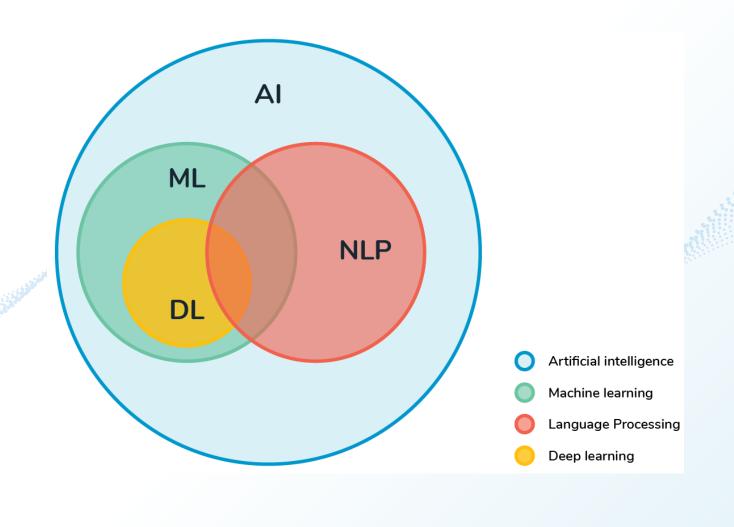
## Procesamiento de Lenguaje Natural

#### Contenido

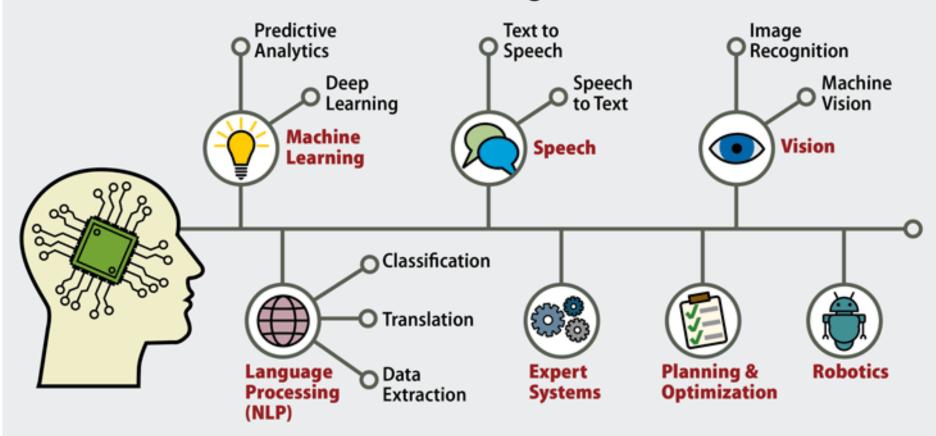
- Conceptos de Inteligencia Artificial / Machine Learning
- Conceptos de Procesamiento de Lenguaje Natural
- Conceptos de Análisis de Sentimientos





La Inteligencia Artificial (IA) es la rama de las ciencias de la computación que se enfoca en brindar a las máquinas o computadoras la capacidad de pensar de manera tan inteligente como los humanos y, en algunos casos, mejor que los humanos, aprendiendo de una gran cantidad de datos.

#### **Artificial Intelligence**



### **Machine Learning**

 El aprendizaje automático es un subcampo de la inteligencia artificial que utiliza algoritmos para aprender automáticamente cómo realizar una tarea determinada sin estar programado explícitamente con reglas.

### **Aplicaciones Machine Learning**

Voice assistants (Alexa, Siri, Google Home etc)



Recommendation Engines



Sentiment Analysis



Face detection and Recognition



Detecting credit card Fraud



Text summarization algorithms



#### Machine Learning: 3 Types of Learning

Background

Types of machine learning

Reinforcement Learning:

feedback to algorithm when it does something

Rewards & right or wrong

Recommendations algorithms

Example: Child gets feedback 'on the job' when it does something right or wrong

Machine Learning Reinforcement Learning

Supervised Learning:

pre-labelled data trains a model to predict new outcomes

> Example: Sorting LEGO blocks by matching them with the colour of the bags



Classification & Regression Algorithms

Unsupervised Learning



**Unsupervised Learning:** 

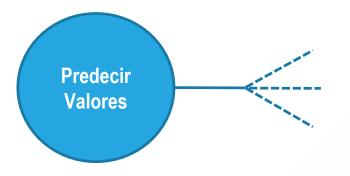
Non-labelled data self organises to predict new outcomes (e.g. clustering)

Clustering Algorithms

#### **Aprendizaje Supervisado**

#### Regresión

Predecir resultados futuros estimando relaciones entre variables



#### **Aplicaciones**

Estimar demanda de productos

Predecir gustos de clientes

Predecir ventas de una tienda

#### Clasificación

Identificar a qué categorías pertenecen nuevos valores



Predecir fraude en uso de Tarjeta de Crédito

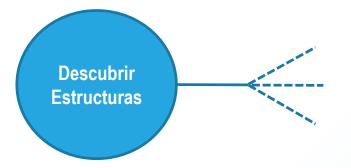
Predecir fallos en dispositivos

Diagnósticos médicos

#### Aprendizaje No Supervisado

#### Clustering

Separar en grupos con características similares



#### **Aplicaciones**

Agrupar clientes

Reconocimiento de Formas

Clasificación de Documentos

#### **Asociaciones**

Detectar relaciones entre las características de una instancia



Relacionar productos de la compra

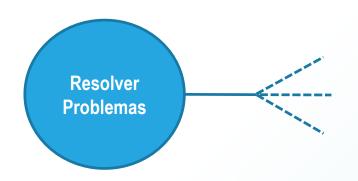
Relacionar gustos de entretenimientos de clientes

Relacionar productos contratados por cliente

#### Aprendizaje por Refuerzo

## Reinforcement Learning

Crea la mejor estrategia para obtener la mayor recompensa



#### **Aplicaciones**

Juegos

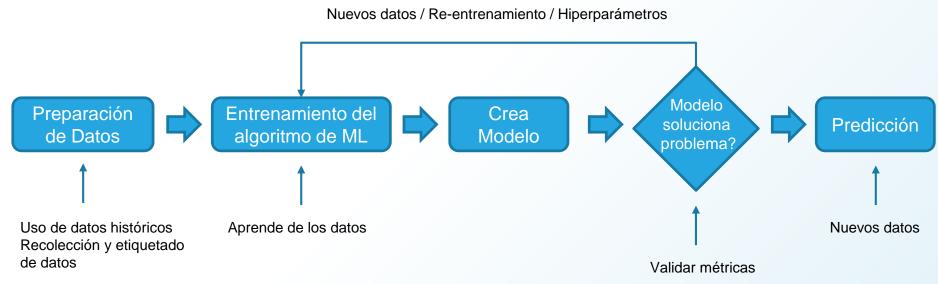
Control de Dispositivos (motores)

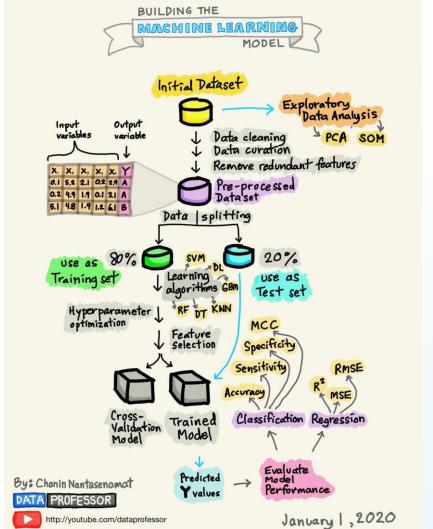
Robótica (control de movimiento)

## Ciclo de Vida – Trabajo de Investigación



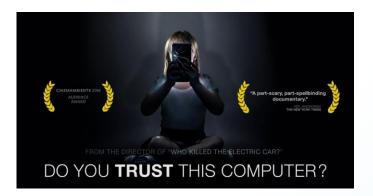
## Ciclo de Vida – Experimentos en Machine Learning





#### **Videos**





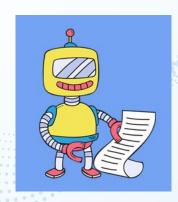


## 1.Definición

Qué es Procesamiento de Lenguaje Natural?

### Lenguaje Natural

- Comunicación por seres humanos
- Ejemplos: Inglés, Español, Portugués
- Procesamiento de Lenguaje Natural es cualquier manipulación computacional del lenguaje natural



### Tecnologías basadas en PLN

- Predicción de texto / Correctores Ortográficos (MS Word/otros editores de texto)
- Búsqueda en la web (Google, Bing, Yahoo)
- Clasificadores de Spam (Todos los servicios de email)
- Traducción automática (Google Translate)
- Asistentes virtuales (Siri, Alexa, Cortana)
- Análisis de sentimientos en tweets y blogs

#### **Técnicas**

- Segmentación de Oraciones
- Segmentación de Palabras (Tokenización)
- Limpieza de datos (stopwords)
- Análisis Lexicográfico (stemming, lematizing)
- Etiquetado Gramatical (POS)
- Chunking
- Reconocimiento de Entidades

#### **NLTK**

• Instalación:

```
import nltk
nltk.download('book')
```

Carga de textos:

```
from nltk.book import *
text9

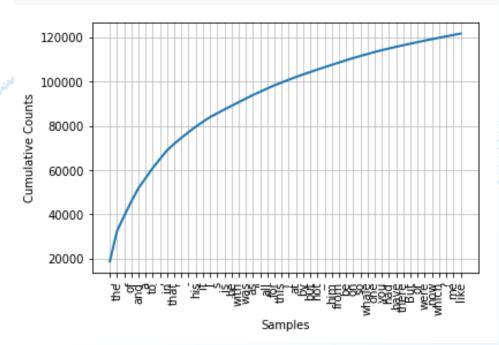
<Text: The Man Who Was Thursday by G . K . Chesterton 1908>
```

#### Contador de Vocabulario

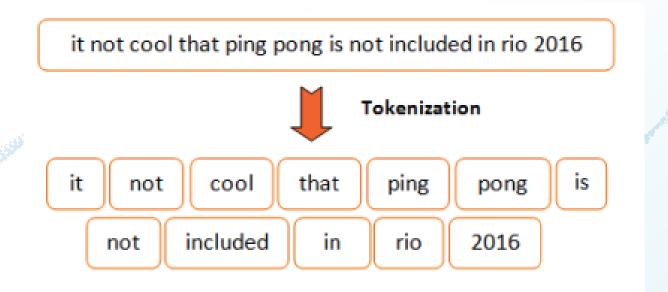
```
print('Cantiddad de símbolos:',len(text3))
Cantiddad de símbolos: 44764
print('Cantiddad de palabras',len(set(text3)))
Cantiddad de palabras 2789
text3.count("smote")
```

### Distribución de Frequencia

```
fdist = FreqDist(text1)
fdist.most_common(50)
fdist.plot(50, cumulative=True)
```



#### **Tokenización**



#### **Tokenization**

```
text = word_tokenize("And now for something completely different")
print(nltk.pos_tag(text))

Tokens: ['And', 'now', 'for', 'something', 'completely', 'different']
```

## Stopwords y Signos de Puntuación



Stop Words: remover palabras comunes pero que no proveen utilidad al descubrimiento del contexto (el, la, de, los, y, etc...)

#### **StopWords**

```
from nltk.corpus import stopwords

def removeStopword(texto):
    stopwordSpanish = stopwords.words('spanish')
    textNew= [w.lower() for w in texto if w not in stopwordSpanish]
    print(textNew)
```

```
['A', 'punto', 'de', 'cumplirse', 'una', 'semana', 'desde', 'que', 'Nicolás', 'Maduro', 'pusiera',
'en', 'marcha', 'las', 'nuevas', 'medidas', 'económicas', 'para', 'Venezuela', ',']
['a', 'punto', 'cumplirse', 'semana', 'nicolás', 'maduro', 'pusiera', 'marcha', 'nuevas',
'medidas', 'económicas', 'venezuela', ',', 'valor', 'bolívar', ',', 'fijado', 'tasa', '60',
'unidades']
```

### **Stemming - Lemmatization**

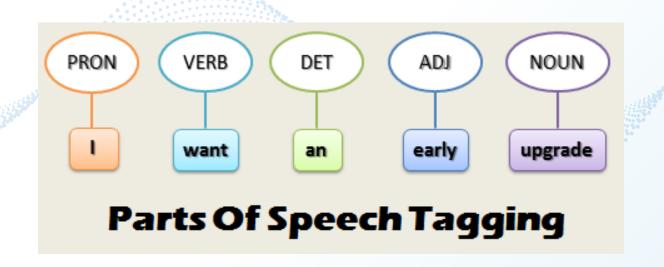
campo --> camp casita --> cas vendedor --> vend panadería --> pan comiendo --> comer limones --> limón corruptas --> corruptos nueces --> nuez

### Stemming y Lemmatization

```
tokens = word_tokenize(raw)
porter = nltk.PorterStemmer()
lancaster = nltk.LancasterStemmer()
print('PorterStemmer',[porter.stem(t) for t in tokens])
print('LancasterStemmer',[lancaster.stem(t) for t in tokens])
wnl = nltk.WordNetLemmatizer()
print('WordNetLematizer',[wnl.lemmatize(t) for t in tokens])
```

```
Porter Stemmer ['la', 'familia', 'suelen', 'escond', 'a', 'lo', 'niño', 'afectado', 'o', 'lo', 'aíslan', 'con', 'lo', 'animal']
Lancaster Stemmer ['las', 'familia', 'suel', 'escond', 'a', 'los', 'niño', 'afectado', 'o', 'los', 'aísl', 'con', 'los', 'anim']
Lemmatization ['Las', 'familias', 'suelen', 'esconder', 'a', 'los', 'niños', 'afectados', 'o', 'los', 'aíslan', 'con', 'los', 'animales']
```

## **Etiquetado Gramatical Tagging Words**

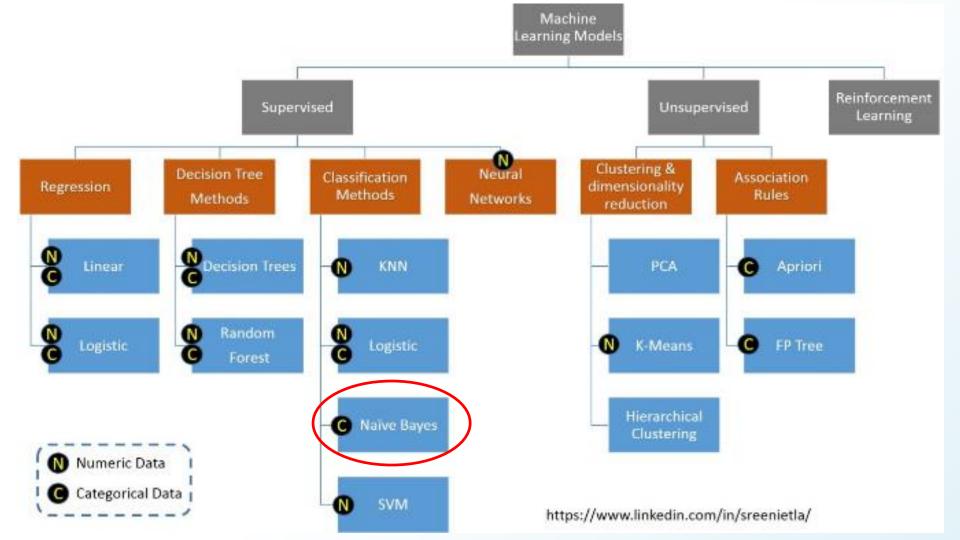


### Etiquetador (tagger)

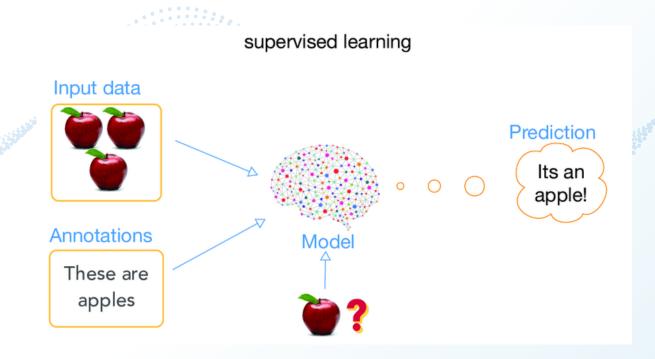
```
text = word tokenize("And now for something completely different")
nltk.pos tag(text)
[('And', 'CC'),
 ('now', 'RB'),
 ('for', 'IN'),
 ('something', 'NN'),
 ('completely', 'RB'),
 ('different', 'JJ')]
nltk.corpus.brown.tagged_words(tagset='universal')
[('The', 'DET'), ('Fulton', 'NOUN'), ...]
```

# 2. Clasificación Supervisada

Y su uso para el Análisis de Sentimientos



## Aprendizaje Automático Clasificación Supervisada

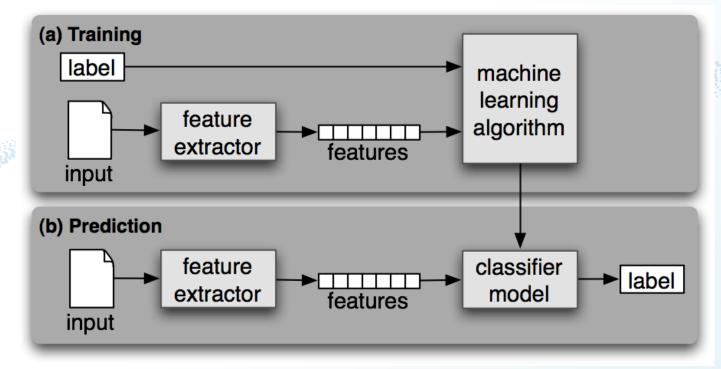


#### **Análisis de Sentimientos**

This is a good book! Postive
This is a awesome book! Postive
This is a bad book! Negative
This is a terrible book Negative



# Aprendizaje Automático Clasificación Supervisada



## Aprendizaje Automático Clasificación Supervisada

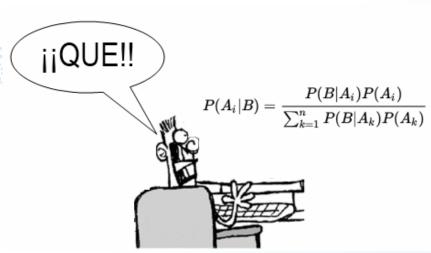
Características (features)



#### **Técnicas**

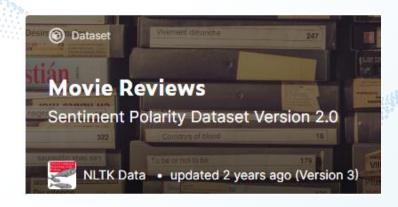
- Tokenización
- Lematización
- POS
- NEF

# **Naive Bayes**



### **Movie Reviews**

1000 archivos
 etiquetados como
 positivos y otros
 1000 etiquetados
 como negativos.



### **Movie Reviews**

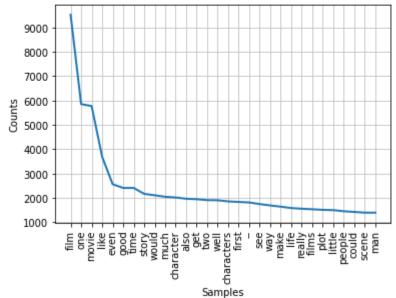
```
def Estadisticas():
    print ('total',len(movie_reviews.fileids()))
    print ('categorias',movie_reviews.categories())
    print ('total positivos',len(movie_reviews.fileids('pos')))
    print ('total negativos',len(movie_reviews.fileids('neg')))

all_words = [word.lower() for word in movie_reviews.words()]
    all_words_frequency = FreqDist(all_words)
    print ('10 palabras más frecuentes',all_words_frequency.most_common(10))
    print ('cantidad de veces que se repite la palabra happy',all_words_frequency['happy'])
```

```
total 2000
categorias ['neg', 'pos']
total positivos 1000
total negativos 1000
10 palabras más frecuentes [(',', 77717), ('the', 76529), ('.', 65876), ('a', 38106),
('and', 35576), ('of', 34123), ('to', 31937), ("'", 30585), ('is', 25195), ('in', 21822)]
cantidad de veces que se repite la palabra happy 215
```

### **Movie Reviews**

```
total positivos 1000
total negativos 1000
1583820
710578
10 palabras más frecuentes [('film', 9517), ('one', 5852);
cantidad de veces que se repite la palabra happy 215
```





- Bag of Words method converts text data into numbers.
- It does this by
  - Creating a vocabulary from the words in all documents
  - · Calculating the occurrences of words:
    - o binary (present or not)
    - o word counts
    - o frequencies

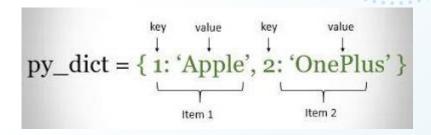
Simple example using word counts:

		a	cat	dog	happy	is	it	my	not	old	wolf
	"It is a dog."	1	0	1	0	1	1	0	0	0	0
	"my cat is old"	0	1	0	0	1	0	1	0	1	0
	"It is not a dog, it a is wolf."	2	0	1	0	2	2	0	1	0	1

```
def bag_of_words(words):
    words_clean = []
    stopwords_english = stopwords.words('english')

for word in words:
    word = word.lower()
    if word not in stopwords_english and word not in string.punctuation:
        words_clean.append(word)

words_dictionary = dict([word, True] for word in words_clean)
    return words_dictionary
```



```
def DatosBOW():
    pos reviews = []
    for fileid in movie reviews.fileids('pos'):
        words = movie reviews.words(fileid)
        pos reviews.append(words)
    neg reviews = []
    for fileid in movie reviews.fileids('neg'):
        words = movie reviews.words(fileid)
        neg reviews.append(words)
    pos_reviews_set = []
    for words in pos reviews:
        pos reviews set.append((bag of words(words), 'pos'))
    neg reviews_set = []
    for words in neg reviews:
        neg reviews set.append((bag of words(words), 'neg'))
    return pos reviews set, neg reviews set
```

POSITIVO [({'plot': True, 'two': True, 'teen': True, 'couples': True, 'go': True, 'church': True, 'party': True, 'drink': True, 'drive': True, 'get': True, 'accident': True, 'one': True, 'guys': True, 'dies': True, 'girlfriend': True, 'continues': True, 'see': True, 'life': True, 'nightmares': True, 'deal': True, 'nightmares': True, 'deal': True, 'watch': True, 'movie': True, 'sorta': True, 'find': True, 'critique': True, 'find': True, 'fuck': True, 'generation': True, 'touches': True, 'cool': True, 'idea': True, 'presents': True,

NEGATIVO [({'plot': True, 'two': True, 'teen': True, 'couples': True, 'go': True, 'church': True, 'party': True, 'drink': True, 'drive': True, 'get': True, 'accident': True, 'one': True, 'guys': True, 'dies': True, 'girlfriend': True, 'continues': True, 'see': True, 'life': True, 'nightmares': True, 'deal': True, 'watch': True, 'movie': True, 'sorta': True, 'find': True, 'critique': True, 'mind': True, 'fuck': True, 'generation': True, 'fuck': True, 'generation': True,

```
def clasificadorBOW(pos_reviews_set,neg_reviews_set):
    size = int(len(pos_reviews_set) * 0.1)
    test_set = pos_reviews_set[:size] + neg_reviews_set[:size]
    train_set = pos_reviews_set[size:] + neg_reviews_set[size:]
    print(len(test_set), len(train_set))
    classifier = NaiveBayesClassifier.train(train_set)
    accuracy = classify.accuracy(classifier, test_set)
    print(accuracy)
    print (classifier.show_most_informative_features(10))
    return classifier
```

```
def pruebaBOW(custom_review,classifier):
    custom_review_tokens = word_tokenize(custom_review)
    custom_review_set = bag_of_words(custom_review_tokens)
    print (custom_review)
    print ('positivo o negativo?',classifier.classify(custom_review_set))
    prob_result = classifier.prob_classify(custom_review_set)
    print ('probabilidad para negativo',prob_result.prob("neg"))
    print ('probabilidad para positivo',prob result.prob("pos"))
```

I hated the film. It was a disaster. Poor direction, bad acting. positivo o negativo? neg probabilidad para negativo 0.8389515850310557 probabilidad para positivo 0.16104841496894456
It was a wonderful and amazing movie. I loved it. Best direction, good acting. positivo o negativo? pos probabilidad para negativo 0.06519158529235963 probabilidad para positivo 0.9348084147076396

### Matriz de Confusión

# Realidad Positivos Negativos TP FP Negativos FN TN

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

Precision = 
$$\frac{TP}{TP + FP}$$

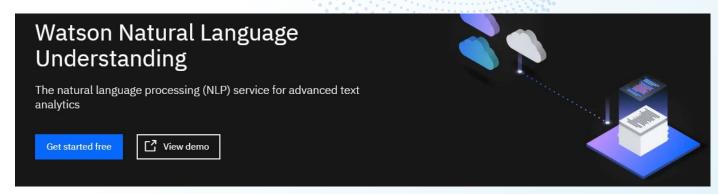
$$Recall = \frac{TP}{TP + FN}$$
 (sensitivity)

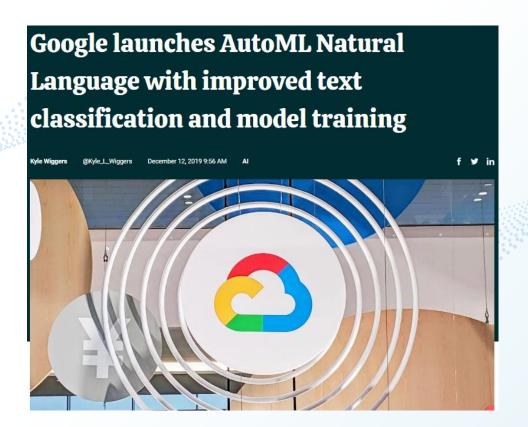
$$F - score = \frac{2 \times (Precision \times Recall)}{(Precision + Recall)}$$

Specificity 
$$= \frac{TN}{TN + FP}$$

Qué puedo usar?









Deep Learning NLP with spaCy











# Gracias!

# Alguna pregunta?

Maria.limaylla@gmail.com