

Week 6 (24 Sep - 30 Sep) Week 7 (1 Oct - 7 Oct)

Project: Building an ETL Pipeline to Transform Data in Data Lake (From Landing Zone to Cleaned Zone)

สามารถเริ่มต้นจากโค้ดที่ <https://github.com/zkan/swu-ds525> ได้

ในโปรเจคนี้เราจะใช้ข้อมูล GitHub event data จาก [API](#) หรือสามารถดาวน์โหลดไฟล์ JSON ได้ที่ [URL](#) นี้

สิ่งที่คาดหวังในโปรเจคนี้

ในโปรเจคนี้เราจะตั้งต้นว่าเรามีไฟล์อยู่บน S3 อยู่แล้ว ดังนั้นทุกคนสามารถที่จะอัปโหลดไฟล์ขึ้นไปใน S3 ได้เลย

1. มีโค้ดที่ใช้ Spark ทำ ETL จากข้อมูล JSON ใน S3 ที่ landing zone และไปสร้างไฟล์ที่ clean แล้วใน cleaned zone ใน data lake (S3)
 2. มีการเขียน documentation อธิบายสิ่งที่ตัวเองทำลงไว้
 3. มี instruction ในการรันโค้ดของตัวเอง
-

หนังสือภาษาไทย Spark Internal

Click <https://github.com/aorjoa/SparkInternals/tree/thai/markdown/thai> link to open resource.

Spark By Examples | Learn Spark Tutorial with Examples

Click <https://sparkbyexamples.com/> link to open resource.

Introduction to PySpark

Click <https://www.datacamp.com/courses/introduction-to-pyspark> link to open resource.

Data Wrangling with Spark

Click <https://github.com/zkan/data-wrangling-with-spark> link to open resource.

Machine Learning with Spark and Zeppelin

Click <https://github.com/zkan/machine-learning-with-spark-and-zeppelin> link to open resource.

Read Nested JSON in Spark DataFrame

Click <https://bigdataprogrammers.com/read-nested-json-in-spark-dataframe/> link to open resource.

How to parse nested JSON objects in spark sql?

Click

<https://stackoverflow.com/questions/29948789/how-to-parse-nested-json-objects-in-spark-sql>

Slides: Week 7

<https://docs.google.com/presentation/d/1gZOfJJ9TzDLQ5ZRAdTLd74ePKzeCSBFVSJsMt7LRx6M/edit?usp=sharing>

Week 6 (Part 1)

<https://youtu.be/mWkGUzsarAQ>

Week 6 (Part 2)

<https://youtu.be/sesxWvbJnzo>

Week 6 Project Walkthru

<https://youtu.be/1sVj7mwAmTE>

```
spark.sql("""  
    select  
        upc  
        , product_size  
        , split(product_size, ' ')[0] as value  
        , split(product_size, ' ')[1] as unit  
        , case  
            when split(product_size, ' ')[1] = 'OZ' then split(product_size, ' ')[0] * 100  
            else  
                split(product_size, ' ')[0] / 100  
            end as is_oz
```

Week 7

<https://youtu.be/c1RtJhyIo9s>

Week 7 Office Hour

https://youtu.be/ZIEPt0OI_w0

Week 6 (24 Sep - 30 Sep)

```
ddd_v1_w_foQ_1383611@runweb62413:~$ cat ~/.aws/credentials
[default]
aws_access_key_id = ASIARC3CKZR2WNZXQLZA
aws_secret_access_key = puIUVfojgB6Ws4jyvv/AecjBP10kqB0W8feUsn7T
aws_session_token = FwogZXIVYXdzECMaDLfAeXU6AITLcZujViLPAYzC1NjTEMc2LtjXg0mEn7nXiv6AUMyTfx5N64EU
uIojF15Qz6lMwDMeJJY4vYQpL0ATv5NaZd0o0B0c+vT2qdtL5eYsLxPjWJvw2moTv3He4SKfZagR2T/p43M+3igpVH3L6/11
T1edcIgTqh8A5brrWrg1GJ+eLvLjPpeUBQ4/vdJo+deai69aTM8gNkb5YCirrzM56kUkCXVUlcw//xklpxZEe3lc89X1kr
uly1LCp05WC1K0a10rqPbq5Gcp0fBfxyfOrmoyjbmlWBIHyi5p7uZbjItbx/VwvVL9mCXoR3+lW3nxGUaoYBM3gNeUNHKFE0X
g/a6Wkvld6Ym4Lm1MN9W
ddd_v1_w_foQ_1383611@runweb62413:~$
```

us-east-1.console.aws.amazon.com/elasticmapreduce/home?region=us-east-1#

poll_form.php Canny edge detect... Simple guide on ho... A Neural Network P... Learner Lab Canny edge detect... Machine Learning u... TensorFlow Lite Co...

aws Services Search for services, blogs, docs, and more [Alt+S] N. Virginia vocabs/user1591029=peeyapak.somvitoon@g.swu.ac.th @ 0748-512...

Amazon EMR

EMR Studio

EMR Serverless New

EMR on EC2

Clusters

Notebooks

Git repositories

Security configurations

Block public access

VPC subnets

Events

EMR on EKS

Virtual clusters

Feedback Looking for language selection? Find it in the new [Unified Settings](#)

Welcome to Amazon Elastic MapReduce

Amazon Elastic MapReduce (Amazon EMR) is a web service that enables businesses, researchers, data analysts, and developers to easily and cost-effectively process vast amounts of data.

You do not appear to have any clusters. Create one now:

Create cluster

How Elastic MapReduce Works

Upload Create Monitor

Additional Information

More about Elastic MapReduce

[EMR overview](#) [FAQs](#) [Pricing](#)

More Help Using Elastic MapReduce

[Forum](#) [Documentation](#) [Developer Guide](#) [API Reference](#) [EMR on GitHub](#) [Help portal](#)

© 2022, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Welcome to AI

Amazon Elastic MapReduce
analysts, and developers to

You do not appear to have a

Create cluster

Create Cluster - Quick Options [Go to advanced options](#)

General Configuration

Cluster name [i](#)

Logging [i](#)
S3 folder [i](#)

Launch mode Cluster [i](#) Step execution [i](#)

Software configuration

Release [i](#)

Applications Core Hadoop: Hadoop 3.2.1 with Hive 3.1.3, Hue 4.10.0, Pig 0.17.0 and Tez 0.9.2
 HBase: HBase 2.4.4 with Hadoop 3.2.1, Hive 3.1.3, Hue 4.10.0, Phoenix 5.1.2, and ZooKeeper 3.5.7

Create Cluster - Advanced Options [Go to quick options](#)

Step 1: Software and Steps

Step 2: Hardware
Step 3: General Cluster Settings
Step 4: Security

Software Configuration

Release [i](#)

| | | |
|---|--|--|
| <input checked="" type="checkbox"/> Hadoop 3.2.1 | <input type="checkbox"/> Zeppelin 0.10.0 | <input type="checkbox"/> Livy 0.7.1 |
| <input type="checkbox"/> JupyterHub 1.4.1 | <input type="checkbox"/> Tez 0.9.2 | <input type="checkbox"/> Flink 1.14.2 |
| <input type="checkbox"/> Ganglia 3.7.2 | <input type="checkbox"/> HBase 2.4.4 | <input checked="" type="checkbox"/> Pig 0.17.0 |
| <input checked="" type="checkbox"/> Hive 3.1.3 | <input type="checkbox"/> Presto 0.272 | <input type="checkbox"/> ZooKeeper 3.5.7 |
| <input type="checkbox"/> JupyterEnterpriseGateway 2.1.0 | <input type="checkbox"/> MXNet 1.8.0 | <input type="checkbox"/> Sqoop 1.4.7 |
| <input checked="" type="checkbox"/> Hue 4.10.0 | <input type="checkbox"/> Phoenix 5.1.2 | <input type="checkbox"/> Trino 378 |
| <input type="checkbox"/> Oozie 5.2.1 | <input type="checkbox"/> Spark 3.2.1 | <input type="checkbox"/> HCatalog 3.1.3 |
| <input type="checkbox"/> TensorFlow 2.4.1 | | |

Multiple master nodes (optional)

Software Configuration

Release [i](#)

| | | |
|--|---|--|
| <input checked="" type="checkbox"/> Hadoop 3.2.1 | <input type="checkbox"/> Zeppelin 0.10.0 | <input type="checkbox"/> Livy 0.7.1 |
| <input type="checkbox"/> JupyterHub 1.4.1 | <input type="checkbox"/> Tez 0.9.2 | <input type="checkbox"/> Flink 1.14.2 |
| <input type="checkbox"/> Ganglia 3.7.2 | <input type="checkbox"/> HBase 2.4.4 | <input checked="" type="checkbox"/> Pig 0.17.0 |
| <input checked="" type="checkbox"/> Hive 3.1.3 | <input type="checkbox"/> Presto 0.272 | <input type="checkbox"/> ZooKeeper 3.5.7 |
| <input checked="" type="checkbox"/> JupyterEnterpriseGateway 2.1.0 | <input type="checkbox"/> MXNet 1.8.0 | <input type="checkbox"/> Sqoop 1.4.7 |
| <input checked="" type="checkbox"/> Hue 4.10.0 | <input type="checkbox"/> Phoenix 5.1.2 | <input type="checkbox"/> Trino 378 |
| <input type="checkbox"/> Oozie 5.2.1 | <input checked="" type="checkbox"/> Spark 3.2.1 | <input type="checkbox"/> HCatalog 3.1.3 |
| <input type="checkbox"/> TensorFlow 2.4.1 | | |

Amazon EMR

EMR Serverless is now GA.
With EMR Serverless, get the benefits of Amazon EMR such as open source compatibility, latest versions and performance optimized runtime for popular frameworks along with easy provisioning, quick job startup, automatic capacity management, and simple cost controls. [Get Started with EMR Serverless.](#)

Configuration details

- Release label: emr-6.7.0
- Hadoop distribution: Amazon 3.2.1
- Applications: Hive 3.1.3, Pig 0.17.0, Hue 4.10.0, JupyterEnterpriseGateway 2.1.0, Spark 3.2.1
- Log URI: s3://aws-logs-074831285365-us-east-1/elastictmapreduce/
- EMRFS consistent view: Disabled
- Custom AMI ID: --
- Amazon Linux Release: 2.0.20220806.1 [Learn more](#)

Application user interfaces

Persistent user interfaces [x]: --

On-cluster user Not Enabled [Enable an SSH Connection](#)
interfaces [x]: --

Network and hardware

- Availability zone: us-east-1b
- Subnet ID: [subnet-04ed0369f48942ed](#)
- Master: Provisioning 1 m5.xlarge
- Core: Provisioning 2 m5.xlarge
- Task: --
- Cluster scaling: Not enabled
- Auto-termination: Terminate if idle for 1 hour

Security and access

- Key name: --
- EC2 instance profile: EMR_EC2_DefaultRole
- EMR role: EMR_DefaultRole
- Auto Scaling role: EMR_AutoScaling_DefaultRole
- Visible to all users: All [Change](#)
- Security groups for Master: [sg-0f6468cd497c78455](#) (ElasticMapReduce-master)
- Security groups for Core & Task: [sg-06317eeab7a22ded0](#) (ElasticMapReduce-Task: slave)

Instances (3) Info

| Name | Instance ID | Instance state | Instance type | Status check | Alarm status | Availability Zone |
|------|---------------------|----------------|---------------|-------------------|--------------|-------------------|
| - | i-07e31bce190575ae0 | Running | m5.xlarge | 2/2 checks passed | No alarms | us-east-1b |
| - | i-0999c0a6ae2c66806 | Running | m5.xlarge | 2/2 checks passed | No alarms | us-east-1b |
| - | i-09bd0a75606ac4842 | Running | m5.xlarge | 2/2 checks passed | No alarms | us-east-1b |

Cluster: My cluster in class Running

Summary

- ID: j-30GA8FX0O1ZL1
- Creation date: 2022-09-24 17:09 (UTC+7)
- Elapsed time: 10 minutes
- After last step completes: Cluster waits
- Termination protection: On [Change](#)
- Tags: -- [View All / Edit](#)
- Master public DNS: [ec2-3-83-128-235.compute-1.amazonaws.com](#) [Connect to the Master Node Using SSH](#)

Configuration details

- Release label: emr-6.7.0
- Hadoop distribution: Amazon 3.2.1
- Applications: Hive 3.1.3, Pig 0.17.0, Hue 4.10.0,

Application user interfaces

Persistent user interfaces [x]: Spark history server, YARN timeline server, Tez UI

On-cluster user Not Enabled [Enable an SSH Connection](#)
interfaces [x]: --

Amazon EMR

EMR Studio

EMR Serverless [New](#)

EMR on EC2

Clusters

Notebooks

Git repositories

Security configurations

Block public access

VPC subnets

Events

EMR on EKS

Virtual clusters

Help

EMR Serverless is now GA.
With EMR Serverless, get the benefits of Amazon EMR such as open source compatibility, latest versions and performance optimized runtime for popular frameworks along with easy provisioning, quick job startup, automatic capacity management, and simple cost controls. [Get Started with EMR Serverless.](#)

Notebooks

Use EMR notebooks based on Jupyter to analyze data interactively with live code, narrative text, visualizations, and more. Create and attach notebooks to Amazon EMR cluster Hadoop, Spark, and Livy. Notebooks run free of charge and are saved in Amazon S3 independently of clusters. Standard billing for clusters and Amazon S3 apply. [Learn more](#)

[Create notebook](#) [View details](#) [Open in JupyterLab](#) [Open In Jupyter](#) [Start](#) [Stop](#) [Delete](#)

Filter: All notebooks [Filter notebooks ...](#) 0 notebooks (all loaded)

| Name | Status | Cluster | Creation time (UTC+7) | Last modified |
|------|--------|---------|-----------------------|---------------|
|------|--------|---------|-----------------------|---------------|

Amazon EMR

EMR Studio

EMR Serverless [New](#)

EMR on EC2

Clusters

Notebooks

Git repositories

Security configurations

Block public access

VPC subnets

Events

EMR on EKS

Virtual clusters

Help

What's new

EMR Serverless is now GA.
With EMR Serverless, get the benefits of Amazon EMR such as open source compatibility, latest versions and performance optimized runtime for popular frameworks along with easy provisioning, quick job startup, automatic capacity management, and simple cost controls. [Get Started with EMR Serverless.](#)

Name your notebook, choose a cluster or create one, and customize configuration options if desired. [Learn more](#)

Notebook name* mynotebook
Names may only contain alphanumeric characters, hyphens (-), or underscores (_).

Description
256 characters max.

Cluster* Choose an existing cluster [Choose](#) My cluster in class **j-30GA8FX0O1ZL1**
 Create a cluster

Security groups Use default security groups
 Choose security groups (vpc-0644f8397501301da)

AWS service role* LabRole
 Make sure this role has the required permissions. [Learn more](#)

aws Services [Search for services, features, blogs, docs, and more](#) [Alt+S] N. Virginia v vocabs/user1591029=peeyapak.somvittoon@g.swu.ac.th @ 0748-312...

Amazon EMR

EMR Studio

EMR Serverless [New](#)

EMR on EC2

Clusters

Notebooks

Git repositories

Security configurations

Block public access

VPC subnets

Events

EMR on EKS

Virtual clusters

Help

What's new

EMR Serverless is now GA.
With EMR Serverless, get the benefits of Amazon EMR such as open source compatibility, latest versions and performance optimized runtime for popular frameworks along with easy provisioning, quick job startup, automatic capacity management, and simple cost controls. [Get Started with EMR Serverless.](#)

Notebook: mynotebook Starting Starting workspace(notebook). Cluster j-30GA8FX0O1ZL1.

[Open in JupyterLab](#) [Open in Jupyter](#) [Stop](#) [Delete](#)

Notebook

Notebook ID: e-EH0DJCMFFK7XD7I09FU1KF5V

Description: --

Last modified: 7 seconds ago

Last modified by: ...assumed-role/voclabs/user1591029=peeyapak.somvittoon@g.swu.ac.th

Created on: 2022-09-24 17:24 (UTC+7)

Created by: ...assumed-role/voclabs/user1591029=peeyapak.somvittoon@g.swu.ac.th

Service IAM role: LabRole

Notebook tags: creatorUserId = AROARC3CKZR23PB73NMSZ:user1591029=peeyapak.somvittoon@g.swu.ac.th [View All / Edit](#)

Notebook location: s3://aws-emr-resources-074831285365-us-east-1/notebooks/

Cluster

us-east-1.console.aws.amazon.com/elasticmapreduce/home?region=us-east-1#notebook-details:e-EH9DJKCMFFK7XD7IO9FU1KF5V

poll_form.php Canny edge detect... Simple guide on ho... A Neural Network P... Learner Lab Canny edge detect... Machine Learning u... TensorFlow Lite Co...

aws Services Search for services, features, blogs, docs, and more [Alt+S] N. Virginia vocabs/user1591029=peeyapak.somvitoon@g.swu.ac.th @ 0748-312...

EMR Serverless is now GA.
With EMR Serverless, get the benefits of Amazon EMR such as open source compatibility, latest versions and performance optimized runtime for popular frameworks along with easy provisioning, quick job startup, automatic capacity management, and simple cost controls. [Get Started with EMR Serverless.](#)

Amazon EMR
EMR Studio
EMR Serverless **New**
EMR on EC2
Clusters
Notebooks
Git repositories
Security configurations
Block public access
VPC subnets
Events
EMR on EKS
Virtual clusters

Help
What's new

Notebook: mynotebook **Ready** Workspace(notebook) is ready to run jobs on cluster j-30GA8FX0O1ZL1.

Open in JupyterLab Open in Jupyter Stop Delete

Notebook

Notebook ID: e-EH9DJKCMFFK7XD7IO9FU1KF5V
Description: --
Last modified: 8 seconds ago
Last modified by: ...assumed-role:vocabs/user1591029=peeyapak.somvitoon@g.swu.ac.th
Created on: 2022-09-24 17:24 (UTC+7)
Created by: ...assumed-role:vocabs/user1591029=peeyapak.somvitoon@g.swu.ac.th
Service IAM role: [LabRole](#)
Security groups for sg-055a0dd20ad1afa6 master instance:
Security groups for sg-07d11a332b983be3 notebook instance:
Notebook tags: creatorUserId = AROARC3CK2R23PB73NMSZ:user1591029=peeyapak.somvitoon@q.swu.ac.th [View All / Edit](#)

e-EH9DJKCMFFK7XD7IO9FU1KF5V.emrnotebooks-prod.us-east-1.amazonaws.com/e-EH9DJKCMFFK7XD7IO9FU1KF5V/lab

poll_form.php Canny edge detect... Simple guide on ho... A Neural Network P... Learner Lab Canny edge detect... Machine Learning u... TensorFlow Lite Co...

File Edit View Run Kernel Git Tabs Settings Help

Launcher

Filter files by name

Name Last Modified

mynotebook.ipynb a minute ago

Notebook

Python 3 PySpark Spark SparkR

Console

Python 3 PySpark Spark SparkR

Other

Launcher

jupyter

Files Running Clusters

Select items to perform actions on them.

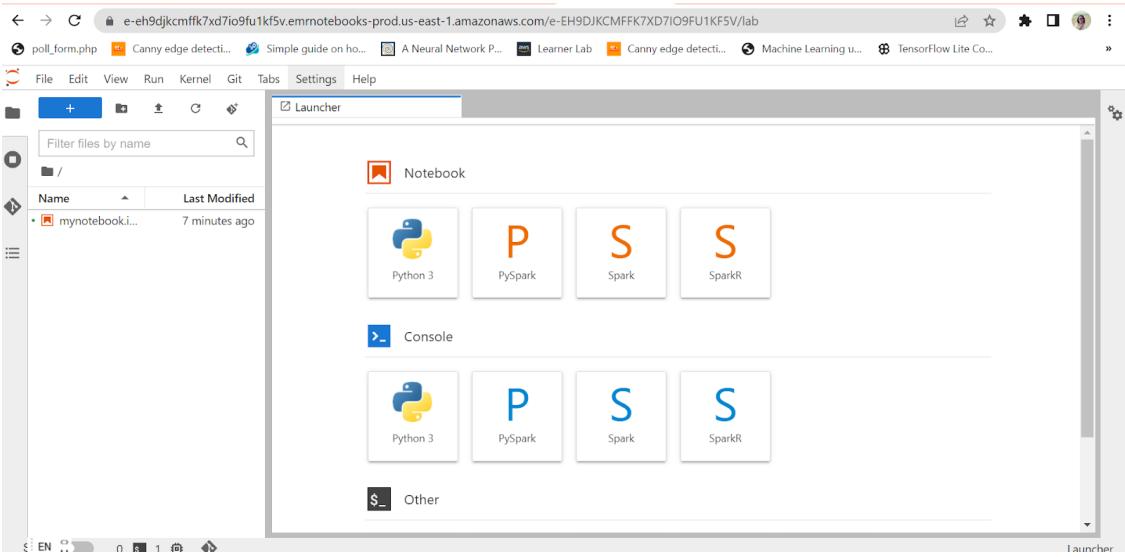
Upload New File size

0 mynotebook.ipynb 2 minutes ago 72 B

A screenshot of a web browser window showing several tabs open. The tabs include: poll.form.php, Canny edge detect..., Simple guide on ho..., A Neural Network P..., Learner Lab, Canny edge detect..., Machine Learning u..., TensorFlow Lite Co..., and poll.form.php. The main content area shows the Jupyter Notebook interface with a 'Running' tab selected. It displays a list of currently running Jupyter processes, terminals, and notebooks. One notebook named 'mynotebook.ipynb' is listed, created by 'Python 3' and shutdown 'seconds ago'.

A screenshot of a Jupyter Notebook interface. The top navigation bar shows 'Files', 'Running', and 'Clusters'. Below it, a message says 'Currently running Jupyter processes'. Under 'Terminals', it says 'There are no terminals running.' Under 'Notebooks', there is one entry: 'mynotebook.ipynb' created by 'Python 3' and shutdown 'seconds ago'.

A screenshot of the AWS IAM (Identity and Access Management) console. The left sidebar shows 'Identity and Access Management (IAM)'. Under 'Access management', 'Roles' is selected, showing 'Policies', 'Identity providers', and 'Account settings'. Under 'Access reports', 'Access analyzer' is selected, showing 'Archive rules', 'Analyzers', 'Settings', and 'Credential report'. The main content area shows the 'LabRole' configuration. It includes a summary table with details like Creation date (September 03, 2022), Last activity (24 minutes ago), ARN (arn:aws:iam::074831285365:role/LabRole), Maximum session duration (1 hour), and Instance profile ARN (arn:aws:iam::074831285365:instance-profile/LabInstanceProfile). Below the summary are tabs for 'Permissions', 'Trust relationships', 'Tags (1)', 'Access Advisor', and 'Revoke sessions'. At the bottom, there is a section for 'Permissions policies (7)' with buttons for 'Edit', 'Delete', 'Simulate', 'Remove', and 'Add permissions'.



jupyter Untitled Last Checkpoint: 3 minutes ago (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help Trusted | PySpark O

In [1]: `from pyspark.sql import SparkSession`

Starting Spark application

| ID | YARN Application ID | Kind | State | Spark UI | Driver log | User | Current session? |
|----|--------------------------------|---------|-------|----------------------|----------------------|------|------------------|
| 0 | application_1664014599155_0001 | pyspark | idle | Link | Link | None | ✓ |

SparkSession available as 'spark'.

In [2]: `spark = SparkSession.builder \
.appName("ETL") \
.getOrCreate()`

In []:

s3.console.aws.amazon.com/s3/upload/kiktitanic?region=us-east-1

Upload succeeded

View details below.

Upload: status

The information below will no longer be available after you navigate away from this page.

Summary

| Destination | Succeeded | Failed |
|-----------------|---------------------------|-------------------|
| s3://kiktitanic | 1 file, 59.8 KB (100.00%) | 0 files, 0 B (0%) |

Files and folders (1 Total, 59.8 KB)

| Name | Folder | Type | Size | Status | Error |
|-------------|--------|----------|---------|-----------|-------|
| titanic.csv | - | text/csv | 59.8 KB | Succeeded | - |

jupyter Untitled Last Checkpoint: 19 minutes ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted PySpark

In [1]: `from pyspark.sql import SparkSession`

Starting Spark application

| ID | YARN Application ID | Kind | State | Spark UI | Driver log | User | Current session? |
|----|--------------------------------|---------|-------|----------------------|----------------------|------|------------------|
| 0 | application_1664014599155_0001 | pyspark | idle | Link | Link | None | ✓ |

SparkSession available as 'spark'.

In [2]: `spark = SparkSession.builder \
 .appName("ETL") \
 .getOrCreate()`

In [3]: `bucket = "s3://kiktitanic"`

In [5]: `df = spark.read.csv(bucket)`

jupyter Untitled Last Checkpoint: 20 minutes ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted | PySpark C

In [3]: `bucket = "s3://kiktitanic"`

In [5]: `df = spark.read.csv(bucket)`

▼ Spark Job Progress

▼ Job [0]: csv at NativeMethodAccesso... Job Progress: 1/1 Tasks...

Stage [ID]: name at [sou... Status Task Progress Elapsed Time (s...) Failed Tas...

Stage [0]: csv at NativeMet... COMP... 1/1 5.601

In [7]: `df = spark.read.option('header', 'true').csv(bucket)`

In [7]: `df = spark.read.option('header', 'true').csv(bucket)`

► Spark Job Progress

In [8]: `df.show()`

► Spark Job Progress

| PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin |
|-------------|----------|--------|-----------------------|--------|-----|-------|-------|-----------|---------|-------|
| 1 | 0 | 3 | Braund, Mr. Owen G. | male | 22 | 1 | 0 | A/5 21171 | 7.25 | n/a |
| 2 | 1 | 1 | Cumings, Mrs. John S. | female | 38 | 1 | 0 | PC 17599 | 71.2833 | |

jupyter Untitled (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help Trusted | PySpark



only showing top 20 rows

In [9]: `df.select('Age', 'Survived')`

DataFrame[Age: string, Survived: string]

In [10]: `df.select('Age', 'Survived').show()`

▶ Spark Job Progress

| Age | Survived |
|------|----------|
| 22 | 0 |
| 38 | 1 |
| 26 | 1 |
| 35 | 1 |
| 35 | 0 |
| null | 0 |
| 54 | 0 |
| 2 | 0 |
| 27 | 1 |

```
In [12]: df.createOrReplaceTempView("titanic")
```

```
In [14]: spark.sql("""
    select
        Age
        , Survived
    from
        titanic
""").show()
```

▶ Spark Job Progress

| Age | Survived |
|------|----------|
| 22 | 0 |
| 38 | 1 |
| 26 | 1 |
| 35 | 1 |
| 35 | 0 |
| null | 0 |
| 54 | 0 |

jupyter Untitled Last Checkpoint: 32 minutes ago (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help Trusted | PySpark

```
+---+-----+
| 2| 0|
| null| 1|
| 31| 0|
| null| 1|
+---+-----+
only showing top 20 rows
```

In [15]: `result = spark.sql("""
 select
 Age
 , Survived
 from
 titanic
""")`

In [16]: `result.show(3)`

▶ Spark Job Progress

```
+---+-----+
|Age|Survived|
+---+-----+
| 2| 0|
| 38| 1|
| 26| 1|
+---+-----+
only showing top 3 rows
```

In [17]: `result.write.mode("overwrite").csv("s3://kiktitanic/cleaned")`

▶ Spark Job Progress

aws Services Search for services, features, blogs, docs, and more [Alt+S] Global vclabs/user1591029=peeyapak.somvittoon@g.swu.ac.th @ 0748-312...

Amazon S3

Buckets

- Access Points
- Object Lambda Access Points
- Multi-Region Access Points
- Batch Operations
- Access analyzer for S3

Block Public Access settings for this account

Storage Lens

- Dashboards
- AWS Organizations settings

Feature spotlight

AWS Marketplace for S3

Amazon S3 > Buckets > kiktitanic > cleaned/

cleaned/

Copy S3 URI

Objects Properties

Objects (2)

Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory [Get 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](#)

| <input type="checkbox"/> | Name | Type | Last modified | Size | Storage class |
|--------------------------|--|------|--|--------|---------------|
| <input type="checkbox"/> | _SUCCESS | - | September 24, 2022, 18:13:28 (UTC+07:00) | 0 B | Standard |
| <input type="checkbox"/> | part-00000-7815f5c1-f264-4755-9ea6-dbef450a12cf-c000.csv | csv | September 24, 2022, 18:13:28 (UTC+07:00) | 4.3 KB | Standard |

jupyter Untitled Last Checkpoint: 37 minutes ago (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help

In [16]: `titanic`

In [16]: `result.show(3)`

▶ Spark Job Progress

```
+---+-----+
|Age|Survived|
+---+-----+
| 22|     0|
| 38|     1|
| 26|     1|
+---+-----+
only showing top 3 rows
```

In [17]: `result.write.mode("overwrite").csv("s3://kiktitanic/cleaned")`

▶ Spark Job Progress

In [18]: `result.write.mode("overwrite").parquet("s3://kiktitanic/cleaned-parquet")`

▶ Spark Job Progress

In []:

s3.console.aws.amazon.com/s3/buckets/kiktitanic?region=us-east-1&tab=objects

poll_form.php Canny edge detect... Simple guide on ho... A Neural Network P... Learner Lab Canny edge detect... Machine Learning u... TensorFlow Lite Co...

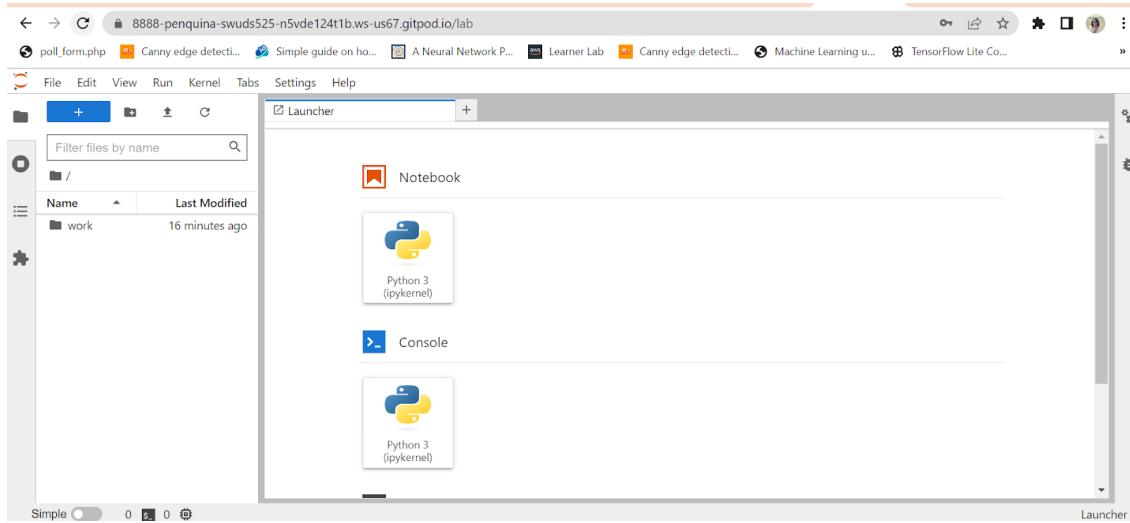
AWS Services Search for services, features, blogs, docs, and more [Alt+S]

Amazon S3 Amazon S3 > Buckets > kiktitanic

kiktitanic Info

Objects (3) 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

| Name | Type | Last modified | Size | Storage class |
|------------------|--------|--|---------|---------------|
| cleaned-parquet/ | Folder | - | - | - |
| cleaned/ | Folder | - | - | - |
| titanic.csv | csv | September 24, 2022, 17:49:40 (UTC+07:00) | 59.8 KB | Standard |



```
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

oken=60fcfc77228bcebae61f235de501111b1685e774ccc1d5134
workshop-pyspark-notebook-1 | [I 2022-10-01 10:35:29.260 ServerApp] Use Control-C to stop this serv
er and shut down all kernels (twice to skip confirmation).
workshop-pyspark-notebook-1 | [C 2022-10-01 10:35:29.264 ServerApp]
workshop-pyspark-notebook-1 |
workshop-pyspark-notebook-1 | To access the server, open this file in a browser:
workshop-pyspark-notebook-1 |     file:///home/jovyan/.local/share/jupyter/runtime/jpserver-25-
open.html
workshop-pyspark-notebook-1 |     Or copy and paste one of these URLs:
workshop-pyspark-notebook-1 |     http://c4749dc8605d:8888/lab?token=60fcfc77228bcebae61f235de50
1111b1685e774ccc1d5134
workshop-pyspark-notebook-1 |     or http://127.0.0.1:8888/lab?token=60fcfc77228bcebae61f235de50111
1b1685e774ccc1d5134
workshop-pyspark-notebook-1 | [I 2022-10-01 10:41:32.420 ServerApp] 302 GET / (192.168.48.138) 0.61
ms
workshop-pyspark-notebook-1 | [I 2022-10-01 10:41:32.624 LabApp] 302 GET /lab? (192.168.48.138) 0.7
2ms
workshop-pyspark-notebook-1 | [I 2022-10-01 10:42:05.615 ServerApp] 302 POST /login?next=%2Flab%3F
(192.168.48.138) 1.19ms
workshop-pyspark-notebook-1 | [I 2022-10-01 10:42:11.971 LabApp] Build is up to date
```

| Port | Address | Actions |
|------|---|---------|
| 3000 | https://3000-penquinawuds525-n5vde124t1b.ws-us67.gitpod.io | |
| 4040 | https://4040-penquinawuds525-n5vde124t1b.ws-us67.gitpod.io | |
| 4041 | https://4041-penquinawuds525-n5vde124t1b.ws-us67.gitpod.io | |
| 8888 | https://8888-penquinawuds525-n5vde124t1b.ws-us67.gitpod.io | |

8888-penquina-swuds525-n5vde124t1b.ws-us67.gitpod.io/lab/tree/work/workshop.ipynb

File Edit View Run Kernel Tabs Settings Help

Launcher workshop.ipynb Python 3 (ipykernel)

Filter files by name

Name Last Modified

- dataset 18 minutes ago
- docker-compose.yml 19 minutes ago
- workshop.ipynb 19 minutes ago

base price and actual selling price; to determine a product's discount, if any

- Promotional support details (e.g., sale tag, in-store display), if applicable for the given product/store/week
- Store information, including size and location, as well as a price tier designation (e.g., upscale vs. value)
- Product information, including UPC, size, and description

To identify outliers, it is suggested to look at

- The ratio of units vs. number of visits
- The ratio of visits vs. number of households
- Some items that may be out-of-stock or discontinued for a store

Source: <https://www.dunnhumby.com/source-files/>

```
[1]: from pyspark.sql import Row, SparkSession
```

```
[2]: spark = SparkSession.builder \
    .appName("breakfast") \
    .getOrCreate()
```

```
[3]: product_data_folder = "dataset/products"
store_data_folder = "dataset/stores"
transaction_data_folder = "dataset/transactions"
```

Launcher workshop.ipynb Python 3 (ipykernel)

1. What is the range of prices offered on products?

2. What is the impact on units/visit of promotions by geographies?

3. Which products would you lower the price to increase sales?

```
[4]: product_df = spark\
    .read\
    .option("header",True) \
    .csv(product_data_folder)
```

```
[5]: product_df.show(1)
```

| UPC | DESCRIPTION | MANUFACTURER | CATEGORY | SUB_CATEGORY | PRODUCT_SIZE |
|------------|----------------------|---------------|------------|--------------|--------------|
| 1111009477 | PL MINI TWIST PRE... | PRIVATE LABEL | BAG SNACKS | PRETZELS | 15 OZ |

only showing top 1 row

```
[ ]:
```

Launcher workshop.ipynb Code Python 3 (ipykernel)

only showing top 20 rows

```
[12]: store_df.show(1)
```

| | STORE_ID | STORE_NAME | ADDRESS_CITY_NAME | ADDRESS_STATE_PROV_CODE | MSA_CODE | SEG_VALUE_NAME | PARKING_SPACE_QTY | SALES_AREA_SIZE_NUM | Avg_WEEKLY_BASKETS |
|-----|----------|------------|-------------------|-------------------------|----------|----------------|-------------------|---------------------|--------------------|
| 408 | 389 | SILVERLAKE | ERLANGER | KY | 17140 | MAINSTREAM | | 46073 | 24767 |

only showing top 1 row

```
product_df.createOrReplaceTempView("products")
transaction_df.createOrReplaceTempView("transactions")
```

```
#Range of price
spark.sql("""
    select
        products.upc
        , min(price)
        , max(price)
        , description
        , category

    from transactions
    join products
    on
        transactions.upc = products.upc
    group by
        1, 4, 5
""").show(3)
```

| upc | min(price) | max(price) | description | category |
|------------|------------|------------|----------------------|------------|
| 1111009477 | 0.89 | 1.83 | PL MINI TWIST PRE... | BAG SNACKS |
| 1111009497 | 0.86 | 1.69 | PL PRETZEL STICKS | BAG SNACKS |
| 1111009507 | 0.8 | 1.69 | PL TWIST PRETZELS | BAG SNACKS |

only showing top 3 rows