# Projet 11

# Réalisez un traitement dans un environnement Big Data sur le Cloud

**Zaccaria Amillou**

# Sommaire

- Contexte

- Données

- EDA

- Modélisation

- RGPD

- Cloud

- Conclusions

# Contexte

Entreprise



Objectif

Mise en place d'une traitement des données sur le cloud

# Données

Fichier Data



Exemple Images

Structure Test

```
        tree -L 1
    .
    ├── apple_6
    ├── apple_braeburn_1
    ├── apple_crimson_snow_1
    ├── apple_golden_1
    ├── apple_golden_2
    ├── apple_golden_3
    ├── apple_granny_smith_1
    ├── apple_hit_1
    ├── apple_pink_lady_1
    ├── apple_red_1
    ├── apple_red_2
    ├── apple_red_3
    ├── apple_red_delicios_1
    ├── apple_red_yellow_1
    ├── apple_rotten_1
    ├── cabbage_white_1
    ├── carrot_1
    ├── cucumber_1
    ├── cucumber_3
    ├── eggplant_violet_1
    ├── pear_1
    ├── pear_3
    ├── zucchini_1
    └── zucchini_dark_1

    25 directories, 0 files
```
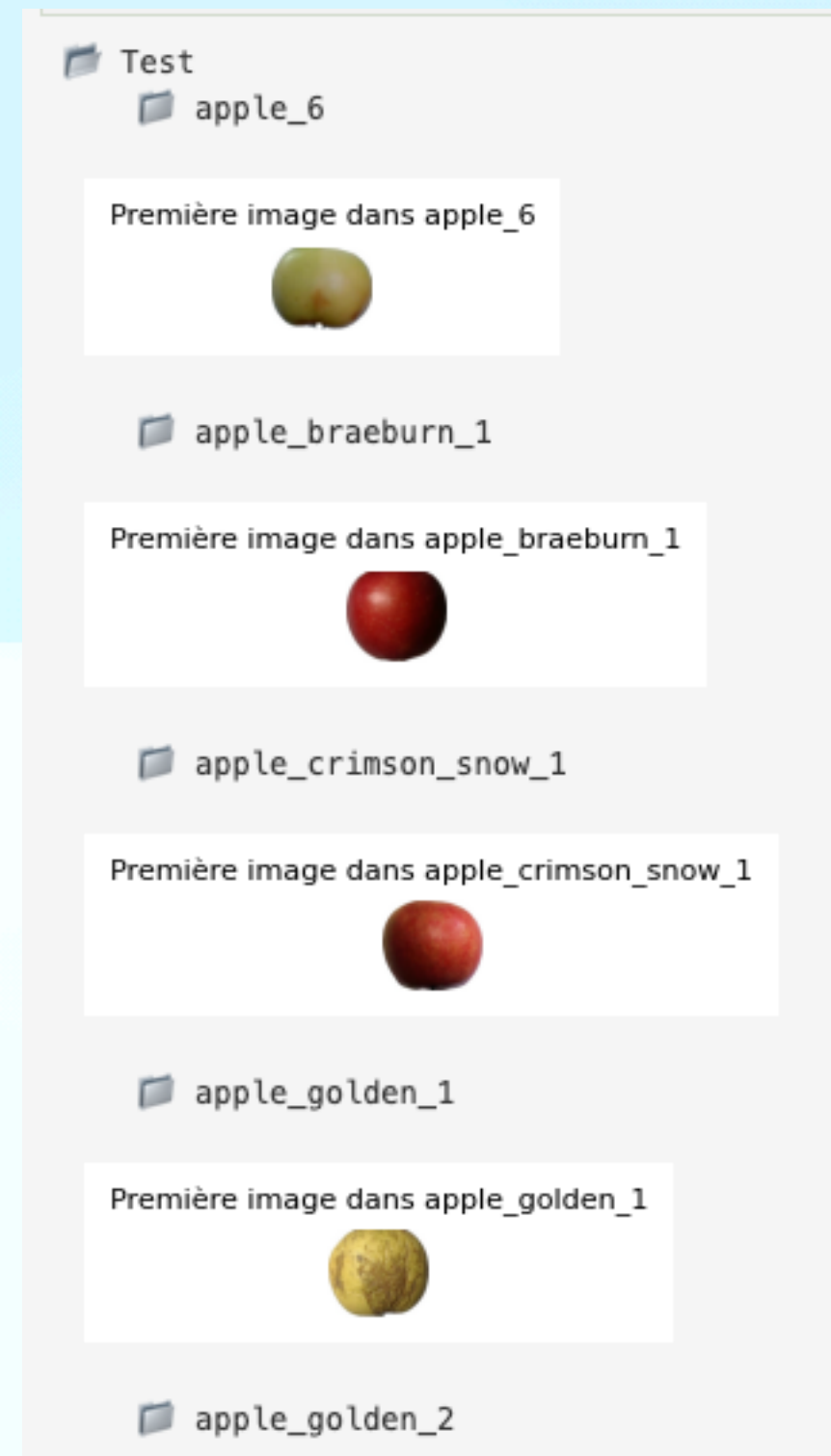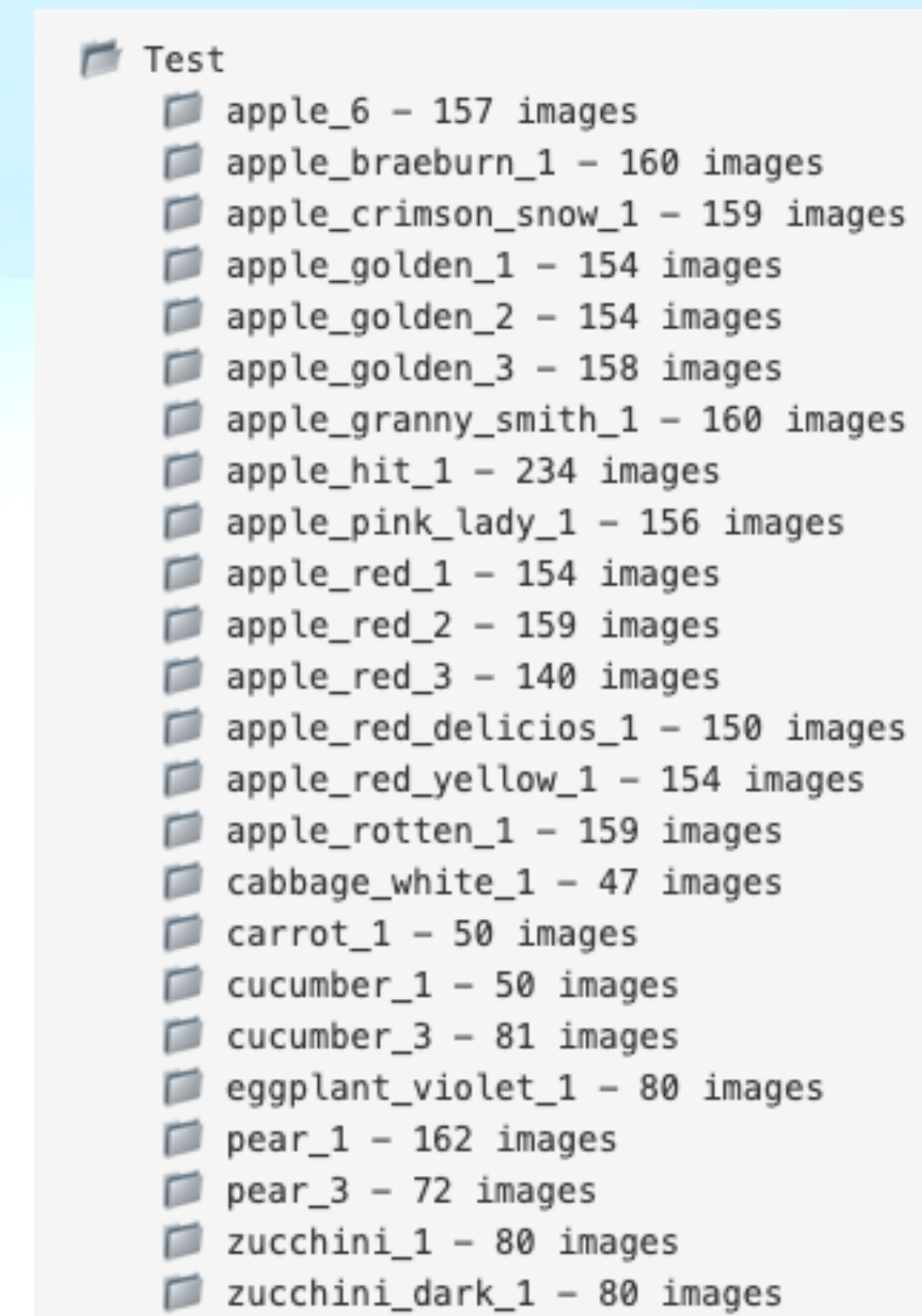
# EDA

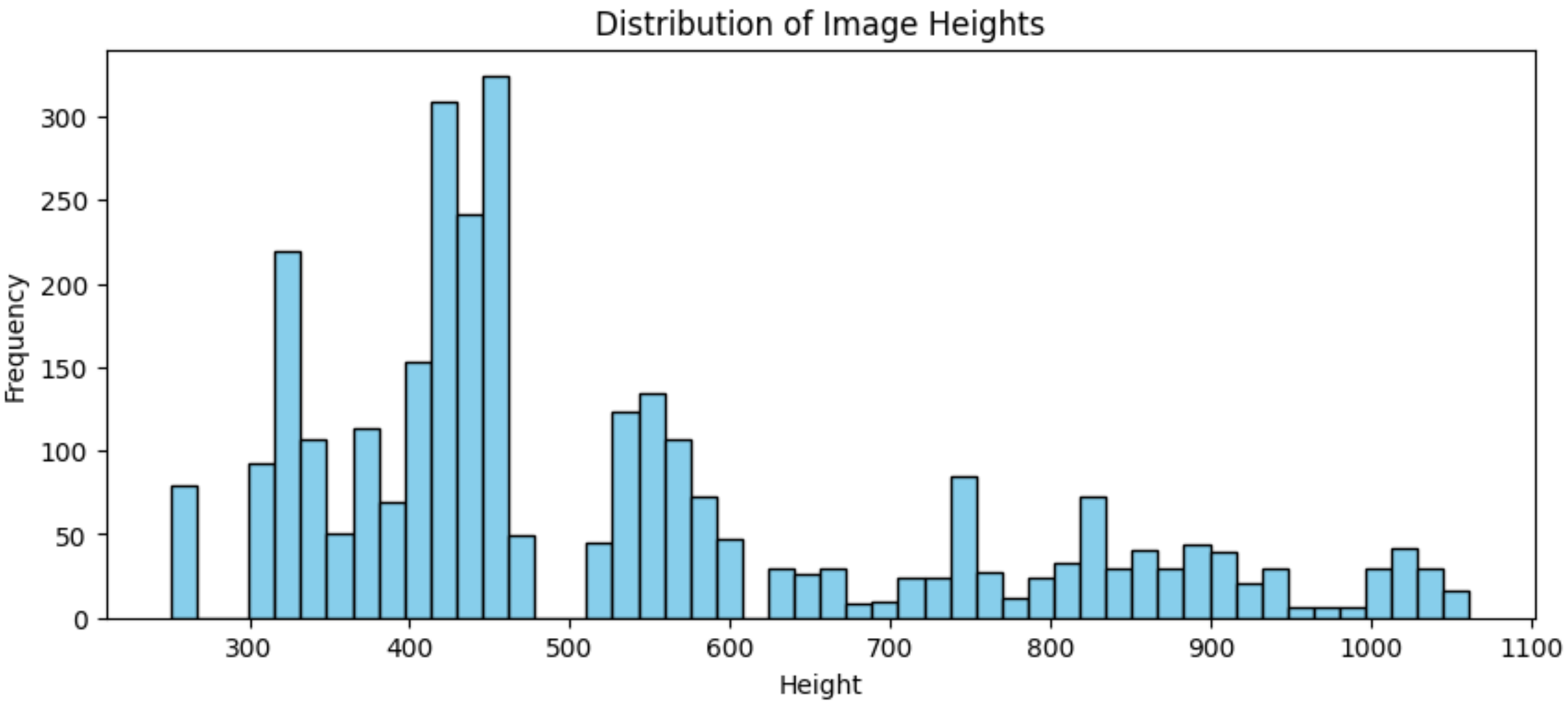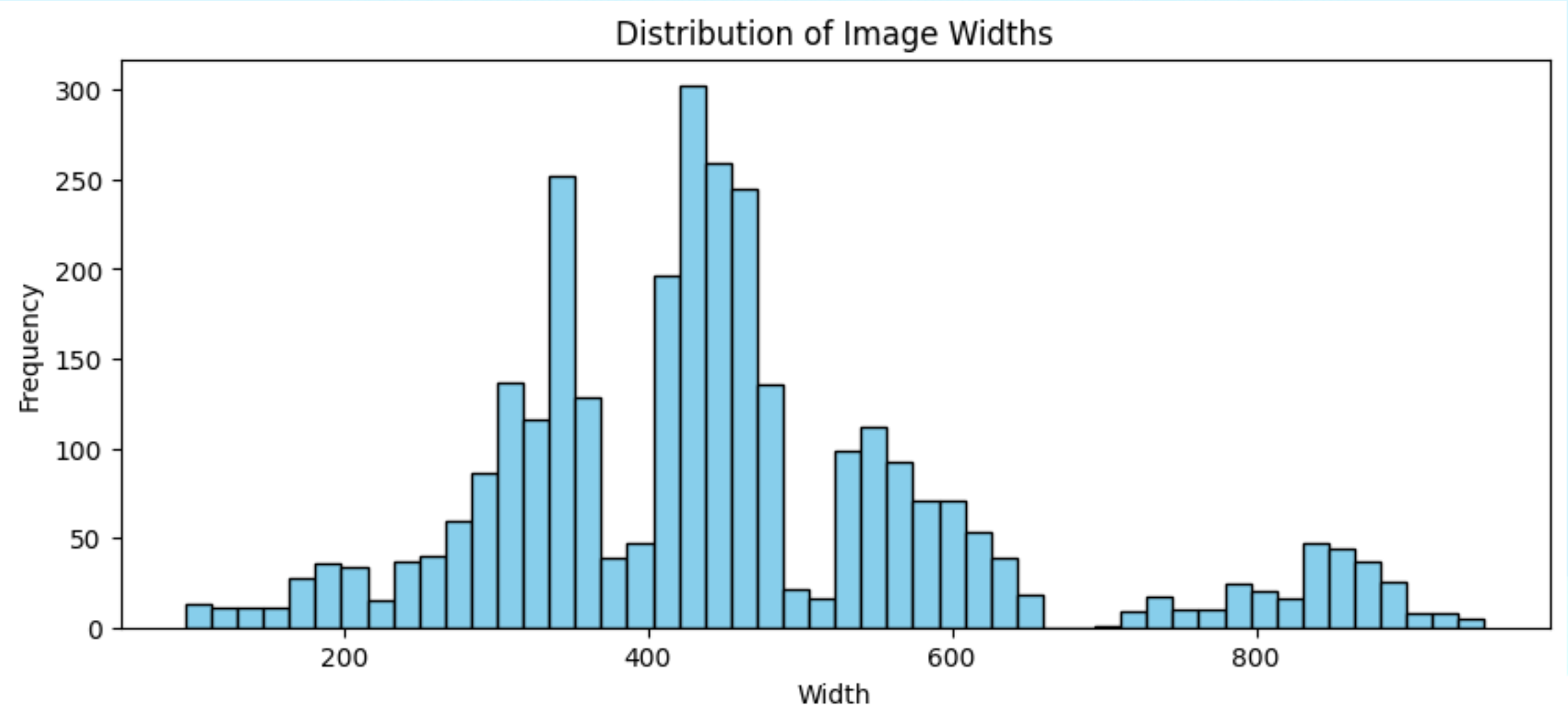## Affichage données



```
Test
    apple_6

Première image dans apple_6


    apple_braeburn_1

Première image dans apple_braeburn_1


    apple_crimson_snow_1

Première image dans apple_crimson_snow_1


    apple_golden_1

Première image dans apple_golden_1


    apple_golden_2
```

## Comptage nombre images

```
Test
    apple_6 – 157 images
    apple_braeburn_1 – 160 images
    apple_crimson_snow_1 – 159 images
    apple_golden_1 – 154 images
    apple_golden_2 – 154 images
    apple_golden_3 – 158 images
    apple_granny_smith_1 – 160 images
    apple_hit_1 – 234 images
    apple_pink_lady_1 – 156 images
    apple_red_1 – 154 images
    apple_red_2 – 159 images
    apple_red_3 – 140 images
    apple_red_delicios_1 – 150 images
    apple_red_yellow_1 – 154 images
    apple_rotten_1 – 159 images
    cabbage_white_1 – 47 images
    carrot_1 – 50 images
    cucumber_1 – 50 images
    cucumber_3 – 81 images
    eggplant_violet_1 – 80 images
    pear_1 – 162 images
    pear_3 – 72 images
    zucchini_1 – 80 images
    zucchini_dark_1 – 80 images
```

# EDA

Graphiques



Class Distribution
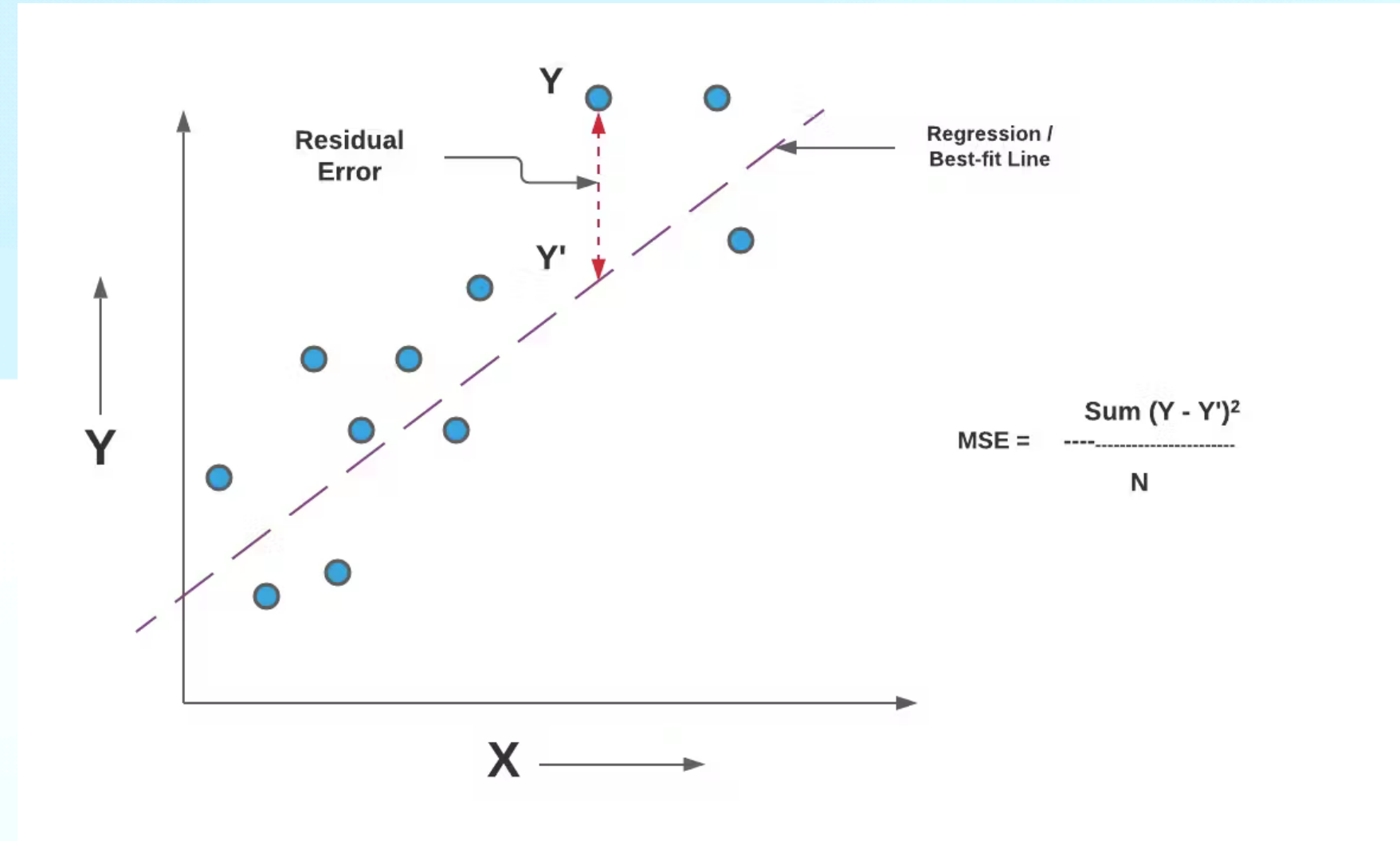


Distribution of Image Widths



Distribution of Image Heights

# EDA

## Calcul difference entre images



```
Average MSE for cucumber_3: 0.5039010431271174
Average MSE for zucchini_1: 0.039984845090060575
Average MSE for eggplant_violet_1: 0.819878872307832
Average MSE for apple_red_yellow_1: 0.2223830180512842
Average MSE for apple_crimson_snow_1: 0.1227337766108603
Average MSE for pear_1: 0.21814232797018862
Average MSE for apple_red_delicios_1: 0.1934729710623934
Average MSE for apple_rotten_1: 0.19269970814960616
Average MSE for apple_golden_3: 0.11211556807380742
Average MSE for apple_golden_2: 0.0825270816466765
Average MSE for apple_red_1: 0.16190522789345568
Average MSE for carrot_1: 0.19650694828175055
Average MSE for apple_granny_smith_1: 0.09207881103573468
Average MSE for apple_braeburn_1: 0.2013861984678091
Average MSE for cabbage_white_1: 0.07342016743440216
Average MSE for cucumber_1: 0.17364077917309909
Average MSE for pear_3: 0.11804350814750816
Average MSE for apple_hit_1: 0.16354016169688232
Average MSE for apple_golden_1: 0.11703786624304234
Average MSE for apple_pink_lady_1: 0.13442306362441345
Average MSE for apple_6: 0.12500830326126317
Average MSE for zucchini_dark_1: 0.3997640414355076
Average MSE for apple_red_2: 0.15895893345860138
Average MSE for apple_red_3: 0.1788846334866389
```

# EDA

Fruits choisi pour l'entrainement

- **apple_granny_smith_1** (nombre élevé et relativement constant d'images : 160)

- **pear_1** (nombre similaire d'images: 162)

- **cucumber_3** (81 images, pour équilibrer l'étude)

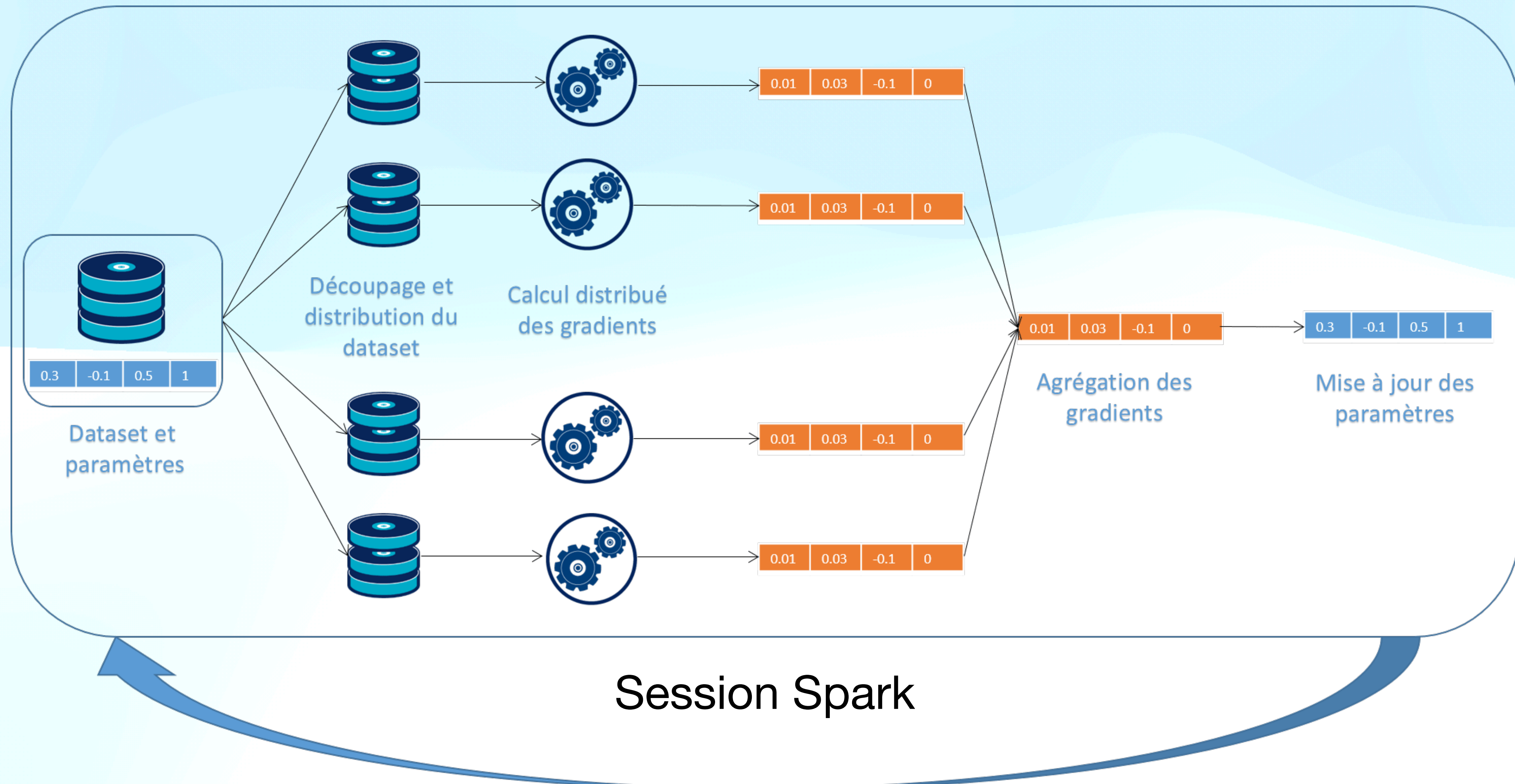- **zucchini_1** (80 images, pour équilibrer l'étude).

# Modélisation

## Architecture utilisé



## Traitement de données distribuées

# Modélisation



Dataset et paramètres
| 0.3 | -0.1 | 0.5 | 1 |

Découpage et distribution du dataset

Calcul distribué des gradients

| 0.01 | 0.03 | -0.1 | 0 |
| 0.01 | 0.03 | -0.1 | 0 |
| 0.01 | 0.03 | -0.1 | 0 |
| 0.01 | 0.03 | -0.1 | 0 |

Agrégation des gradients
| 0.01 | 0.03 | -0.1 | 0 |

Mise à jour des paramètres
| 0.3 | -0.1 | 0.5 | 1 |

Session Spark

# Modélisation

## Session Spark

SparkSession - in-memory

**SparkContext**

Spark UI

Version
v3.5.1

Master
local

AppName
P11

## Initialisation modèle

```python
1  model = MobileNetV2(weights='imagenet',
2                      include_top=True,
3                      input_shape=(224, 224, 3))
```

## Diffusion poids du modèle

```python
1  brodcast_weights = sc.broadcast(new_model.get_weights())
```

# Modélisation

T-Sne

```python
# Assuming df is your DataFrame
features = df['features'].apply(lambda x: np.array(x))

# Convert features into a 2D array
features_2d = np.array(features.tolist())

# Initialize t-SNE
tsne = TSNE(n_components=2, random_state=0)

# Apply t-SNE to the data
tsne_results = tsne.fit_transform(features_2d)

# Create a DataFrame with t-SNE results and labels
tsne_df = pd.DataFrame({'X': tsne_results[:, 0], 'Y': tsne_results[:, 1], 'label': df['label']})

# Plot the results with labels as hue
plt.figure(figsize=(10, 10))
sns.scatterplot(x='X', y='Y', hue='label', data=tsne_df)
plt.show()
```

# Modélisation

## PCA

```python
# Define a UDF to convert array to vector
list_to_vector_udf = udf(lambda l: Vectors.dense(l), VectorUDT())

# Load the data
df = spark.read.parquet(PATH_Result)

# Convert the array of floats to a vector
df = df.withColumn("features_vec", list_to_vector_udf(df["features"]))

# Apply PCA
pca = PCA(k=2, inputCol="features_vec", outputCol="pcaFeatures")
model = pca.fit(df)
result = model.transform(df)
result = result.drop('features_vec')
result = result.drop('features')
# Convert the Spark DataFrame to a Pandas DataFrame
result_pd = result.toPandas()
```

## Résultats

|   | path | label | pcaFeatures |
|---|------|-------|-------------|
| 0 | file:/Users/zaccaria/Documents/Progetti/ocia/o... | cucumber_3 | [11.308647269291315, -4.089936426123957] |
| 1 | file:/Users/zaccaria/Documents/Progetti/ocia/o... | cucumber_3 | [10.564156308055871, -4.0341143125233065] |
| 2 | file:/Users/zaccaria/Documents/Progetti/ocia/o... | cucumber_3 | [12.559707640014853, -3.3085778039897167] |
| 3 | file:/Users/zaccaria/Documents/Progetti/ocia/o... | cucumber_3 | [10.112701970729733, -2.994906086358103] |
| 4 | file:/Users/zaccaria/Documents/Progetti/ocia/o... | cucumber_3 | [12.475473290124103, -5.459863035549889] |

# Cloud

3 Approches

- Local

- Script

- Notebook

Script

Notebook

Compute

# RGPD

Article 3 – Territorial Scope

 ➔ Europe (Paris) eu-west-3

 ➔ FranceCentral

# Cloud

## Cluster EMR

## Sommaire

# Cloud

IAM

Utilisateur

Policies



zaccaria-p11 Info

## Summary

ARN
⧉ arn:aws:iam::975049997739:user/zaccaria-p11

Console access
⚠ Enabled without MFA

Created
June 17, 2024, 09:02 (UTC+02:00)

Last console sign-in
ⓘ Never

Permissions | Groups (1) | Tags (1) | Security credentials | Access Advisor

## Permissions policies (2)

Permissions are defined by policies attached to the user directly or through groups.

🔍 Search

Filter by Type
All types

| | Policy name ↗ | ▲ | Type |
|---|---|---|---|
| ☐ | ⊞ 📦 AdministratorAccess | | AWS managed - job function |
| ☐ | ⊞ 📦 IAMUserChangePassword | | AWS managed |

17

# Cloud
## Databricks

### Deployment

### Initialisation cluster

| State | Name | Runtime | Active memory | Active cores | Active DBU / h | Source | Creator | N |
|-------|------|---------|---------------|--------------|----------------|--------|---------|---|
| ↻ | Zaccaria Amillou's Cluster | 14.3 | - | - | 1.5 | UI | zaccaria.amillou@gmail.com | |

**⋯ Deployment is in progress**

Deployment name : p11-oc_databricks_oc
Subscription       : Azure subscription 1
Resource group  : p11-oc

∨ Deployment details

**Resource**

🔄 databricks_oc

### Import des données

**1.2 Import données**

```
▶    ✓ 01:10 PM (<1s)                          8

PATH = "dbfs:/mnt/p11_mount/data"
PATH_Data = PATH+'/Test'
PATH_Result = PATH+'/Results'

print('PATH:         '+\
      PATH+'\nPATH_Data:    '+\
      PATH_Data+'\nPATH_Result: '+PATH_Result)
```

```
PATH:        dbfs:/mnt/p11_mount/data
PATH_Data:   dbfs:/mnt/p11_mount/data/Test
PATH_Result: dbfs:/mnt/p11_mount/data/Results
```

18

# Cloud

## Résultats

# Conclusions

Analyse des données

Mise en place du modèle

Traitement des données sur le cloud

Respect des normes RGPD

# Merci pour votre attention