



UNIVERSIDAD DE LOS ANDES
FACULTAD DE INGENIERIA
INGENIERÍA DE SISTEMAS Y COMPUTACIÓN

 **Universidad de
los Andes**
Facultad de Ingeniería

PROYECTO ENTREGA GRUPAL

Presentado a
Haydemar Nuñez

Presentado por
Tomas Acosta Bernal - 202011237
Diego Andres Rubio - 201816492
David Alvarez Polo - 201911318

BOGOTÁ DC - COLOMBIA
29 de Octubre 2023

Índice general

1. Segunda Parte	2
1.1. Descripción del proceso e implementación realizados por el ingeniero de datos	2
1.1.1. Implementación	2
1.1.2. Código	3
1.2. Desarrollo de la aplicación y justificación	5
1.2.1. Descripción del usuario/rol de la organización que va a utilizar la aplicación	5
1.2.2. Conexión entre aplicación y el proceso de negocio	5
1.2.3. Importancia de la aplicación para el rol	5
1.3. Resultados-Aplicación	5
1.4. Conclusiones	7
1.5. Trabajo en equipo	7
1.6. Referencias	8

Capítulo 1

Segunda Parte

1.1. Descripción del proceso e implementación realizados por el ingeniero de datos

El ingeniero de datos llevó a cabo el siguiente procedimiento para la parte de automatización-preparación de datos y modelamiento para la aplicación:

1. **Carga de datos:**

Se cargaron los datos desde un archivo de Excel en un DataFrame de Pandas.

2. **Preparación de datos:**

Se realizaron las siguientes operaciones de limpieza y preparación de los datos:

- Limpieza de caracteres no ASCII.
- Conversión de caracteres a minúsculas.
- Eliminación de signos de puntuación.
- Eliminación de stop words.
- Reemplazo de special coders por caracteres con tildes.
- Reemplazo de números por su representación textual.
- Tokenización.
- Normalización.

3. **Entrenamiento de modelos:**

Se entrenaron tres modelos de clasificación de textos:

- Regresión logística.
- Naive Bayes.
- Random Forest.

1.1.1. Implementación

Para esta parte se realizaron las siguientes tareas:

1. **Pipeline:** Se implementó un pipeline que integraba las funciones necesarias para llevar a cabo las transformaciones esenciales de los datos. Este pipeline permitió gestionar eficazmente el flujo de datos y aplicar las preparaciones requeridas de manera sistemática.

```

from sklearn.pipeline import Pipeline
from sklearn.base import BaseEstimator, TransformerMixin

class TextPreprocessor2(BaseEstimator, TransformerMixin):
    def __init__(self):
        pass

    def fit(self, X, y=None):
        return self

    def transform(self, X):
        p = engine()
        preprocessed_data = []

        for sentence in X:
            words = sentence.split()
            words = to_lowercase(words)
            words = replace_numbers(words, p)
            words = remove_punctuation(words)
            words = remove_non_ascii(words)
            words = remove_stopwords(words)
            words = remove_specialCoders(words)
            preprocessed_sentence = ' '.join(words)
            preprocessed_data.append(preprocessed_sentence)

        return pd.Series(preprocessed_data)

# Crear un pipeline
preprocessing_pipeline = Pipeline([
    ('text_preprocessor', TextPreprocessor2()),
    ('vectorizer', CountVectorizer(max_features=3000))
])

```

Figura 1 Pipeline implementado

2. **Framework para desarrollo de API:** Se utilizó Flask para desarrollar una API REST que permitiera acceder a los resultados de los modelos.

1.1.2. Código

Para la creación de nuestra aplicación es relevante destacar algunos fragmentos de código que se presentan a continuación:

1. **Importación de las librerías:** Realizamos la importación de todas las librerías necesarias para desplegar nuestra aplicación. Entre estas se destaca **Streamlit**. Esta nos permite hacer nuestra aplicación web y conectarla con nuestro notebook de manera sencilla.

```

import streamlit as st
import pandas as pd
import matplotlib.pyplot as plt
import proyecto_grupo13
from pandas_profiling import ProfileReport
import plotly.express as px

```

Figura 2 Librerías importadas

2. **Carga de datos y creación de gráficos:** Se realiza la respectiva carga de datos y luego se hace lo siguiente:
 - Se crea un diccionario llamado dataGraph que contiene los datos del gráfico de barras. La clave Modelo contiene los nombres de los modelos de clasificación. La clave Exactitud contiene las exactitudes de los modelos de clasificación.
 - Se crea un DataFrame de Pandas llamado df a partir del diccionario dataGraph
 - Crea un gráfico de barras con la biblioteca plotly.express. El gráfico muestra la exactitud de cada modelo de clasificación.

- Se utiliza la función `st.plotly_chart()` para mostrar el gráfico de barras en la aplicación web.

```
data = pd.read_excel('cat_345.xlsx')

dataGraph = {
    'Modelo' : ['Regresión Logística', 'Naive Bayes', 'Random Forest'],
    'Exactitud' : [proyecto_grupo13.accuracy2, proyecto_grupo13.accuracy3, proyecto_grupo13.accuracy4]
}

df = pd.DataFrame(dataGraph)
```

Figura 3 Carga y gráficos

3. **Modelado en la aplicación:** A continuación se muestra el código para modelar uno de los algoritmos usados en la anterior etapa en la aplicación web:

```
## Describir brevemente cada modelo y razones de elección
## Mostrar matrices de confusión creadas
st.title('Resultados')
st.write('_Nota: Las exactitudes son valores entre 0 y 1_')

## Gráficos sobre resultados
st.title('Gráfico comparativo sobre modelos')
fig = px.bar(df, x='Modelo', y='Exactitud', title='Gráfico de Barras')
st.plotly_chart(fig)

st.header('Regresión Logística Multivariada')
st.markdown('<div style="text-align: justify;">Este modelo nos permite analizar opiniones y relacionarlas con l
st.text('')
st.write('Exactitud del modelo de Regresión Logística:', proyecto_grupo13.accuracy2)
st.write('**Razones de elección:**')
st.markdown("""
- Capacidad para manejar problemas de clasificación multiclase
- Interpretabilidad
- Suposiciones
""")
st.subheader('Matriz de confusión - Regresión Logística')
st.image('./confusion_matrix_logist.png')
st.subheader('Predicciones Modelo Test Regresión Logística')
st.write(proyecto_grupo13.dataFramePredictedLogit)
```

Figura 4 Modelamiento en la aplicación

1.2. Desarrollo de la aplicación y justificación

1.2.1. Descripción del usuario/rol de la organización que va a utilizar la aplicación

:

El rol que va a utilizar la aplicación es el analista de datos. Este va a ser el responsable de realizar análisis de datos para la organización y para ello necesitará de la aplicación web para visualizar los resultados y así realizar dicho análisis de manera exitosa.

1.2.2. Conexión entre aplicación y el proceso de negocio

La aplicación web desarrollada en este proyecto apoyará el proceso de análisis de datos de la organización. El analista utilizará la aplicación para mostrar los resultados de sus análisis, lo que le permitirá tomar mejores decisiones a futuro y que le beneficien a la empresa. En este caso, sobre los mejores modelos o las situaciones en las que es más indicado usar uno en específico.

1.2.3. Importancia de la aplicación para el rol

Esta aplicación es de gran importancia para el analista por los siguientes motivos:

- Le permitirá compartir sus análisis con otros usuarios de la organización.
- Podrá obtener retroalimentación de otros usuarios sobre sus análisis.
- Logrará identificar tendencias y oportunidades en los datos.

1.3. Resultados-Applicación

A continuación se presenta nuestra aplicación implementada:



Figura 5 Modelamiento en la aplicación



Figura 6 Visualización de la aplicación

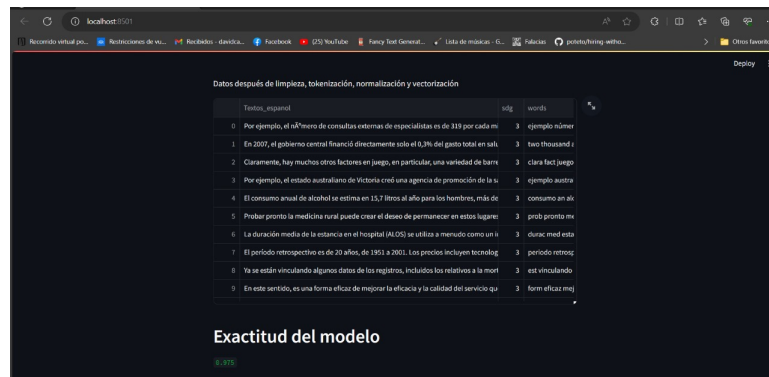


Figura 7 Visualización de la aplicación



Figura 8 Visualización de la aplicación

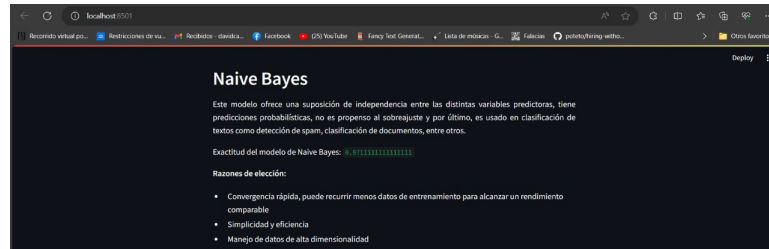


Figura 9 Visualización de la aplicación

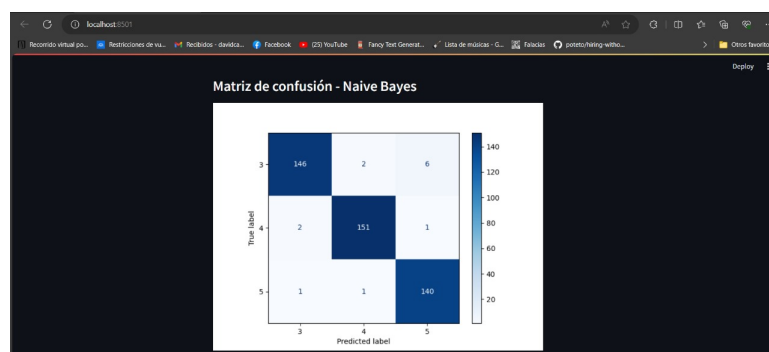


Figura 10 Visualización de la aplicación

Texto_ingresado	Predicciones
0 Se basa en los nueve años de educación básica (seis de primaria y tres de secundaria	3
1 En la última década, y en particular desde 2010, el número de años por cápita ha tem	4
2 Aún, por ejemplo, el estigma asociado a los beneficiarios de prestaciones sociales o a	3
3 Hay muchos profesores no cualificados en las escuelas, ya que es difícil contratar a p	3
4 A raíz de su preocupación por el hecho de que los médicos de todo el sistema sanitari	3
5 Hay que eliminar las prácticas discriminatorias de las instituciones financieras, com	4
6 Los centros de educación básica alternativa (EBA) ofrecen tres años con un plan de e	5
7 El uso de estadísticas de género puede proporcionar una comprensión más completa,	3
8 Cabe preguntarse, sin embargo, si los alumnos con más talento del país asisten a este	4
9 A pesar de que hay más pensionistas mujeres que hombres, los ingresos totales de la	5

Figura 11 Visualización de la aplicación

Texto_ingresado	Predicciones
0 Se basa en los nueve años de educación básica (seis de primaria y tres de secundaria	4
1 En la última década, y en particular desde 2010, el número de años por cápita ha tem	4
2 Aún, por ejemplo, el estigma asociado a los beneficiarios de prestaciones sociales o a	3
3 Hay muchos profesores no cualificados en las escuelas, ya que es difícil contratar a p	3
4 A raíz de su preocupación por el hecho de que los médicos de todo el sistema sanitari	3
5 Hay que eliminar las prácticas discriminatorias de las instituciones financieras, com	4
6 Los centros de educación básica alternativa (EBA) ofrecen tres años con un plan de e	3
7 El uso de estadísticas de género puede proporcionar una comprensión más completa,	3
8 Cabe preguntarse, sin embargo, si los alumnos con más talento del país asisten a este	3
9 A pesar de que hay más pensionistas mujeres que hombres, los ingresos totales de la	3

Conclusión

Se recomendó el uso del modelo de Naive Bayes, este obtuvo la exactitud más alta con un valor de 97.1%. No obstante, esta no es la única razón por la que se recomendó su uso sino también porque es fácilmente escalable, lo que ayuda a manejar cantidades más grandes de datos. Asimismo, su simplicidad para operar en este tipo de tareas de clasificación de textos.

Figura 12 Visualización de la aplicación

1.4. Conclusiones

A partir de esto podemos concluir que con la aplicación realizada el analista de datos podrá observar los resultados que arrojó cada algoritmo y a partir de esto escoger alguno dependiendo el caso o problema que se le presente a la hora de clasificar textos. Encontramos que si se desea mejorar la clasificación de la categoría 4 es mejor usar el algoritmo de regresión logística multivariada, mientras que para la categoría 5 es mejor Naive Bayes. En este caso no es recomendable usar Random Forest ya que mostró un mayor margen de error de predicción en comparación con los otros algoritmos. La implementación de Naive Bayes en particular podría brindar ventajas sustanciales a la UNFPA, puede automatizar la clasificación de opiniones de los ciudadanos sobre temas sociales y económicos. Este enfoque respalda a la organización de dos maneras: Permite identificar de manera rápida la categoría de ODS y reduce la necesidad de clasificar manualmente los comentarios

1.5. Trabajo en equipo

Tomás Acosta

Roles: Líder de proyecto, líder de datos

Tareas realizadas: Entendimiento de los datos, preparación de los datos, creación de modelo, elaboración del informe

Tiempo: 20 horas

Retos: Preparar datos para mantenerlo en lenguaje natural, aprender herramientas y conceptos nuevos para la preparación de los datos y obtención de palabras más útiles

Soluciones: Definición de funciones dentro de pipeline para remover caracteres que no sean ASCII, pasar todas las palabras a minúsculas, remover signos de puntuación, convertir caracteres especiales, implementar código para eliminar prefijos y sufijos, reemplazar números por palabras, realizar tokenización y lematización de los datos (normalización de los datos mediante el uso de lems y stems)

Puntaje: 33.3

Diego Rubio

Rol: Líder de negocio

Tareas realizadas: Entendimiento del negocio, entendimiento de los datos, preparación de los datos, documento/informe

Tiempo: 20 horas

Retos: Exploración de temáticas correspondientes para cada ODS, identificar el rol beneficiado y su conexión con la aplicación.

Soluciones: Investigar sobre el negocio a manera mas detallada y sobre los roles beneficiados por la aplicación.

Puntaje: 33.3

David Polo

Rol: Líder de analítica

Tareas realizadas: Implementación de la aplicación y wiki

Tiempo: 20 horas

Retos: Encontrar soluciones para implementar la aplicación y desplegarla

Soluciones: Usar librerías de python para poder crear la aplicación web

Puntaje: 33.3

1.6. Referencias

1. Scikit-learn developers. (2021). Scikit-learn: Machine Learning in Python. Recuperado de:
<https://scikit-learn.org/stable/index.html>
2. SciELO. (2014). Clasificación automática de textos usando redes de palabras. Recuperado de:
https://www.scielo.cl/scielo.php?script=sci_arttext&pid=S0718-09342014000300001.