

# Homework Lab

DP-203 Microsoft Azure Data  
engineering

Salaheddine HOUKMI  
Clément ROUANET

# Sujet

## Projet :

Notre initiative vise à implémenter et optimiser des solutions de stockage et de traitement de données sur le cloud, axées sur un objectif principal : **Trouver le "produit gagnant" et les emplacements stratégiques pour les publicités à partir du dataset de la plateforme de e-commerce brésilienne Olist.**

Pour démarrer un business en ligne, surtout dans le e-commerce, on commence par chercher une niche rentable en essayant de vendre plusieurs produits et en observant les résultats de chacun. Une fois le produit gagnant identifié, on lance nos publicités dans une région bien déterminée, en procédant également par essais. Toutes ces démarches nécessitent un grand budget.

## Objectif :

Notre objectif dans ce projet est de créer une solution de machine learning capable de prédire le produit gagnant en se basant sur plusieurs indicateurs tels que les ventes par catégories, les moyens de paiement, les avis et retours clients, etc., ainsi que de déterminer la région cible pour nos publicités, le tout sans nécessiter un grand budget. De plus, pour la création de notre boutique e-commerce, nous devons spécifier les moyens de paiement, ce qui représente un grand défi pour les professionnels du domaine. C'est pourquoi nous avons également décidé d'identifier les moyens de paiement les plus utilisés par région.

## Jeux de données :

Le périmètre de notre projet sera le marché brésilien. L'ensemble des données a été généreusement fourni par Olist, le plus grand marché en ligne du Brésil. Olist connecte les petites entreprises de tout le pays aux grandes chaînes de distribution sans complications, et avec un seul contrat. Ces commerçants peuvent vendre leurs produits via la boutique Olist et les expédier directement aux clients en utilisant les partenaires logistiques d'Olist.

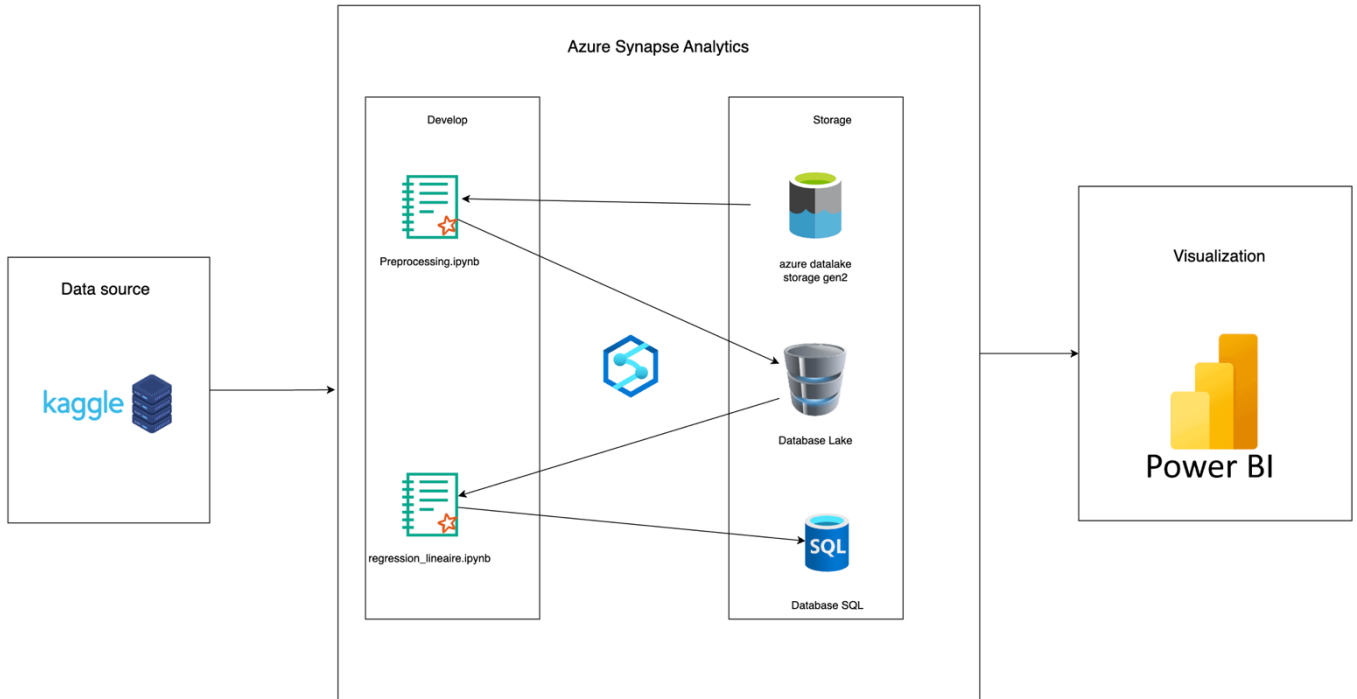
(Source de données : <https://www.kaggle.com/datasetsolistbr/brazilian-ecommerce>)

Les données sont divisées en plusieurs ensembles de données, pour une meilleure compréhension, voici un modèle conceptuel de données :



## Architecture et traitements :

### 1. Schéma d'architecture



### 2. Collecte et stockage des données ( Dataset ) :

Nous opterons pour Azure, un fournisseur de services cloud renommé, comme notre infrastructure principale et Azure synapse comme service analytique pour nos traitements. Les données seront récupérées et stockées dans un conteneur sur data lake storage gen2 d'azure, garantissant ainsi une gestion sécurisée et évolutive de nos jeux de données.

Notre gestion de stockage sera comme suite :

- *Dossier raw\_data* : Dossier contenant les données brutes sous forme de fichiers CSV.
- *Dossier cleaned\_data ( fichiers parquet )* : Dossier contenant les données nettoyés ( suppression des lignes dupliquées ,traitement de données aberrantes et nulles,) et transformées ( traduction de la data du portugais à l'anglais, changement des types de variables , ajout de colonnes opérationnels ,changement du type des fichiers du csv en parquet ) .
- *Dossier Linear\_regression* : Dossier contenant les données de prédiction issues de la régression linéaire.

▲ Azure Data Lake Storage Gen22

▲ cloud-workspace (Primary - datalak...

▲ projetcloud (Primary)

▲ (Attached Containers)

Nom	^	Dernière modification	Type de contenu
cleaned_data		18/06/2024 12:55:00	Folder
linear_regression		19/06/2024 18:43:15	Folder
raw_data		18/06/2024 12:54:53	Folder
synapse		18/06/2024 12:49:32	Folder

▲ Azure Data Lake Storage Gen22

▲ cloud-workspace (Primary - datalak...

▲ projetcloud (Primary)

▲ (Attached Containers)

Nom	^	Dernière modification
customers.csv		18/06/2024 13:36:20
geolocation.csv		18/06/2024 13:36:37
order_items.csv		18/06/2024 13:36:23
order_payments.csv		18/06/2024 13:36:24
order_reviews.csv		18/06/2024 13:36:32
orders.csv		18/06/2024 13:36:33
product_category_name_translation.csv		18/06/2024 13:36:34
products.csv		19/06/2024 16:53:40
sellers.csv		18/06/2024 13:36:35
states_name.csv		18/06/2024 14:06:07

← → ∨ ↑

projetcloud > cleaned\_data

Nom	^	Dernière modification	Type de contenu
customers		19/06/2024 16:56:11	Folder
geolocation		19/06/2024 16:56:25	Folder
order_items		19/06/2024 16:56:16	Folder
order_payments		19/06/2024 16:56:17	Folder
order_reviews		19/06/2024 16:56:19	Folder
orders		19/06/2024 16:56:14	Folder
products		19/06/2024 16:56:21	Folder
products_translated		19/06/2024 16:56:22	Folder
sellers		19/06/2024 16:56:13	Folder
states_name		19/06/2024 16:56:27	Folder

### 1. Traitement de données :

### 1.1. Nettoyage et transformation de données ( Preprocessing ) :

Le traitement de données est effectuée par le framework PySpark.

On a procédé comme suite :

- Nettoyage de données : Traitements des valeurs aberrantes et nulles , suppression des lignes dupliquées
- Transformation de données : Traduction de la data du portugais à l'anglais, changement des types de variables , ajout de colonnes opérationnels ,changement du type des fichiers du csv au parquet.

```
1 # Suppression des valeurs nulles de "order_reviews"
2 dataframes["order_reviews"] = dataframes["order_reviews"] \
3     .filter(F.col("review_creation_date").isNotNull() & F.col("review_answer_timestamp").
4
5 # Vérification que toutes les valeurs nulles de order_reviews ont été supprimées
6 null_counts = count_nulls(dataframes["order_reviews"])
7 print(f"Null values in order_reviews DataFrame :")
8 null_counts.show(truncate=False)
```

✓ - Commande exécutée en 1 s 59 ms le 4:56:09 PM, 6/19/24

Null values in order reviews DataFrame :

```

1 # Remplacement des valeurs nulles de "products" : String -> "Other" & Int -> -1
2 dataframes["products"] = dataframes["products"] \
3     .withColumn("product_category_name", F.when(F.col("product_category_name").isNull(), "Other"))
4     .withColumn("product_name_length", F.when(F.col("product_name_length").isNull(), -1))
5     .withColumn("product_description_length", F.when(F.col("product_description_length").isNull(), -1))
6     .withColumn("product_photos_qty", F.when(F.col("product_photos_qty").isNull(), -1))
7     .withColumn("product_weight_g", F.when(F.col("product_weight_g").isNull(), -1))
8     .withColumn("product_length_cm", F.when(F.col("product_length_cm").isNull(), -1))
9     .withColumn("product_height_cm", F.when(F.col("product_height_cm").isNull(), -1))
10    .withColumn("product_width_cm", F.when(F.col("product_width_cm").isNull(), -1))
11
12 # Vérification que toutes les valeurs nulles de products ont été remplacée
13 null_counts = count_nulls(dataframes["products"])
14 print(f"Null values in products DataFrame :")
15 null_counts.show(truncate=False)

```

# Ajout champs

```
1 ataframes["order_items"] = dataframes["order_items"] \
2     .withColumn("total_items_value", F.col("price") * F.col("order_item_id")) \
3     .withColumn("total_freight_value", F.col("freight_value") * F.col("order_item_id")) \
4     .withColumn("total_order_value", F.col("total_items_value") + F.col("total_freight_val
5
6 ataframes["order_items"].show(10)
```

✓ - Commande exécutée en 501 ms le 4:56:11 PM, 6/19/24

---

```

1  # Jointure entre products et products_translated sur la colonne product_category_name
2  joined_df = dataframes["products"].alias("p") \
3      .join(
4          dataframes["products_translated"].alias("pt"),
5          F.col("p.product_category_name") == F.col("pt.product_category_name"),
6          "left"
7      )
8
9  # Suppression du champs avec le nom du produit en portugais
10 joined_df = joined_df.drop("product_category_name")
11
12 # Renommage de la colonne product_category_name_english
13 joined_df = joined_df.withColumnRenamed("product_category_name_english", "product_category_name")
14
15 # Réagencement des colonnes du DataFrame
16 new_columns = [joined_df.columns[0]] + ["product_category_name"] + [col for col in joined_df.columns if col != "product_category_name"]
17 dataframes["products"] = joined_df.select(new_columns)
18
19 # Mettre à jour le DataFrame dans le dictionnaire
20 dataframes["products"].show(5)
21

```

✓ - Commande exécutée en 1 s 106 ms le 4:56:00 PM, 6/19/24

## 1.2. Machine learning ( régression linéaire ) :

### 1.2.1. Traitements

Afin de répondre à notre sujet et d'arriver à notre objectif de trouver le niche gagnant des produits pour le e-commerce, nous avons opté pour un apprentissage supervisé qui se représente dans la régression linéaire pour prédire les ventes futures de chaque catégorie de produits dans les années 2019, 2020 et 2021.

# Régression linéaire



```

1  # Regression linéaire : prévision des prix
2  schema = StructType([
3      StructField("category", StringType(), False),
4      StructField("sales", IntegerType(), False),
5      StructField("year", IntegerType(), False),
6      StructField("month", IntegerType(), False),
7      StructField("date", TimestampType(), True),
8  ])
9
10 price_predictions = spark.createDataFrame(sc.emptyRDD(), schema)
11
12 price_predictions.show()

```

[43] ✓ - Commande exécutée en 2 s 776 ms le 9:26:47 PM, 6/19/24

```

...
+-----+-----+-----+-----+-----+
|category|sales|year|month|date|
+-----+-----+-----+-----+-----+

```

```

2 def regression_lineaire(df_train, future_date, category) :
3     df_train = df_train.groupBy("year", "month") \
4         .agg(F.sum("sales").alias("sales"))
5
6     # Assembler les fonctionnalités
7     assembler = VectorAssembler(
8         inputCols=["year", "month"],
9         outputCol="features"
10    )
11
12    # Transformer les données en utilisant l'assembler
13    sales_data = assembler.transform(df_train)
14    sales_data = sales_data.select("features", "sales")
15
16    # Créer le modèle de régression linéaire
17    lr = LinearRegression(featuresCol="features", labelCol="sales")
18
19    # Ajuster le modèle aux données d'entraînement
20    lr_model = lr.fit(sales_data)
21
22    # Transformer les données en utilisant l'assembler
23    df_predicted = spark.createDataFrame(future_date, ["year", "month"])
24    df_predicted = assembler.transform(df_predicted).select("features")
25
26    # Prédire les prix pour les dates futures
27    df_predicted = lr_model.transform(df_predicted)
28
29    # Assemblage des DataFrames train et predicted en RDD
30    rdd = sc.parallelize(sales_data.union(df_predicted).rdd.collect()) \
31        .map(lambda row: Row(features=row.features.toArray().tolist(), sales=row.sales)) \
32        .map(lambda x: (int(x.features[0]), int(x.features[1]), x.sales))
33
34    # Créer un DataFrame à partir de l'RDD avec les noms de colonnes appropriés
35    df = rdd.toDF(["year", "month", "sales"]) \
36        .withColumn("category", F.lit(category)) \
37        .withColumn("date", F.to_date(F.concat(F.col("year"), F.lit("-"), F.col("month"), F.lit("-01")))) \
38        .select("category", "sales", "year", "month", "date")
39
40    return df

```

De plus, pour avoir une bonne visualisation et faciliter l'analyse des données nous avons généré également en plus, un dataframe avec les 10 meilleures catégories prédites au niveau de ventes.

```

1  n = 10
2
3  # Récupération de la plus grande date de prédiction
4  max_date_row = price_predictions.agg(F.max("date")).head()
5  max_date = max_date_row[0]
6
7  # Récupération des n catégories ayant les meilleures ventes
8  best_categories_predicted = price_predictions.select("category") \
9      .filter(F.col("date") == max_date) \
10     .orderBy(F.col("sales").desc()) \
11     .head(n)
12
13  # Mise sous forme de liste
14  best_categories_predicted = [c.category for c in best_categories_predicted]
15  best_categories_predicted

```

✓ - Commande exécutée en 51 s 634 ms le 9:36:28 PM, 6/19/24

```

['health_beauty',
'watches_gifts',
'bed_bath_table',
'housewares',
'computers_accessories',
'sports_leisure',
'furniture_decor',
'auto',
'baby',
'telephony']

```

```

1  # Création du dataframe avec uniquement les n meilleures catégories
2  price_predictions_top_category = price_predictions \
3      .filter(F.col("category").isin(best_categories_predicted))
4
5  price_predictions_top_category.show()

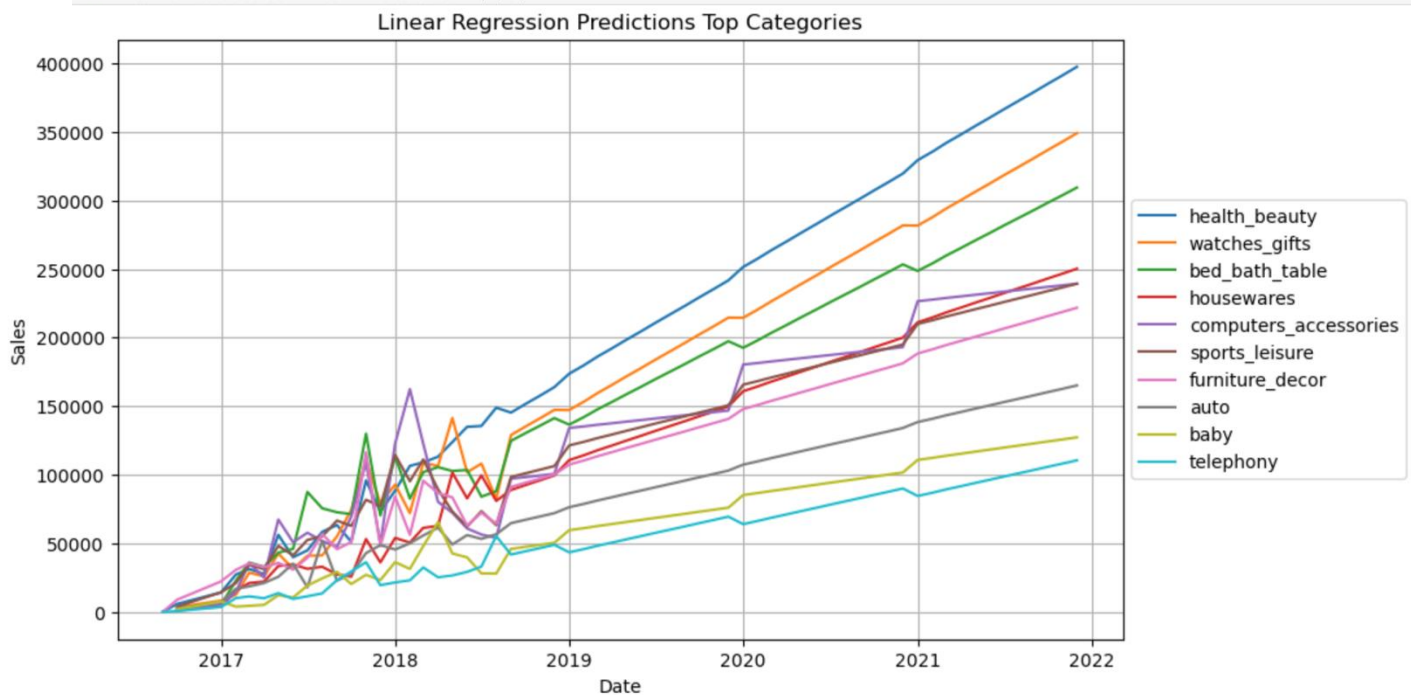
```



### 1.2.2. Représentation graphique des prédictions. ( Matplotlib )

```
1 # Initialiser la figure
2 plt.figure(figsize=(10, 6))
3
4 for category in best_categories_predicted :
5     # Choisir la bonne catégorie et trier les données par date
6     predictions_pd = price_predictions_top_category.filter(F.col("category") == category).orderBy("date").toPandas()
7
8     # Tracer les données
9     if not predictions_pd.empty:
10         plt.plot(predictions_pd['date'], predictions_pd['sales'], label=category)
11
12
13 plt.xlabel('Date')
14 plt.ylabel('Sales')
15 plt.title('Linear Regression Predictions Top Categories')
16 plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
17 plt.grid(True)
18 plt.show()
```

✓ - Commande exécutée en 12 s 827 ms le 9:37:58 PM, 6/19/24



### 1.3. Mise en place d'un Datwarehouse.

Après le traitement de la machine learning, tout les données ainsi que les données de prédictions sont sauvegardés dans la zone de staging ( data lake storage gen2 ) sous forme de fichiers parquet.

```
1 # Chemin de destination dans ADLS Gen2
2 predictions_path = "abfss://projetcloud@datalakecloud.dfs.core.windows.net/linear_regression/price_predictions"
3 top_category_predictions_path = "abfss://projetcloud@datalakecloud.dfs.core.windows.net/linear_regression/price_predictions_top_category"
4
5 # Écrire les DataFrame au format Parquet
6 price_predictions.write.mode("overwrite").parquet(predictions_path)
7 price_predictions_top_category.write.mode("overwrite").parquet(top_category_predictions_path)
```

Ensuite, pour avoir une bonne structuration de données, nous avons créé un datawarehouse ( base de données SQL ) qu'on a alimenté par nos fichiers parquets groupés par catégories sous forme de tables externes (dbo.customers, etc ).

```
1 SELECT TOP (100) [customer_id]
2 , [customer_unique_id]
3 , [customer_zip_code_prefix]
4 , [customer_city]
5 , [customer_state]
6 FROM [dbo].[customers]
```

Résultats Messages

Afficher

Table

Graphique

↳ Exporter des résultats

Rechercher

customer_id	customer_unique_id	customer_zip_code_...	customer_city	customer_state
06b8999e2fba1a1fbc88172c00ba...	861eff4711a542e4b93843c6dd7f...	14409	franca	SP
18955e83d337fd6b2def6b18a42...	290c77bc529b7ac935b93aa66c3...	9790	sao bernardo d...	SP
4e7b3e00288586ebd08712fdd03...	060e732b5b29e8181a18229c7b0...	1151	sao paulo	SP
b2b6027bc5c5109e529d4dc6358...	259dac757896d24d7702b9acbbf...	8775	mogi das cruze...	SP

Data

Workspace

Linked

Filtrer les ressources par nom

Database\_ECommerce (SQL)

Tables externes

dbo.customers

Colonnes

customer\_id (varchar(256), null)

customer\_unique\_id (varchar(256), null)

customer\_zip\_code\_prefix (int, null)

customer\_city (varchar(256), null)

customer\_state (varchar(256), null)

dbo.geolocation

dbo.order\_items

dbo.order\_payments

dbo.order\_reviews

dbo.orders

dbo.price\_predictions

dbo.price\_predictions\_top\_categories

dbo.products

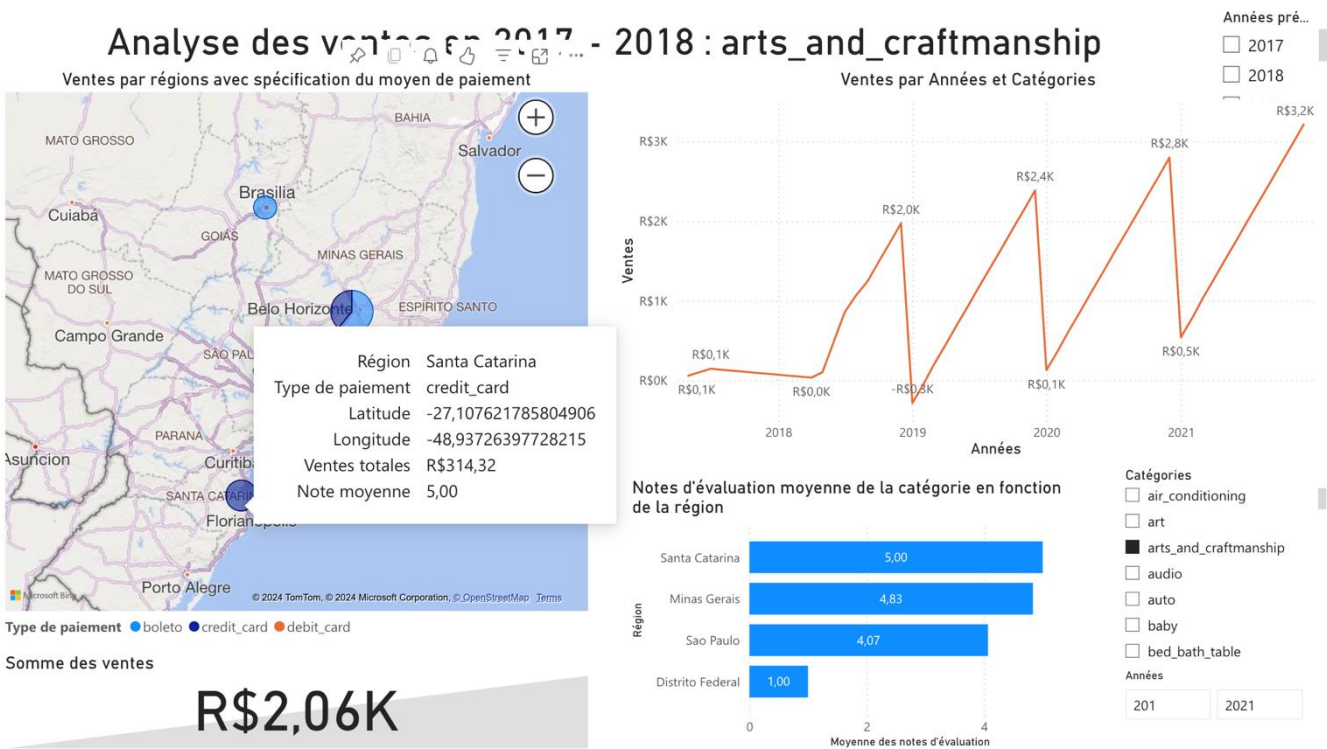
dbo.products\_translated

dbo.sellers

dbo.state\_name

À partir des données stockés dans le datawarehouse sur azure synapse, nous avons créé un rapport PowerBI, pour analyser les résultats de nos traitements sous formes de plusieurs visualisations :

- Carte : représente les régions avec les grandes ventes par catégories avec spécification du moyen de paiement.
- Courbe : représente les prédictions de ventes par année et catégorie
- Box KPI : représente la somme des vente par catégorie.
- Histogramme : représente les notes d'évaluation moyenne de la catégorie en fonction de région.
- Filtres : catégories et années.
- Titres dynamiques et visualisation dynamique.



Conclusion :

Nous avons conclu à partir de cette analyse que la catégorie "health\_beauty" est la plus rentable et présente un grand potentiel pour l'avenir, avec un gain de 0,40M \$R en décembre 2021, soit à peu près la moitié des gains cumulés de 2017 à 2019. Le moyen de paiement le plus utilisé est la carte de crédit, donc nous opterons pour une passerelle de paiement de type Stripe ou similaire pour notre boutique. Les produits de cette catégorie ont de bonnes notes, supérieures à 4/5, et sont vendus dans toutes les régions du Brésil, avec une majorité des ventes à São Paulo.