

Proyecto 1

Inteligencia De Negocios

Jairo Nicolás Gómez

Paola Campiño

Felipe Duque

Análisis de Sentimientos en Películas

Abril 2 de 2023

Contenido

1. Entendimiento del negocio y enfoque analítico	3
2. Modelado y evaluación:.....	3
a. Regresión logística (Felipe Duque):	3
b. Resultados: Regresión logística (Felipe Duque):	4
c. Árbol de decisión (Paola Andrea Campiño):	4
d. Resultados: Árbol de decisión:	5
e. SVM(Jairo Nicolás Gómez):	5
f. Resultados SVM (Jairo Nicolás Gómez):	6
3. Resultados en base a todos los modelos	6
4. Acta de reunión con la experta en estadística:.....	7
5. Trabajo en Grupo:.....	7

1. Entendimiento del negocio y enfoque analítico

Oportunidad/problema Negocio	La oportunidad que vemos en este proyecto es poder ver y analizar las opiniones de los espectadores para mejorar los factores que son más criticados en las películas y mantener los aspectos más positivos y que son mejor resaltados por los espectadores.
Enfoque analítico (Descripción del requerimiento desde el punto de vista de aprendizaje automático)	Llegar a obtener un modelo con una alta precisión para analizar las reseñas sobre las películas y lograr clasificar de la mejor manera si son negativas o positivas.
Organización y rol dentro de ella que se beneficia con la oportunidad definida	<p>Este modelo podría llegar a tener un beneficio en empresas como lo puede ser cineColombia, Netflix y plataformas de streaming de películas o incluso páginas web de reseñas. Pues de esta forma muy rápidamente se podría llegar a una idea de cómo se ha recibido las películas entre el público para tomar decisiones de negocio como quitarla de taquilla, no recomendarla.</p> <p>En un cierto punto se podría hacer un análisis detallado de como la opinión de las personas podría cambiar en cuanto a una película al paso del tiempo.</p>
Técnicas y algoritmos a utilizar	Utilizaremos una técnica de clasificación y los algoritmos de regresión logística, arboles de decisión y SVM (máquinas de vectores de soporte).

2. Modelado y evaluación:

a. Regresión logística (Felipe Duque):

El modelo de regresión logística para análisis de textos es una técnica de aprendizaje supervisado que se utiliza para predecir la probabilidad de que un texto pertenezca a una o más categorías predefinidas. Este modelo es especialmente útil para el análisis de texto, ya que permite modelar la relación entre las palabras o características de un texto y la probabilidad de que ese texto pertenezca a una determinada categoría.

El modelo de regresión logística utiliza una función logística para modelar la relación entre las características de entrada del texto y la variable de salida categórica. Esta función logística transforma una variable continua en una variable binaria, es decir, la probabilidad de que el texto pertenezca a una categoría determinada.

El modelo de regresión logística se utiliza en análisis de textos porque permite la clasificación de grandes volúmenes de datos de texto en diferentes categorías de

manera rápida y precisa. Por ejemplo, en este caso se puede utilizar para analizar opiniones de los clientes sobre una película.

b. Resultados: Regresión logística (Felipe Duque):

Luego de entrenar el modelo de regresión logística con el 70% (3500 registros) de los datos proporcionados fue testeado con el 30% restante (1500 registros) y lo primero que generamos es una matriz de confusión que arroja los siguientes datos:

- El modelo predijo 594 registros positivos correctamente.
- El modelo predijo 668 registros negativos correctamente.
- El modelo predijo 143 registros positivo erróneamente.
- El modelo predijo 95 registros negativos erróneamente.

A partir de lo anterior encontramos las siguientes métricas cuantitativas para la evaluación de la efectividad del modelo:

- Para clasificar comentarios negativos obtuvimos:
 - Precisión: 86%
 - Recall: 81%
 - F1-score: 83%
 - 737 registros
- Para clasificar comentarios positivos obtuvimos:
 - Precisión: 82%
 - Recall: 88%
 - F1-score: 85%
 - 763 registros

Finalmente, en promedio el modelo tiene una precisión del 84%, un recall del 84% y un f1-score del 84% lo cual es indicativo de que se ha logrado un buen modelo de clasificación para categorizar los sentimientos sobre las películas entre positivo y negativo.

Informe	de				clasificación:
	precision		recall	f1-score	support
0	0.86		0.81	0.83	737
1	0.82		0.88	0.85	763
accuracy				0.84	1500
macro avg	0.84		0.84	0.84	1500
weighted avg	0.84	0.84	0.84	1500	

c. Árbol de decisión (Paola Andrea Campiño):

Los arboles de decisión es un algoritmo de clasificación en un aprendizaje supervisado no paramétrico. La solución plantea un árbol, que consta de nodo raíz, ramas, nodos y hojas. Los nodos hojas representan los resultados posibles dentro del

conjunto de datos, donde hay un numero de reseñas y un numero de entropía. Cabe resaltar que entre menor sea dicho número mejor son la clasificación de las reseñas.

d. Resultados: Árbol de decisión:

Para la implementación del árbol de decisión por este algoritmo se planteó 3 diferentes modelos. El primero con un 0.691 de precisión en la que se define el número mínimo de elementos para generar un nodo de 2 elementos, usando como métrica la entropía y una altura por default. Las métricas relacionadas a este modelo se encuentran a continuación:

	precision	recall	f1-score	support
0	0.68	0.73	0.70	498
1	0.71	0.65	0.68	502
accuracy			0.69	1000
macro avg	0.69	0.69	0.69	1000
weighted avg	0.69	0.69	0.69	1000

Para el segundo modelo se intentó encontrar el mejor árbol posible del que salió que los parámetros deberían ser que la métrica sería entropía, la altura sería de 28 y que el número mínimo de elementos debía estar 4 del que si se llegó a las siguientes métricas:

	precision	recall	f1-score	support
0	0.66	0.80	0.72	498
1	0.75	0.59	0.66	502
accuracy			0.69	1000
macro avg	0.70	0.69	0.69	1000
weighted avg	0.70	0.69	0.69	1000

En base a lo anterior se puede decir que no es muy bueno usar árbol de decisión como modelo para hacer la clasificación, pues el porcentaje de error es de un 30% que considerablemente alto.

Y Finalmente para el 3r modelo se hizo una prueba con todos los parámetros por default los cuales, no generaron muchos cambios en términos del porcentaje de error. En si se considera que al hacer un modelo con árboles de decisión se puede notar que la precisión sigue estando en un 69% y un porcege de error del casi 31%.

e. SVM(Jairo Nicolás Gómez):

El algoritmo de máquinas de vectores de soporte es un algoritmo de aprendizaje supervisado, el cual se usó para la clasificación binaria y poder separar de la mejor forma posible dos clases diferentes de puntos de datos. Lo anterior con el objetivo de negocio en cuenta para mirar las películas con comentarios positivos y negativos, y ver cómo mejorar para un futuro con base en la retroalimentación de los espectadores. Después de hacer todo el proceso de entendimiento de datos y

limpieza, y crear la matriz de palabras con el countvectorizer, se realizó un modelo el cual en la primera iteración arrojo resultados muy buenos, tomando 80% de registros como entrenamiento y 20% de registros de prueba.

f. Resultados SVM (Jairo Nicolás Gómez):

Lo anterior arrojo un 82% de precisión por lo que no fue necesario hacer un segundo modelo

Classification				Report:
	precision	recall	f1-score	support
0	0.83	0.80	0.81	485
1	0.82	0.84	0.83	515
accuracy			0.82	1000
macro avg	0.82	0.82	0.82	1000
weighted avg	0.82	0.82	0.82	1000

En esta tabla se observa que:

Precisión: En ambas clases es del 82%, lo que indica que el modelo tiene buena capacidad para clasificar las instancias

Recall: El recall promedio es de 82%, lo que indica que el modelo tiene buena capacidad para detectar correctamente las instancias positivas

F1-Score: El modelo tiene un buen equilibrio entre la precisión y el recall pues es del 82% también.

3. Resultados en base a todos los modelos

En base a los resultados de los modelos anteriores podemos concluir que lo mejor sería usar el modelo de regresión logística, ya presenta un menor porcentaje de error cerca del 16% y una precisión de casi del 84%.

Cabe resaltar que las pruebas se hicieron con una muestra 1000 reseñas sobre la población, por lo que es muy posible que teniendo un mayor número de datos podría llegar a una precisión mayor. Es muy posible que sea necesario mejorar la limpieza de datos de manera que podamos llegar un error menor lo cual garantizaría una mayor confiabilidad sobre el modelo planteado.

Si bien el porcentaje de error es considerablemente alto, es importante enfatizar que se puede mejorar el modelo y a futuro podría llegar a ser de valor dentro de la organización de manera que se pueda tomar decisiones en cuanto a la recepción de una película nueva, incluso podría llegar a ser de valor para el descubrimiento de películas que podrían llegar a ser buenas.

4. Acta de reunión con la experta en estadística:

Acta de Reunión	
Fecha: 1/04/2023 Hora Inicio: 6:51pm	Hora de Finalización: 7:16pm
Asistentes: Paola Andrea Campiño (Grupo 26) Jairo Nicolas Gomez Mendoza (Grupo 26) Valery Sharith Fonseca Pana (Experta en estadística)	Orden Reunión: 1. Presentación 2. Contexto 3. Presentar modelos 4. Presentación de métricas 5. Evaluación de métricas.
Observaciones/Recomendación:	
<ul style="list-style-type: none">- Hay que mejorar la interpretación tabla confusión, expresarlo en % de error.- En este caso se habla de una muestra y no una población. (Dentro de una población se seleccionó 5000 muestras)- Hablar del margen de error.	
Firma: Paola Andrea Campiño (Grupo 26) Jairo Nicolas Gomez Mendoza (Grupo 26) Valery Sharith Fonseca Pana (Experta en estadística)	

5. Trabajo en Grupo: