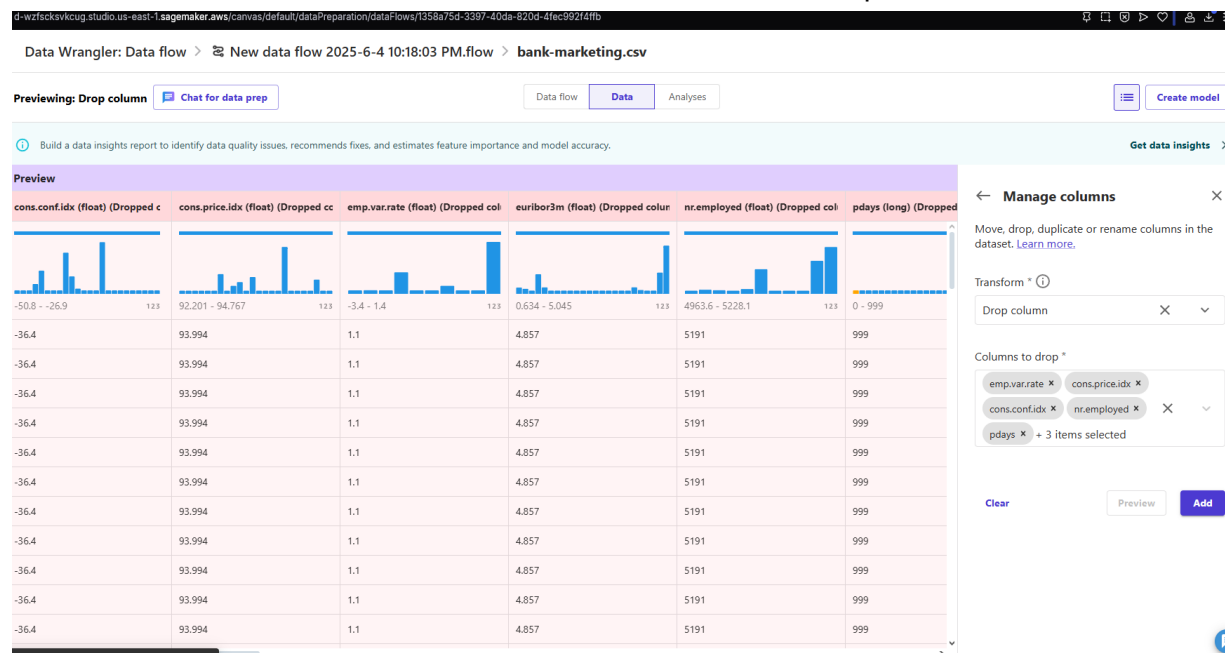


Actividad 9: refuerzo con AWS Sagemaker Autopilot

Antes de empezar con los puntos de la actividad voy a mostrar el proceso que seguí la primera vez que hice la tarea, menos detalladamente y que nos llevó a que nos pasara el mismo problema en el apartado del inference que cuando hicimos el ejemplo en clase el día 20 de mayo.

En el proceso de seguir las pautas para agilizar los pasos si hacemos el modelo sobre el dataset inicial sin meter el transform nos saldrán estos dos problemas:



Aquí ya nos avisa de la posibilidad de que de problema realizar el transform. Después en la vez que realizamos el modelo sobre el dataset ya modificado no nos dará este problema.

Select dataset for predictions

To make predictions on a dataset, select it or import it. The dataset that you select must have the same number of feature columns as the training dataset.

Search datasets in Canvas

Name	Columns	Rows	Cells	Created	Status	
bank	V1	21	41,188	864,943	06/04/2025 11:35 PM	Ready
bank-inference	V1	12	100	1,200		Incompatible
canvas-sample-product-descriptions.csv	V1	5	120	600		Incompatible
canvas-sample-shipping-logs.csv	V1	12	1,000	12,000		Incompatible
canvas-sample-housing.csv	V1	10	1,000	10,000		Incompatible
canvas-sample-databricks-dolly-15k.csv	V1	2	2,000	4,000		Incompatible
canvas-sample-loans-part-2.csv	V1	5	1,000	5,000		Incompatible
canvas-sample-retail-electronics-forecasting.csv	V1	6	40,500	243,000	06/04/2025 10:10 PM	Incompatible
canvas-sample-maintenance.csv	V1	9	1,000	9,000	06/04/2025 10:10 PM	Incompatible
canvas-sample-loans-part-1.csv	V1	19	1,000	19,000	06/04/2025 10:10 PM	Incompatible
canvas-sample-diabetic-readmission.csv	V1	16	1,000	16,000	06/04/2025 10:10 PM	Incompatible

Close

Generate predictions

Como podemos observar nos dice que faltan columnas para realizar las predicciones y no nos deja activar el dataset de inference por eso mismo. Con la finalidad de poder hacerlo de esta manera por si no encontraba otra solución, hice un edit del dataset de inference en pandas y le añadí esas columnas, rellenando sus datos con el valor medio y en el caso de las categoricas con el resultado más común de la columna.

Select dataset for predictions

To make predictions on a dataset, select it or import it. The dataset that you select must have the same number of feature columns as the training dataset. [?](#) [+ Create dataset](#)

Search datasets in Canvas

	Name	Columns	Rows	Cells	Created	Status
<input checked="" type="radio"/>	bank-inference-modificado	20	100	2000	06/05/2025 4:09 PM	Ready
<input type="radio"/>	bank	21	41,188	864,948	06/04/2025 11:35 PM	Ready
<input type="radio"/>	canvas-sample-product-descriptions.csv	5	120	600	06/04/2025 10:10 PM	Incompatible ?
<input type="radio"/>	canvas-sample-shipping-logs.csv	12	1000	12,000	06/04/2025 10:10 PM	Incompatible ?
<input type="radio"/>	canvas-sample-housing.csv	10	1000	10,000	06/04/2025 10:10 PM	Incompatible ?
<input type="radio"/>	canvas-sample-databricks-dolly-15k.csv	2	2000	4000	06/04/2025 10:10 PM	Incompatible ?
<input type="radio"/>	canvas-sample-loans-part-2.csv	5	1000	5000	06/04/2025 10:10 PM	Incompatible ?
<input type="radio"/>	canvas-sample-retail-electronics-forecasting.csv	6	40,500	243,000	06/04/2025 10:10 PM	Incompatible ?
<input type="radio"/>	canvas-sample-maintenance.csv	9	1000	9000	06/04/2025 10:10 PM	Incompatible ?
<input type="radio"/>	canvas-sample-loans-part-1.csv	19	1000	19,000	06/04/2025 10:10 PM	Incompatible ?
<input type="radio"/>	canvas-sample-diabetic-readmission.csv	16	1000	16,000	06/04/2025 10:10 PM	Incompatible ?

[Close](#) [Generate predictions](#)

Aquí observamos que ya nos permite seleccionar el dataset de inference.

Cargamos ese dataset que hemos generado dentro de la pestaña de canvas de datasets para poder ver que se hizo la predicción correctamente.

Datasets > datasetPrediction [V1](#) [+ Create a data flow](#) [Update dataset](#) [+ Create a model](#) [Dataset details](#)

[Data](#) [Version history](#) [Auto update](#)

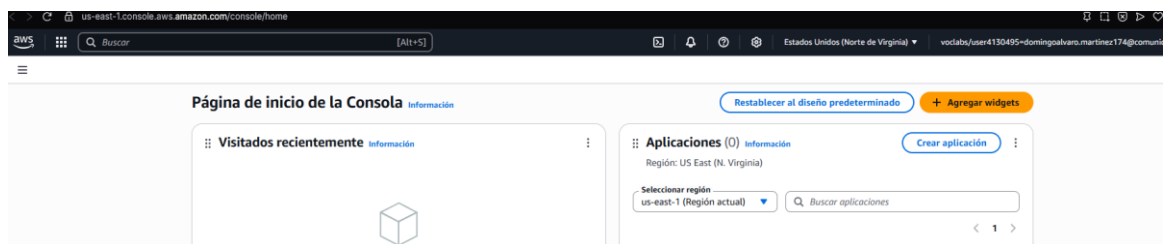
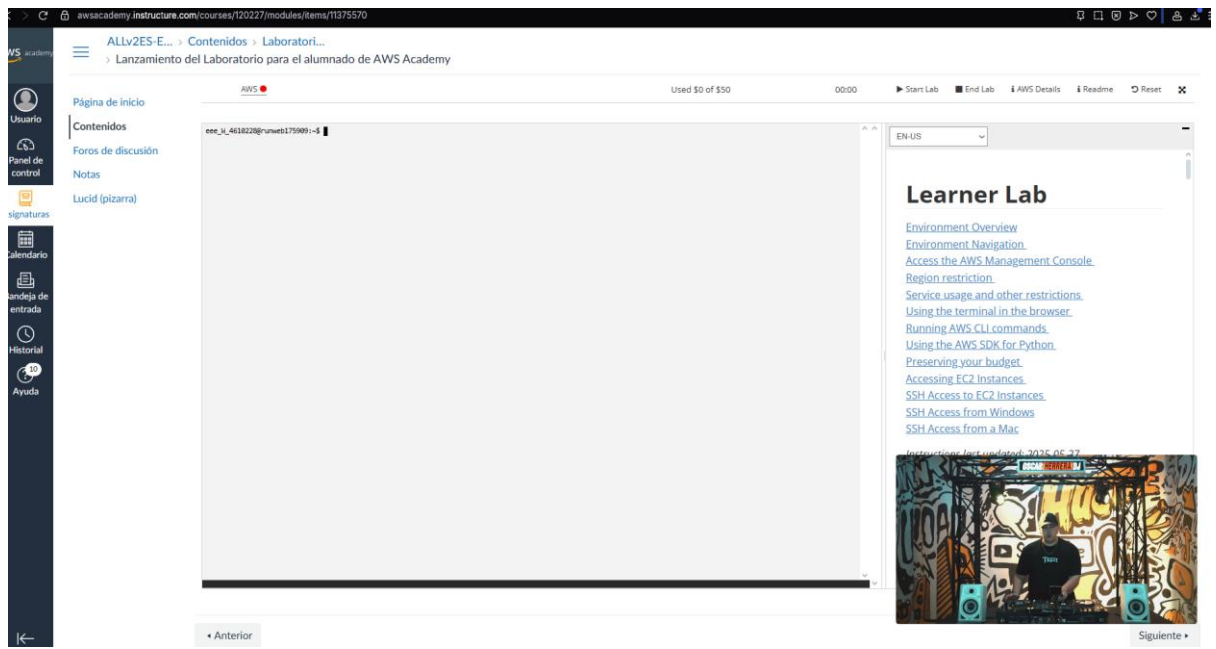
Previewing up to the first 100 rows of datasetPrediction

y	probability	age	job	marital	education	default	housing	loan
no	0.9998318553	0.6226415094	technician	married	high.school	no	no	yes
no	0.999846518	0.5849056604	unknown	married	unknown	unknown	yes	no
no	0.999524653	0.1698113208	blue-collar	married	basic.9y	no	no	no
no	0.9998098612	0.2264150943	admin.	married	high.school	no	no	no
no	0.9995096922	0.0566037736	housemaid	married	high.school	no	yes	no
no	0.9993556738	0.641509434	retired	married	professional.course	no	yes	yes
no	0.999589068	0.4528301887	services	married	high.school	unknown	yes	no
no	0.9998573065	0.5094339623	admin.	divorced	university.degree	unknown	yes	no
no	0.9997653365	0	entrepreneur	married	university.degree	no	yes	yes
no	0.9998453259	0.2264150943	technician	divorced	professional.course	no	yes	yes
no	0.9997922778	0.1886792453	blue-collar	married	basic.9y	no	no	no
no	0.9995871186	0.3396226415	blue-collar	single	basic.4y	no	yes	no

Dataset type: Tabular Total dataset cells (columns x rows): 2200 (22 x 100) Data source: Local

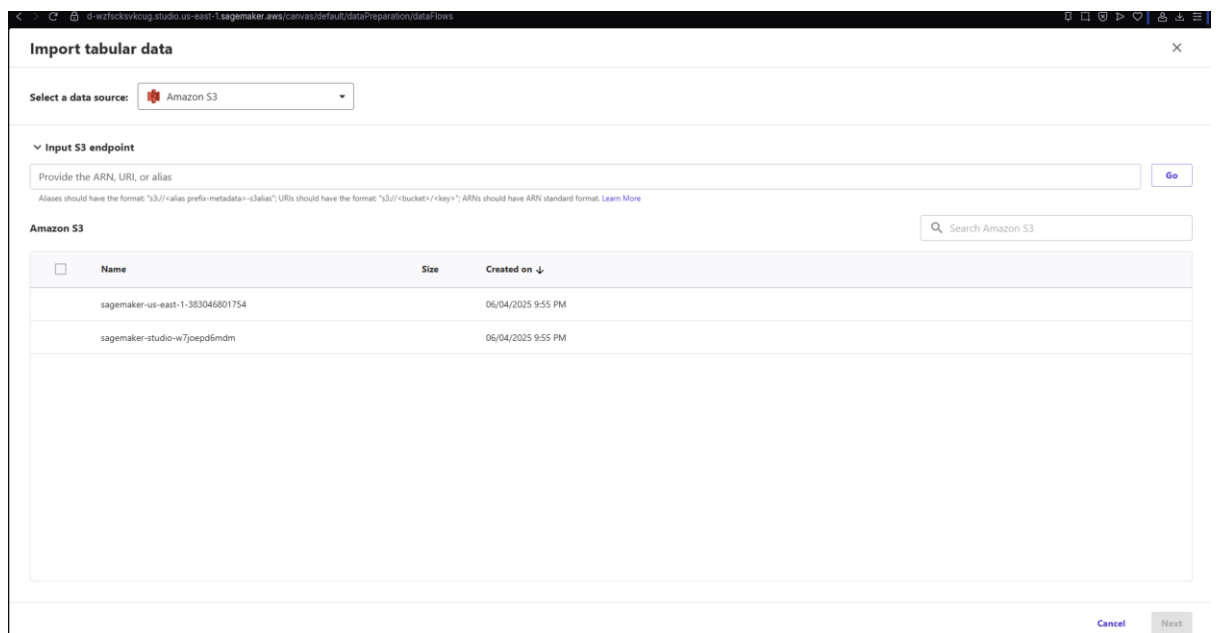
Vemos que ha generado bien el dataset de predicción y vamos ya con el ejercicio con su proceso para que el inference no de problemas a la hora de cargar ese dataset.

1. Empezamos con la entrada a nuestro laboratorio de AWS.

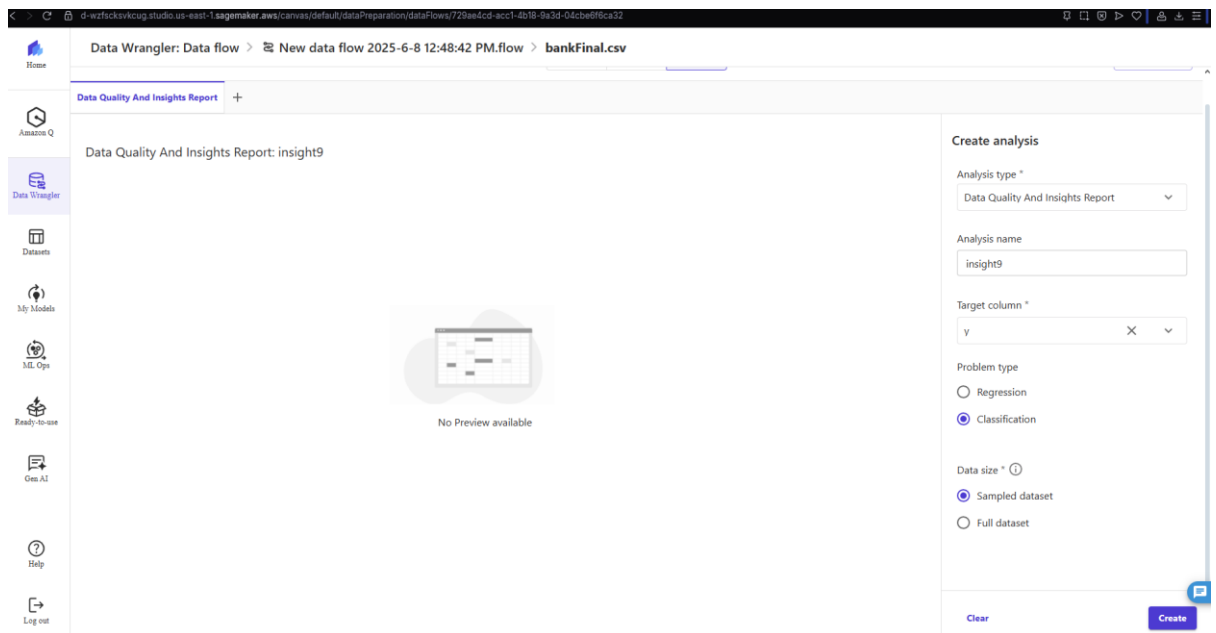


Ya tenemos la cuenta cargada, procedemos a entrar a canvas.

2. Carga del dataset.



3. Generamos data insight.



4. Analizamos el resultado de estos.

Data Wrangler: Data flow > New data flow 2025-6-8 12:48:42 PM.flow > bankFinal.csv

Step 2: Data types Data flow Data Analyses Create model

Data Quality And Insights Report +

Data quality and insights report: insight9

Target column	Type	Dataset	Date
y	Classification	bankFinal.csv	8 de junio de 2025, 12:54 CEST

Summary

Dataset statistics

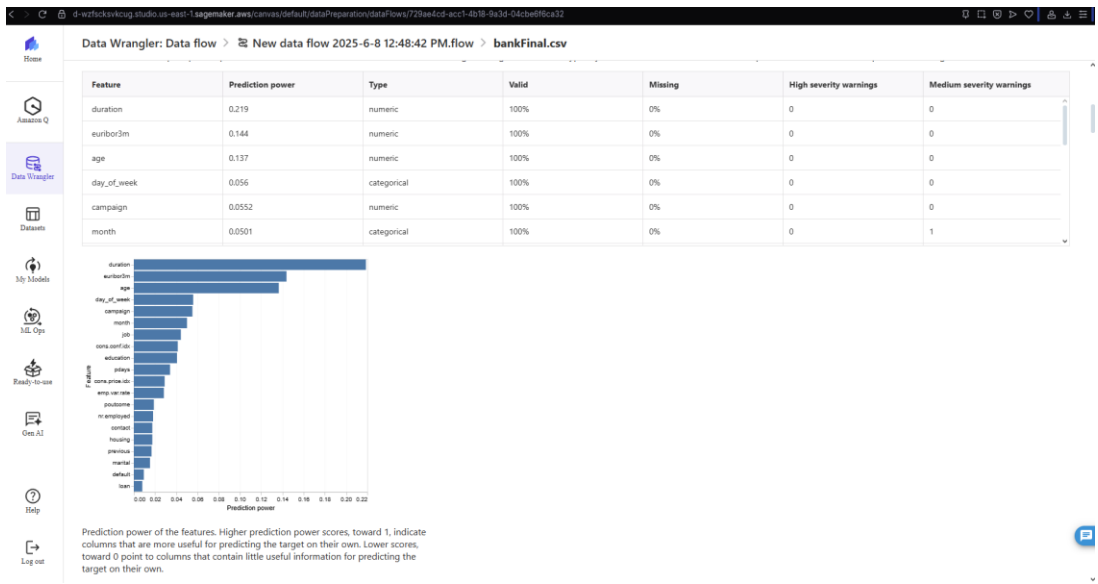
Key	Value
Number of features	21
Number of rows	41188
Missing	0%
Valid	100%
Duplicate rows	0.0583%

Feature type	Count
numeric	10
categorical	9
text	0
datetime	0
binary	1
unknown	0

High Priority Warnings

No high severity warnings were detected in the data.

y and insight report finished generating. X



Home

Amazon Q

Data Wrangler

Datasets

3rd Models

MLOps

Ready-to-use

Gen AI

Help

Log out

Data Wrangler: Data flow

New data flow 2025-6-8 12:48:42 PM.flow > bankFinal.csv

Canvas detected that 0.0583% of the data are duplicate. Some data sources could include valid duplicates. Other data sources could have duplicates that point to problems in data collection. Duplicate samples that result from faulty data collection could interfere with machine learning processes that rely on splitting the data into independent training and validation folds.

The following are examples of processes that can suffer from duplicated samples:

- Quick model analysis
- Prediction power estimation
- Hyper-parameters optimization

The most common duplicate rows are presented below. The number of occurrences of a row is given in left most column named Duplicate count. You can remove duplicate samples from the dataset using the Drop duplicates transform under Manage rows.

Note: features of type vector are ignored in duplicate rows detection

Duplicate count	age	job	marital	education	default	housing
2	24	services	single	highschool	no	yes
2	27	technician	single	professional.course	no	no
2	32	technician	single	professional.course	no	yes
2	35	admin.	married	university.degree	no	yes
2	36	retired	married	unknown	no	no
2	39	admin.	married	university.degree	no	no

Duplicate rows Medium

We found that 0.0583% of the data are duplicate. Some data sources could include valid duplicates and in other cases these duplicates could point to problems in data collection. Duplicate samples resulting from faulty data collection, could derail machine learning processes that rely on splitting to independent training and validation folds. For example quick model scores, prediction power estimation and automatic hyper parameter tuning. Duplicate samples could be removed from the dataset using the Drop duplicates transform under Manage rows.

Anomalous samples

Canvas detects anomalous samples using the Isolation forest algorithm after basic preprocessing. The isolation forest associates an anomaly score to each sample (row) of the dataset.

Home

Amazon Q

Data Wrangler

Datasets

3rd Models

MLOps

Ready-to-use

Gen AI

Help

Log out

Data Wrangler: Data flow

New data flow 2025-6-8 12:48:42 PM.flow > bankFinal.csv

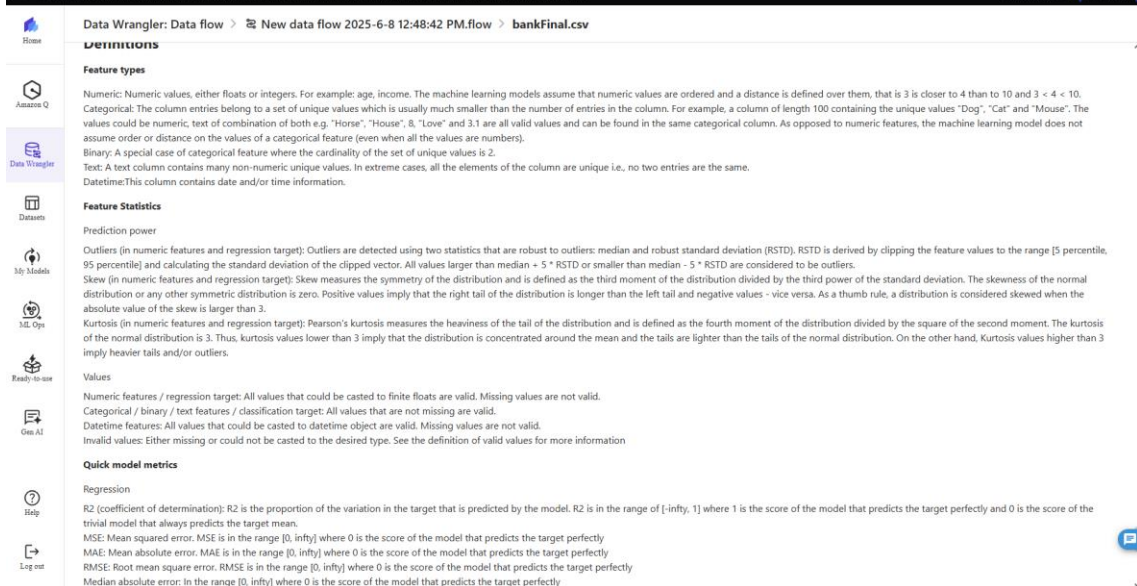
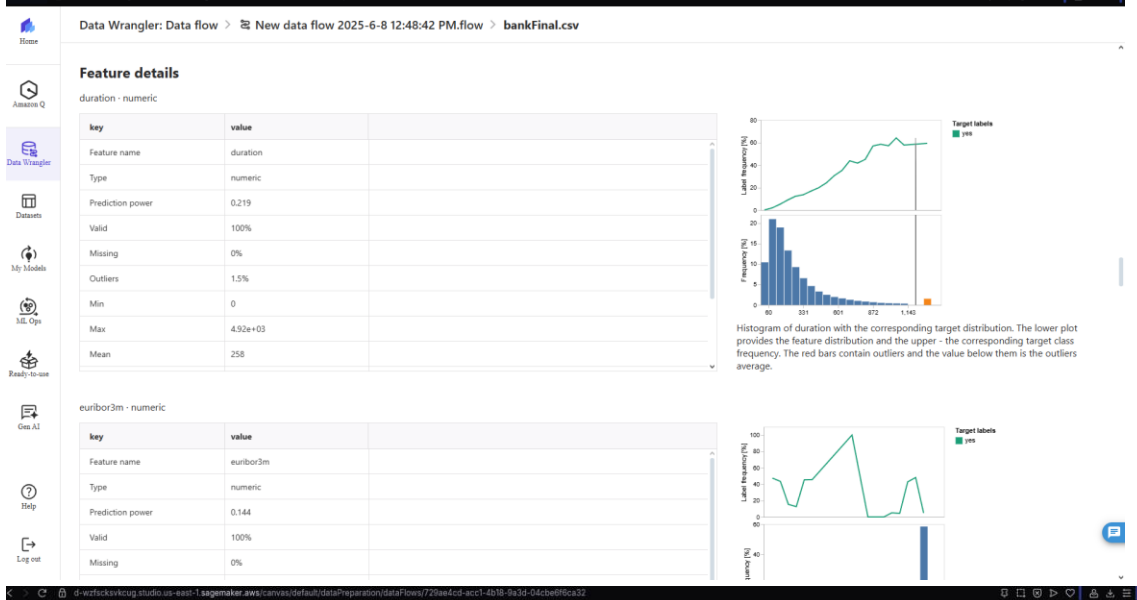
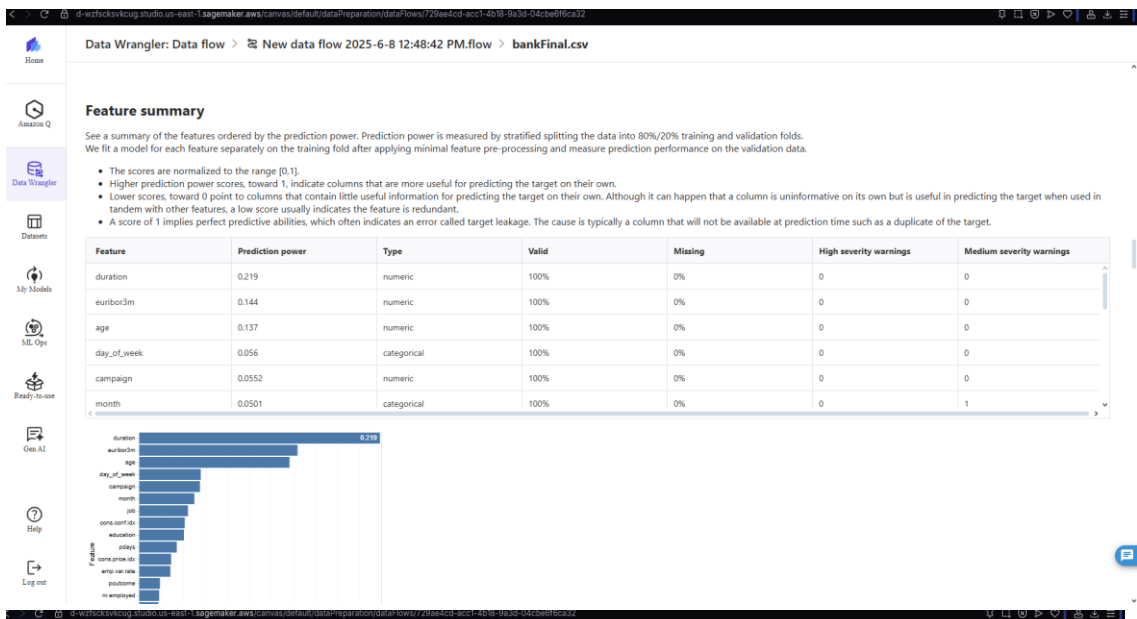
Quick model

Quick model provides an estimate of the expected predicted quality of a model that you train on your data.

The data is split into training and validation folds where Canvas uses 80% of the samples for training and 20% of the values for validation. For classification the sample is stratified split. For a stratified split, each data partition has the same ratio of labels. For classification problems, it's important to have the same ratio of labels between the training and classification folds. Canvas trains the XGBoost model with the default hyper-parameters. It applies early stopping on the validation data and performs minimal feature pre-processing.

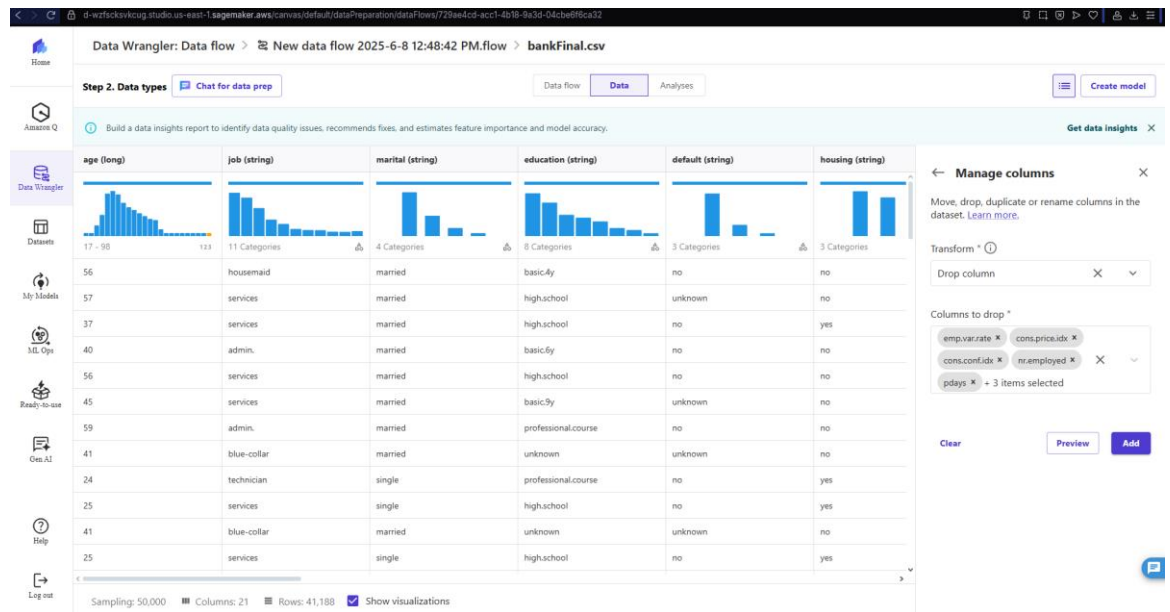
Metric	Validation scores	Train scores
Accuracy	0.916	0.932
Balanced accuracy	0.753	0.793
ROC-AUC	0.949	0.964
F1	0.594	0.669
Precision	0.656	0.736
Recall	0.543	0.614

class	precision	recall	f1-score	support
no	0.9432396251673361	0.96385089192086	0.9534506089309879	7310.0
yes	0.65625	0.5431034482758621	0.5943396256415094	928.0

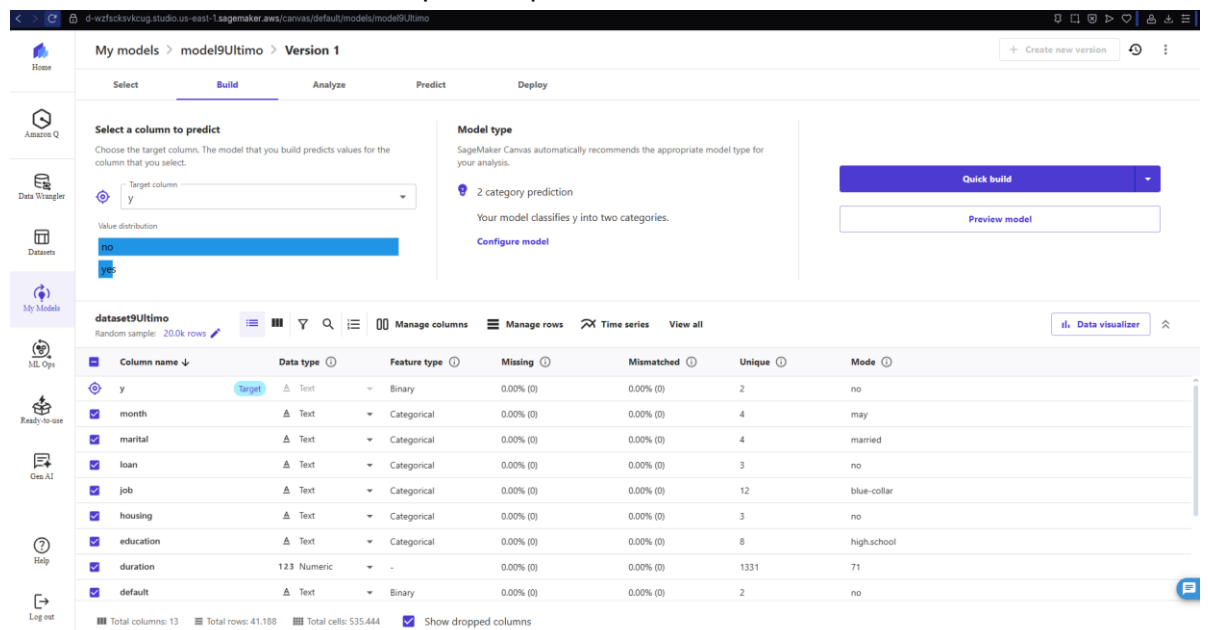


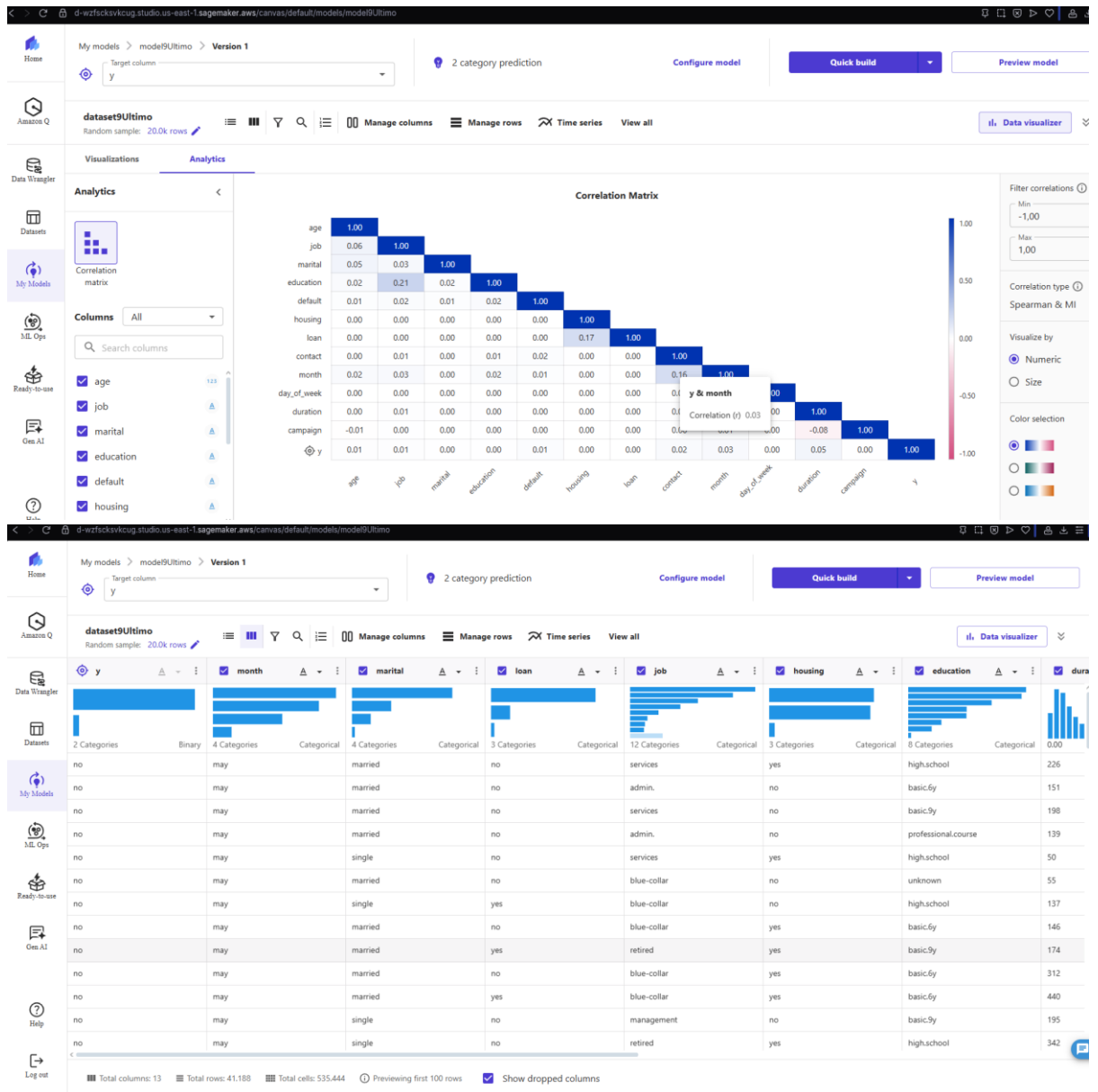
Vemos en el análisis que al final del todo nos da una descripción detallada del dataset. En las primeras capturas lo más reseñable es el porcentaje bajo de datos duplicados. Podemos seguir con una imagen de una estimación que daría un quick model. Seguimos con una captura de las características ordenadas por poder de predicción. Tenemos los datos normalizados de 0 a 1 para un análisis más fácil. En la última captura vemos una explicación de las cosas más importantes relacionadas con el EDA.

5. Aplicación del transform al dataset.



Le eliminamos las columnas que no aportan nada al entrenamiento del modelo.

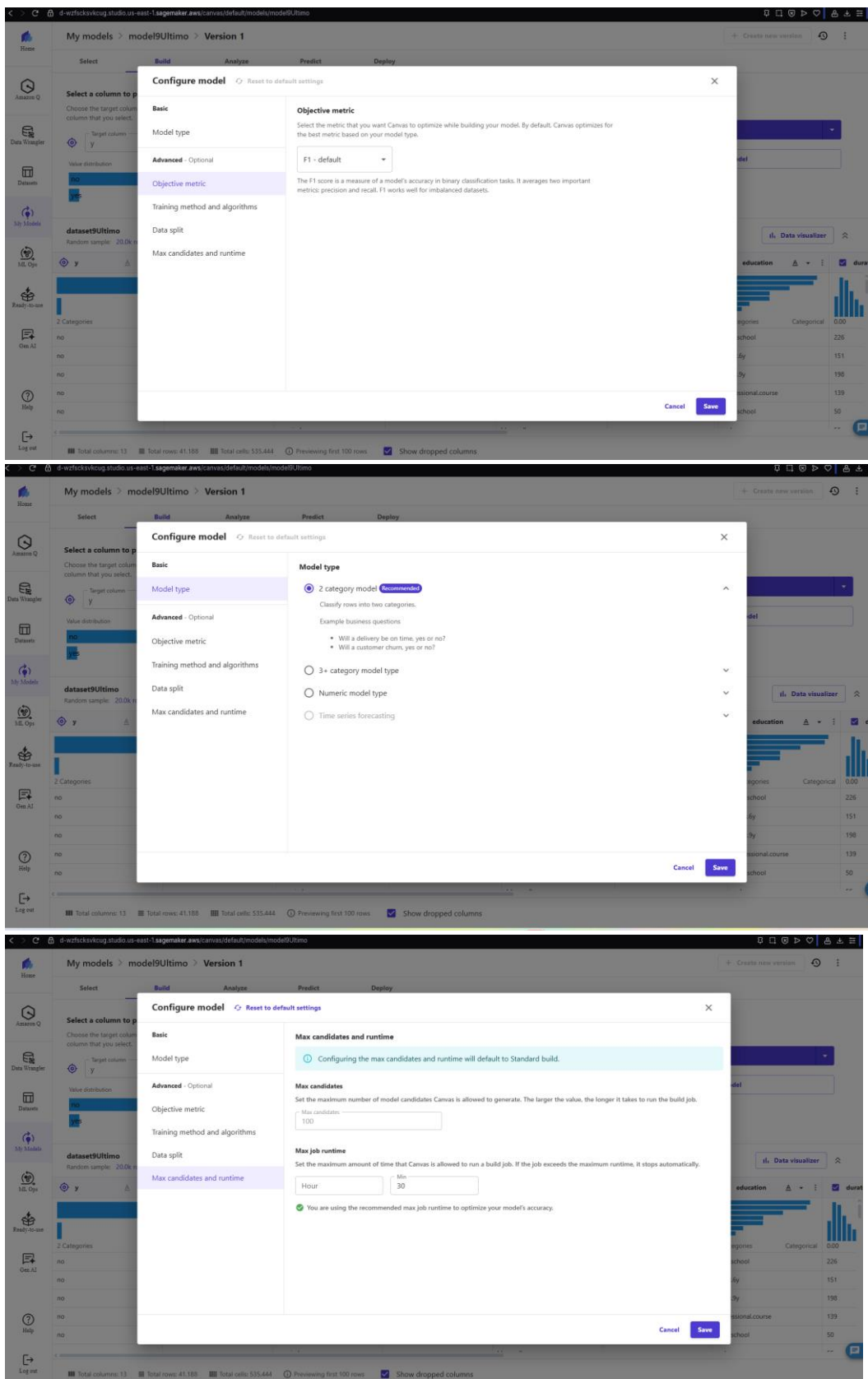




En la pestaña de build del modelo vemos como ha quedado el dataset y algunas graficas relevantes de este.

6. Entrenamiento del modelo

Ahora cambiaremos en una sola ventana al darle a la configuración del modelo.



En la casilla min del apartado de max job runtime cambiamos a 30 lo que hará automáticamente que cambie de quick build a standard build. Esto hará que el entrenamiento del modelo sea más largo.

Home

Assess Q

Data Visualizer

Datasets

My Models

ML Ops

Ready-to-use

Gen AI

Help

Log out

My models > model9Ultimo > Version 1

Create new version

Select

Build

Analyze

Predict

Deploy

Preprocessing your dataset. This can take a few minutes. You can navigate away from this page. This won't interrupt the process. [Learn more](#)

Select a column to predict

Choose the target column. The model that you build predicts values for the column that you select.

Target column

y

Value distribution

no

yes

Model type

SageMaker Canvas automatically recommends the appropriate model type for your analysis.

2 category prediction

Your model classifies y into two categories.

Configure model

Validating your data...

Validating your data...

dataset9Ultimo

Random sample: 20.0k rows

Data visualizer

y	month	marital	loan	job	housing	education	duration
no	may	married	no	services	yes	highschool	226
no	may	married	no	admin.	no	basic.5y	151
no	may	married	no	services	no	basic.3y	198
no	may	married	no	admin.	no	professional.course	139

Total columns: 13Total rows: 41,188Total cells: 535,444Previewing first 100 rowsShow dropped columns

Home

Assess Q

Data Visualizer

Datasets

My Models

ML Ops

Ready-to-use

Gen AI

Help

Log out

My models

Search models

New model

Grid

List

Filter by problem type: 2 category prediction

Last viewed

Ready

model9Ultimo

Versions1

Targety

Problem type2 category prediction

Updated2025-6-8 1:47:06 PM

View

Resources

Tutorials

Documentation

What's new

Home

Assess Q

Data Visualizer

Datasets

My Models

ML Ops

Ready-to-use

Gen AI

Help

Log out

My models > model9Ultimo > Version 1

Create new version

Select

Build

Analyze

Predict

Deploy

Model status

Standard build

Accuracy89.597%

F10.603

The model predicts the correct Y 89.597% of the time.

Predict

Deploy

Overview

Scoring

Advanced metrics

Model leaderboard

Column impact

Search columns...

1duration55.121%

2contact10.972%

3month9.451%

4education4.504%

5campaign4.344%

6default4.026%

7job3.839%

Impact of duration on prediction of y

no

Impact on prediction

yes

duration

dataset9Ultimo

Total columns: 13

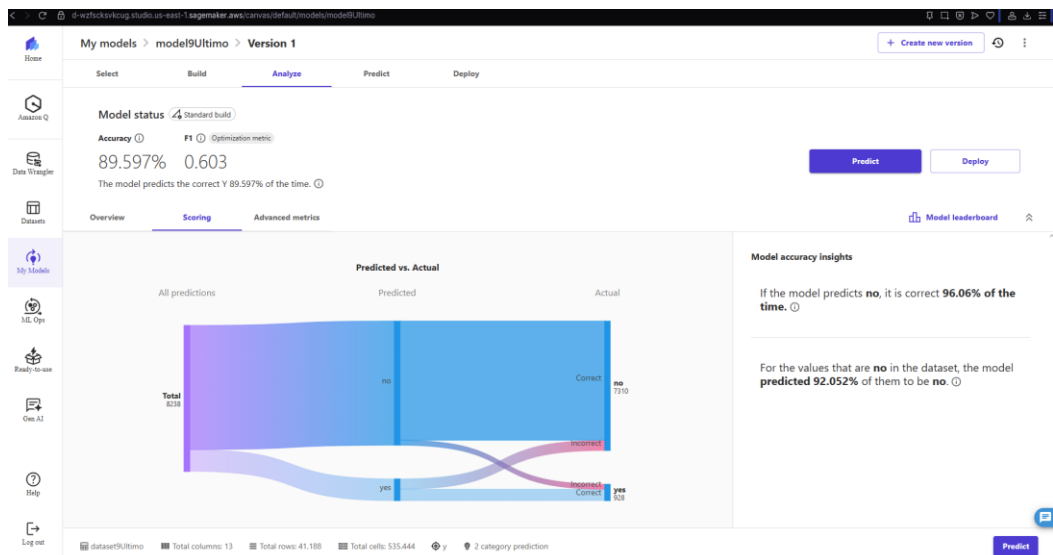
Total rows: 41,188

Total cells: 535,444

y

2 category prediction

Predict



My models > model9Ultimo > Version 1

Model status Standard build

Accuracy F1 Optimization metric

89.597% 0.603

The model predicts the correct Y 89.597% of the time.

Predict Deploy

Overview Scoring **Advanced metrics**

Model leaderboard

Positive Class	F1	Accuracy	Precision	Recall	AUC-ROC
yes no	60.342%	89.597%	52.879%	70.259%	0.924

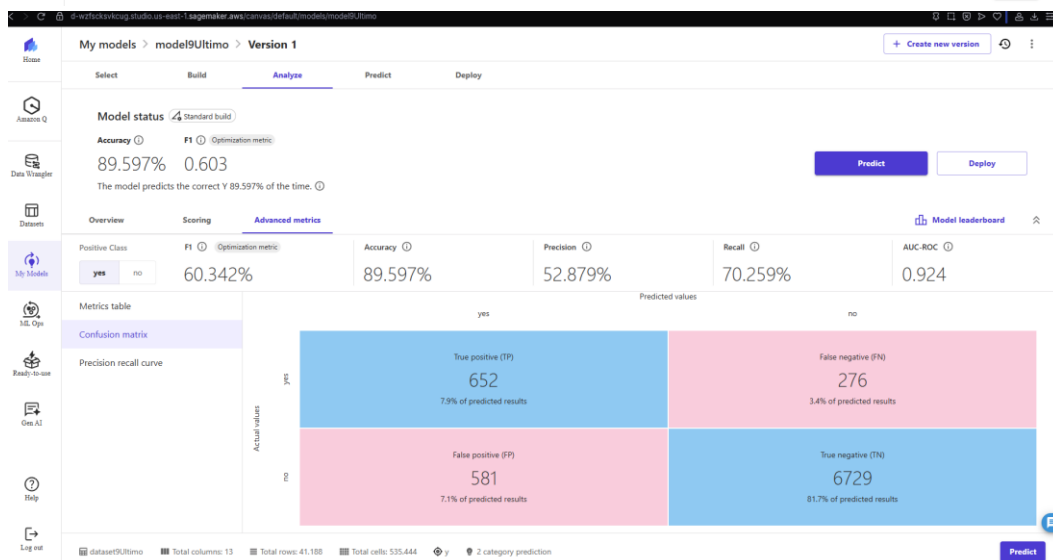
Metrics table

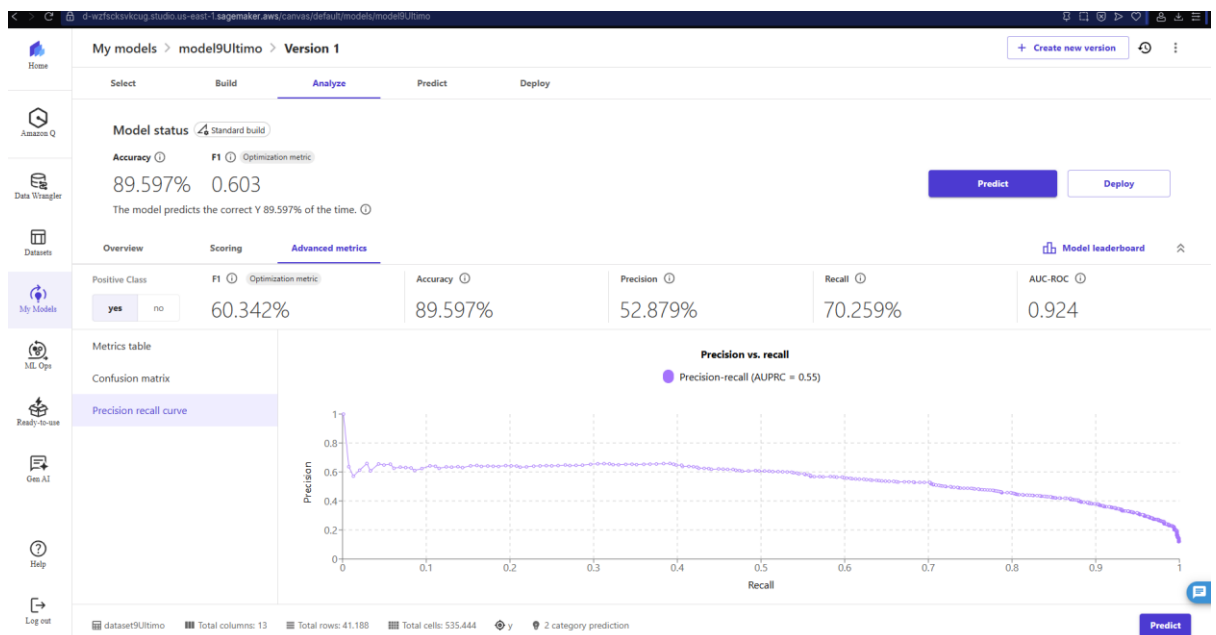
Confusion matrix

Precision recall curve

Metric name	Value
precision	0.529
recall	0.703
accuracy	0.896
f1	0.603
auc	0.924

dataset9Ultimo Total columns: 13 Total rows: 41,188 Total cells: 535,444 y 2 category prediction





A continuación, vemos la captura del proceso de entrenamiento del modelo y el modelo ya listo. Después podemos observar las métricas más importantes representadas en gráficas y múltiples apartados.

Matriz de confusión, curva de precisión del recall, AUC-ROC , accuracy , f1....

Model leaderboard									
Search leaderboard									
Model name	F1	Optimization	Accuracy	AUC	Balanced Accuracy	Precision	Recall	Log Loss	Inference latency (seconds)
FULL-13383046801754Canvas1749381630834	60.342%	Optimized model	89.597%	0.924	81.155%	52.879%	70.259%	0.232	0.167
FULL-14383046801754Canvas1749381630834	54.454%		87.958%	0.902	77.457%	47.440%	63.901%	0.267	0.266
FULL-18383046801754Canvas1749381630834	60.157%		89.548%	0.924	81.034%	52.717%	70.043%	0.234	0.351
FULL-17383046801754Canvas1749381630834	58.591%		87.800%	0.920	82.918%	47.432%	76.616%	0.265	0.184
FULL-16383046801754Canvas1749381630834	58.296%		86.975%	0.923	84.288%	45.593%	80.819%	0.289	0.132
FULL-15383046801754Canvas1749381630834	58.296%		86.975%	0.923	84.288%	45.593%	80.819%	0.289	0.129
FULL-14383046801754Canvas1749381630834	60.157%		89.548%	0.924	81.034%	52.717%	70.043%	0.234	0.381
FULL-12383046801754Canvas1749381630834	58.591%		87.800%	0.920	82.918%	47.432%	76.616%	0.265	0.175
FULL-11383046801754Canvas1749381630834	60.157%		89.548%	0.924	81.034%	52.717%	70.043%	0.234	0.361
FULL-110383046801754Canvas1749381630834	59.015%		90.204%	0.920	78.158%	55.812%	62.608%	0.250	0.177

Mostramos la tabla de modelos y sus resultados.

7. Pasamos ya al apartado del inference con el dataset descargado de la página de amazon.

Select dataset for predictions

To make predictions on a dataset, select it or import it. The dataset that you select must have the same number of feature columns as the training dataset. [?](#) [+ Create dataset](#)

Search datasets in Canvas

	Name		Columns	Rows	Cells	Created	Status
<input checked="" type="radio"/>	bank-inferenceFinal	VI	12	100	1200	06/08/2025 1:57 PM	Ready
<input type="radio"/>	dataset9Ultimo	VI	13	41.188	535.444	06/08/2025 1:13 PM	Ready
<input type="radio"/>	datasetFinal	VI	21	41.188	864.948	06/08/2025 12:46 PM	Ready
<input type="radio"/>	datasetPrediction	VI	22	100	2200	06/05/2025 4:14 PM	Ready
<input type="radio"/>	bank-inference-modificado	VI	20	100	2000	06/05/2025 4:09 PM	Ready
<input type="radio"/>	bank	VI	21	41.188	864.948	06/04/2025 11:35 PM	Ready
<input type="radio"/>	canvas-sample-product-descriptions.csv	VI	5	120	600	06/04/2025 10:10 PM	Incompatible ?
<input type="radio"/>	canvas-sample-shipping-logs.csv	VI	12	1000	12.000	06/04/2025 10:10 PM	Incompatible ?
<input type="radio"/>	canvas-sample-housing.csv	VI	10	1000	10.000	06/04/2025 10:10 PM	Incompatible ?
<input type="radio"/>	canvas-sample-databricks-dolly-15k.csv	VI	2	2000	4000	06/04/2025 10:10 PM	Incompatible ?
<input type="radio"/>	canvas-sample-loans-part-2.csv	VI	5	1000	5000	06/04/2025 10:10 PM	Incompatible ?
<input type="radio"/>	canvas-sample-retail-electronics-forecasting.csv	VI	6	40.500	243.000	06/04/2025 10:10 PM	Incompatible ?

[Close](#) [Generate predictions](#)

Vemos que el dataset de inference ya no nos muestra ninguna incompatibilidad como en el primer intento.

[Home](#) [Datasets](#) [finalModePredict9](#) [VI](#) [+ Create a data flow](#) [Update dataset](#) [+ Create a model](#) [Dataset details](#)

[Data](#) [Version history](#) [Auto update](#)

Previewing up to the first 100 rows of finalModePredict9

y	probability	age	job	marital	education	default	housing	loan
no	0.9997381568	0.6226415094	technician	married	high.school	no	no	yes
no	0.9999905229	0.5849056604	unknown	married	unknown	unknown	yes	no
no	0.9997378588	0.1698113208	blue-collar	married	basic.9y	no	no	no
no	0.9998580217	0.2264150943	admin.	married	high.school	no	no	no
no	0.9993964434	0.0566037736	housemaid	married	high.school	no	yes	no
no	0.9985105991	0.641509434	retired	married	professional.course	no	yes	yes
no	0.9999926686	0.4528301887	services	married	high.school	unknown	yes	no
no	0.9997512698	0.5094339623	admin.	divorced	university.degree	unknown	yes	no
no	0.9999595881	0	entrepreneur	married	university.degree	no	yes	yes
no	0.9997615218	0.2264150943	technician	divorced	professional.course	no	yes	yes
no	0.9997927547	0.1886792453	blue-collar	married	basic.9y	no	no	no
no	0.9996770024	0.3396226415	blue-collar	single	basic.4y	no	yes	no
no	0.9993048511	0.358490566	blue-collar	single	basic.9y	no	yes	no

[Dataset type: Tabular](#) [Total dataset cells \(columns x rows\): 1400 \(14 x 100\)](#) [Data source: Local](#)

Cargamos el dataset resultante del predict y vemos las 100 primeras columnas.

8. Realizar el deploy del modelo.

Model status Standard build

Accuracy 89.597% F1 0.603
The model predicts the correct Y 89.597% of the time.

Column impact

Column	Impact
1 duration	55.121%
2 contact	10.972%
3 month	9.451%
4 education	4.504%
5 campaign	4.344%
6 default	4.026%
7 job	3.839%

Impact of duration on prediction of y

Impact on prediction

duration

Create Deployment

Deploy your model to a SageMaker endpoint so that you can make predictions from outside of the Canvas application, test and monitor your model to proactively detect issues such as model drift.

Selected model version

model9Ultimo

Ready Created: 06-08-2025-2:00 PM

Deployment type

Real-time

Deployment name

finalDeploy

Instance type

ml.m5.12xlarge

Instance count

2

Cancel Deploy

Cambiamos el nombre al deploy.

Deployments

Filter by status: In service Failed Creating

Deployment name Status Deployment URL Created

Missing permissions to call AI services

Canvas can't create the endpoint because you don't have the necessary permissions. Contact your administrator to grant you access and try again. If you're an administrator or an individual user, go to the IAM console and check that the IAM role has the AmazonSageMakerFullAccess and AmazonSageMakerCanvasDirectDeployAccess policies attached. [Click here](#) to learn more about attaching AmazonSageMakerFullAccess and/or AmazonSageMakerCanvasDirectDeployAccess to an IAM role.

If you continue to see this issue, contact AWS support and provide the following code: 40ebe670-49a5-4f14-a4e5-d0e37a7ef6e2 to resolve the issue.

Close

Y vemos la captura que nos da canvas que como bien vimos en la clase del día 20, esta capada esta parte del modelado.