

# Du RAW au schéma en étoile (DW) : génération, staging, dimensions, fait, agrégations

Programme Data Engineer

ECE

23 octobre 2025

# Objectifs du Mini-TP et plan

## Objectifs :

- Simuler des données **RAW** (clients, produits, commandes).
- Nettoyer/normaliser (*staging*).
- Construire des **dimensions** (date, produit, client) avec SK.
- Construire la **table de faits** fact\_sales et réaliser des agrégations.
- Produire un **graphique** des ventes mensuelles.

## Plan (12 diapositives) : Contexte

→ RAW → Staging → Dimensions → Fait → Persist/Analytics → Graphique → TP.

# Contexte et setup (imports, reproductibilité)

```
1 import numpy as np, pandas as pd
2 from datetime import datetime, timedelta
3 from pathlib import Path
4 import matplotlib.pyplot as plt
5
6 # Reproductibilite et sortie
7 rng = np.random.default_rng(42)
8 out_dir = Path("dw_example"); out_dir.mkdir(exist_ok=True)
```

# Génération RAW : clients et produits

```
1 # Clients (avec un peu de "saleté")
2 countries = ["FR","DE","ES","IT","UK","US"]
3 raw_customers = pd.DataFrame({
4     "customer_id": range(1,121),
5     "name": [f"Cust {i}" for i in range(1,121)],
6     "email": [f"user{i}@Example.com" for i in range(1,121)],
7     "country": rng.choice(countries, 120)
8 })
9
10 # Produits
11 categories = ["Electronics","Home","Sports","Toys"]
12 raw_products = pd.DataFrame({
13     "product_id": range(1,61),
14     "product_name": [f"Item {i}" for i in range(1,61)],
15     "category": rng.choice(categories, 60),
16     "brand": rng.choice(["Acme","Globex","Umbrella","Initech","Soylent
17     "], 60),
18     "unit_price": np.round(rng.uniform(10,300,60),2)
```

# Génération RAW : lignes de commande (order lines)

```
1 start, end = datetime(2025,1,1), datetime(2025,9,30)
2 N = 2500
3 raw_orders = pd.DataFrame({
4     "order_id": 10000 + rng.integers(0,10000, N),
5     "order_line_id": range(1, N+1),
6     "order_ts": [start + timedelta(days=int(rng.integers(0,(end-start)
7         .days+1))) for _ in range(N)],
8     "customer_id": rng.integers(1,121, N),
9     "product_id": rng.integers(1,61, N),
10    "quantity": rng.integers(1,6, N)
11 })
12 # Persist RAW (optionnel)
13 raw_customers.to_csv(out_dir/"raw_customers.csv", index=False)
14 raw_products.to_csv(out_dir/"raw_products.csv", index=False)
15 raw_orders.to_csv(out_dir/"raw_orders.csv", index=False)
```

# Staging : nettoyage et standardisation

```
1 stg_customers = raw_customers.assign(  
2     email = raw_customers["email"].str.strip().str.lower(),  
3     country = raw_customers["country"].str.upper()  
4 )  
5  
6 stg_products = raw_products.assign(  
7     category = raw_products["category"].str.title(),  
8     brand     = raw_products["brand"].str.title()  
9 )  
10  
11 stg_orders = raw_orders.assign(order_ts = pd.to_datetime(raw_orders[  
    "order_ts"])))
```

## Dimensions : dim\_date, dim\_product, dim\_customer

```
1 # dim_date (SK sur calendrier des commandes)
2 dim_date = (pd.DataFrame({"date": pd.to_datetime(stg_orders["
    order_ts"].unique())})
3     .sort_values("date").reset_index(drop=True))
4 dim_date["date_sk"] = range(1, len(dim_date)+1)
5 dim_date = dim_date.assign(
6     day = dim_date["date"].dt.day,
7     month = dim_date["date"].dt.month,
8     quarter = dim_date["date"].dt.quarter,
9     year = dim_date["date"].dt.year
10 )
11
12 # dim_product & dim_customer (SK)
13 dim_product = stg_products.drop_duplicates("product_id").
14     reset_index(drop=True)
15 dim_product["product_sk"] = range(1, len(dim_product)+1)
16
17 dim_customer = stg_customers.drop_duplicates("customer_id").
```

## Table de faits : fact\_sales (quantité, montant, coût, marge)

```
1 fact = (stg_orders
2   .merge(dim_customer[["customer_id","customer_sk"]], on="customer_id"
3   .merge(dim_product[["product_id","product_sk","unit_price","
4     unit_cost"]], on="product_id")
5   .merge(dim_date[["date","date_sk"]], left_on="order_ts", right_on="
6     date"))
7
8 fact["amount"] = fact["quantity"] * fact["unit_price"]
9 fact["cost"]   = fact["quantity"] * fact["unit_cost"]
10 fact["margin"] = fact["amount"] - fact["cost"]
11
12 fact_sales = fact[[
13   "order_id","order_line_id","date_sk","product_sk","customer_sk",
14   "quantity","unit_price","amount","cost","margin"
15 ]].sort_values(["order_id","order_line_id"])
```



# Persistence du schéma en étoile (CSV) et aperçu

```
1 # Sauvegarde des dimensions et du fait
2 dim_date.to_csv(out_dir/"dim_date.csv", index=False)
3 dim_product[["product_sk","product_id","product_name","category",
4              "brand",
5              "unit_price","unit_cost"]].to_csv(out_dir/"dim_product.
6              csv", index=False)
7 dim_customer[["customer_sk","customer_id","name","email","country"]]
8 .to_csv(out_dir/"dim_customer.csv", index=False)
9 fact_sales.to_csv(out_dir/"fact_sales.csv", index=False)
```

**Fichiers produits :** dim\_date.csv, dim\_product.csv, dim\_customer.csv,  
fact\_sales.csv

## Agrégations : ventes mensuelles, top 10 produits

```
1 # Ventes mensuelles
2 sales_by_month = (fact_sales
3     .merge(dim_date[["date_sk", "month", "year"]], on="date_sk")
4     .groupby(["year", "month"], as_index=False)["amount"].sum()
5     .sort_values(["year", "month"]))
6
7 # Top 10 produits par CA
8 top_products = (fact_sales
9     .merge(dim_product[["product_sk", "product_name"]], on="product_sk"
10         )
11     .groupby(["product_sk", "product_name"], as_index=False)["amount"].
12         sum()
13     .sort_values("amount", ascending=False).head(10))
```

## Graphique : chiffre d'affaires mensuel (matplotlib)

```
1 import matplotlib.pyplot as plt
2 x = range(len(sales_by_month))
3 labels = [f"{int(y)}-{int(m):02d}" for y,m in zip(sales_by_month["
    year"], sales_by_month["month"])]
4 plt.figure(); plt.plot(x, sales_by_month["amount"], marker="o")
5 plt.title("Chiffre d'affaires mensuel (exemple)")
6 plt.xlabel("Mois"); plt.ylabel("Montant")
7 plt.xticks(x[::max(1, len(x)//10)], [labels[i] for i in range(0,len(
    labels), max(1,len(x)//10))],
8         rotation=45, ha="right")
9 plt.tight_layout()
10 plt.savefig(out_dir/"monthly_sales.png")
11 plt.show()
```

# TP : exercices proposés

- Ajouter une dimension **Promotion** et lier au fait.
- Créer une table **fact\_inventory** (stock journalier) et comparer CA vs stock.
- Mettre en place une **dimension SCD2** sur dim\_customer (valid\_from/to, is\_current).
- Construire une vue matérialisée des **ventes hebdomadaires**.
- Ajouter des **règles DQ** (emails valides, valeurs non nulles) et un rapport d'anomalies.

*Astuce* : séparer RAW, STAGING, DIM, FACT en modules pour réutiliser le code.

# Exemple flocon : dimension Géographie

**Objectif** : réduire la redondance des attributs de localisation (ville, région, pays) en normalisant la dimension Magasin. **Tables proposées** :

- `dim_magasin(sk_magasin, id_ville, nom_magasin, type_magasin, ...)`
- `dim_ville(id_ville, ville, id_region)`
- `dim_region(id_region, region, id_pays)`
- `dim_pays(id_pays, pays, code_iso2)`

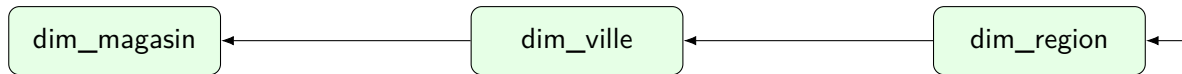
**Avantages** : référentiels centralisés (pays/régions), mise à jour cohérente, moindre stockage.

**Inconvénients** : plus de jointures pour la BI, modèle moins intuitif qu'une étoile plate.

# DDL (extrait) — flocon Géographie

```
1 CREATE TABLE dim_pays (  
2     id_pays INT PRIMARY KEY,  
3     pays VARCHAR(80) NOT NULL,  
4     code_iso2 CHAR(2) UNIQUE  
5 );  
6  
7 CREATE TABLE dim_region (  
8     id_region INT PRIMARY KEY,  
9     region VARCHAR(80) NOT NULL,  
10    id_pays INT NOT NULL REFERENCES dim_pays(id_pays)  
11 );  
12  
13 CREATE TABLE dim_ville (  
14     id_ville INT PRIMARY KEY,  
15     ville VARCHAR(80) NOT NULL,  
16     id_region INT NOT NULL REFERENCES dim_region(id_region)  
17 );  
18
```

# Schéma en flocon — vue Géographie



*Remarque :* pour la consommation BI, on peut matérialiser une vue aplatie `dim_magasin_flat` (magasin+ville+region+pays) pour limiter les jointures.