

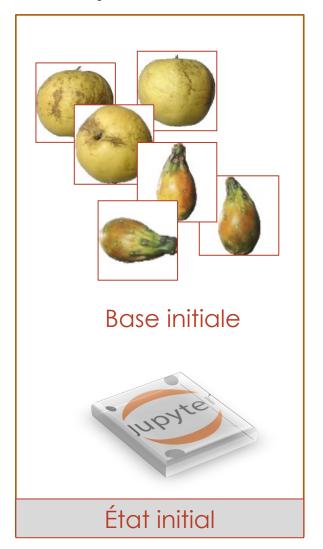
Préparation à la reconnaissance d'image

Déploiement du modèle dans le cloud

Introduction

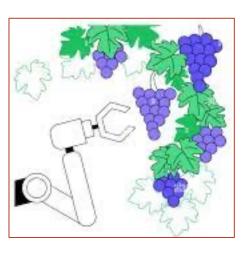
2/15

Objectifs:









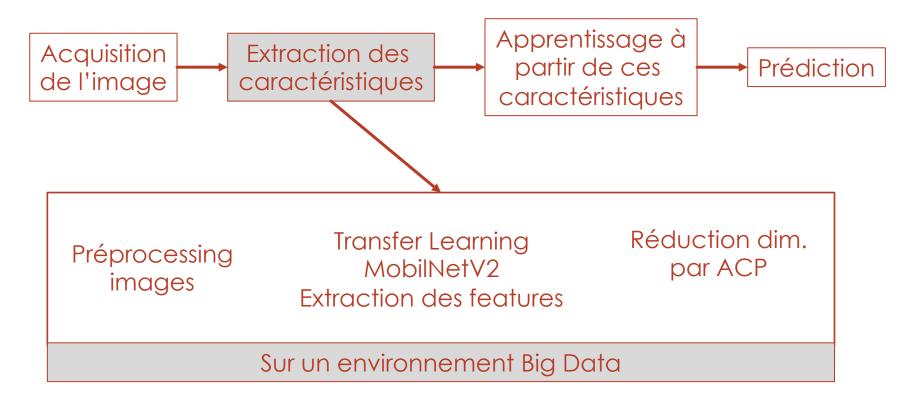
Utilisation finale





Introduction

Reconnaissance d'image :







Briques de l'environnement Big Data sur Amazon Web Service



► IAM: Gestion des utilisateurs et de leurs droits



EC2 : Instanciation des serveurs virtuels pour le calcul Création de la paire de clé rsa pour le tunnel SSH Ouverture du port 22 vers l'extérieur pour le tunnel SSH (groupe de sécurité ElasticMapReduce-master)



S3 : Création du bucket elisepoupi-p8
 Stockage du Notebook, du fichier de bootstrapping et des données



EMR: Création du cluster pour les calculs distribués par PySpark sur des serveurs EC2 Accès à JupyterHub par le tunnel SSH vers EC2 pour lancer les calculs du Notebook stocké sur S3



Fruits!

► IAM : Deux utilisateurs créés



Services	Elise	Lecteur
\$3	Contrôle total	Lecture seule
EC2	Contrôle total	Lecture seule
EMR	Contrôle total	Lecture seule
IAM	Contrôle total	



Nom d'utilisateur	Mot de passe	URL de connexion
Lecteur]%1m*A5\$	https://105083061881.signin.aws.a mazon.com/console

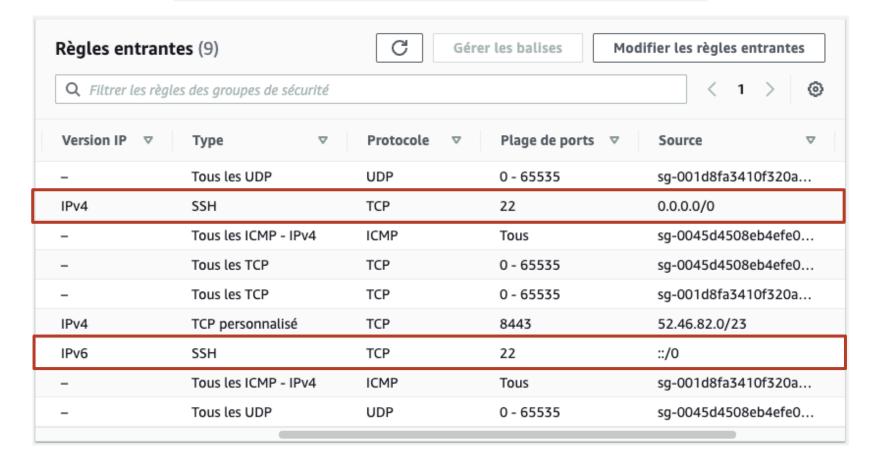




EC2:

sg-0045d4508eb4efe04 - ElasticMapReduce-master



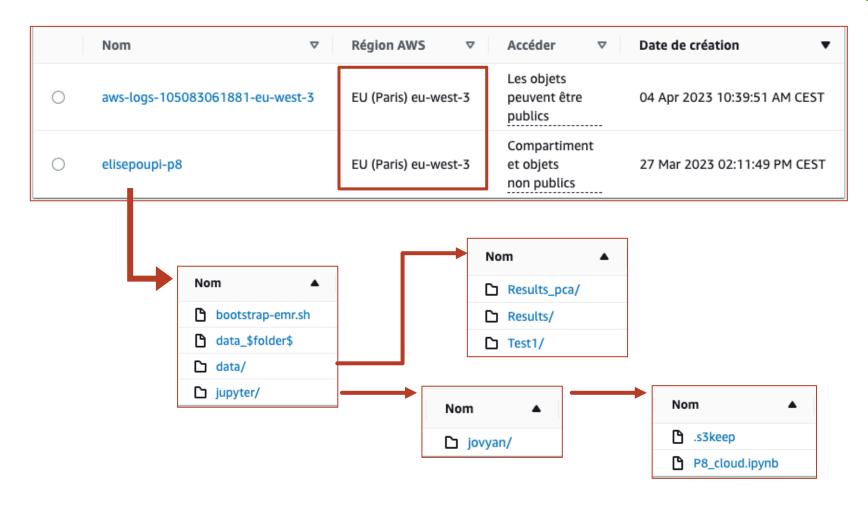






S3:









Retrait de Tensorflow suite bug

et insertion de Tensorflow dans le

fichier bootstrap-emr.sh

m5 : adapté pour les charges de travail EMR

m5.xlarge: 4 vCPU - mémoire 16 Gio

Détails de configuration EMR:

Étiquette de version: emr-6.9.0

Distribution Hadoop: Amazon 3.3.3

Applications: JupyterHub 1.4.1, Spark 3.3.0

URI de connexion : s3://aws-logs-105083061881-eu-west-

3/elasticmapreduce/

Vue cohérente EMRFS: Désactivé

ID d'AMI personnalisée : --

Version d'Amazon Linux : 2.0.20230404.0 En savoir plus [2]

Réseau et matériel

Zone de disponibilité : eu-west-3a

ID de sous-réseau (subnet): subnet-03c843f28b79cc8d7

Maître: Action d'amorçage 1 m5.xlarge

Principal: Action d'amorçage 2 m5.xlarge

Tâche: --

Cluster scaling: Custom policy

Résiliation automatique : Arrêter en cas d'inactivité pour 1 heure, 30 minutes

Sécurité et accès

Nom de clé: elise-ec2

Profil d'instance EC2 : EMR_EC2_DefaultRole

Rôle EMR: EMR DefaultRole

Auto Scaling role: EMR AutoScaling DefaultRole

















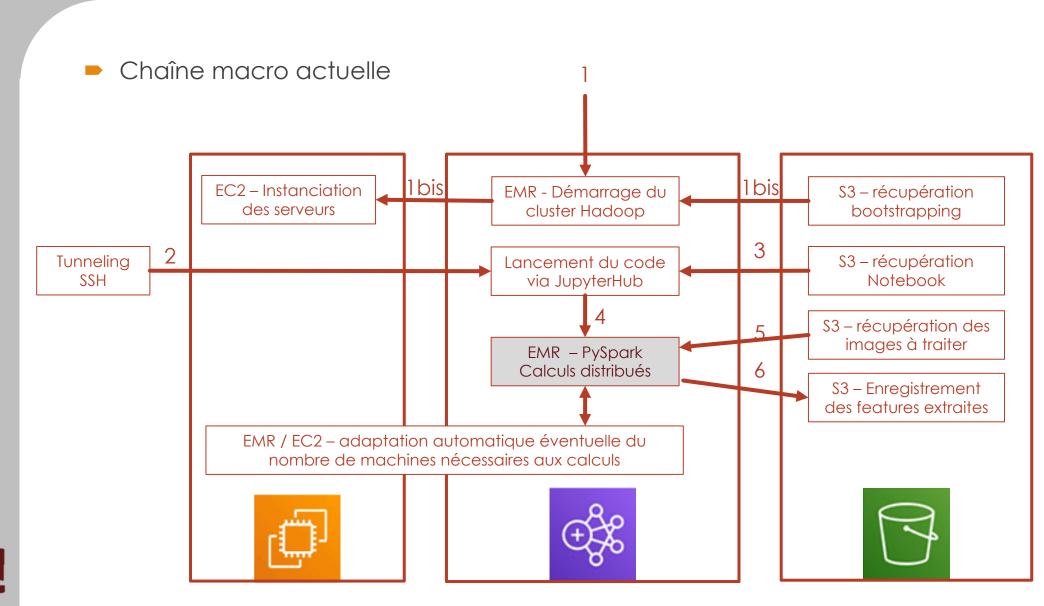
- ► Latence réduite au maximum pour les utilisateurs français
- Données conservées en France, région conforme aux normes ISO 27001, ISO 27017, ISO 27018, SOC 1 (SAS 70), SOC2 et SOC3 en matière de sécurité et disponibilité, PCI DSS niveau 1

Région	Coût de l'instance m5.xlarge
Europe (Stockholm) eu-north-1	0,204 USD/h -8,9%
Europe (Irlande) eu-west-1	0,214 USD/h -4,5%
Europe (Espagne) eu-south-2	0,214 USD/h -4,5%
Europe (Londres) eu-west-2	0,222 USD/h -0,9%
Europe (Paris) eu-west-3	0,224 USD/h -
Europe (Milan) eu-south-1	0,224 USD/h -
Europe (Francfort) eu-central-1	0,230 USD/h +2,7%
Europe (Zurich) eu-central-2	0,253 USD/h +12,9%





De l'environnement au traitement

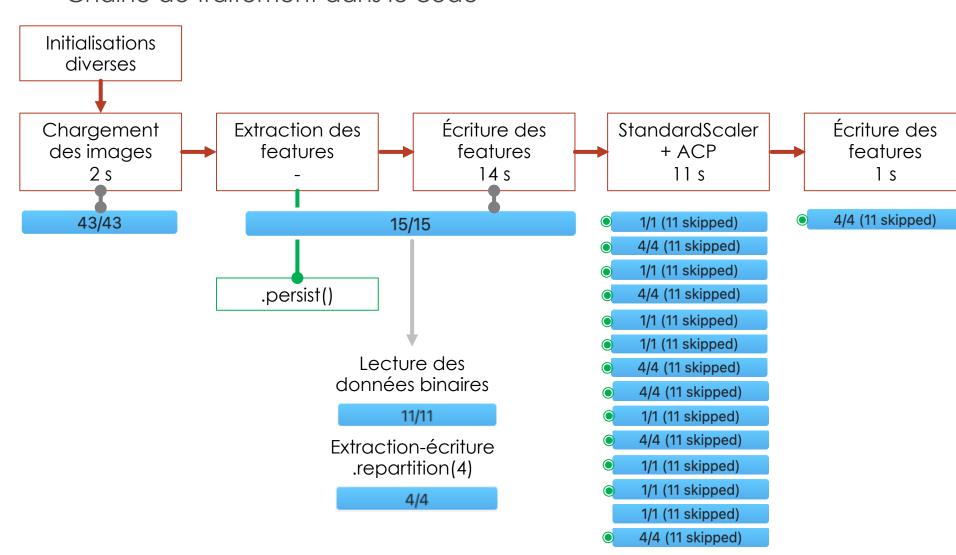




Fruits!

Chaîne de traitement des images

Chaine de traitement dans le code

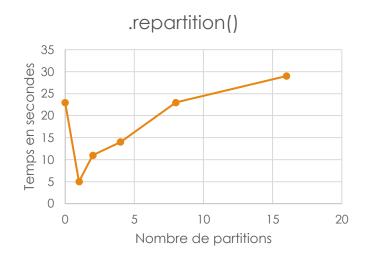






Chaîne de traitement des images

- Focus sur .repartition(n)
 - Opération dite par PySpark coûteuse en ressources et à éviter dans le cadre de l'optimisation des performances



 Ici elle reste intéressante avec peu de données, mais à surveiller par la suite et voir s'il faut laisser PySpark gérer le découpage



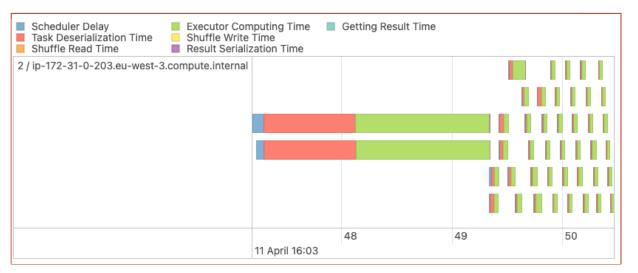


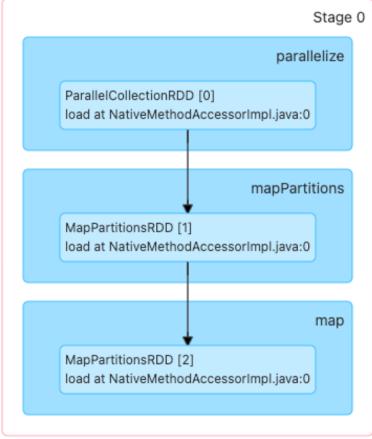
Chaîne de traitement des images

Chargement des images (43 dossiers différents ici)

en local: 0,3 secondes sur AWS: 4,0 secondes

```
images = spark.read.format("binaryFile") \
coption("pathGlobFilter", "*.jpg") \
coption("recursiveFileLookup", "true") \
load(PATH_Data)
```









Conclusion

- Solution AWS avec EC2, EMR et S3 :
 - → Tolérance aux pannes : Cloud et PySpark
 - → Maintenabilité: EMR et sa pré-configuration, séparation EMR et S3
 - → Coût faible: Cloud, EC2 et instanciations à la demande, possibilité de savings plan sur un an
 - ➤ ✓ Augmentation du volume de données : passage à l'échelle horizontal permis par Spark et EMR/EC2, grande rapidité d'actions avec retrait et ajout de machines possibles en quelques minutes en fonction des besoins







Merci pour votre attention

Place aux questions!

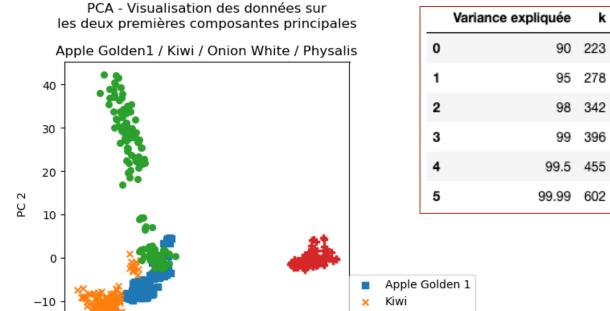
Hors propos

PCA: sur 4 types de fruits

-20

-10

602 variables pour une variance expliquée à 99,99%

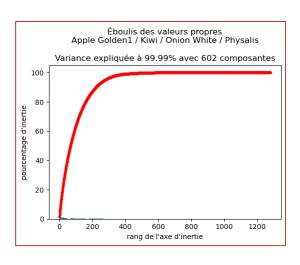


20

10

PC 1

Onion White Physalis



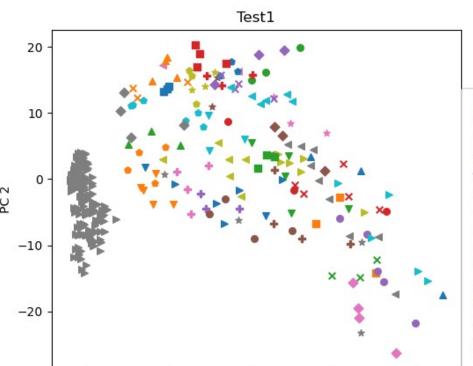




Hors propos

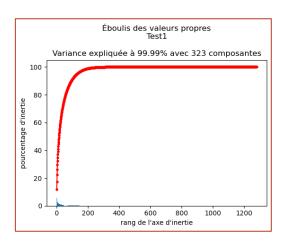
- PCA: sur 43 types de fruits
- 323 variables pour une variance expliquée à 99,99%

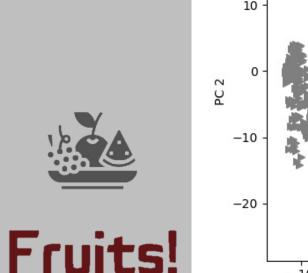
PCA - Visualisation des données sur les deux premières composantes principales



PC 1

94
131
174
204
234
323





-10

