

# Proyecto 2

PEA DATA SCIENCE & IA

CARLOS ARREGUI Y JUAN PABLO DELZO

## Tabla de contenido

1.	Introducción.....	4
2.	Exploración y limpieza de los datos .....	5
3.	Base de datos.....	8
3.1.	Creación de la base de datos .....	8
3.2.	Carga de datos a MySQL.....	9
4.	Modelos de Machine Learning.....	14
4.1.	BigML.....	14
4.2.	Vertex.ai .....	17
4.3.	Watson.x.....	19
5.	Resultados de los modelos .....	24
5.1.	BigML.....	24
5.2.	Vertex.ai .....	27
5.3.	Watson.x.....	29
6.	PowerBI.....	31
7.	Conclusiones .....	32

## Tabla de figuras

Figura 1. Boxplot de los precios de alquiler por distrito .....	5
Figura 2. Código generado en Python para el ETL en el preprocesamiento de datos .....	6
Figura 3. Tabla en csv de los valores preprocesados sin la columna distrito transformada .....	7
Figura 4. Tabla en csv de los valores preprocesados con la columna de distritos transformados a variables binarias .....	7
Figura 5. Creación de la cuenta en Aiven .....	8
Figura 6. Credenciales asignadas para acceder a la base de datos en Aiven .....	9
Figura 7. Ajustes de conexión a la base de datos de Aiven con DBeaver.....	9
Figura 8. Carga de datos hacia la base de datos en Aiven.....	10
Figura 9. Selección de importación de datos mediante csv en DBeaver .....	10
Figura 10. Configuración de parámetros para una lectura correcta de los datos .....	11
Figura 11. Comprobación de lectura de datos .....	11
Figura 12. Configuración de parámetros para importación de datos .....	12
Figura 13. Paso de confirmación para la importación de los datos .....	12
Figura 14. Imagen de demostración de la carga correcta de los datos del csv a la base de datos MySQL .....	13
Figura 15. Creación de conexión de BigML a la base de datos MySQL de Aiven .....	15
Figura 16. Creación del dataset en BigML.....	15
Figura 17. Comprobación del dataset y asignación de la columna objetivo en BigML .....	16
Figura 18. Dataset creado con los modelos de ML generados en verde a su derecha .....	16
Figura 19. Carga de datos en Vertex.ai.....	17
Figura 20. Ajuste para el entrenamiento de nuevos modelos en Vertex.ai .....	17
Figura 21. Secuencia de pasos para ajustar los parámetros de tiempo de ejecución y la métrica en la columna objetivo .....	19
Figura 22. Selección del objetivo de optimización .....	19
Figura 23. Creación de instancia en Watson.x .....	20
Figura 24. Creación de proyecto en Watson.x.....	20
Figura 25. Importación de archivos en Watson.x.....	21
Figura 26. Elección de crear un modelo de aprendizaje automático con los datos cargados .....	21
Figura 27. Asociación del servicio de WatsonMachineLearning en Watson.x .....	22
Figura 28. Creación del modelo en Watson.x.....	22
Figura 29. Asignación del archivo csv en Watson.x.....	23
Figura 30. Resultados obtenidos por el modelo Random Forest .....	24
Figura 31. Resultados obtenidos por el modelo Ensemble .....	24
Figura 32. Resultados obtenidos por el modelo Linear Regression .....	25
Figura 33. Resultados obtenidos por el modelo DeepNet.....	25
Figura 34. Comparación de resultados obtenidos mediante optimización por RMSE y MAE .....	27
Figura 35. Comparación de la diferencia en importancias de cada atributo según el ajuste .....	27
Figura 36. Mapa de relaciones obtenida por Watson.x .....	29
Figura 37. Tabla de interconexiones del modelo de Watson.x.....	29

Figura 38. Gráfico de métricas del modelo de Watson.x.....	29
Figura 39. Minería de datos de la tabla procesada con WindSurf .....	31

# 1. Introducción

El principal objetivo es que una vez ya desarrollado el primer proyecto de forma tradicional y gracias al feedback recibido, el siguiente paso es utilizar herramientas en IA y otras herramientas ya conocidas que no requieran programación.

El proyecto consta de 4 apartados:

- El primero está desarrollado con *WindSurf*, con el fin de crear código con el mínimo esfuerzo en Python (solo se tuvo que corregir o añadir puntualmente).
- El segundo apartado está desarrollado con *MySQL*, *DBeaver* y *Aiven*.
- El tercer apartado está desarrollado con *BigML*, *Vertex AI* e *IBM Watsonx*.
- El último apartado está desarrollado en PowerBI con Python y PowerQuery en el backend, y *DAX* en el frontend.

Preprocesamiento de datos:

- Análisis exploratorio de datos.
- Llenado o eliminación de datos nulos.
- Eliminación de atípicos.
- Escalado de datos.
- Transformación de variables.

Carga de datos:

- Creación de base de datos sin herramientas de programación

Procesamiento de datos y métodos supervisados:

- Herramientas ML en IA

Explotación de datos:

- Visualización interactiva de los datos.
- Visualización de modelos.

## 2. Exploración y limpieza de los datos

Mediante *WindSurf* se realizaron la siguiente serie de peticiones:

Con Python, Pandas y Matplotlib con fondo fivethirtyeight, con el archivo adjunto en input:

1. Elimina la columna Unnamed: 0.
2. Analiza qué columnas tienen valores nulos.
3. Si la columna Neighborhood no tiene valores nulos, reemplaza los valores nulos en función de la moda en su columna de cada valor de Neighborhood.
4. En las columnas Rooms, Bathrom y Square\_meters, reemplaza a valor int redondeando los valores quedándote con la parte entera.
5. Chequea si  $\text{Square\_meters\_price} * \text{Square\_meters}$  difieren con price con una tolerancia de 1000 euros.
6. Elimina la columna Square\_meters\_price. Detecta si hay outliers con el método de IQR con k igual a 3 en Price y apártalos del dataframe.
7. Analiza si hay una relación lineal entre Price promedio y Neighborhood. Si tal valor es menor de 0.1, con get\_dummies de Pandas, transforma la columna Neighborhood con drop\_first True y dtype a int.
8. Guarda dos archivos con extensión csv, uno con la tabla transformada con get\_dummies y el otro sin transformarlo.

El código, gráfico y tablas generados tardaron menos de dos minutos, fueron los siguientes:



Figura 1. Boxplot de llos precios de alquiler por distrito

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np
5 from scipy import stats
6
7 # Set the style for plots
8 plt.style.use('fivethirtyeight')
9
10 # Read the data
11 df = pd.read_csv('input/DATA_Barcelona_Fotocasa_HousingPrices_Augmented.csv')
12 print("\n")
13 print(df.dtypes)
14
15 # Remove Unnamed: 0 column
16 df = df.drop('Unnamed: 0', axis=1)
17
18 # Check null values
19 print("\nNull values in each column:")
20 print(df.isnull().sum())
21
22 # If neighborhood has no null values, fill nulls based on mode per neighborhood
23 if df['neighborhood'].isnull().sum() == 0:
24     for column in df.columns:
25         if column != 'neighborhood' and df[column].isnull().any():
26             df[column] = df.groupby('neighborhood')[column].transform(lambda x: x.fillna(x.mode().iloc[0] if not x.mode().empty else x.mean()))
27
28 # Convert rooms, bathroom and square_meters to int
29 df['rooms'] = df['rooms'].round().astype(int)
30 df['bathroom'] = df['bathroom'].round().astype(int)
31 df['square_meters'] = df['square_meters'].round().astype(int)
32
33 # Check price consistency
34 price_diff = abs(df['square_meters_price'] * df['square_meters'] - df['price'])
35 inconsistent_prices = price_diff > 1000
36 print("\nNumber of inconsistent prices:", inconsistent_prices.sum())
37
38 # Remove square_meters_price column
39 df = df.drop('square_meters_price', axis=1)
40
41 # Detect outliers using IQR method
42 Q1 = df['price'].quantile(0.25)
43 Q3 = df['price'].quantile(0.75)
44 IQR = Q3 - Q1
45 k = 2.5
46 lower_bound = Q1 - k * IQR
47 upper_bound = Q3 + k * IQR
48
49 # Remove outliers
50 df_no_outliers = df[(df['price'] >= lower_bound) & (df['price'] <= upper_bound)]
51 print(f"\nRemoved {len(df) - len(df_no_outliers)} outliers")
52
53 # Create boxplot
54 plt.figure(figsize=(15, 8))
55 sns.boxplot(data=df_no_outliers, x='neighborhood', y='price')
56 plt.xticks(rotation=45, ha='right')
57 plt.title('Price Distribution by Neighborhood')
58 plt.tight_layout()
59 plt.savefig('output/price_neighborhood_boxplot.png')
60 plt.close()
61
62 # Calculate correlation using one-way ANOVA
63 neighborhoods = df_no_outliers['neighborhood'].unique()
64 price_groups = [df_no_outliers[df_no_outliers['neighborhood'] == n]['price'] for n in neighborhoods]
65 f_statistic, p_value = stats.f_oneway(*price_groups)
66 correlation = np.sqrt(f_statistic / (f_statistic + len(df_no_outliers) - len(neighborhoods)))
67 print(f"\nCorrelation strength (based on ANOVA): {correlation:.4f}")
68
69 # If correlation is less than 0.25, create dummy variables
70 if correlation < 0.25:
71     df_dummy = pd.get_dummies(df_no_outliers, columns=['neighborhood'], drop_first=True, dtype=int)
72     # Save both versions
73     df_dummy.to_csv('output/housing_prices_dummy.csv', index=False)
74     df_no_outliers.to_csv('output/housing_prices_processed.csv', index=False)
75 else:
76     df_no_outliers.to_csv('output/housing_prices_processed.csv', index=False)
77     print("\nCorrelation is >= 0.25, dummy variables were not created")

```

Figura 2. Código generado en Python para el ETL en el preprocesamiento de datos

housing\_prices\_processed.csv

output > housing\_prices\_processed.csv > data

	price	rooms	bathroom	Lift	terrace	square_meters	real_state	neighborhood
1	750	3	1	True	False	60	flat	Horta- Guinardo
2	770	2	1	True	False	59	flat	Sant Andreu
3	1300	1	1	True	True	30	flat	Gràcia
4	2800	1	1	True	True	70	flat	Ciutat Vella
5	720	2	1	True	False	44	flat	Sant Andreu
6	1100	3	2	False	False	118	attic	Horta- Guinardo
7	1350	3	3	True	False	84	flat	Sarria-Sant Gervasi
8	900	2	1	True	False	60	flat	Sarria-Sant Gervasi
9	1165	4	2	True	False	106	flat	Horta- Guinardo
10	1050	2	1	True	False	70	flat	Les Corts
11	2500	4	2	True	False	161	flat	Sarria-Sant Gervasi
12	1335	2	1	True	False	55	flat	Sant Martí
13	2200	4	2	True	True	175	flat	Eixample
14	2500	3	2	True	True	145	flat	Sarria-Sant Gervasi
15	980	2	2	True	False	82	flat	Eixample
16								

Figura 3. Tabla en csv de los valores preprocesados sin la columna distrito transformada

housing\_prices\_dummy.csv

output > housing\_prices\_dummy.csv > data

	price	rooms	bathroom	Lift	terrace	square_meters	real_state	neighborhood_Eixample	neighborhood_Gràcia	neighborhood_Horta- Guinardo	neighborhood_Les Corts	neighborhood_Nou Barris	neighborhood_Sant Andreu	neighborhood_Sant Martí	neighborhood_Sants-Montjuic	neighborhood_Sarria-Sant Gervasi
1	750	3	1	True	False	60	flat	0	0	1	0	0	0	0	0	0
2	770	2	1	True	False	59	flat	0	0	0	0	0	1	0	0	0
3	1300	1	1	True	True	30	flat	0	1	0	0	0	0	0	0	0
4	2800	1	1	True	True	70	flat	0	0	0	0	0	0	0	0	0
5	720	2	1	True	False	44	flat	0	0	0	0	0	1	0	0	0
6	1100	3	2	False	False	118	attic	0	0	1	0	0	0	0	0	0
7	1350	3	3	True	False	84	flat	0	0	0	0	0	0	0	0	1
8	900	2	1	True	False	60	flat	0	0	0	0	0	0	0	0	1
9	1165	4	2	True	False	106	flat	0	0	1	0	0	0	0	0	0
10	1050	2	1	True	False	70	flat	0	0	0	1	0	0	0	0	0
11	2500	4	2	True	False	161	flat	0	0	0	0	0	0	0	0	1
12	1335	2	1	True	False	55	flat	0	0	0	0	0	0	0	1	0
13	2200	4	2	True	True	175	flat	1	0	0	0	0	0	0	0	0
14																

Figura 4. Tabla en csv de los valores preprocesados con la columna de distritos transformados a variables binarias



## 3. Base de datos

### 3.1. Creación de la base de datos

Se crea una base de datos *MySQL* en la plataforma *Aiven*. Es una plataforma de gestión de bases de datos y servicios en la nube que permite a las empresas desplegar, operar y escalar fácilmente tecnologías de código abierto como *PostgreSQL*, *Kafka*, *OpenSearch*, *Redis*, *MySQL*, entre otras. A continuación, se detallan los pasos para la creación de la base de datos:

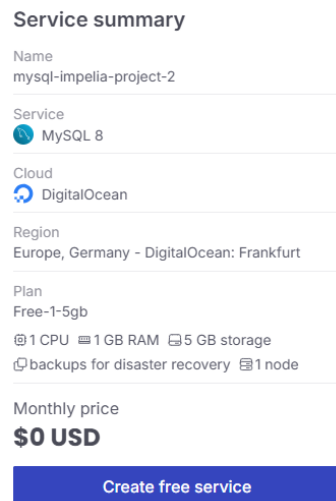


Figura 5. Creación de la cuenta en Aiven

Para la creación de la base de datos, se ha escogido un plan gratuito en el que se incluyen un servidor en Ámsterdam de 1 CPU, con 1 GB de memoria RAM y 5 GB de almacenamiento de datos. Los datos de conexión a la base de datos son los siguientes:

- Nombre defaultdb
- Host mysql-impelia-project-2-impelia-project-2.h.aivencloud.com
- Puerto 13961
- Usuario avnadmin
- Contraseña AVNS\_VjEvOz2Yre8wHS8cDiw

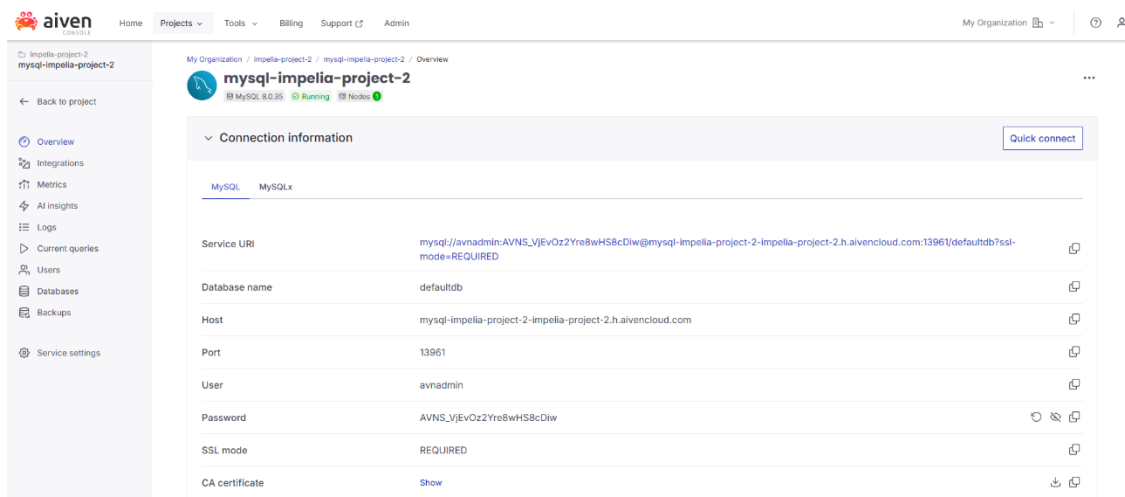


Figura 6. Credenciales asignadas para acceder a la base de datos en Aiven

### 3.2. Carga de datos a MySQL

Para cargar los datos en la base de datos *MySQL* de *Aiven*, se utilizó la plataforma *DBeaver*, un software gratuito especializado en la gestión de base de datos. Para poder subir los datos, primero nos conectamos a ella introduciendo la información mencionada en el apartado anterior. Además, también se autorizó usar el protocolo SSL y se cargó el certificado CA, descargado directamente desde *Aiven*.

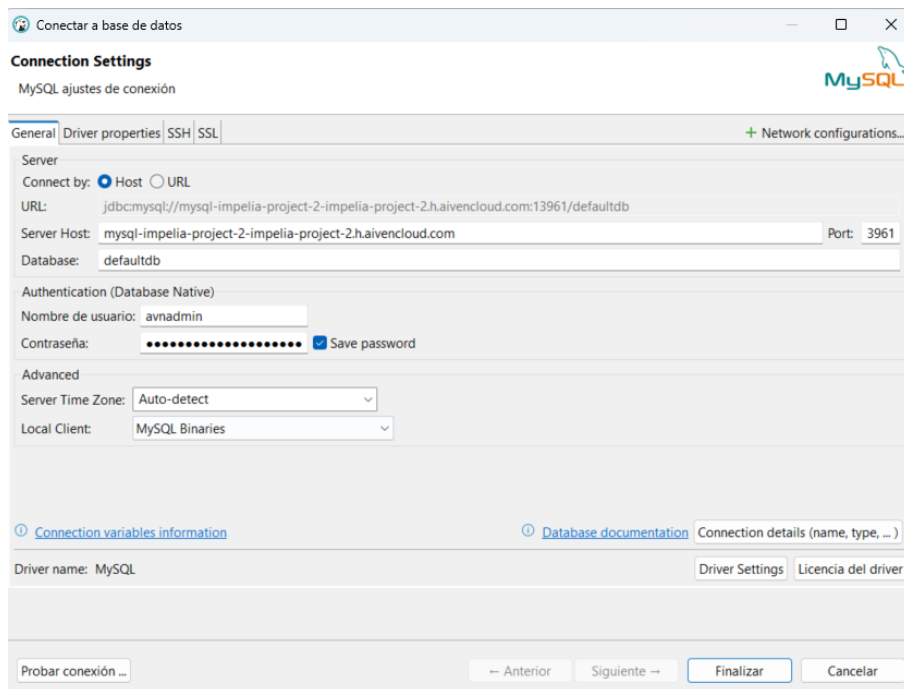


Figura 7. Ajustes de conexión a la base de datos de Aiven con DBeaver

Una vez terminada la conexión a la base de datos, se cargaron los datos del .csv posterior al ETL. Se adjuntan capturas de pantalla del proceso:

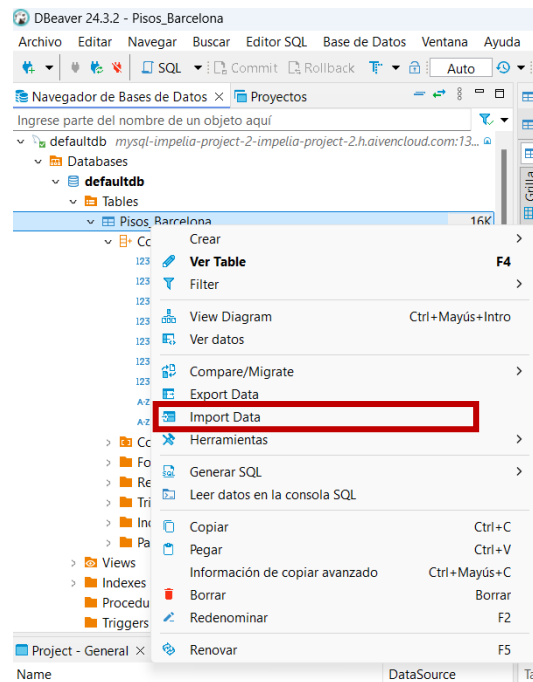


Figura 8. Carga de datos hacia la base de datos en Aiven

En primer lugar, se hace clic con el botón derecho sobre la tabla creada, y se va a *Import Data*. Acto seguido en la pantalla que se abre, se selecciona CSV como fuente de la importación.

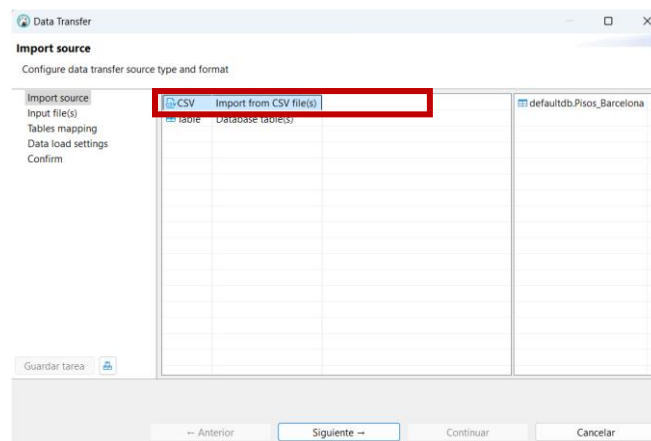


Figura 9. Selección de importación de datos mediante csv en DBeaver

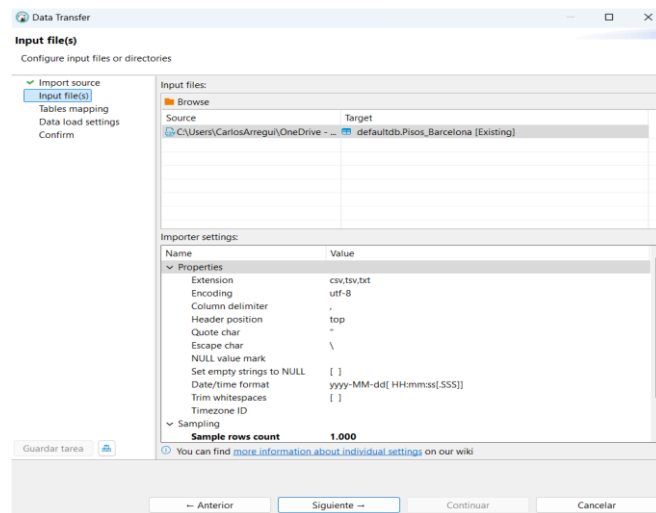
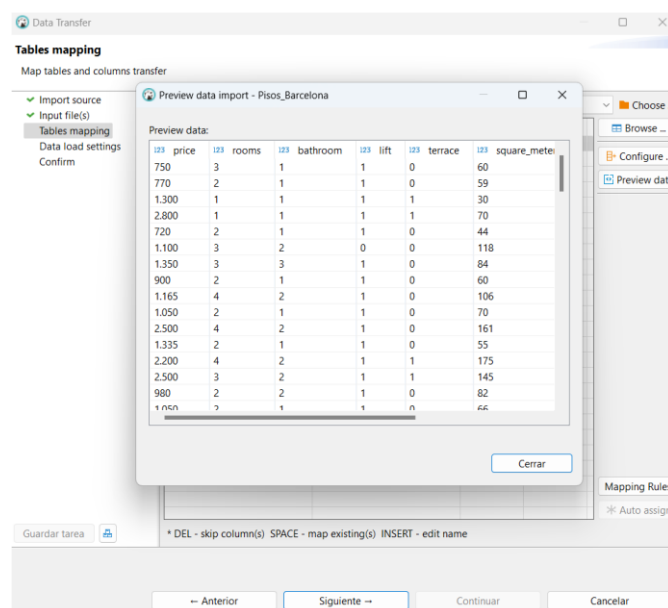


Figura 10. Configuración de parámetros para una lectura correcta de los datos

Se selecciona el archivo de la carpeta local para importar el CSV y se mapean los datos para comprobar que los está leyendo bien y que las columnas tienen el nombre que necesitamos.



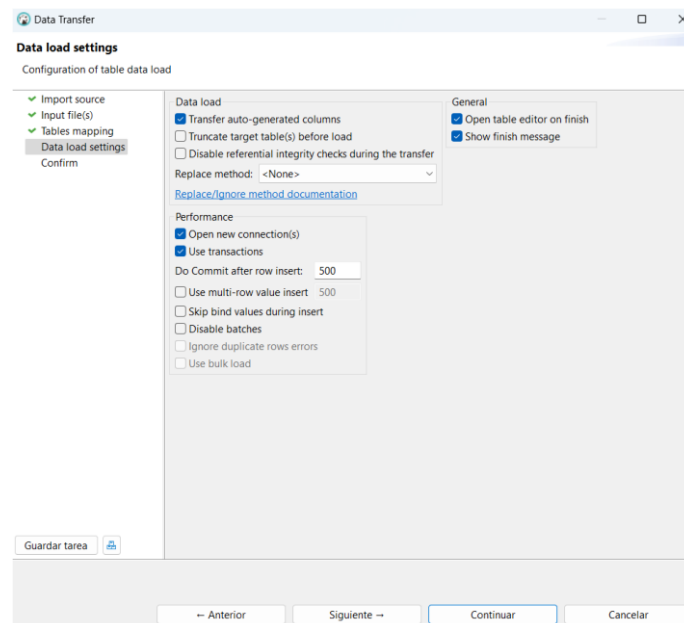


Figura 12. Configuración de parámetros para importación de datos

En la configuración de la carga de datos se deja la predeterminada, subiendo únicamente el *commit* cada 500 filas insertadas en vez de cada 100. Se clicca sobre *Continuar* y se confirman todos los pasos anteriores para empezar la carga.

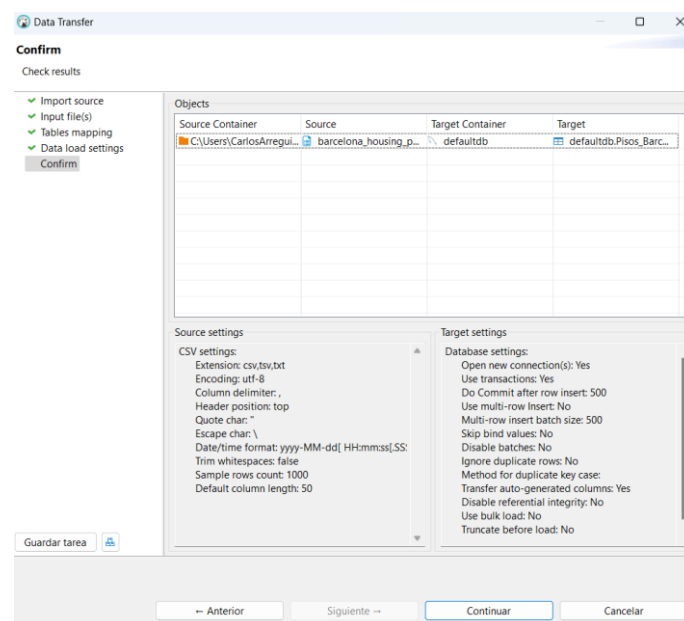


Figura 13. Paso de confirmación para la importación de los datos

Una vez terminada la carga de todas las filas del CSV, se puede comprobar que la carga ha sido realizada yendo a ver la tabla creada *Pisos\_Barcelona*. Tal y como puede verse en la imagen siguiente, se ha creado automáticamente una columna que contiene un

identificador para cada piso y luego, todas las columnas que contiene el CSV separadas por una coma.

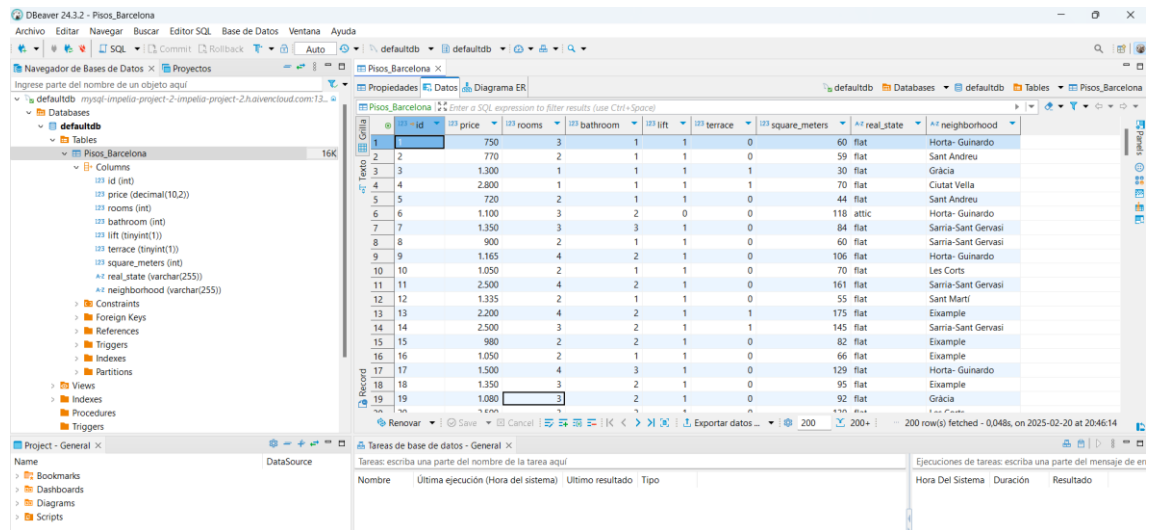


Figura 14. Imagen de demostración de la carga correcta de los datos del csv a la base de datos MySQL

## 4. Modelos de Machine Learning

En los últimos años, el auge del *Machine Learning* ha revolucionado múltiples sectores, incluyendo el mercado inmobiliario. La capacidad de analizar grandes volúmenes de datos y extraer patrones complejos permite mejorar la precisión en la estimación de precios de viviendas, optimizando tanto la toma de decisiones de compradores y vendedores como la estrategia de inversión de empresas del sector. En este contexto, el presente trabajo se centra en la predicción de precios de pisos en Barcelona mediante la implementación de modelos de *Machine Learning* en tres plataformas avanzadas: **BigML**, **Vertex AI** y **Watsonx**.

El objetivo principal de este estudio es construir modelos predictivos que permitan estimar el valor de un inmueble a partir de diversas características, como ubicación, superficie, número de habitaciones, antigüedad del edificio, entre otros factores relevantes. Para ello, se han utilizado diferentes técnicas de aprendizaje supervisado, explorando la efectividad de cada plataforma en términos de facilidad de uso, rendimiento del modelo y capacidad de escalabilidad.

BigML, con su enfoque centrado en la simplicidad y automatización, ha permitido una rápida implementación de modelos sin necesidad de una infraestructura compleja. Por otro lado, Vertex AI, la plataforma de Google ha proporcionado herramientas avanzadas para el entrenamiento y despliegue de modelos en la nube, facilitando la integración con grandes volúmenes de datos. Finalmente, Watsonx de IBM ha ofrecido capacidades de inteligencia artificial explicable, permitiendo analizar la interpretabilidad de los modelos y mejorar su transparencia en la toma de decisiones.

A lo largo del trabajo, se explicarán los procesos, paso a paso, que se han utilizado para crear estos modelos y se comparan los resultados obtenidos en cada plataforma, evaluando métricas de rendimiento como el error cuadrático medio (MSE), el coeficiente de determinación ( $R^2$ ) y la eficiencia en el manejo de datos. Además, se analizan las ventajas y desafíos que presenta cada solución en el desarrollo de modelos para el sector inmobiliario.

Este estudio no solo busca determinar cuál de estas herramientas proporciona las predicciones más precisas, sino también aportar una visión práctica sobre su aplicabilidad en proyectos del mundo real, contribuyendo al desarrollo de modelos de estimación de precios más eficientes y accesibles para el mercado inmobiliario de Barcelona.

### 4.1. BigML

BigML es una plataforma de inteligencia artificial basada en la nube que se especializa en la creación de modelos de aprendizaje automático de manera accesible y simplificada. Su diseño intuitivo permite a los usuarios, independientemente de su nivel técnico, construir y desplegar modelos predictivos con facilidad.

El proceso general es sencillo y la aplicación es muy intuitiva. Primero nos conectamos a la base de datos de Aiven introduciendo los valores de Host, Puerto, Usuario, Contraseña y Nombre de la base de datos. Seguidamente se crea un dataset usando la misma fuente de datos, y finalmente se entrenan los tipos de modelos que se consideren necesarios.

Create new connector

MySQL

Host: mysql-impelia-project-2-impelia-project-2.h.aivencloud.c Port: 13961

Username: avnadmin Password: .....

Database: defaultdb ☒ Use SSL ☐ Verify certificates

Cancel Create

Figura 15. Creación de conexión de BigML a la base de datos MySQL de Aiven

Para crear el dataset, solo es necesario clicar a la derecha de la fuente de datos y seleccionar *1-click dataset* para crearlo.

Sources

Type	Name	Size
CSV	Pisos_Barcelona	9min, 647.7 KB
ZIP	hot-dog-or-not.zip	13min, 2.2 MB
ZIP	firetruck.zip	13min, 1.5 MB
ZIP	grape-strawberry.zip	13min, 1.9 MB
CSV	Country Stats Mashup	13min, 12.0 KB
CSV	Fictional Wine Sales	13min, 51.9 KB
CSV	Titanic Survival	13min, 78.0 KB
CSV	US Car Accidents in 2011	13min, 685.5 KB
CSV	Premier League 2011-2012 Season	13min, 24.7 KB
CSV	Churn in the Telecom Industry	13min, 270.4 KB

1 to 10 of 14 sources

Figura 16. Creación del dataset en BigML

Una vez creado el dataset de los datos limpios de la base de datos de Aiven, se comprueba que se ha conectado correctamente y que se leen bien todos los datos. También se establece como columna objetivo el precio



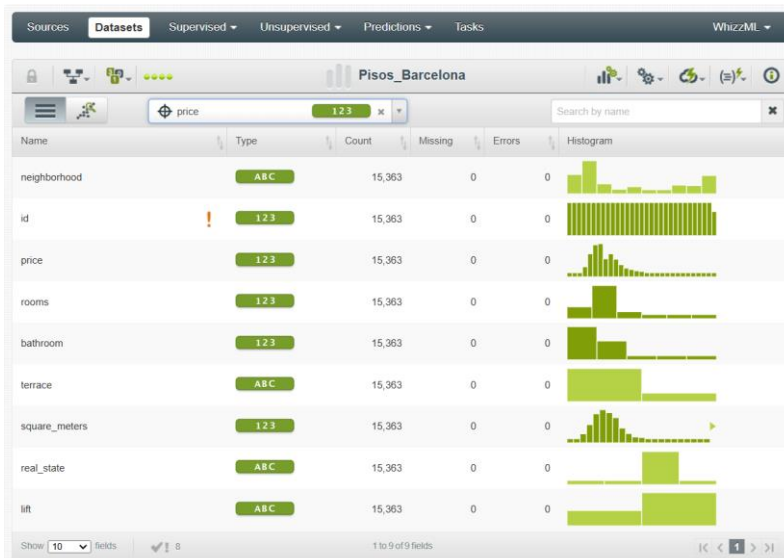


Figura 17. Comprobación del dataset y asignación de la columna objetivo en BigML

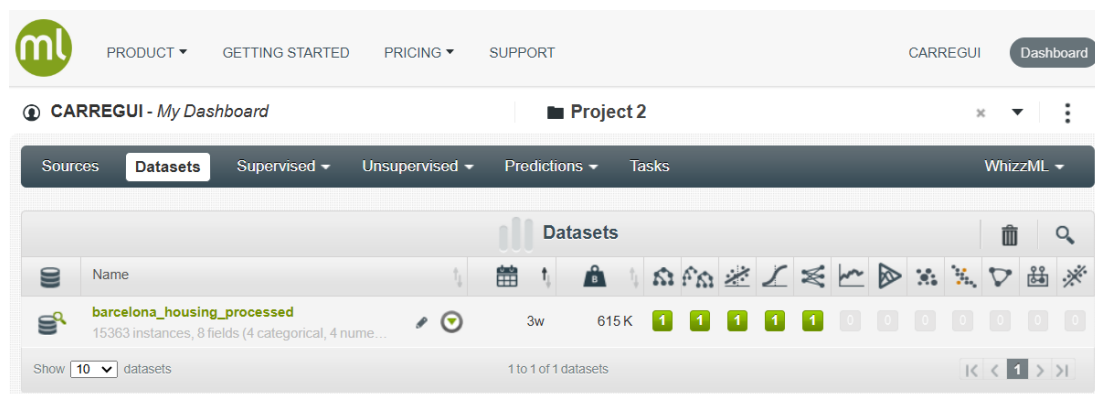


Figura 18. Dataset creado con los modelos de ML generados en verde a su derecha

Como se observa en la imagen anterior, una vez creado el dataset con el conjunto de datos, a la derecha se nos ofrecen todas las opciones de modelos disponibles (marcadas en verde). Para nuestro proyecto, se realiza la comparativa de varias opciones:

1. Random Forest
2. Ensemble
3. Linear Regression
4. DeepNet

## 4.2. Vertex.ai

Vertex AI es la plataforma de Google diseñada para simplificar el desarrollo, entrenamiento, implementación y monitoreo de modelos de inteligencia artificial (IA) y aprendizaje automático. Integra diversas herramientas y servicios en una sola interfaz, permitiendo a los desarrolladores y científicos de datos gestionar el ciclo de vida completo de los modelos de IA de manera eficiente.

El proceso para poder llegar a entrenar un modelo de ML a través de Vertex.ai empieza por cargar el conjunto de datos limpios. Para ello se crea un conjunto de datos y se adjunta el CSV. El nombre asignado para nuestro caso es *datos-impelia*.

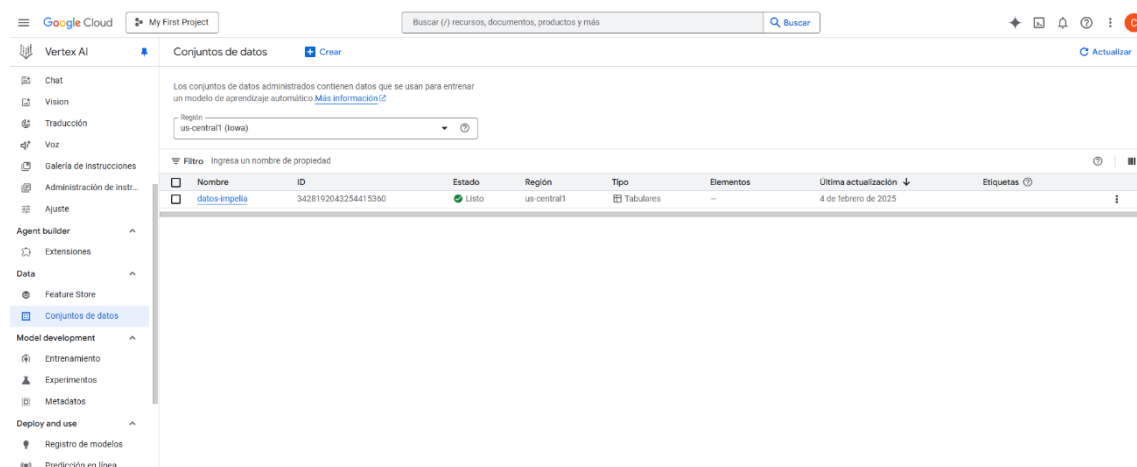


Figura 19. Carga de datos en Vertex.ai

Luego se clicla en el apartado *Entrenamiento* para entrenar un modelo nuevo.

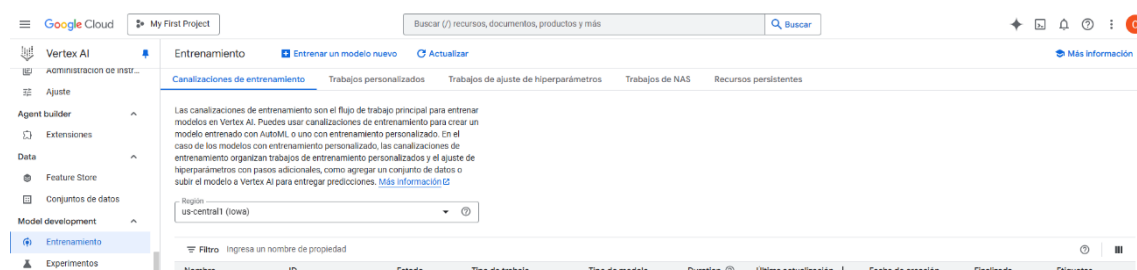


Figura 20. Ajuste para el entrenamiento de nuevos modelos en Vertex.ai

A continuación, se abre una pantalla que nos va guiando en todo el proceso de creación del modelo. En la primera pantalla se asigna el conjunto de datos que queremos utilizar el objetivo, en el cual marcaremos *Regresión*.

En la segunda pantalla se asigna el nombre del modelo (*regresion-proyecto-2-impelia*), una breve descripción del uso del modelo y la columna de los datos que se va a querer predecir, en nuestro caso *Price*. La tercera pantalla no la usaremos para este caso.

En la cuarta pantalla se realiza un análisis previo de los datos, pudiendo cambiar el tipo de datos de cada columna en caso de que no se haya especificado bien de forma automática, contar los valores distintos que hay en cada columna y la correlación con la columna de destino. Este paso sirve de comprobación de los datos en el caso de que no se haya hecho una ETL previa, o para comprobar que la ETL está bien hecha. También nos da la opción de eliminar alguna columna del entrenamiento del modelo si así lo consideramos. Para eliminarle peso a los pisos demasiado caros o que tengan alguna característica atípica, se selecciona el objetivo de optimización MAE, que es el error absoluto medio de modelo.

Para terminar la configuración, se ingresa un presupuesto de horas estimado, que es el tiempo estimado que tardará en finalizar. Se selecciona también la opción *Habilitar la interrupción anticipada*. Esta opción es importante porque si el modelo tarda menos en de lo esperado en entrenarse, se parará antes de consumir todas las horas establecidas como presupuesto, lo que puede ahorrar dinero y tiempo en la creación de un modelo de ML.

Nombre de la columna	Transformación	Porcentaje faltante (recuento)	Valores distintos	Correlación con el destino
bathroom	Categorica	-	7	0.247
lift	Categorica	-	2	0.074
neighborhood	Categorica	-	10	0.031
price	Destino	-	726	-
real_state	Categorica	-	4	0.111
rooms	Numerica	-	11	0.22
square_meters	Numerica	-	250	0.44
terrace	Categorica	-	2	0.033

El entrenamiento incluye 8 columnas de atributos en total

Figura 21. Secuencia de pasos para ajustar los parámetros de tiempo de ejecución y la métrica en la columna objetivo

Figura 22. Selección del objetivo de optimización

Para comprobar la diferencia de modelos, se han creado dos modelos con objetivos de optimización distintos. El primer modelo se optimiza con la raíz cuadrada del promedio de los errores al cuadrado (RMSE), que va a permitirnos poder introducir y predecir pisos más caros o con características más atípicas de un modo más exacto. El segundo modelo se optimiza con el promedio de los valores absolutos de los errores (MAE), el cual se centra más en predecir mejor los pisos más comunes y deja los pisos extremos como valores atípicos.

### 4.3. Watson.x

Watsonx es la plataforma de inteligencia artificial empresarial de IBM, diseñada para facilitar la creación, entrenamiento, implementación y gestión de modelos de Machine Learning e IA generativa a escala.

El primer paso para crear el modelo de ML es crear una instancia. En nuestro caso, es una instancia gratuita en un servidor de Sídney con las características que se observan en la imagen:

Crear

Acerca de

Tipo

Servicio

Proveedor

IBM

Última actualización

02/06/2025

Categoría

IA / Aprendizaje automático

Conformidad

Habilitado para HIPAA

Habilitado para IAM

Ubicación

Sidney (au-syd)

Frankfurt (eu-de)

Londres (eu-gb)

Tokio (jp-tok)

Dallas (us-south)

Toronto (ca-tor)

Enlaces relacionados

Documentos

Condiciones

Active Promotion ya ha aplicado un código promocional que le proporciona \$200 sobre este servicio. <b></b>

X

Seleccione una ubicación

Sidney (au-syd)

Seleccione un plan de precios

Los precios mostrados no incluyen impuestos. Los precios mensuales que se muestran son para el país o ubicación: España

Plan	Características y prestaciones	Tarifas
Lite	<div>1 usuario autorizado</div> <div>Límite mensual de 10 horas de unidad de capacidad</div> <div>Entorno = Número de unidades de capacidad necesarias por hora</div> <ul style="list-style-type: none"> <li>1 vCPU + 4 GB RAM = 0,5</li> <li>2 vCPU + 8 GB RAM = 1</li> <li>4 vCPU + 16 GB RAM = 2</li> </ul> <div>Decision Optimization + Watson NLP = Entorno + 6</div> <div>Generador de datos sintéticos, 2 vCPU + 8 GB RAM = 7 (requiere Watson Machine Learning)</div>	Gratis

El plan Lite para Watson Studio le ofrece todo lo que necesita para convertirse en un mejor científico de datos o experto en dominios en un entorno de colaboración.

Los servicios del plan Lite se suprimen tras 30 días de inactividad.

Figura 23. Creación de instancia en Watson.x

A continuación, se crea un proyecto. En este proyecto se adjuntarán los datos del CSV al que ya se le ha realizado el ETL previamente. Para cargarlos, es suficiente con arrastrar el archivo deseado a la pantalla *Importar archivos de datos*.

## Crear un proyecto

Empiece con un proyecto nuevo en blanco o seleccione desde donde importar un proyecto existente.

+ Nuevo

Archivo local

Ejemplo

Definir detalles

Nombre

Proyecto\_2

Descripción (opcional)

Predicción del precio de la vivienda en Barcelona

Etiquetas (opcional)

Añadir etiquetas

Añada etiquetas para que los proyectos sean más fáciles de encontrar. Para añadir etiquetas, sepárelas con comas y pulse Intro.

Almacenamiento

CloudObjectStorage

El proyecto incluye la integración con [Cloud Object Storage](#) para almacenar activos de proyecto.

Valores avanzados

Figura 24. Creación de proyecto en Watson.x

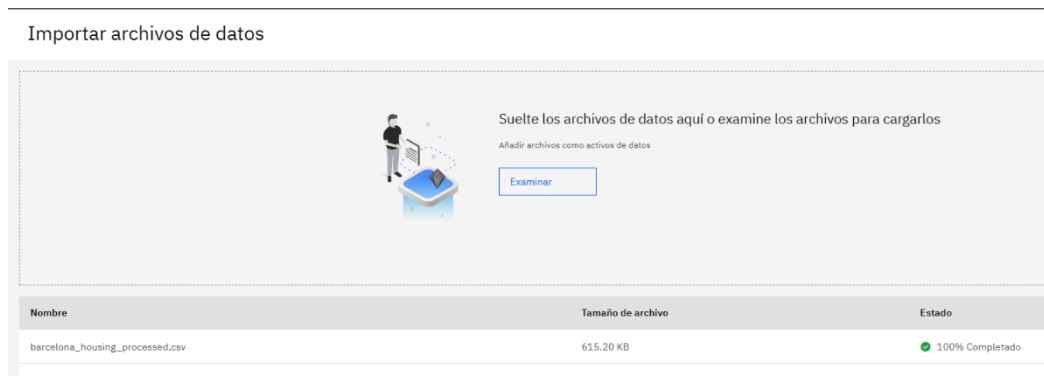
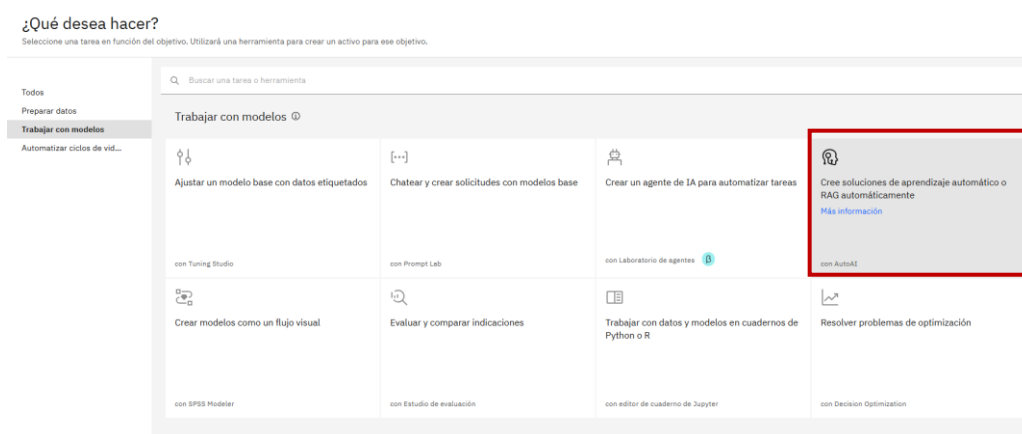


Figura 25. Importación de archivos en Watson.x

Con los datos cargados, debemos decidir qué queremos hacer con ellos. Para este proyecto, se pide crear un modelo de aprendizaje automático:



## Crear modelos de aprendizaje automático automáticamente

Define the details to create an AutoAI experiment asset using machine learning models or retrieval-augmented generation patterns and open it in the AutoAI tool.

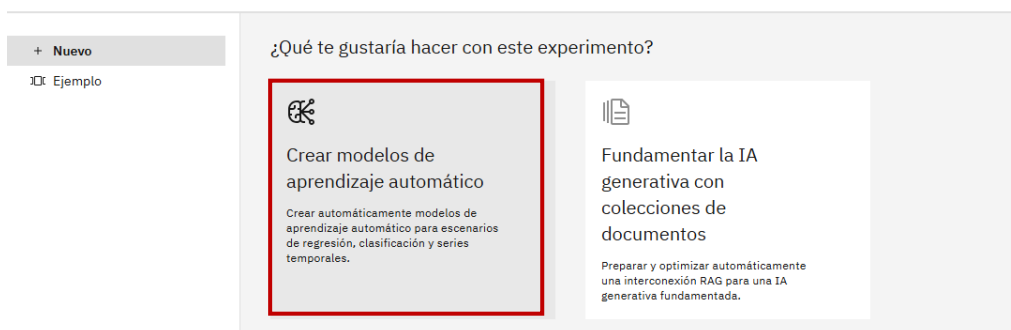


Figura 26. Elección de crear un modelo de aprendizaje automático con los datos cargados

Para poder crear el modelo, hay que asociar un servicio que sea capaz de realizar esa tarea. Se asocia un servicio predeterminado llamado *WatsonMachineLearning* que se encuentra en Sídney, la misma ubicación donde hemos creado la instancia.

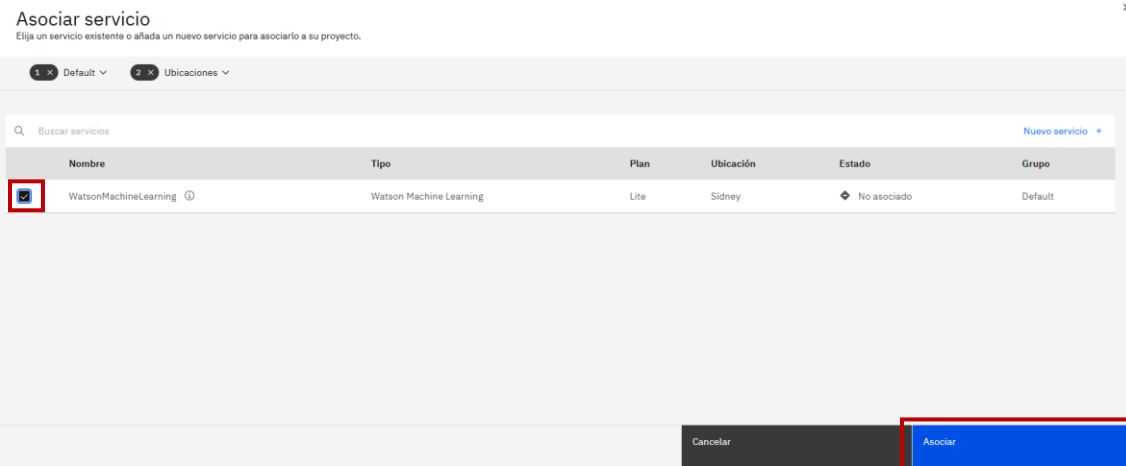


Figura 27. Asociación del servicio de WatsonMachineLearning en Watson.x

El siguiente paso es crear el modelo. Nuestro modelo lo llamaremos Proyecto\_2\_v1.0. Este nombre permitirá hacer un seguimiento de mejoras y cambios en el modelo gracias a poder cambiar la versión del nombre.

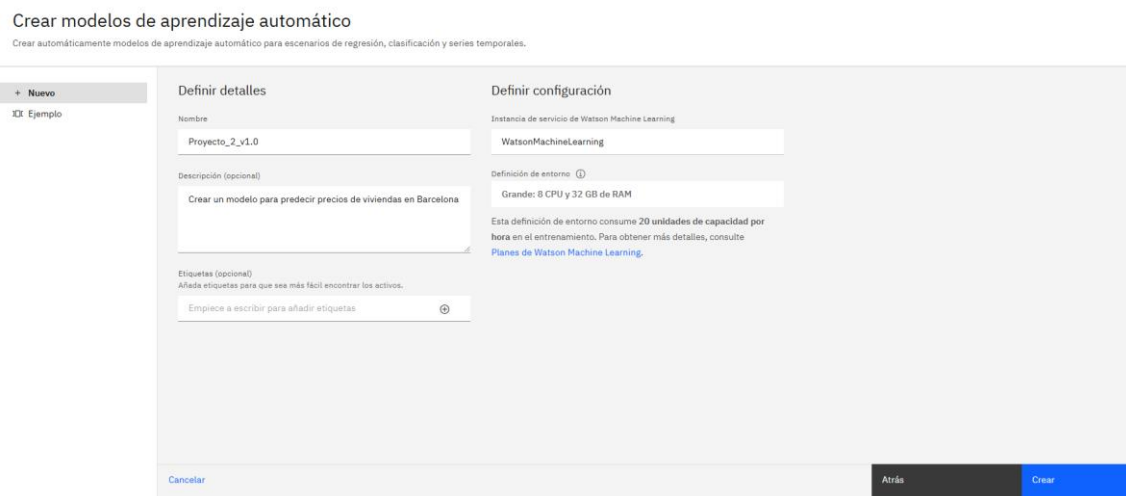


Figura 28. Creación del modelo en Watson.x

Se asigna el archivo CSV cargado previamente como origen de datos, se marca que no queremos un análisis de serie temporal y se indica que el objetivo de predicción debe ser la columna *Price*. Automáticamente el sistema asigna el tipo de predicción *Regresión* optimizado para *RMSE & tiempo de ejecución*.

Configurar experimento AutoAI  
Proyecto\_2\_v1.0 

Autoguardado: 18:12:47

**Añade origen de datos**

Añada archivos como, por ejemplo, datos tabulares (CSV).

Examinar

Seleccionar desde proyecto

barcelona\_housing\_processed.csv

Tamaño: 615.2 KB

Columnas: 8

**Configurar detalles**

¿Quiere crear un análisis de serie temporal?

Habilite esta opción para pronosticar la actividad futura en un rango de fecha/hora especificado. Los datos deben estar estructurados y ser secuenciales. [Información adicional](#)

Si

No

¿Qué desea pronosticar?

Columna de predicción: price

Columna de predicción: price

CUH restante: 20 CUH

TIPO DE PREDICCIÓN

OPTIMIZADO PARA

Regresión

RMSE & tiempo de ejecución

Valores de experimento

Ejecutar experimento

Figura 29. Asignación del archivo csv en Watson.x

23



## 5. Resultados de los modelos

### 5.1. BigML

Para comparar los resultados de los modelos escogidos en BigML, se escogerá un piso con unas características determinadas y se observarán las predicciones de todos ellos para todos los distintos barrios. Para nuestro caso, escogeremos un apartamento de 3 habitaciones con 2 baños, con ascensor y sin terraza. Los gráficos adjuntos muestran el precio de la vivienda en función de los metros cuadrados (eje X) y del barrio (eje Y). Se muestran los mapas de calor del precio para los distintos modelos:

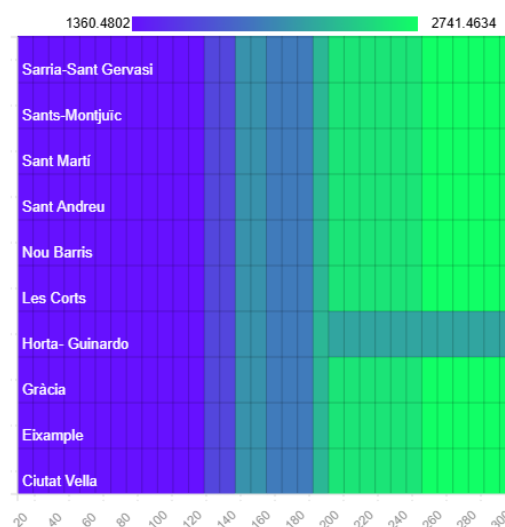


Figura 30. Resultados obtenidos por el modelo Random Forest

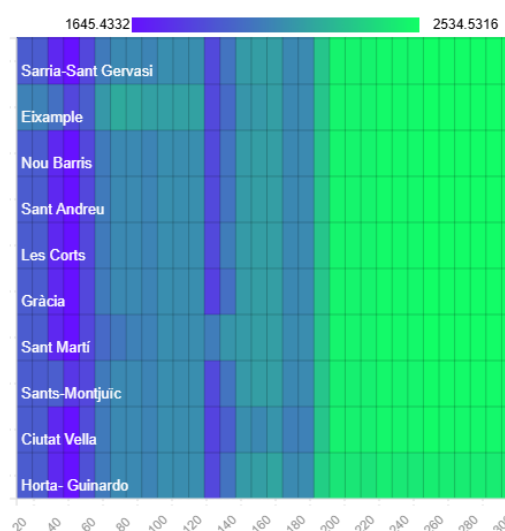


Figura 31. Resultados obtenidos por el modelo Ensemble

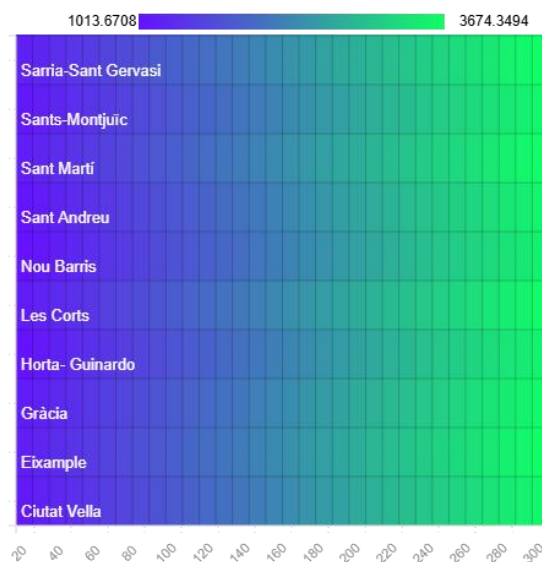


Figura 32. Resultados obtenidos por el modelo Linear Regression

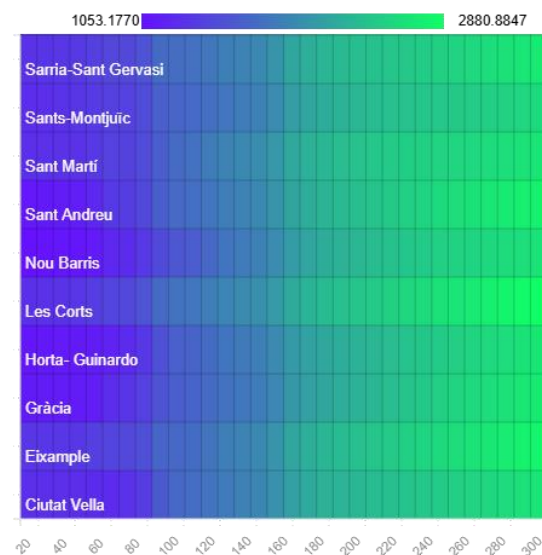


Figura 33. Resultados obtenidos por el modelo DeepNet

Tabla 1. Resumen de precios mínimos y máximos por cada modelo creado en BigML

Modelo	Precio mínimo	Precio máximo
Random Forest	1360 €	2741 €
Ensemble	1645 €	2535 €
Linear Regression	1014 €	3674 €
DeepNet	1053 €	2880 €

La tabla muestra una comparación de los precios mínimos y máximos estimados por diferentes modelos de predicción.

La diferencia en los rangos de precios estimados por los modelos puede deberse a varios factores relacionados con la forma en que cada modelo maneja los datos y los patrones en la información. Aquí hay algunas posibles explicaciones:

### **Dispersión y robustez del modelo**

- Linear Regression muestra el rango más amplio (1014 € - 3674 €), lo que sugiere que el modelo puede estar capturando relaciones lineales pero no necesariamente las no lineales o interacciones complejas entre variables. Esto podría explicar la mayor dispersión en las predicciones.
- Ensemble, en cambio, tiene el rango más ajustado (1645 € - 2535 €), lo que indica que la combinación de múltiples modelos ayuda a estabilizar la predicción, reduciendo la variabilidad.

### **Capacidad para capturar no linealidades**

- Random Forest y DeepNet pueden capturar relaciones no lineales mejor que Linear Regression, lo que podría hacer que sus predicciones sean más representativas de los datos reales. Sin embargo, DeepNet aún presenta una mayor dispersión (1053 € - 2880 €) en comparación con Random Forest (1360 € - 2741 €), posiblemente debido a la sensibilidad de redes neuronales a la cantidad y calidad de los datos disponibles.

### **Efecto de outliers y ruido en los datos**

- Linear Regression es más sensible a valores atípicos, lo que podría explicar por qué su rango de predicción es tan amplio. Unos pocos valores extremos en los datos de entrenamiento pueden influir significativamente en la predicción del modelo.
- Modelos basados en árboles (como Random Forest) suelen ser más robustos a outliers, lo que ayuda a reducir la variabilidad en las predicciones.

### **Sobreajuste vs. Generalización**

- DeepNet, a pesar de ser un modelo complejo, puede estar sobreajustando los datos de entrenamiento, lo que lleva a una mayor dispersión en las predicciones.
- Ensemble logra un mejor equilibrio al combinar modelos, reduciendo el sobreajuste y proporcionando predicciones más estables.

### **Conclusión**

El modelo Linear Regression presenta la mayor variabilidad debido a su simplicidad y sensibilidad a valores extremos. Ensemble ofrece predicciones más consistentes gracias a la combinación de modelos, mientras que Random Forest y DeepNet capturan mejor las relaciones no lineales pero con diferentes niveles de estabilidad.

## 5.2. Vertex.ai

En Vertex.ai se han comparado dos modelos con criterios de optimización distintos: uno minimizando el RMSE y el otro minimizando el MAE. Se presentan recortes de la comparativa de los resultados de ambos modelos y de la diferencia de importancia de los atributos entre ellos:

Versión del modelo	Evaluación	MAE	MAPE	RMSE	RMSLE	R <sup>2</sup>
<input checked="" type="checkbox"/> <a href="#">regresion-proyecto-2-impelia &gt; Versión 1</a>	<a href="#">untitled_1214977826490190968</a>	219.321	17.798	323.792	0.23	0.576
<input checked="" type="checkbox"/> <a href="#">regresion-impelia-proyecto-2-MAE &gt; Versión 1</a>	<a href="#">untitled_7281355709692816704</a>	213.667	16.364	333.746	0.232	0.564

Figura 34. Comparación de resultados obtenidos mediante optimización por RMSE y MAE

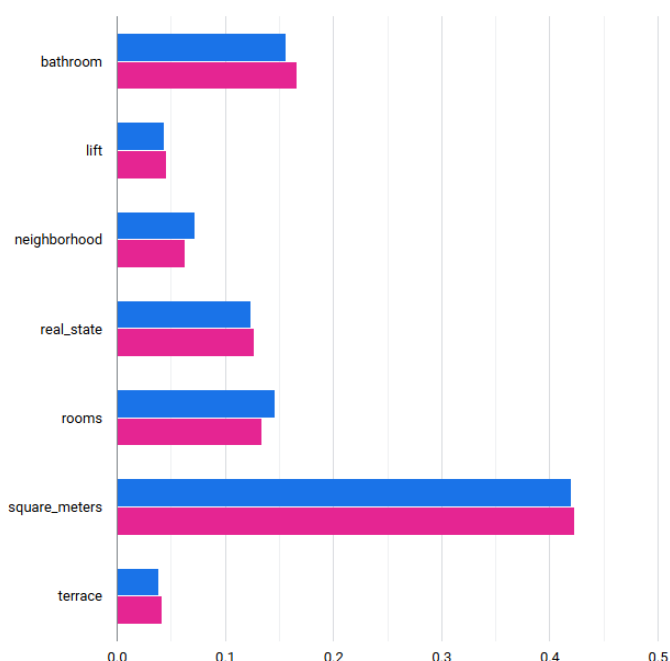


Figura 35. Comparación de la diferencia en importancias de cada atributo según el ajuste

Las pequeñas diferencias entre ambos modelos pueden deberse a distintos motivos. Algunos de los más importantes se listan a continuación:

### Optimización y su impacto en las métricas

- El modelo optimizado con RMSE tiene un menor RMSE (323.792 vs. 333.746), lo que era esperado ya que este modelo fue ajustado para minimizar esta métrica.
- El modelo optimizado con MAE tiene un menor MAE (213.667 vs. 219.321), lo cual también es coherente con su criterio de optimización.

### **Diferencias en los errores**

- MAE: Indica el error promedio absoluto. El modelo optimizado con MAE tiene un menor valor, lo que sugiere que, en promedio, sus predicciones están más cerca de los valores reales.
- RMSE: Penaliza más los errores grandes. El modelo optimizado con RMSE tiene un menor valor aquí, lo que indica que es mejor reduciendo grandes desviaciones en los datos.

### **$R^2$ y ajuste general**

- El modelo optimizado con RMSE tiene un mejor  $R^2$  (0.576 vs. 0.564), lo que sugiere que explica mejor la variabilidad en los datos en comparación con el modelo optimizado con MAE.
- Sin embargo, la diferencia en  $R^2$  no es muy grande, por lo que ambos modelos tienen un desempeño similar en términos generales.

### **Impacto en la interpretación y aplicación**

- Si se busca estabilidad y predicciones más consistentes, el modelo optimizado con MAE puede ser preferible, ya que minimiza los errores absolutos en promedio.
- Si es crítico reducir grandes errores y mejorar el ajuste global, el modelo optimizado con RMSE es mejor, ya que penaliza más los errores altos y ofrece un  $R^2$  ligeramente superior.

### **Conclusión:**

- El modelo optimizado con MAE es mejor si se quiere minimizar el error medio absoluto, lo cual es útil en situaciones donde errores grandes y pequeños tienen la misma importancia.
- El modelo optimizado con RMSE es mejor si se quiere evitar grandes desviaciones, ya que castiga más los errores grandes y mejora la capacidad explicativa del modelo.
- La elección final dependerá del contexto de uso: Si hay valores atípicos que se quieren minimizar, RMSE es mejor. Si se busca precisión promedio sin importar valores extremos, MAE es la mejor opción.

### 5.3. Watson.x

Se muestran capturas de pantalla de los resultados obtenidos del modelo. Los resultados que se muestran son el mapa de relaciones, la tabla de clasificación de interconexiones y el gráfico de métricas.

#### Mapa de relaciones ⓘ

Columna de predicción: price

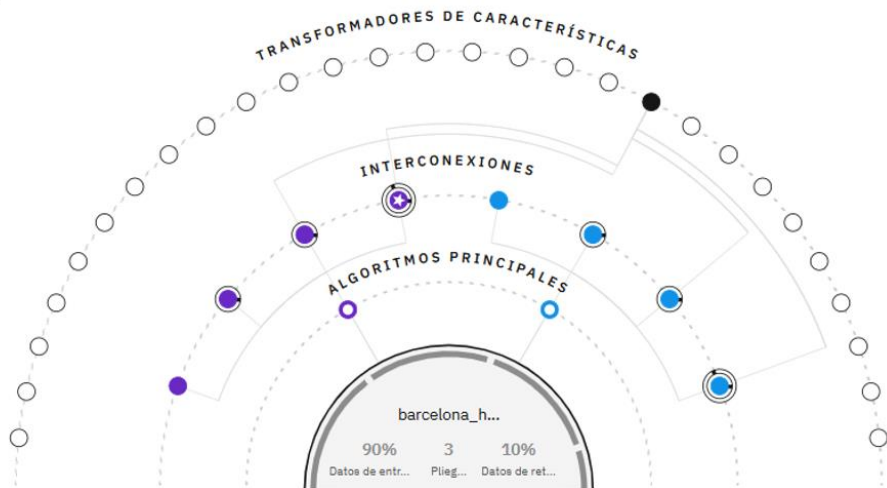


Figura 36. Mapa de relaciones obtenia por Watson.x

	Clasificación	↑	Nombre	Algoritmo	RMSE (Optimizado) Validación Cruzada	Mejoras	Tiempo de creación
★	1		4 Interconexiones	Regresor de bosque aleatorio de instantáneas	343.651	HPO-1 FE HPO-2	00:00:33
	2		3 Interconexiones	Regresor de bosque aleatorio de instantáneas	343.651	HPO-1 FE	00:00:25
	3		2 Interconexiones	Regresor de bosque aleatorio de instantáneas	344.709	HPO-1	00:00:06
	4		1 Interconexiones	Regresor de bosque aleatorio de instantáneas	344.709	Ninguna	00:00:02
	5		8 Interconexiones	Regresor de árbol de decisiones	358.274	HPO-1 FE HPO-2	00:00:29
	6		7 Interconexiones	Regresor de árbol de decisiones	359.580	HPO-1 FE	00:00:25
	7		6 Interconexiones	Regresor de árbol de decisiones	362.946	HPO-1	00:00:03
	8		5 Interconexiones	Regresor de árbol de decisiones	376.624	Ninguna	00:00:01

Figura 37. Tabla de interconexiones del modelo de Watson.x

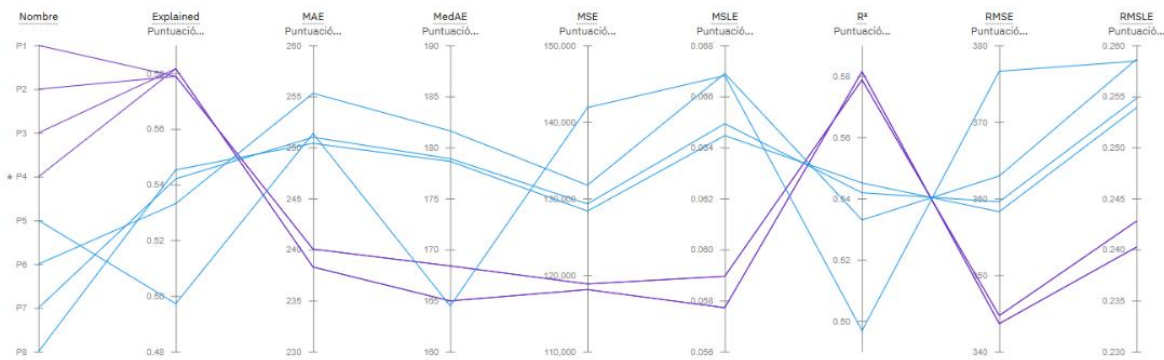


Figura 38. Gráfico de métricas del modelo de Watson.x

Analizando la tabla de métricas, se pueden observar varias líneas que representan diferentes modelos de regresión (P1 hasta P8) evaluados con diversas métricas de rendimiento:

- Varianza explicada: mide cuánta variabilidad de los datos es capturada por el modelo.
- **MAE** (Error Absoluto Medio): promedio de los errores absolutos.
- **MedAE** (Error Absoluto Mediano): la mediana de los errores absolutos.
- **MSE** (Error Cuadrático Medio): promedio de los errores al cuadrado.
- **MSLE** (Error Cuadrático Logarítmico Medio): similar al MSE, pero usando logaritmos.
- **R<sup>2</sup>**: Coeficiente de determinación.
- **RMSE** (Raíz del Error Cuadrático Medio): raíz cuadrada del MSE.
- **RMSLE** (Raíz del Error Cuadrático Logarítmico Medio): Raíz cuadrada del MSLE.

Analizamos los resultados de los 8 modelos:

### **Modelos más destacados**

P1 (línea púrpura) muestra el mejor rendimiento en varias métricas. Tiene la varianza explicada más alta (cercana a 1), un R<sup>2</sup> elevado, y valores bajos en métricas de error (MSE, MSLE). P3 y P4: También muestran un buen rendimiento en general.

### **Modelos con menor rendimiento**

P8 (línea azul clara) presenta valores más altos en métricas de error y valores más bajos en varianza explicada y R<sup>2</sup>. P7 también muestra rendimiento inferior en la mayoría de las métricas.

### **Comportamiento inconsistente**

Se observa que algunos modelos tienen buen rendimiento en ciertas métricas, pero peor en otras, lo que podría indicar diferentes fortalezas según el tipo de error que se quiera minimizar.

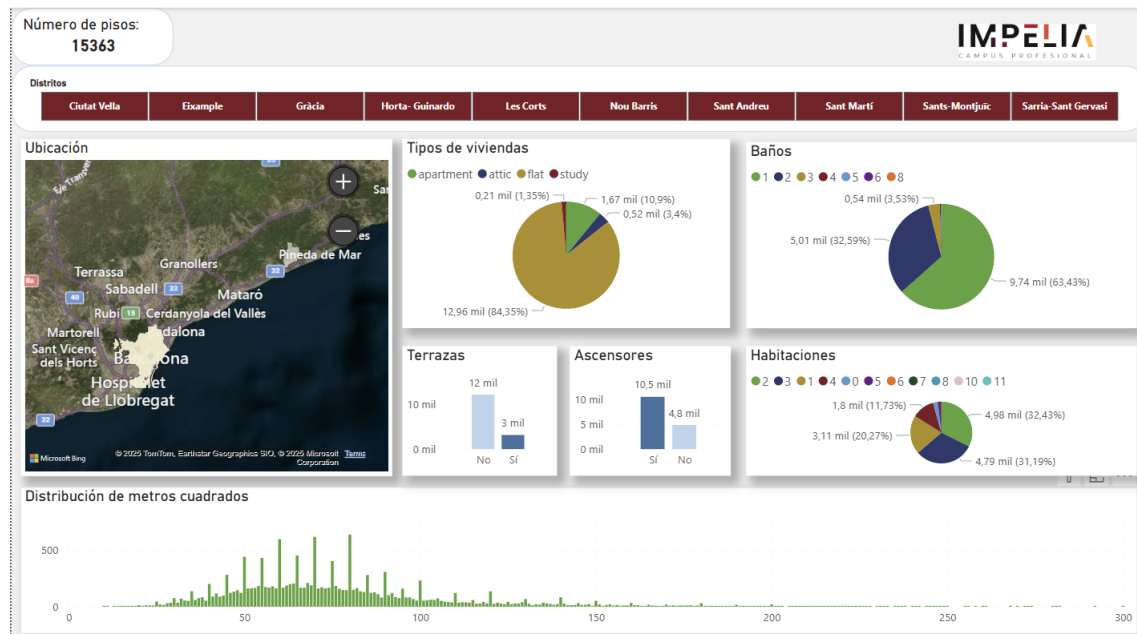
### **Conclusión**

P1 parece ser el modelo con mejor rendimiento general, con la mejor combinación de alta varianza explicada, alto R<sup>2</sup> y bajos valores de error.

Hay compensaciones evidentes entre los diferentes modelos porque algunos tienen mejor MAE, pero peor MSE, lo que indica que responden de manera diferente a errores grandes versus pequeños.

## 6. PowerBI

A modo de introducción, se diseñaron dos dashboards en función del distrito. El primero viene a ser el estudio, columna por columna donde informa la distribución de variables. El segundo es la variación del precio de alquiler y las principales columnas que dependen de ello.



cds

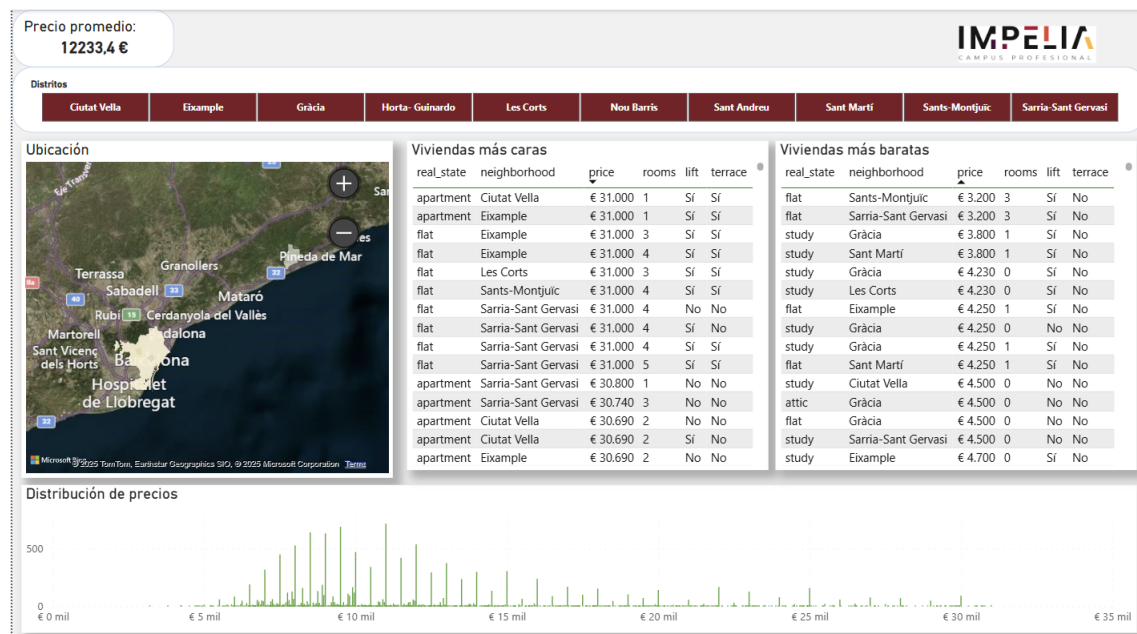


Figura 39. Minería de datos de la tabla procesada con WindSurf



## 7. Conclusiones

Después de realizar este proyecto, como equipo hemos llegado a las siguientes conclusiones:

1. Se comprobó que, mediante la conjunción de los métodos tradicionales de programación junto con las herramientas de IA como WindSurf, realizar el preprocesamiento de los datos ha sido una tarea mucho más rápida, sencilla y asistida.
2. Se han utilizado tres herramientas No-Code para la generación de modelos de aprendizaje automático: BigML, Vertex.ai (Google) y Watson.x (IBM). Todas ellas son de pago, y se ha sacado provecho de los periodos de prueba gratuitos que ofrecen para probar las aplicaciones.
3. Las  $R^2$  obtenidas en los modelos de Vertex.ai y Watson.x no son valores demasiado altos (0.56 y 0.57, respectivamente), lo que indicaría que el modelo no representa con fidelidad los datos. Viendo que ambas plataformas han dado una  $R^2$  similar, es probable que haya algún error en el criterio del preprocesamiento de datos que el equipo no haya sabido resolver.