



# Procesamiento de Lenguaje Natural

## Clasificación

Mauricio Toledo-Acosta  
[mauricio.toledo@unison.mx](mailto:mauricio.toledo@unison.mx)

Departamento de Matemáticas  
Universidad de Sonora



## Section 1

# Introducción



- **Chapter 1.2, 1.4.** Eisenstein, J. (2018). Natural language processing. Jacob Eisenstein.
- **Chapter 4.** Jurafsky, D., Martin, J. H. (2019). Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition.



# Clasificación

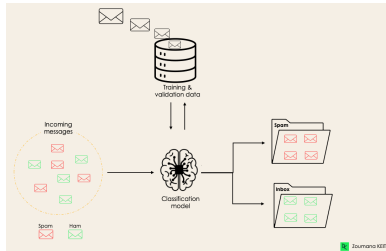
Procesamiento  
de Lenguaje  
Natural

Introducción

Clasificación  
en el PLN

## Clasificación

La **clasificación** consiste en asignar categorías a objetos. En el ML supervisado un modelo intenta predecir la etiqueta correcta de datos de entrada.



Source



- **Clasificación Binaria.** Cada dato tiene sólo una de dos posibles etiquetas.
- **Clasificación MultiClase.** Cada dato tiene sólo una de varias posibles etiquetas.
- **Clasificación MultiEtiqueta.** Cada dato tiene una o más de varias posibles etiquetas.

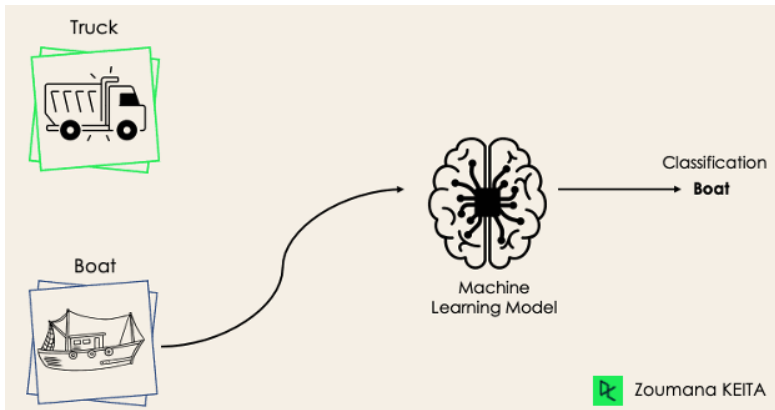


# Tipos de Clasificación: ejemplos

Procesamiento  
de Lenguaje  
Natural

Introducción

Clasificación  
en el PLN



Source

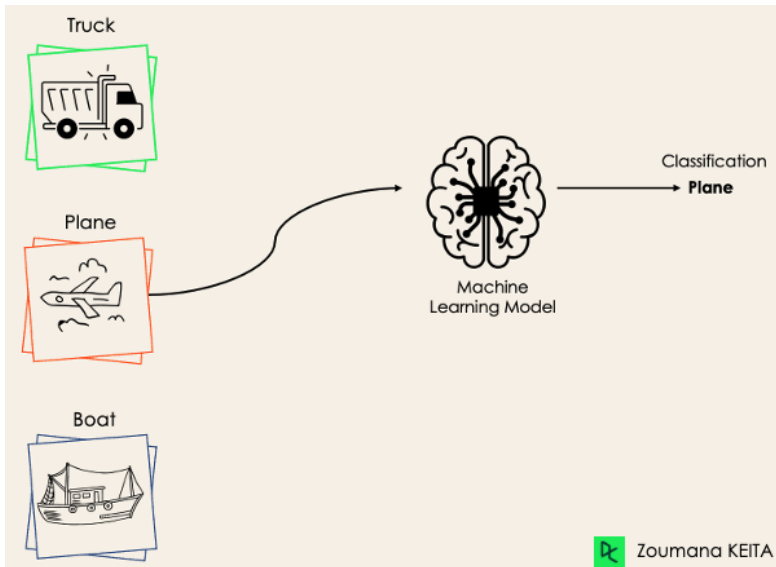


# Tipos de Clasificación: ejemplos

Procesamiento  
de Lenguaje  
Natural

Introducción

Clasificación  
en el PLN



Zoumana KEITA

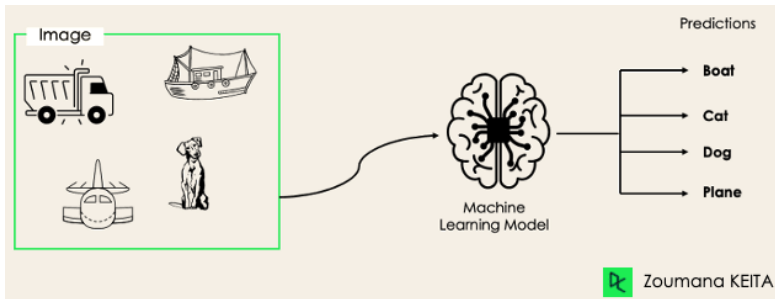


# Tipos de Clasificación: ejemplos

Procesamiento  
de Lenguaje  
Natural

Introducción

Clasificación  
en el PLN



Source





# Ejemplos de clasificación

- Clasificación de imágenes (identificar objetos en fotos).
- Diagnóstico médico (clasificar si un tumor es benigno o maligno). Esto puede ser por medio de imágenes, mediciones, etc.
- Reconocimiento de voz (identificar palabras habladas).
- Detección de fraude (identificar transacciones fraudulentas).
- Análisis de sentimientos (Identificar el sentimiento detrás de un texto).





# Algoritmos de Clasificación

Procesamiento  
de Lenguaje  
Natural

Introducción

Clasificación  
en el PLN

- Clasificación Lineal
  - SVM
  - Regresión Logística
  - Perceptrón
  - Clasificador de Mínimos Cuadrados
- Clasificación No lineal
  - Árboles de Decisión
  - Clasificadores de Ensamble
  - K-NN
  - Redes Neuronales
  - **Naive Bayes**

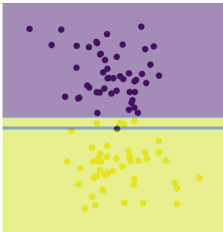


# Procesamiento de Lenguaje Natural

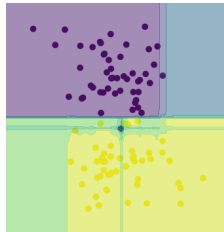
## Introducción

### Clasificación en el PLN

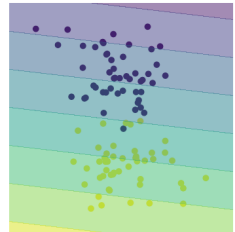
DecisionTreeClassifier



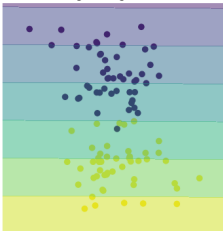
RandomForestClassifier



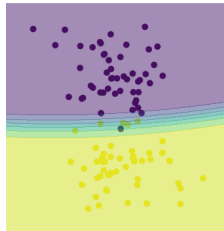
SVC



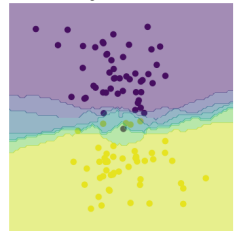
LogisticRegression



GaussianNB



KNeighborsClassifier







# Matriz de Confusión Binaria

Procesamiento  
de Lenguaje  
Natural

Introducción

Clasificación  
en el PLN

		Predicted condition	
		Positive (PP)	Negative (PN)
Actual condition	Positive (P)	True positive (TP)	False negative (FN)
	Negative (N)	False positive (FP)	True negative (TN)





14/27







# ROC-AUC Score

Procesamiento  
de Lenguaje  
Natural

Introducción

Clasificación  
en el PLN

## ROC-AUC Score

La curva paramétrica ROC (Receiver Operating Characteristic) muestra los valores FPR y TPR en varios valores de umbral de probabilidad. El **score AUC** es el area bajo la curva ROC, es una medida de rendimiento para los problemas de clasificación que representa el grado o medida de separabilidad. Indica la capacidad del modelo para distinguir entre clases.



# ROC-AUC Score

## ROC-AUC Score

La curva paramétrica ROC (Receiver Operating Characteristic) muestra los valores FPR y TPR en varios valores de umbral de probabilidad. El **score AUC** es el area bajo la curva ROC, es una medida de rendimiento para los problemas de clasificación que representa el grado o medida de separabilidad. Indica la capacidad del modelo para distinguir entre clases.

$$\text{FPR} = \frac{\text{FP}}{\text{TN} + \text{FP}}$$

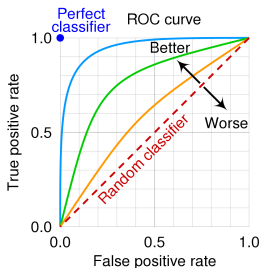
$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$



La curva paramétrica ROC (Receiver Operating Characteristic) muestra los valores FPR y TPR en varios valores de umbral de probabilidad. El **score AUC** es el area bajo la curva ROC, es una medida de rendimiento para los problemas de clasificación que representa el grado o medida de separabilidad. Indica la capacidad del modelo para distinguir entre clases.

$$\text{FPR} = \frac{\text{FP}}{\text{TN} + \text{FP}}$$

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}} \leftarrow \text{Recall}$$



El valor ROC-AUC es un número  $0 \leq s \leq 1$ . Entre más grande es  $s$ , el clasificador es mejor.

- Si  $s = 1$ , el clasificador es perfecto.
- Si  $s = \frac{1}{2}$ , el clasificador es aleatorio.
- Si  $s = 0$ , el clasificador predice perfectamente las clases *al revés*.



Umbral: 0.5



# Un ejemplo

Procesamiento  
de Lenguaje  
Natural

Introducción

Clasificación  
en el PLN

Umbral: 0.2

y_test	y_pred	probabilidades
0	0	0.048
0	0	0.145
1	1	0.905
0	1	0.24
1	1	0.215
0	1	0.231
0	0	0.116
1	1	0.551
1	0	0.172
1	1	0.803

$$\begin{pmatrix} 3 & 2 \\ 1 & 4 \end{pmatrix}, \quad TPR = 0.8, \quad FPR = 0.4$$



Umbral: 0.75



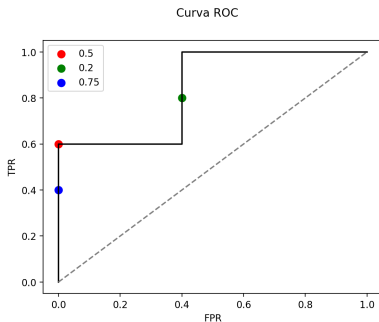


# Un ejemplo

Procesamiento  
de Lenguaje  
Natural

Introducción

Clasificación  
en el PLN



El área bajo la curva es 0.84.







# Clasificación en PLN

Procesamiento  
de Lenguaje  
Natural

Introducción

Clasificación  
en el PLN

- Detección de SPAM
- Identificación de idioma
- Atribución de autoría
- Detección de tópicos o temática
- Predicción de la siguiente palabra



# ¿Cómo hacemos la clasificación?

- **Métodos basados en reglas.** Las reglas suelen generarse por expertos en el dominio y pueden incluir coincidencias de palabras o secuencias específicas, patrones sintácticos, patrones léxicos (longitud de palabras, frecuencias).
  - Las reglas suelen ser frágiles.
  - Las reglas pueden cambiar con el tiempo.
  - + Eficientes y no requieren muchos datos.
  - + Transparentes y explicables.
- **Métodos de aprendizaje automático.**
- **Métodos Híbridos:** Lexicon-based, ...



Algunas palabras tienen una polaridad globalmente reconocida, por ejemplo *bueno*, *malo*, *perfecto*, *feo*. La simple presencia de una de estas palabras en un texto puede ser una pista importante sobre el sentimiento expresado.

*"A good tool that works perfectly"*

"I had an horrible experience"

En la clasificación Lexicon-based<sup>1</sup>, se crea una lista de palabras para cada etiqueta y se clasifica cada documento en función de cuántas palabras de cada lista están presentes.

<sup>1</sup>Taboada, M., J. Brooke, M. Tofiloski, K. Voll, and M. Stede (2011).  
Lexicon-based methods for sentiment analysis. *Computational linguistics* 37(2),  
267–307.



Existen diccionarios ya recopilados que asocian sentimientos a palabras:

- MPQA
- The General Inquirer lexicon
- SentiWordNet
- Appraisal lexicon

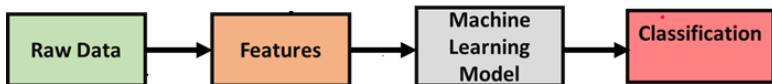


# Clasificación en el Machine Learning

Procesamiento  
de Lenguaje  
Natural

Introducción

Clasificación  
en el PLN



¿Quiénes son las *features* en el NLP?

- Métodos clásicos: BOW, TF-IDF, ...
- Embeddings: Redes Neuronales, LLMs.





- Features
- Raw Data
- Evaluación
- Rule-based / Lexicon-based / Machine Learning