



# Analyse Exploratoire des Données - Fashion Store Sales

## 1. Structure du fichier

Colonnes et types :

Nombre total de colonnes : 19

Nombre total de lignes : 2253

Colonne	Type estimé	Description
<b>sale_date</b>	String	Date de vente
<b>item_id</b>	Integer	Identifiant unique de l'article vendu
<b>sale_id</b>	Integer	Identifiant de la vente
<b>product_id</b>	Integer	Identifiant du produit
<b>quantity</b>	Integer	Quantité vendue
<b>original_price</b>	Float	Prix original
<b>unit_price</b>	Float	Prix unitaire après remise
<b>discount_applied</b>	Float	Montant de la remise
<b>discount_percent</b>	String	Pourcentage de remise (format XX.XX%)
<b>discounted</b>	Integer	Indicateur binaire (0/1)
<b>item_total</b>	Float	Total pour l'article
<b>channel</b>	String	Canal de vente
<b>channel_campaigns</b>	String	Campagne marketing
<b>total_amount</b>	Float	Montant total de la transaction
<b>product_name</b>	String	Nom du produit
<b>category</b>	String	Catégorie produit
<b>brand</b>	String	Marque
<b>color</b>	String	Couleur
<b>size</b>	String	Taille
<b>catalog_price</b>	Float	Prix catalogue

<b>cost_price</b>	Float	Prix de revient
<b>customer_id</b>	Integer	ID client
<b>gender</b>	String	Genre
<b>age_range</b>	String	Tranche d'âge
<b>signup_date</b>	Date	Date d'inscription
<b>first_name</b>	String	Prénom
<b>last_name</b>	String	Nom
<b>email</b>	String	Email
<b>country</b>	String	Pays

## 2. Cardinalité des colonnes principales

Données catégorielles :

category : 8 valeurs uniques

- Dresses (21.393697 %)
- Shoes (20.949845 %)
- Pants (15.623613 %)
- T-Shirts (21.837550 %)
- Sleepwear (20.195295 %)

brand : 1 valeur unique

- Tiva (100%) → Marque exclusive

gender : 1 valeur unique

- Female (100%) → Boutique féminine

country : 6 pays

- Germany (23.834887 %)
- France (22.103862 %)
- Italy (18.419885 %)
- Portugal (8.921438 %)

- Netherlands (14.469596 %)

- Spain (12.250333 %)

channel : 2 valeurs

- App Mobile (51.930759 %)

- E-commerce (48.069241 %)

age\_range : 5 tranches

- 16-25 (20.195295 %)

- 36-45 (21.171771%)

- 46-55 (19.218819 %)

- 56-65 (18.863737 %)

- 26-35 (20.550377 %)

### 3. Entités métier principales identifiées

#### ENTITÉ 1 : Produit

- **Clé :** *product\_id*
- **Attributs :** *product\_name, category, brand, color, size, catalog\_price, cost\_price*
- **Relations :** Apparaît dans plusieurs ventes

#### ENTITÉ 2 : Client

- **Clé :** *customer\_id*
- **Attributs :** *first\_name, last\_name, email, gender, age\_range, country, signup\_date*
- **Observations :** Certains clients récurrents (même *customer\_id* sur plusieurs lignes)

#### ENTITÉ 3 : Vente

- **Clé composite :** *sale\_id*
- **Attributs :** *sale\_date, channel, channel\_campaigns, total\_amount*
- **Relations :** Contient plusieurs articles (ventes multi-articles)

#### ENTITÉ 4 : Transaction

- **Clé :** *item\_id* (un par ligne de vente)
- **Attributs :** *quantity, original\_price, unit\_price, discount\_applied, item\_total*
- **Relations :** Liée à une vente et un produit

## 4. Redondances identifiées

### Redondance 1 : Prix tripliqués

original\_price, unit\_price, catalog\_price

Analyse :

- original\_price = Prix avant remise
- unit\_price = Prix après remise (original\_price - discount\_applied)
- catalog\_price = Prix catalogue (souvent égal à original\_price)

Problème : catalog\_price semble redondant avec original\_price dans 95% des cas

### Redondance 2 : Calculs dérivés

discount\_percent peut être calculé à partir de discount\_applied / original\_price

discounted (0/1) peut être déduit de discount\_applied > 0

item\_total = unit\_price × quantity (cohérent)

### Redondance 3 : Relations prix

Pour les ventes sans remise :

- original\_price = unit\_price = catalog\_price
- discount\_applied = 0
- discount\_percent = "0.00%"
- discounted = 0

## 5. Anomalies potentielles

### Anomalie 1 : Valeurs manquantes

first\_name : 5.14 % manquants

last\_name : 2.84 % manquants

email : 9.94 % manquants

total\_amount : 9.98 % manquants

### Anomalie 2 : Incohérences temporelles

signup\_date POSTERIEURE à sale\_date dans 8.21% des cas

Ex: Vente 2025-04-06 mais inscription 2025-05-21

### Anomalie 3 : Discounts incohérents

Cas où  $\text{discount\_percent} \neq (\text{discount\_applied} / \text{original\_price}) \times 100$

Ex: original\_price=66.87, discount\_applied=6.69, discount\_percent="10.00%"

$66.87 \times 0.10 = 6.687$  (arrondi à 6.69) → OK

Mais d'autres cas montrent des écarts dus aux arrondis

## 6. Relations métier-clés

### 7. Modélisation

Normalisation :

**Dim\_channel** (

channel\_id,

    channel\_name

)

**Dim\_campaign** (

campaign\_id,

    #channel\_id,

    campaign\_name,

);

**Dim\_category** (

category\_id,

    category\_name

);

**Dim\_brand** (

brand\_id,

    brand\_name

);

**Dim\_color** (

color\_id,

    color\_name

);

**Dim\_size** (

size\_id,

    size\_label

);

**Dim\_country** (

country\_id,

    country\_name

);

**Dim\_customer** (

customer\_id,

    first\_name,

    last\_name,

    email,

    gender,

    age\_range,

    signup\_date,

    #country\_id

);

**Dim\_product** (

product\_id,

    product\_name,

```
#category_id,  
#brand_id,  
#color_id,  
#size_id ,  
cost_price,  
original_price  
)
```

```
Fact_sale (  
    sale_id ,  
    sale_date,  
    total_amount,  
    #customer_id ,  
    #campaign_id  
)
```

```
Fact_sale_item (  
    item_id,  
    #sale_id ,  
    #product_id ,  
    quantity,  
    discount_percent  
)
```