



# Fruits!

**EA Stephen**

Data scientist

**Déployer un modèle dans le cloud**



# Fruits!

## SOMMAIRE

- 1- Introduction
- 2- Scalabilité
- 3- L'architecture et environnement
- 4- Les étapes de la chaîne de traitement
- 5- Conclusion



# Fruits!

## 1. Introduction

### 1.1 Mission

### 1.2 Le jeu de données

Fruits est une start-up de l'Agritech souhaitant mettre à disposition une application mobile permettant au grand public d'avoir un moteur de classification d'images de fruits/légumes.

Développer dans un environnement Big data une première chaîne de traitement des données (preprocessing des images et réduction dimension)

Il faut prendre en compte que le volume des données va augmenter rapidement après la livraison de ce projet.

Il faut donc adapter une architecture cloud en fonction du volume et du script pyspark pour la première chaîne de traitement

Le nombre total d'images : 90483.

**Taille de l'ensemble d'entraînement : 67 692  
images (un fruit ou un légume par image).**

Taille de l'ensemble de test : 22 688 images (un  
fruit ou un légume par image).

Le nombre de classes : 131 (fruits et légumes).

Taille de l'image : 100 x 100 pixels.





# Fruits!

## 2. Scalabilité

**2.1 Evaluation du besoin de FRUITS par rapport au 3V's**

**2.2 Evolution des architectures proposées pour FRUITS**

### • The 3 V's of Big Data

#### Volume

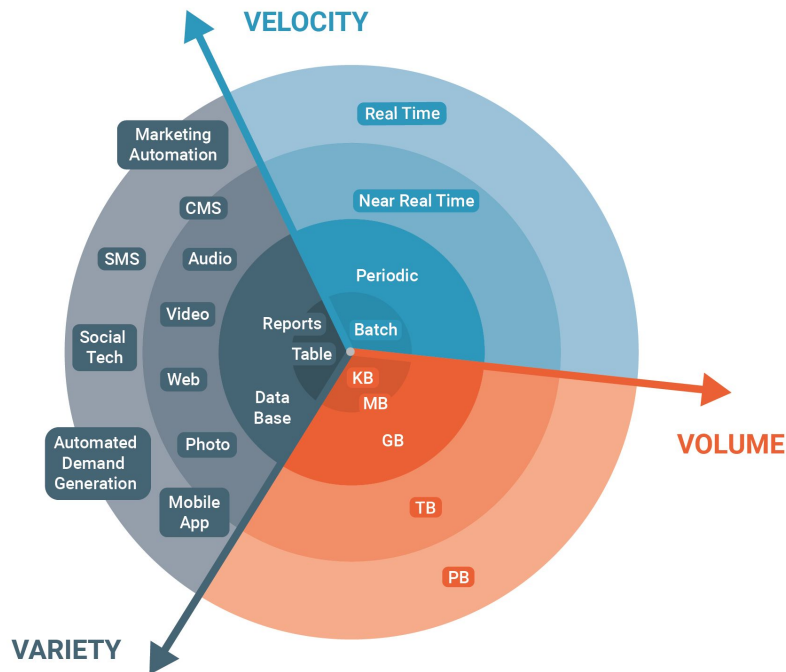
La quantité de données astronomiques générées par les entreprises est en constante augmentation

#### Vélocité

La vélocité fait référence à la vitesse à laquelle les données sont générées. Les données à grande vitesse sont générées à un rythme tel qu'elles nécessitent des techniques de traitement distribuées (hadoop, Spark).

#### Variété

Les données ont des types différents structurés (databases), semi-structurés (Csv, Json, Xml), non structurés (images, textes, vidéo, audios)



### single instance



**Scaling Up:** Augmentation du type d'instance (besoins de calcul, de mémoire, de réseaux ou de stockage.)



**Avantage :** Le coût



**Désavantage:** Possible latence



**Désavantage :** pas tolérance au pannes, distribution local spark, Il ne nécessite aucun gestionnaire de ressources. ( Spark s'exécute localement avec autant de threads de travail que de cœurs logiques sur votre machine.)

### Cluster EC2



**Scaling out :** Augmentation cluster besoin de calcul, de mémoire, de réseaux ou de stockage.



**Avantages :** Le coût comparé à serverless



tolérance au pannes , moins de latence, distribution des calculs sur plusieurs machines (hadoop, spark), répliquions des données partitionnées



**Désavantages :** Mis en oeuvre plus complexe (docker swarm, Kubernetes)

### EMR (serverless)



**Scaling out :** Augmentation cluster besoins de calcul, de mémoire, réseaux ou de stockage.



**Désavantages :** le coût



**Avantages:** moins de latence, distribution des calculs sur plusieurs machines (hadoop, spark)



**Tolérance au pannes ,** répliquions des données partitionnées



**Préconfiguré** (pas de gestion serveurs, pas d'installation , pas configuration, pas de maintenance)





# Fruits!

## 3. Architecture et environnement

- 3.1 EC2 single
- 3.2 EC2 cluster
- 3.3 EMR



# Fruits!

## 3.1 EC2 single

3.1.1 L'instance

3.1.2 Environnement containeriser

3.1.3 La configuration de l'environnement big data

## 3.1.1 L'instance

### EC2 description



- ❑ AMI (Ubuntu server 22.04)



Amazon EC2

- ❑ Type de l'instance (t2.xlarge)
  - 4 vCPU
  - 16 Gio Mémoire
  - 30 Gio stockage



- ❑ Paire de clés (connexion SSH)



- ❑ Groupe de sécurité port
  - Port 22 (SSH)
  - Port 8888 (Jupyterlab)
  - Protocole TCP
  - Source IP 0.0.0.0/0 (entrante, sortante)



- ❑ Rôle IAM
  - Full Access Aws service et ressource



Elastic IP

- ❑ Non
  - IP dynamique



- ❑ Région
  - eu-west-3

## 3.1.2 Environnement containeriser

1

```
1 # Copyright (c) Jupyter Development Team.
2 # Distributed under the terms of the Modified BSD License.
3 ARG OWNER=jupyter
4 ARG BASE_CONTAINER=$OWNER/scipy-notebook
5 FROM $BASE_CONTAINER
6
7 LABEL maintainer="Jupyter Project <jupyter@googlegroups.com>"
8
9 # Fix DL4006
10 SHELL ["bin/bash", "-o", "pipefail", "-c"]
11
12 USER root
13
14 # Spark dependencies
15 # Default values can be overridden at build time
16 # (ARGS are in lower case to distinguish them from ENV)
17 ARG spark_version="3.2.1"
18 ARG hadoop_version="3.2"
19 ARG spark_checksum="145ADACF189FECF05FBA369841D2084DD66546B1D14FC181AC49D89F3C85E4FEC09B25F56F0A8767155419CD430838FB651992AEB37D3A6F91E7E0B9D1F9AE"
20 ARG openjdk_version="8"
21
22 ENV APACHE_SPARK_VERSION="${spark_version}" \
23     HADOOP_VERSION="${hadoop_version}"
24 ENV JUPYTER_ENABLE_LAB=yes
25 RUN apt-get update --yes && \
26     apt-get install --yes --no-install-recommends \
27     "openjdk-${openjdk_version}-jre-headless" \
28     ca-certificates-java && \
29     apt-get clean && rm -rf /var/lib/apt/lists/*
30
31 # Spark installation
32 WORKDIR /tmp
33 RUN wget -q "https://archive.apache.org/dist/spark/spark-${APACHE_SPARK_VERSION}/spark-${APACHE_SPARK_VERSION}-bin-hadoop${HADOOP_VERSION}.tgz" && \
34     echo "${spark_checksum} *spark-${APACHE_SPARK_VERSION}-bin-hadoop${HADOOP_VERSION}.tgz" | sha512sum -c - && \
35     tar xzf "spark-${APACHE_SPARK_VERSION}-bin-hadoop${HADOOP_VERSION}.tgz" -C /usr/local --owner root --group root --no-same-owner && \
36     rm "spark-${APACHE_SPARK_VERSION}-bin-hadoop${HADOOP_VERSION}.tgz"
37
38 WORKDIR /usr/local
```

**Image jupyter-notebook**

**Spécifier version spark, hadoop, Installer java (compatible dépendances)**

**Installer spark**

2

```
# Configure Spark
ENV SPARK_HOME=/usr/local/spark
ENV SPARK_OPTS="--driver-java-options=-Xms1024M --driver-java-options=-Xmx4096M --driver-java-options=-Dlog4j.logLevel=info" \
    PATH="${PATH}:${SPARK_HOME}/bin"

RUN ln -s "spark-${APACHE_SPARK_VERSION}-bin-hadoop${HADOOP_VERSION}" spark && \
    # Add a link in the before_notebook hook in order to source automatically PYTHONPATH
    mkdir -p /usr/local/bin/before-notebook.d && \
    ln -s "${SPARK_HOME}/sbin/spark-config.sh" /usr/local/bin/before-notebook.d/spark-config.sh

# Fix Spark installation for Java 11 and Apache Arrow library
# see: https://github.com/apache/spark/pull/27356, https://spark.apache.org/docs/latest/#downloading
RUN cp -p "${SPARK_HOME}/conf/spark-defaults.conf.template" "${SPARK_HOME}/conf/spark-defaults.conf" && \
    echo 'spark.driver.extraJavaOptions -Dio.netty.tryReflectionSetAccessible=true' >> "${SPARK_HOME}/conf/spark-defaults.conf" && \
    echo 'spark.executor.extraJavaOptions -Dio.netty.tryReflectionSetAccessible=true' >> "${SPARK_HOME}/conf/spark-defaults.conf"

WORKDIR /jars
ENV SPARK_JARS=/usr/local/spark/jars

RUN wget https://repo1.maven.org/maven2/com/amazonaws/aws-java-sdk-bundle/1.12.226/aws-java-sdk-bundle-1.12.226.jar

# Get Hadoop-AWS Jar
RUN wget https://repo1.maven.org/maven2/org/apache/hadoop/hadoop-aws/3.3.1/hadoop-aws-3.3.1.jar

# Get jets3t JAR
RUN wget https://repo1.maven.org/maven2/net/java/dev/jets3t/jets3t/0.9.4/jets3t-0.9.4.jar

USER ${NB_UID}
```

**Configuration spark**

**Télécharger fichier jar pour accéder à S3 avec les fonctionnalités de spark (aws-sdk-bundle, hadoop-aws, jets3t)**

3

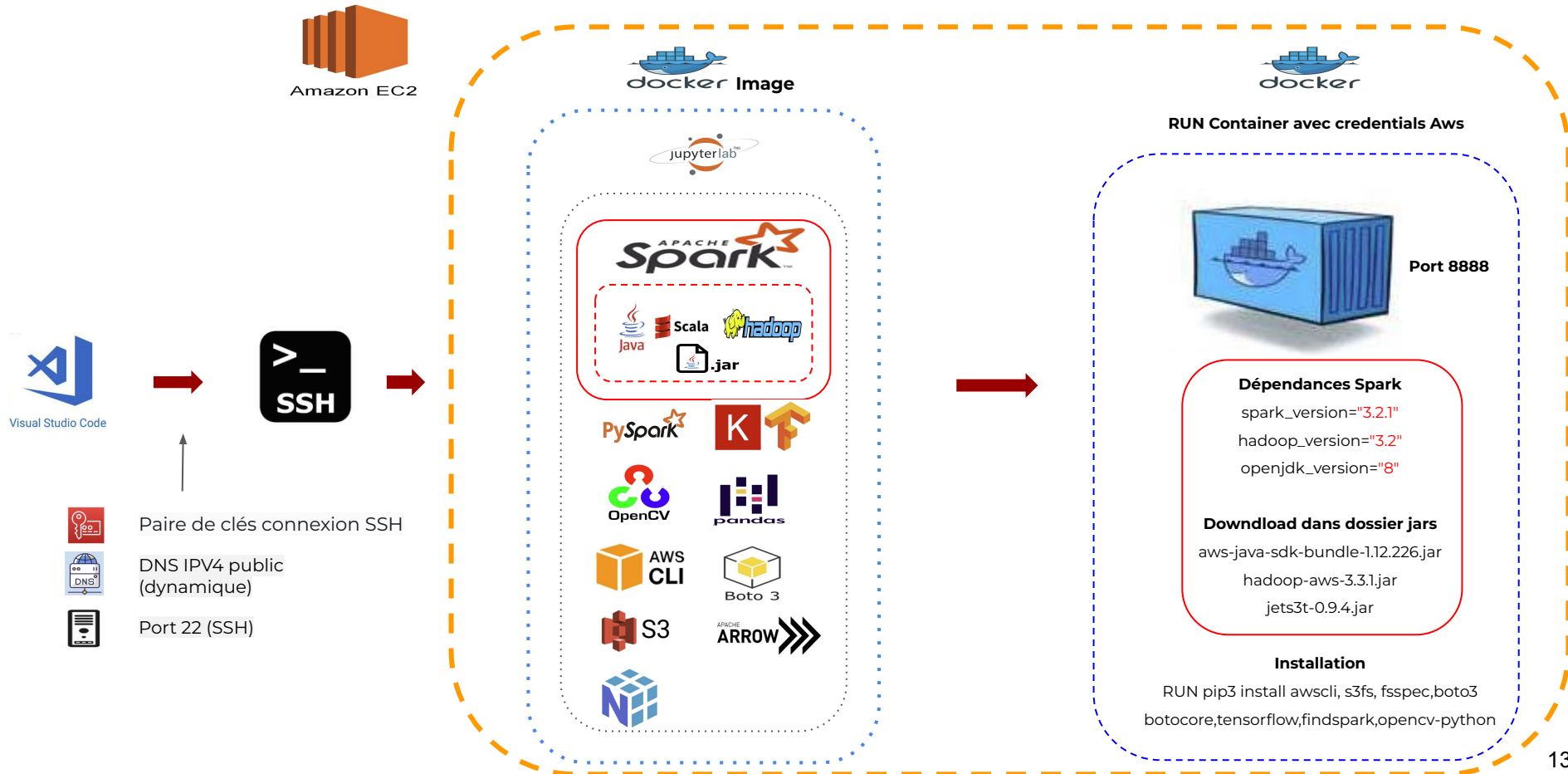
```
# Install pyarrow
RUN mamba install --quiet --yes \
    'pyarrow' && \
    mamba clean --all -f -y && \
    fix-permissions "${CONDA_DIR}" && \
    fix-permissions "/home/${NB_USER}"

RUN pip3 install awscli
RUN pip3 install s3fs
RUN pip3 install fsspec
RUN pip3 install boto3
RUN pip3 install botocore

WORKDIR "${HOME}"
```

**Install autres packages**

## 3.1.3 La configuration de l'environnement big data





# Fruits!

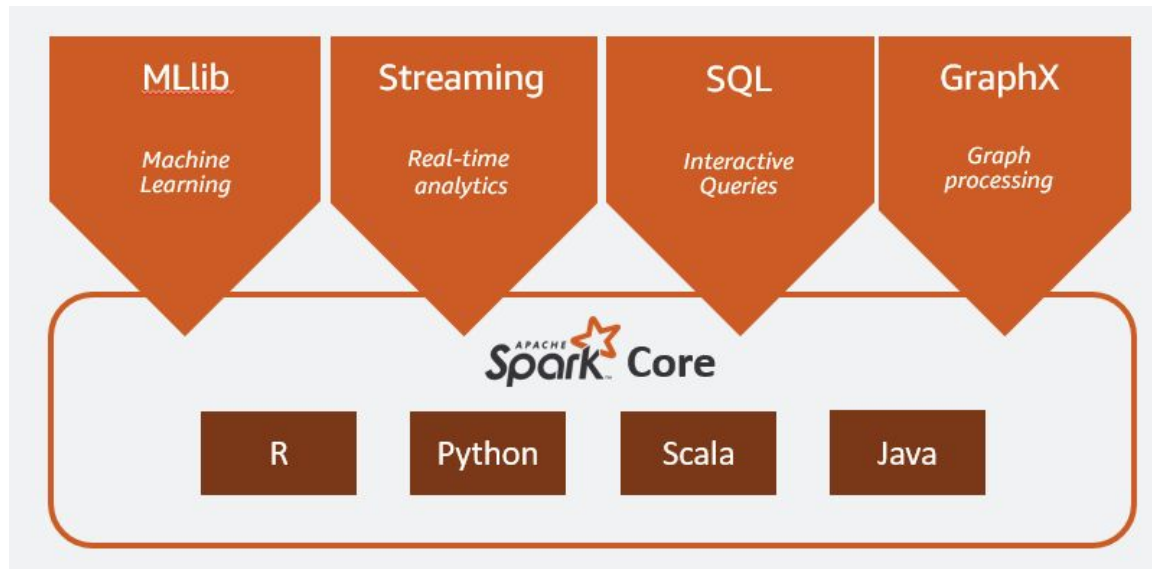
## 3.2 EC2 cluster

- 3.2.1 Apache Spark eco system
- 3.2.2 Les structures de données dans Spark
- 3.2.3 Comment fonctionne Spark
- 3.2.4 L'architecture big data

## 3.2.1 Apache Spark eco system

Le framework Spark comprend :

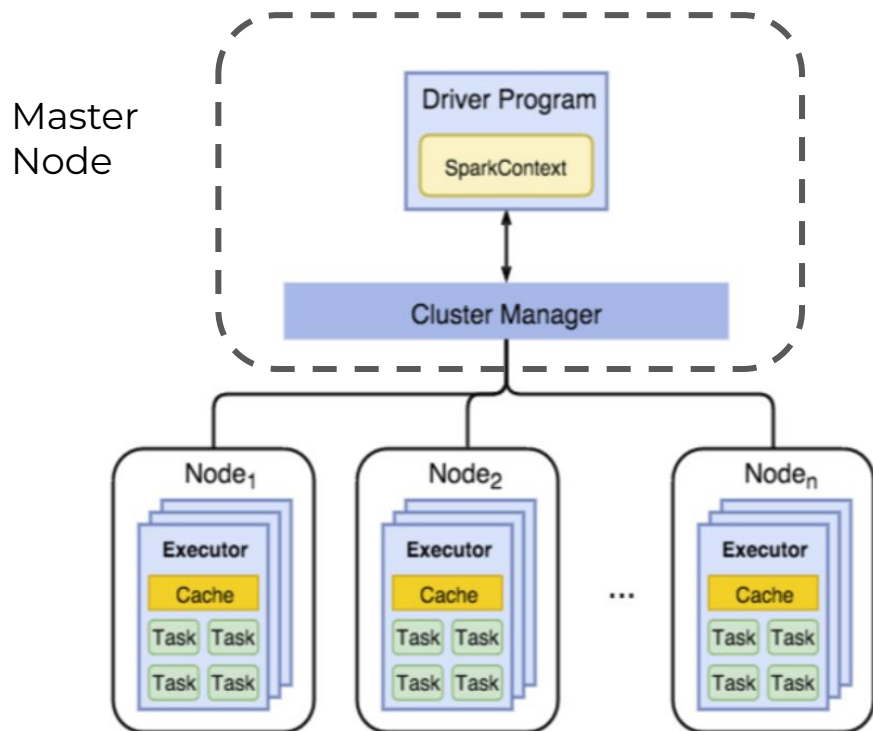
- Spark Core comme base de la plateforme
- Spark SQL pour les requêtes interactives
- Spark Streaming pour des analyses en temps réel
- Spark MLlib pour l'apprentissage automatique
- Spark GraphX pour le traitement de graphes



## 3.2.2 Les structures de données dans Spark

	RDD	Dataframe	Dataset
<b>Immutable</b>	Oui	Oui	Oui
<b>Tolérance panne</b>	Oui	Oui	Oui
<b>Type safe</b>	Oui	Non	Oui
<b>Optimizer</b>	Non	Catalyst Engine	Catalyst Engine
<b>Memory usage</b>	High	Faible en raison de la mémoire hors tas binaire Tungsten	Faible en raison de la mémoire hors tas binaire Tungsten
<b>Exécution</b>	Lente	Augmenté grâce à des plans d'exécution optimisés à l'aide d'un catalyseur Engine	Augmenté grâce à des plans d'exécution optimisés à l'aide d'un catalyseur Engine
<b>Lazy evaluation</b>	Oui	Oui	Oui
<b>Schéma</b>	Définir manuellement schéma	Automatique	Automatique





### Driver program (pilote)

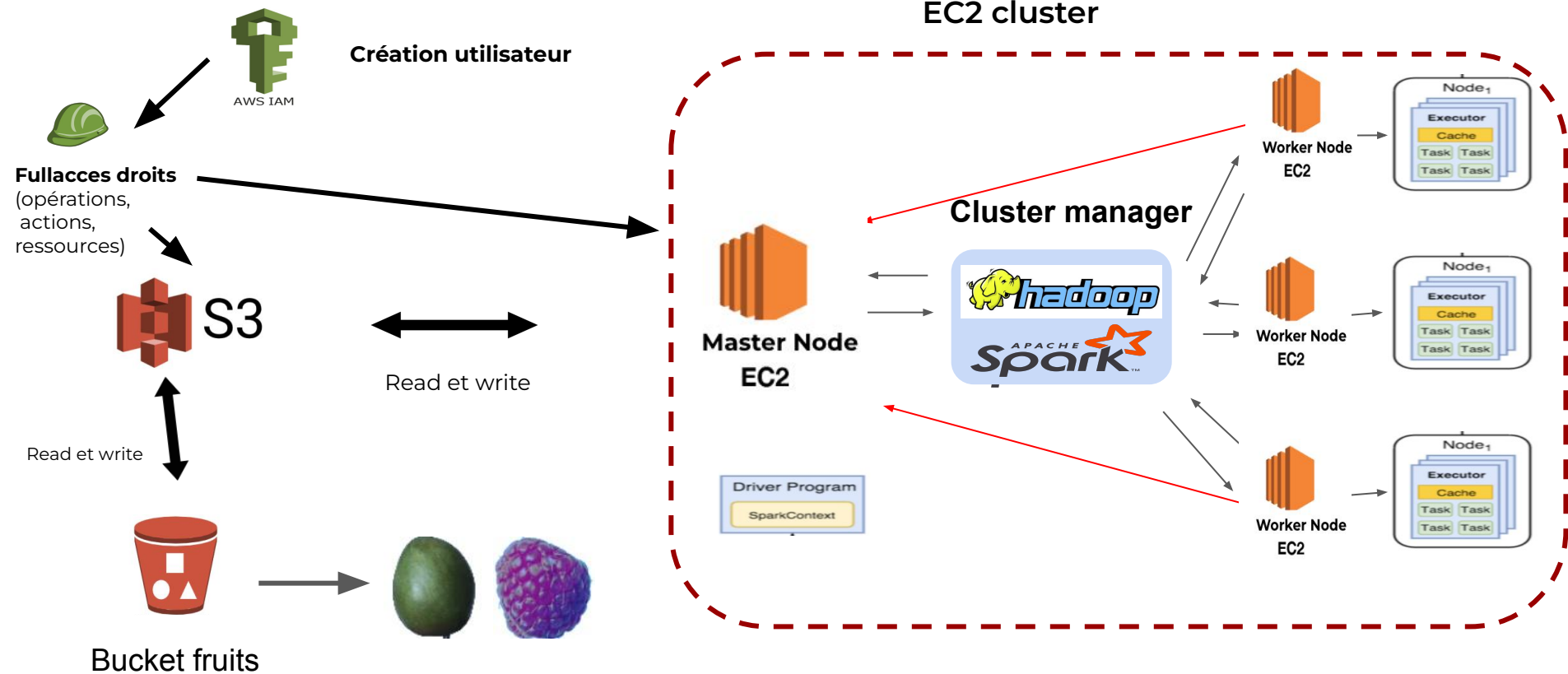
- Le pilote est composé du programme, et une session Spark (entrée vers les fonctionnalités de spark). Il est le responsable de l'exécution du code sur les différentes machines.

### Cluster manager

- Gère l'allocation de ressources, la division des programmes, l'exécution du programme (Standalone, Yarn, Mesos)

### Workers nodes

- Chaque exécuteur, ou nodes Worker, reçoit une tâche du pilote et exécute cette tâche.
- Les données sont répliquées et conservées en mémoire vive (cache) afin que l'application soit tolérante au panne( Spark s'exécute 100 fois plus vite en mémoire et 10 fois plus vite sur disque que hadoop)





# Fruits!

## 3.3 EMR

3.3.1 L'architecture big data EMR

3.3.2 Auto scaling EMR

### Détails de configuration

Étiquette de version : emr-6.7.0

Distribution Hadoop : Amazon

Applications : TensorFlow 2.4.1, Spark 3.2.1,  
JupyterEnterpriseGateway 2.1.0, JupyterHub 1.4.1

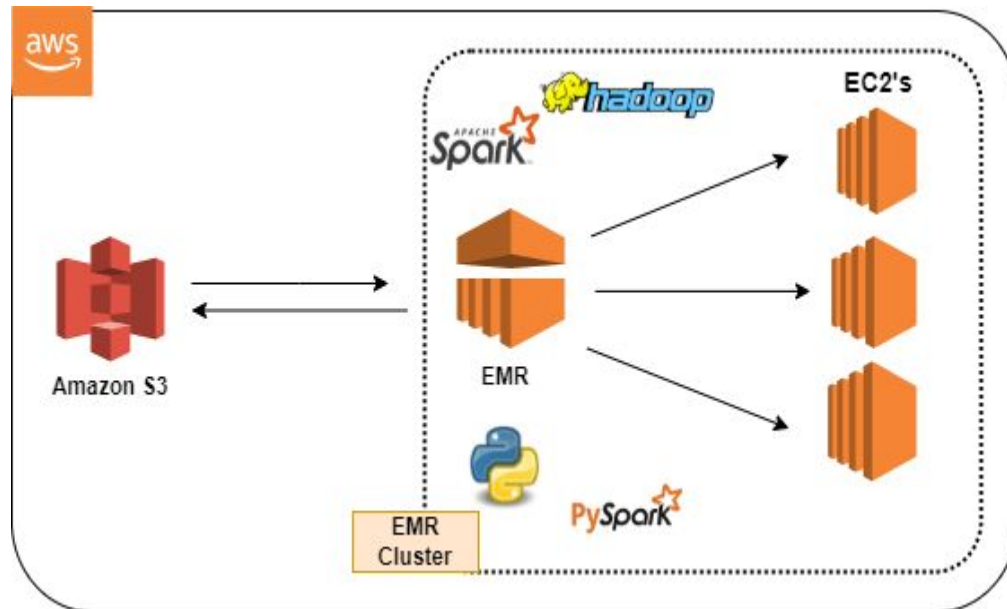
URI de connexion : s3://aws-logs-933635246190-eu-west-3/elasticmapreduce/ 

Vue cohérente EMRFS : Désactivé

ID d'AMI personnalisée : --

Version d'Amazon Linux : 2.0.20220719.0 [En savoir plus](#) 

Nom et type de nœud	Type d'instance
<b>MASTER</b>	<b>m5.xlarge</b>
Groupe d'instances maître - 1	4 Cœurs virtuels, 16 Gio de mémoire, stockage EBS uniquement Stockage sur EBS : 64 Gio
<b>CORE</b>	<b>m5.xlarge</b>
Groupe d'instances principal - 2	4 Cœurs virtuels, 16 Gio de mémoire, stockage EBS uniquement Stockage sur EBS : 64 Gio



### Configuration du matériel

Type d'instance  Le type d'instance sélectionné ajoute un volume EBS GP2 par défaut de 64 Gio par instance. [En savoir plus](#)

Nombre d'instances  (1 nœud maître et 2 nœuds principaux)

Mise à l'échelle du cluster ☒ scale cluster nodes based on workload

#### EMR-managed scaling

EMR will automatically increase and decrease the number of instances in core and task nodes based on workload. Set a minimum and maximum limit of the number of instances for the cluster nodes. Master nodes do not scale. [Learn more](#)

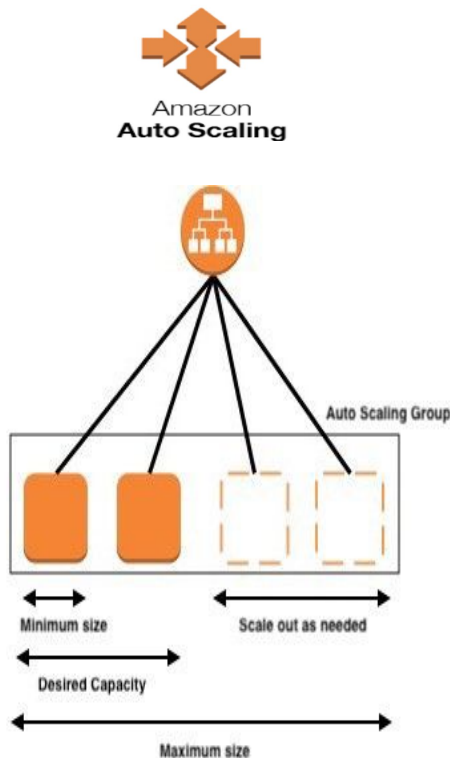
#### Unités principales et unités de tâches

Minimum :

Maximum :

Résiliation automatique ☒ Activer la résiliation automatique [En savoir plus](#)

Arrêter le cluster lorsqu'il est inactif après  heures  minutes



**EMR augmentera et diminuera automatiquement le nombre d'instances dans les nœuds principaux et de tâches en fonction de la charge de travail. Définissez une limite minimale et maximale du nombre d'instances pour les nœuds de cluster. Les nœuds maîtres ne sont pas mis à l'échelle.**



# Fruits!

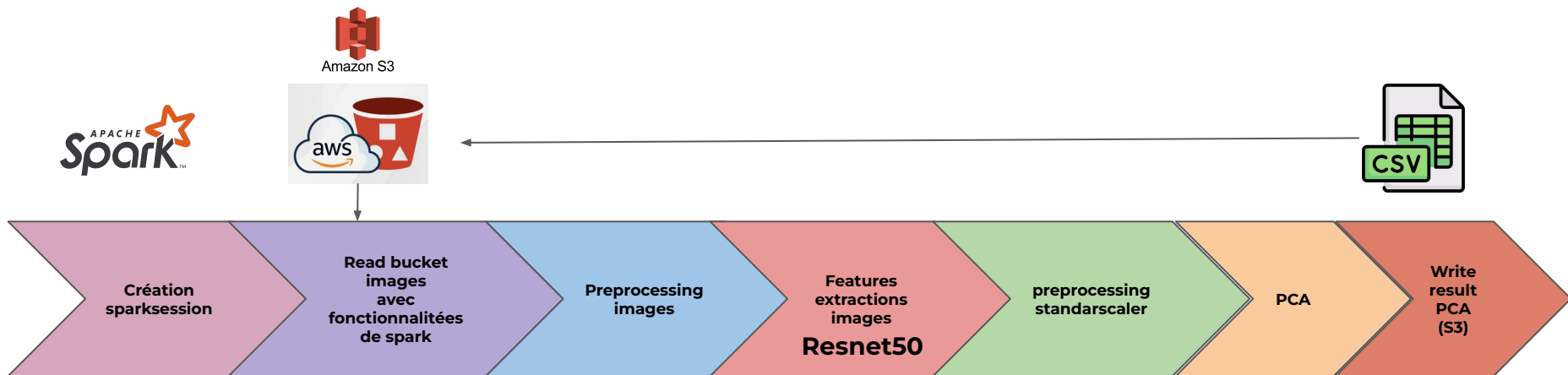
## 4. Chaîne de traitement

4.1 Les étapes de la chaîne de traitement

4.2 Features extraction Resnet50

4.3 Visualisation PCA

4.4 Write résultat PCA dans data lake (S3)



IDENTIFIANT	ID d'application YARN	Général	État	Interface utilisateur Spark	Journal du conducteur	Utilisateur	Séance actuelle?
0	application_1665586751956_0001	cvsark	inactif	<a href="#">lien</a>	<a href="#">lien</a>	Auon	✓

```
def read_S3_buckets_img(s3path):
    s3path = s3path

    # Read images avec le format binaryFile
    df_images = spark.read.format("binaryFile") \
        .option("pathGlobFilter", "*.jpg") \
        .load(s3path)
    df_images.printSchema()

    return df_images
```

Resize

100x100



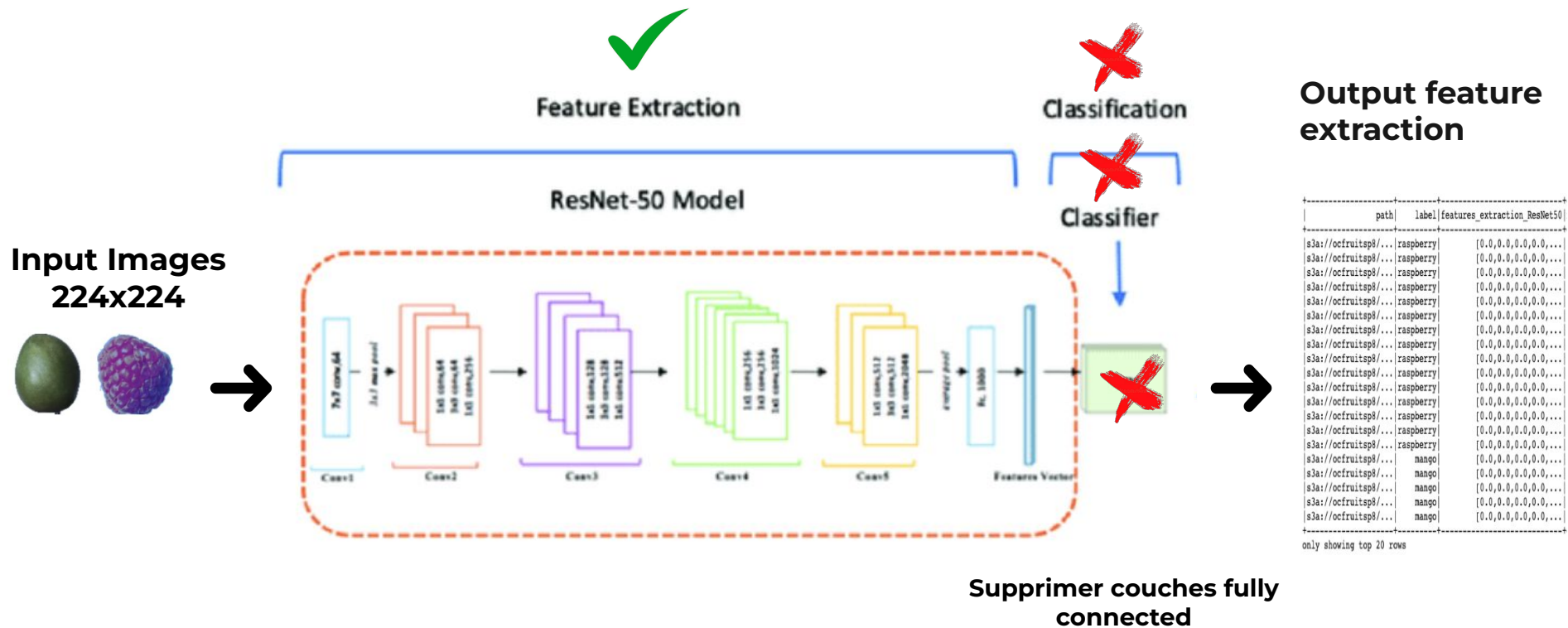
224x224



	path	label	features_extraction_ResNet50	features_scaled_StandardScaler	features_pca
s3a://ocfruitsp8/...	raspberr...		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[98.3656270320505...
s3a://ocfruitsp8/...	raspberr...		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[88.7803466332628...
s3a://ocfruitsp8/...	raspberr...		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[95.7223947966706...
s3a://ocfruitsp8/...	raspberr...		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[104.156610510264...
s3a://ocfruitsp8/...	raspberr...		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[87.7651299642067...
s3a://ocfruitsp8/...	raspberr...		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[90.9680860429685...
s3a://ocfruitsp8/...	raspberr...		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[89.8776923367681...
s3a://ocfruitsp8/...	raspberr...		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[92.0226157636022...
s3a://ocfruitsp8/...	raspberr...		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[89.4623227934633...
s3a://ocfruitsp8/...	raspberr...		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[95.6694779208975...
s3a://ocfruitsp8/...	raspberr...		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[89.8530285023048...
s3a://ocfruitsp8/...	raspberr...		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[91.8247326805036...
s3a://ocfruitsp8/...	raspberr...		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[94.3949962773765...
s3a://ocfruitsp8/...	raspberr...		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[93.8360389146309...
s3a://ocfruitsp8/...	raspberr...		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[91.4965047227552...
s3a://ocfruitsp8/...	mango		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[-84.965653169756...
s3a://ocfruitsp8/...	mango		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[-86.810272337467...
s3a://ocfruitsp8/...	mango		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[-84.428132693828...
s3a://ocfruitsp8/...	mango		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[-88.516856670406...
s3a://ocfruitsp8/...	mango		[0.0,0.0,0.0,0.0,...	[0.0,0.0,0.0,0.0,...	[-106.13608345547...

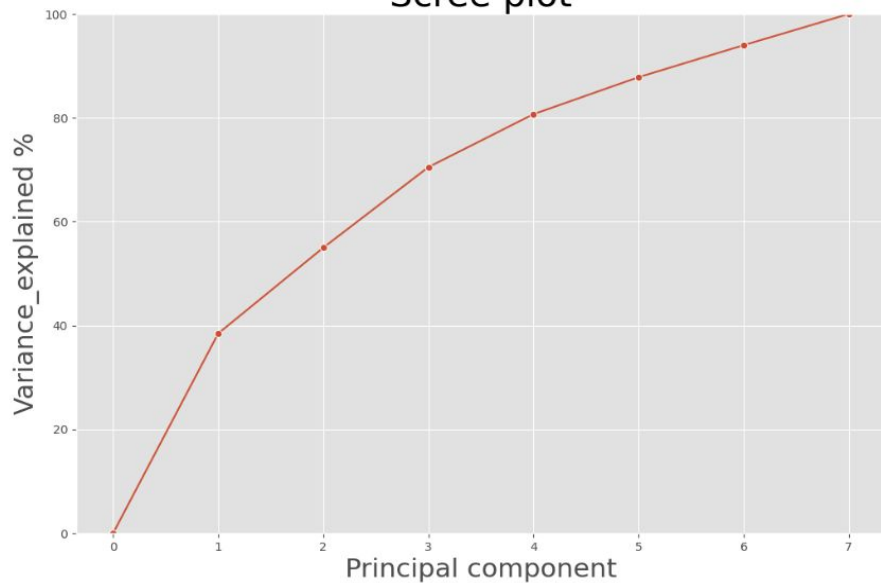
only showing top 20 rows

## 4.2 Features extractions Resnet50

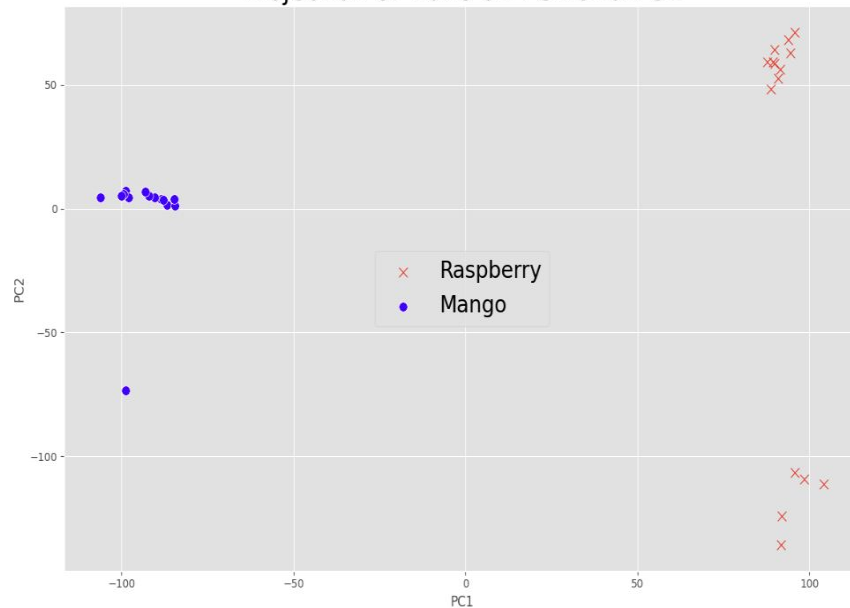




Scree plot



Projection of fruits on PC1 and PC2



La variance expliquée sur un plan 2d (2 components) est de 55%

La variance pour chaque composante : [0.3847107 0.54991919 0.70502387 0.80705559 0.87827463 0.93989402 1. ]

# 4.4 Résultat sauvegarder dans S3

Amazon S3 > Buckets > bucketsavepca > pca\_result/

pca\_result/

Objects

Properties

Objects (1)

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

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	pca_result.csv	csv	September 19, 2022, 14:55:01 (UTC+02:00)	5.7 KB	Standard

Copy S3 URI



pca\_result

path	label	features_pca
s3a://odnfruits8/fruits_em/raspberry/12_100.jpg	raspberr	[98.3656270205053,-109.32981336579955,-25.63056476127914,-68.18995732607699,-7.01318457234865,62.647423967992836,103.5232903032754]
s3a://odnfruits8/fruits_em/raspberry/122_100.jpg	raspberr	[88.78034663326284,48.41495427042843,12.025070453467311,-6.88188333098825,104.49574646740214,-4.299763485124397,6.4878242624415]
s3a://odnfruits8/fruits_em/raspberry/11_100.jpg	raspberr	[56.722947966707,-106.50930570287791,-25.27343200363455,-7.17201881688206,-12.438785003255878,28.67947328428462,85.72690022580195]
s3a://odnfruits8/fruits_em/raspberry/12_100.jpg	raspberr	[104.15661051028427,-111.18585027980257,-22.20723782392402,-150.67960274396002,-9.63176217430504,-76.2809070245423,-114.5989803445892]
s3a://odnfruits8/fruits_em/raspberry/120_100.jpg	raspberr	[87.7651299420676,59.188563021846644,13.871435447545327,-1.483275149657166,58.945045033181914,-3.703671685342076,1.32046132046446]
s3a://odnfruits8/fruits_em/raspberry/121_100.jpg	raspberr	[90.989806429855,52.5259967513743,13.29833658954394,5.368020961562979,108.4091349842089,4.32942008525379,5.789148071491941]
s3a://odnfruits8/fruits_em/raspberry/119_100.jpg	raspberr	[89.8778923367681,58.557588381805196,13.195911830530626,1.3439625199923875,28.716901102350287,-3.6201719244778023,1.64341889730397596]
s3a://odnfruits8/fruits_em/raspberry/1_100.jpg	raspberr	[92.022617636022,-124.20349047817432,-26.19578273248906,125.79891026327515,2.5893727816005176,-16.51603926273897,-27.822336296461174]
s3a://odnfruits8/fruits_em/raspberry/11_100.jpg	raspberr	[89.4023276340339,59.124612809117,12.84632242075805,1.854033026854868,25.408641988334693,-1.448372121813576,-5.894189463168727]
s3a://odnfruits8/fruits_em/raspberry/116_100.jpg	raspberr	[95.66947792088759,71.29659793219349,16.331522331316666,9.869535411027455,-47.74885588218335,-44.15364903830798,-4.876447870613723]
s3a://odnfruits8/fruits_em/raspberry/115_100.jpg	raspberr	[85.8302850230487,64.1254051461187,14.74341952721627,6.004771193650225,-50.08865884474984,-1.836967814175747,-3.542635623338967]
s3a://odnfruits8/fruits_em/raspberry/9_100.jpg	raspberr	[91.8247326850362,-135.8760729178806,-28.887140401170694,140.79775293509056,9.239178624470745,-14.24826911109806,-13.97203261942791]
s3a://odnfruits8/fruits_em/raspberry/113_100.jpg	raspberr	[94.39496627737656,62.7064476259832,16.2487633447116,9.688047607062863,-43.5005588363063,3.415348214221824,0.2796223796069265]
s3a://odnfruits8/fruits_em/raspberry/117_100.jpg	raspberr	[93.83602891463093,68.203684118533,15.851379443304014,12.988050602233505,-41.7610982508149,1.043342049202695,-4.086516133725723]
s3a://odnfruits8/fruits_em/raspberry/114_100.jpg	raspberr	[91.496042272526,56.1481951128432,12.786688949736488,11.043653034859405,-79.90311151688904,4.31505084749294,1.483058766527961]
s3a://odnfruits8/fruits_em/mango/102_100.jpg	mango	[54.9665316975651,2.2728141959192973,-17.088924674561445,-1.1585650931654646,2.7406430262579872,40.15126116135094,-23.985927862847895]
s3a://odnfruits8/fruits_em/mango/101_100.jpg	mango	[96.8102723746703,1.43982379670555,-14.885021404476592,-1.638437895488705,4.833057184301963,54.0599393039577,-14.4008949010208]
s3a://odnfruits8/fruits_em/mango/100_100.jpg	mango	[48.4201329808206,1.1577484254513455,-19.01421489794532,-1.58803894484869,3.913034500561563,55.2420960443834,33.999973373081]
s3a://odnfruits8/fruits_em/mango/104_100.jpg	mango	[68.516886670464,3.80017050593372,-22.79558304575613,-0.943584287913335,2.7948817179813487,45.87427425508929,-29.73309337965399]
s3a://odnfruits8/fruits_em/mango/103_100.jpg	mango	[106.1360534547664,4.571158242155822,-19.95949607429373,1.567190461177862,-4.822269388045405,45.9162542314374337,56.03591966547]
s3a://odnfruits8/fruits_em/mango/106_100.jpg	mango	[90.3278415693244,4.52871797276062,-23.310891972214384,-0.954668706479843,0.2307813112692195,20.242302831903362,-12.32171066246975]
s3a://odnfruits8/fruits_em/mango/105_100.jpg	mango	[67.8844869433849,3.64392121320385,-21.015356579753883,-1.942031413978144,2.335021186972143,46.0552097138329,-26.613833687171534]
s3a://odnfruits8/fruits_em/mango/105_100.jpg	mango	[84.7170763045981,3.9265760814619604,-25.780870653796554,-1.2827643379178668,0.5680234227407708,26.42229291890231,-16.85342910530895]
s3a://odnfruits8/fruits_em/mango/110_100.jpg	mango	[98.79319590623418,7.127576697238092,-27.20246878858103,-0.4359849188741834,-4.183013176375659,-43.21017530473461,27.74478916502565]
s3a://odnfruits8/fruits_em/mango/108_100.jpg	mango	[98.0479809205212,4.434196387566206,-22.4694414166402,-0.02916599989471264,-3.1425991523201036,-18.587834404402812,8.878038183033763]
s3a://odnfruits8/fruits_em/mango/111_100.jpg	mango	[99.43484973255165,9.47361492154818,-25.96097557867991,-1.6376531368119718,-4.57142602654242,-50.7805563377839,31.74698260449984]
s3a://odnfruits8/fruits_em/mango/107_100.jpg	mango	[92.16145902164784,5.20330689919983,-26.835914416224,-0.3047497459274135,-0.5414487988153,8.42222746480348,-6.498935116371755]
s3a://odnfruits8/fruits_em/mango/109_100.jpg	mango	[93.1275482081308,6.793301698732098,-28.13785841064593,-0.958644028389452,-1.3387659944280218,-11.411796471910387,6.103476235564138]
s3a://odnfruits8/fruits_em/mango/112_100.jpg	mango	[100.4656055937708,5.29641543250045,-22.08169437591003,-0.7297394493350813,-0.812300378831283,-66.3887002053838,34.843820659904]
s3a://odnfruits8/fruits_em/mango/13_100.jpg	mango	[98.7977458383773,-73.335206626269,30.58247323032184,-1.5404411505780955,-0.8341115051902498,-0.6824948924523592,2.477181525708118]

### La mission

#### Les architectures utilisées

- EC2 single (mode Local spark)
- EC2 cluster (mode standalone spark)
- EMR (préconfiguré Yarn)

#### Preprocessing et PCA

- Les images sont séparées lors de la projection en 2D

### Amélioration possible

- Mise à l'échelle autoscaling EMR
- Preprocessing images
- Tester d'autres modèles features extractions et hyperparamètres
- Entraîner un modèle de classification avec les features extractions
- Approfondir connaissance Spark et service AWS