# CLASIFICACIÓN Y REGRESIÓN EN DATASET DE PELÍCULAS

Marco
Pérez González

marco.perez
@alumnos.upm.es

David
Cano Rosillo

david.cano.rosillo
@alumnos.upm.es

José Javier
Román Camacho

josejavier.roman
@alumnos.upm.es

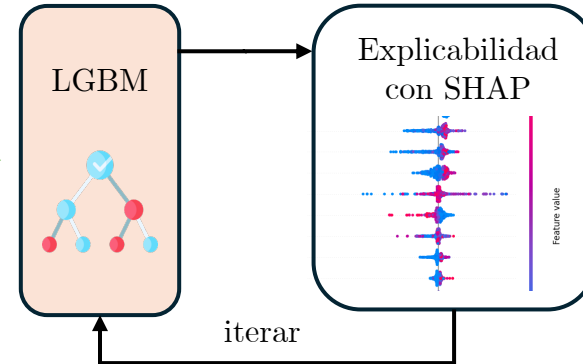# CLASIFICACIÓN Y REGRESIÓN EN DATASET DE PELÍCULAS

| Datos | Procesado | Métricas |
|---|---|---|



**Datos tabulares**

2896 películas

**numéricos**
duration, content_rating.

**categóricos**
color, language, country...

**texto**
keywords, overview,

LGBM

Explicabilidad con SHAP

iterar

output

### Predicción de score

| | CV | TEST |
|---|---|---|
| MAE | **0.24** | **0.22** |
| MSE | **0.13** | **0.09** |

**Pósteres**

2894 x .jpg

**CONCATENAR**
(keywords, overview, plot_keywords)

frecuencia

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Vector DB

Regresión logística

$\mathcal{P}(war) = 0.6$
$\mathcal{P}(crime) = 0.3$
$\mathcal{P}(drama) = 0.2$

LGBM

OPTUNA ⚡⚡

### Predicción de género

| | CV | TEST |
|---|---|---|
| ACC | 0.92 | 0.27 |

**Imágen**
matrix 3 x 224 x 244

$\mathcal{P}(war) = 0.2$
$\mathcal{P}(crime) = 0.4$
$\mathcal{P}(drama) = 0.2$

Pre-trained resnet + finetuning

Predicción de score

| | CV | TEST |
|---|---|---|
| MAE | 0.24 | 0.22 |
| MSE | 0.13 | 0.09 |

numéricos
duration, content_rating.

categóricos
color, language, country...

LGBM

Explicabilidad con SHAP

iterar

output

Datos tabulares
2896 películas

# Ingeniería de características

- Variables asociadas a los directores (LOO target encoding)
- Variables asociadas a longitudes (del reparto, descripciones, etc)
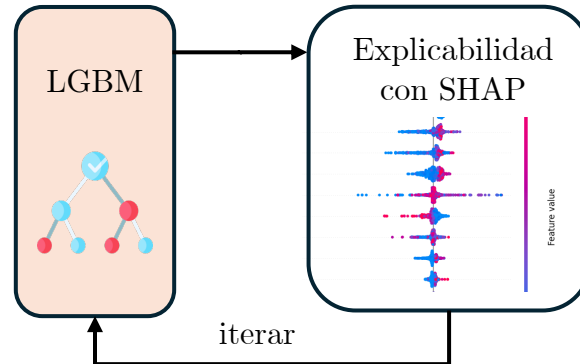- Proporciones de likes en facebook (película, actores, etc)

**Predicción de score**

| | CV | TEST |
|---|---|---|
| MAE | 0.24 | 0.22 |
| MSE | 0.13 | 0.09 |

numéricos
duration, content_rating.

categóricos
color, language, country...

LGBM

Explicabilidad con SHAP

output

iterar

Datos tabulares
2896 películas

color | duration | vote averga | keywords

# Ingeniería de características

- Variables asociadas a los directores (LOO target encoding)
- Variables asociadas a longitudes (del reparto, descripciones, etc)
- Proporciones de likes en facebook (película, actores, etc)

Training Sets | Test Set

Iteration 1 → $Error_1$
Iteration 2 → $Error_2$
Iteration 3 → $Error_3$
Iteration 4 → $Error_4$
Iteration 5 → $Error_5$

$$Error = \frac{1}{5}\sum_{i=1}^{5} Error_i$$

vote_average
num_critic_for_reviews
num_voted_users
title_year
gross
num_user_for_reviews
content_rating
len_keywords
budget
loo_director_imdb_score
popularity

High
Low

Feature value

SHAP value (impact on model output)

# CLASIFICACIÓN Y REGRESIÓN EN DATASET DE PELÍCULAS

| Datos | Procesado | Métricas |
|---|---|---|



**Datos tabulares**

2896 películas

**numéricos**
duration, content_rating.

**categóricos**
color, language, country...

**texto**
keywords, overview,

LGBM

Explicabilidad con SHAP

iterar

output

**Predicción de score**

| | CV | TEST |
|---|---|---|
| MAE | **0.24** | **0.22** |
| MSE | **0.13** | **0.09** |

**Pósteres**

2894 x .jpg

**Imágen**
matrix 3 x 224 x 244

**CONCATENAR**
(keywords, overview, plot_keywords)

frecuencia

$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

Vector DB

Regresión logística

$\mathcal{P}(\text{war}) = 0.6$
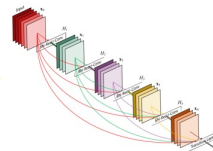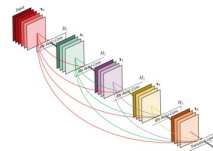$\mathcal{P}(\text{crime}) = 0.3$
$\mathcal{P}(\text{drama}) = 0.2$

$\mathcal{P}(\text{war}) = 0.2$
$\mathcal{P}(\text{crime}) = 0.4$
$\mathcal{P}(\text{drama}) = 0.2$

Pre-trained resnet + finetuning

LGBM

OPTUNA

**Predicción de género**

| | CV | TEST |
|---|---|---|
| ACC | 0.92 | 0.27 |

# CLASIFICACIÓN Y REGRESIÓN EN DATASET DE PELÍCULAS

| Datos | Procesado | Métricas |
|---|---|---|

**numéricos**
duration, content_rating.

**categóricos**
color, language, country...

**texto**
keywords, overview,

Datos tabulares
2896 películas

Pósteres
2894 x .jpg

**Imágen**
matrix 3 x 224 x 244

**CONCATENAR**
(keywords, overview, plot_keywords)

frecuencia

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Vector DB

Regresión logística

$\mathcal{P}(\text{war}) = 0.6$
$\mathcal{P}(\text{crime}) = 0.3$
$\mathcal{P}(\text{drama}) = 0.2$

Pre-trained resnet + finetuning

$\mathcal{P}(\text{war}) = 0.2$
$\mathcal{P}(\text{crime}) = 0.4$
$\mathcal{P}(\text{drama}) = 0.2$

LGBM

OPTUNA
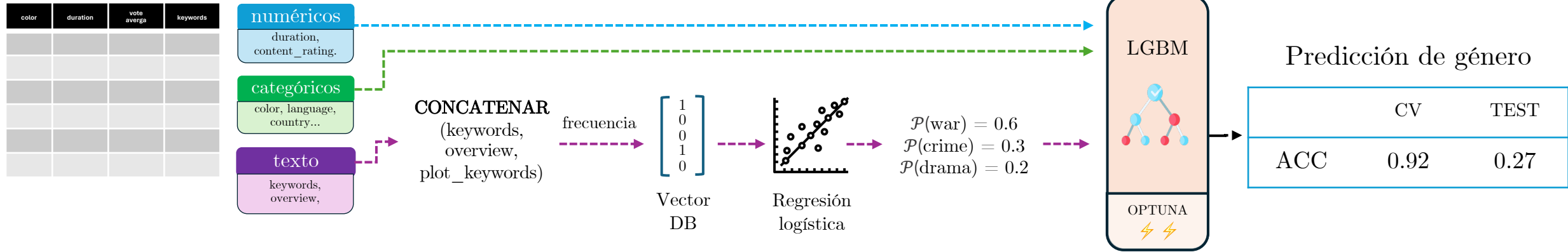⚡⚡

## Predicción de género

| | CV | TEST |
|---|---|---|
| ACC | 0.92 | 0.27 |

**Datos**

**Procesado**

**Métricas**

| color | duration | vote averga | keywords |

**numéricos**
duration, content_rating.

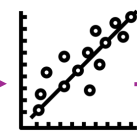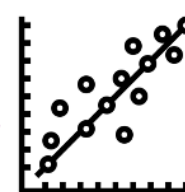**categóricos**
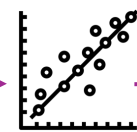color, language, country...

**texto**
keywords, overview,

CONCATENAR
(keywords, overview, plot_keywords)

frecuencia

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Vector DB

Regresión logística

$\mathcal{P}(\text{war}) = 0.6$
$\mathcal{P}(\text{crime}) = 0.3$
$\mathcal{P}(\text{drama}) = 0.2$
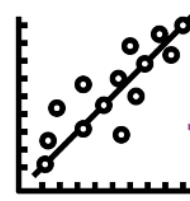
LGBM

OPTUNA
⚡⚡

**Predicción de género**

| | CV | TEST |
|---|---|---|
| ACC | 0.92 | 0.27 |

CONCATENAR
(keywords, overview, plot_keywords)

frecuencia

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Vector DB

Regresión logística

$\mathcal{P}(\text{war}) = 0.6$
$\mathcal{P}(\text{crime}) = 0.3$
$\mathcal{P}(\text{drama}) = 0.2$

| color | duration | vote averga | keywords |

**numéricos**
duration, content_rating.

**categóricos**
color, language, country...

**texto**
keywords, overview,

**CONCATENAR**
(keywords, overview, plot_keywords)

frecuencia

$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

Vector DB

Regresión logística

$\mathcal{P}(\text{war}) = 0.6$
$\mathcal{P}(\text{crime}) = 0.3$
$\mathcal{P}(\text{drama}) = 0.2$

LGBM

OPTUNA
⚡⚡

Predicción de género

|  | CV | TEST |
|------|------|------|
| ACC | 0.92 | 0.27 |

**CONCATENAR**
(keywords, overview, plot_keywords)

frecuencia

$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

Vector DB

Regresión logística

$\mathcal{P}(\text{war}) = 0.6$
$\mathcal{P}(\text{crime}) = 0.3$
$\mathcal{P}(\text{drama}) = 0.2$

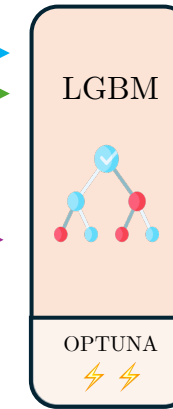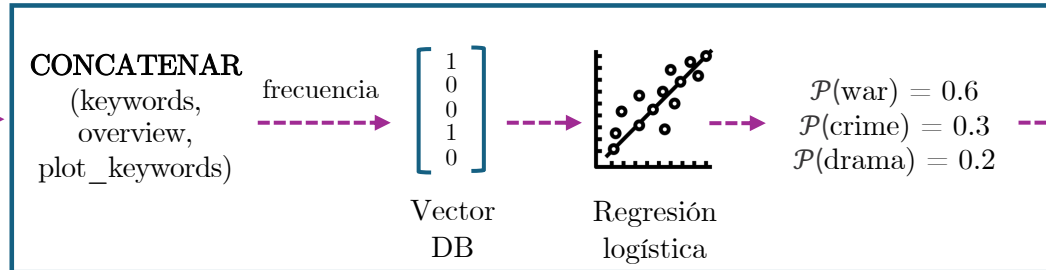## Cosas que no funcionaron

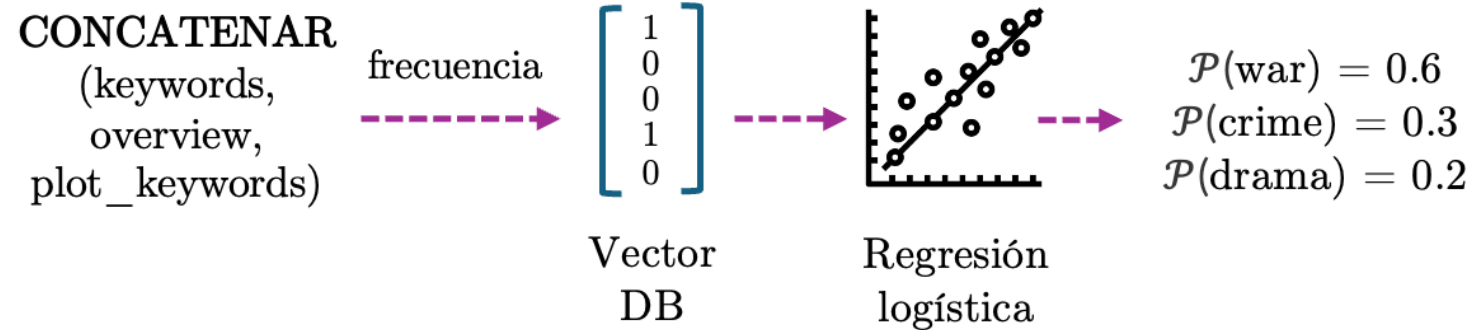- BERT
- CountVectorizer

**Datos**

**Procesado**

**Métricas**

| color | duration | vote averga | keywords |

**numéricos**
duration, content_rating.

**categóricos**
color, language, country...

**texto**
keywords, overview,

LGBM

OPTUNA

CONCATENAR
(keywords, overview, plot_keywords)

frecuencia

$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

Vector DB

Regresión logística

$\mathcal{P}(\text{war}) = 0.6$
$\mathcal{P}(\text{crime}) = 0.3$
$\mathcal{P}(\text{drama}) = 0.2$

**Predicción de género**

| | CV | TEST |
|---|---|---|
| ACC | 0.92 | 0.27 |

CONCATENAR
(keywords, overview, plot_keywords)

frecuencia

$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

Vector DB

Regresión logística

$\mathcal{P}(\text{war}) = 0.6$
$\mathcal{P}(\text{crime}) = 0.3$
$\mathcal{P}(\text{drama}) = 0.2$

| GENRE | AUC |
|---|---|
| Western | 0.990389 |
| War | 0.985437 |
| Horror | 0.956172 |
| Documentary | 0.954239 |
| Music | 0.947796 |
| Science Fiction | 0.946767 |
| History | 0.929853 |
| Animation | 0.923654 |
| Family | 0.917338 |

| GENRE | AUC |
|---|---|
| Crime | 0.910368 |
| Fantasy | 0.890410 |
| Action | 0.888968 |
| Romance | 0.867276 |
| Adventure | 0.866570 |
| Thriller | 0.865806 |
| Comedy | 0.854768 |
| Mystery | 0.824911 |
| Foreign | 0.490550 |

| color | duration | vote averga | keywords |

numéricos
duration, content_rating.

categóricos
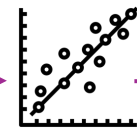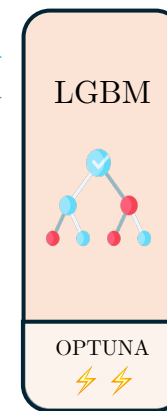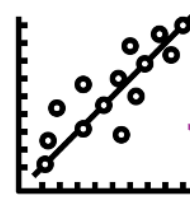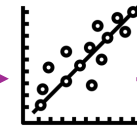color, language, country...

texto
keywords, overview,

**CONCATENAR**
(keywords, overview, plot_keywords)

frecuencia

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Vector DB

Regresión logística

$\mathcal{P}(\text{war}) = 0.6$
$\mathcal{P}(\text{crime}) = 0.3$
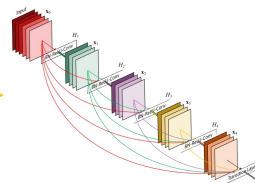$\mathcal{P}(\text{drama}) = 0.2$

LGBM

OPTUNA
⚡⚡

Predicción de género

| | CV | TEST |
|---|---|---|
| ACC | 0.92 | 0.27 |

**CONCATENAR**
(keywords, overview, plot_keywords)

frecuencia

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Vector DB

Regresión logística

$\mathcal{P}(\text{war}) = 0.6$
$\mathcal{P}(\text{crime}) = 0.3$
$\mathcal{P}(\text{drama}) = 0.2$

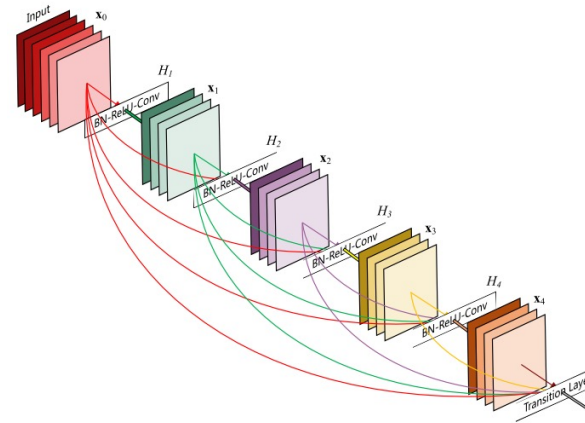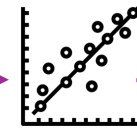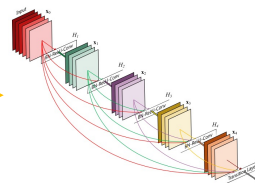" Nuestro objetivo es generar las mejores *features* para LGBM "

# Pre-trained resnet + **finetuning**



Imágen

matrix 3 x 224 x 244

$\mathcal{P}(\text{war}) = 0.2$
$\mathcal{P}(\text{crime}) = 0.4$
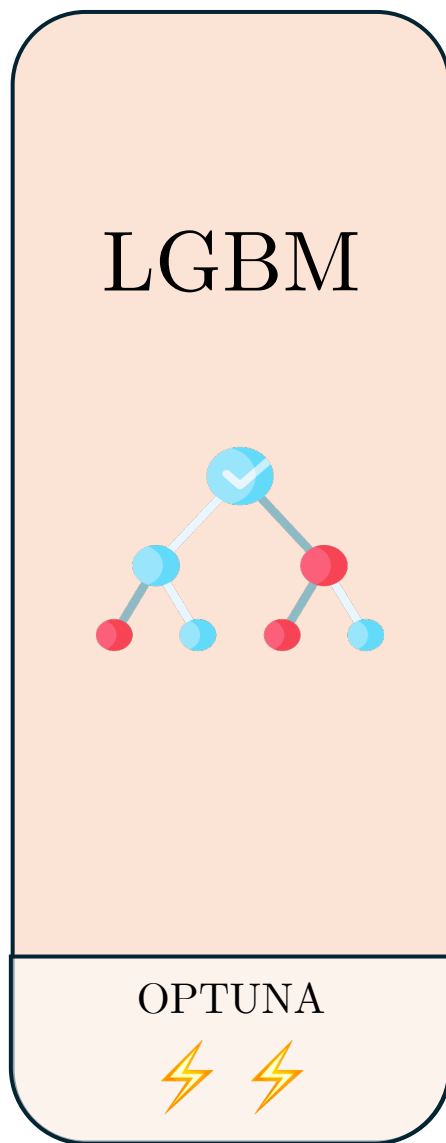$\mathcal{P}(\text{drama}) = 0.2$

LGBM

OPTUNA

$\mathcal{P}(\text{war}) = 0.6$
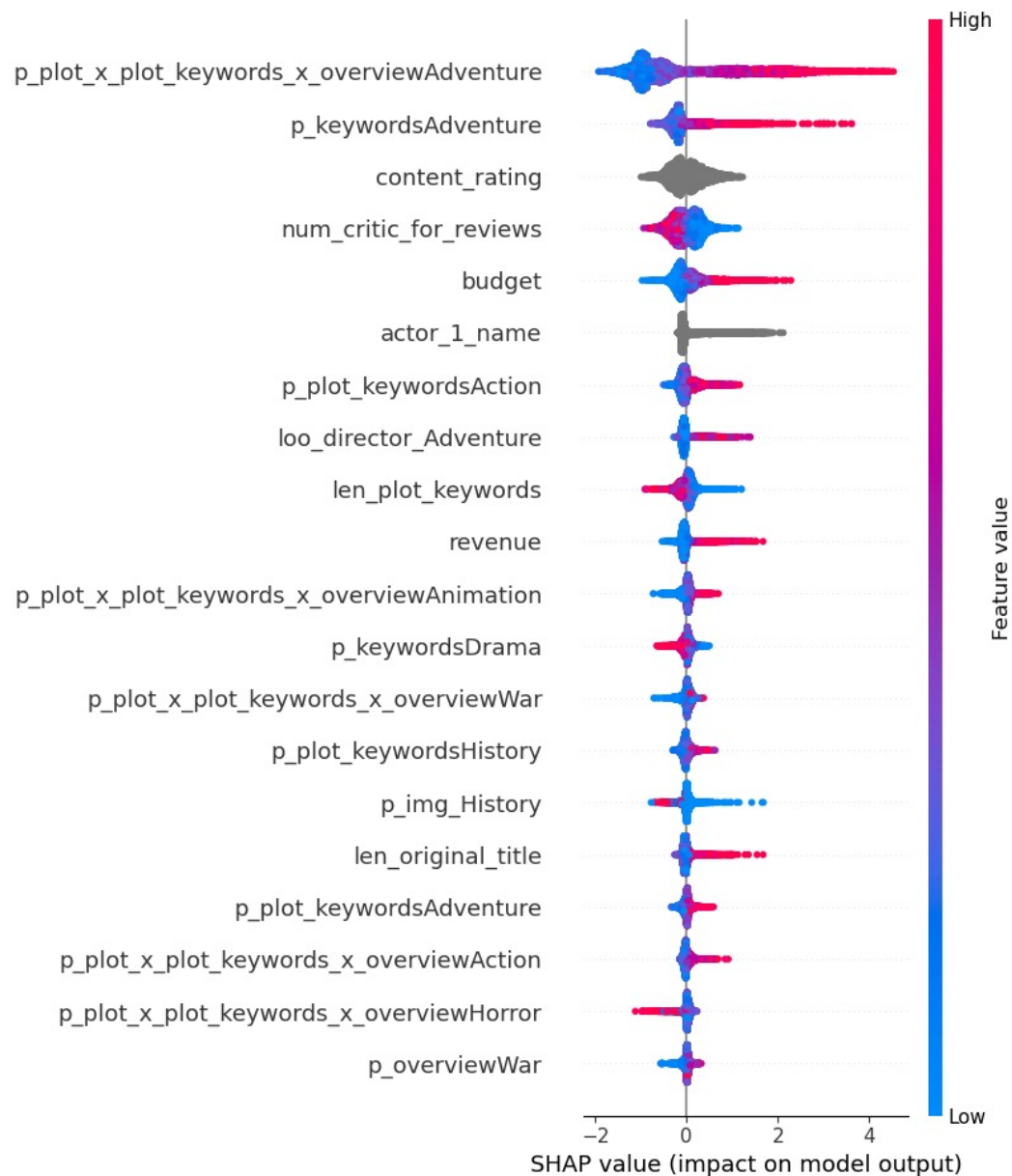$\mathcal{P}(\text{crime}) = 0.3$
$\mathcal{P}(\text{drama}) = 0.2$

$\mathcal{P}(\text{war}) = 0.2$
$\mathcal{P}(\text{crime}) = 0.4$
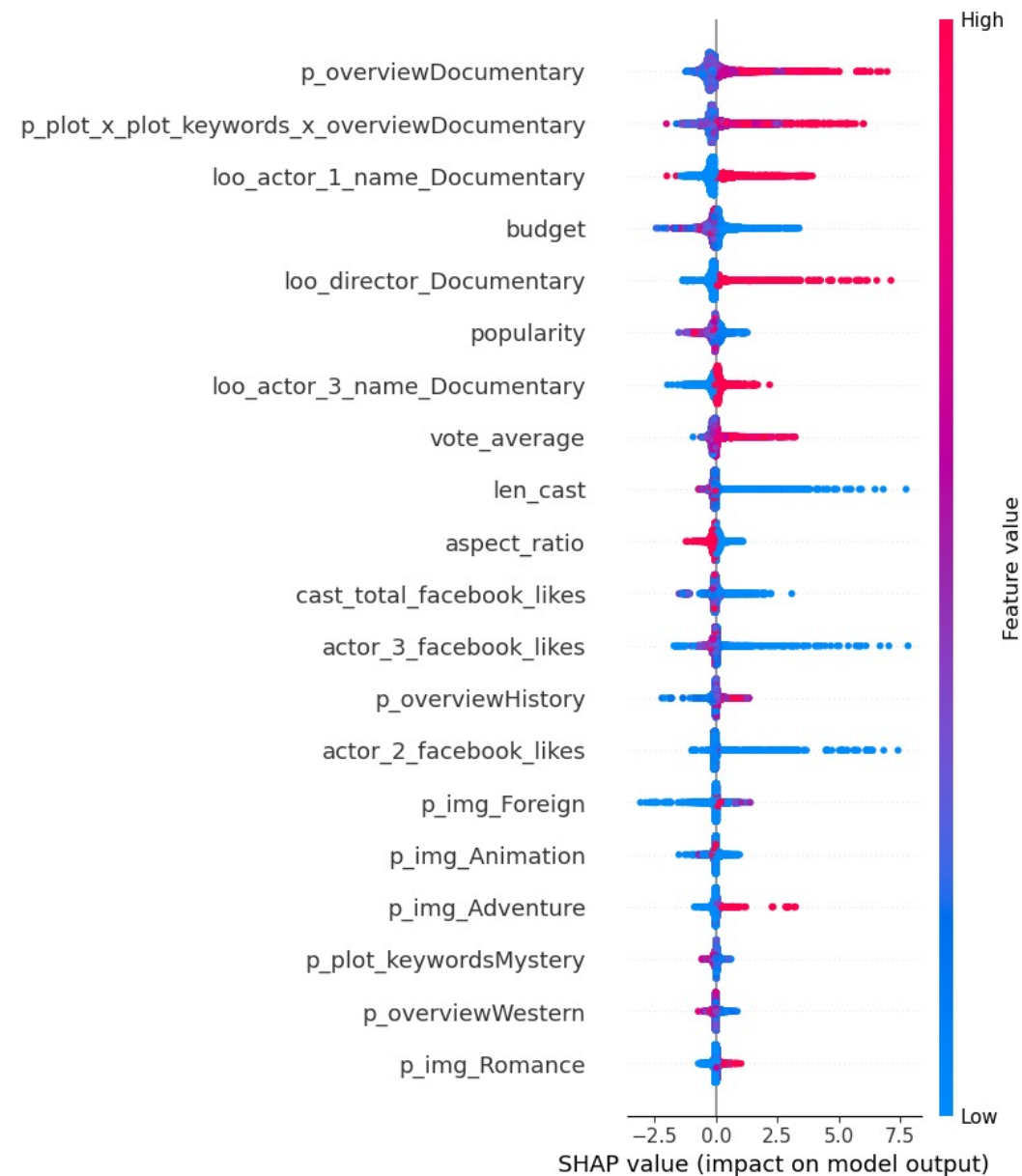$\mathcal{P}(\text{drama}) = 0.2$
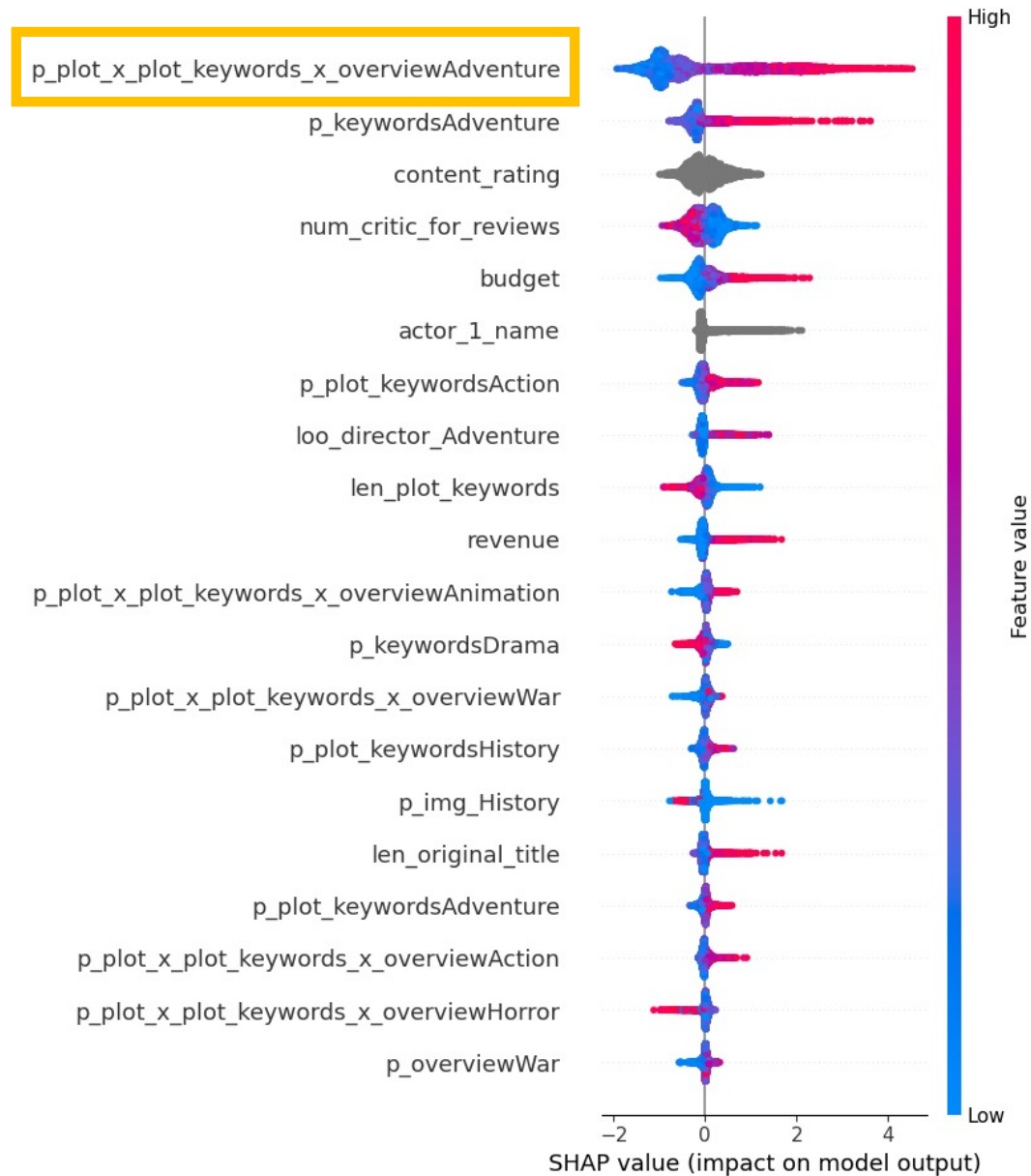
Predicción de género

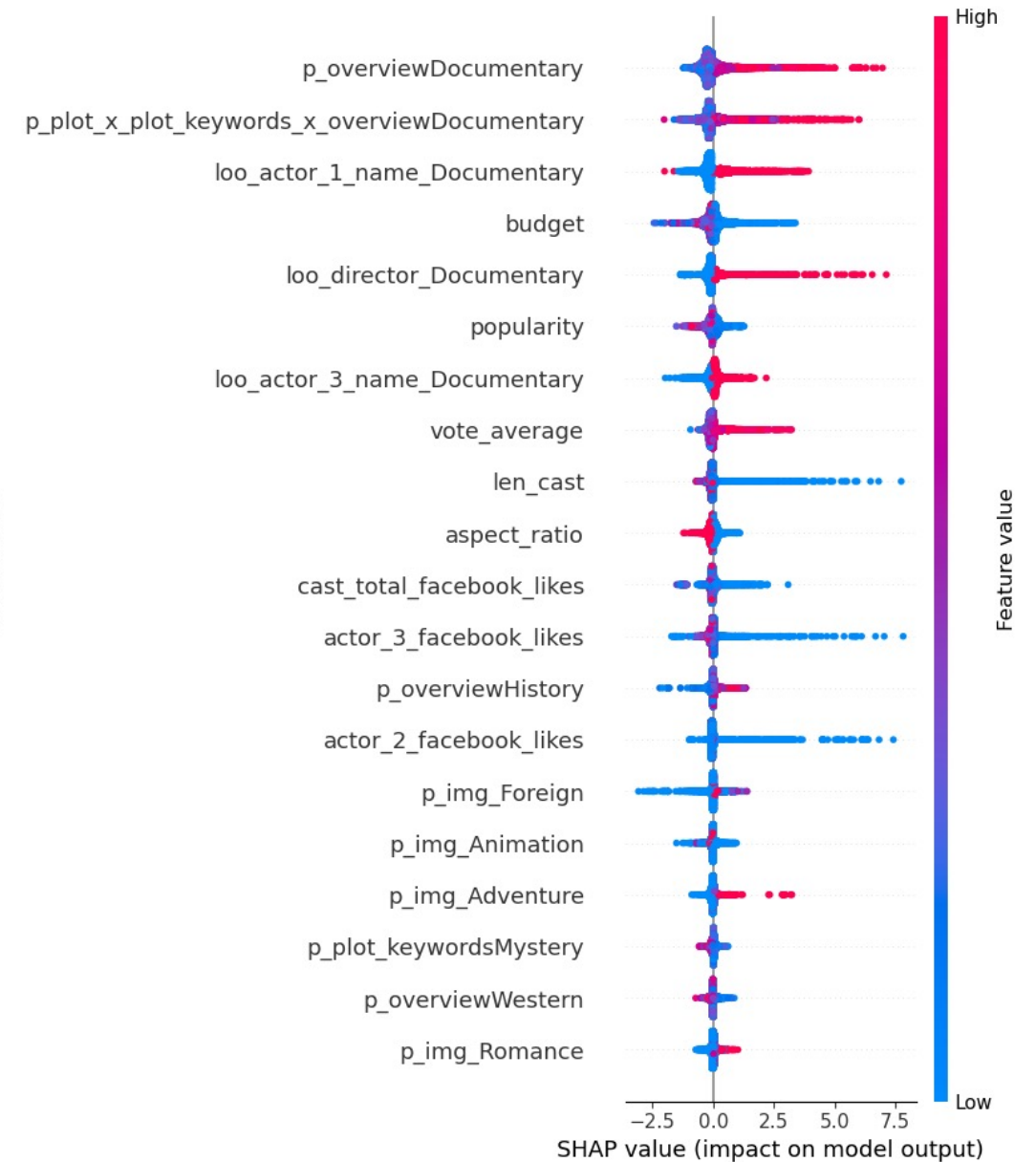|  | CV | TEST |
| --- | --- | --- |
| ACC | 0.92 | 0.27 |

EL MODELO FINAL

Adventure

Documentary

EL MODELO FINAL

Adventure

| p_plot_x_plot_keywords_x_overviewAdventure |
| p_keywordsAdventure |
| content_rating |
| num_critic_for_reviews |
| budget |
| actor_1_name |
| p_plot_keywordsAction |
| loo_director_Adventure |
| len_plot_keywords |
| revenue |
| p_plot_x_plot_keywords_x_overviewAnimation |
| p_keywordsDrama |
| p_plot_x_plot_keywords_x_overviewWar |
| p_plot_keywordsHistory |
| p_img_History |
| len_original_title |
| p_plot_keywordsAdventure |
| p_plot_x_plot_keywords_x_overviewAction |
| p_plot_x_plot_keywords_x_overviewHorror |
| p_overviewWar |

SHAP value (impact on model output)

Documentary

| p_overviewDocumentary |
| p_plot_x_plot_keywords_x_overviewDocumentary |
| loo_actor_1_name_Documentary |
| budget |
| loo_director_Documentary |
| popularity |
| loo_actor_3_name_Documentary |
| vote_average |
| len_cast |
| aspect_ratio |
| cast_total_facebook_likes |
| actor_3_facebook_likes |
| p_overviewHistory |
| actor_2_facebook_likes |
| p_img_Foreign |
| p_img_Animation |
| p_img_Adventure |
| p_plot_keywordsMystery |
| p_overviewWestern |
| p_img_Romance |

SHAP value (impact on model output)

EL MODELO FINAL

Adventure
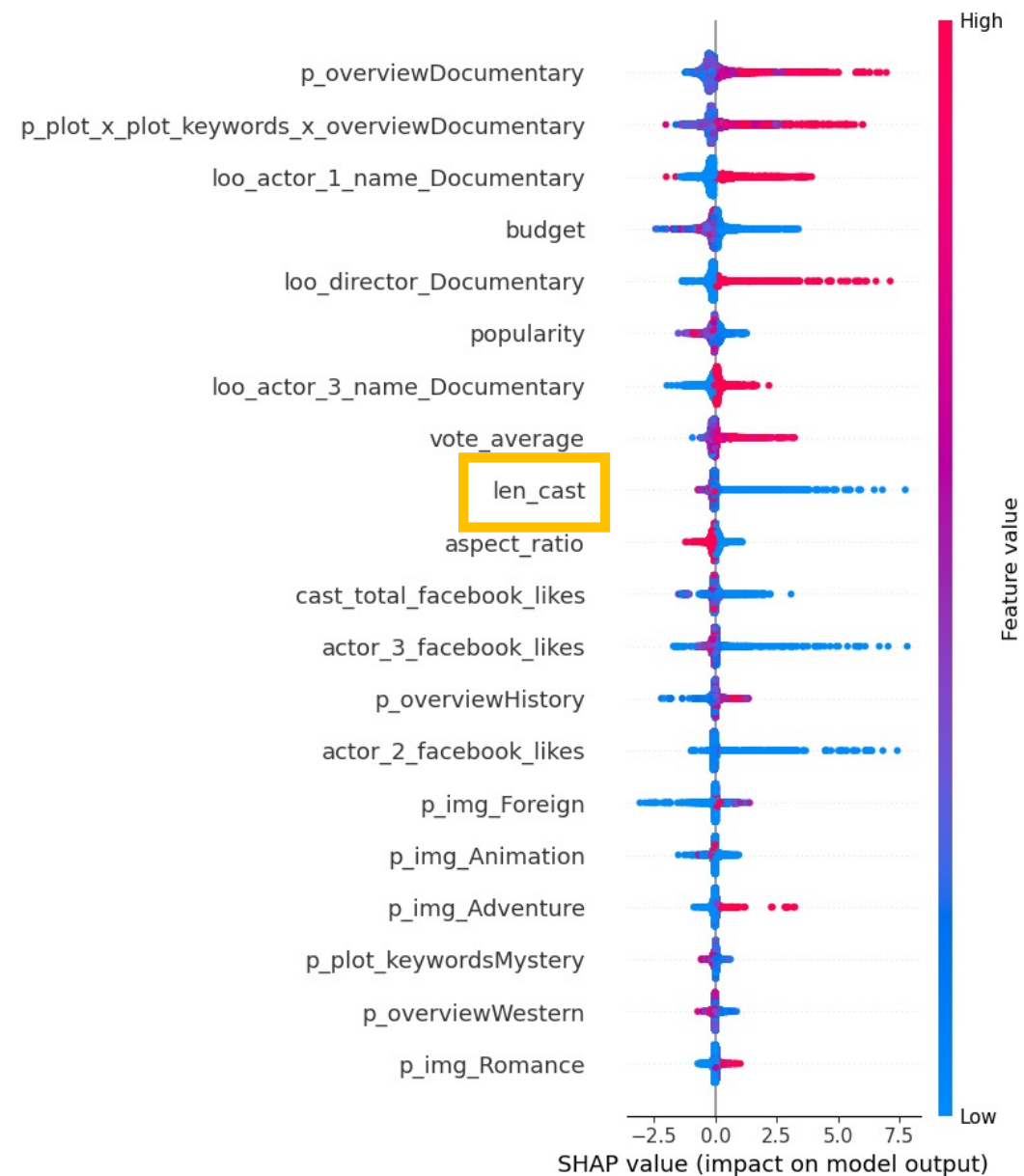
Documentary

EL MODELO FINAL

Adventure

Documentary

# CLASIFICACIÓN Y REGRESIÓN EN DATASET DE PELÍCULAS

| Datos | Procesado | Métricas |
|---|---|---|



**Datos tabulares**

2896 películas

**numéricos**
duration, content_rating.

**categóricos**
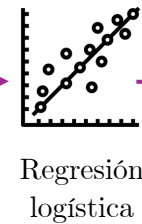color, language, country...

**texto**
keywords, overview,

LGBM

Explicabilidad con SHAP

iterar

output

### Predicción de score

| | CV | TEST |
|---|---|---|
| MAE | **0.24** | **0.22** |
| MSE | **0.13** | **0.09** |

**Pósteres**

2894 x .jpg

CONCATENAR
(keywords, overview, plot_keywords)

frecuencia

$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

Vector DB

Regresión logística

$\mathcal{P}(\text{war}) = 0.6$
$\mathcal{P}(\text{crime}) = 0.3$
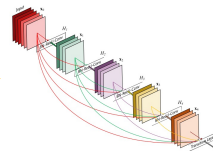$\mathcal{P}(\text{drama}) = 0.2$

**Imágen**
matrix 3 x 224 x 244

Pre-trained resnet + finetuning

$\mathcal{P}(\text{war}) = 0.2$
$\mathcal{P}(\text{crime}) = 0.4$
$\mathcal{P}(\text{drama}) = 0.2$

LGBM

OPTUNA

### Predicción de género

| | CV | TEST |
|---|---|---|
| ACC | 0.92 | 0.27 |