

Machine Learning aplicado al mundo del cine

¿Puede un algoritmo predecir la nota de una película?

Rubén Peinado

Objetivos del trabajo

- ★ Lanzamiento de nuevas películas.
- ★ Programación de la cartelera de cine en una sala de proyección de películas:
 - Queremos proyectar películas que dispongan de una determinada nota mínima.
 - Predecir la nota que puede tener una película.
 - Fidelizar a un público en concreto para nuestra sala de cine.
 - Destacarnos de otras salas de cine convencional.

Estructura del trabajo



Fuente de datos



Filmaffinity

Sistema recomendador de cine con una base de datos en la que se encuentra la ficha completa (técnica y artística) de gran cantidad de películas.

Listado de películas

Más de 10.000 películas

Ficha técnica y artística

10 categorías

Premios y Festivales

47 festivales y 278 categorías

Limpieza de datos - features



Numéricas



2/10

Año y duración.



Categóricas



8/10

Título, pais, director, guión, música, fotografía, género y actores.

Preprocesamiento de datos









Film
Total premios y
premios Razzie



Géneros get_dummies



Duración

Discretización y
eliminación outliers



País
Discretización por frecuencia

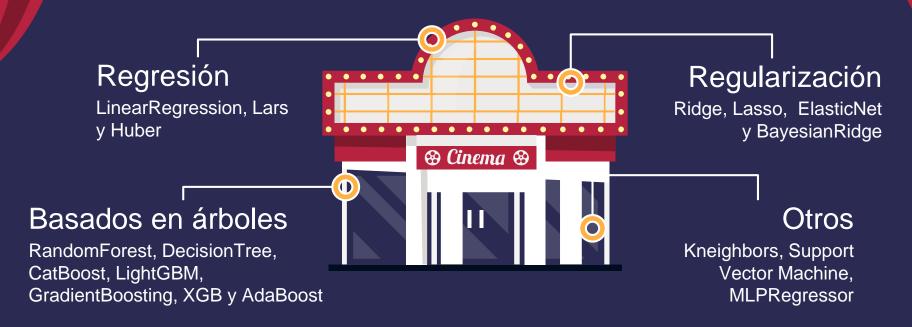
Película: The Artist



Actriz: Meryl Streep



Selección de algoritmos



Elección del modelo

Paso 1

Paso 2

Paso 2



Comparativa de resultados (métricas) 14 modelos diferentes



Voting con el mejor de cada grupo Ridge, LinearRegressor, Huber, CatBoost, SVM y MLP



Combinación y Ponderación

Combinar/ ponderar varios modelos

Evaluación del modelo



Métricas

Utilizaremos las métricas más apropiadas para la decisión de eficacia del modelo según la fase del trabajo.

R2 Score

Medida de la performance del modelo.

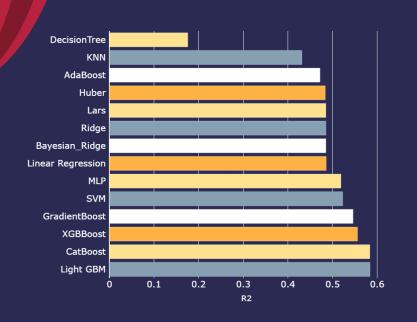
MAE

Medida del error absoluto en las predicciones.

RMSE

Medida que señala un menor error promedio.

Comparativa de modelos -Inicio



R2 Score

CatBoost SVM

0.58 0.52

Mejor MAE

Bayesian Ridge Linear Regressor

0.48 0.48

Mejores modelos

CatBoost

Basados en árboles

Support Vector Machine

Otros

Bayesian Ridge

Regularización

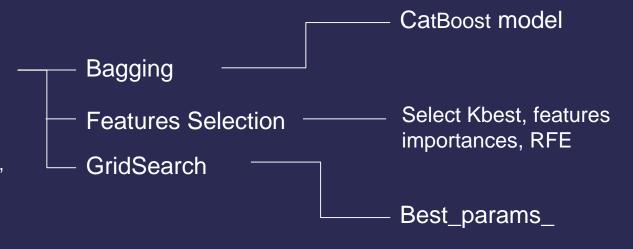
Linear Regressor

Regresión

Optimización del modelo

CatBoost

Potente y flexible, procesamiento eficiente de grandes volúmenes de datos, regularización incorporada y herramientas de interpretación.



Comparativa de modelos - Optimizados



MAE

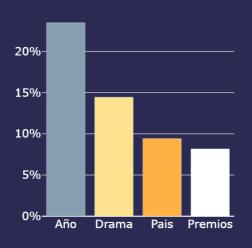
Bagging CatBoost + LGBM

0.559 0.561

GridSearch Light GBM

0.563 0.565

Features Importances



23%

Año

Preferencia de cine clásico

9%

Pais

Feature desbalanceada

14%

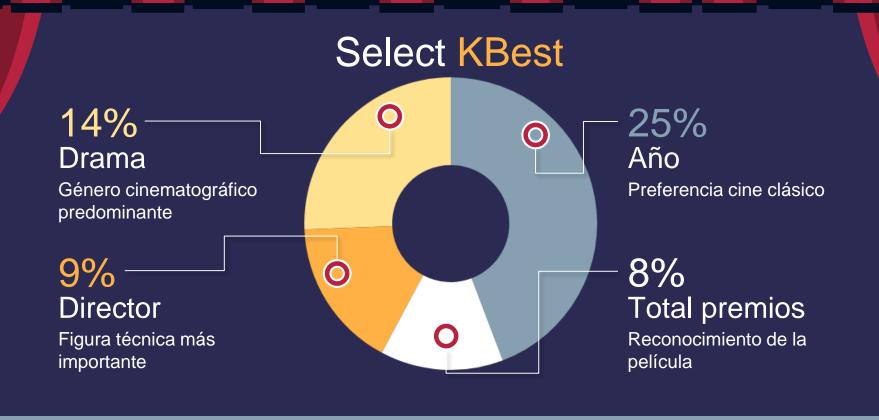
Drama

Género cinematográfico predominante

8%

Total premios

Reconocimiento de la película



Recursive Feature Elimination



Reconocimiento de la película



Drama



Género cinematográfico predominante





Feature desbalanceada



Año

Preferencia cine clásico

Comparativa de modelos – Features Select.



MAE

Bagging CatBoost + LGBM

0.5593 0.5614

Select KBest Grid Search

0.5618 0.5633

Comparativa de modelos – Features Select.



RMSE

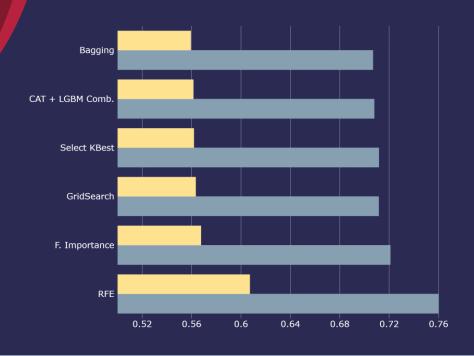
Bagging CatBoost + LGBM

0.7068 0.7079

Select KBest GridSearch

0.7116 0.7115

Comparativa de modelos – Features Select.



RMSE/ MAE

Bagging

CatBoost + LGBM

Select KBest

GridSearch

Modelos ganadores - RMSE



Conclusiones del trabajo

	Resumen	Mejoras/ próximos pasos
Data Analysis	Una única fuente facilita el proceso pero limita la obtención de nuevas features.	Ampliar fuentes de datos para mayor información. Enriquecer el dataframe con nuevas fuentes.
Feature Engineering	Preprocesado mixto (automático/ manual). Favorece el trabajo del modelo.	Establecer nuevas estrategias para sacar partido a las features iniciales.
Machine Learning	Se mejora la performance con la optimización. Ponderación de modelos, bagging o feature selection.	Establecer nuevo punto de vista de resolución del problema. Clasificación multiclase.



