

UACM

Universidad Autónoma
de la Ciudad de México

NADA HUMANO ME ES AJENO

Redes Neuronales

Sabino Miranda

Clasificación de textos y Modelado del texto

Clasificación de texto

Clasificación de texto

El propósito es la **clasificación** de **documentos** en alguna categoría previamente definida.

Clasificación de texto

Clasificación de texto

El propósito es la **clasificación** de **documentos** en alguna categoría previamente definida.

Tareas

- Identificación de polaridad
(positivo, negativo, neutro)
- Clasificación de emociones
(felicidad, tristeza, ira, alegría, etc.)
- Clasificación de emociones
(intensidad de la emoción: 0, 1, 2, 3)
- Detección de agresividad
(agresivo, no agresivo)

Clasificación de texto

El propósito es la **clasificación** de **documentos** en alguna categoría previamente definida.

Tareas (cont.)

- Detección de comentarios tóxicos
(tóxico, obsceno, amenazante, insulto)
- Perfilado de autores
(hombre, mujer, edad)
- Variedad de leguaje
(Español: Argentina, Colombia, México, Perú, España)
- ...

Procedimiento

Esquema general de un clasificador

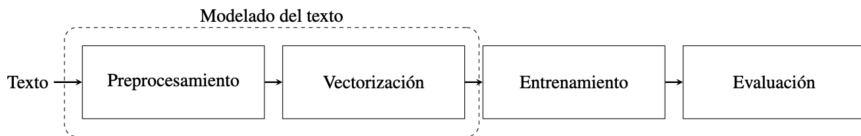


Figura 1: Diagrama a bloques del clasificador

Conjunto de entrenamiento: emociones

Texto	Clase
Pesa mas la #rabia q el #cemento	tristeza
Estoy cansada de que me prometan cosas y después no lo cumplan	tristeza
⋮	⋮
Estos pulmones asmáticos no me sirven para nada	enojo
⋮	⋮
Que te duela la panza de tanto reírte, eso no se paga con nada	alegría
⋮	⋮

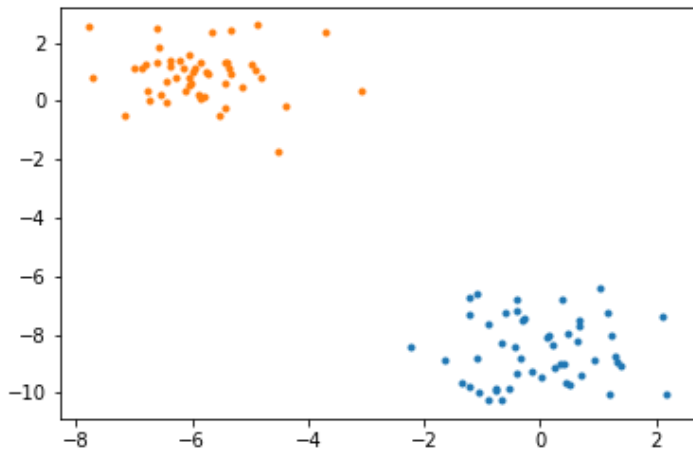
Conjunto de entrenamiento: polaridad

Texto	Clase
Que bien se siente volver a tener mi celular :D	positivo
@Dani3lS bueno en el galaxy s2 no es tan dificil ademas tengo dedos delgados y puntiagudos jajaja	positivo
:	:
@MovistarMX tengo un galaxy ace pero no puedo conectarme a internet de datos	negativo
:	:
Mi celular esta fallando): Dios concedeme otro celular :D	negativo
:	:

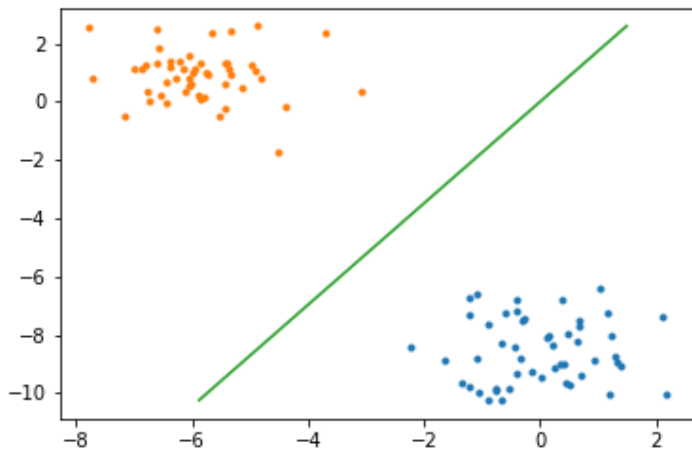
Aprendizaje supervisado



Aprendizaje supervisado



Función de decisión



Clases

- $\mathbf{x} \in \mathbb{R}^2$
- Naranja $\rightarrow negativo(-1)$
- Azul $\rightarrow positivo(1)$

Clasificador

$$f(x) = \text{sign}\left(\sum_i \alpha_i \mathbf{x}_i + \alpha_0\right)$$

Modelos de texto

Modelo de texto

Función que transforma el texto en un vector, i.e.,

$$m : \text{text} \rightarrow \mathbb{R}^d$$

Modelo de texto

Función que transforma el texto en un vector, i.e.,

$$m : \text{text} \rightarrow \mathbb{R}^d$$

Bolsa de palabras

- Modelo simple
- Empieza con un diccionario
- Cada palabra es una dimensión
- Coeficientes diferente de cero indican presencia

Ejemplo de Bolsa de Palabras

Corpus

Karla Ximena dijo el lunes que se siente como
Pero si Milan se siente como herido
Chris Evans dijo exactamente ...

Ejemplo de Bolsa de Palabras

Corpus

Karla Ximena dijo el lunes que se siente como
Pero si Milan se siente como herido
Chris Evans dijo exactamente ...

Representación vectorial

Karla	Ximena	dijo	lunes	que	se	siente	como	Chris	...
1	1	1	1	1	1	1	1	0	...
0	0	0	0	0	1	1	1	0	...
0	0	1	0	0	0	0	0	1	...

Características

- Simple de implementar
- Competitivo
 - Optimizar parámetros
 - Modelo de pesados (TF, TF-IDF, etc.)
 - Vectores dispersos

Características

- Simple de implementar
- Competitivo
 - Optimizar parámetros
 - Modelo de pesados (TF, TF-IDF, etc.)
 - Vectores dispersos

Limitantes - Orden

- Ana ama a Pablo
- Pablo ama a Ana

Características

- Simple de implementar
- Competitivo
 - Optimizar parámetros
 - Modelo de pesados (TF, TF-IDF, etc.)
 - Vectores dispersos

Limitantes - Orden

- Ana ama a Pablo
- Pablo ama a Ana

Limitantes - Partes de la Oración

- Sustantivos
- Verbos

Vectores densos (*Embeddings*)

Características

- Capturan la semántica y sintaxis
- Basados en predicción
- Vectores de menor dimensionalidad
- Vectores preentrenados con millones de datos
 - Word2Vec, Glove, FastTextm y Tipo BERT: RoBERTA, BETO, MarIA

Vectores densos (*Embeddings*)

Características

- Capturan la semántica y sintaxis
- Basados en predicción
- Vectores de menor dimensionalidad
- Vectores preentrenados con millones de datos
 - Word2Vec, Glove, FastTextm y Tipo BERT: RoBERTA, BETO, MarIA

Ejemplo

- “Que te duela la panza de tanto reírte, eso no se paga con nada”
- 300 dimensiones
[-0.02781,-0.03400,-0.01853,-0.00697,0.00401,0.007483 ...,]

- Máquinas de vectores de soporte (SVM), es un clasificador muy usado en identificación de polaridad[Tellez et al., 2017]. En Política, se usa UMAP para reducción de dimensionalidad y SVM [Cabrera-Pineda et al., 2023].

- Máquinas de vectores de soporte (SVM), es un clasificador muy usado en identificación de polaridad[Tellez et al., 2017]. En Política, se usa UMAP para reducción de dimensionalidad y SVM [Cabrera-Pineda et al., 2023].
- Modelos de redes neuronales es común su uso en tareas de perfilado y clasificación.
 - Identificación de ideología [Ahuir et al., 2023] combina modelos pre-entrenados BERT: BETO y MarIA; más esquema de votación (sumar predicciones parciales) para la predicción final.
 - En detección de odio, [Indurthi et al., 2019], usa representación de texto en *Universal Sentence Encoder* (basado en Transformers) y clasificador SVM con kernel RBF.

- Modelado del texto según la frecuencia del término (TF) y la frecuencia del término–Frecuencia inversa del documento (TF-IDF)

Ejercicio. Modelado del texto con pesos TF y TF-IDF

Consultar el notebook: *08_Modelado_del_texto.ipynb*

Red neuronal como clasificador de textos

Consultar el notebook: *09_Entrenamiento_Evaluacion.ipynb*

Referencias

- V. Ahuir, L. F. Hurtado, F. García-Granada, and E. Sanchis. ELiRF-VRain at PoliticES-IberLEF2023: Dealing with Long Texts in Transformer-based Systems for User Profiling. In *Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2023)*, CEUR Workshop Proceedings. CEUR-WS, 2023.
- H. Cabrera-Pineda, E. S. Tellez, and S. Miranda. Infotec-labd at politices-iberlef2023: Explainable non-linear low-dimensional projections. In *Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2023) co-located with the 39th Conference of the Spanish Society for Natural Language Processing (SEPLN 2023)*, Jaén, Spain, September, 2023, volume 3496 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2023.
- V. Indurthi, B. Syed, M. Shrivastava, N. Chakravartula, M. Gupta, and V. Varma. FERMI at SemEval-2019 task 5: Using sentence embeddings to identify hate speech against immigrants and women in Twitter. In J. May, E. Shutova, A. Herbelot, X. Zhu, M. Apidianaki, and S. M. Mohammad, editors, *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 70–74, Minneapolis, Minnesota, USA, June 2019. Association for Computational Linguistics.
- E. S. Tellez, S. Miranda-Jiménez, M. Graff, D. Moctezuma, O. S. Siordia, and E. A. Villaseñor. A case study of spanish text transformations for twitter sentiment analysis. *Expert Systems with Applications*, 81:457 – 471, 2017. ISSN 0957-4174.

Referencias (1)

- ❶ Deep Learning. Ian Goodfellow, Yoshua Bengio, and Aaron Courville. MIT Press, 2016.
<http://www.deeplearningbook.org>
- ❷ Dive into Deep Learning. Aston Zhang, Zachary C. Lipton, Mu li, and Alexander J. Smola. Cambridge University Press, 2023. <https://d2l.ai>
- ❸ Neural Networks and Deep Learning A Textbook (2nd Edition). Charu C. Aggarwal. Springer, 2023.
<https://doi.org/10.1007/978-3-031-29642-0>
- ❹ Deep Learning: Foundations and Concepts. Christopher M. Bishop and Hugh Bishop. Springer, 2024.
<https://doi.org/10.1007/978-3-031-45468-4>

Referencias (2)

- ⑤ PyTorch documentation.
`https://pytorch.org`
- ⑥ Numpy documentation.
`https://numpy.org`
- ⑦ Python documentation.
`https://www.python.org`