nous regardons la station qui a le même mail que cette dernière, pour prendre son numéro. (Nous corrigeons par rapport aux doublons)

### 5.2.3 Localisation :

Nous avons donc opté pour la solution qui consiste à remplacer la valeur nulle par la localisation d'une mesure de la même table faite **à quelques minutes près (la différence des deux timestamp ne dépasse pas les 60min)** par le même capteur sur un autre polluant.

Si on trouve aucune localisation remplaçante, on supprime le tuple concerné de notre table.

## 5.3 Détection et élimination de doublons :

Cette amélioration doit être exécuter après l'amélioration de la complétude et la conformité, pour qu'on puisse détecter tous les doublons possibles (Pour éviter le cas (t1(1, Null, Null, 3) et t2(2, a, b, Null)).

- Au niveau de l'union des stations 1 et 2, On va d'abord faire la jointure entre les deux tables s1 et s2 puis regarder si dans le tableau on a des doublons pour ensuite les supprimer. Nous avons comparé chaque valeur des deux tables afin de garder un seul tuple avec les valeurs les plus cohérents possibles.

Le critère pour détecter les doublons de prendre les tuples qui ont soit la même adresse mail ou le même numéro de téléphone. Pour cela nous avons utilisé le composant **tUniqRow**.

- Si par exemple, un tuple a son champ nul pour un attribut donné et l'autre tuple détecté comme son doublon a une valeur numérique non nulle pour le même attribut, nous allons garder la valeur non nulle.

Et si les deux champs sont non nuls, on garde celui qui est conforme à la règle de constitution des valeurs de l'attribut. Et si les deux ont respecté les règles de conformité, on privilégie arbitrairement la valeur du tuple provenant de la source S1.

## 5.4 Hétérogénéité des échelles, de la granularité :

### 5.4.1 Hétérogénéité des échelles :

Au niveau des sources S1.mesure, S2.mesure et S4, lorsqu'on détecte qu'une valeur de Taux\_releve est en dehors de **l'intervalle (qui a été défini pour chaque polluant comme hypothèse de base)**, on va la corriger par :

- $((\text{MAX}(\text{Intervalle})) + (\text{MIN}(\text{Intervalle}))) / 2.$

### 5.4.2 Granularité :

Au niveau de la source S5, nous avons la majorité des valeurs sont en dehors de leurs intervalle, ce qui nous permet de dire que y a un problème de granularité. Pour mettre sur la même échelle les valeurs de S5 nous avons utilisé la méthode suivante :

- $\text{Facteur} = \text{AVG}(\text{S4.mesure}) / \text{AVG}(\text{S5.mesure})$

Pour obtenir un facteur F qui sera utilisé pour la mise à l'échelle ( ou on va multiplier les valeurs de s5 par ce facteur).