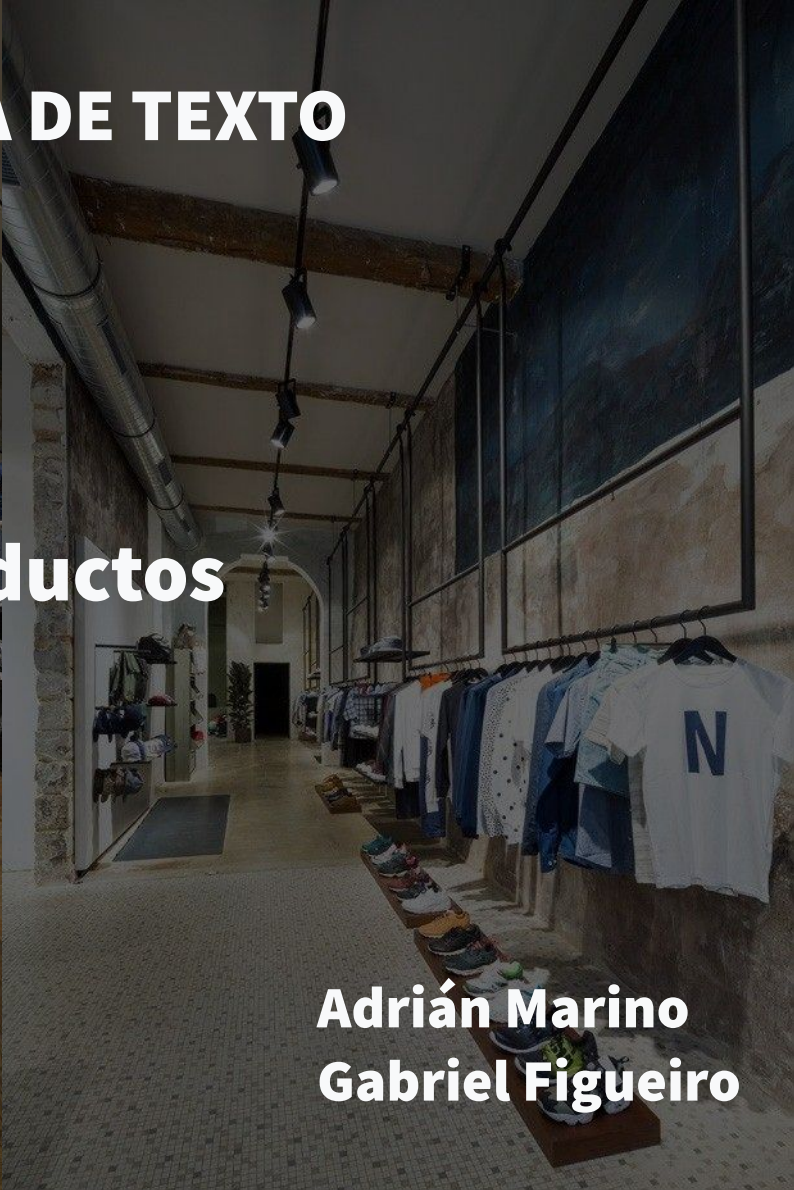
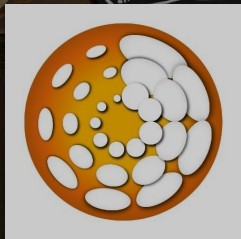


RECUPERACIÓN Y MINERÍA DE TEXTO

Clasificación de Productos



Adrián Marino
Gabriel Figueiro

Pregunta

- **Caso de Uso**
 - El usuario realiza el alta de un nuevo producto en un e-commerce para su posterior venta.
- **Pregunta**
 - ¿Es posible sugerir la categoría del producto?
- **Contexto**
 - Existe un árbol de categorías.
 - Cada nodo de hoja contiene productos.



Caso de uso: Mercado Libre

←

1

¿Qué querés publicar?

Productos

Vehículos

Inmuebles

Servicios

←

2

Indicá tu producto, marca y modelo

Tableta gráfica Wacom Intuos Small black

Evita incluir condiciones de venta. 41/60

Continuar

Tableta gráfica Wacom I...

1 2 3 4 5 6 7 8 9 0

Q W E R T Y U I O P

A S D F G H J K L

↑ Z X C V B N M ↵

!#1 ? < English (US) > . ↵

←

3

Confirmá la clasificación de tu producto

Computación

Periféricos de PC

Mouses y Teclados

Tabletas Digitalizadoras

Wacom

Intuos Small CTL-4100

Black

Tu publicación

Está bien clasificado

Quiero clasificarlo yo

←

¿Cuál es su condición?

Nuevo

Usado

Reacondicionado

←

¿Es este tu producto?

1 Revisá que este producto sea el que vendés para evitar cancelaciones que afecten tu reputación.

Tableta gráfica Wacom Intuos Small black

Usado

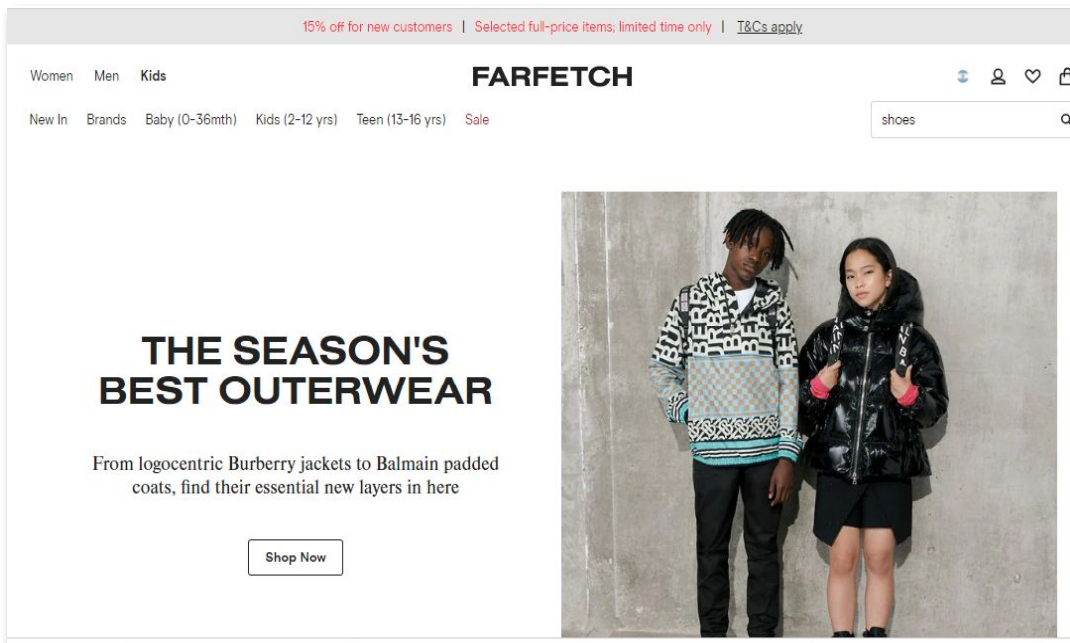
Ficha técnica

Marca	Wacom	Línea	Intuos
Modelo	Small	Modelo alfanumérico	CTL-4100
Color	Black	Largo del área de trabajo	152 mm
Ancho del área de trabajo	95 mm	Resolución de la imagen	2540 lpi
Niveles de sensibilidad de ...	4096	Es multi-touch	No

Continuar

Dataset: Farfetch

FARFETCH



Farfetch is a British-Portuguese online luxury fashion retail platform that sells products from over 700 boutiques and brands from around the world. The company was founded in 2007 by the Portuguese entrepreneur José Neves with its headquarters in London and main branches in Lisbon and Porto.

Wikipedia

Dataset disponible en:

<https://eval.ai/web/challenges/challenge-page/1721/overview>

Dataset: Farfetch

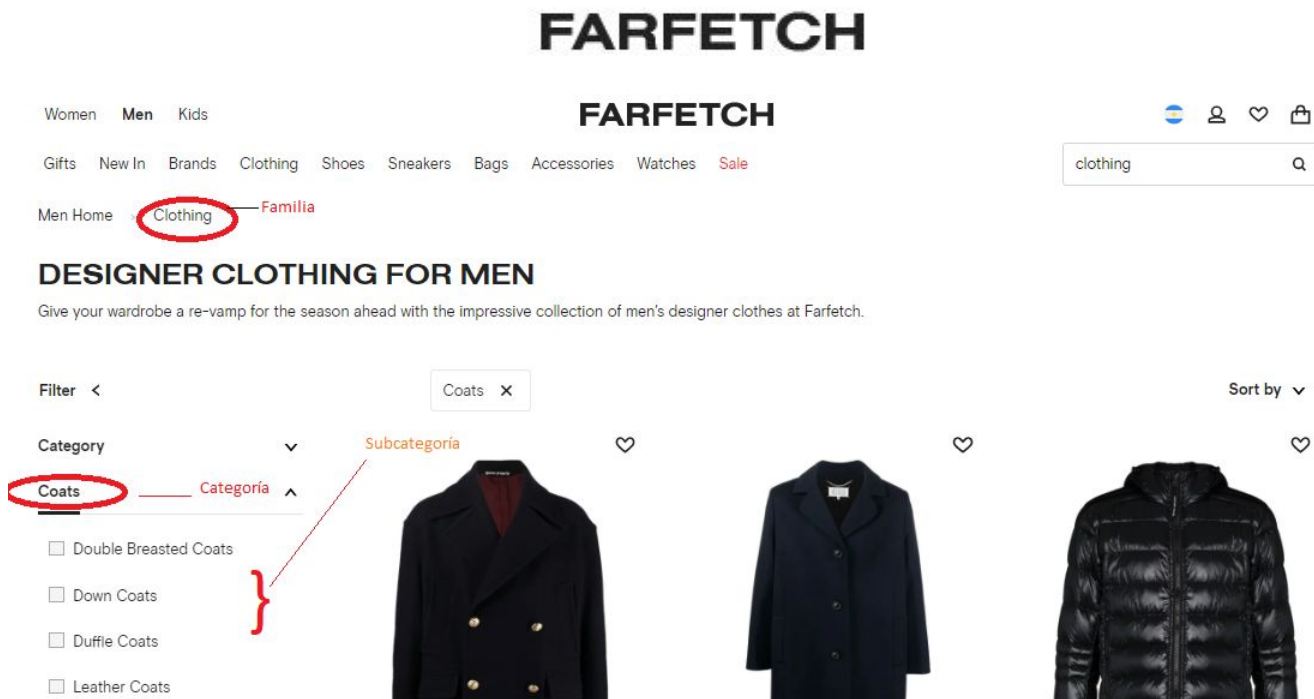
FARFETCH

Features

- product.id
- **product.gender**
- product.main_colour
- product.second_color
- **product.brand**
- **product.materials**
- **product.description**
- product.attributes
- **product.highlights**

Target

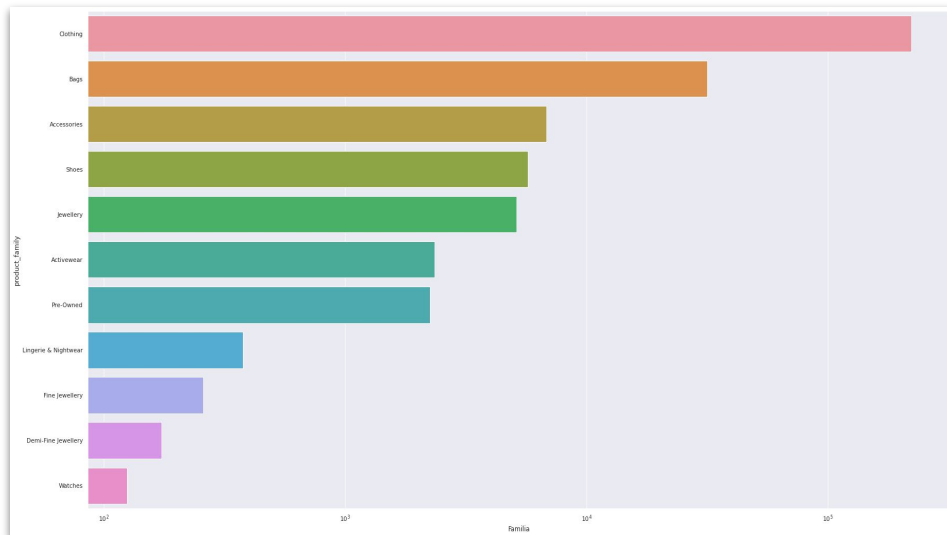
- Branch
 - family
 - category
 - subcategory



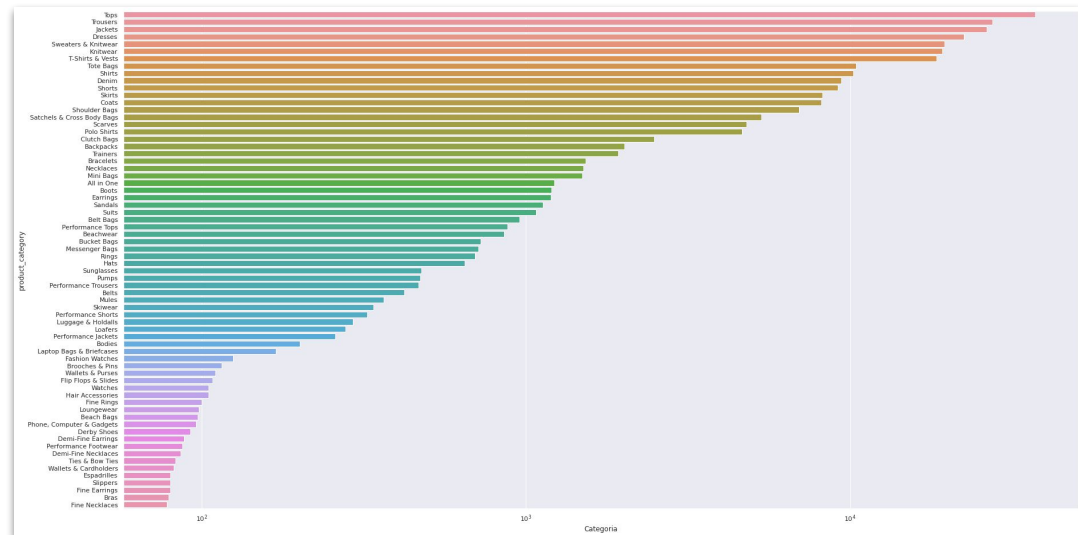
Dataset: Farfetch

FARFETCH

Familias de productos



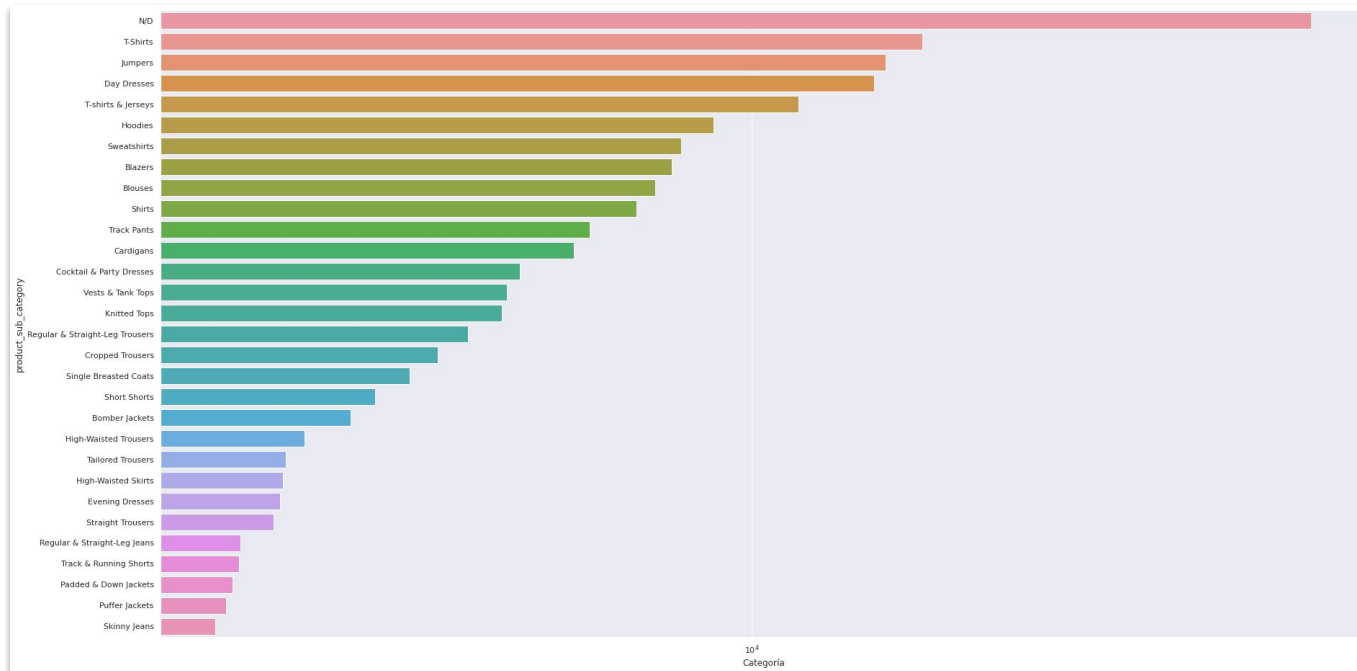
Categorías de productos



Dataset: Farfetch

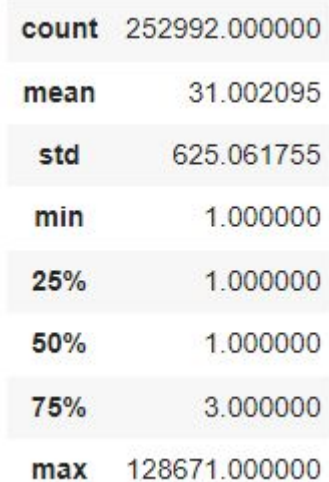
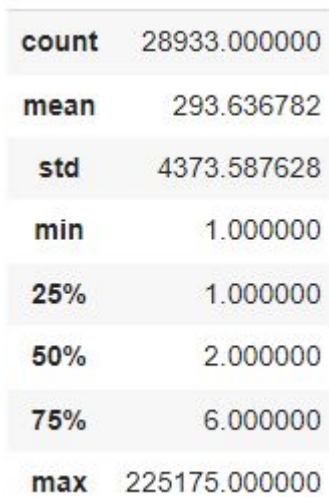
FARFETCH

Subcategorías de productos



FARFETCH

Bigramas



Dataset: Farfetch

FARFETCH

Nube de palabras por familia de productos



Lingerie & Nightwear



Underwear & Socks



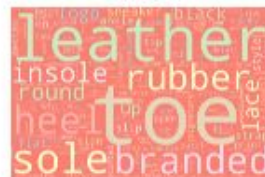
Homeware



Bags



Shoes



Fine Jewellery



Jewellery



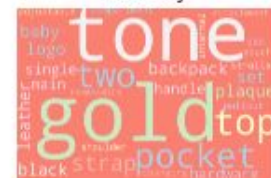
Demi-Fine Jewellery



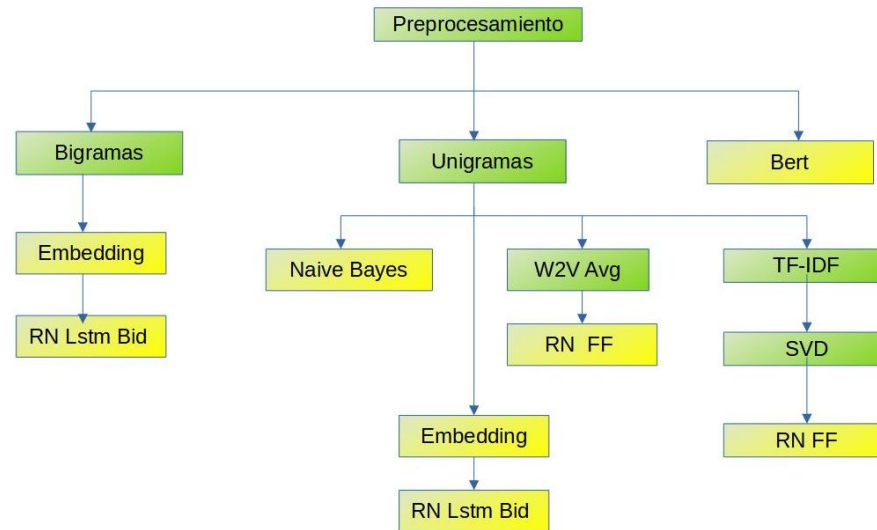
Accessories



Nursery



Modelos



Preprocesamiento

Concatenación de variables description + highlights + gender

Eliminación de:

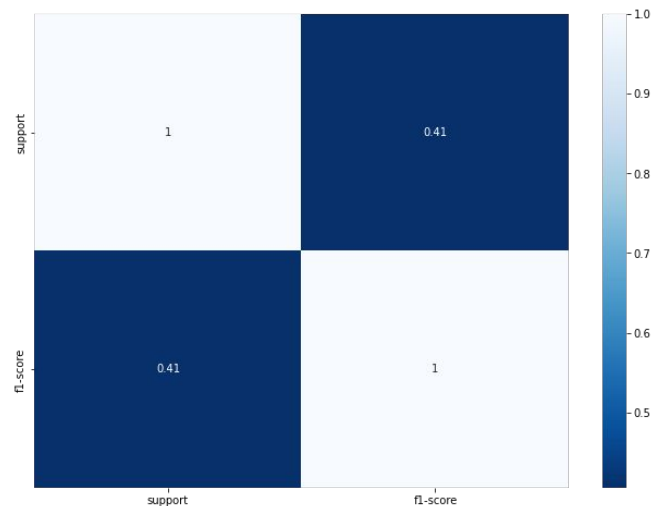
- Caracteres especiales
- Stop words en inglés
- Números
- Signos de puntuación

Modelos

- **Naive Bayes**

- Features: short description, gender, highlights
- Tokenizado con unigramas
- Frecuencia mínima de tokens: 10

F1-score macro: 0.53

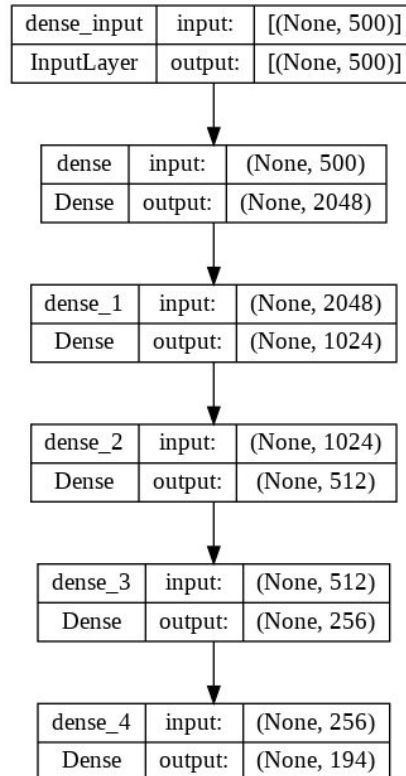


Modelos: Clasificador RN Feed Forward

- **Red Neuronal Feed Forward TF-IDF SVD**
 - Features: short description, gender, highlights
 - Tokenizado con unigramas
 - Frecuencia mínima del token: 10
 - TF-IDF
 - SVD de 500 componentes

Modelos: Clasificador RN Feed Forward

- Red Neuronal Feed Foward TFIDF SVD



F1-score: 0.64

Modelos: Clasificador RN Feed Forward

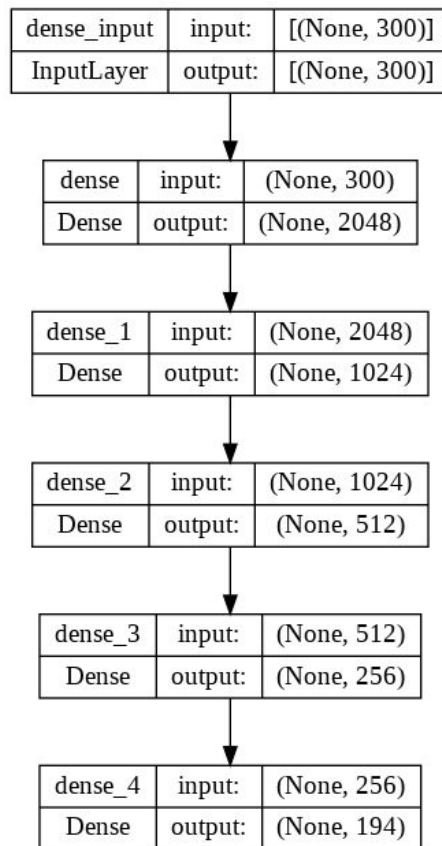
- **word2Vect Averange**

- Features: short description, gender, highlights
- Tokenizado con unigramas
- Frecuencia mínima del token: 20
- vocabulario de 2.800 tokens
- word2vect: ventana de 5 palabras y 30 iteraciones 300 dimensiones
- Promediamos los embeddings para obtener un sentence embedding

Modelos: Clasificador RN Feed Forward

- Red Neuronal Feed Foward

F1-score: 0.61

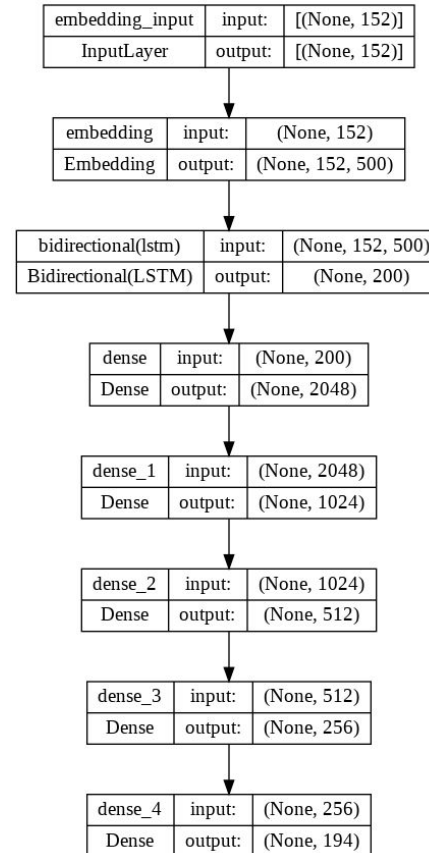


Modelos: Clasificador RNN LSTM Bidireccional Tokenizada con Unigramas

- **Red Neuronal LSTM Bidireccional con unigramas**
 - Features: short description, gender, highlights
 - Tokenizado con unigramas
 - Frecuencia mínima del token: 20
 - vocabulario de 2.800 tokens
 - Largo máximo de descripción de 152 tokens

Modelos: Clasificador RNN LSTM Bidireccional Tokenizada con Unigramas

F1-score: 0.66



Modelos: Clasificador RNN LSTM Bidireccional

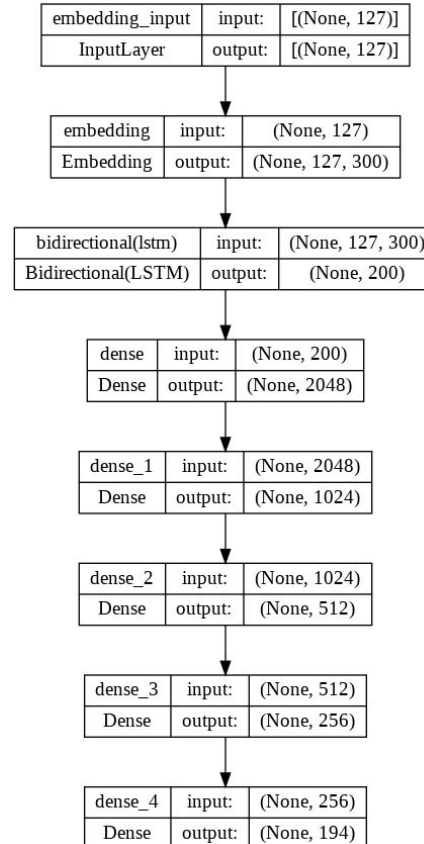
Tokenizada con Bigramas

- **Red Neuronal LSTM Bidireccional**
 - Features: short description, gender, highlights
 - Tokenizado con bigramas
 - Frecuencia mínima del token: 20
 - vocabulario de 8.541 tokens
 - Largo máximo de descripción de 152 tokens

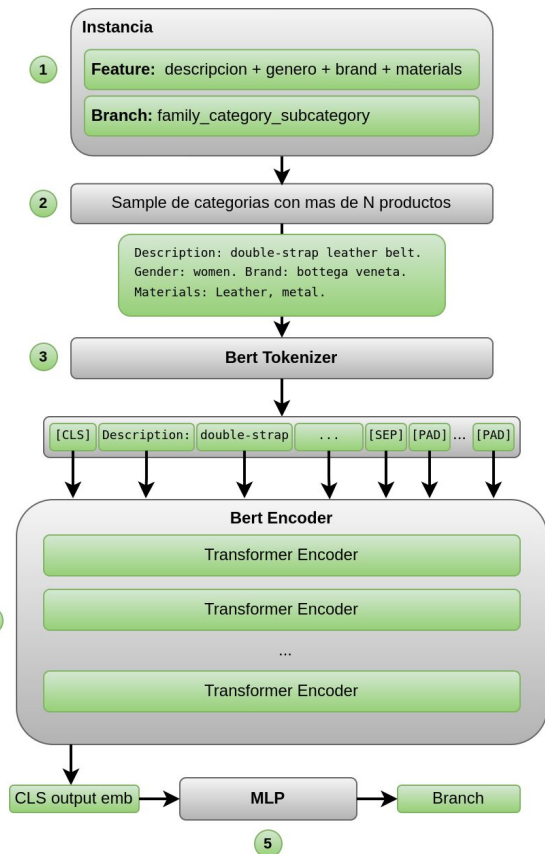
Modelos: Clasificador RNN LSTM Bidireccional

Tokenizada con bigramas

F1-score: 0.64



Modelos: Clasificador BERT



1. Variables

- Feature: Es la variable input del modelo.
 - Es una concatenación de variables.
 - Género y branch aumentaron el f1-score en un 10% y 1%.
- Branch: Es la clase a predecir.

2. Sampling: Para balancear el dataset sampleamos las categorías por encima de 1129 productos (Decil 0.7).

3. Tokenizer

- Generamos las secuencias con Bert Tokenizer.
- El tokenizer wrapa los tokens con CLS y SEP y completa con token PAD al máximo de secuencia definido (Max len: 39).
- Tamaño máximo de la secuencia de 512.

4. Encoder: Probamos un Encoder BERT pre-entrenado con los siguientes pesos :

- bert-base-cased (*)
- bert-base-uncased (*)
- distilbert-base-uncased

5. MLP

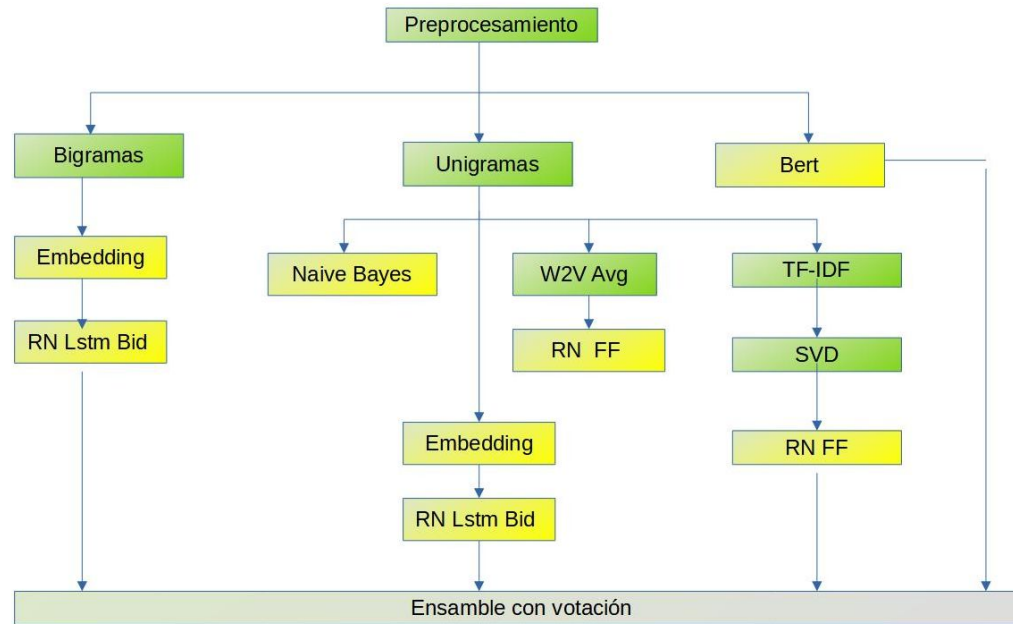
- Tomamos el embedding de la salida correspondiente al token CLS.
- Embedding de 768 dimensiones (en el caso del modelo de estos modelos).
- Este embedding es la entrada de una red MPL (Capa densa).
- La salida de la MLP pasa por una Softmax con igual número de salidas que clases(branch) en nuestro dataset.

Modelos: Clasificador BERT

Métricas (avg)

Por Nivel	Precision	Recall	F1-Score	Accuracy
Familia	0.85	0.88	0.86	0.99
Categoría	0.84	0.86	0.84	0.93
Subcategoría	0.69	0.71	0.70	0.81

Modelos: Ensamble por votación



F1-score: 0.70

Resultados Generales

F0.5-score

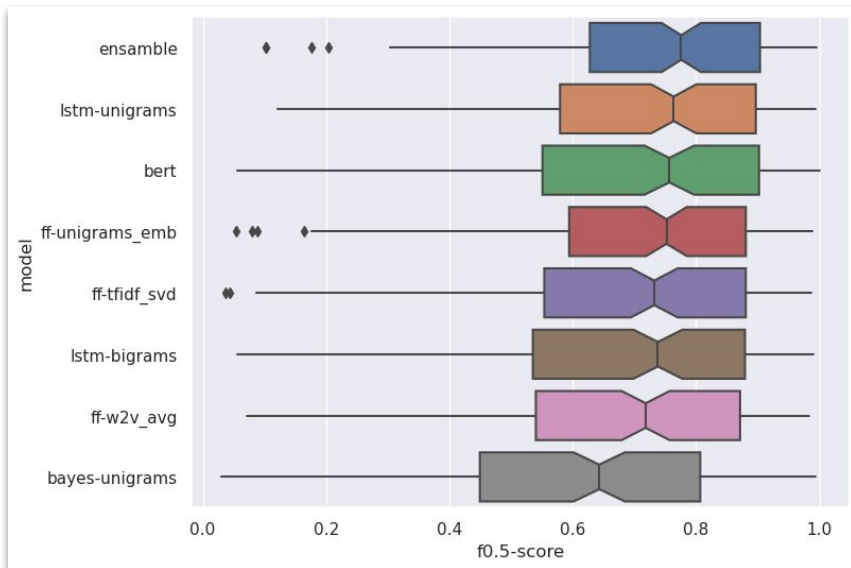


Tabla Comparativa

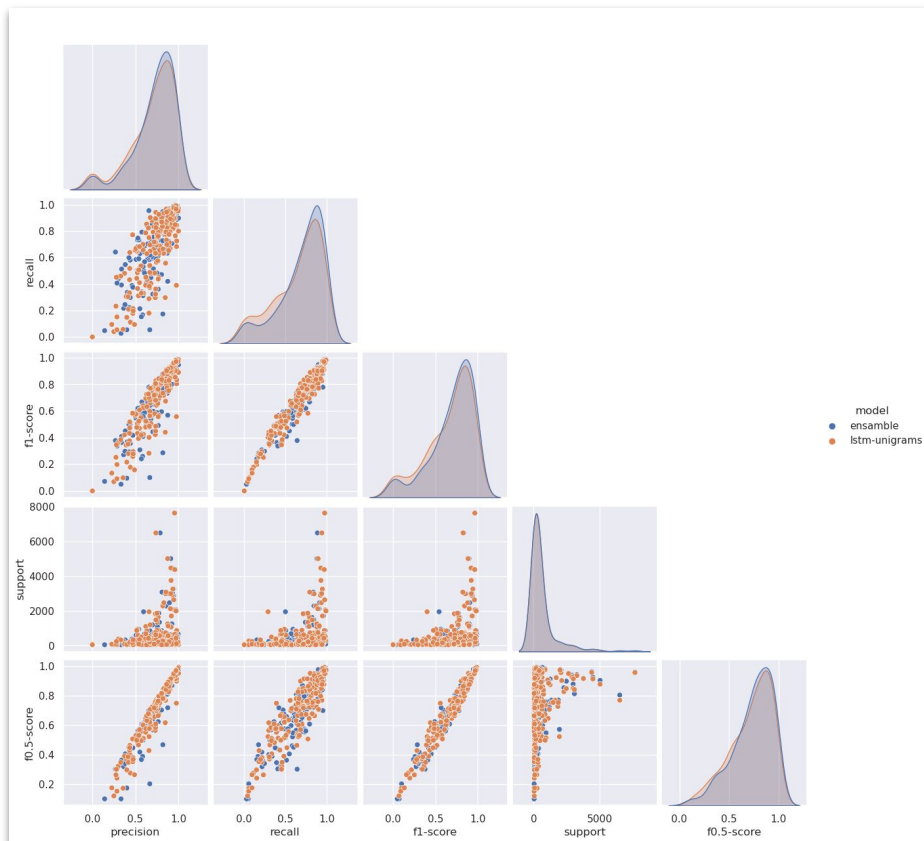
	precision	recall	f1-score	f0.5-score
model				
ensemble	0.73	0.70	0.70	0.74
lstm-unigrams	0.71	0.65	0.67	0.72
bert	0.69	0.73	0.70	0.71
ff-unigrams_emb	0.70	0.66	0.66	0.71
ff-tfidf_svd	0.67	0.63	0.64	0.70
lstm-bigrams	0.68	0.63	0.64	0.69
ff-w2v_avg	0.66	0.59	0.61	0.69
bayes-unigrams	0.63	0.51	0.53	0.60

Resultados



- **Distribución de F0.5-score por clase**
 - **Ensamble**
 - Está sesgada a la derecha.
 - Mayor cantidad de clases con f-score alto.
 - Pocas clases hasta 0.3 aprox.
 - **Bayes**
 - Menos sesgada a la derecha.
 - Menor concentración de clases para valores **alto** de f-score.
 - Mayor concentración de clases para valores **bajos** de f-score.
- **Precision**
 - **Ensamble**
 - Mayor número de clases con alta precision.
 - **Bayes**
 - Se distribuye en valores más bajos de precisión.
- **Recall**
 - **Bayes**
 - Tiende a capturar pocos ejemplos positivos por clase.
 - Mayor probabilidad en valores bajos de recall.

Resultados



- **F1-score**

- **Lstm-unigrams**

- Concentra más categorías en valores bajo de f0.5-score.

- **F0.5-score**

- **Lstm-unigrams**

- Concentra más categorías en valores bajo de f0.5-score.

- **Ensamble**

- Aumento importante en precision, recall en menor medida.

- **En resumen**

- Ambos modelos son muy similares en términos de precisión.
 - Bert gana en Recall.

	precision	recall	f1-score	f0.5-score
model				
ensemble	0.73	0.70	0.70	0.74
lstm-unigrams	0.71	0.65	0.67	0.72
bert	0.69	0.73	0.70	0.71

Proximos Pasos

- Analizar en detalle ejemplos mal clasificados.
- Probar con modelos bert-large-cased/uncased.
- Ensamble utilizando Probabilidades.
- Nuevos Features
 - main y secondary colors.
 - Imágenes de productos.
- Detectar productos mal asignados a una categoría.
- Fusionar ramas.
 - Categorías sinónimo.
 - Contienen el mismo tipo de productos.



¿Preguntas?