

Projet n°8 : Déployez un modèle dans le cloud

Parcours Data Scientist - OpenClassrooms

Jérémy Vangansberg - octobre 2021



Fruits!

Introduction (1)

- La start-up **Fruits!** cherche à proposer des solutions innovantes pour la récolte de fruits.
- L'entreprise souhaite dans un premier temps se faire connaître en mettant à disposition du grand public une application qui pourra permettre de prendre en photo un fruit et avoir des informations sur celui-ci. Cette application implique de mettre en place un **moteur de classification d'image**.
- L'application doit être conçu dans un environnement Big Data pour faire face à une potentielle évolution du nombre d'utilisateurs.

Introduction (2)



- Le but du projet est de mettre en place le pipeline de traitement des images.
- Il n'est pas **demandé d'effectuer la modélisation.**

Sommaire



- A. Les données
- B. Le Big Data et son architecture
- C. Les solutions envisagées
- D. Databricks : fonctionnement et paramétrage
- E. La chaîne de traitement
- F. Conclusion

Les données

Le jeu de données brut



- Données d'entraînement : 67692 images (un fruit ou un légume par image)
- Données de test : 22688 images (un fruit ou un légume par image)
- Nombre de classes : 131 (fruits ou légumes)
- Taille des images : 100 x 100 pixels



Fruits!

Le jeu de données

- Parmi ce jeu de données, j'ai sélectionné **8 images** pour faire une démonstration du pipeline.
- Classes sélectionnées :
 - Banane
 - Cerise
 - Aubergine
 - Kiwi
- Note : je propose également un pipeline alternatif où je traite l'intégralité des images.

Le Big Data et son architecture

Le Big Data



- Quelques chiffres en 2020:
 - 40 zétaoctets (10^{21} octets) de données générées.
 - 150 millions d'e-mails envoyés par minute.
- Qu'est ce que le Big Data ?
 - C'est lorsque le **volume de données est trop important** pour utiliser des **méthodes traditionnelles**.
- Quand commence le Big Data ?
 - **Ça dépend**, par exemple 100 Go de données à traiter peut paraître être un volume important mais à l'échelle de l'un des acteurs des GAFAM, c'est une goutte d'eau. D'une manière générale, on retient que lorsque les données **sont trop grandes pour être stockées en RAM**, on est confronté à un problème de Big Data.



Fruits!

Les 3V du Big Data

- J'ai évoqué précédemment que le **volume** de données est ce qui caractérise le Big Data. En effet, il est l'un des 3 V du Big Data mais il y a également :
- La **vélocité**, il faut être capable de traiter les données dans des délais acceptables. Par exemple, imaginons qu'une société mette une heure à traiter 1 Go de données de log pour faire un rapport journalier. Si la société produit plus de 24 Go de données par jour, elle n'est déjà plus capable de traiter les données de manière journalière.
- La **variété**, le big data se caractérise aussi par l'hétérogénéité des formats. Ils peuvent être structurés (JSON), semi-structurés (logs), non-structurés (texte et image).

Le cloud : l'émergence du Big Data



- **Avant l'ère du cloud**, les entreprises devaient s'affranchir de frais considérables pour acheter et maintenir des serveurs.
- Désormais avec des acteurs comme **AWS, Microsoft Azure et Google Cloud** des capacités de calcul importantes sont à la disposition de tout le monde à des coûts acceptables.

3 briques du Big Data



- Map Reduce
- Hadoop Map Reduce
- Spark

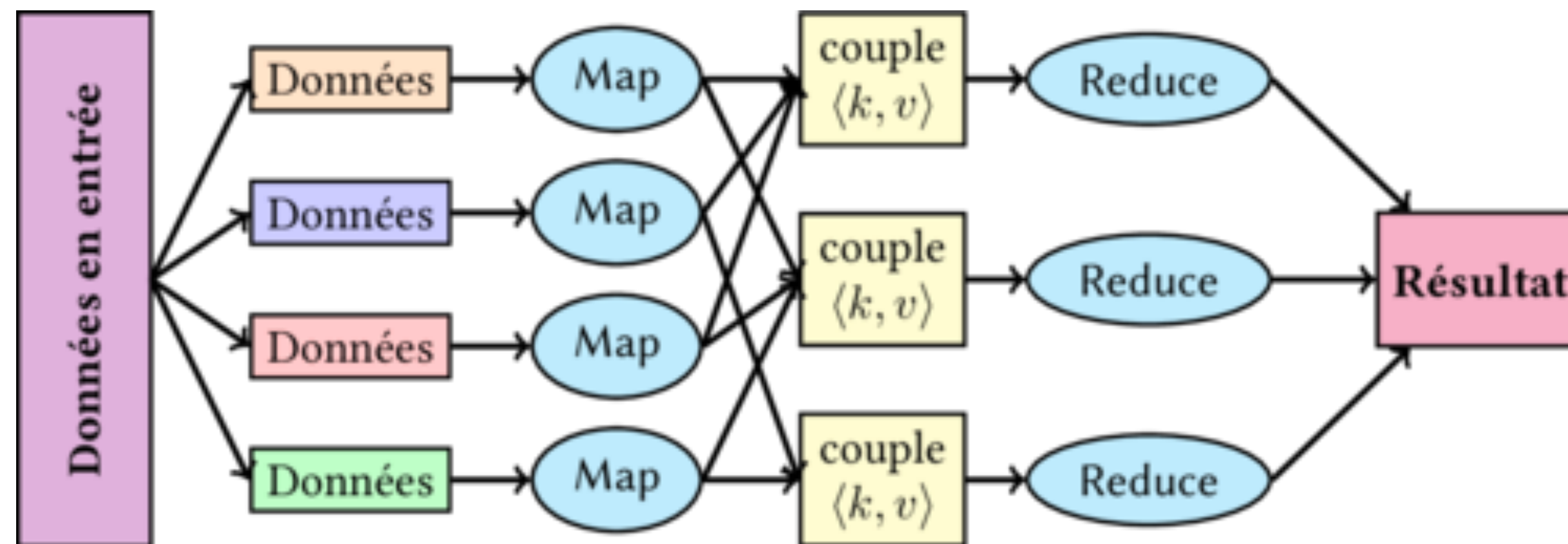
Map Reduce



Fruits!

- L'algorithme de Map Reduce, est le composant fondamental du Big Data. Il permet de **paralléliser** et **distribuer** des calculs.

ARCHITECTURE MAP REDUCE



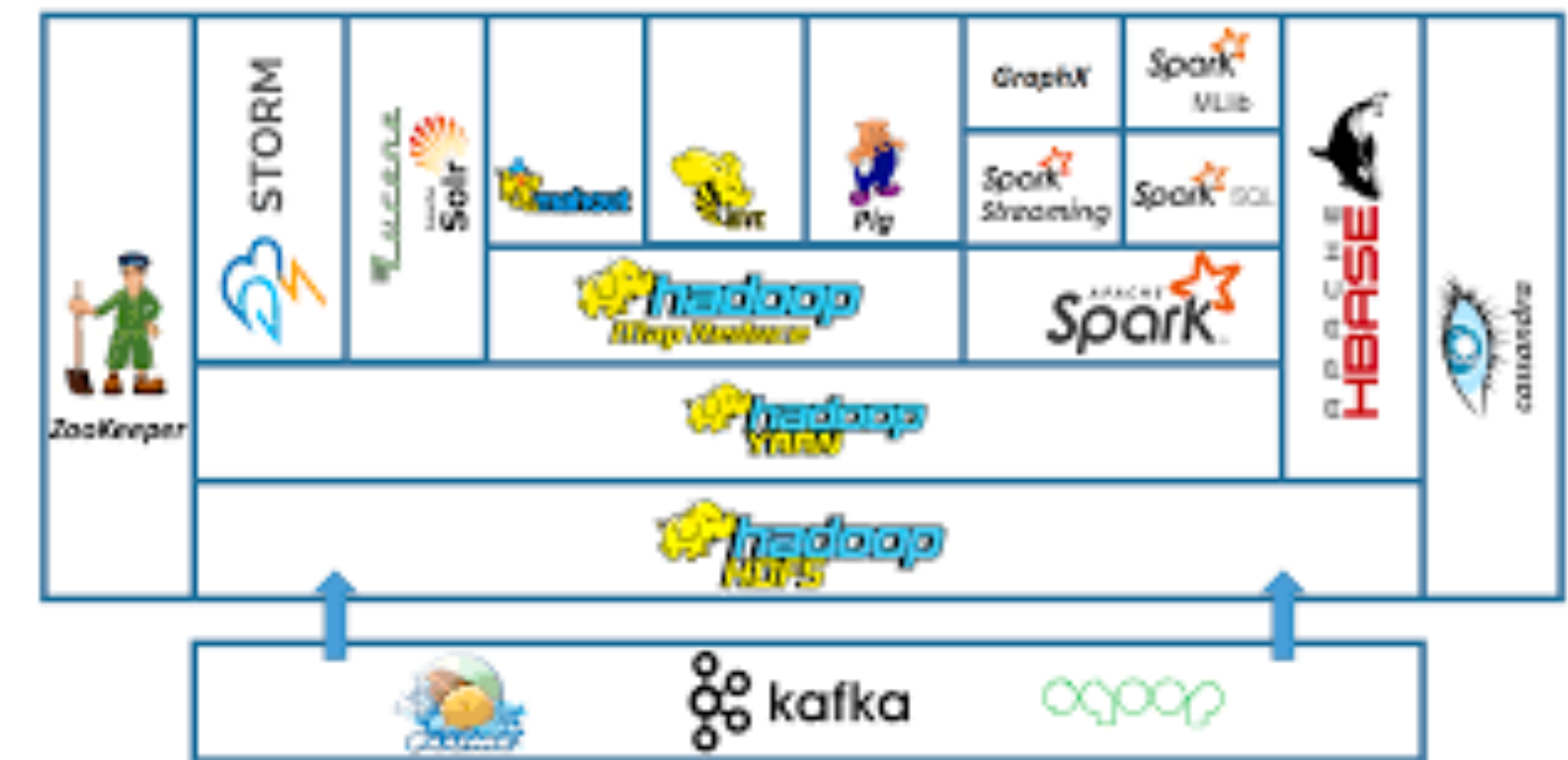
- Cette architecture **développée par Google** au début des années 2000 s'articule autour de deux opérations : Map et Reduce. La première permet **d'appliquer une fonction** à des données qui ont été au préalable partitionnées. La seconde permet **d'agréger les résultats**.



Fruits!

- C'est un **projet open-source** inspiré des travaux de Google qui permet d'appliquer l'architecture Map Reduce de manière **générique**. Le projet a rejoint la fondation Apache en 2008.
- Aujourd'hui tout un Écosystème existe autour d'Hadoop.

ÉCOSYSTÈME HADOOP





- L'un des inconvénients majeurs lorsque l'on veut travailler avec Hadoop MapReduce, est que l'on doit **traduire un problème en un enchaînement d'opération *map* et *reduce***. De plus Hadoop MapReduce écrit les données sur disque entre chaque opération.
- Spark permet de **s'abstraire** de MapReduce car il s'utilise comme un langage de haut niveau.
- Spark va stocker les données en **RAM**, ce qui peut accélérer le traitement des données par un facteur entre 10 et 100.



Architecture Spark : les composantes

- Spark dépend de plusieurs éléments :
 - Hadoop
 - Java
 - Scala
- Pour pouvoir utiliser le langage python avec spark, il faut se servir de la librairie **Pyspark**.
- De plus, pour pouvoir mettre en place le pipeline de traitement des images, j'aurai également besoin de :
 - pandas, PIL, tensorflow, os, io, numpy, boto3

Les solutions envisagées



Fruits!

Les solutions envisagées



Amazon EC2

Service web qui propose des capacités de calcul dans le cloud



Amazon SageMaker

Service clé en main qui permet de faire du machine learning



Amazon EMR

Semblable à EC2 mais spécialisé dans le traitement de big data



databricks

Entreprise qui travaille avec AWS. Elle permet de lancer des clusters et des notebooks pré-configurés pour le big data



Fruits!

Difficultés rencontrées

- La principale difficulté a été liée aux conflits entre les différentes versions.
- En effet les **versions évoluent rapidement** et les tutoriels deviennent rapidement obsolètes.
- Pour pouvoir me concentrer sur le pipeline de traitement des données et de résoudre les problèmes liés à la configuration, j'ai décidé de me tourner vers **Databricks**.

Databricks

Databricks



Avantage :

- Databricks permet de réduire le temps de déploiement d'un modèle de big data

Inconvénient :

- Solution payante



Databricks - marche à suivre

1. Créer un compte Databricks et profiter des 14 jours d'essai gratuit
2. Associer son compte avec AWS en mode *Cross Account Role*
 - <https://docs.databricks.com/administration-guide/account-api/iam-role.html>
3. Définir un *bucket* S3 où stocker les informations de Databricks
4. Paramétrer les politiques de sécurité du *bucket* S3 où le jeu de données est stocké
 - <https://docs.databricks.com/administration-guide/cloud-configurations/aws/instance-profiles.html>
5. Déployer un *cluster* (avec l'*instance profile* défini au point précédent)
6. Lancer un *notebook* associé à ce *cluster*

Databricks : cluster



Create Cluster

New Cluster DBU / hour: 3 - 9 2-8 Workers:61-244 GB Memory, 8-32 Cores
1 Driver:30.5 GB Memory, 4 Cores

Cluster Name Please enter a cluster name

Cluster Mode ?
Standard

Databricks Runtime Version ?
Runtime: 8.3 (Scala 2.12, Spark 3.1.1)

Note Databricks Runtime 8.x and later use Delta Lake as the default table format. [Learn more](#)

Autopilot Options
☒ Enable autoscaling ?
☐ Enable autoscaling local storage ?
☒ Terminate after minutes of inactivity ?

Worker Type ? Min Workers Max Workers
i3.xlarge 30.5 GB Memory, 4 Cores

New Configure separate pools for workers and drivers for flexibility. [Learn more](#)

Driver Type
Same as worker 30.5 GB Memory, 4 Cores

DBU / hour: 3 - 9 ?

► Advanced Options

tarifs pour la configuration actuelle

- Databricks nous permet de lancer un cluster qui va lancer une ou plusieurs instances ec2

choix des versions de scala et spark

choix du type d'instance EC2 et leur nombre

Databricks : notebook



P8_01_notebook (Python)

Cluster-single-node-ML | File | Edit | View: Standard | Permissions | Run All | Clear

Comments | Experiment | Revision history

Cmd 1

Modules import

Cmd 2

```
# Standard modules
import os
import io
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from PIL import Image
import boto3

# tensorflow modules
from tensorflow.keras.applications.resnet50 import ResNet50, preprocess_input
from tensorflow.keras.preprocessing.image import img_to_array

# pyspark modules
import pyspark
from pyspark.sql import SparkSession
from pyspark import SparkContext, SparkConf
from pyspark.sql.functions import col, pandas_udf, PandasUDFType, split
from pyspark.ml.linalg import Vectors, VectorUDT
from pyspark.sql.functions import udf
from pyspark.ml.feature import StringIndexer, StandardScaler
from pyspark.ml.feature import PCA
```

Command took 34.23 seconds -- by jeremy.vangansberg@gmail.com at 14/10/2021, 10:48:41 on Cluster-single-node-ML

Cmd 3

On associe le notebook avec le cluster

Les principales librairies sont déjà installées

Les paramètres du cluster



Clusters / cluster-multi-i3-2

cluster-multi-i3-2 [Edit] [Start] [Clone] [Delete]

Configuration | Notebooks (0) | Libraries | Event Log | Spark UI | Driver Logs | Metrics | Apps | Spark Cluster UI - Master

Cluster Mode ?

Standard

Databricks Runtime Version

8.4 ML (includes Apache Spark 3.1.2, Scala 2.12)

Autopilot Options

☒ Enable autoscaling ?

☐ Enable autoscaling local storage ?

☒ Terminate after 120 minutes of inactivity ?

Worker Type ?

i3.xlarge 30.5 GB Memory, 4 Cores

Min Workers: 2 Max Workers: 8

Driver Type

i3.xlarge 30.5 GB Memory, 4 Cores

DBU / hour: 3 - 9 ? i3.xlarge

► Advanced Options

Mode standard (plusieurs noeuds)

Version ML avec les principaux modules

type d'instance, 30GB de mémoire et 4 cores. 2 à 8 noeuds

coût en DBU par heure

La chaîne de traitement

Stratégie



Fruits!

- Afin de réduire les coûts associés aux clusters, j'ai décidé de limiter le jeu de données à **8 images** de 4 classes différentes.
- Pour extraire les features des images, j'ai effectué du *transfer learning* en utilisant le modèle **ResNet50**.
- Utiliser une série de fonctions pour transformer les données binaires en une liste de vecteur.
- Réduire la dimension de ce vecteur à 16 dimensions grâce à une PCA.

Étapes (1)



Fruits!

1. INITIALISATION D'UNE SESSION SPARK

```
Cmd 4

# Instanciation of spark with credentials

spark = (SparkSession.builder
  .appName('jv-p8')
  .config('spark.hadoop.fs.s3a.access.key', '██████████')
  .config('spark.hadoop.fs.s3a.secret.key', '████████████████████████████████████████')
  .config('spark.hadoop.fs.s3a.impl', 'org.apache.hadoop.fs.s3a.S3AFileSystem')
  .getOrCreate()
)

sc = spark.sparkContext
sc.setSystemProperty('com.amazonaws.services.s3a.enableV4', 'true')
sc._jsc.hadoopConfiguration().set("fs.s3a.endpoint", "s3.eu-west-1.amazonaws.com")

Command took 0.08 seconds -- by jeremy.vangansberg@gmail.com at 14/10/2021, 10:48:41 on Cluster-single-node-ML
```

2. CHARGEMENT DES DONNÉES

```
Cmd 7

# Loading the data

s3_url = "s3a://s3-jv-ocr/reduced_dataset/*"
image_df = image_df = spark.read.format("binaryfile").load(s3_url)
image_df.printSchema()

▶ image_df: pyspark.sql.dataframe.DataFrame = [path: string, modificationTime: timestamp ... 2 more fields]

root
 |-- path: string (nullable = true)
 |-- modificationTime: timestamp (nullable = true)
 |-- length: long (nullable = true)
 |-- content: binary (nullable = true)

Command took 2.36 seconds -- by jeremy.vangansberg@gmail.com at 14/10/2021, 10:48:41 on Cluster-single-node-ML
```

Étapes (2)



Fruits!

3. MODIFICATION DU JEU DE DONNÉES

```
# Data restructuring
```

```
image_df = image_df.withColumn('label', split(col('path'), '/').getItem(4))
image_df = image_df.select('path', 'content', 'label')
image_df.show()
```

▶ (2) Spark Jobs

▶  image_df: pyspark.sql.dataframe.DataFrame = [path: string, content: binary ... 1 more field]

```
+-----+-----+-----+
|           path|           content|   label|
+-----+-----+-----+
|s3a://s3-jv-ocr/r...|[FF D8 FF E0 00 1...|    Kiwi|
|s3a://s3-jv-ocr/r...|[FF D8 FF E0 00 1...|    Kiwi|
|s3a://s3-jv-ocr/r...|[FF D8 FF E0 00 1...|Cherry_1|
|s3a://s3-jv-ocr/r...|[FF D8 FF E0 00 1...|Cherry_1|
|s3a://s3-jv-ocr/r...|[FF D8 FF E0 00 1...|  Banana|
|s3a://s3-jv-ocr/r...|[FF D8 FF E0 00 1...|  Banana|
|s3a://s3-jv-ocr/r...|[FF D8 FF E0 00 1...|Eggplant|
|s3a://s3-jv-ocr/r...|[FF D8 FF E0 00 1...|Eggplant|
+-----+-----+-----+
```

Command took 2.57 seconds -- by jeremy.vangansberg@gmail.com at 14/10/2021, 10:48:41 on Cluster-single-node-ML

Étapes (3)



Fruits!

4. CHARGEMENT DU MODÈLE

Cmd 13

```
# Deeplearning model instantiation

model = ResNet50(include_top=False)
model.summary() # verify that the top layer is removed
```

conv5_block3_1_relu (Activation)	(None, None, None, 5 0	conv5_block3_1_bn[0][0]

conv5_block3_2_conv (Conv2D)	(None, None, None, 5 2359808	conv5_block3_1_relu[0][0]

conv5_block3_2_bn (BatchNormali	(None, None, None, 5 2048	conv5_block3_2_conv[0][0]

conv5_block3_2_relu (Activation)	(None, None, None, 5 0	conv5_block3_2_bn[0][0]

conv5_block3_3_conv (Conv2D)	(None, None, None, 2 1050624	conv5_block3_2_relu[0][0]

conv5_block3_3_bn (BatchNormali	(None, None, None, 2 8192	conv5_block3_3_conv[0][0]

conv5_block3_add (Add)	(None, None, None, 2 0	conv5_block2_out[0][0] conv5_block3_3_bn[0][0]

conv5_block3_out (Activation)	(None, None, None, 2 0	conv5_block3_add[0][0]
=====		
Total params: 23,587,712		
Trainable params: 23,534,592		
Non-trainable params: 53,120		

Command took 5.41 seconds -- by jeremy.vangansberg@gmail.com at 14/10/2021, 10:48:41 on Cluster-single-node-ML

Étapes (4)



Fruits!

5. DÉFINIR UNE SÉRIE DE FONCTION POUR TRANSFORMER LES DONNÉES EN VECTEUR

```
bc_model_weights = spark.sparkContext.broadcast(model.get_weights())
def model_fn():
    """
    Returns a ResNet50 model with top layer removed and broadcasted pretrained weights.
    """
    model = ResNet50(weights=None, include_top=False)
    model.set_weights(bc_model_weights.value)
    return model

def preprocess(content):
    """
    Preprocesses raw image bytes for prediction.
    """
    img = Image.open(io.BytesIO(content)).resize([50, 50])
    arr = img_to_array(img)
    return preprocess_input(arr)

def featurize_series(model, content_series):
    """
    Featurize a pd.Series of raw images using the input model.
    :return: a pd.Series of image features
    """

    input = np.stack(content_series.map(preprocess))
    preds = model.predict(input)
    output = [p.flatten() for p in preds]
    return pd.Series(output)

@pandas_udf('array<float>', PandasUDFType.SCALAR_ITER)
def featurize_udf(content_series_iter):
    """
    This method is a Scalar Iterator pandas UDF wrapping our featurization function.
    The decorator specifies that this returns a Spark DataFrame column of type ArrayType(FloatType).

    :param content_series_iter: This argument is an iterator over batches of data, where each batch
                                is a pandas Series of image data.
    """
    model = model_fn()
    for content_series in content_series_iter:
        yield featurize_series(model, content_series)
```

Étapes (5)



Fruits!

6. RÉDUCTION DE LA DIMENSION

Cmd 18

Dimensions reduction with a PCA

Cmd 19

```
# standard scaler before doing the PCA

scaler = StandardScaler(withMean=True, withStd=True,
                        inputCol='features',
                        outputCol='features_scaled')
scaler_fitted = scaler.fit(features_df)
features_scaled_df = scaler_fitted.transform(features_df)
```

▸ (2) Spark Jobs

▸  features_scaled_df: pyspark.sql.dataframe.DataFrame = [path: string, label: string ... 2 more fields]

Command took 8.99 seconds -- by jeremy.vangansberg@gmail.com at 19/10/2021, 13:36:10 on cluster-multi-i3-2

Cmd 20

```
# I choose a number of dimensions pretty low because the cluster is on a single node to reduce the

pca = PCA(k=16, inputCol="features_scaled", outputCol="features_pca")
model = pca.fit(features_scaled_df)
df_transformed = model.transform(features_scaled_df)
```

▸ (5) Spark Jobs

▸  df_transformed: pyspark.sql.dataframe.DataFrame = [path: string, label: string ... 3 more fields]

Command took 34.16 minutes -- by jeremy.vangansberg@gmail.com at 19/10/2021, 13:36:10 on cluster-multi-i3-2

Étapes (6)



Fruits!

7. EXPORT DES DONNÉES VERS S3

Cmd 22

Export the data scaled to S3

Cmd 23

```
# Data selection and transformation to pandas

df_pandas = df_transformed.select('path', 'label', 'features_pca').toPandas()
```

► (1) Spark Jobs

/databricks/spark/python/pyspark/sql/pandas/conversion.py:92: UserWarning: toPandas attempted Arrow optim
Unable to convert the field features_pca. If this column is not necessary, you may consider dropping it
Direct cause: Unsupported type in conversion to Arrow: VectorUDT
Attempting non-optimization as 'spark.sql.execution.arrow.pyspark.fallback.enabled' is set to true.
warnings.warn(msg)

Command took 7.16 seconds -- by jeremy.vangansberg@gmail.com at 19/10/2021, 13:36:10 on cluster-multi-i3-2

Cmd 24

```
# Export to S3 bucket

from io import StringIO

bucket = "s3-jv-ocr"

csv_buffer = StringIO()
df_pandas.to_csv(csv_buffer)

s3_resource = boto3.resource('s3')

s3_resource.Object(bucket, 'df_reduced_pca_16.csv').put(Body=csv_buffer.getvalue())
```

Les données dans le S3



BUCKET S3-JV-OCR

LES PARTITIONS DE FICHIERS CSV

Amazon S3 > s3-jv-ocr

s3-jv-ocr

Info

Objects

Properties

Permissions

Metrics

Management

Access Points

Objects (4)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in y

Copy S3 URI

Copy URL

Download

Open

Delete

Find objects by prefix

	Name	Type	Last modifi
	alternative_full_data_set/	Folder	-
	dataset/	Folder	-
	df_reduced_pca_16.csv	csv	October 19
	reduced_dataset/	Folder	-

To enable sorting in the table below, use the search to reduce the size of the list to 999 objects or fewer.

Objects (999+)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Copy S3 URI

Copy URL

Download

Open

Delete

Actions

Create folder

Upload

Find objects by prefix

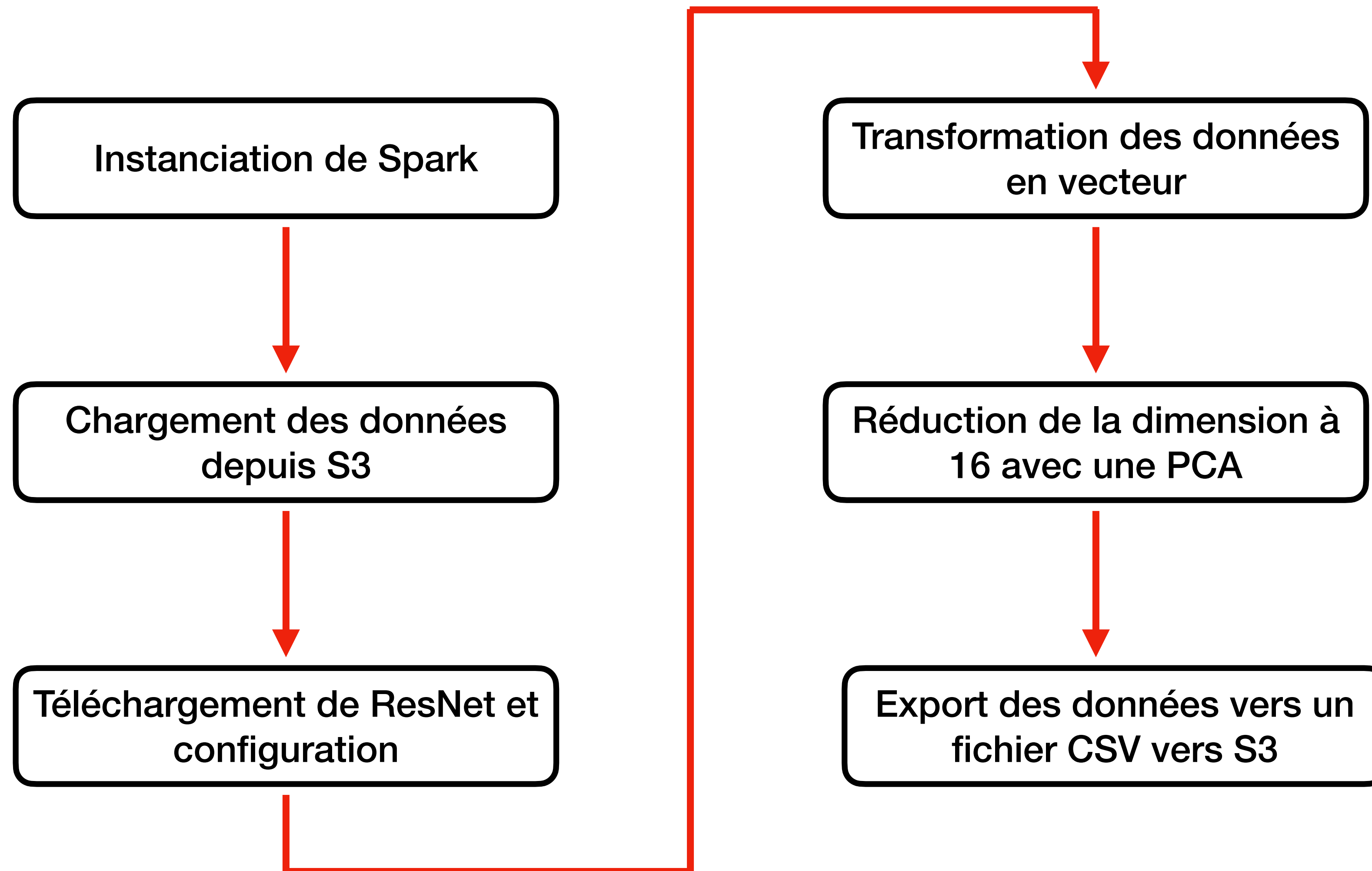
	Name	Type	Last modified	Size	Storage class
	_SUCCESS	-	October 19, 2021, 16:00:44 (UTC+02:00)	0 B	Standard
	_committed_2017549523352621764	-	October 19, 2021, 16:00:44 (UTC+02:00)	160.1 KB	Standard
	_started_2017549523352621764	-	October 19, 2021, 15:50:18 (UTC+02:00)	0 B	Standard
	part-00000-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-172-1-c000.csv	csv	October 19, 2021, 15:50:44 (UTC+02:00)	621.8 KB	Standard
	part-00001-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-173-1-c000.csv	csv	October 19, 2021, 15:50:47 (UTC+02:00)	630.5 KB	Standard
	part-00002-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-174-1-c000.csv	csv	October 19, 2021, 15:50:27 (UTC+02:00)	647.2 KB	Standard
	part-00003-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-175-1-c000.csv	csv	October 19, 2021, 15:50:27 (UTC+02:00)	677.3 KB	Standard
	part-00004-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-176-1-c000.csv	csv	October 19, 2021, 15:50:27 (UTC+02:00)	676.5 KB	Standard
	part-00005-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-177-1-c000.csv	csv	October 19, 2021, 15:50:47 (UTC+02:00)	659.6 KB	Standard
	part-00006-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-178-1-c000.csv	csv	October 19, 2021, 15:50:45 (UTC+02:00)	658.7 KB	Standard
	part-00007-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-179-1-c000.csv	csv	October 19, 2021, 15:50:55 (UTC+02:00)	660.0 KB	Standard
	part-00008-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-180-1-c000.csv	csv	October 19, 2021, 15:50:44 (UTC+02:00)	640.8 KB	Standard
	part-00009-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-181-1-c000.csv	csv	October 19, 2021, 15:50:47 (UTC+02:00)	634.5 KB	Standard
	part-00010-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-182-1-c000.csv	csv	October 19, 2021, 15:50:27 (UTC+02:00)	629.5 KB	Standard
	part-00011-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-183-1-c000.csv	csv	October 19, 2021, 15:50:27 (UTC+02:00)	629.3 KB	Standard
	part-00012-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-184-1-c000.csv	csv	October 19, 2021, 15:50:27 (UTC+02:00)	627.0 KB	Standard
	part-00013-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-185-1-c000.csv	csv	October 19, 2021, 15:50:47 (UTC+02:00)	641.1 KB	Standard
	part-00014-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-186-1-c000.csv	csv	October 19, 2021, 15:50:45 (UTC+02:00)	642.4 KB	Standard
	part-00015-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-187-1-c000.csv	csv	October 19, 2021, 15:50:55 (UTC+02:00)	641.1 KB	Standard
	part-00016-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-188-1-c000.csv	csv	October 19, 2021, 15:50:34 (UTC+02:00)	643.3 KB	Standard
	part-00017-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-189-1-c000.csv	csv	October 19, 2021, 15:50:34 (UTC+02:00)	652.6 KB	Standard
	part-00018-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-190-1-c000.csv	csv	October 19, 2021, 15:50:34 (UTC+02:00)	650.4 KB	Standard
	part-00019-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-191-1-c000.csv	csv	October 19, 2021, 15:50:34 (UTC+02:00)	652.6 KB	Standard
	part-00020-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-192-1-c000.csv	csv	October 19, 2021, 15:50:34 (UTC+02:00)	646.6 KB	Standard
	part-00021-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-193-1-c000.csv	csv	October 19, 2021, 15:50:35 (UTC+02:00)	637.5 KB	Standard
	part-00022-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-194-1-c000.csv	csv	October 19, 2021, 15:50:41 (UTC+02:00)	649.4 KB	Standard
	part-00023-tid-2017549523352621764-d416d19e-0299-4910-b35e-66778819af0b-195-1-c000.csv	csv	October 19, 2021, 15:50:41 (UTC+02:00)	649.5 KB	Standard

34

Récapitulatif de la chaîne de traitement



Fruits!





Fruits!

Version alternative (1)

- **Le projet imposait d'effectuer une PCA pour réduire la dimension,** cependant cette méthode a été une contrainte à plusieurs niveaux :
- Premièrement, j'ai dû réduire considérablement la taille du jeu de données de 67000 images à 8
- La PCA met environ 35 minutes à s'exécuter
- En prenant du recul, je suis arrivé à la conclusion que pour réduire la dimension des images, il était plus efficient d'utiliser une fonction *resize* de PIL



Fruits!

Version alternative (2)

- La fonction *resize* permet de **redimensionner** une image avec la méthode *Nearest-neighbor interpolation*.
- En procédant de cette manière, j'ai réussi à réduire les dimensions des 67 000 images en 20 x 20 en quelques secondes.



Fruits!

Conclusion

- Configurer un serveur manuellement pour utiliser Spark est un **challenge** important à cause de l'évolution rapide des versions.
- Les solutions comme Databricks permettent de déployer des modèles plus rapidement.
- La réduction de dimension avec une **PCA est gourmande en ressources de calcul**, même en distribuant les calculs.
- En utilisant une fonction de resizing uniquement, il est possible de traiter l'intégralité du jeu de données en quelques secondes.



Fruits!

Recommandations : Data Architect vs Data Scientist

Selon moi :

- Le paramétrage d'une instance pour utiliser Spark est le **coeur d'activité d'un poste de Data Architect**.
- Les Data Scientists doivent plutôt s'orienter vers des solutions comme Databricks pour **gagner du temps sur le déploiement** de solution Big Data afin de **se concentrer sur le pipeline** de traitement des données.