Designed for:

Designed by:

Date:

## TAREA DE APRENDIZAJE



### **DECISIONES**

Los resultados del modelo se convierten en

decisiones procesables al servir como una

función de los ODS. Por ejemplo, si un texto

pobreza), esta información puede orientar a

un analista o institución a priorizarlo dentro

reducción de pobreza. De manera similar,

un texto etiquetado como ODS 3 (Salud y

programas de salud pública, mientras que

destinaría a proyectos educativos. De esta

mismo, sino que permite organizar, filtrar y

tomar decisiones estratégicas sobre a qué

área dirigir recursos, políticas o acciones

concretas. En el contexto del usuario final,

traduce cada texto en una categoría que

seguimiento de tendencias relacionadas

con los Objetivos de Desarrollo Sostenible.

facilita la asignación de prioridades,

generación de reportes temáticos o

estas decisiones son inmediatas: el modelo

uno de ODS 4 (Educación de calidad) se

forma, la clasificación no es un fin en sí

clasificación automatizada de textos en

es clasificado como ODS 1 (Fin de la

de iniciativas sociales o políticas de

bienestar) podría canalizarse hacia



PROPUESTA DE VALOR



El beneficiario final de este modelo son principalmente organizaciones, instituciones públicas, ONGs y centros de investigación que trabajan con grandes volúmenes de información vinculada a los Objetivos de Desarrollo Sostenible (ODS). La propuesta de valor radica en automatizar la clasificación de textos en función de su temática ODS, lo cual facilita la gestión, análisis y priorización de documentos, reportes o comunicaciones relacionadas con pobreza, salud y educación. Con ello se aborda el problema específico de la sobrecarga de información y la falta de categorización eficiente, que dificulta tomar decisiones rápidas y fundamentadas. El modelo permite ahorrar tiempo, mejorar la precisión en la asignación de textos a categorías relevantes y apoyar la trazabilidad de los avances en los ODS. Sin embargo, el uso del modelo también implica riesgos: por ejemplo, clasificaciones erróneas podrían desviar recursos a áreas menos prioritarias, generar interpretaciones equivocadas o invisibilizar información importante. Por eso, es clave acompañar el modelo con validaciones humanas y métricas de desempeño claras para mitigar estos riesgos.

## RECOLECCIÓN DE **DATOS**



### **FUENTES DE DATOS**



Las fuentes de datos provienen principalmente de corpus textuales en español que contienen documentos, artículos, reportes y fragmentos de texto vinculados a los ODS seleccionados (1, 3 y 4). En este caso, se han utilizado datasets ya estructurados en tablas internas, con campos como textos y labels (etiquetas correspondientes al ODS). También se cuenta con un dataset de prueba sin etiquetas, que permite evaluar la capacidad del modelo para generalizar y predecir sobre nuevos textos. En un escenario productivo, podrían integrarse datos externos a través de APIs públicas de Naciones Unidas, observatorios de sostenibilidad o bases abiertas de organizaciones internacionales, enriqueciendo así el entrenamiento y la cobertura temática. Estas fuentes resultan apropiadas para el objetivo del análisis, ya que ofrecen textos suficientemente variados y etiquetados, lo cual posibilita entrenar un modelo supervisado que aprenda las diferencias lingüísticas y contextuales entre cada ODS. No obstante, es importante mantener un equilibrio entre las clases y la calidad de los textos, para evitar sesgos y asegurar que el modelo sea robusto y aplicable a distintos contextos reales

En este caso, la tarea de aprendizaie corresponde a un aprendizaje supervisado de tipo clasificación de texto, ya que el modelo se entrena con ejemplos previamente etiquetados en tres categorías de ODS: ODS 1 (Fin de la pobreza), ODS 3 (Salud y bienestar) y ODS 4 (Educación de calidad). Lo que se predice es la clase a la que pertenece un nuevo texto en español, basándose en sus características lingüísticas y en el vocabulario procesado durante el entrenamiento. Los posibles resultados de la tarea de aprendizaje son, por lo tanto, las tres etiquetas de salida asociadas a cada texto. Los resultados se observan de manera inmediata, en el momento en que un texto nuevo es ingresado al modelo, lo que permite una clasificación instantánea y útil en escenarios de análisis documental o monitoreo de información. Esto implica que, a diferencia de modelos predictivos que anticipan eventos en horizontes temporales,

aquí la salida se obtiene en tiempo real al

procesar el texto.

## SIMULACIÓN DE IMPACTO



La simulación de impacto en este caso debe considerar tanto los beneficios de las predicciones correctas como los costos asociados a errores de clasificación. Cuando el modelo asigna correctamente un texto al ODS correspondiente, permite a la empresa u organización clasificar información de forma eficiente y automatizada, ahorrando costos de tiempo y recursos humanos, y generando valor en la aestión de conocimiento. Sin embargo, una clasificación incorrecta puede tener costos significativos: por ejemplo, ubicar un texto de salud (ODS 3) dentro de educación (ODS 4) podría llevar a errores en reportes estratégicos, indicadores de sostenibilidad o políticas públicas, reduciendo la credibilidad del sistema y

Los criterios de éxito del modelo se definen en términos de métricas como precisión, recall, F1-score y balance entre clases, además de su capacidad de mantener un rendimiento estable en el dataset de prueba. Para el despliegue, el modelo debe alcanzar un umbral mínimo de desempeño (por ejemplo, >80% en F1 promedio ponderado) que garantice que los resultados son confiables en un entorno real. En cuanto a restricciones de equidad, el modelo debe asegurar que todas las clases ODS tengan representación adecuada, evitando sesgos hacia la clase mayoritaria.

afectando la toma de decisiones.

# APRENDIZAJE (USO DEL MODELO)

El uso del modelo se plantea



principalmente en modo por lotes, ya que el flujo de trabajo consiste en recibir un conjunto de textos (por ejemplo, reportes, documentos institucionales o mensajes recopilados de distintas fuentes) y clasificarlos de manera masiva según el ODS correspondiente. Esto permite procesar grandes volúmenes de información de forma periódica, optimizando recursos. Sin embargo, el modelo también puede adaptarse a un uso en tiempo real, por ejemplo, integrándose en una plataforma donde, al ingresar un nuevo texto, se devuelva inmediatamente la predicción de su categoría ODS. La frecuencia de uso dependerá del contexto de aplicación: en entornos institucionales podría ejecutarse mensualmente o trimestralmente para la consolidación de reportes de sostenibilidad. En cambio, en sistemas de monitoreo social o académico podría usarse a diario o incluso en línea para dar soporte a la clasificación de información en tiempo real.

# CONSTRUCCIÓN DE MODELOS



En la ingeniería de características de este modelo se parte de los textos en español y se generan variables numéricas que permiten entrenar el clasificador. Las características clave incluyen:

**INGENIERÍA DE** 

**CARACTERÍSTICAS** 

de clasificación supervisada multiclase, capaz de distinguir entre los textos que pertenecen a ODS 1, ODS 3 u ODS 4. Sin embargo, en la práctica pueden explorarse varias arquitecturas y enfoques (por ejemplo, modelos de regresión logística, árboles de decisión, SVM) y luego seleccionar el que ofrezca mejor desempeño según métricas como FI-score o precisión balanceada. En ese sentido, aunque el producto final es un modelo, durante la construcción será necesario experimentar con múltiples variantes.

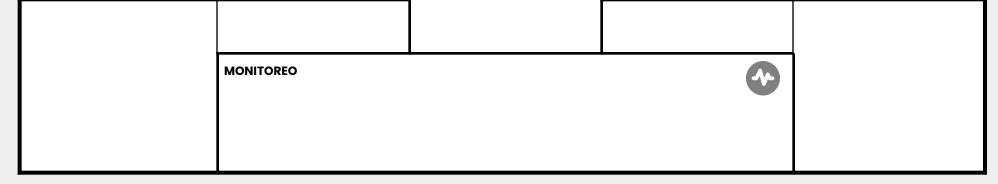
Para este caso se necesita un único modelo

Respecto a la actualización, el modelo debería reentrenarse periódicamente una o dos veces al año, o cuando se disponga de un volumen considerable de nuevos datos que reflejen cambios en la forma en que se redactan los textos sobre ODS. Esto es importante porque el vocabulario y los contextos de uso evolucionan.

En cuanto al tiempo disponible, considerando la ingeniería de características, limpieza de datos, tokenización, anállsis exploratorio, entrenamiento y validación, un ciclo inicial de construcción puede requerir entre 2 y 4 semanas de trabajo intensivo. Este plazo incluye tanto la preparación de los datos como la comparación de modelos candidatos y la selección del óptimo. Una vez en operación, los reentrenamientos periódicos serían más rápidos, ya que gran parte de la infraestructura estaría lista.

- Características estructurales: longitud en caracteres, número de palabras, número de oraciones y número de tokens después de la limpieza, lo que refleja la complejidad y densidad de cada texto.
- Preprocesamiento y normalización: conversión a minúsculas, eliminación de puntuación, números y espacios múltiples, además de la eliminación de stopwords en español (tanto de NLTK como de una lista manual adaptada al dominio ODS).
- Representación BoW (Bag of Words): los textos se transformarán en vectores según la frecuencia de aparición de cada palabra en el vocabulario, capturando así las palabras más frecuentes y discriminativas sin necesidad de ponderaciones adicionales.
- Agregaciones estadísticas: se evalúa la distribución de tokens por texto, la proporción de palabras únicas y se detectan outliers que pudieran alterar el comportamiento del modelo.
- Vocabulario por clase ODS: se construyen subconjuntos de palabras frecuentes dentro de cada categoría, permitiendo identificar términos más representativos de cada objetivo de desarrollo sostenible.

De esta forma, el modelo se nutrirá tanto de variables cuantitativas sencillas como de la representación BoW, lo que asegura una base coherente con el análisis exploratorio que realizamos y suficiente para capturar patrones de lenguaje en los textos.









Version 1.2. Created by Louis Dorard, Ph.D. Licensed under a <u>Creative Commons Attribution-ShareAlike 4.0 International License</u>. Please keep this mention and the link to ownml.co when sharing.

**OWNML.CO**