# Sistemas Recomendadores IIC-3633

## Explicabilidad en Sistemas Recomendadores

# Esta clase

1. Explainable AI
2. Explicabilidad en sistemas de recomendación

# Motivación de XAI

Salud / Tamizaje
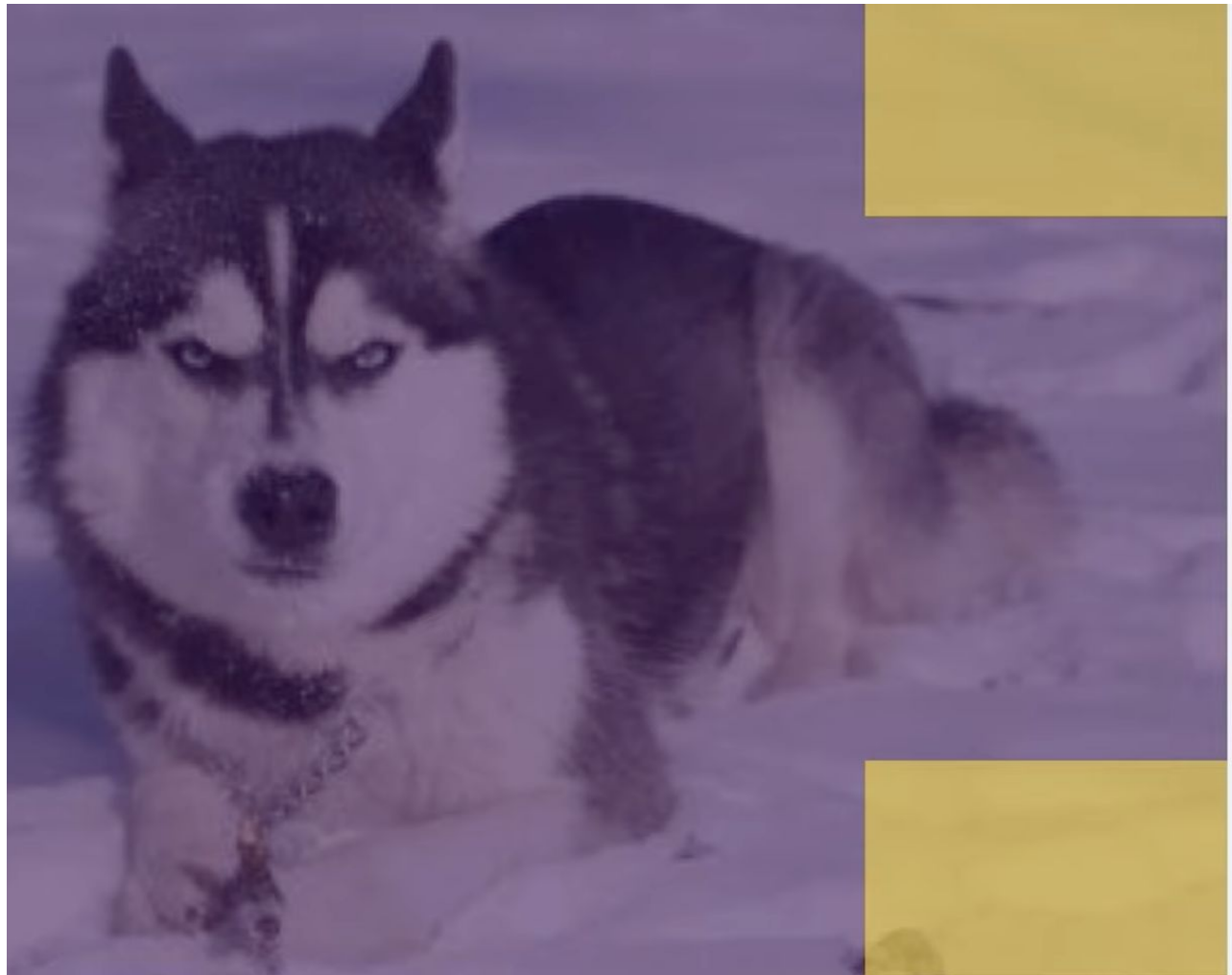
Predicción de a quien dar un crédito bancario y por qué?

Husky o Lobo?

Área de la imagen en la que se fijó el modelo para clasificarlo como lobo.

Toma la decisión correcta pero basado en malas razones (la nieve)

El modelo debería prestar atención a esta zona para determinar si es un lobo o no.
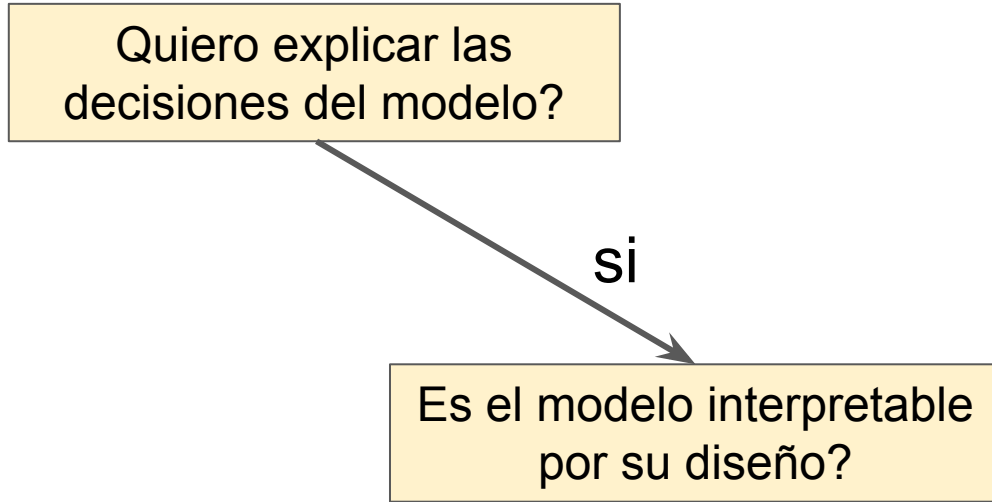
inicio

Quiero explicar las decisiones del modelo?

inicio

Quiero explicar las decisiones del modelo?

si

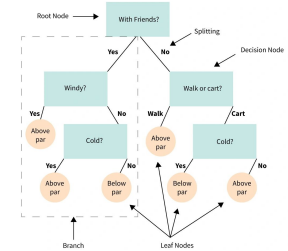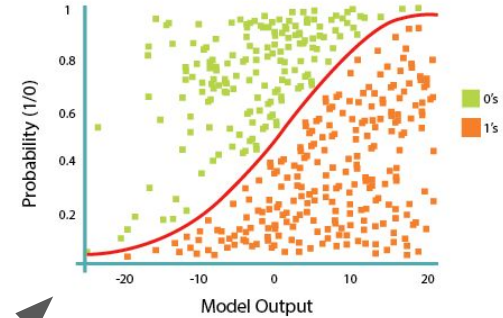Es el modelo interpretable por su diseño?
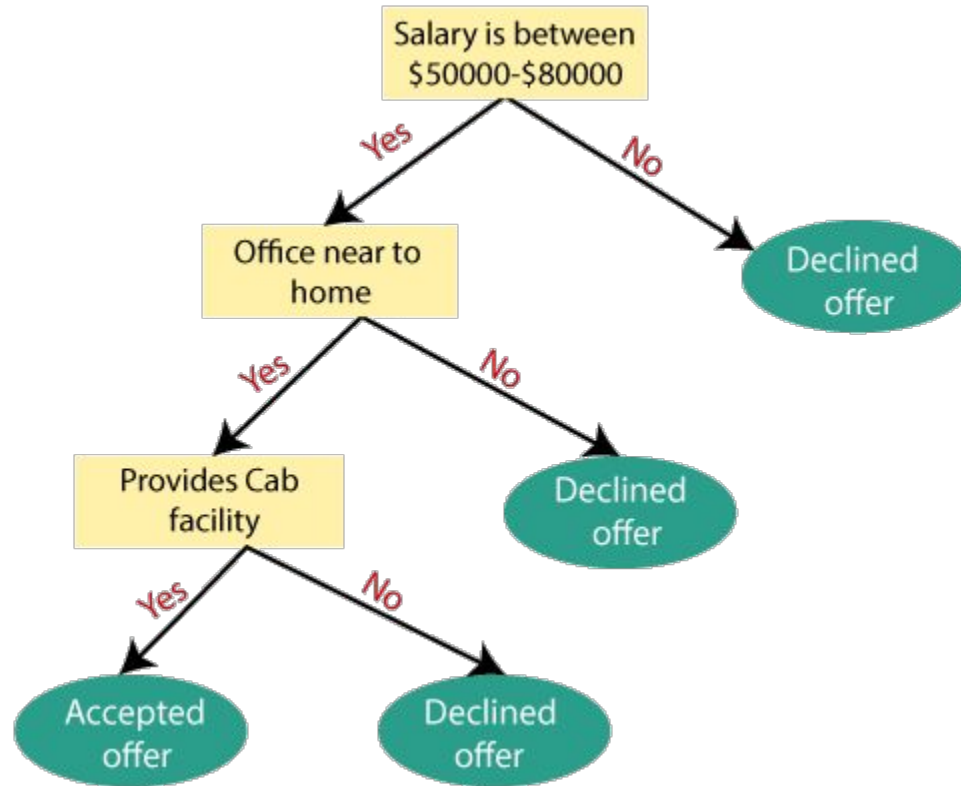
inicio

Quiero explicar las decisiones del modelo?

KNN , LR , DECISION TREE LINEAR

si

si

Es el modelo interpretable por su diseño?

Muy bueno para hacerse una idea de los features más importantes a la hora de hacer inferencia !!
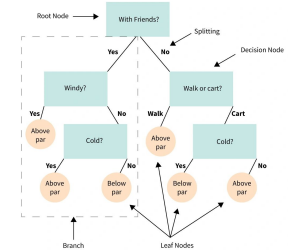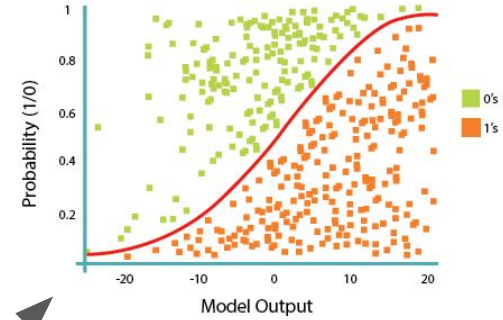
inicio

Quiero explicar las decisiones del modelo?

KNN , LR , DECISION TREE LINEAR

si

si

Es el modelo interpretable por su diseño?

no

¿Necesita de otro modelo para interpretarlo?

inicio

Quiero explicar las decisiones del modelo?

**si**

KNN , LR ,
DECISION
TREE
LINEAR
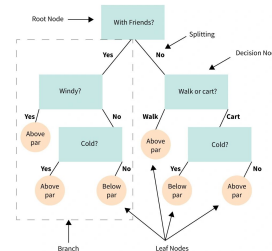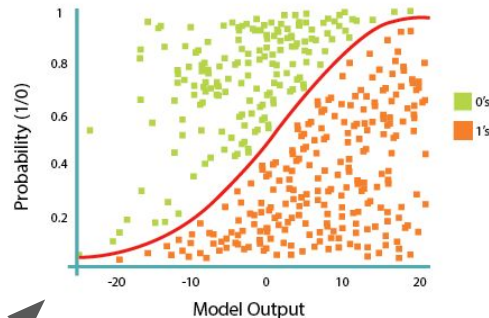


Es el modelo interpretable por su diseño?

**si**



**no**

¿Necesita de otro modelo para interpretarlo?

**si**

Model Agnostic Explanations
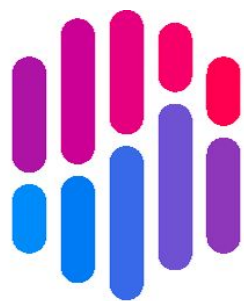
Example based
(Counterfactual, Adversarial,
Influence)

# Model-agnostic explanations

# LIME



Figure 3: Toy example to present intuition for LIME. The black-box model's complex decision function $f$ (unknown to LIME) is represented by the blue/pink background, which cannot be approximated well by a linear model. The bold red cross is the instance being explained. LIME samples instances, gets predictions using $f$, and weighs them by the proximity to the instance being explained (represented here by size). The dashed line is the learned explanation that is locally (but not globally) faithful.

# LIME



(a) Original Image    (b) Explaining *Electric guitar*    (c) Explaining *Acoustic guitar*    (d) Explaining *Labrador*

Modelo predijo 3 clases:
- Guitarra electrica (**p**=0.32) → pixeles del mástil de la guitarra
- Guitarra acustica (**p**=0.24) → pixeles del cuerpo de la guitarra
- Labrador (**p**=0.21) → cara del perro labrador.

# Example-based explanations

# Counterfactual

¿Qué tendría que cambiar para obtener un resultado diferente?

**Ejemplo:**

Si tu ingreso anual hubiera sido $X más alto, tu solicitud de crédito habría sido aprobada.

Así podemos encontrar si hay sesgos del modelo por ciertas características.

stop sign
Confidence: 0.9153

Adversarial perturbation

flowerpot
Confidence: 0.8374

Ejemplos adversarios:

entradas perturbadas imperceptibles que induzcan a errores en un modelo

stop sign: 99%

sports ball: 80%

# Funciones de influencia (Influence functions)

Las funciones de influencia se utilizan para analizar la sensibilidad de un estimador (por ejemplo, la media o la mediana) ante perturbaciones en el conjunto de datos.

inicio

Quiero explicar las decisiones del modelo?

KNN , LR , DECISION TREE LINEAR



si

Es el modelo interpretable por su diseño?

si



no

- attention
- gradient saliency
- integrated gradients

¿Necesita de otro modelo para interpretarlo?

si

Model Agnostic Explanations

SHAP

no

# Attention

Estudiar la atencion del modelo al generar inferencia sobre nuevos ejemplos.

# Gradient saliency maps

Estudiar cuáles partes de la entrada (ej. imagen) tienen mayor influencia en el gradiente de la red neuronal.

# Integrated gradients

Se elige un punto de referencia neutro, como una imagen en blanco o negro.

Luego, se traza un camino desde ese punto hasta la entrada real.

A lo largo de este camino, se suman los gradientes del modelo para determinar la importancia de cada característica.

|  | Gradient | | Integrated Gradients | |
|---|---|---|---|---|
|  | Standard training | Adv trained L2=4 | Standard training | Adv trained L2=4 |
| Melanoma | | | | |
| Benign keratosis | | | | |
| Melanocytic Nevi | | | | |

inicio

Quiero explicar las decisiones del modelo?

KNN , LR , DECISION TREE LINEAR



**no**

entender cómo aprende el modelo
ej. feature visualization.

**si**

Es el modelo interpretable por su diseño?

**si**



- attention
- gradient saliency (feat importance)
- integrated gradients

**no**

**no**

¿Necesita de otro modelo para interpretarlo?

**si**

Model Agnostic Explanations

SHAP

**si**

Example-based

counterfactual adversarial

# Feature visualization



| Edges | Textures | Patterns | Parts | Objects |

Explorar que aprende cada una de las neuronas.

fuente: https://distill.pub/2017/feature-visualization/

# Explainable Recommendation

La recomendación explicable busca desarrollar modelos que no solo generen recomendaciones de alta calidad, sino también explicaciones intuitivas.

Las explicaciones pueden ser

- post-hoc
- provenir directamente de un modelo interpretable o transparente

Aborda el problema del por qué, ayudando a los humanos (usuarios o diseñadores de sistemas) a entender por qué ciertos ítems son recomendados por el algoritmo.

# Beneficios

- Mejora la transparencia de los sistemas de recomendación.
- Aumenta la persuasión y efectividad de las recomendaciones.
- Fomenta la confianza y satisfacción del usuario.
- Facilita a los diseñadores del sistema una mejor depuración.

**Bob**

⭐⭐⭐⭐⭐ Not a bad price and it works. Can't do much more than that. We'll see how long it takes me to wear it out.

⭐⭐⭐⭐⭐ For the price, I would definitely buy again. It's sturdily constructed, bright, and feels good in hand.

⭐⭐⭐⭐⭐ I love this camera. It is amazing. It gives professional quality. I am still learning all the excellent features. The more I learn, the better I love this camera.

**Helen**

⭐⭐⭐⭐⭐ Love the watch, but the delivery was WAY after the original predicted date and that was disappointing since it was a prime item and should have been two day max.

⭐☆☆☆☆ Continuously shuts down. Numerous errors and issues.

**William**

⭐⭐⭐⭐⭐ Pretty good performance in night, especially when recording videos. Good choice for a starter.

⭐⭐⭐⭐⭐ Great lens for its price sure it doesn't have phase detection but the a7ii has it built in and with the update the autofocus is pretty fast once you calibrate it to your camera

**Fred**

⭐⭐⭐⭐⭐ I have been wanting a new camera for years. This one is great for me

⭐⭐⭐⭐⭐ This lens is great and perfectly functional. Be sure to take off both protective films from the front and back of the lens. The focus doesn't pull as cleanly as the canon brand 50mm but it is a perfect intermediate lens when you are on a budget.

⭐⭐⭐⭐⭐ For the price, I would definitely buy again. It's sturdily constructed, bright, and feels good in hand.

---

**Traditional explanation**

Recommend →

**User based explanation**

The lens is recommended to you, because your similar user William and Fred have bought this item before.

**Item based explanation**

The lens is recommended to you, because you bought a camera before.

---

**Textual explanation**

Recommend →

**Feature-level explanation**

| Feature | likeness |
|---|---|
| color | 0.87 |
| quality | 0.54 |
| Focal Length | 0.66 |
| Focus Type | 0.71 |

**Sentence-level explanation**

Structured: You might be interested in [feature] (can be quality, color, etc), on which this product performs well.

Unstructured: Great and deserve the price.

---

**Visual explanation**

Recommend →

**Visual explanation**

(a)

(b)

**Your rating for similar movies**

**Your neighbors' rating for this movie**

| Rating | Number of Neighbors |
|---|---|
| ★★★★★ | 2 |
| ★★★★ | 3 |
| ★★★ | 4 |
| ★★ | 0 |
| ★ | 0 |

(a)

Your recommendation is based on how Movielens thinks you like the following aspects

**Relevance**　　　　**Your Preference**

| | | |
|---|---|---|
| Wes Anderson | | ★★★★★ |
| Deadpan | | ★★★★ |
| Quirky | | ★★★★ |
| Witty | | ★★★★ |
| Off-beat comedy | | ★★★★ |
| Notable soundtrack | | ★★★★★ |
| Stylized | | ★★★★ |

(b)

**Amazon Books**

**Recommended**

This is well written book with a very good detail of a person that love his dog but didn't restrain his freedom. I think that I relive the joy of my experiences with my dogs, a Labrador and a Siberian Husky. Both were rescued, one from the shelter and the other from the street. After 4 months with me, the owner of the husky appeared and I returned the dog. Two weeks later the dog escaped and returned to my house. He decided who will be his owner. The author described with details the relationship of them, concerns, disappointments and health issues. The final chapter was a surprise that I am still enjoying.

Who are similar to you [u1]

Similarity to u1 | Friend? | Likes

u2 — Yes — c

u3 — No — b

u4 — No — b

u5 — No — a

We recommend these.

Why is u5 similar to you?

- Likes the things you[u1] like:

  d  e  g

- Dislikes the things you dislike:

  h  i

Why is u2 similar to you?

- Likes the things you[u1] like:

  j  k

- Common friends:

  u6

**Behavior triplets**

(User A, Item A, buy)
(User B, Item C, buy)
(User A, Word A, mention_word)
(User A, Word B, mention_word)
(Item D, Word B, mention_word)
(User B, Word A, mention_word)
(Item D, Word C, mention_word)
(Item C, Word C, mention_word)
(Item A, Item B, also_bought)
(Item B, Item C, also_view)
(Item A, Category A, belong_to_category)
(Item B, Category A, belong_to_category)
(Item C, Category B, belong_to_category)
(Item B, Brand A, belong_to_brand)
(Item C, Brand B, belong_to_brand)

**User behavior Graph**

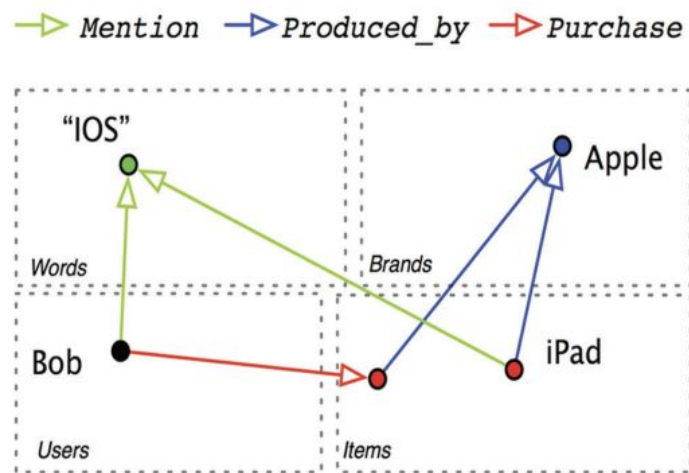| | |
|---|---|
| → | buy |
| ·····► | mention |
| → | also_view |
| → | belong_to_category |
| ·····► | belong_to_brand |
| → | also_bought |

(a) A knowledge graph of users and items
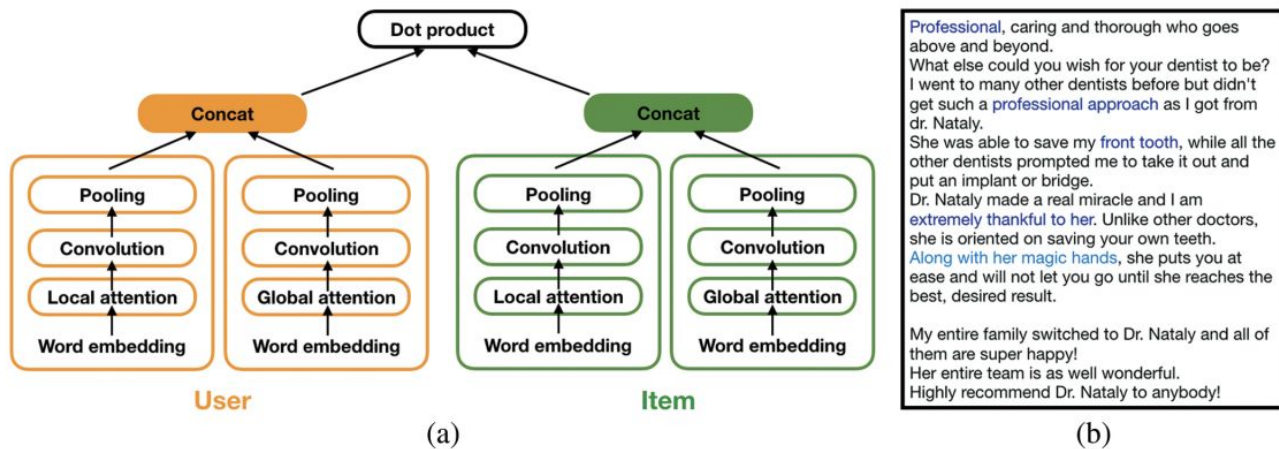
(b) Extracting explanation path

**Figure 3.8:** (a) The dual-attention architecture to extract user and item representations. A user document and an item document are fed into the user network (left) and item network (right). (b) The model generates attention scores for each review and highlights the high attention words as explanations (Seo *et al.*, 2017).