

# Inteligencia artificial explicable (xai) con lime y shap: interpretación del modelo del TITANIC

Alejandro Lara Caballero

Erik Reyes Reyes

<sup>1</sup>Universidad Autónoma Metropolitana  
Unidad Cuajimalpa Departamento de Matemáticas Aplicadas y Sistemas

[alarac@cua.uam.mx](mailto:alarac@cua.uam.mx)

[ereyes@cua.uam.mx](mailto:ereyes@cua.uam.mx)

Referencia de este artículo [1].

## RESUMEN

La adopción y uso de modelos de aprendizaje automático para resolver diversos problemas en ciencia e ingeniería se han incrementado en años recientes. Aunque muchos de estos modelos han alcanzado una gran precisión, se puede considerar que funcionan como una “caja negra”, lo que dificulta una transparencia en el proceso de toma de decisiones. Ante la creciente necesidad de explicar el porqué de los resultados obtenidos por estos métodos, surge la Inteligencia Artificial Explicable (XAI, por sus siglas en inglés). Mediante bibliotecas en Python y usando el conjunto de datos del Titanic, en este trabajo se aplican dos técnicas populares: LIME y SHAP. Estos métodos permiten a los ingenieros validar la lógica detrás de variables críticas como el género, la edad y la clase socioeconómica. Estas técnicas pueden trasladarse fácilmente a otros ámbitos de estudio.

## 1. Introducción

El aprendizaje automático ha permitido desarrollar modelos de gran capacidad predictiva, como las redes neuronales o los ensambles, sin embargo, estas estrategias se pueden considerar opacas, ya que no es fácil interpretar su funcionamiento interno. En campos críticos como la medicina, las finanzas o la ingeniería forense, no solo es crucial que un modelo acierte, sino también conocer bajo qué criterios lo hace para generar confianza.

En este contexto surge la Inteligencia Artificial Explicable (XAI), que engloba técnicas y métodos que busca transformar los modelos de ‘caja negra’ en sistemas más transparentes. Dos herramientas de uso extendido son LIME y SHAP. Estas técnicas se aplicaron en un caso de estudio clásico: los datos de supervivencia del Titanic. El objetivo no es solo predecir quién sobrevivió al naufragio, sino entender qué factores influyeron en esa probabilidad y verificar si el modelo está aprendiendo patrones reales de la tragedia o simplemente sesgos injustificados.

Cuando hablamos de XAI, existen dos categorías principales de métodos. Por un lado, están los

modelos transparentes por construcción, como la regresión lineal, cuya lógica es clara desde el inicio. Por otro, encontramos los métodos denominados post-hoc, que actúan como traductores externos para explicar modelos complejos de 'caja negra'. En esta última clasificación destacan herramientas como LIME y [1].

## 2. Metodología

### 2.1 Datos, preprocesamiento y modelo

Para aplicar ambos métodos se utilizó el conjunto de datos del Titanic, el cual reúne información descriptiva de algunos pasajeros del famoso trasatlántico y un valor binario que indica si sobrevivieron o no al hundimiento. Este conjunto incluye un identificador único para cada persona (PassengerId), el resultado de supervivencia (Survived), la clase socioeconómica asociada al boleto (Pclass), así como datos personales como nombre, género y edad. También incorpora el número de boleto (Ticket), la cuota pagada (Fare), el número de camarote (Cabin) y el puerto donde cada pasajero embarcó, entre otros. Los datos se obtuvieron del repositorio público de Kaggle en [2].

Para cargar el conjunto de datos y comenzar el análisis, fue necesario preparar primero el entorno de trabajo en Python. Para ello se emplearon bibliotecas ampliamente utilizadas: pandas y numpy para manipular la información, matplotlib y seaborn para crear visualizaciones de los datos, y scikit-learn para aplicar clasificar los datos. Este proceso inicial se ilustra en la Figura 1.

```
1 # Suprimir Advertencias: Evita que aparezcan mensajes de advertencia que distraen.
2 import warnings
3 warnings.filterwarnings("ignore")
4
5 # Matplotlib Inline: asegura que los gráficos generados por matplotlib se muestren dentro del notebook.
6 %matplotlib inline
7
8 # Importar Librerías de Visualización:
9 import matplotlib.pyplot as plt # Para gráficos básicos.
10 import seaborn as sns # Para visualizaciones mejoradas y estadísticas.
11 sns.set() # Establece un estilo predeterminado atractivo para los gráficos de seaborn.
12
13 # Importar la Librería Principal de Aprendizaje Automático (scikit-learn):
14 import sklearn
15 from sklearn import preprocessing
16 from sklearn import metrics
17 from sklearn import model_selection
18 from sklearn.preprocessing import MinMaxScaler, StandardScaler, LabelEncoder, OneHotEncoder,
19     LabelBinarizer
20 from sklearn.metrics import roc_auc_score, roc_curve, accuracy_score, recall_score, precision_score,
21     f1_score, precision_recall_curve
22 from sklearn.model_selection import KFold, StratifiedKFold, GridSearchCV, train_test_split
23 from sklearn import linear_model
24 from sklearn.linear_model import LogisticRegression
25
26 # Importar Utilidades de Fecha y Hora:
27 from datetime import datetime, date, time, timedelta
28
29 # Importar Librerías Numéricas y de Manejo de Datos:
30 import numpy as np
31 import pandas as pd
32
33 # Importar Gestión de Sistema y Memoria:
34 import os, gc
35
36 # Cargar el dataset del Titanic:
37 df_train = pd.read_csv('https://raw.githubusercontent.com/alarac64/ColabRepo/rnfs/heads/main/train.csv')
38 df_test = pd.read_csv('https://raw.githubusercontent.com/alarac64/ColabRepo/rnfs/heads/main/test.csv')
39
40 # Mostrar las primeras filas del dataframe
41 print(df_train.head())
42 # Proporciona un resumen del dataframe, incluyendo tipos de datos y valores faltantes.
43 display(df_train.info())
```

Figura 1: Configuración del entorno y carga de datos

## 2.2 Construcción y Entrenamiento del Modelo

Una vez completada la configuración del entorno y la carga de los datos, el siguiente paso consiste en construir y entrenar el modelo de aprendizaje automático. En términos simples, un modelo es una herramienta matemática que aprende patrones a partir de los datos para poder hacer predicciones. En esta etapa utilizamos configuraciones predeterminadas que suelen ofrecer un desempeño adecuado. Para la demostración se emplea un clasificador basado en LightGBM [3], que funciona de forma apropiada en el conjunto de datos del Titanic. Este proceso se observa en la Figura 2.

```

1 ! pip uninstall -y lightgbm
2 ! pip cache purge
3 ! pip install --no-cache-dir lightgbm==4.5.0 # Force install 4.5.0
4 import lightgbm as lgb
5
6 exclude_columns = [
7     'Name',
8     'Ticket',
9     'PassengerId',
10    'Survived'
11 ] # Columnas a excluir
12
13 features = [c for c in df_train_fe.columns if c not in exclude_columns] # Columnas de características
14 target = df_train_fe['Survived'] # Variable objetivo
15 print(len(target)) # Imprimir la longitud del objetivo
16
17 gc.collect() # Liberación de memoria
18
19 X_train, X_test, y_train, y_test = train_test_split(df_train_fe[features], target, test_size=0.2,
20                                                  random_state=42) # Dividir los datos
21
22 param = { # Parámetros de LightGBM
23     'boost': 'gbdt',
24     'learning_rate': 0.008,
25     'feature_fraction': 0.38,
26     'bagging_freq': 1,
27     'bagging_fraction': 1,
28     'max_depth': -1,
29     'num_leaves': 17,
30     'lambda_l2': 0.9,
31     'lambda_l1': 0.9,
32     'max_bin': 200,
33     'metric': ['auc', 'binary_logloss'],
34     'tree_learner': 'serial',
35     'objective': 'binary',
36     'verbosity': 1,
37     'seed': 42
38 }
39
40 num_round = 10000 # Número de rondas de Boosting
41 callbacks = [lgb.early_stopping(stopping_rounds=1000, verbose=1)] # Parada temprana
42 evals_result = {}
43
44 lgb_train = lgb.Dataset(X_train, y_train) # Datos de entrenamiento
45 lgb_valid = lgb.Dataset(X_test, y_test, reference=lgb_train) # Datos de validación
46
47 clf = lgb.train(param, lgb_train, num_round,
48               valid_sets=[lgb_train, lgb_valid],
49               callbacks=callbacks)
50
51 oof = clf.predict(X_test, num_iteration=clf.best_iteration) # Hacer predicciones OOF
52 print(f'OOF AUC: {roc_auc_score(y_test, oof):.4f}') # Imprimir OOF AUC
53
54 # Convertir probabilidades a predicciones binarias (0 o 1)
55 oof_binary = [1 if p >= 0.5 else 0 for p in oof]
56
57 # Calcular e imprimir la precisión
58 oof_accuracy = accuracy_score(y_test, oof_binary)
59 print(f'OOF Accuracy: {oof_accuracy:.4f}')
60
61 predictions = clf.predict(df_test_fe[features], num_iteration=clf.best_iteration) # Hacer predicciones en
62 # los datos de prueba
63 print("Predictions on test data (first 5):", predictions[:5]) # Imprimir las primeras 5 predicciones

```

Figura 2: Entrenamiento del modelo.

### 2.3 LIME

LIME [4] es el acrónimo de Local Interpretable Model-agnostic Explanations. Es una técnica que permite comprender cómo funciona un modelo de aprendizaje a partir de una predicción o ejemplo concreto. LIME construye una versión simplificada del modelo alrededor del ejemplo que se desea analizar.

Para generar esta interpretación, LIME crea variaciones del dato de entrada original y observa cómo cambia la salida del modelo. Posteriormente, simplifica el modelo aproximándolo a una regresión lineal. Este proceso puede expresarse matemáticamente como la búsqueda del modelo interpretable  $g$  que mejor se aproxima al modelo original  $f$  en torno al punto  $x$ :

$$\xi(x) = \underset{g \in G}{\operatorname{argm}^{\prime} \operatorname{in}} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

donde:

- $\mathcal{L}(f, g, \pi_x)$  mide qué tan bien el modelo sencillo  $g$  imita al modelo real  $f$  cerca del punto de interés  $x$ ,  $\pi_x$  es una función que asigna más peso a las muestras que están más cerca de  $x$ ,
- $\Omega(g)$  penaliza la complejidad del modelo  $g$  para asegurar que sea fácil de interpretar.

Es decir, LIME construye un acercamiento o "zoom" sobre un caso particular, creando un modelo simple y más fácil de entender que el modelo más complejo que tomó esa decisión.

En el bloque de código de la Figura 3 se muestra cómo utilizar la biblioteca LIME en Python para generar explicaciones sobre predicciones individuales en el conjunto de datos del Titanic. Para ello se emplea el módulo `lime_tabular`, que se emplea para interpretar modelos que reciben datos tabulares como entrada. Los comentarios incluidos en el propio código guían paso a paso cada parte del proceso.

```

1 import lime # Obtener la libreria LIME
2 import lime.lime_tabular # Obtener la parte de LIME para tablas de datos
3
4 def predict_fn(x):
5     """
6     Esta función toma datos y le pide a nuestro modelo (clf) que prediga la probabilidad
7     de supervivencia para ellos. Devuelve estas probabilidades tanto para 'no sobrevivió' como para
8     'sobrevivió'.
9     """
10    preds = clf.predict(x, num_iteration=clf.best_iteration).reshape(-1,1)
11    # Obtener la predicción del modelo (probabilidad de sobrevivir)
12    p0 = 1 - preds
13    # Calcular la probabilidad de no sobrevivir
14    return np.hstack((p0, preds))
15
16 explainerLine = lime.lime_tabular.LimeTabularExplainer(
17     X_train.values,
18     mode='classification',
19     feature_names=features,
20     class_names=["NotSurvived", "Survived"],
21     verbose=True
22 )
23 # Esto configura LIME para explicar nuestros datos tabulares y nuestras predicciones de supervivencia.
24 # - X_train.values: Los datos de los que aprendió el modelo.
25 # - mode='classification': Estamos prediciendo categorías (sobrevivió o no).
26 # - feature_names=features: Los nombres de las columnas en nuestros datos (como 'Age', 'Sex').
27 # - class_names=["NotSurvived", "Survived"]: Los nombres de las cosas que estamos prediciendo.
28 # - verbose=True: Dice a LIME lo que está haciendo.
29
30 np.random.seed(1)
31 # Asegurarse de que LIME haga la misma cada vez que lo ejecutamos para consistencia.
32
33 i = 0
34 # Queremos explicar la predicción para la primera persona en nuestros datos de entrenamiento.
35
36 exp = explainerLine.explain_instance(
37     X_train[features].values[i],
38     predict_fn,
39     num_features=10
40 )
41 # Pídele a LIME que explique la predicción para la persona elegida.
42 # - X_train[features].values[i]: Los datos de esa primera persona.
43 #

```

Figura 3: Código de Uso de LIME en Python.

## 2.4 SHAP

Por su parte, la técnica SHAP [5] es el acrónimo de SHapley Additive exPlanations. Es una técnica que permite interpretar las predicciones de un modelo calculando cuánto aporta cada característica o atributo al resultado final. Está inspirado en la teoría de juegos cooperativos, la cual utiliza los valores de Shapley para definir de manera justa cuánto aporta cada jugador al resultado final del grupo. En el contexto de un modelo de aprendizaje, los “jugadores” representan las variables, y SHAP calcula cuánto suma o resta cada una para obtener la predicción o resultado. La cantidad SHAP para una característica  $i$  se calcula a través de:

$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} [f(S \cup \{i\}) - f(S)]$$

donde:

- $F$  representa el conjunto total de características,
- $S$  es un subconjunto de características que no incluye a  $i$ ,  $f(S)$  simboliza la predicción del modelo usando solo las características en  $S$ ,
- $\phi_i$  mide la contribución promedio de la característica  $i$  en todas las combinaciones posibles.

El valor de Shapley descompone cualquier predicción  $f(x)$  como una suma aditiva:

$$f(x) = \phi_0 + \sum_{i=1}^M \phi_i$$

donde:

- $\phi_0$  es el valor base (la predicción promedio del modelo),  $\phi_i$  indica cuánto aporta
- cada característica para alejarse de ese valor base.

De esta forma, SHAP muestra qué variables influyeron más en una predicción específica y en qué dirección.

En la Figura 4 se muestra cómo emplear la biblioteca shap para generar explicaciones basadas en valores de Shapley. En particular, la función `shap.force_plot` permite visualizar de manera intuitiva cómo cada característica influye en la predicción de un caso específico, mostrando cómo cada característica contribuye a empujar la predicción desde el valor base (predicción promedio) hasta la predicción real.

```

1 import shap # Obtener SHAP para explicar las predicciones del modelo
2 shap.initjs() # Configurar SHAP para mostrar gráficos en el navegador
3
4 # Crear un explicador SHAP para nuestro modelo LightGBM ('clf')
5 # Le dice a SHAP cómo nuestro modelo basado en árboles toma decisiones
6 # 'raw' output significa que queremos la puntuación directa del modelo
7 explainer = shap.TreeExplainer(clf, model_output='raw')
8
9 # Calcular cuánto contribuyó cada característica a cada predicción en nuestros datos de entrenamiento
10 shap_values = explainer.shap_values(X_train)
11 instance_ind=0
12 # Mostrar un gráfico explicando la predicción para la primera persona en nuestros datos de entrenamiento
13 # Muestra cómo el valor de cada característica empujó la predicción hacia arriba o hacia abajo desde la
    predicción promedio
14 shap.force_plot(explainer.expected_value, shap_values[0], X_train.iloc[instance_ind,:])

```

Figura 4: Código para generar gráfica de fuerza SHAP en Python.

A diferencia de LIME, SHAP permite obtener una visión global sobre qué características son más importantes para el modelo en todo el conjunto de datos y, además, muestra cómo se distribuye su impacto en las predicciones. Para ello puede utilizarse la función `shap.summary_plot`, que resume visualmente la influencia de cada variable de manera clara. El código que genera esta gráfica se presenta en la Figura 5.

```

1 import shap # Para explicar las predicciones del modelo
2
3 # Mostrar un resumen de qué características fueron más importantes y cómo afectaron las predicciones
4 shap.summary_plot(shap_values, X_train)
5
6 # Mostrar un gráfico de barras más simple de la importancia de las características
7 shap.summary_plot(shap_values, X_train, plot_type='bar')

```

Figura 5: Código para generar Gráfica de resumen de características SHAP en Python.

### 3. Resultados y discusión

En este apartado se presentan algunas gráficas obtenidas al aplicar las técnicas LIME y SHAP al conjunto de datos del Titanic, utilizando los códigos en Python descritos previamente.

La Figura 6 presenta una explicación local generada con la técnica LIME. En este caso, el modelo estima una probabilidad del 71 % de que el pasajero sobreviviera, y LIME identifica que factores importantes como ser mujer, pertenecer a una clase alta y tener camarote en cierta zona contribuyen positivamente a esta conclusión.



Figura 6: Gráfica LIME tabular.

La Figura 7 muestra una gráfica de fuerza SHAP que explica las razones por las que un pasajero al azar del Titanic tuvo altas posibilidades de sobrevivir. La gráfica parte de un valor promedio del modelo y muestra cómo cada característica del pasajero empuja la predicción hacia arriba (en rojo) o hacia abajo (en azul). En este caso, variables como *Cabin* = 5, *Title* = 3, *Sex* = 0 y *Pclass* = 1 contribuyen claramente a aumentar la probabilidad estimada de supervivencia, hasta llegar a un valor final de  $f(x) = 0,88$ .



Figura 7: Gráfica de fuerza SHAP.

Por último, la Figura 8 muestra un SHAP summary plot que explica cómo cada característica del conjunto influye en la predicción de supervivencia. Las barras indican la importancia global de las variables, destacando el nivel social, el género, la ubicación del camarote y el tipo del boleto como las más relevantes. El diagrama superior muestra, para cada pasajero, si el valor de una característica aumenta o disminuye su probabilidad de sobrevivir: los puntos rojos representan valores altos y los azules valores bajos. Así, ser mujer o viajar en primera clase suele desplazar los puntos hacia la derecha (mayor probabilidad de supervivencia), mientras que ser hombre o pertenecer a clases bajas los desplaza hacia la izquierda (menor probabilidad).

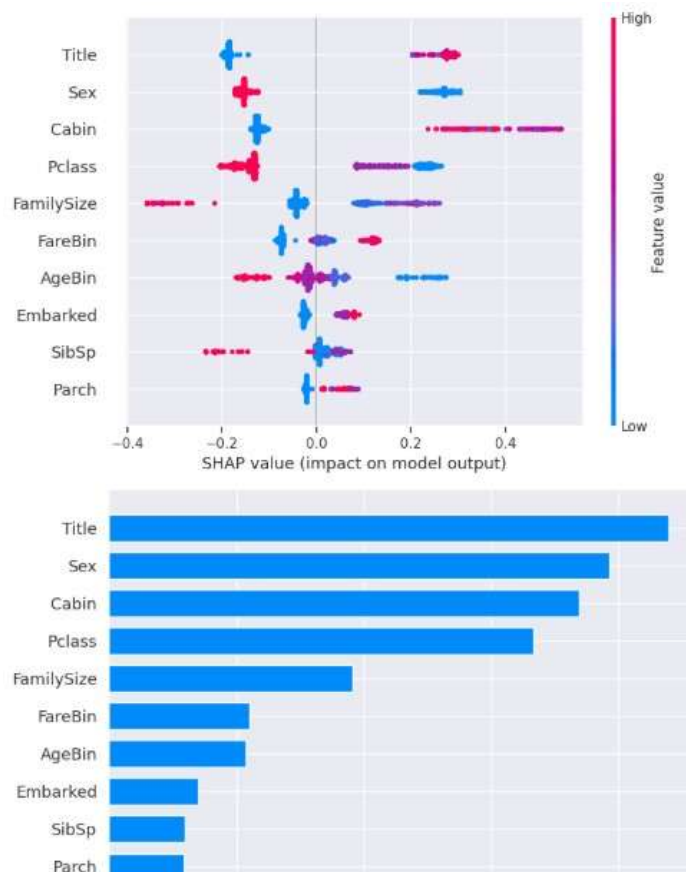


Figura 8: Gráfica de resumen de características SHAP.

#### 4. Conclusiones

En este trabajo se aplicaron las técnicas LIME y SHAP para interpretar el conjunto de datos del Titanic. Ambas herramientas permiten aportar mayor transparencia a los modelos de aprendizaje automático considerados como “cajas negras”. Las dos coinciden en señalar que el sexo, la clase social y la edad fueron factores determinantes en la supervivencia. LIME destaca por su claridad y rapidez para ofrecer explicaciones locales, mientras que SHAP proporciona una base matemática más sólida y permite obtener también una visión global del modelo. Es importante señalar que estas técnicas pueden trasladarse a otros contextos, pero sus explicaciones no deben interpretarse como relaciones causales.

#### Referencias

[1] Adadi, A., & Berrada, M. (2018). Peeking inside the black-box: A survey on explainable artificial

- intelligence (XAI). IEEE Access, 6, 52138–52160. <https://doi.org/10.1109/ACCESS.2018.2870052>
- [2] Yasser, H. (s.f.). Titanic Dataset [Conjunto de datos]. Kaggle. Recuperado de <https://www.kaggle.com/datasets/yasserh/titanic-dataset>
- [3] Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., & Liu, T.-Y. (2017). LightGBM: A highly efficient gradient boosting decision tree. In Advances in Neural Information Processing Systems 30 (pp. 3146– 3154). Curran Associates, Inc.
- [4] LIME: Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). “Why should I trust you?": Explaining the predictions of any classifier. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (pp. 1135–1144). ACM. <https://doi.org/10.1145/2939672.2939778>
- [5] SHAP: Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, & R. Garnett (Eds.), Advances in Neural Information Processing Systems 30 (pp. 4765–4774). Curran Associates, Inc.

## Referencia

Lara, A. & Reyes, E. (marzo - abril, 2026) Inteligencia artificial explicable (xai) con lime y shap: interpretación del modelo del TITANIC. *Boletín UPIITA. año 20, (11.) 2026*  
<https://www.boletin.upiita.ipn.mx/index.php/ciencia/1105-cyt-numero-113/2484-inteligencia-artificial-explicable-xai-con-lime-y-shap-interpretacion-del-modelo-del-titanic>