

# **PREDICCIÓN RESEÑAS AMAZON**

# Índice

**01**

---

**Recopilación y  
lectura de datos**

**02**

---

**Preprocesamiento  
básico**

**03**

---

**Modelo inicial**

**04**

---

**Mejoras  
preprocesamiento**

**05**

---

**Mejoras modelo**

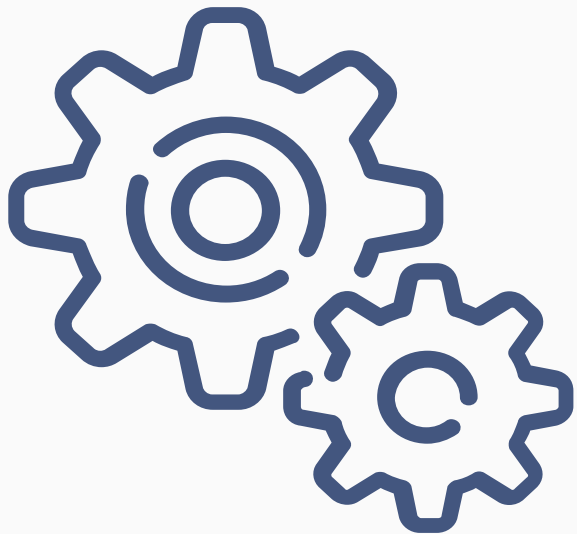
**06**

---

**Conclusiones**

# PRODUCTOS ELEGIDOS





# **PREPROCESAMIENTO BÁSICO**

# **PREPROCESAMIENTO BÁSICO**

**1. LETRAS MINÚSCULAS**

**2. ELIMINACIÓN  
SIGNOS PUNTUACIÓN**

**3. ELIMINACIÓN  
STOPWORDS**

**4. TOKENIZACIÓN**

# VECTORIZACIÓN TF - IDF

**BOW**

**VS**

**TF - IDF**

**VS**

**TF / IDF**

El más simple y rápido

Clasificación textos pequeños

Bueno cuando la semántica no es relevante

Simple, eficiente, rápido y equilibrado

Reduce palabras irrelevantes

Destaca palabras importantes

No tiene en cuenta la relación semántica

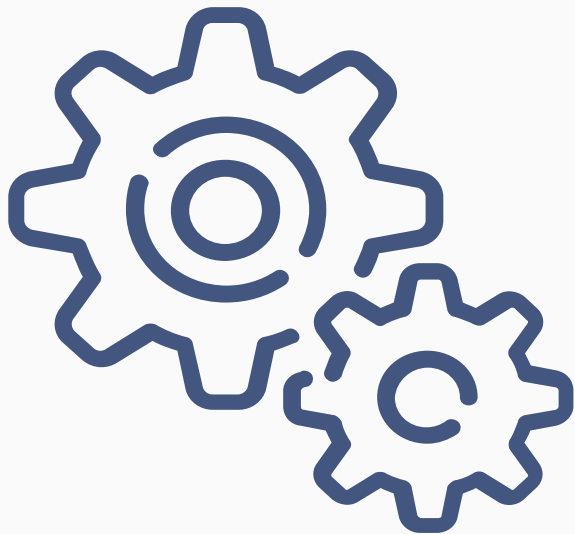
Mide importancia de palabras

Reduce palabras irrelevantes

No distingue palabras comunes y clave

No considera la frecuencia

No suficientes por sí solos



# ELECCIÓN MODELO INICIAL

# MODELOS SELECCIONADOS

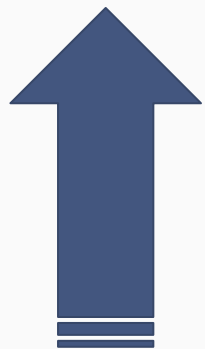
NAIVE BAYES  
MULTINOMIAL

**50%**

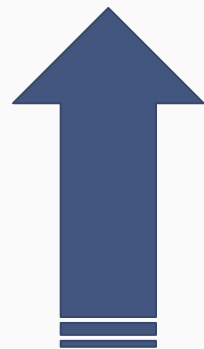
REGRESIÓN  
LOGÍSTICA

**58%**





**MEJORAS BASADAS EN  
EL MODELO BASE**



# Mejoras en cuanto a preprocesamiento

---

Lemming / Stemming

# Recordando...

**Análisis  
lingüístico para  
sacar el lema o  
forma canónica  
de cada palabra**

*Running -> Run  
Better -> Good*



**Se basa en recorte  
de sufijos. No  
devuelve palabras  
con sentido siempre**

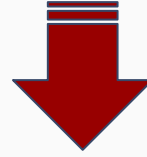
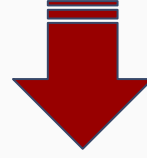
*Running -> Run  
Better -> Bet*



# Intuímos

Rendimiento

Tiempo Empleado



**Lemming**

**Stemming**

# Resultados



# Conclusiones

**Idea de simpleza  
durante el resto del  
proyecto**



**Modelos sencillos no  
funcionan con datos  
complejos**

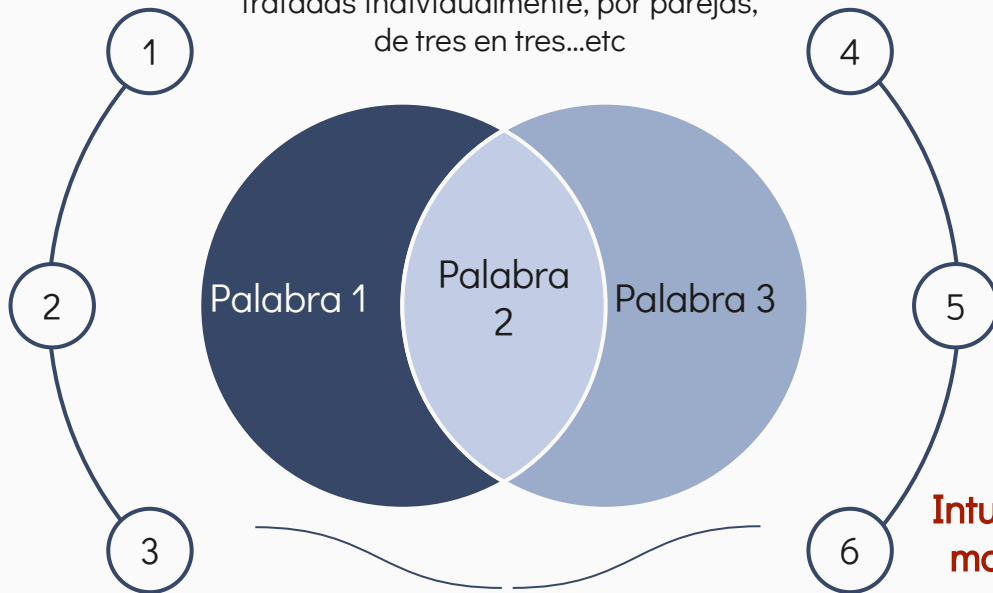
# Mejoras en cuanto al tratamiento aplicado

---

Vectorizador TF-IDF de Sklearn

# N-Gramas

Las distintas palabras pueden ser tratadas individualmente, por parejas, de tres en tres...etc



Tratados como  
un solo token

**Intuimos que un valor de N  
mayor conseguirá mayor  
contexto y por lo tanto  
mejor resultado**

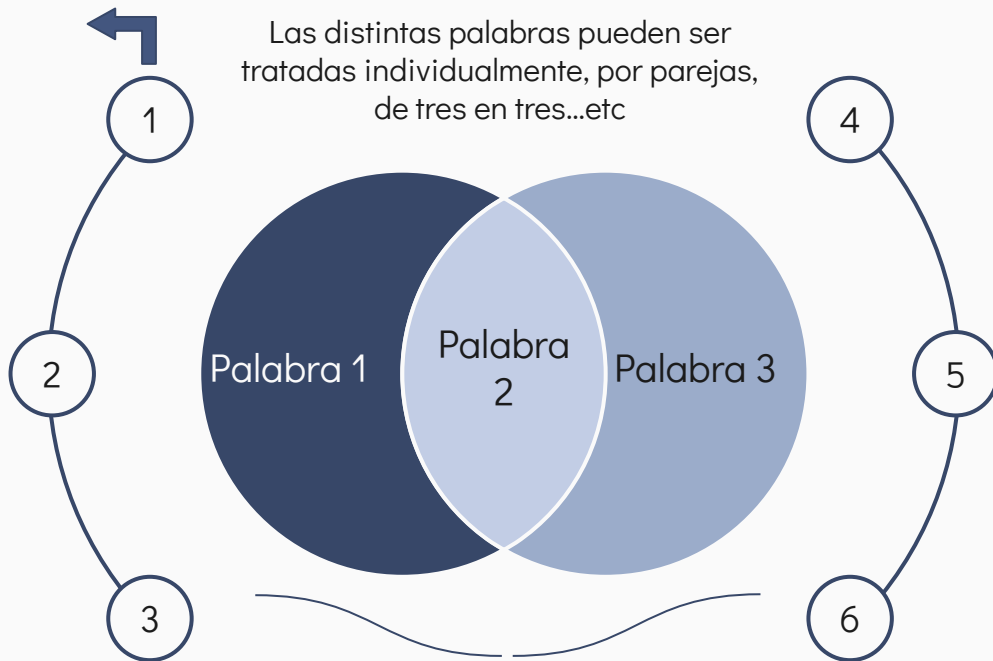


Una vez más el mejor  
resultado es el más  
sencillo

Implementamos incluso  
rangos en los que  
generaba combinaciones,  
ampliando el espacio de  
vectores con lo que  
aprendían los modelos

# N-Gramas

Las distintas palabras pueden ser  
tratadas individualmente, por parejas,  
de tres en tres...etc



Tratados como  
un solo token

# N-Gramas

Vimos que el resultado era realmente malo para Naive Bayes

$$P(c | x) = \frac{P(x | c)P(c)}{P(x)}$$

Diagram illustrating the components of the Naive Bayes formula:

- $P(c | x)$  is labeled as **Posterior Probability**.
- $P(x | c)$  is labeled as **Likelihood**.
- $P(c)$  is labeled as **Class Prior Probability**.
- $P(x)$  is labeled as **Predictor Prior Probability**.

$$P(c | X) = P(x_1 | c) \times P(x_2 | c) \times \dots \times P(x_n | c) \times P(c)$$

Con diferentes pruebas manuales y discutiendo con otros grupos, dimos con el porqué. Al usar conteo de palabras para calcular las probabilidades, si como palabra teníamos en cuenta agrupaciones de más de 3, las posibilidades de que se repitieran eran muy bajas por lo que no tenía ejemplos suficientes como para generar un buen entrenamiento

# Normalización

## ¿Para qué?

Normalizar los textos de entradas resulta importante para equiparar la relevancia de reviews con distinto conteo de palabras



Si no lo normalizamos :  
Las reviews con más palabras conseguirán solo por ese efecto valores más altos

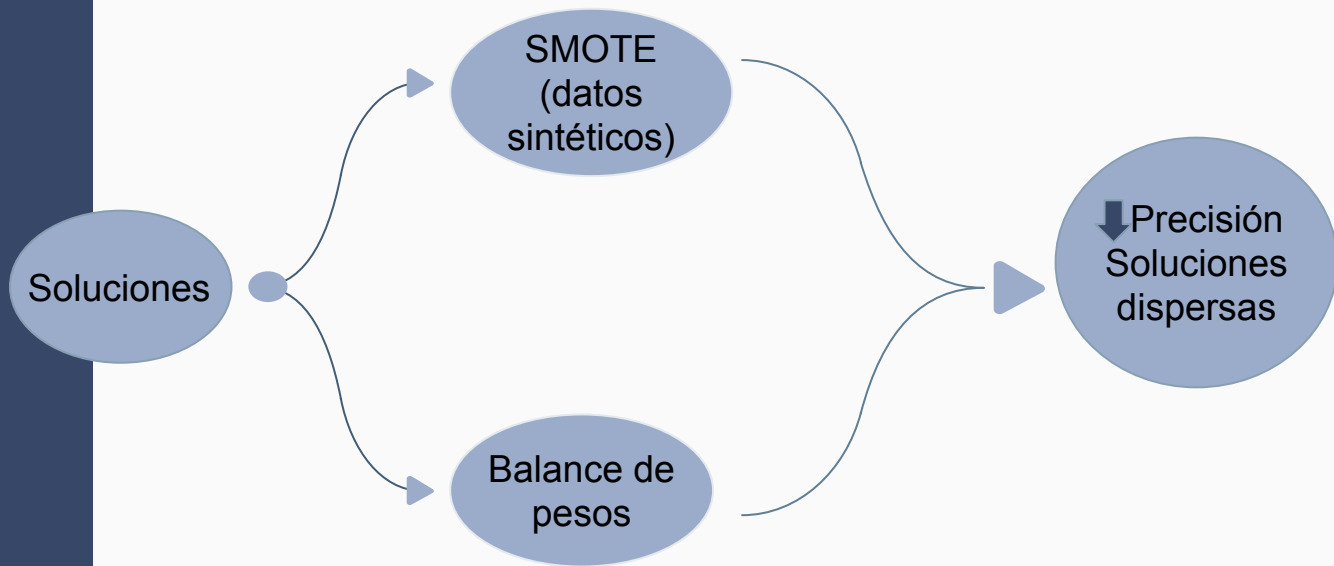
# Normalización

## ¿Resultados?

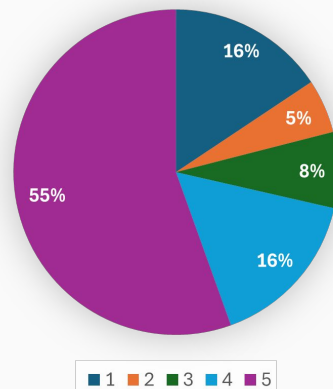
No se apreciaban cambios notables en los modelos con los que evaluamos, pero lo mantuvimos en mente para más adelante



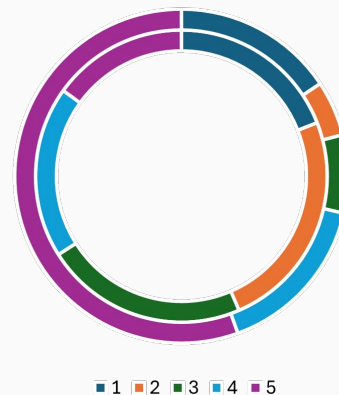
# Balanceo de clases



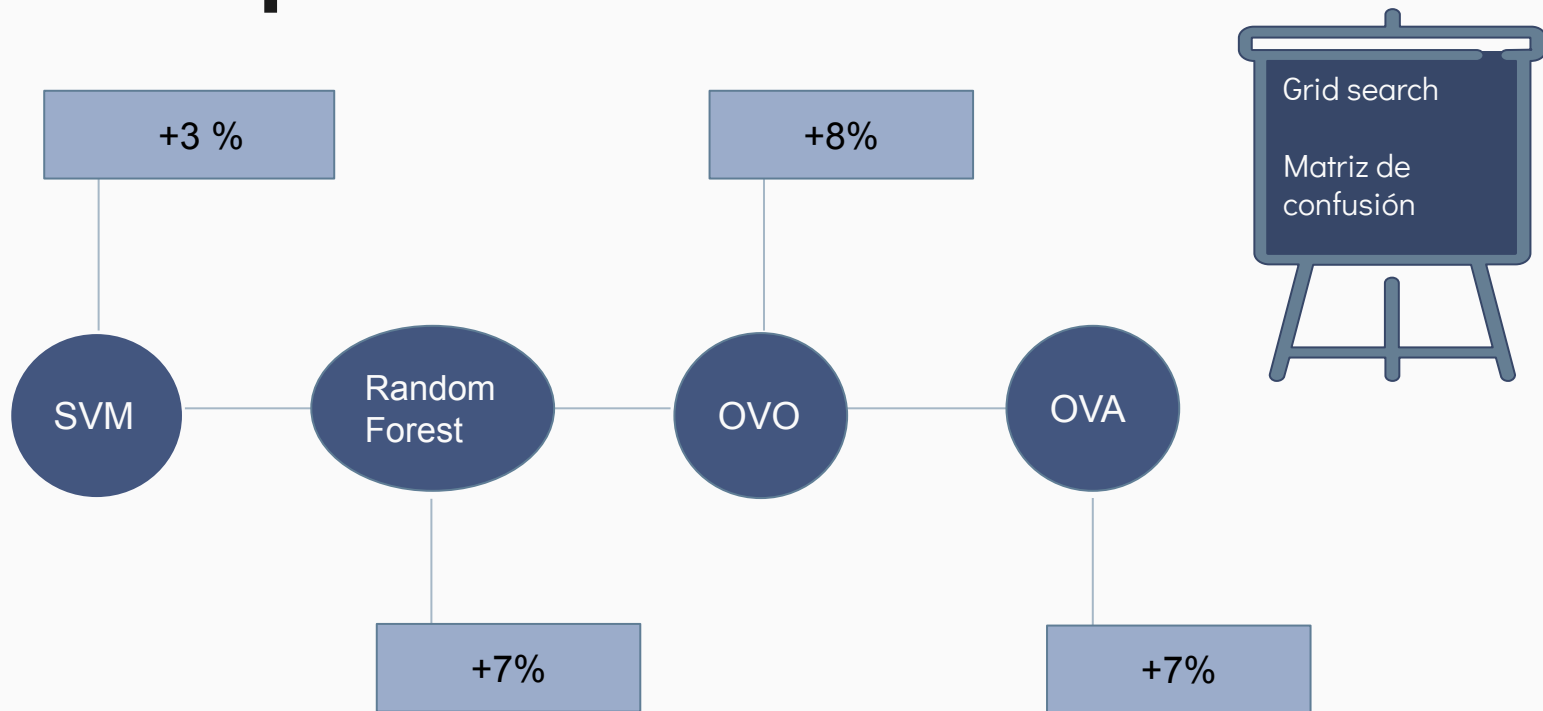
Repeticiones



Relación frecuencia - pesos



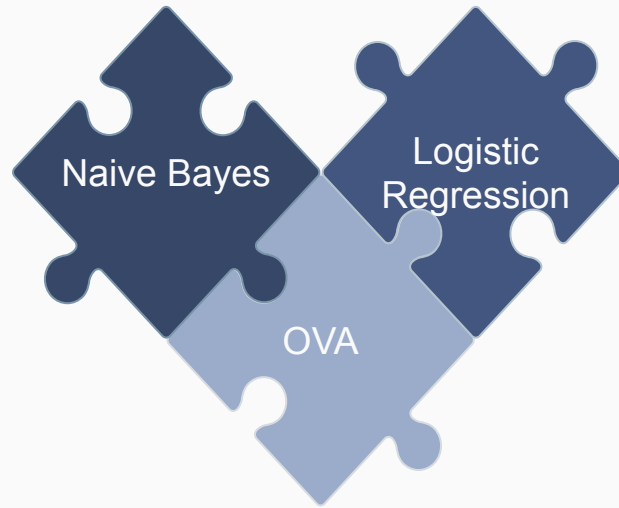
# Comprobación de modelos



# Ensemble



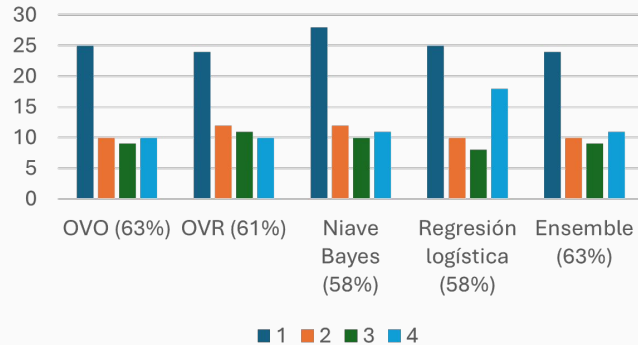
- Equilibrar clasificadores
- Distintos pesos



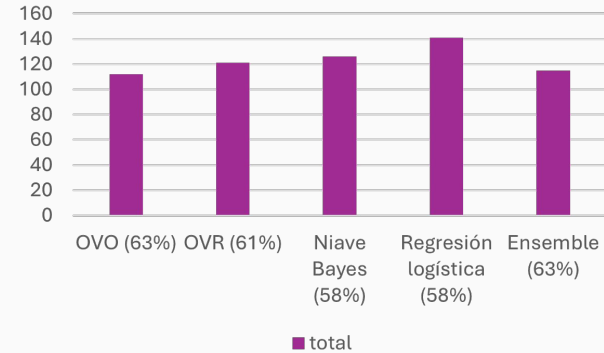
# Nueva métrica: Desviación de estrellas

total =  $\sum$  (nº de desviaciones \* valor de la desviación)

Desviaciones de estrellas



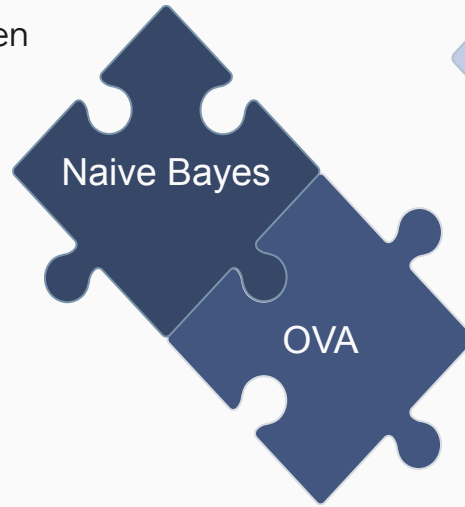
Desviaciones totales





# Nueva métrica: Desviación de estrellas

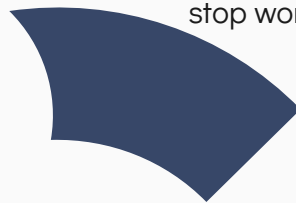
- Mayor precisión (64%)
- Menor desviación en predicciones
- OVA con apoyo de Naive Bayes



## **Conclusiones**

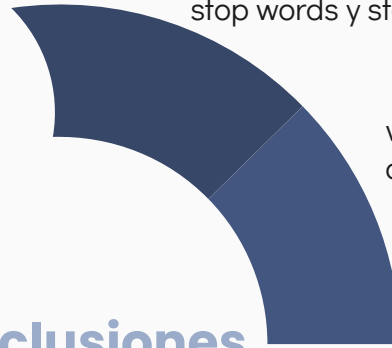
# Preprocesamiento

stop words y stemming



Conclusiones

# Preprocesamiento



stop words y stemming

vectorización TF-IDF sin n-gramas y  
con normalización L2

**Conclusiones**

# Preprocesamiento



# Preprocesamiento

stop words y stemming

vectorización TF-IDF sin n-gramas y  
con normalización L2

Conclusiones

**No balancear clases**

Afecta negativamente  
la precisión.

**Métricas**

Permite un análisis  
más detallado

## Preprocesamiento

stop words y stemming

vectorización TF-IDF sin n-gramas y  
con normalización L2

## Conclusiones

## No balancear clases

Afecta negativamente  
la precisión.

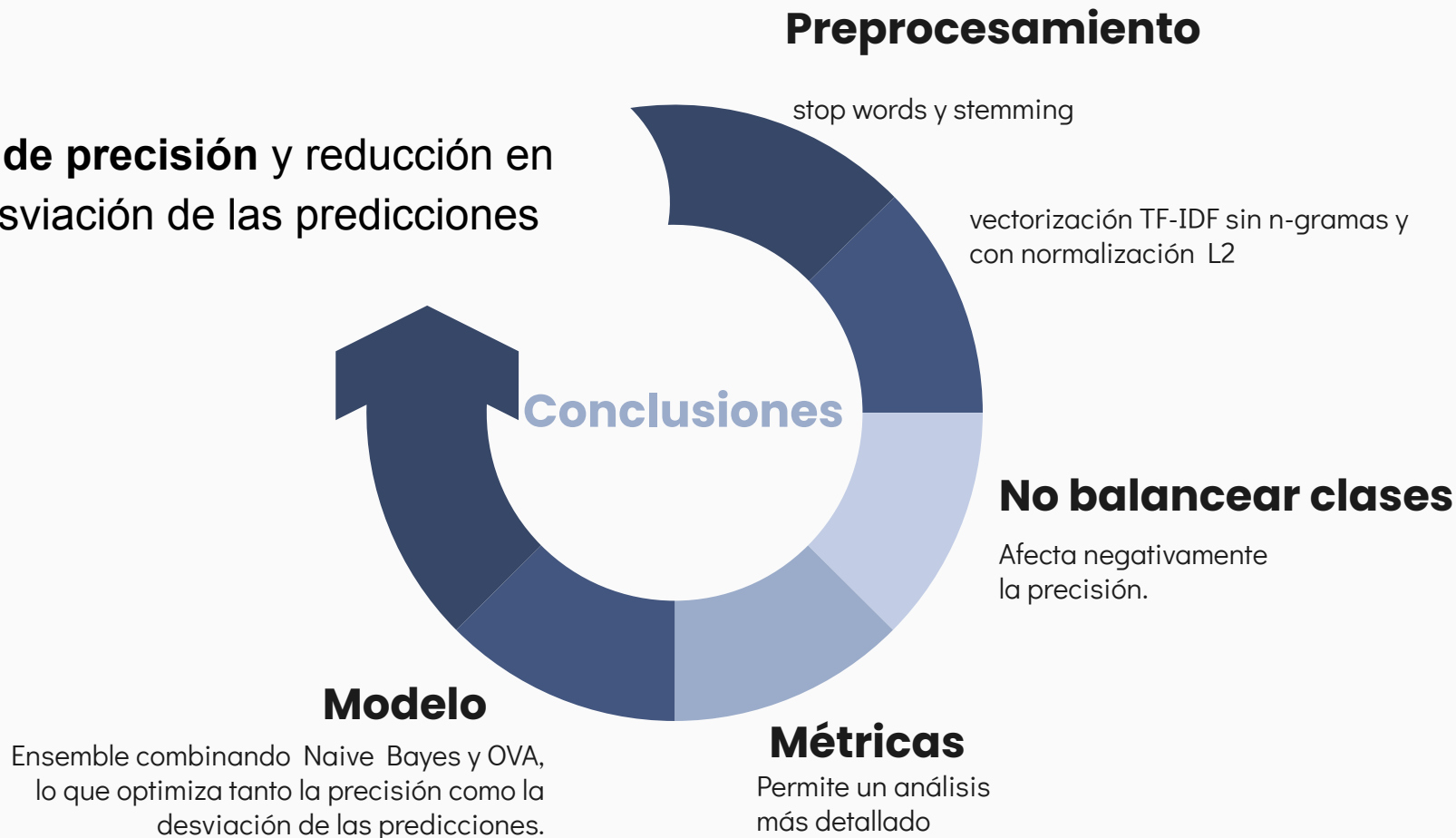
## Modelo

Ensemble combinando Naive Bayes y OVA,  
lo que optimiza tanto la precisión como la  
desviación de las predicciones.

## Métricas

Permite un análisis  
más detallado

**64% de precisión** y reducción en la desviación de las predicciones





# ¡Gracias!

---

¿Alguna pregunta?

