

Práctica de Azure

OPTIMIZACIÓN DE GRANDES VOLÚMENES DE DATOS
2023/2024

Carlos Oliva López y Christian Graf Aray
MÁSTER DE APRENDIZAJE AUTOMÁTICO Y DATOS MASIVOS

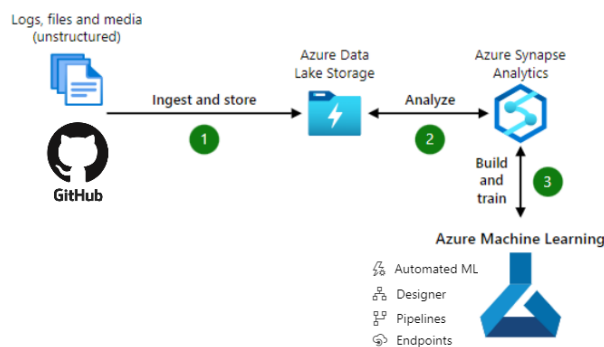
Contenido

1. Introducción	1
2. Azure Synapse Analytics.....	1
2.1. Creación del pool de Spark	1
2.2. Ingesta de datos	2
2.3. Lectura y procesamiento de datos.....	4
2.5. Visualización de los datos	5
3. Azure ML Studio	7
3.1. Creación del clúster.....	7
3.2. Importación del dataset	8
3.3. AutoML: encontrar el mejor modelo de clasificación	10
3.4. Pipelines de Designer: un entrenamiento más personalizado	14
3.4.1. Pipeline de entrenamiento	14
3.4.2. Pipeline de inferencia y despliegue.....	16
4. Costes	18
5. Referencias.....	21

1. Introducción

En esta práctica analizaremos las posibilidades y ventajas que nos aportan los servicios de Azure como Synapse Analytics y ML Studio. Para probar nuestro punto, utilizaremos estas dos herramientas con un [dataset](#) que recoge más de 10.000 enlaces a diferentes endpoints y sitios web, clasificándolos como phishing o legítimos. En total, por cada uno de esos enlaces se recogen 87 features, que van desde el propio enlace (la URL) hasta el número de ciertos tipos de caracteres especiales. Al ser un dataset de clasificación, la etiqueta a predecir será la columna “status”.

Por lo tanto, en la primera parte de la práctica utilizaremos Azure Synapse Analytics para realizar un análisis exploratorio y tratamiento de los datos anteriormente descritos. Una vez guardado esos datos curados, los utilizaremos en Azure ML Studio para realizar su clasificación, examinando las distintas características que nos ofrece Azure. El diagrama de la arquitectura final quedará así:



El repositorio del grupo se puede encontrar [aquí](#), y se han dado permisos a la [suscripción](#) y los workspaces creados para que se puedan consultar los enlaces que van apareciendo a lo largo de la memoria.

2. Azure Synapse Analytics

En primer lugar, crearemos el [workspace](#) de Synapse Analytics que nos permitirá realizar el EDA y tratamiento de los datos.

2.1. Creación del pool de Spark

Para poder trabajar en el workspace, empezaremos creando [un pool de Spark](#) que nos permitirá lanzar notebooks utilizando PySpark. Debido a que solo se necesita tratar un dataset pequeño, escogemos pocos vCores como máximo y de la familia “Memory Optimized” para ahorrar costes.

New Apache Spark pool

Basics * Additional settings * Tags Review + create

Create an Synapse Analytics Apache Spark pool with your preferred configurations. Complete the Basics tab then go to Review + Create to provision with smart defaults, or visit each tab to customize.

Apache Spark pool details

Name your Apache Spark pool and choose its initial settings.

Apache Spark pool name *	<input type="text" value="sparkpool"/>
Node size family *	<input type="text" value="Memory Optimized"/>
Node size *	<input type="text" value="Small (4 vCores / 32 GB)"/>
Autoscale * ⓘ	<input checked="" type="radio"/> Enabled <input type="radio"/> Disabled
Number of nodes *	<input type="text" value="3"/> <input type="range" value="3"/> <input type="text" value="6"/>
Estimated price ⓘ	Est. cost per hour 1.84 to 3.68 USD View pricing details
Dynamically allocate executors * ⓘ	<input type="radio"/> Enabled <input checked="" type="radio"/> Disabled

2.2. Ingesta de datos

Posteriormente, importaremos los datos con los que vamos a trabajar. Estos se han descargado de Kaggle y se han subido al repositorio de GitHub, situándose [aquí](#). Los pasos para la ingesta se muestran a continuación. Para que funcione correctamente la ingesta mediante HTTP, es importante utilizar el [enlace tipo raw](#) de GitHub.

Copy Data tool

Properties Source Dataset Configuration Destination Settings Review and finish

Source data store

Specify the source data store for the copy task. You can use an existing data store connection or specify a new data store.

Source type

Connection * [Edit](#) [+ New connection](#)

Integration runtime * [Edit](#)

Base URL
<https://raw.githubusercontent.com/ETSISI-OGVD/practicaogvd23-2>

Relative URL ⓘ

Request method ⓘ

Additional headers ⓘ

Binary copy ⓘ ☐

Request timeout ⓘ

Max concurrent connections ⓘ

Copy Data tool

Properties Source Dataset Configuration Destination Settings Review and finish

File format settings

File format [Detect text format](#) [Preview data](#)

Column delimiter [Edit](#)

Row delimiter [Edit](#)

☒ First row as header ⓘ

Advanced

Skip line count ⓘ

Quote character ⓘ
 [Edit](#)

Escape character ⓘ
 [Edit](#)

Null value ⓘ

Encoding ⓘ

Preview data

Linked service: PhisingDatasetConnection

Object: https://raw.githubusercontent.com/ETSIISI-OGVD/practicaogvd23-24-grupo_christiangraf_carlosoliva/main/dataset_1...

PreviewSchema

Prop_0	url
0	http://www.progarchives.com/album.asp?id=61737
1	http://signin.eday.co.uk/ws.edayjsapi.dll/sign.inusingsslpuseridcopartnerid2siteid.zdfsx949xyss1prnbh0soabfdzgdh2kpx
2	http://www.avevaconstruction.com/blesstool/image.htm
3	http://www.jp519.com/
4	https://www.velocidrone.com/
5	https://support-appleid.com.secureupdate.duilaweryork.com/ap/b5aed586dda5d21/?cmd=_update&dispatch=b5aed586dda5d219f&locale=_US
6	https://www.authpro.com/auth/ubabankng/?action=reg
7	http://litlee.com.au/alibaba/login.alibaba.com.php
8	http://www.tutorialspoint.com/dbms/

Copy Data tool

Properties

Source

Destination

Dataset

Configuration

Settings

Review and finish

Destination data store

Specify the destination data store for the copy task. You can use an existing data store connection or specify a new data store.

Destination type

Azure Data Lake Storage Gen2

Connection *

synapse-practica-ws-WorkspaceDefault

Edit

New connection

Integration runtime *

AutoResolveIntegrationRuntime

Edit

Folder path

If the identity you use to access the data store only has permission to subdirectory instead of the entire account, specify the path to browse.

practica-fs/

Browse

File name

phising.csv

Copy behavior ⓘ

Select...

Max concurrent connections ⓘ

Block size (MB) ⓘ

Metadata ⓘ

+ New

Copy Data tool

Properties

Source

Destination

Dataset

Configuration

Settings

Review and finish

File format settings

File format

DelimitedText

Column delimiter

Comma (,)

Edit

Row delimiter

Default (\r\n or \n)

Edit

Add header to file ⓘ

☒

Advanced

Quote character ⓘ

No quote character

Edit

Escape character ⓘ

Backslash (\)

Edit

Copy Data tool

Properties

Source

Destination

Settings

Review and finish

Settings

Enter name and description for the copy data task, more options for data movement

Task name *

CopyPipelineValid

Task description

Fault tolerance ⓘ

Enable logging ⓘ

☒

Logging settings

Storage connection name * ⓘ

synapse-practica-ws-WorkspaceDefault

Test connection

Edit

New

Integration runtime *

AutoResolveIntegrationRuntime

Edit

Logging level ⓘ

Info

Logging mode ⓘ

Reliable

Best effort

Folder path ⓘ

practica-fs

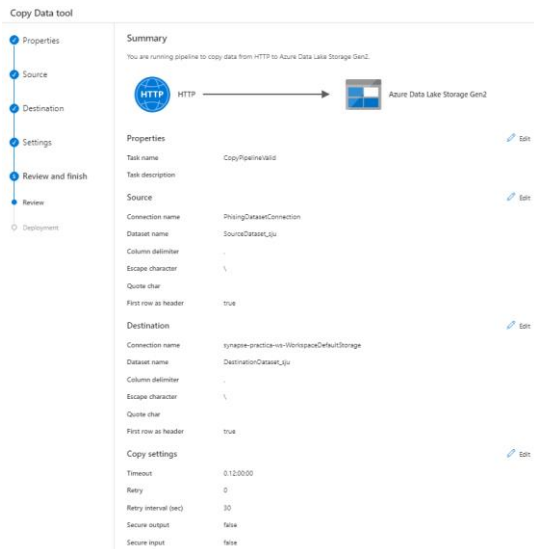
Browse

Enable staging ⓘ

☐

Advanced

3



2.3. Lectura y procesamiento de datos

Una vez creado el [notebook](#) y asignado el pool al que creamos en el apartado 2.1, podemos comenzar con el procesamiento. Se ha utilizado la función de `spark.read` para leer los datos previamente guardados:

```
1 path = "abfss://practica-fs@roberwido.dfs.core.windows.net/phishing.csv"
2 phishing_df = spark.read.option("header", "true").option("inferSchema", "true").csv(path)
```

✓ - Command executed in 31 sec 52 ms on 4:16:28 PM, 3/19/24

Utilizando el atributo `dtypes` podemos verificar que la mayoría de las columnas de los datos son ya numéricas

```
1 # vemos los tipos de datos cargados
2 phishing_df.dtypes
```

✓ - Command executed in 183 ms on 4:03:22 AM, 3/19/24

```
('domain_with_copyright', 'string'),
('whois_registered_domain', 'int'),
('domain_registration_length', 'int'),
('domain_age', 'int'),
('web_traffic', 'int'),
('dns_record', 'int'),
('google_index', 'int'),
('page_rank', 'int'),
('status', 'string')]
```

Las columnas “domain_with_copyright” y “status” son de tipo string, por lo que vamos a convertirlas en números enteros.

```
1 # Vamos a comprobar los valores únicos para la columna 'domain_with_copyright'
2 unique_values = phishing_df.select('domain_with_copyright').distinct().collect()
3 unique_values
```

```
[Row(domain_with_copyright='One'),
 Row(domain_with_copyright='0'),
 Row(domain_with_copyright='one'),
 Row(domain_with_copyright='zero'),
 Row(domain_with_copyright='Zero'),
 Row(domain_with_copyright='1')]
```

Se puede verificar que los valores de la columna son ahora 1 o 0. Luego hacemos lo mismo con la columna “status”:

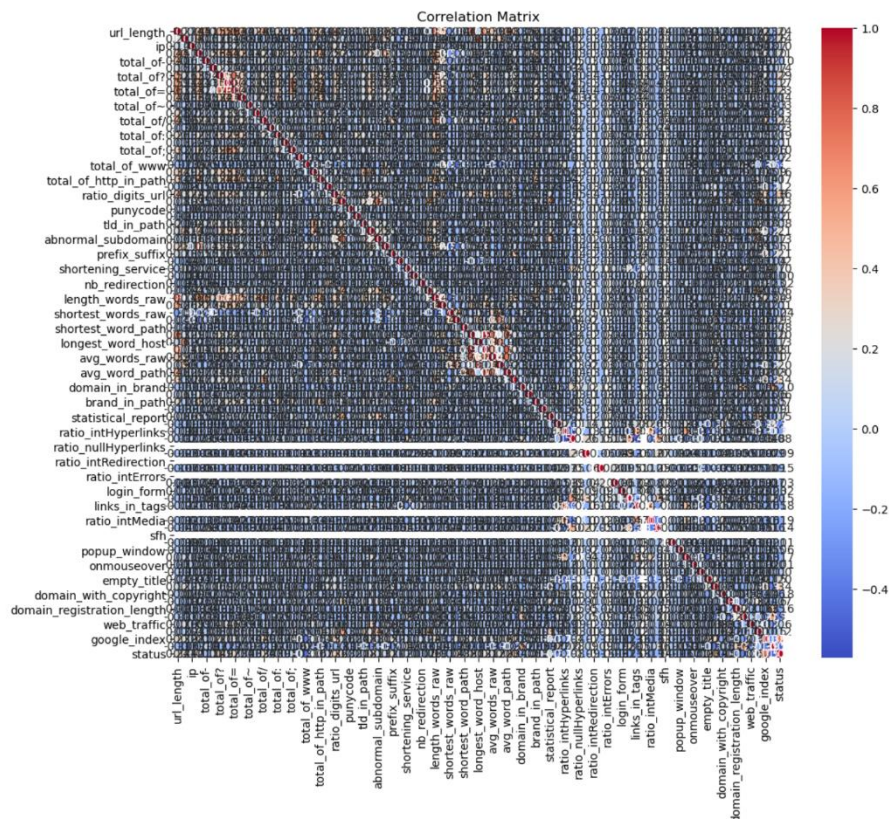
2.5. Visualización de los datos

Se generará una matriz de correlación para buscar columnas con un alto potencial de ser predictivas.

5

```
1 plt.figure(figsize=(12, 10))
2 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
3 plt.title('Correlation Matrix')
4 plt.show()
```

[18] ✓ - Command executed in 8 sec 715 ms on 4/03/49 AM, 3/19/24

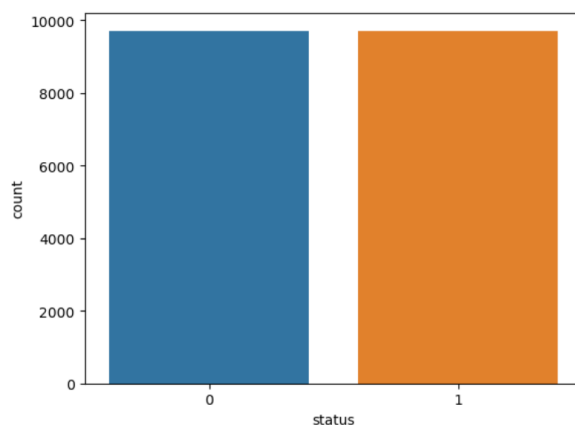


Podemos ver, por ejemplo, que la variable “google_index” tiene un alto nivel de correlación con la columna status. Esta columna será una de las que posteriormente se elegirán como las más relevantes en el siguiente apartado.

Podemos también verificar que tenemos un dataset balanceado pues hay un número similar de ejemplos para cada opción de “status”.

```
1 sns.countplot(data=df, x='status')
2 plt.show()
```

[22] ✓ - Command executed in 145 ms on 4:03:51 AM, 3/19/24



3. Azure ML Studio

Como se comentaba en la introducción, la idea es entrenar el modelo para ver si es capaz de predecir si un enlace es phishing o no. Para explorar las características que nos ofrece Azure ML Studio, vamos a:

- Lanzar un Automated ML para determinar el mejor modelo para clasificar el dataset.
- Crear en Designer una Pipeline que encuentre los mejores parámetros para ese tipo de modelo, personalizando más el entrenamiento.
- Finalmente, con este último modelo, crear en Designer una Pipeline de inferencia y desplegarlo para ver como funcionaría si tuviéramos una API.

3.1. Creación del clúster

En primer lugar se ha creado un clúster de computación para ejecutar los pasos descritos anteriormente. En la sección de costes se discutirá si el clúster es adecuado para las necesidades, pero se ha decidido elegir una Standard_DS3_v2, configurando el rango de nodos de 0 a 4. En las siguientes capturas se puede ver los pasos de creación.

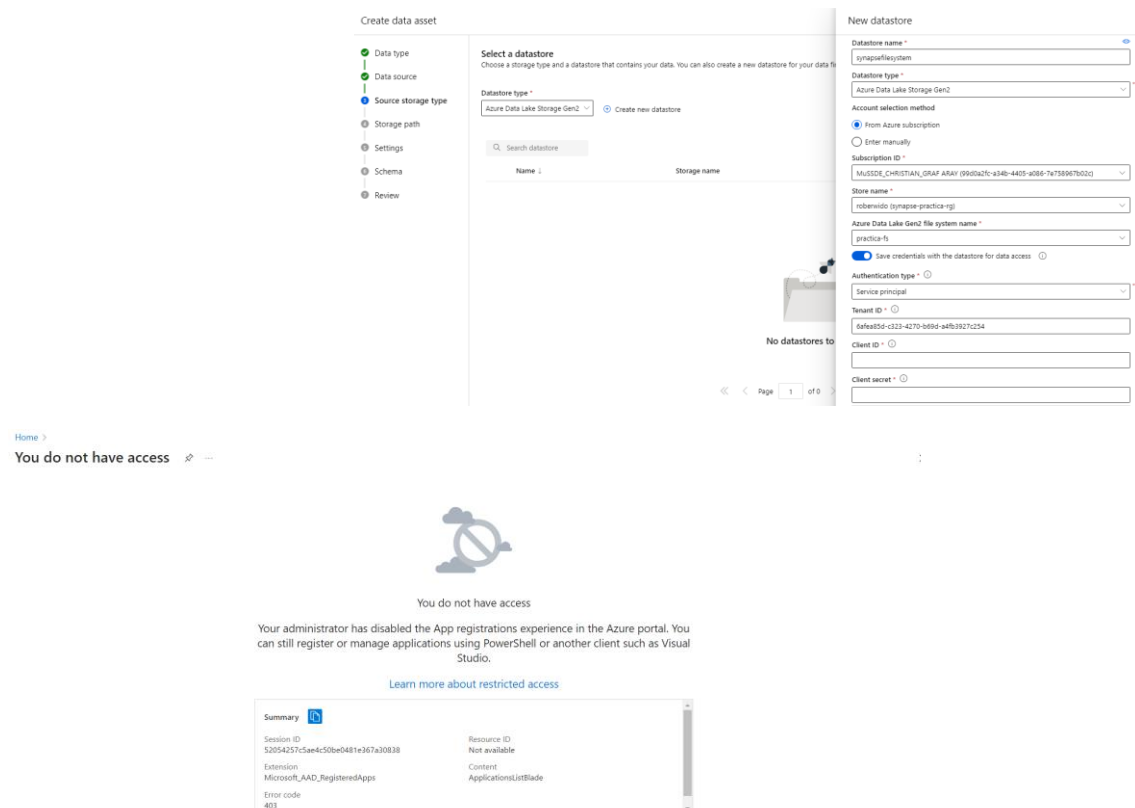
The first screenshot shows the 'Create compute cluster' wizard in Azure ML Studio. The 'Virtual Machine' tab is selected. The 'Location' is set to 'West Europe'. The 'Virtual machine tier' is set to 'Dedicated'. The 'Virtual machine type' is set to 'CPU'. The 'Virtual machine size' is set to 'Standard_DS3_v2'. The table below shows the available options:

Name	Category	Workload types	Available quota	Cost
Standard_DS1_v2 2 cores, 14GB RAM, 28GB storage	Memory optimized	Development on Notebooks (or other IDE) and light weight testing	24 cores	\$0.19/hr
Standard_DS3_v2 4 cores, 14GB RAM, 28GB storage	General purpose	Classical ML model training on small datasets	24 cores	\$0.27/hr
Standard_F4ls_v4 4 cores, 32GB RAM, 150GB storage	Memory optimized	Data manipulation and training on medium-sized datasets (1-10GB)	350 cores	\$0.35/hr
Standard_F4s_v2 4 cores, 8GB RAM, 32GB storage	Compute optimized	Data manipulation and training on large datasets (>10 GB)	6 cores	\$0.19/hr

The second screenshot shows the 'Configure Settings' step. The 'Standard_DS3_v2' is selected. The 'Compute name' is 'cpucluster'. The 'Minimum number of nodes' is 0 and the 'Maximum number of nodes' is 4. The 'Idle seconds before scale down' is 120. The 'Enable SSH access' checkbox is checked. The 'Add tags' section is empty.

3.2. Importación del dataset

Nuestra primera idea era poder importar directamente desde el filesystem creado para el primer apartado el dataset curado. Para ello, había que crear un datastore de tipo Data Lake Gen 2 y enlazar con el filesystem del usuario utilizado anteriormente. Como se ve en las siguientes dos capturas, nos pedían 2 campos especiales: Client ID y Client secret. Lamentablemente, la única forma de acceder a ellos está restringida por la organización (UPM) a nuestras cuentas de usuarios, y al no poder acceder al servicio “[App Registrations](#)” para obtener esos campos, hemos tenido que subir los datos desde GitHub como en el apartado anterior.



Así que, en las siguientes capturas se puede apreciar cómo se importa el [dataset curado desde GitHub](#). Para que funcione correctamente en Azure, es importante utilizar el enlace que ofrece GitHub de [tipo raw](#). Así, el formato y los tipos de datos se detectan automáticamente.

Create data asset

1 Data type
2 Data source

Set the name and type for your data asset

Name *

Description

Type *

Create data asset

✔ Data type

✔ Data source

● Web URL

● Settings

● Schema

● Review

Choose a source for your data asset

Choose the data source you want to create your asset from. A data source can be from a local storage location on your computer, from an attached datastore, from Azure publicly available web location.

From Azure storage

Create a data asset from registered data storage services including Azure Blob Storage, Azure file share, and Azure Data Lake.

From local files

Create a data asset by uploading files from your local drive.

From SQL databases

Create a dataset from Azure SQL database and Azure PostgreSQL database.

From web files

Create a data asset from a single file located at a public web URL.

From Azure Open Datasets

Create a dataset with one-click from pre-made data sets. These data sets are created by the general public and published as Azure Open Datasets.

Back

Next

Create data asset

✔ Data type

✔ Data source

● Web URL

● Settings

● Schema

● Review

Enter a web URL

Specify the URL of a public web page you want your data retrieved from.

Web URL *

https://raw.githubusercontent.com/ETSI-OGVD/practicaogvd3-24-grupo_christiangraf_carlosoliva/main/phishing_curated.csv

Skip data validation

If you choose to skip validation, we will not validate your data path, or try to access your data for preview and schema.

● Skip data validation

Back

Next

Create data asset

✔ Data type

✔ Data source

✔ Web URL

● Settings

● Schema

● Review

Settings

These settings determine how the data is parsed. The initial settings are automatically detected; you can change them as needed to reparse the data.

File format

Delimiter

Example

Encoding

Delimited

Comma

Field1,Field2,Field3

UTF-8

Column headers

Skip rows

All files have same headers

None

Dataset contains multi-line data

Note: Processing tabular files with multi-line data is slower because multiple CPU cores cannot be used to ingest the data in parallel. Checking this option may result in slower processing times.

Data preview

url	url_len...	hostna...	ip	total_of...	total_of...	total_o...	total_of...	total_of...	total_of...	total_of...	total_of...	total_of...	total_o...
http://w...	46	20	0	3	0	0	1	0	1	0	0	0	3
http://st...	128	120	0	10	0	0	0	0	0	0	0	0	3
http://w...	52	25	0	3	0	0	0	0	0	0	0	0	4
http://w...	21	13	0	2	0	0	0	0	0	0	0	0	3
https://...	28	19	0	2	0	0	0	0	0	0	0	0	3
https://s...	128	50	1	4	1	0	1	2	3	2	0	0	5
https://...	50	15	0	2	0	0	1	0	1	0	0	0	5
http://st...	51	14	0	5	0	0	0	0	0	0	0	0	4
http://w...	35	22	0	2	0	0	0	0	0	0	0	0	4

Back

Next

Review

Cancel

9

Create data asset

Schema

Column types are auto-detected based on the initial subset of the data and can be updated here. Values not aligning with the specified column type will fail conversion and would be either nullified or replaced with error value. Any conversions preview errors are non-blocking and you can proceed.

Search column name

Include	Column name	Type	Example values	Date format	Properties
<input type="checkbox"/>	Path	String		Not applicable to selected type	Not applicable to se...
<input checked="" type="checkbox"/>	url	String	http://www.progerchives.com/album.a...	Not applicable to selected type	Not applicable to se...
<input checked="" type="checkbox"/>	url_length	Integer	46, 126, 52	Not applicable to selected type	Not applicable to se...
<input checked="" type="checkbox"/>	hostname_length	Integer	20, 120, 25	Not applicable to selected type	Not applicable to se...
<input checked="" type="checkbox"/>	ip	Integer	0, 0, 0	Not applicable to selected type	Not applicable to se...
<input checked="" type="checkbox"/>	total_of	Integer	3, 10, 3	Not applicable to selected type	Not applicable to se...
<input checked="" type="checkbox"/>	total_ph	Integer	0, 0, 0	Not applicable to selected type	Not applicable to se...
<input checked="" type="checkbox"/>	total_ph0	Integer	0, 0, 0	Not applicable to selected type	Not applicable to se...
<input checked="" type="checkbox"/>	total_ph1	Integer	1, 0, 0	Not applicable to selected type	Not applicable to se...
<input checked="" type="checkbox"/>	total_ph2	Integer	0, 0, 0	Not applicable to selected type	Not applicable to se...

Back Next Cancel

Create data asset

Review

Review the settings for your data asset and make any changes as needed.

Data type

Name: phishing

Description: URLs that can be classified into phishing or not

Type: tabular

Data source

Type: WebURL

Web URL

Web URL: https://raw.githubusercontent.com/ETSI-OGV/practicaogv23-24-group-christiangraf_carlosoliva/main/phishing_curated.csv

Skip data validation: false

Settings

Delimiter: Comma

Encoding: UTF-8

File format

Back Create

Schema

Column name	Type
url	String
url_length	Integer
hostname_length	Integer
ip	Integer
total_of	Integer

(showing 5 of 87 columns)

3.3. AutoML: encontrar el mejor modelo de clasificación

Una vez importados los datos, lanzaremos un proceso de AutoML que puede consultarse [aquí](#). El proceso seguido se puede ver en las siguientes capturas, se ha permitido que se prueben todos los modelos de machine learning disponibles, con límites de tiempo de 15 minutos cada uno y 120 en total. Además, se ha permitido que se utilicen todos los nodos del clúster.

Universidad Politécnica de Madrid > ml-practica-ws > Training job

Submit an Automated ML job

Training method

Basic settings

Let's start with some basic information about your training job.

Job name *

experiment-phishing

Experiment name *

Select existing Create new

New experiment name *

experiment-automl-phishing

Description

Tags

Name Value Add

Submit an Automated ML job

PREVIEW

Training method

Basic settings

Task type & data

Task settings

Compute

Review

Task type & data

Choose the type of task that you would like your model to perform and the data to use for training. [Learn more](#)

Select task type *

Classification

Select data

Make sure your data is preprocessed into a supported format.

Create

Refresh

Show supported data assets only

Search

Name	Type	Created on	Modified on
phishing	Table	Mar 18, 2024 11:50 PM	Mar 18, 2024 11:50 PM

<<

<

Page 1

>

>>

25/Page

Submit an Automated ML job

PREVIEW

Training method

Basic settings

Task type & data

Task settings

Compute

Review

Task settings

Task type

Classification

Data

phishing [View data](#)

Target column *

status (Integer)

Classification settings

Enable deep learning

View additional configuration settings

View featurization settings

Limits

Max trials

4

Max concurrent trials

4

Max nodes

4

Universidad Politécnica de Madrid

mi-practica-ws

Training job

Submit an Automated ML job

PREVIEW

Training method

Basic settings

Task type & data

Task settings

Compute

Review

Max concurrent trials

4

Max nodes

4

Metric score threshold

Enter metric score threshold

Experiment timeout (minutes)

120

Iteration timeout (minutes)

15

Enable early termination

Validate and test

You can choose a validation type and select test data as an optional step.

Validation type

Train-validation split

Percentage validation of data *

15

Automated ML recommends that between 10 and 30 percent of data is held out for validation

Test data

None

Back

Next

11

Additional configuration

Primary metric ⓘ

AUCWeighted

☒ Explain best model ⓘ

☐ Enable ensemble stacking ⓘ

☒ Use all supported models

Blocked models ⓘ

A list of models that Automated ML will not use during training.

Positive class label ⓘ

1

Universidad Politécnica de Madrid > ml-practica-ws > Training job

Submit an Automated ML job REVIEW

- Training method
- Basic settings
- Task type & data
- Task settings
- Compute**
- Review

Compute

Select and configure the compute resource for executing your training job.

Select compute type

Compute cluster

Select Azure ML compute cluster *

cpucluster

[+ New](#)

Submit an Automated ML job REVIEW

- Training method
- Basic settings
- Task type & data
- Task settings
- Compute
- Review**

Review

Review or make changes to your job before submission.

Basic settings

Name: automl-phishing

Experiment name: experiment-automl-phishing

Description: ...

Timeout (hours): ...

Tags: ...

Task type & data

Task type: Classification

Data: phishing

Task settings

Target column: status

Limits: ✓

Max trials: 4

Max concurrent trials: 4

Max nodes: 4

Metric score threshold: --

Experiment timeout (minutes): 120

Iteration timeout (minutes): 15

Enable deep learning: No

Validate type: Train-validation split

Percentage validation of data: 15

Compute settings

Compute type: Azure ML compute cluster

Selected Azure ML compute cluster: cpucluster

Una vez finalizado, observamos que [el mejor modelo](#) ha sido un XGBoostClassifier. La ventaja de Azure es que nos permite activar la explicabilidad, que en la última captura nos devuelve las 4 features más importantes. Además, también explica que tratamiento previo de los datos ha hecho (incluso descartando features).

Universidad Politécnica de Madrid > ml-practica-ws > Jobs > experiment-automl-phishing > automl-phishing

automl-phishing Completed

Overview Data guardrails Models + child jobs Outputs + logs Child jobs

[Refresh](#) [Edit and submit \(preview\)](#) [Register model](#) [Cancel](#) [Delete](#) [Compare \(preview\)](#)

Properties

Status: Completed

Created on: Mar 18, 2024 11:58 PM

Start time: Mar 18, 2024 11:58 PM

Duration: 13m 16.22s

Compute duration: 13m 16.22s

Compute target: cpucluster

Name: automl-phishing

Script name: --

Created by: CARLOS OLIVA LOPEZ

Job type: Automated ML

Experiment: experiment-automl-phishing

Arguments: None

See all properties: [Kan SCIN](#)

See YAML job definition: [Job YAML](#)

Inputs

Input name: training_data

Data asset: phishing1

Asset URI: [experiment-phishing1](#)

Outputs

Output name: best_model

Model: azureml(automl-phishing_1_output_mlflow_log_model_264064485:1

Asset URI: [experiment-automl-phishing_1_output_mlflow_log_model_264064485:1](#)

Best model summary

Algorithm name: MaxAbsScaler, XGBoostClassifier

Hyperparameters: [See View hyperparameters](#)

AUC weighted: 0.99937 [View all other metrics](#)

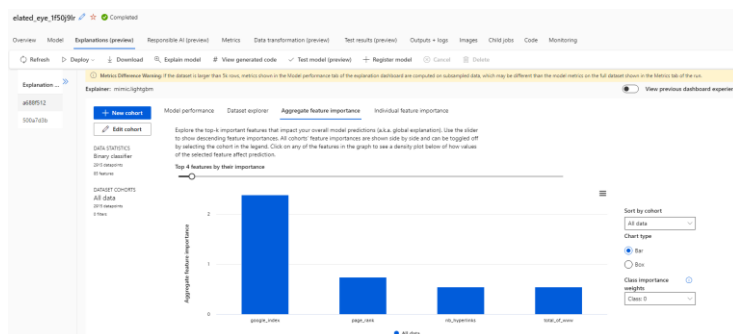
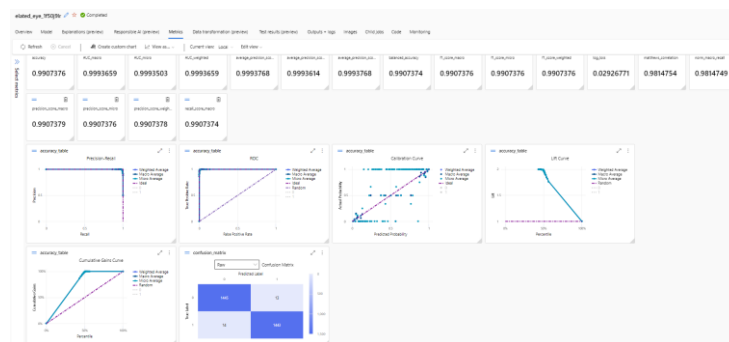
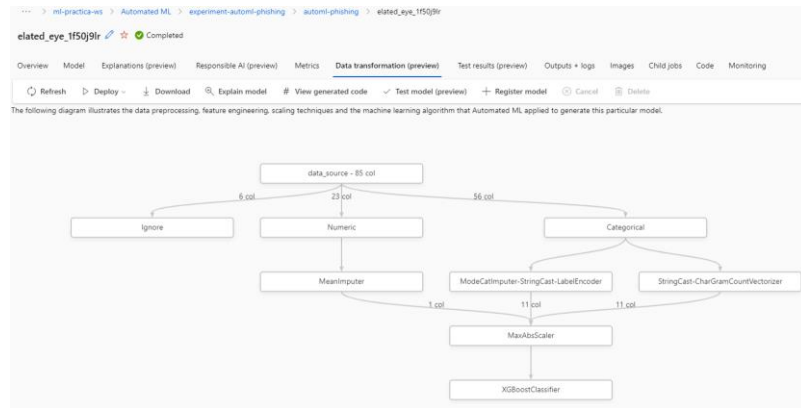
Sampling: 100.00 %

Registered models: No registration yet

Tags

ml_experiment_000: 0.3768881.602850.307929.2 iteration_000: 0.1:2.3

Práctica de Azure – Carlos Oliva López y Christian Graf Aray



3.4. Pipelines de Designer: un entrenamiento más personalizado

Como aclaración antes de empezar este punto, hemos encontrado que la forma de procesar un dataset de entrada no es igual cuando se utiliza AutomatedML que Designer. En AutoML el dataset curado funciona perfecto, pero para Designer hemos tenido que realizar una segunda curación (se encuentra en el [notebook del EDA](#)). Lo que se ha hecho es sustituir los nombres de columnas que tenían caracteres especiales como “total_of?” o “total_of\$” por, respectivamente, “total_of_questionmarks” y “total_of_dollars” (y así con todas las columnas que cumplieran esta condición). Por lo tanto, el nuevo dataset curado para pipelines se puede encontrar [aquí](#). Adjuntamos una captura del error que ocurría al utilizar el anterior dataset, que claramente no era soportado por esos caracteres.

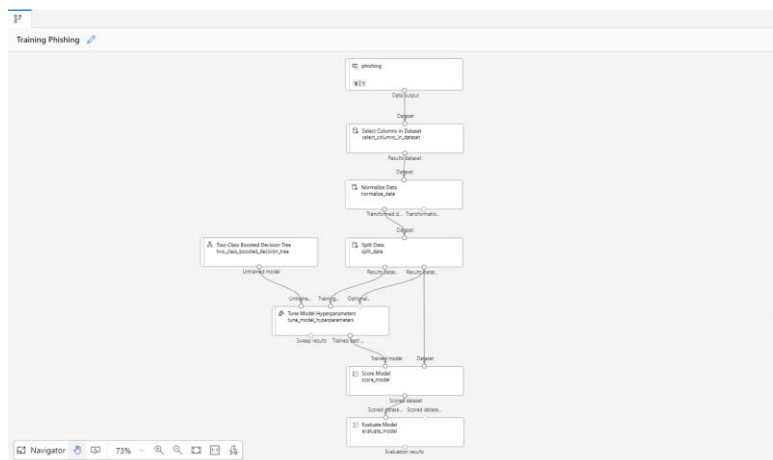
```

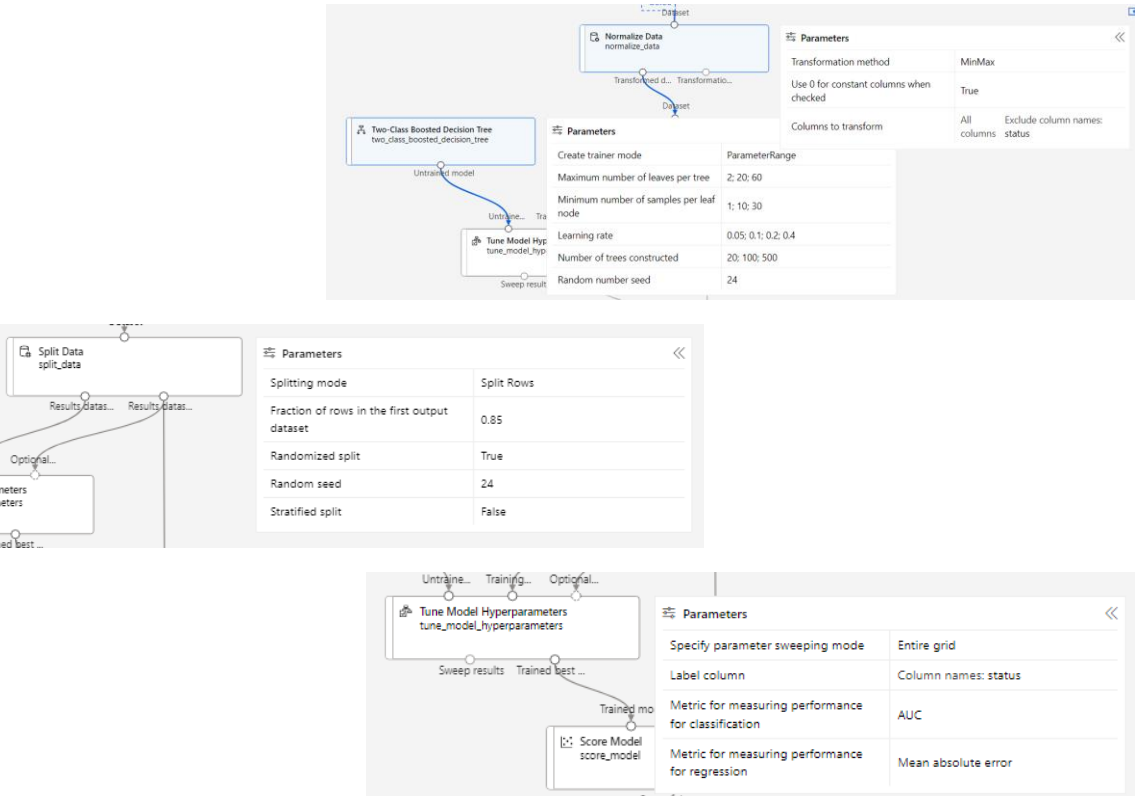
172 return self.reader.read_all(column_indices=column_indices,
173 file "pyarrow_parquet.py", line 1136, in pyarrow_parquet.ParquetReader.read_all
174 File "pyarrow/error.pxi", line 94, in pyarrow.lib.check_status
175 pyarrow.lib.ArrowInvalid: Column 4 named total_of_ expected length 19431 but got length 38862
176
177 The above exception was the direct cause of the following exception:
178
179 Traceback (most recent call last):
180 File "/azureml-envs/azureml_8f317849db35f281450cf74333640b98/lib/python3.8/site-packages/azureml/studio
181 execute(sys.argv)
182 File "/azureml-envs/azureml_8f317849db35f281450cf74333640b98/lib/python3.8/site-packages/azureml/studio
183 return execute_with_cli(original_args)
184 File "/azureml-envs/azureml_8f317849db35f281450cf74333640b98/lib/python3.8/site-packages/azureml/studio
185 ret = func(*args, **kwargs)
186 File "/azureml-envs/azureml_8f317849db35f281450cf74333640b98/lib/python3.8/site-packages/azureml/studio
187 do_execute_with_env(parser, FolderRuntimeEnv())
188 File "/azureml-envs/azureml_8f317849db35f281450cf74333640b98/lib/python3.8/site-packages/azureml/studio
189 ModuleReflector(parser.module_entry, env).exec()
190 File "/azureml-envs/azureml_8f317849db35f281450cf74333640b98/lib/python3.8/site-packages/azureml/studio
191 self._handle_exception(ex)
192 File "/azureml-envs/azureml_8f317849db35f281450cf74333640b98/lib/python3.8/site-packages/azureml/studio
193 raise exception
194 File "/azureml-envs/azureml_8f317849db35f281450cf74333640b98/lib/python3.8/site-packages/azureml/studio
195 reflected_input_ports = self._reflect_input_ports(input_ports)
196 File "/azureml-envs/azureml_8f317849db35f281450cf74333640b98/lib/python3.8/site-packages/azureml/studio
197 value = self._env.handle_input_port(annotation, input_value)
198 File "/azureml-envs/azureml_8f317849db35f281450cf74333640b98/lib/python3.8/site-packages/azureml/studio
199 ErrorHappening.rerthrow(e, InvalidDatasetError(
200 File "/azureml-envs/azureml_8f317849db35f281450cf74333640b98/lib/python3.8/site-packages/azureml/studio
201 raise err from e
202 azureml.studio.common.error.InvalidDatasetError: Dataset contains invalid data, failed to load dataset

```

3.4.1. Pipeline de entrenamiento

A continuación, realizaremos la última prueba con Azure ML Studio. Diseñaremos [una Pipeline de entrenamiento](#) con la herramienta gráfica Designer. En este caso, hemos aprovechado para definir cómo queremos que sea el entrenamiento. En las siguientes capturas se aprecia la configuración de cada uno de los bloques del pipeline. A destacar, la normalización aplicada y el uso de un ParameterRange a la hora de entrenar, lo que permite hacer una búsqueda de hiperparámetros en forma de grid.





Set up pipeline job

1 Basics

2 Inputs & outputs

3 Runtime settings

4 Review + Submit

Basics

Experiment name

☐ Select existing ☒ Create new

New experiment name *

experiment-train-phishing

Job display name

Training Phishing

Job description

Pipeline to train on the phishing dataset

Job tags

Name : Value Add

Set up pipeline job

1 Basics

2 Inputs & outputs

3 Runtime settings

4 Review + Submit

Runtime settings

Default compute ⓘ

Select compute type

Compute cluster

Select Azure ML compute cluster

cpucluster

Create Azure ML compute cluster Refresh Compute

Default datastore ⓘ

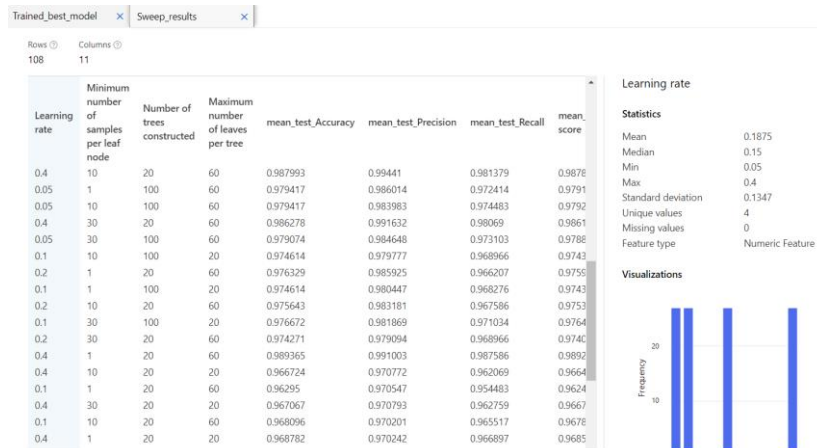
Select datastore *

workspaceblobstore

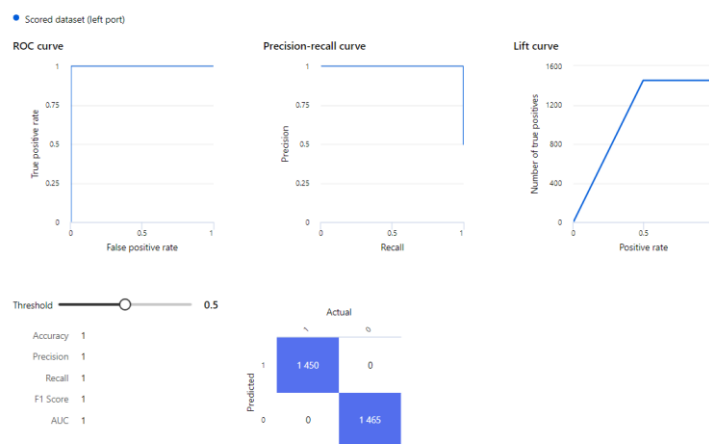
Advanced settings

☒ Continue on step failure ⓘ

Una vez finalizado [el experimento](#), podemos observar el valor de las distintas métricas para cada conjunto de hiperparámetros.

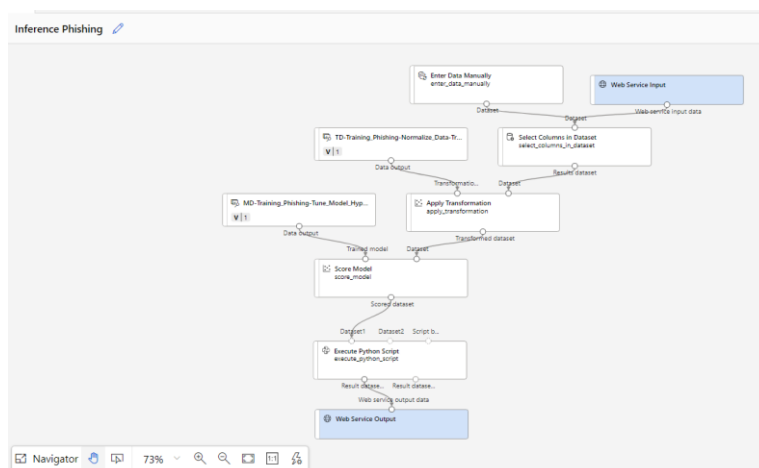


Para el mejor modelo, a pesar de haber utilizado un split del conjunto total como validación durante el entrenamiento, se ha producido overfitting ya que como vemos, los resultados son del 100%.



3.4.2. Pipeline de inferencia y despliegue

Por último, crearemos [una Pipeline de inferencia](#) en la que introduciremos [unos datos manualmente](#) y comprobaremos la salida.



A destacar, la configuración del script de Python que se encuentra [aquí](#). Simplemente renombra las columnas de salida. Si comprobamos las salidas, y comparamos con el [dataset original](#), veremos que ha acertado (se ha utilizado desde la línea 16 a la línea 24).

Result_Dataset		
Rows	Columns	
9	3	
url	PhishingPrediction	Probability
		
https://www.youtube.com/channel/UCLSj5w7gU48j0iaesjfwjQ	0	0.000002
		
https://fieldstonerpmyscharisma.com/bz/p/danfiles/E275H1wQSDpRtP4DIeCt5H8Xo4sRqJgc9j807wUYV7FgIw=Kt959K	1	1
		
http://www.webopedia.com/TERM/C/CI.html	0	0
		
http://www.allmenus.com/n/c/charlotte/89187-mimosagrill/menu/	0	0
		
http://www.makeuseof.com/tag/p2p-peer-peer-file-sharing-works/	0	0
		
http://bdo-onlineverify.xyz/bdoverificati on/secureity/verify/login.php	1	1
		
http://www.picfront.org/	0	0
		
http://kam-net.ci202584619/verification.n.php	1	1
		
http://www.bedennews.nl/	0	0

Por último, desplegaremos este modelo en un endpoint para que, hipotéticamente, alguien pudiera utilizarlo como API con sus propios datos. Se puede observar que, con estos [datos en formato JSON](#), las salidas son correctas y el servicio responde.

Set up real-time endpoint

☒ Deploy new real-time endpoint

☐ Replace an existing real-time endpoint

Name *

predict-phishing

Description

Classify if a url is a phishing or a legitimate one

Compute type *

Azure Container Instance

> Advanced

Deploy

Cancel

[Universidad Politécnica de Madrid](#)
[mi-practica-vo](#)
[Endpoints](#)
[predict-phishing](#)

[predict-phishing](#)

[Details](#)
[Test](#)
[Console](#)
[Logs](#)

Input data to test endpoint

```

{
  "inputs": {
    "inputs": {
      "url": "https://www.youtube.com/channel/UCJ5jwGj0d8gJwYfjwF10",
      "url_length": 0,
      "hostname_length": 15,
      "ip": 0,
      "total_of_periods": 2,
      "total_of_hyphens": 0,
      "total_of_ats": 0,
      "total_of_questionmarks": 0,
      "total_of_underscores": 0,
      "total_of_equals": 0,
      "total_of_underscores": 0,
      "total_of_tildes": 0,
      "total_of_percentages": 0,
      "total_of_slashes": 0,
      "total_of_stars": 0,
      "total_of_colons": 1,
      "total_of_commas": 0,
      "total_of_semicolon": 0,
      "total_of_dollars": 0,
      "total_of_ams": 0,
      "total_of_com": 0
    }
  }
}

```

Test result

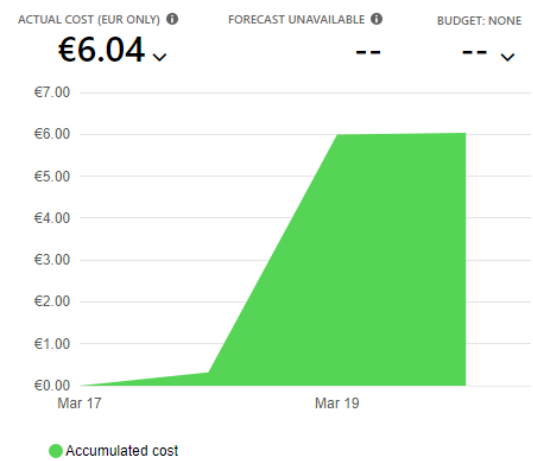
```

{
  "results": {
    "isMaliciousOutput": {
      0: {
        "url": "https://www.youtube.com/channel/UCJ5jwGj0d8gJwYfjwF10",
        "probability": 0.00000034100192206458
      }
    }
  }
}

```

4. Costes

En esta última sección, analizaremos los costes de las operaciones realizadas en los últimos dos apartados. Lo primero es echar un vistazo a los costes acumulados a lo largo de la realización de la práctica. En total, un gasto de algo más de 6€ en 3 días.



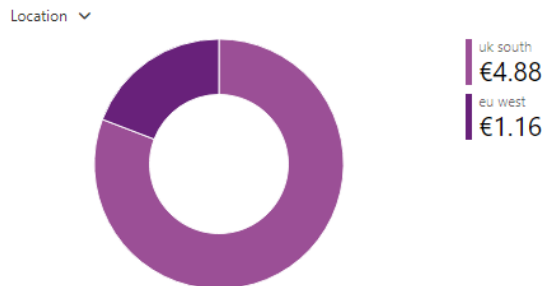
Analizando el coste diario, observamos que el grueso del coste fue el día 19 y causado por Synapse Analytics.



Si analizamos por recurso, el gasto de Synapse Analytics es única y exclusivamente por el uso de los vCores, es decir, por cada vez que se ha arrancado el pool de PySpark. El gasto por el filesystem del usuario creado es ínfimo (<0,01€). En cambio, los gastos de ML Studio están más repartidos entre varios componentes. Los primeros 0,49€ que aparecen son los equivalentes al clúster creado, mientras que los 0,17€ son equivalentes al endpoint del modelo desplegado. El resto de gastos, como por ejemplo el load balancer, se debe a los gastos asociados cada vez que se lanzaba un pipeline y se requería escalar el número de nodos del clúster creado. Por lo tanto, el propio reescalamiento vemos que lleva un coste asociado.

Resource	Resource type	Location	Resource group name	Tags	Cost
synapse-practica-rg / sparkpool	Synapse workspace	uk south	synapse-practica-rg	--	€4.86
Service name: Azure Synapse Analytics					Cost: €4.86
Meter: vCore					
ml-practica-rg	Azure Machine Learning workspace	eu west	ml-rg	amlresourceprovisionerbatch computepeserver...	€1.02
Service name: Virtual Machines					Cost: €0.49
Meter: D3 v2/D3S v2					
Service name: Virtual Machines					Cost: €0.17
Meter: B4ds v4					
Service name: Storage					Cost: €0.13
Meter: P10 LRS Disk					
Service name: Load Balancer					Cost: €0.13
Meter: Standard Included LB Rules and Outbound Rules					
Service name: Load Balancer					Cost: €0.07
Meter: Standard Data Processed					
Service name: Virtual Network					Cost: €0.03
Meter: Standard IPv4 Static Public IP					
Service name: Bandwidth					Cost: <€0.01
Meter: Intra Continent Data Transfer Out					
Service name: Bandwidth					Cost: €0
Meter: Standard Data Transfer Out					

Por ultimo a nivel general, podemos distinguir los gastos de Synapse Analytics y ML Studio porque, debido a limitaciones de Azure en el momento de intentar crear el primero, no nos permitía la región West Europe. Así, dado que la latencia no era un problema, escogimos la región más cercana (South UK) y los gastos han quedado repartidos en dos zonas.



Pasando a analizar más concretamente los gastos de Synapse, podemos apreciar en el propio workspace todas las ejecuciones realizadas y el tiempo que duraron. En este caso, fueron 4.

Application name	Submitter	Submit time	Status	Pool	Type	Attempts	Livy ID	Running duration
EDA_Phising_sparkpool_1710...	c.graf@alumnos.upm.es	3/19/2024, 11:17:18 PM	Stopped	sparkpool	Spark session	All Attempts	3	3m 59s
EDA_Phising_sparkpool_1...	c.oliva@alumnos.upm.es	3/19/2024, 4:12:44 PM	Stopped	sparkpool	Spark session	All Attempts	2	7m 11s
EDA_Phising_sparkpool_1710...	c.graf@alumnos.upm.es	3/19/2024, 3:59:21 AM	Stopped	sparkpool	Spark session	All Attempts	1	6m 9s
EDA_Phising_sparkpool_1710...	c.graf@alumnos.upm.es	3/19/2024, 2:52:36 AM	Stopped (session time	sparkpool	Spark session	All Attempts	0	1h 44m

Y los vCores utilizados por día, que encajan perfectamente con las ejecuciones de arriba.



Teniendo en cuenta que el vCore por hora es de 0,56€ en la familia del pool de PySpark seleccionado, podríamos calcular de la primera sesión que se ejecutó, que fue la más larga de todas con una duración de 1 hora y 44 minutos.

$$GastoSesion1\ vCore = \left(1h + \frac{44}{60}h\right) * 0,56 \frac{\text{€}}{h} = 0,97\text{€ } vCore$$

Como el pool puede utilizar entre 3 y 6 vCores, el gasto de esa primera sesión está comprendido entre 2,91€ y 5,82€. Si analizamos la grafica de costes diarios veremos que el coste estuvo cerca de los 5€. Por lo tanto, habría que analizar si es rentable el uso de Synapse Analytics. ¿Vale la pena ese gasto para un EDA de un dataset verdaderamente pequeño? Lo cierto es que esa sesión debería haberse acortado, ya que la mayoría del tiempo se utilizó para programar código y probar en tiempo real, mientras que el paradigma a adoptar en estos casos (si se quiere ahorrar) es programar todo de antemano

y solo ejecutar cuando sea necesario hacer una prueba. Es decir, que el tiempo durante el cual el pool esté activo sea única y exclusivamente utilizado para ejecución. Sería interesante poder comparar con el gasto en luz equivalente a hacer estas ejecuciones en un ordenador local, aunque lo cierto es que no se tendría la misma infraestructura ni disponibilidad de cores. Por lo que hemos visto, el gasto puede dispararse muy rápidamente si se descuidan los tiempos de ejecución, lo que podría suponer gastos innecesarios para una empresa o incluso un problema económico para pymes o pequeñas empresas emergentes.

Por último, aunque los gastos de ML Studio han sido mínimos (1,18€) conviene analizarlos también. El coste aproximado por nodo utilizado es de 0,25€ en el clúster escogido, un DS3 v2. Si recordamos lo dicho anteriormente, el gasto de todas las ejecuciones en este clúster fue de 0,49€, al que habría que sumar los costes de 0,13€ por reescalamiento automático.

Resource properties

Virtual machine size

Standard_DS3_v2 (4 cores, 14 GB RAM, 28 GB disk)

Processing unit

CPU - General purpose

Estimated cost

\$0.27/hr per node

Teniendo en cuenta los experimentos lanzados (pueden verse abajo), tanto los intentos acertados como fallidos, es cierto que no parece un gasto muy alto y hasta parece rentable. De hecho, hay que tener en cuenta que con ese gasto se han entrenado varios modelos (aunque nos hayamos quedado con los mejores solo) y se ha realizado una pipeline de inferencia también.

Display name (5 visualized)	Parent job name	Experiment	Status	Created on	Duration	Created by	Compute target
affable_key_d3tfv6vs		prepare_image	Completed	Mar 19, 2024 5:47 PM	18m 12s	CARLOS OLIVA LO...	Serverless
> Inference Phishing (5)		experiment-inference-phishing	Completed	Mar 19, 2024 5:26 PM	9m 10s	CARLOS OLIVA LO...	
> Training Phishing (7)		experiment-train-phishing	Completed	Mar 19, 2024 4:54 PM	14m 32s	CARLOS OLIVA LO...	
> Training Phishing (2)		experiment-train-phishing-1	Failed	Mar 19, 2024 12:30 PM	10m 33s	CARLOS OLIVA LO...	
> automi-phishing (9)		experiment-automi-phishing	Completed	Mar 18, 2024 11:58 PM	13m 16s	CARLOS OLIVA LO...	cpucluster

Se podrían haber ahorrado gastos si se hubiera escogido un clúster más humilde, ya que lo cierto es que los entrenamientos realizados eran de machine learning y no de deep learning con redes neuronales. Sería interesante analizar la relación de bajar los recursos de computo para ver si compensa, ya que habría que tener en cuenta que el tiempo de ejecución aumentaría. Por último, lo que sí que parece caro es levantar un endpoint con el modelo para poder utilizarlo como API. Teniendo en cuenta que estuvo levantado tan solo 18 minutos (como se aprecia en el primer job de la captura de arriba), el consumo de 0,17€ nos hace una idea de cómo de rápido subiría el precio, aproximadamente 14€ por día levantado. Al igual que se comentaba antes, también sería interesante analizar qué otros tipos de recursos endpoint ofrece Azure, ya que la inferencia de un modelo como el que hemos entrenado (básicamente un árbol de decisión) requiere muy poco cómputo.

5. Referencias

- Dataset - <https://www.kaggle.com/datasets/winson13/dataset-for-link-phishing-detection>
- Repositorio de GitHub - https://github.com/ETSISI-OGVD/practicaogvd23-24-grupo_christiangraf_carlosoliva/tree/main
- Suscripción de Azure - <https://portal.azure.com/#@upm365.onmicrosoft.com/resource/subscriptions/99d0a2fc-a34b-4405-a086-7e758967b02c/overview>
- Workspace de Synapse Analytics - <https://web.azuresynapse.net/?workspace=%2fsubscriptions%2f99d0a2fc-a34b-4405-a086-7e758967b02c%2fresourceGroups%2fsynapse-practica-rg%2fproviders%2fMicrosoft.Synapse%2fworkspaces%2fsynapse-practica-ws>
- Workspace de ML Studio - <https://ml.azure.com/?tid=6afea85d-c323-4270-b69d-a4fb3927c254&wsid=/subscriptions/99d0a2fc-a34b-4405-a086-7e758967b02c/resourceGroups/ml-rg/providers/Microsoft.MachineLearningServices/workspaces/ml-practica-ws>
- Get Started with Azure Synapse Analytics - <https://learn.microsoft.com/en-us/azure/synapse-analytics/get-started>
- What is automated machine learning (AutoML)? - <https://learn.microsoft.com/en-us/azure/machine-learning/concept-automated-ml?view=azureml-api-2>
- Tutorial: Train a classification model with no-code AutoML in the Azure Machine Learning studio - <https://learn.microsoft.com/en-us/azure/machine-learning/tutorial-first-experiment-automated-ml?view=azureml-api-2>
- Explore classification with Azure Machine Learning Designer - <https://microsoftlearning.github.io/AI-900-AIFundamentals/instructions/02b-create-classification-model.html>
- El resto de enlaces que han ido apareciendo a lo largo de esta memoria.