

Revelando los secretos de la caja negra: Una inmersión en las profundidades de la inteligencia artificial



The diagram consists of three concentric circles. The outermost circle is dark blue and contains the text 'ARTIFICIAL INTELLIGENCE' and its definition. The middle circle is a medium blue and contains the text 'MACHINE LEARNING' and its definition. The innermost circle is a light blue and contains the text 'DEEP LEARNING' and its definition. The circles are nested, indicating that Deep Learning is a subset of Machine Learning, which is a subset of Artificial Intelligence.

ARTIFICIAL INTELLIGENCE

A program that can sense, reason,
act, and adapt

MACHINE LEARNING

Algorithms whose performance improve
as they are exposed to more data over time

DEEP LEARNING

Subset of machine learning in
which multilayered neural
networks learn from
vast amounts of data

IA en nuestro cotidiano



(sin asunto)  Recibidos 



C

para mí

0:58 [Ver detalles](#)

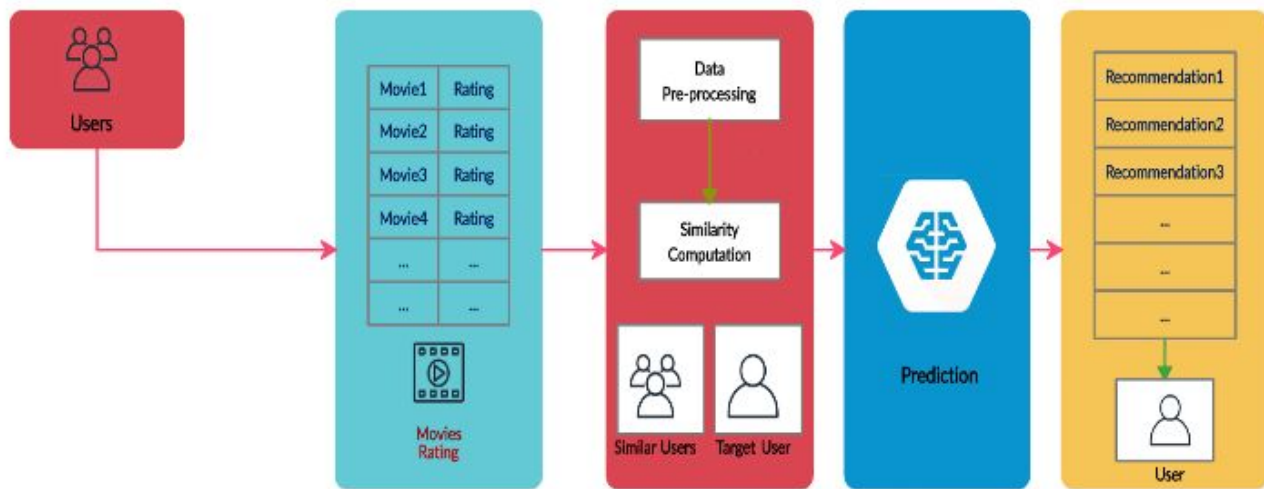


¿Vas a venir mañana?

Sí.

¿A qué hora?

No puedo.



NETFLIX

ChatGPT



What are you?

I'm a large language model trained by OpenAI. I'm a form of artificial intelligence that has been designed to process and generate human-like language.



Are you human?

I'm not a human and I don't have the ability to think or feel in the same way that a person does.





We have removed your post from Facebook

Post unavailable

We are unable to show content that goes against our Community Standards on hate speech.

To protect others' privacy and keep everyone safe from harmful content, we do not show content that violates Community Standards on hate speech.

What happened

Your appeal was reviewed and your post does not follow our Community Standards for hate speech.

What you can do

You can go to [Account status](#) to see how content violations can affect you.

Close

TRADITIONAL RULE-BASED APPROACH



MACHINE LEARNING APPROACH

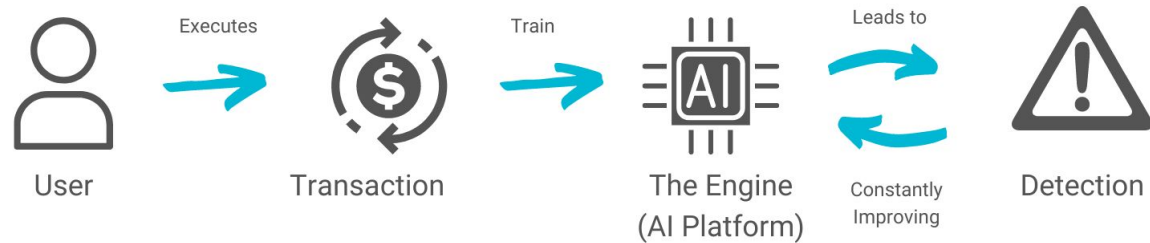
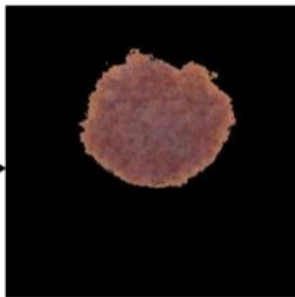


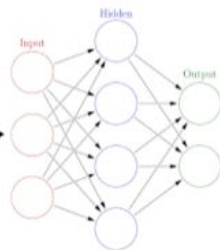
Image of skin lesion



Segmented lesion

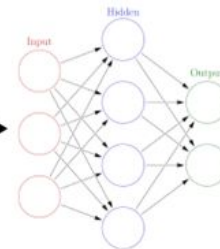
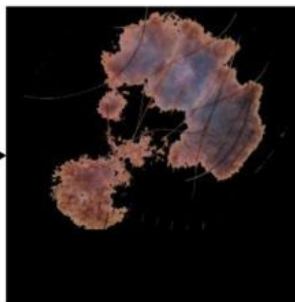
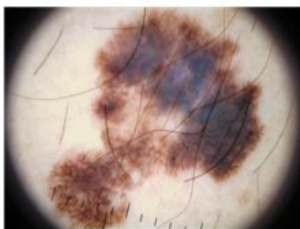


Learning Model

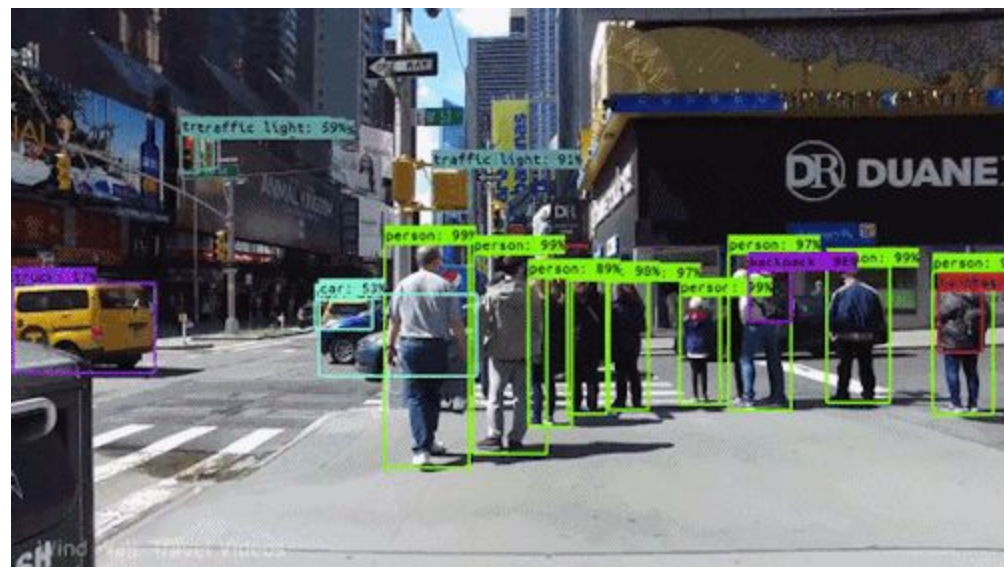


Prediction

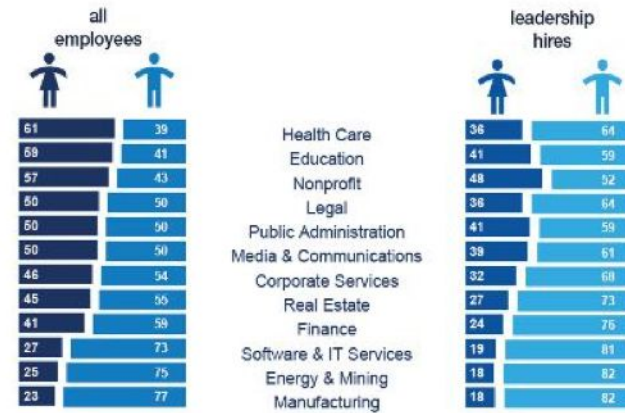
Benign



Malignant



Limitaciones y riesgos de la IA



Source: LinkedIn data featured in the
Global Gender Gap Report 2017, World Economic Forum

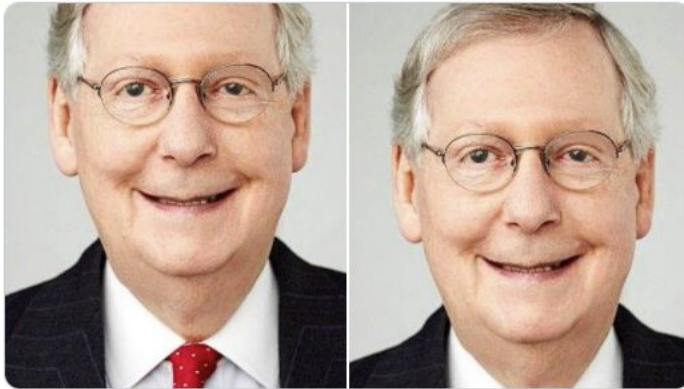


Tony "Abolish ICE" Arcieri 🦀
@bascule



Trying a horrible experiment...

Which will the Twitter algorithm pick: Mitch McConnell or Barack Obama?



8:05 AM · Sep 20, 2020



♥ 193.5K 💬 2.8K ↗ Share this Tweet



Explore Images Related to Crime Generated by Stable Diffusion

A color photograph of a **drug dealer**



Explore Images Related to Crime Generated by Stable Diffusion

A color photograph of an inmate



Explore Images Related to Crime Generated by Stable Diffusion

A color photograph of a **terrorist**



Como evitar esto?

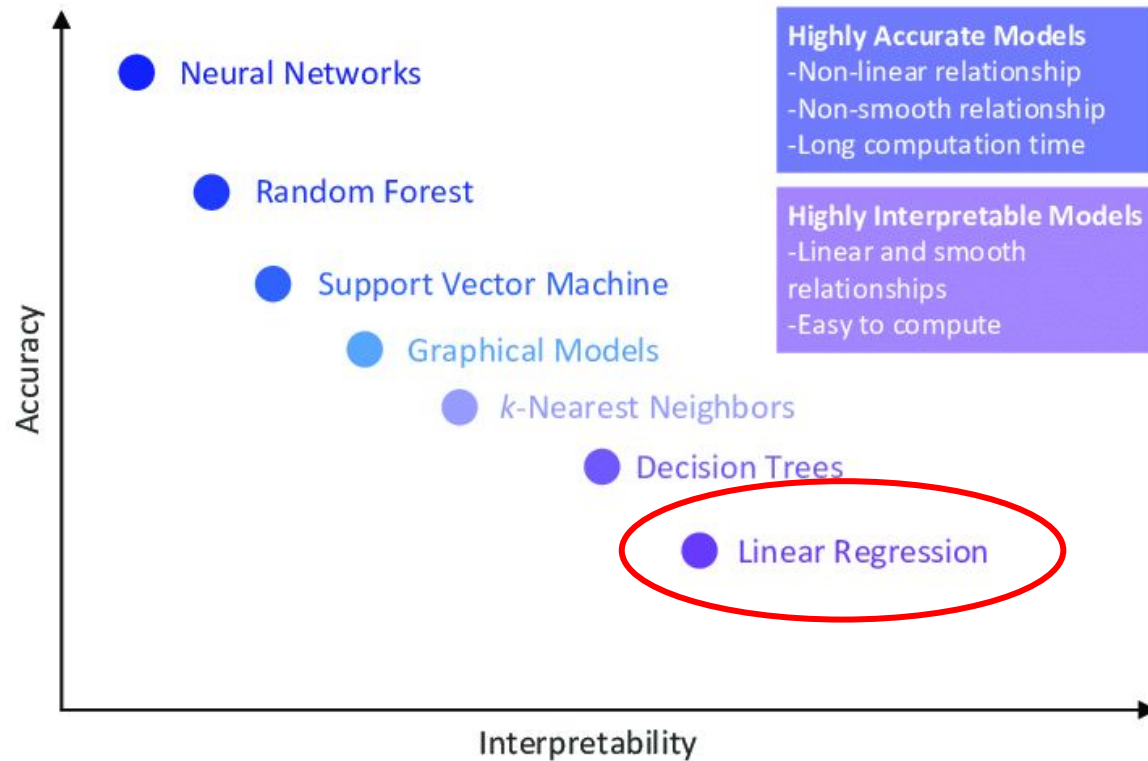
Interpretabilidad de Modelos

Modelos Paramétricos

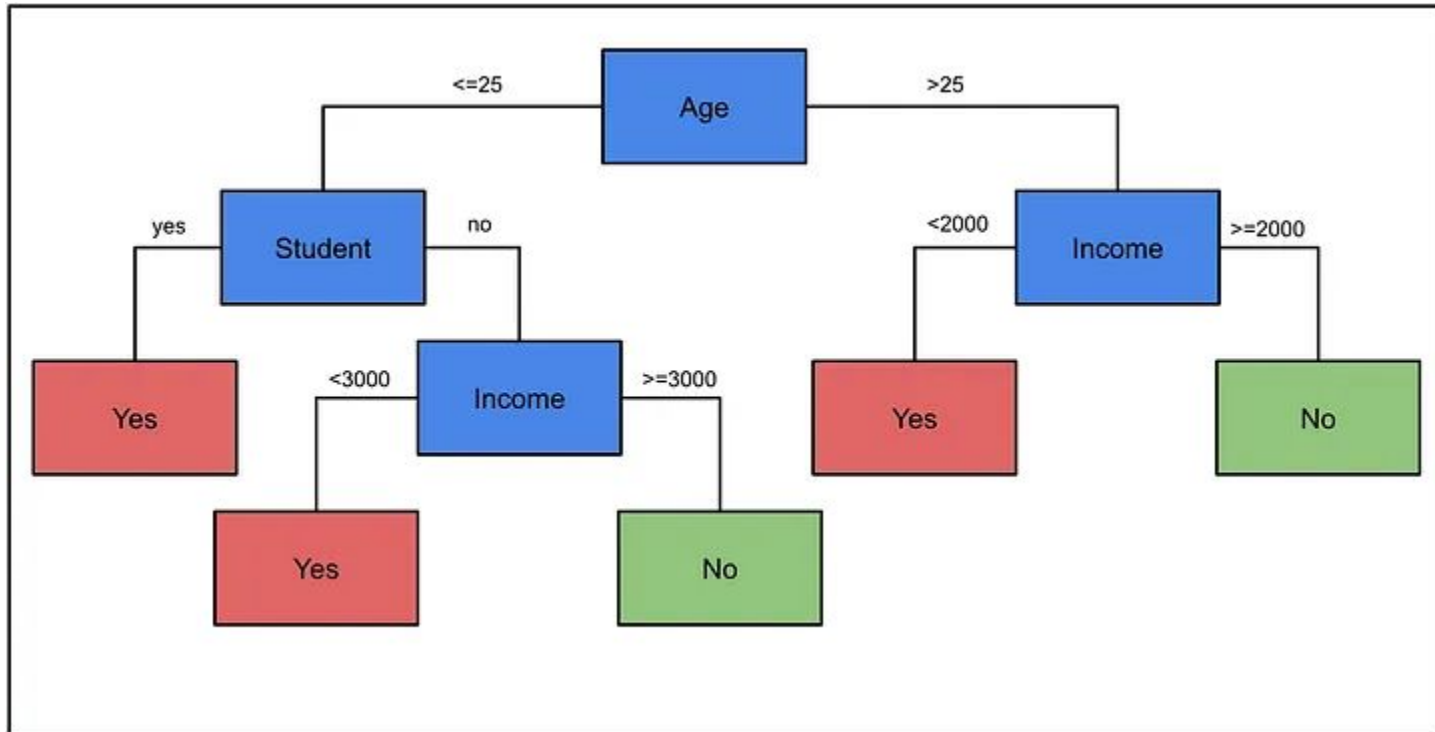
- La Regresión Lineal es un modelo paramétrico, lo que significa que la hipótesis se describe en términos de coeficientes que ajustamos para mejorar la precisión del modelo.

$$Y = 100*\mathbf{age} + 10*\mathbf{income} + 200$$

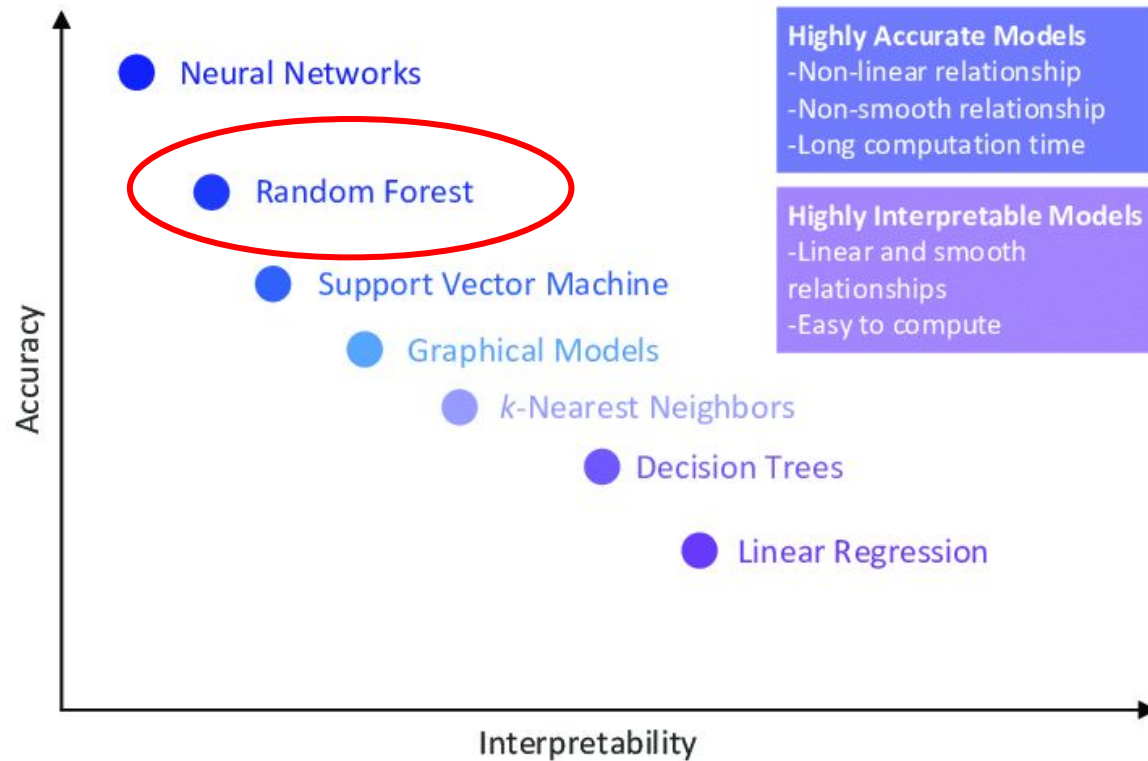
$$Y = X_0 + B_1 X_1 + B_2 X_2 + \dots + B_n X_n$$



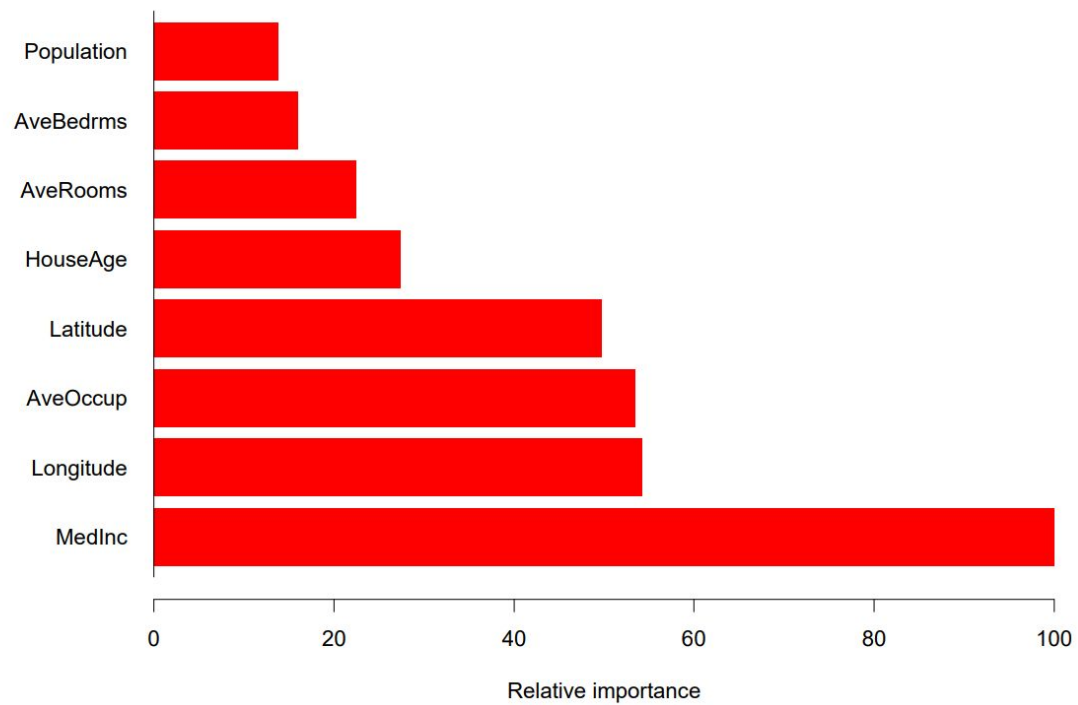
Modelos No Paramétricos

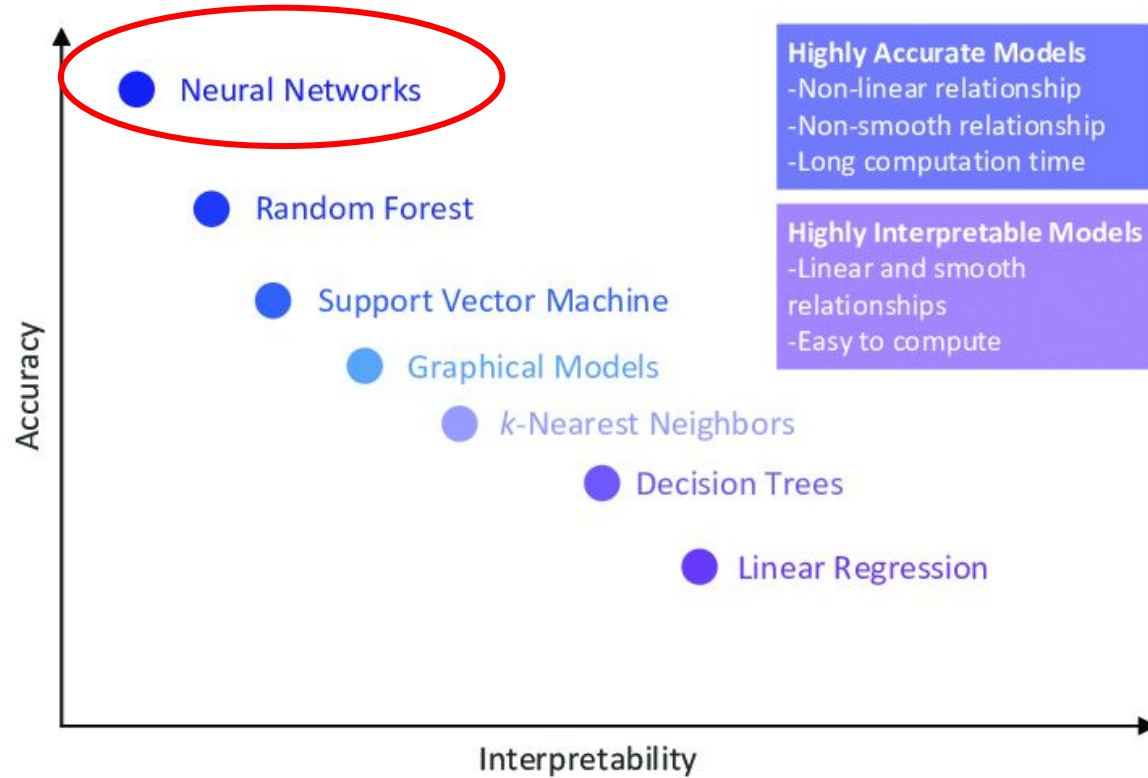


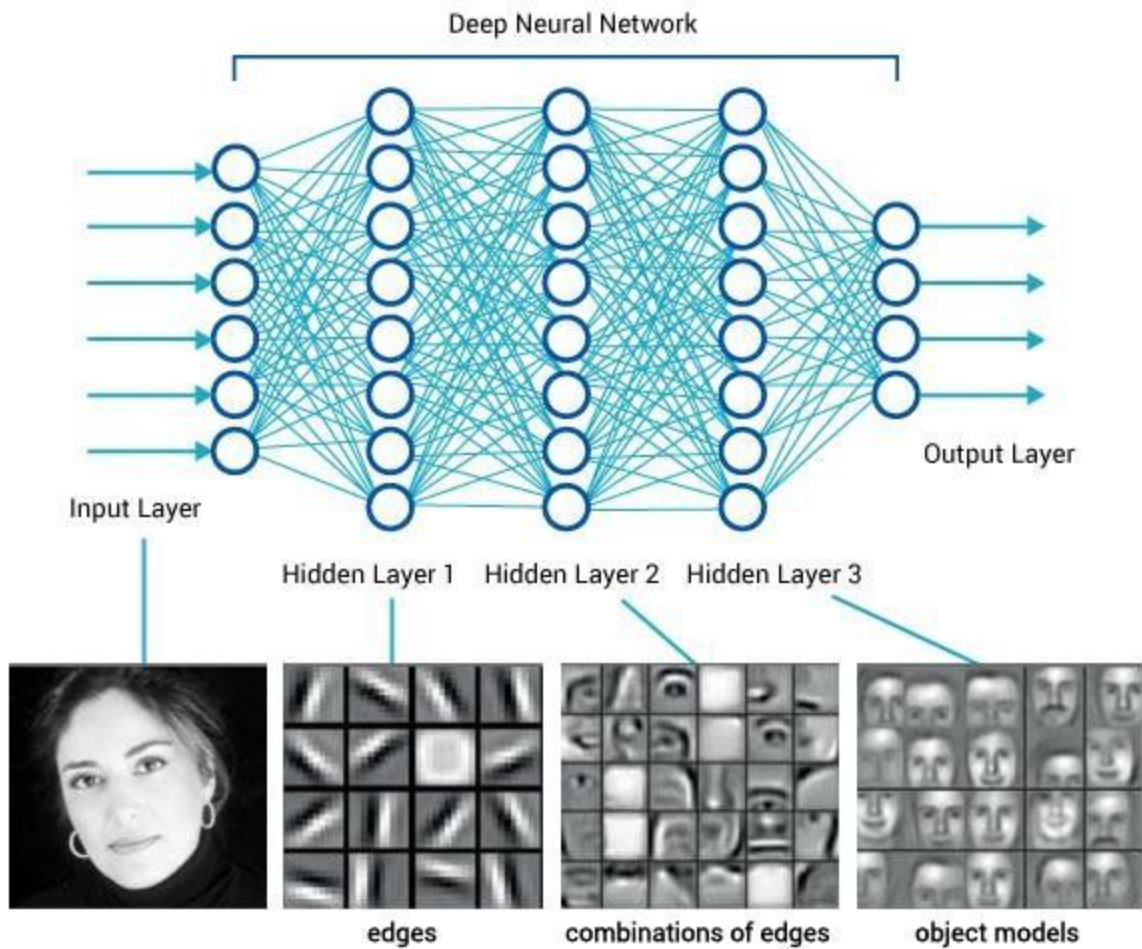
- Modelos como **KNN** o **Naive Bayes** no tienen parámetros que permitan entender la importancia de las variables ni su interpretación.
- Los modelos basados en árboles (**Random Forest, XGBOOST**) son otro ejemplo, no tenemos coeficientes. Sin embargo aún nos puede interesar saber cuáles features son más importantes.

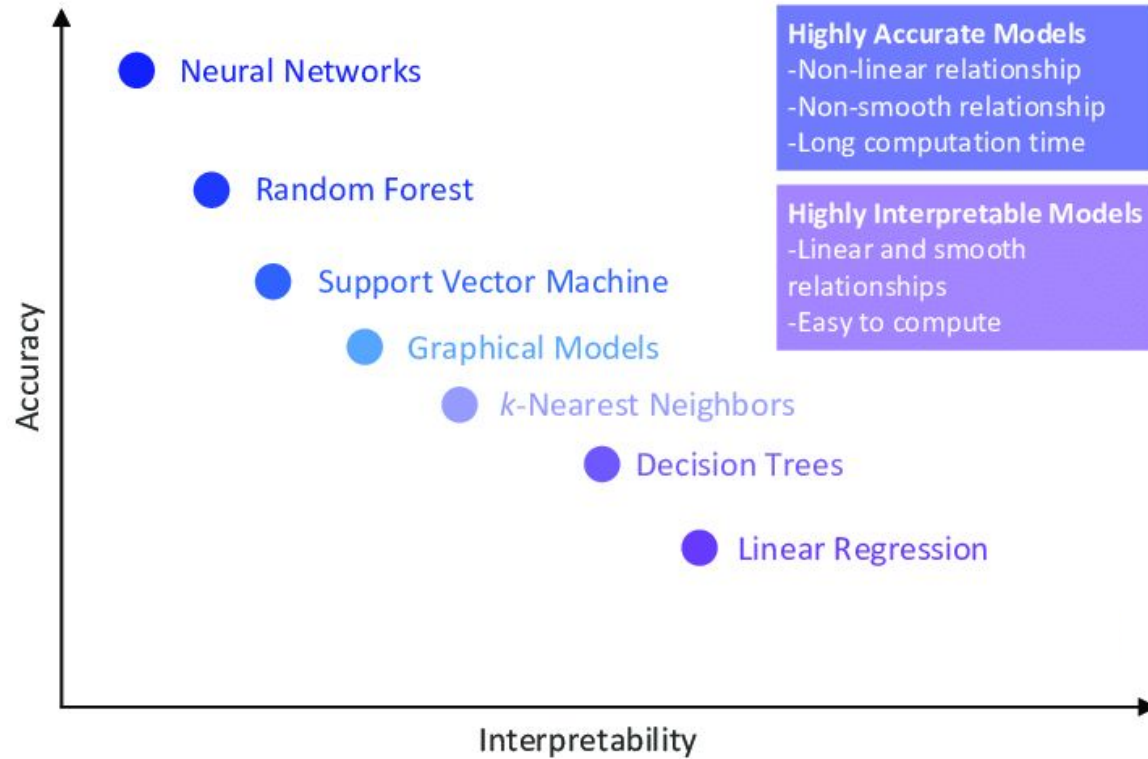


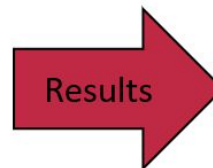
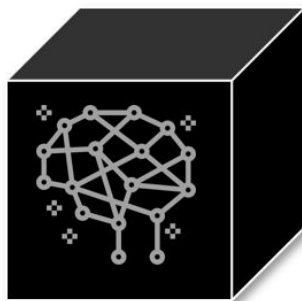
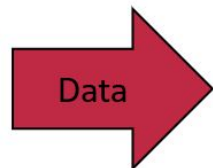
- Un random forest es un conjunto de árboles entrenados en muestras aleatorias y subconjuntos aleatorios de features.
- Podemos **calcular la importancia de los features de cada árbol y luego promediar las importancias en todo el ensamble.**



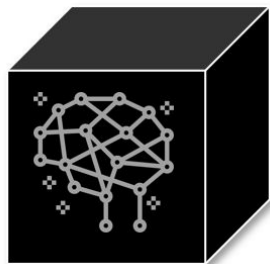




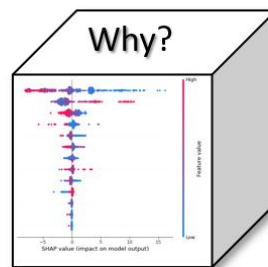
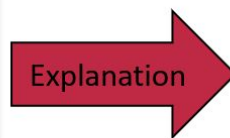
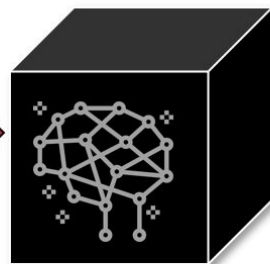
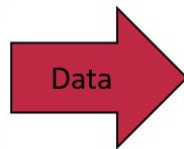




This is an
insect!



This is an insect!



Because it has 6 legs



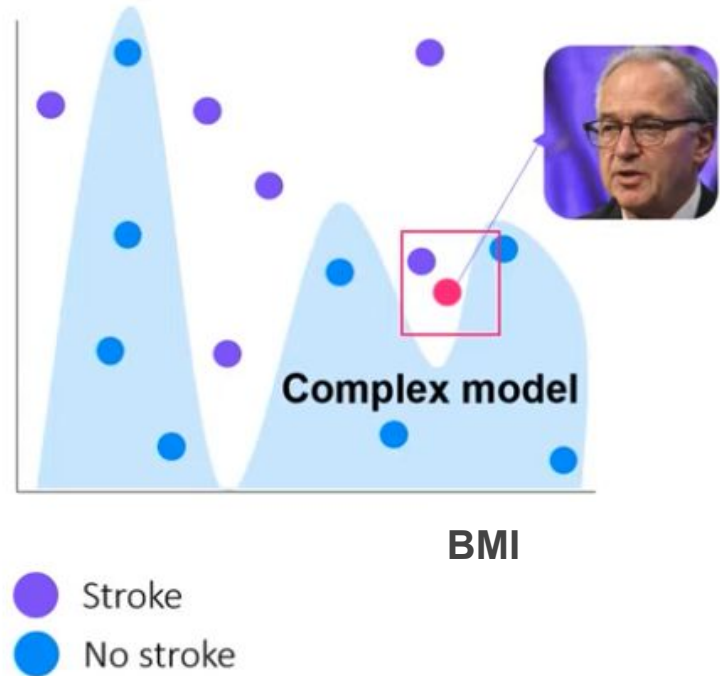
This is an insect!

Explicabilidad Local con LIME

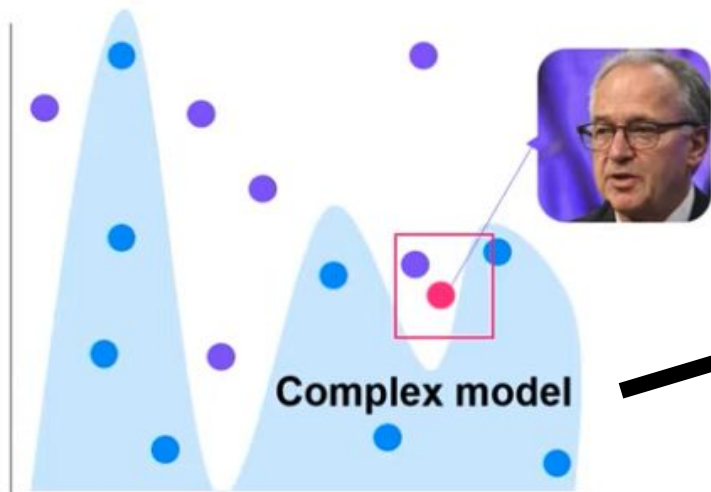
Complejidad: Nuestro principal problema es la alta no linealidad de nuestros modelos, como no tenemos parámetros, es muy difícil darle una interpretación general al modelo.

<https://github.com/marcotcr/lime>

Edad



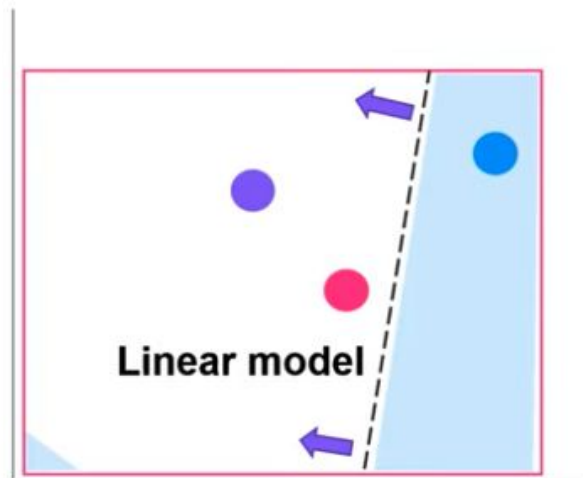
Edad



BMI

- Stroke
- No stroke

Edad

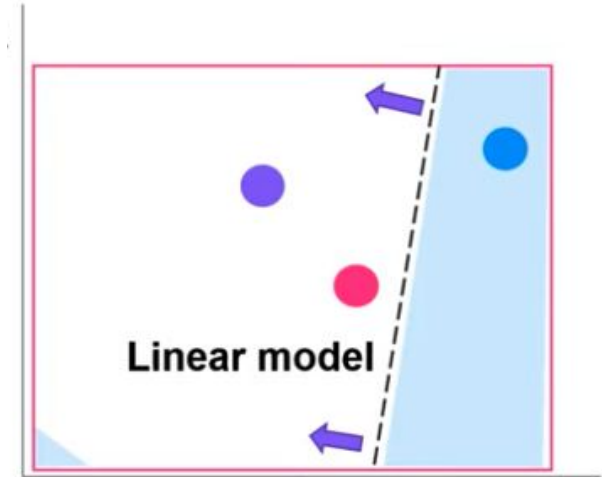


BMI

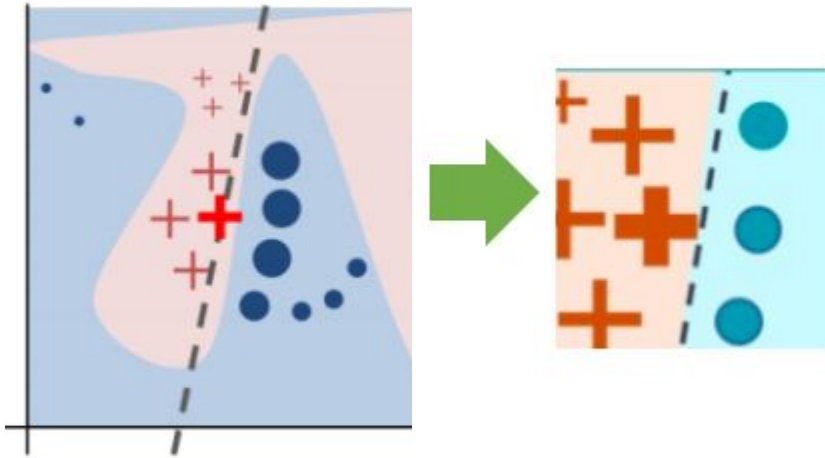
Local Interpretable
Model-Agnostic Explanations
(LIME)

Localmente: Alrededor de la predicción que buscamos explicar podemos recuperar un comportamiento lineal.

Edad



BMI

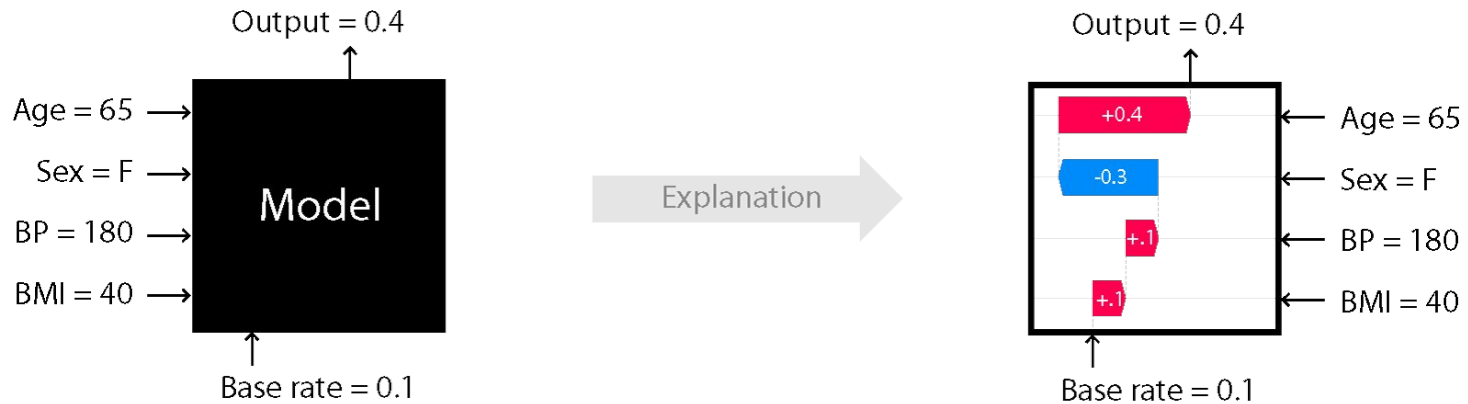


Confianza en una predicción: No podemos aceptar a ciegas la predicción de un modelo, especialmente cuando estamos decidiendo sobre la vida de otras personas..

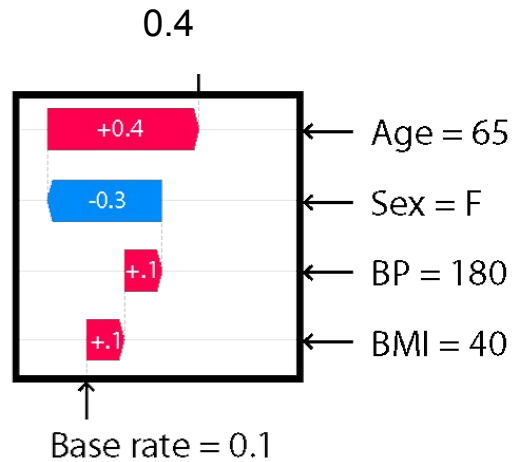
Confianza en el modelo: Más allá de las métricas, la confianza e interpretabilidad en las predicciones nos ayudan a confiar en el modelo para ponerlo en producción con datos reales.

Otro enfoque: SHAP

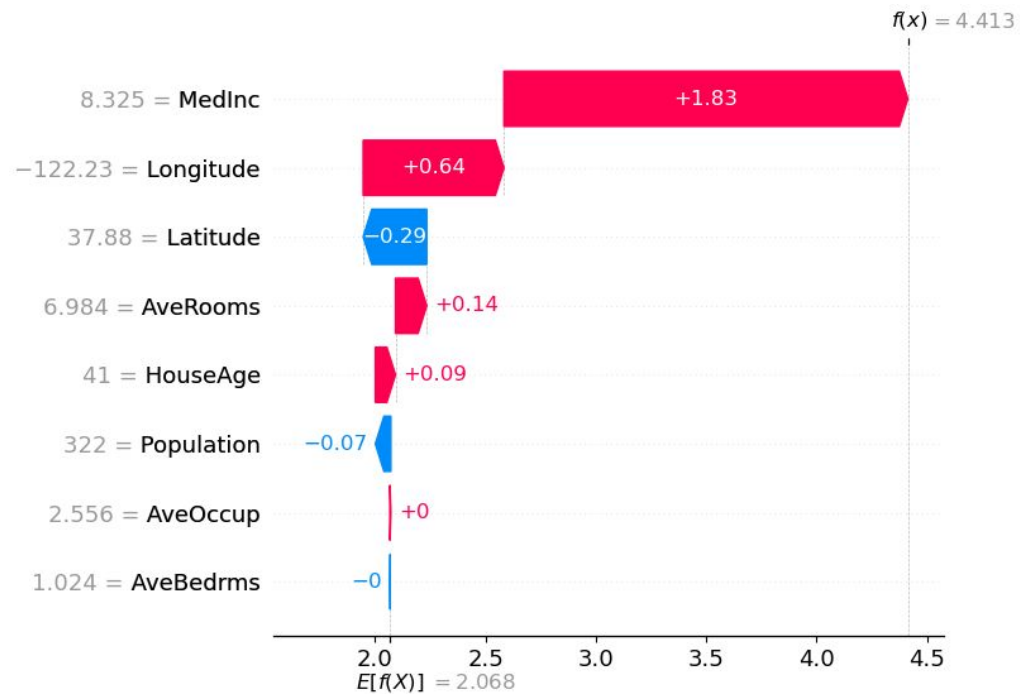
SHAP (SHapley Additive exPlanations)



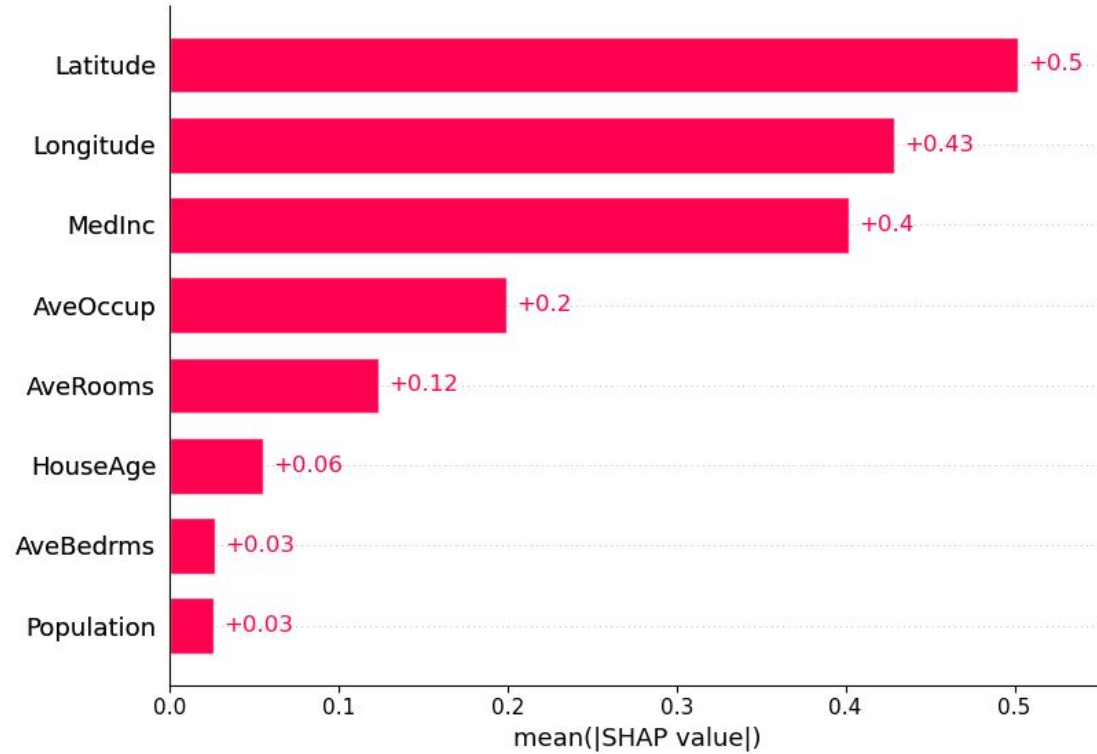
<https://github.com/slundberg/shap>



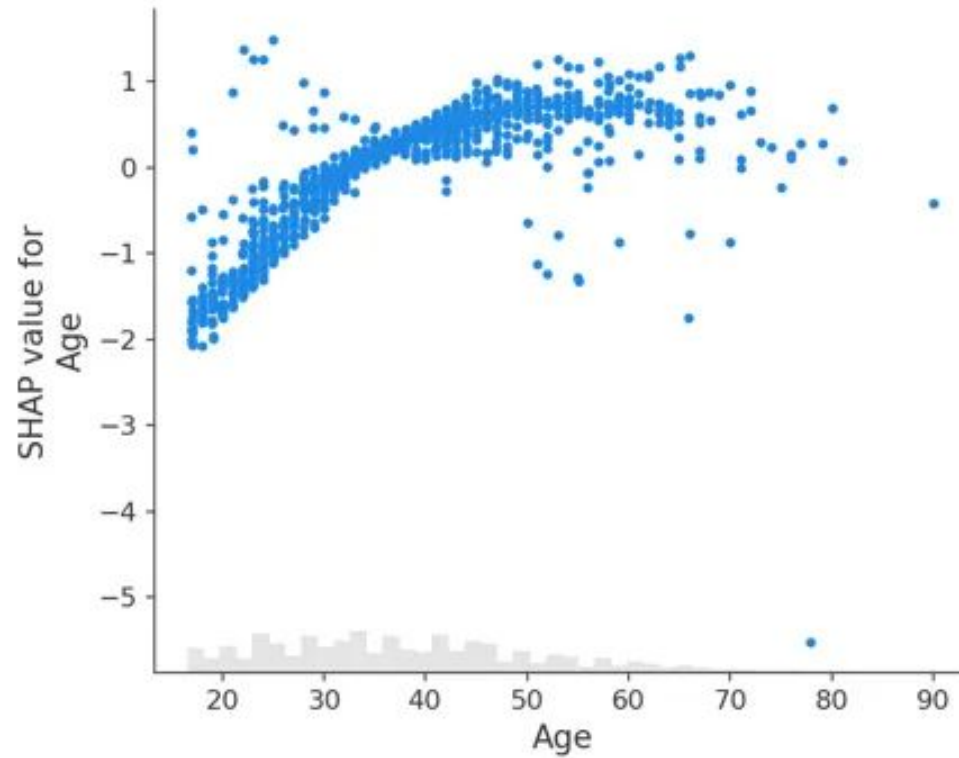
- Un método más robusto basado en el concepto de shapley value (teoría de juegos).
- Calcula la importancia para todas las permutaciones de variables y hace un promedio, esto retiene una mayor complejidad que LIME.

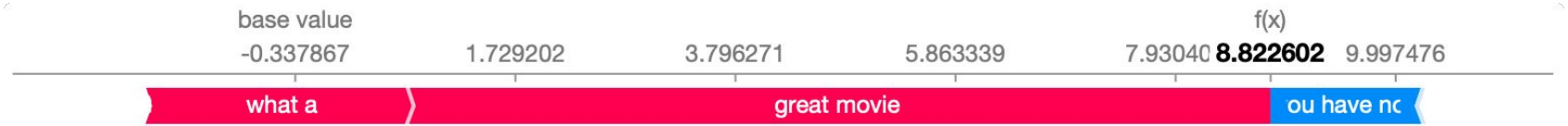


SHAP (SHapley Additive exPlanations) GLOBAL



Income vs. Age

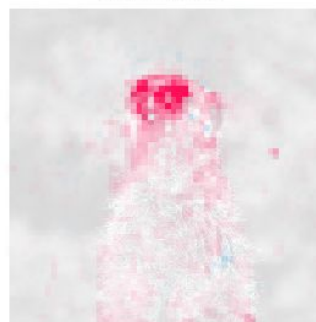




what a great movie! ... if you have no taste .



suricata



mangosta

