# SUPERSTORE

François SOK

# 1. Collecte et préparation des données

**Légende:**

- L'exposition du problème rencontré

→ En vert, les solutions proposées

# A. Détection des anomalies



Anomalies détectées dans les colonnes suivantes:
- location_id,
- state,
- postal_code,
- region et
- col_23.

# i. Valeurs invalides

- Colonnes concernées: product_name, category, sub_category, sales_team, sales_team_manager, city, state, postal_code, region et col_23
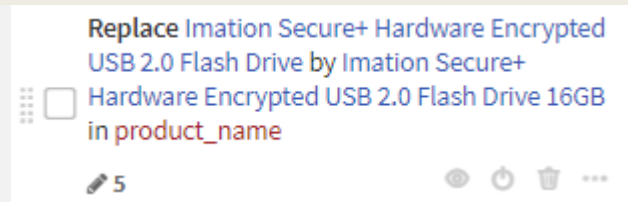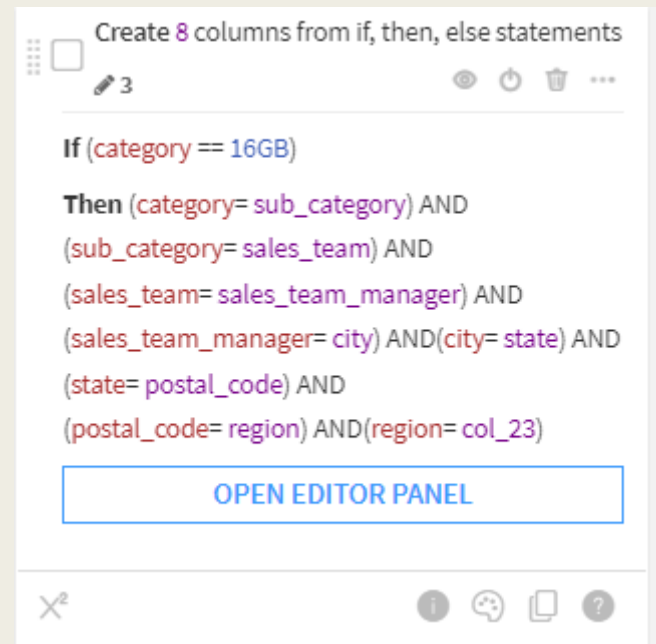  - *Concerne uniquement les lignes dont la « Category » est « 16GB ».*

# i. Valeurs invalides - suite

- Constat: il y a un décalage des colonnes.

→ product_name: **remplacer** « Imation Secure+ Hardware Encrypted USB 2.0 Flash Drive » par « Imation Secure+ Hardware Encrypted USB 2.0 Flash Drive 16GB ». En gros, on rajoute « 16GB » à la fin du texte.
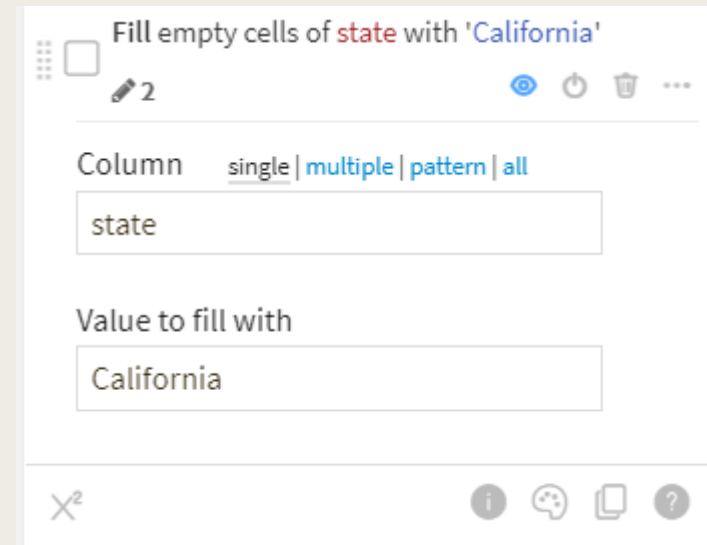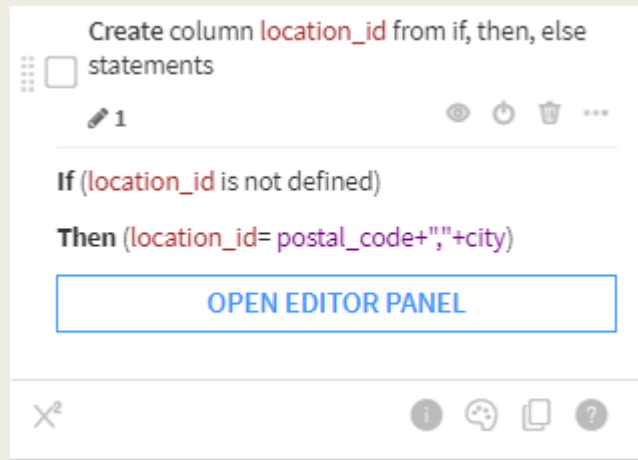
Replace Imation Secure+ Hardware Encrypted USB 2.0 Flash Drive by Imation Secure+ Hardware Encrypted USB 2.0 Flash Drive 16GB in product_name

5

→ Pour corriger ce décalage, on pourrait faire un « shift ». Seulement, la fonction n'existe pas dans Dataiku donc on utilise la méthode suivante:

Create 8 columns from if, then, else statements

3

If (category == 16GB)

Then (category= sub_category) AND
(sub_category= sales_team) AND
(sales_team= sales_team_manager) AND
(sales_team_manager= city) AND(city= state) AND
(state= postal_code) AND
(postal_code= region) AND(region= col_23)
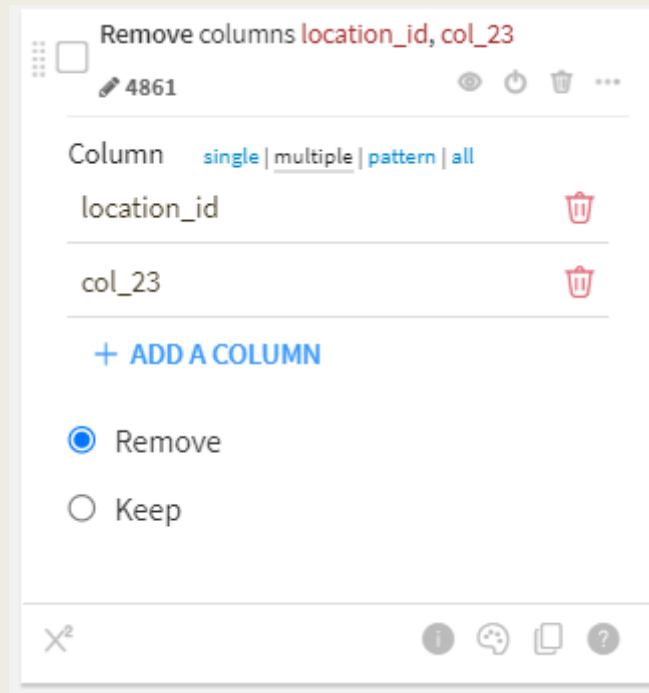
OPEN EDITOR PANEL

# ii. Valeurs manquantes

- Colonnes concernées: location_id et state

→ location_id: corrigés via la formule de la capture d'écran ci-dessous. Etape facultative (voir « iii. Doublons »

→ state: : corrigés en remplaçant par les valeurs trouvées dans « location_id » et en faisant le lien avec « postal_code »

# iii. Doublons

- Colonnes concernées: col_23 et location_id

→ Supprimer « col_23 » car a servi uniquement pour remplacer les chiffres dans « region »

→ Supprimer « location_id » car doublon avec les colonnes: city, state, postal_code et region

# Résumé des étapes dans Dataiku

# iv. Valeurs aberrantes

■ Colonnes concernées: sales et profit

■ sales et profit: les valeurs min. et max. sont bizarres, il faut investiguer.



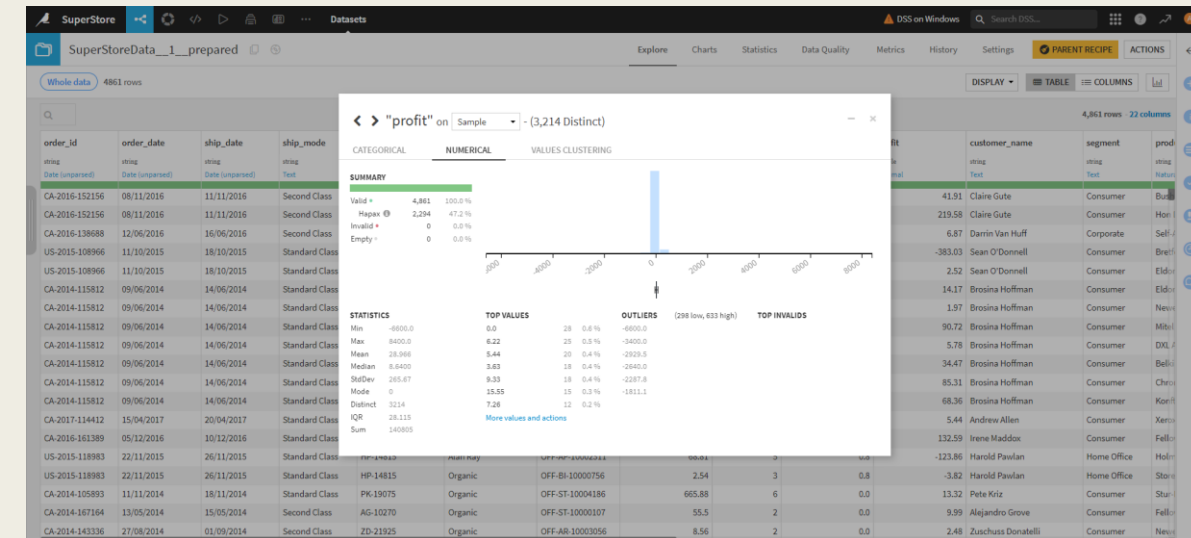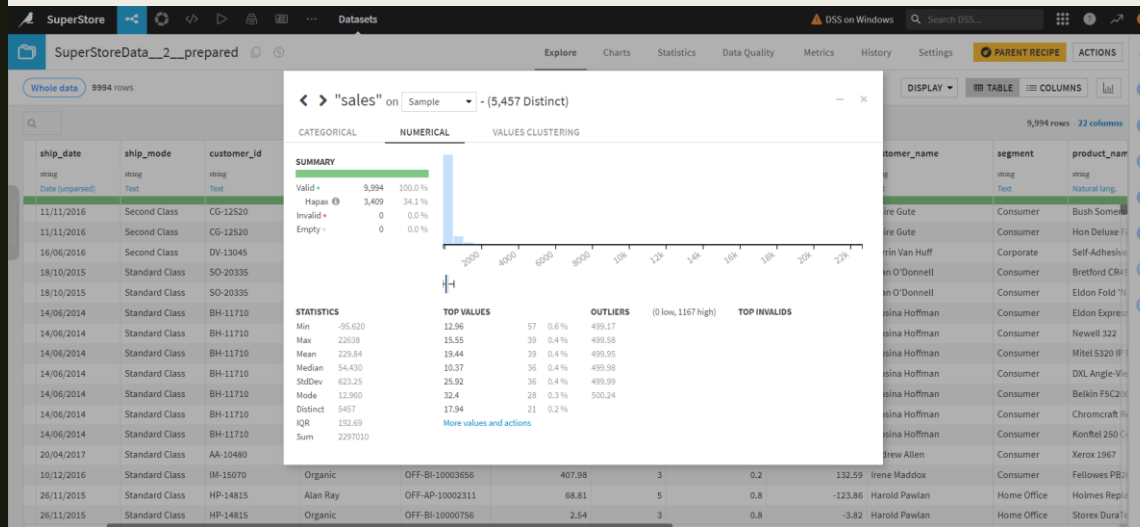→ Il y a un doublon, on le supprime rapidement.

→ Il y a une valeur aberrante dans « sales »: le « -95,62 ».

→ *On constate qu'il s'agit d'une erreur de frappe donc on enlève le signe « - ».*

# 2. EDA

Dataiku