

Clasificador de Películas

Lester Andrés García Aquino
Carnet: 1003115
Universidad Rafael Landívar
Guatemala

Oskar Majus De Paz
Carnet: 1034711
Universidad Rafael Landívar
Guatemala

Max Fernando Díaz
Carnet: 1145916
Universidad Rafael Landívar
Guatemala

Ángel Javier Jiménez Morales
Carnet: 1032517
Universidad Rafael Landívar
Guatemala

Walter Gerardo Acabal Matias
Carnet: 1152418
Universidad Rafael Landívar
Guatemala

ABSTRACT

En este proyecto, presentamos el desarrollo e implementación de un sistema de clasificación de películas basado en Naive Bayes, utilizando reseñas del sitio web Rotten Tomatoes. Nuestro sistema aprovecha las bases probabilísticas del clasificador Naive Bayes para categorizar las reseñas de películas como positivas o negativas. La implementación se realizó utilizando Python para el modelo de aprendizaje automático y Spring Boot para los servicios backend y frontend basados en Java, con Flask facilitando la integración de ambos. El preprocesamiento de datos incluyó la limpieza de texto y la tokenización. El modelo de Naive Bayes entrenado demostró una alta precisión y eficiencia en la clasificación de reseñas de películas. Nuestros resultados resaltan la efectividad de Naive Bayes para tareas de procesamiento de lenguaje natural y muestran la interoperabilidad fluida entre las tecnologías de Python y Java. Este trabajo subraya el potencial para el despliegue escalable y robusto de modelos de aprendizaje automático en aplicaciones web.

CONCEPTOS CCS

• Clasificación de la información • Algoritmos y métodos: Aprendizaje Automático, Teorema Bayes • Lenguajes de Programación: Python, Java, • Bibliotecas de Software: Flask, Pandas. • Tecnologías: SpringBoot

PALABRAS CLAVE

Clasificador Naive Bayes, Teorema de Bayes, Python, Java, SpringBoot

1 INTRODUCCIÓN

La clasificación de películas es una tarea importante en la industria cinematográfica y en los sistemas de recomendación de películas. La clasificación de películas puede utilizarse para mejorar la búsqueda de películas, la recomendación de películas y la organización de colecciones de película.

Existen diversos métodos para la clasificación de películas, incluyendo métodos basados en reglas, métodos basados en aprendizaje automático y métodos basados en redes neuronales. Los métodos basados en el teorema de Bayes son un tipo de método de aprendizaje automático que se ha utilizado con éxito para la clasificación de películas.

Este estudio describe la implementación de un clasificador de películas utilizando el Teorema de Bayes, con una arquitectura que combina Spring Boot para el desarrollo del sistema web y la librería de Python Flask para ejecutar el Teorema de Bayes dentro del lenguaje de Java.

2 FUNDAMENTOS

2.1 Teorema de Bayes

El Teorema de Bayes, también conocido como la regla de Bayes, es un principio fundamental en la teoría de la probabilidad y la estadística. Nombrado en honor a Thomas Bayes, este teorema describe la probabilidad de un evento basado en el conocimiento previo de las condiciones que podrían estar relacionadas con el evento.

El Teorema de Bayes proporciona una forma matemática de actualizar las probabilidades de hipótesis, dada la evidencia. La fórmula del teorema de Bayes es la siguiente:

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

Donde:

- **$P(A|B)$** es la probabilidad posterior de la hipótesis **A** dado el dato **B**.

- $P(B|A)$ es la probabilidad de observar el dato **B** dado que la hipótesis **A** es verdadera.
- $P(A)$ es la probabilidad priora de la hipótesis **A**.
- $P(B)$ es la probabilidad total de observar el dato (**B**)

El teorema de Bayes tiene una amplia gama de aplicaciones y se utiliza en diversos campos, incluyendo la estadística, la informática, la medicina y la ingeniería.

En el contexto del aprendizaje automático, el teorema de Bayes se utiliza a menudo en la inferencia bayesiana y los modelos probabilísticos. Permite actualizar nuestras creencias sobre una hipótesis basándonos en nuevas evidencias.

Aplicaciones del Teorema de Bayes en el aprendizaje automático:

- Clasificador Naive Bayes
- Optimización Bayesiana
- Redes de Creencia Bayesiana

2.1 Clasificador Naive Bayes

En el ámbito del aprendizaje automático y la minería de datos, el clasificador es un algoritmo probabilístico utilizado para tareas de clasificación. Basado en el teorema de Bayes, este método es popular debido a la simplicidad, eficiencia y efectividad, especialmente con conjuntos de datos de gran tamaño.

El clasificador funciona calculando la probabilidad de que la instancia pertenezca a cada clase y asigna la etiqueta de la clase con mayor probabilidad.

Pasos del algoritmo:

1. Entrenar el modelo:
 - Se recibe un conjunto de datos de entrenamiento con ejemplos ya etiquetados.
 - Para cada clase, se calcula la probabilidad a priori y la probabilidad condicional de cada característica.
2. Clasificar una nueva instancia
 - Se recibe una nueva instancia de datos sin etiqueta.
 - Se calcula la probabilidad posterior de cada clase ($P(A|B)$) utilizando el teorema de Bayes:
 - $$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$
 - Donde $P(B)$ es la probabilidad marginal de la instancia

Aplicaciones:

- Clasificación de correos electrónicos
- Filtrado de documentos.
- Detección de anomalías
- Recomendaciones

Ventajas:

- Facilidad de Implementación.
- Buena Precisión
- Manejo de datos dispersos

- Escalabilidad

Desventajas:

- Suposición de independencia
- Sensibilidad de datos ruidosos
- Problemas con clases desequilibradas.

3 METODOLOGIA

Para realizar este clasificador de películas se tomó la decisión de utilizar el clasificador de Naive Bayes sobre el perceptrón simple, se tomó esta decisión ya que los datos eran textuales y categóricos.

- Por lo tanto, Naive Bayes es particularmente efectivo con datos textuales debido a su capacidad para manejar grandes cantidades de palabras en un texto.
- Permite clasificar de manera más efectiva cuando se trata de texto con múltiples características independientes. En el caso de las películas, las palabras en una reseña pueden ser tratadas como características independientes que contribuyen a la probabilidad de una clasificación positiva o negativa.
- Bayes es rápido de entrenar y eficiente, lo que es útil cuando se trabaja con grandes conjuntos de datos.
- Bayes puede funcionar mejor que el perceptrón simple con conjuntos de datos más pequeños, lo cual es útil cuando no se tiene una gran cantidad de críticas para entrenar el modelo.

3.1 Construcción e Implementación del Clasificador

Para construir el clasificador se recolectaron críticas de películas de la página de Rotten Tomatoes, estas críticas fueron descargadas en archivo de formato .CSV. El clasificador fue desarrollado en Python 3 utilizando la librería de Pandas, ya que Python es un lenguaje muy popular para el aprendizaje automático. El clasificador se conectó por la librería Flask a la aplicación creada en el lenguaje de programación Java.

El conjunto de datos de Rotten Tomatoes se carga y limpia, y el modelo naive Bayes se entrena con estos datos limpios. Todo esto se realiza al iniciar el programa.

Pasos por seguir:

- El proceso comienza con la carga del archivo .CSV utilizando la librería de pandas.
- Luego se limpian los datos del archivo:

- más posible sobre las características y patrones presentes en los datos.
- 30% para prueba:
 - Este conjunto de datos se utiliza para probar la capacidad del modelo para generalizar a nuevos datos. No se utiliza durante el entrenamiento y sirve para evaluar la precisión.
- Variaciones:
 - División 80/20 o 60/40 dependiendo de la cantidad total de datos disponibles y de lo complejo que sea el problema
- Validación Cruzada:
 - Para una evaluación más robusta, se puede utilizar la validación cruzada, especialmente en conjuntos de datos más pequeños o variados. Esto implica dividir los datos en k subconjuntos y utilizar cada uno de ellos como un conjunto de prueba en diferentes iteraciones, mientras que el resto actúa como conjunto de entrenamiento.

CONCLUSIONES

- Para la tarea de clasificación de películas, el modelo bayesiano, demuestra ser más conveniente que el perceptrón simple, porque proporciona una metodología robusta para la clasificación de texto.
- La combinación de Python y Java mediante Flask y Spring Boot ofrece una solución eficiente y escalable para aplicaciones de clasificación en tiempo real.
- Al reservar un conjunto de datos de prueba que el modelo nunca ha visto durante su entrenamiento, se puede verificar si el modelo simplemente está memorizando los datos o si realmente está aprendiendo patrones generalizables.

REFERENCIAS

- [1] Autores: GeekforGeeks. 2024. Naïve Bayes Classifiers. Disponible en: [Naive Bayes Classifiers - GeeksforGeeks](#)
- [2] GeekforGeeks.(2024). Bayes theorem in machine learning. Disponible en: [Bayes Theorem in Machine learning - GeeksforGeeks](#)
- [3] IBM. (s.f). Clasificadores Naïve Bayes. IBM. Disponible en: <https://www.ibm.com/topics/naive-bayes>
- [4] Barrios, J. (2020). Inteligencia artificial y aprendizaje automático para todos. Disponible en: <https://www.juanbarrios.com/inteligencia-artificial-y-machine-learning-para-todos/>

Price:\$15.00

Repositorio de GitHub

https://github.com/LesterAGarciaA97/RottenTomatoes_MovieClassifier