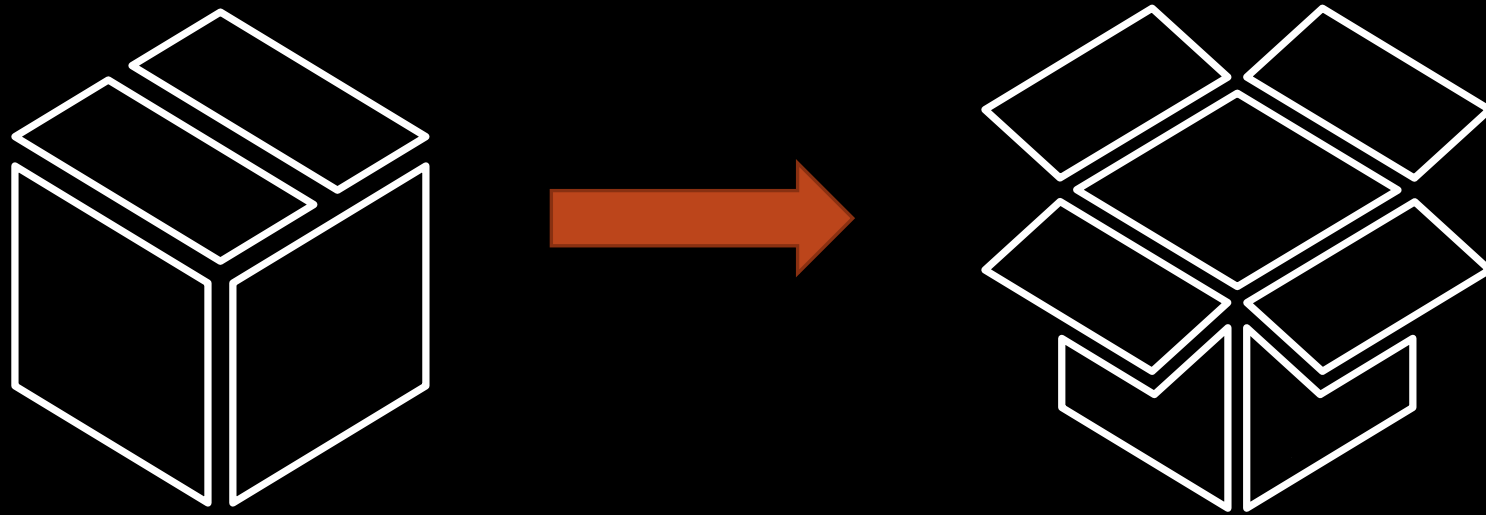


# Explicando modelos de Aprendizaje Automático



Ing. Daniel Hernández Mota, Científico de datos.

Agosto, 2021

## Daniel Hernández Mota:

Ingeniero en Nanotecnología  
(Computación Cuántica)



Científico de datos en Kueski  
(Modelo de fraude)



Mentor en SaturdaysAi  
(Guadalajara 2 y 3, LATAM 1 y 2)



Casi un IronHacker



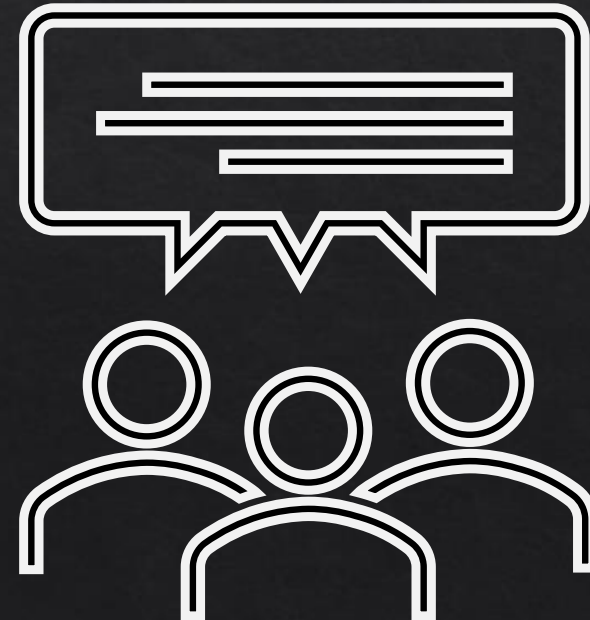
dhdzmota

**¿Por qué es importante  
explicar un modelo de  
Aprendizaje Automático?**

# Confianza

- Sentido.
- Lógica.
- Apego a la realidad.

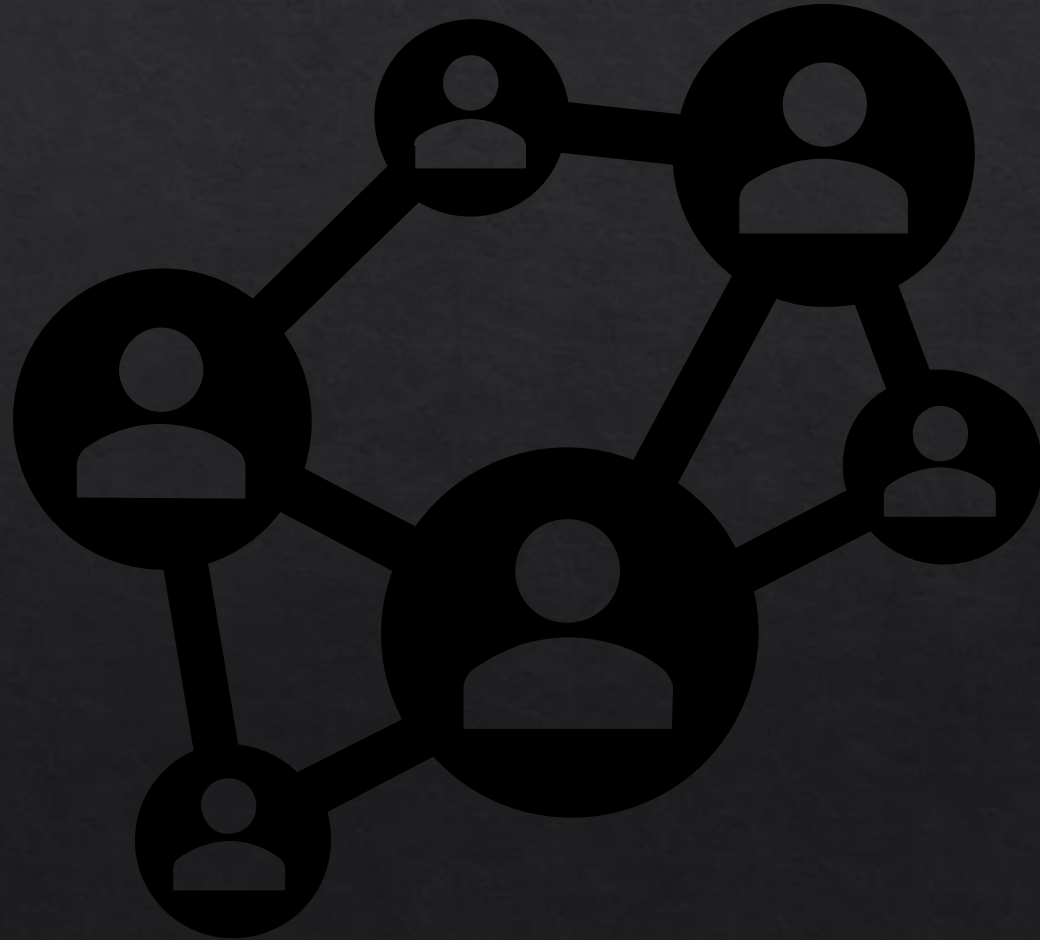
Si los usuarios no confían en el desempeño del modelo o en la predicción, no lo usarán.





# La IA se usa en todos lados...

- Identificación de objetos
- Seguros
- Redes Sociales
- Medicina
- Bancos
- Préstamos
- Medio ambiente
- Terrorismo
- Prisión



# ¿Cómo sabemos que podemos confiar en un modelo?

Feature Leaking

Target Leaking

Sesgo

Sin interpretación clara

Modelo con **muy buenas métricas**

[Incluso en el conjunto de prueba más riguroso]

# ¿Cómo sabemos que podemos confiar en un modelo?

Modelo de clasificación: ¿Husky o Lobo?



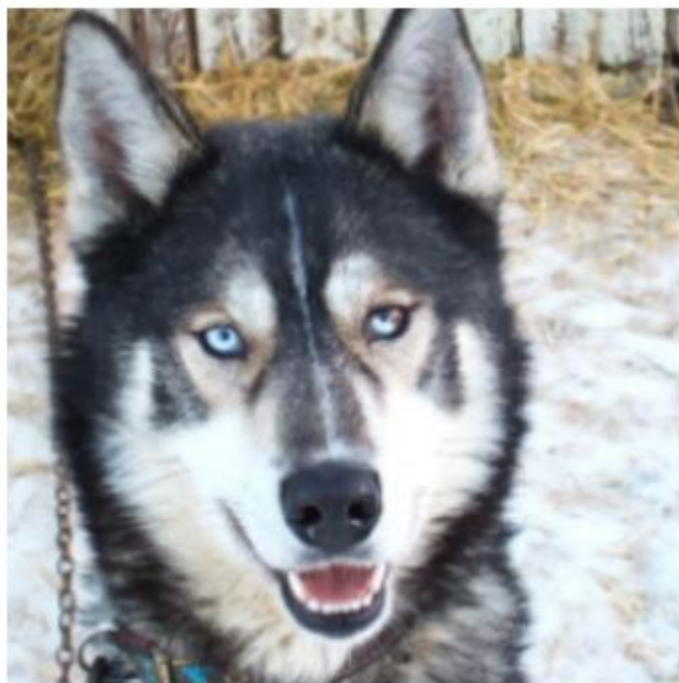
Resultados sorprendentes:

- Métricas de evaluación (precisión, exactitud, exhaustividad, etc.) casi perfectas.
- Todo corroborado en el conjunto de prueba.
- Modelo funcionaba de maravilla... pero comenzó a fallar rápidamente...

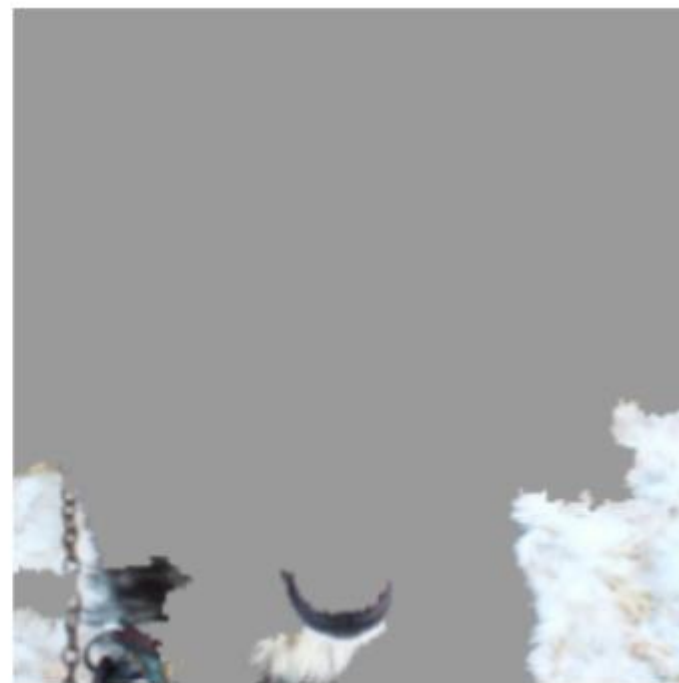


# ¿Cómo sabemos que podemos confiar en un modelo?

Falsos positivos...



(a) Husky classified as wolf



(b) Explanation

# Algunas técnicas...

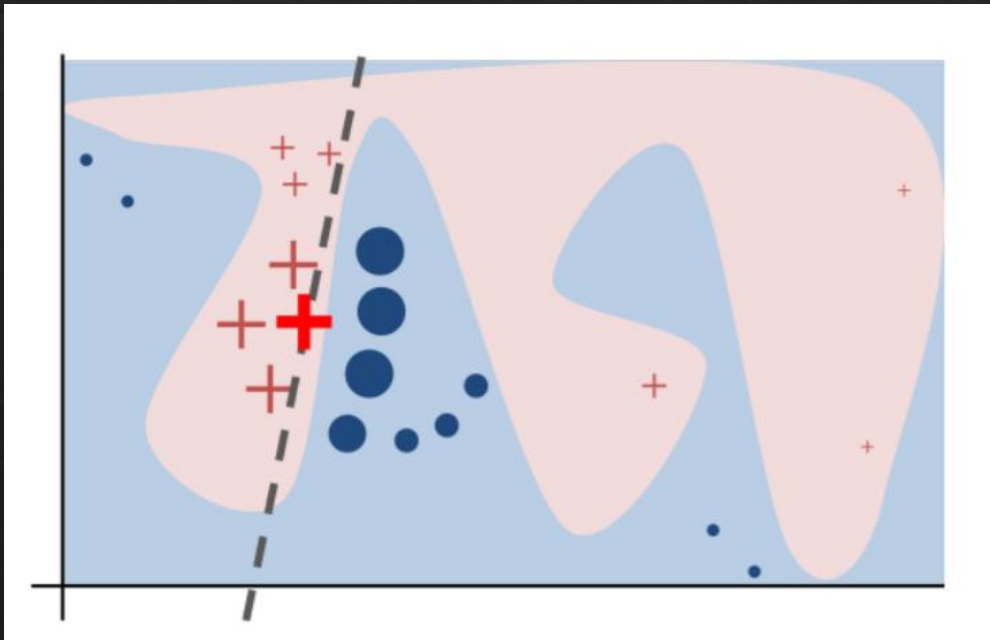
- **LIME** (Explicación local de modelos agnósticos interpretables)
- **ANCHORS** (Explicaciones de alta precisión de modelos agnósticos)
- **SHAP** (Explicaciones aditivas de Shapley)

Buscan explicar el comportamiento de los modelos a través de las variables.

LIME

# LIME (modelo sustituto local)

- LIME: aproxima un modelo interpretable (regresión lineal) de manera local alrededor de una predicción.



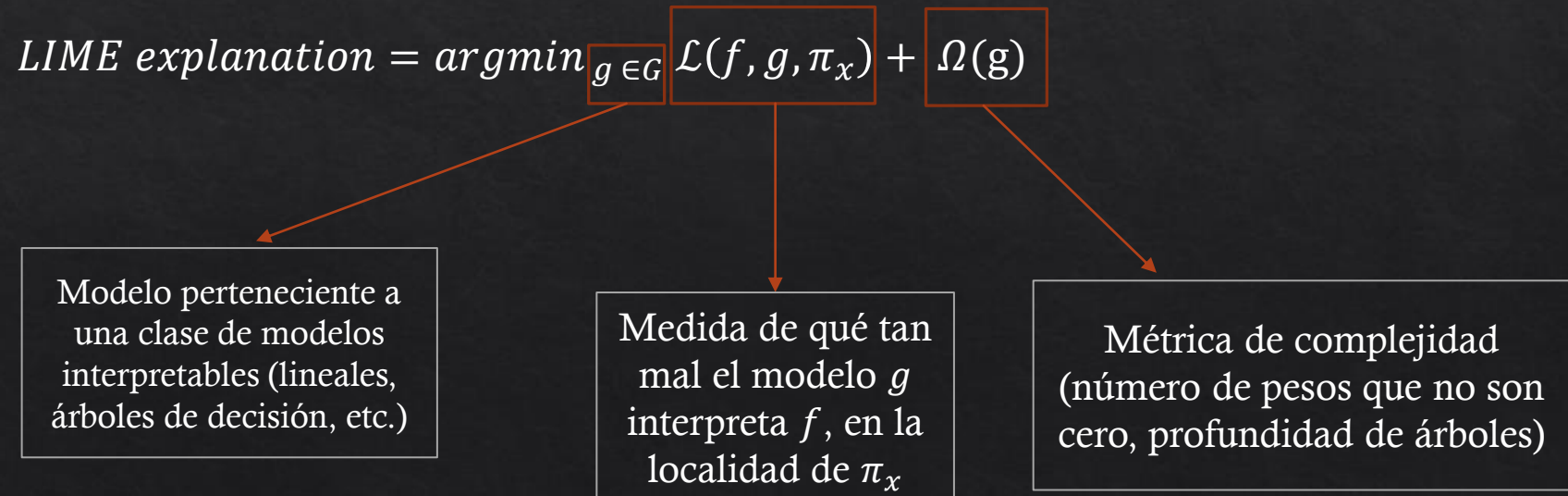
Variables que son importantes de manera local puede que no sean importantes de manera global (vice versa)

$$LIME\ explanation = argmin_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$



# LIME (modelo sustituto local)

**Objetivo formal:** Identificar un modelo interpretable en el espacio de variables interpretación representativa, que es confiable de manera local, para un clasificador.  
Asegurando interpretabilidad y fidelidad local



# LIME (modelo sustituto local)

$$\mathcal{L}(f, g, \pi_x) = \sum_{z, z' \in Z} \pi_x(z) (f(z) - g(z'))^2$$

$$\pi_x(z) = e^{-D(x, z)^2 / (\sigma^2)}$$

La receta...

- Selecciona la instancia de interés: datos de interés para la cual se desea obtener una explicación.
- Perturbar esa instancia y hacer que el modelo haga predicciones para esas nuevas instancias.
- Ponderar las nuevas instancias de acuerdo a la proximidad que tienen con la instancia de interés.
- Entrenar un modelo interpretable en estos conjunto de datos generado a partir de las instancias.
- **Explicar la predicción interpretando el modelo local.**

# LIME (modelo sustituto local)



		edible	poisonous	Feature	Value
gill-size=broad	odor=foul		0.26	odor=foul	True
		0.13		gill-size=broad	True
	stalk-surface-abo...		0.11	stalk-surface-above-ring=silky	True
	spore-print-color=...		0.08	spore-print-color=chocolate	True
	stalk-surface-bel...		0.06	stalk-surface-below-ring=silky	True

		atheism	christian	Text with highlighted words
Posting	Host	0.15		From: johncbad@triton.unm.edu (jchadwic)
	Host	0.14		Subject: Another request for Darwin Fish
	NNTP	0.11		Organization: University of New Mexico, Albuquerque
	edu	0.04		Lines: 11
	have	0.01		NNTP-Posting-Host: triton.unm.edu
	There	0.01		
				Hello Gang,
				There have been some notes recently asking where to obtain the DARWIN fish.
				This is the same question I have and I have not seen an answer on the net. If anyone has a contact please post on the net or email me.

# LIME (modelo sustituto local)

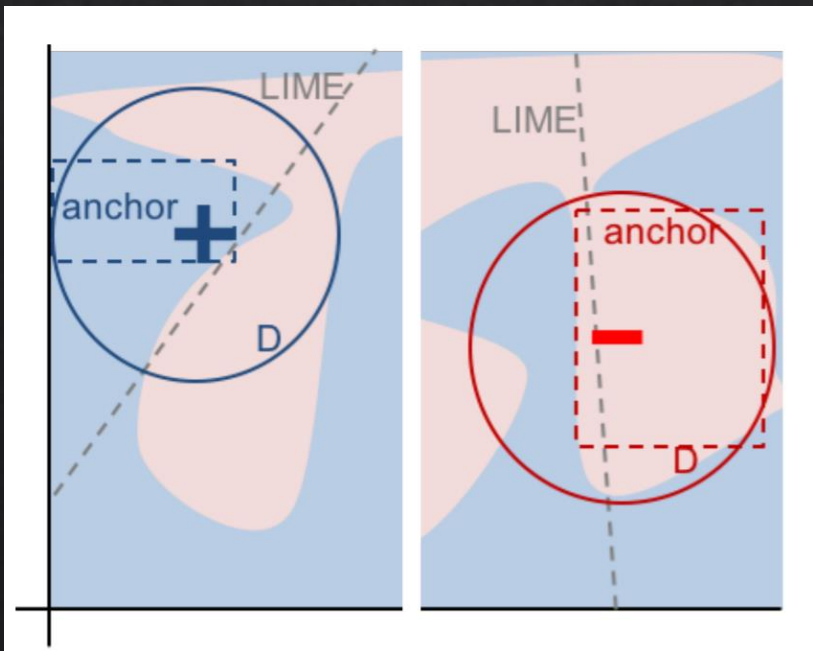
Ventajas	Desventajas
<p>Sirve para debuggear modelos de Aprendizaje Máquina</p> <p>LIME es uno de los pocos métodos que utilizan datos tabulares, texto, e imágenes.</p> <p>Es extremadamente fácil de usar (librerías en R y Python)</p>	<p>La definición de vecindad de perturbación no es tan sencilla de hacer.</p> <p>Inestabilidad de explicación: la confianza no es tan buena.</p> <p>Es algo lento (en particular para imágenes)</p>



Anchor

# Anchor

- Esta técnica explica un modelo generando una regla de decisión (anchor), tal que, dado a que se cumple esa regla, la predicción no se verá modificada sustancialmente.



```
IF (country == 'United_states' AND  
    capital_loss < 1000 AND  
    race == 'white' AND  
    relationship == 'Husband' AND  
    28 < Age < 37 AND  
    sex == 'male' AND  
    occupation == 'Data Scientist')
```

```
THEN PREDICT salary > $50K
```

# Anchor

**Definición formal:**  $A$  es una regla (conjunto de sentencias) que actúa sobre variables tal que  $A(x) = 1$  si todas las sentencias se cumplen para la instancia  $x$ . Entonces se puede definir que  $A$  es un ancla si  $A(x) = 1$  y  $A$  es una condición suficiente para  $f(x)$ .

$$\max_{A \text{ s.t. } P(\text{prec}(A) \geq \tau) \geq 1 - \delta} \text{cov}(A).$$

$$\text{prec}(A) = \mathbb{E}_{\mathcal{D}(z|A)} [\mathbb{1}_{f(x)=f(z)}]$$

Precisión: Proporción de valores correctos, de las instancias generadas, dentro del espacio del ancla.

$$\text{cov}(A) = \mathbb{E}_{\mathcal{D}(z)} [A(z)]$$

Cobertura: Proporción de valores dentro de la población de perturbación a los cuales les aplica el ancla.

# Anchor

Ventajas	Desventajas
<p>Reglas permiten identificar las predicciones de instancias no vistas.</p> <p>Explicaciones son confiables por diseño.</p> <p>Las explicaciones especifican para cuales instancias son válidas.</p> <p>Es extremadamente fácil de usar (librerías en R y Python)</p>	<p>Dado el gran espacio de variables, más de un ancla puede aplicar a la misma instancia (esto podría generar menor interpretabilidad)</p> <p>Generar la distribución de perturbación es difícil.</p> <p>Sigue siendo una explicación local (aunque en realidad, te dice la cobertura)</p>



SHAP

# SHAP (Shapley Additive exPlanations)

- Técnica que emplea teoría de juegos para explicar el comportamiento del modelo.

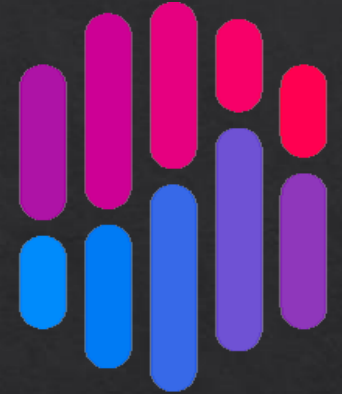
**Valores Shapley:** Método de teoría de juegos (de coaliciones) que te permite determinar cómo repartir la “paga” de manera justa entre todos los jugadores.



Una predicción se puede explicar asumiendo que cada variable es un JUGADOR en un JUEGO (tarea de predecir) donde la predicción es la “paga”.

# SHAP (Shapley Additive exPlanations)

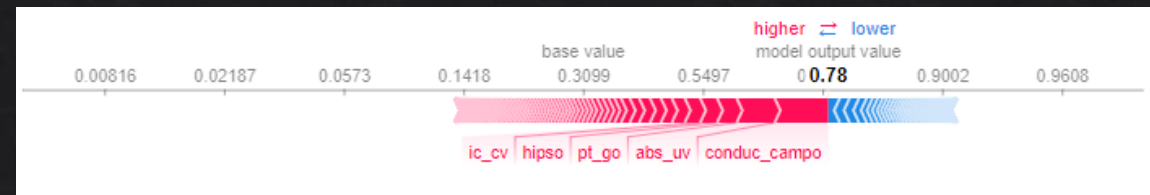
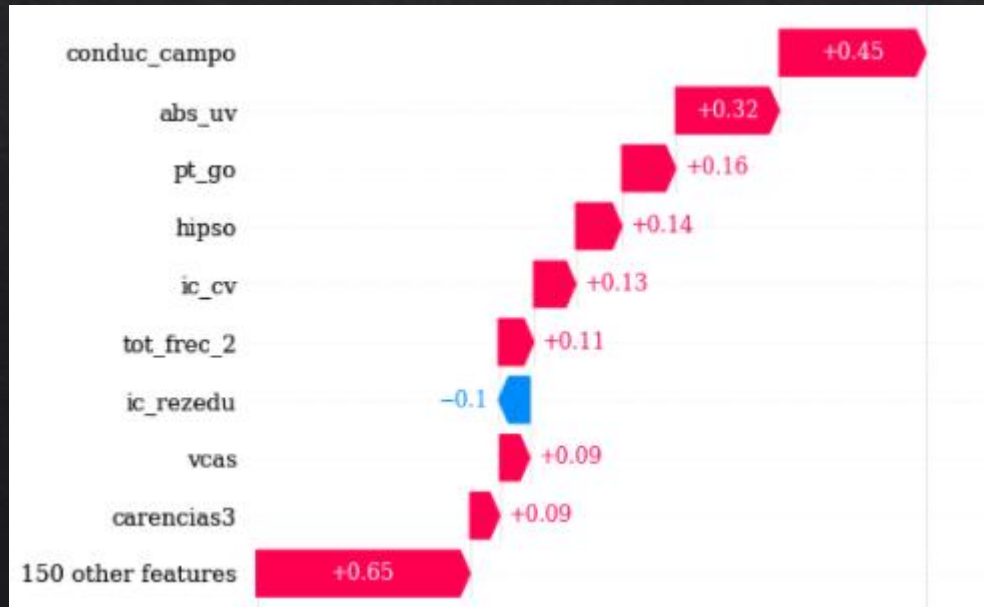
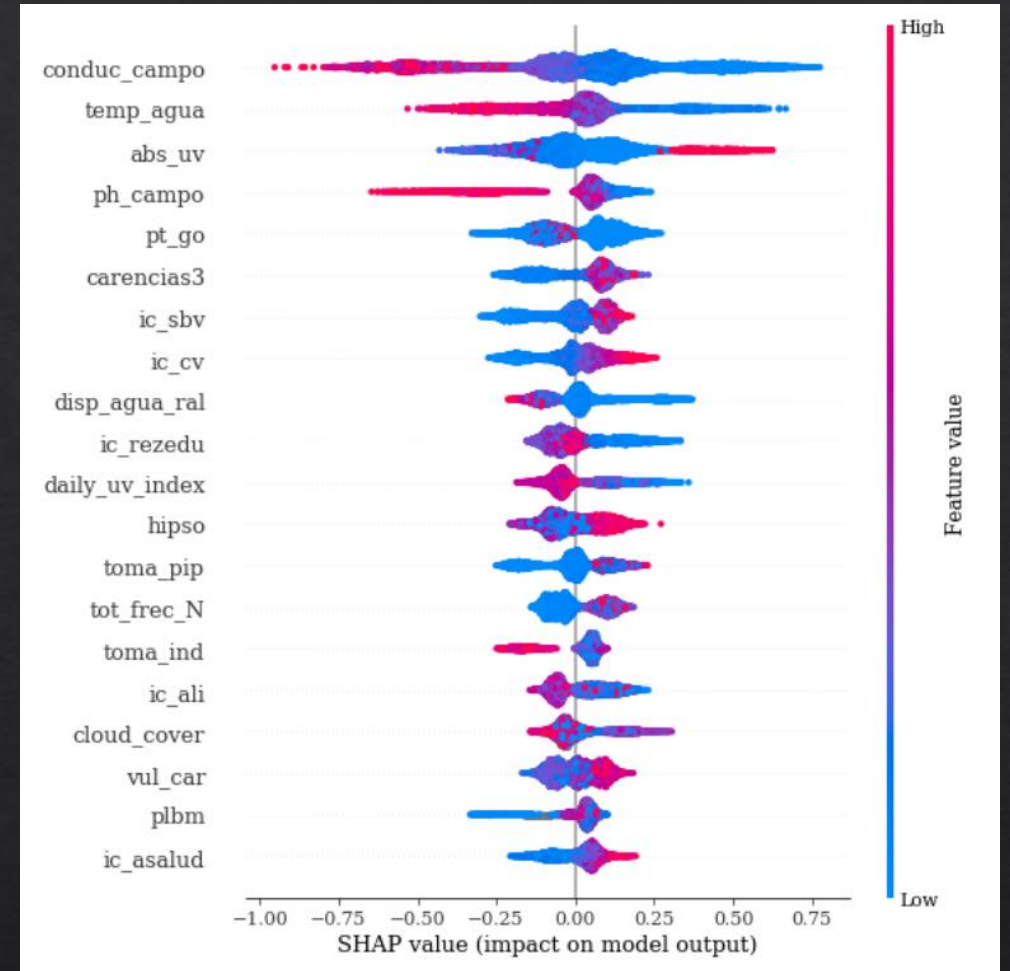
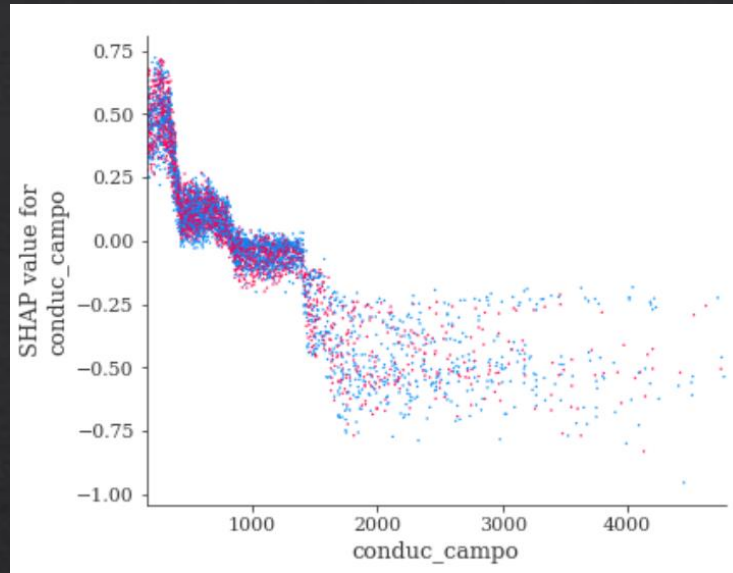
SHAP se basa en la magnitud de las contribuciones marginales promedio de las variables a lo largo de muchas muestras.



$$\phi_j(val) = \sum_{S \subseteq \{x_1, \dots, x_p\} \setminus \{x_j\}} \frac{|S|! (p - |S| - 1)!}{p!} (val(S \cup \{x_j\}) - val(S))$$

En realidad es la contribución a la paga (sumado y ponderado) entre todas las posibles combinaciones, para así detectar la contribución marginal de las variables.

# SHAP





# SHAP

Ventajas	Desventajas
Distribución Equitativa por justicia (solución única demostrada por teoremas)	Obtener valores Shapley requiere mucho tiempo computacional
Explicación completa y comprensiva por cada variable	Se deben de usar aproximaciones para atenuar ese tiempo.
Explicaciones locales y globales	Siempre se requiere información para calcular el valor shapley (no hay perturbación)
Implementación optimizada para árboles (XGBOOST, CatBoost LGBM)	Explicabilidad no es multivariada.
Es extremadamente fácil de usar (librerías en R y Python)	

# CONCLUSIÓN:

Explicar un modelo de Aprendizaje Automático brinda confianza, seguridad, e interpretación a las decisiones que está realizando el modelo para predecir.

Ayuda a encontrar:

- Patrones y sesgos en los datos.
- Importancia en las variables (target/feature leaking).
- La razón de la predicción.

# Vamos a programar...

◆ [https://github.com/dhdzmota/explaining\\_ml\\_ironhack](https://github.com/dhdzmota/explaining_ml_ironhack)

# Referencias

Tulio, M., Singh, S., y Guestrin, C. (2016) Why should I trust you? Explaining the predictions of any classifier. Recuperado el 17 de octubre del 2020 de: <https://arxiv.org/abs/1602.04938>

Tulio, M. (2020) Lime: Explaining the predictions of any machine learning classifier. GitHub. Recuperado el 17 de octubre del 2020 de : <https://github.com/marcotcr/lime>

Tulio, M. (2016) Lime 0.2.0.1. GitHub. Recuperado el 01 de Noviembre del 2020 de : <https://pypi.org/project/lime/#history>

Sharma, A. (2018) Decrypting your Machine Learning model using LIME. Recuperado el 01 de noviembre del 2020 de: <https://towardsdatascience.com/decrypting-your-machine-learning-model-using-lime-5adc035109b5>

Korobov, M. y Lopuhin, K., (2016) ELI5. Recuperado el 17 de octubre del 2020 de: <https://eli5.readthedocs.io/en/latest/overview.html>

Korobov, M. y Lopuhin, K., (2016) ELI5. Recuperado el 17 de octubre del 2020 de: <https://github.com/TeamHG-Memex/eli5>

Tulio, M., Singh, S y Guestrin C. (2018) Anchors: High-Precision Model-Agnostic Explanations. Recuperado el 17 de octubre del 2020 de: <https://homes.cs.washington.edu/~marcotcr/aaai18.pdf>

Tulio, M. (2019) Anchor. Github. Recuperado el 17 de octubre del 2020 de: <https://github.com/marcotcr/anchor>

Lundber, S. Su-In, L. (2017) A unified approach to interpreting model predictions. Advances in Neural Information Processing Systems. Recuperado el 17 de octubre.

Shapley, L., (1953) A value for n-person games. Contributions to the Theory of Games 2.28:

Molnar, C. (2020) Interpretable Machine Learning. Recuperado el 18 de septiembre del 2020 de: <https://christophm.github.io/interpretable-ml-book/shap.html>