



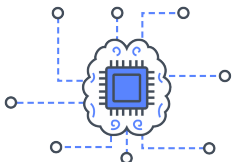
# Text Mining

Marcelo Mendoza

<http://www.inf.utfsm.cl/~mmendoza>

[mmendoza@inf.utfsm.cl](mailto:mmendoza@inf.utfsm.cl)

A 131, Campus San Joaquín - UTFSM

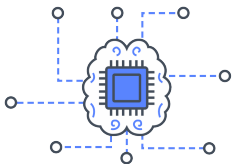




## Conceptos de NLP

### POS tagging

- ▶ Etiquetar cada término de acuerdo a la función que este cumple en el texto.
- ▶ Puede ayudarnos en tareas como detección de estilo, parsing, detección de colocaciones.
- ▶ Tarea importante en NLP.





## Conceptos de NLP

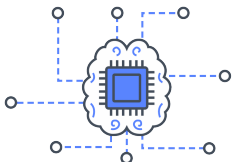
### POS tagging

- ▶ Etiquetar cada término de acuerdo a la función que este cumple en el texto.
- ▶ Puede ayudarnos en tareas como detección de estilo, parsing, detección de colocaciones.
- ▶ Tarea importante en NLP.

Text:

John likes the blue house at the end of the street .

Adjective	Determiner	Preposition
Adverb	Noun	Pronoun
Conjunction	Number	Verb

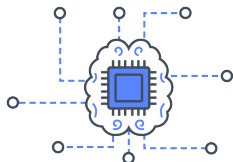




# Conceptos de NLP

## POS tagging

	Tag	Description	Example
Open Class	<b>ADJ</b>	Adjective: noun modifiers describing properties	<i>red, young, awesome</i>
	<b>ADV</b>	Adverb: verb modifiers of time, place, manner	<i>very, slowly, home, yesterday</i>
	<b>NOUN</b>	words for persons, places, things, etc.	<i>algorithm, cat, mango, beauty</i>
	<b>VERB</b>	words for actions and processes	<i>draw, provide, go</i>
	<b>PROPN</b>	Proper noun: name of a person, organization, place, etc..	<i>Regina, IBM, Colorado</i>
	<b>INTJ</b>	Interjection: exclamation, greeting, yes/no response, etc.	<i>oh, um, yes, hello</i>
Closed Class Words	<b>ADP</b>	Adposition (Preposition/Postposition): marks a noun's spacial, temporal, or other relation	<i>in, on, by under</i>
	<b>AUX</b>	Auxiliary: helping verb marking tense, aspect, mood, etc.,	<i>can, may, should, are</i>
	<b>CCONJ</b>	Coordinating Conjunction: joins two phrases/clauses	<i>and, or, but</i>
	<b>DET</b>	Determiner: marks noun phrase properties	<i>a, an, the, this</i>
	<b>NUM</b>	Numeral	<i>one, two, first, second</i>
	<b>PART</b>	Particle: a preposition-like form used together with a verb	<i>up, down, on, off, in, out, at, by</i>
	<b>PRON</b>	Pronoun: a shorthand for referring to an entity or event	<i>she, who, I, others</i>
	<b>SCONJ</b>	Subordinating Conjunction: joins a main clause with a subordinate clause such as a sentential complement	<i>that, which</i>
Other	<b>PUNCT</b>	Punctuation	<i>, , ()</i>
	<b>SYM</b>	Symbols like \$ or emoji	<i>\$, %</i>
	<b>X</b>	Other	<i>asdf, qwfg</i>



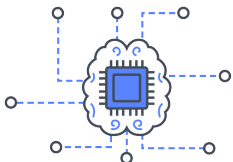


## Conceptos de NLP

### POS tagging en NLTK

```
>>> 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')]
```

```
>>> text = word_tokenize("They refuse to permit us to obtain the refuse permit")
>>> nltk.pos_tag(text)
[('They', 'PRP'), ('refuse', 'VBP'), ('to', 'TO'), ('permit', 'VB'), ('us', 'PRP'),
 ('to', 'TO'), ('obtain', 'VB'), ('the', 'DT'), ('refuse', 'NN'), ('permit', 'NN')]
```





# Conceptos de NLP

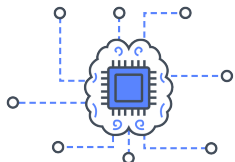
## POS tagging en Spacy (Español)

```
!python -m spacy download es_core_news_sm  
!python -m spacy download es_core_news_md
```

```
import spacy  
import es_core_news_sm  
import es_core_news_md  
  
nlp = es_core_news_sm.load()  
doc = nlp('Un desastroso espíritu posee tu tierra:  
donde la tribu unida blandió sus mazas,  
hoy se enciende entre hermanos perpetua guerra,  
se hieren y destrazan las mismas razas.')
```

```
for token in doc: print(token.text, "|", token.lemma_, '|', token.pos_)
```





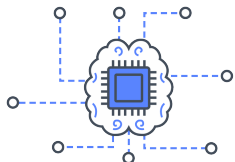
# Conceptos de NLP

## POS tagging en Spacy (Español)

```
!python -m spacy download es_core_news_sm  
!python -m spacy download es_core_news_md
```

```
import spacy  
import es_core_news_sm  
import es_core_news_md  
  
nlp = es_core_news_sm.load()  
doc = nlp('''Un desastroso espíritu posee tu tierra:  
donde la tribu unida blandió sus mazas,  
hoy se enciende entre hermanos perpetua guerra,  
se hieren y destrozan las mismas razas.''' )  
  
for token in doc: print(token.text, "|", token.lemma_, '|', token.pos_)
```

Un | Un | DET  
desastroso | desastroso | NOUN  
espíritu | espíritu | PROPN  
posee | poseer | VERB  
tu | tu | DET  
tierra | tierra | NOUN  
: | : | PUNCT

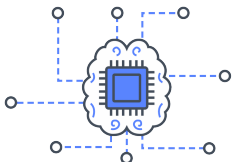




## Conceptos de NLP

### POS tagging ¿Cómo funciona?

- ▶ Se dispone de un corpus etiquetado.
- ▶ La secuencia de tags es interpretada como una cadena de Markov:  
 $P(x_{t+1} \mid x_t, \dots, x_1) = P(x_{t+1} \mid x_t)$ ,  $x_1, \dots, x_{t+1}$  representan tags
- ▶ Usamos un modelo generativo para términos, con tags como estados ocultos:  $P(t \mid x_1, \dots, x_{t+1}) = P(t \mid x_{t+1})$

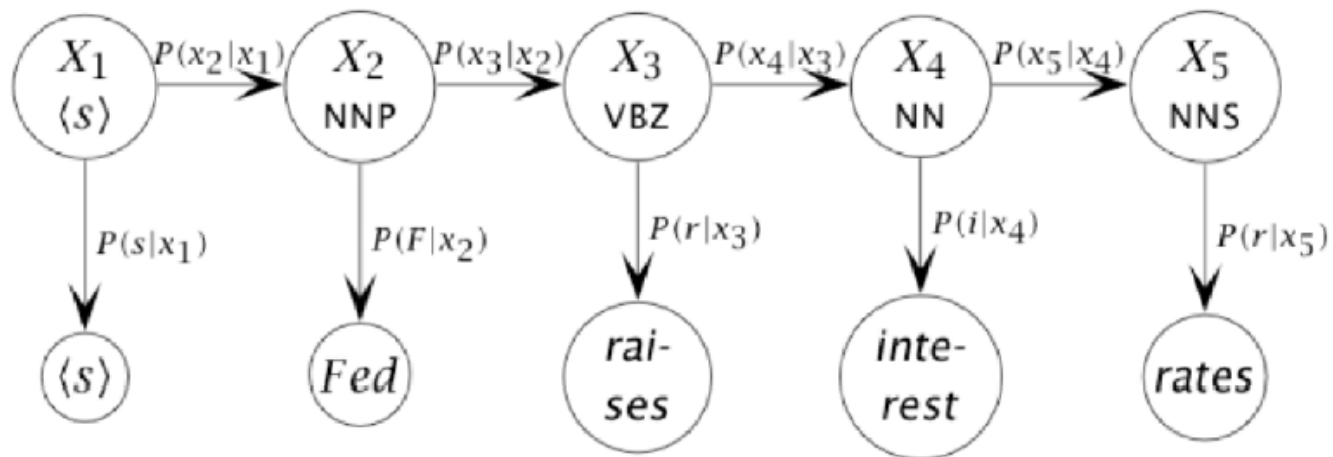




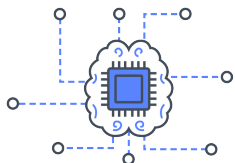


## Conceptos de NLP

### POS tagging ¿Cómo funciona?



- En general muestran buena precisión (sobre 90 %).



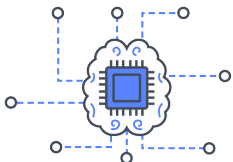


## Conceptos de NLP

HMM POS tagging:

¿Cómo defino la cadena de transiciones?

$Q = q_1 q_2 \dots q_N$	a set of $N$ <b>states</b>
$A = a_{11} \dots a_{ij} \dots a_{NN}$	a <b>transition probability matrix</b> $A$ , each $a_{ij}$ representing the probability of moving from state $i$ to state $j$ , s.t. $\sum_{j=1}^N a_{ij} = 1 \quad \forall i$
$O = o_1 o_2 \dots o_T$	a sequence of $T$ <b>observations</b> , each one drawn from a vocabulary $V = v_1, v_2, \dots, v_V$
$B = b_i(o_t)$	a sequence of <b>observation likelihoods</b> , also called <b>emission probabilities</b> , each expressing the probability of an observation $o_t$ being generated from a state $q_i$
$\pi = \pi_1, \pi_2, \dots, \pi_N$	an <b>initial probability distribution</b> over states. $\pi_i$ is the probability that the Markov chain will start in state $i$ . Some states $j$ may have $\pi_j = 0$ , meaning that they cannot be initial states. Also, $\sum_{i=1}^n \pi_i = 1$





## Conceptos de NLP

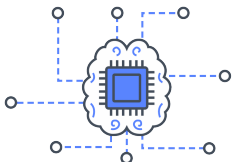
HMM POS tagging:

¿Cómo defino la cadena de transiciones?

$Q = q_1 q_2 \dots q_N$	a set of $N$ <b>states</b>
$A = a_{11} \dots a_{ij} \dots a_{NN}$	a <b>transition probability matrix</b> $A$ , each $a_{ij}$ representing the probability of moving from state $i$ to state $j$ , s.t. $\sum_{j=1}^N a_{ij} = 1 \quad \forall i$
$O = o_1 o_2 \dots o_T$	a sequence of $T$ <b>observations</b> , each one drawn from a vocabulary $V = v_1, v_2, \dots, v_V$
$B = b_i(o_t)$	a sequence of <b>observation likelihoods</b> , also called <b>emission probabilities</b> , each expressing the probability of an observation $o_t$ being generated from a state $q_i$
$\pi = \pi_1, \pi_2, \dots, \pi_N$	an <b>initial probability distribution</b> over states. $\pi_i$ is the probability that the Markov chain will start in state $i$ . Some states $j$ may have $\pi_j = 0$ , meaning that they cannot be initial states. Also, $\sum_{i=1}^n \pi_i = 1$

Objetivo del modelo:

$$\hat{t}_{1:n} = \underset{t_1 \dots t_n}{\operatorname{argmax}} P(t_1 \dots t_n | w_1 \dots w_n)$$





## Conceptos de NLP

HMM POS tagging:  $\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} P(t_1 \dots t_n | w_1 \dots w_n)$

Aplicamos la regla de Bayes:  $\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} \frac{P(w_1 \dots w_n | t_1 \dots t_n) P(t_1 \dots t_n)}{P(w_1 \dots w_n)}$

y simplificamos el denominador:  $\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} P(w_1 \dots w_n | t_1 \dots t_n) P(t_1 \dots t_n)$

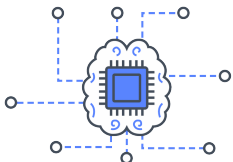
Se asume que la probabilidad de la palabra depende sólo de su tag:

$$P(w_1 \dots w_n | t_1 \dots t_n) \approx \prod_{i=1}^n P(w_i | t_i)$$

y se asume que el tag sólo depende del tag previo:

$$P(t_1 \dots t_n) \approx \prod_{i=1}^n P(t_i | t_{i-1})$$

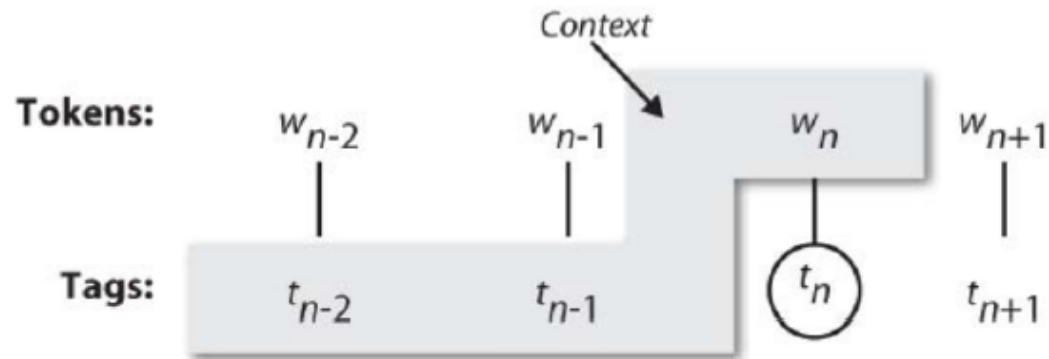
Finalmente:  $\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} P(t_1 \dots t_n | w_1 \dots w_n) \approx \operatorname{argmax}_{t_1 \dots t_n} \prod_{i=1}^n \overbrace{P(w_i | t_i)}^{\text{emission}} \overbrace{P(t_i | t_{i-1})}^{\text{transition}}$



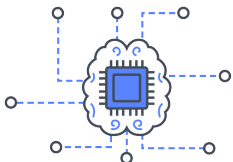


## Conceptos de NLP

### POS tagging ¿Cómo funciona?



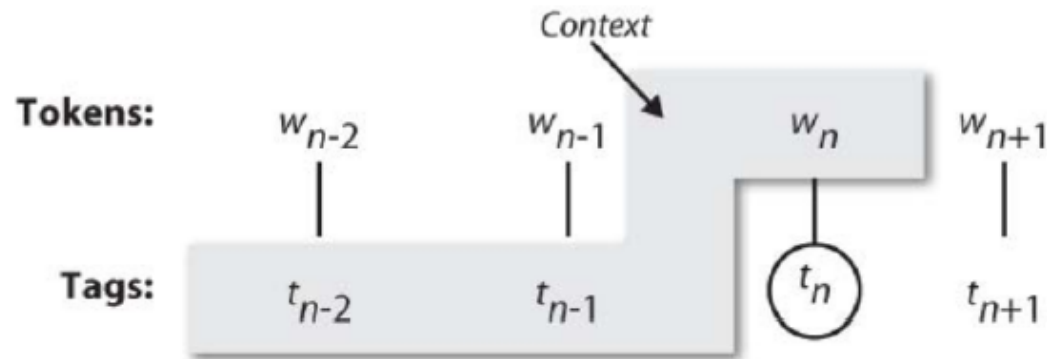
- ▶ Considera los tags de las dos palabras precedentes.
- ▶ En general muestra mejor precisión que HMM.



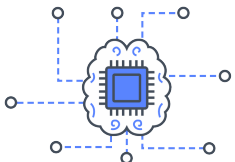


## Conceptos de NLP

### POS tagging ¿Cómo funciona?



- ▶ Considera los tags de las dos palabras precedentes.
- ▶ En general muestra mejor precisión que HMM.



Hinrich Schütze, [Yoram Singer](#): Part-of-Speech Tagging using a Variable Memory Markov Model. [ACL1994](#): 181-187



## Conceptos de NLP

### POS tagging ¿Cuáles datos usan?

Treebanks: [Penn treebank](#) (más famoso), [UAM Spanish Treebank](#), ...

[Treebank viewer](#):

Sentence Count: 317 Displayed Tree (Sentence): 287

((I (PRP)) (am (VBP)) (proud (JJ)) (that (IN)) (John (NNP))  
((John (NNP)) (is (VBZ)) (tough (JJ)) (to (TO)) (please (VI)  
((It (PRP)) (is (VBZ)) (tough (JJ)) (to (TO)) (please (VB)) (I  
((John (NNP)) (is (VBZ)) (likely (JJ)) (to (TO)) (win (VB)) (I  
((Who (PRP)) (does (VBZ)) (his (PRP\$)) (mother (NN)) (I  
((His (PRP\$)) (mother (NN)) (loves (VBZ)) (everyone (NI  
((Every (RB)) (man (NN)) (likes (NNS)) (some (DT)) (syn  
((Every (RB)) (man (NN)) (asked (VBD)) (some (DT)) (a  
((Someone (NN)) (gave (VBD)) (every (JJ)) (actress (NN)  
((Which (NNP)) (man (NN)) (liked (VBD)) (which (WDT))  
((John (NNP)) (seems (VBZ)) (that (IN)) (he (PRP)) (likes  
((John (NNP)) (asked (VBD)) (Mary (NNP)) (about (IN)) (I  
((We (PRP)) (like (IN)) (myself (PRP)))  
((John (NNP)) (seems (VBZ)) (is (VBZ)) (crazy (NN)))  
((John (NNP)) (tried (VBD)) (Bill (NNP)) (to (TO)) (seem  
((John (NNP)) (tried (VBD)) (to (TO)) (be (VB)) (arrested  
((Mary (NNP)) (is (VBZ)) (proud (VBN)) (of (IN)) (Bill (NN  
((John (NNP)) (seems (VBZ)) (that (IN)) (he (PRP)) (is (V  
((The (DT)) (destruction (NN)) (Rome (NNP)) (worked (V  
((The (DT)) (destruction (NN)) (of (IN)) (Rome (NNP)) (I  
((The (DT)) (belief (NN)) (John (NNP)) (to (TO)) (be (VB  
((The (DT)) (belief (NN)) (of (IN)) (John (NNP)) (to (TO))  
((The (DT)) (belief (NN)) (that (IN)) (John (NNP)) (is (VB  
((John (NNP)) ('s (POS)) (belief (NN)) (to (TO)) (be (VB))  
((John (NNP)) (seems (VBZ)) (that (IN)) (his (PRP\$)) (be

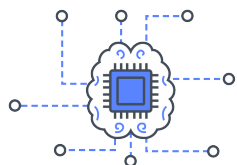
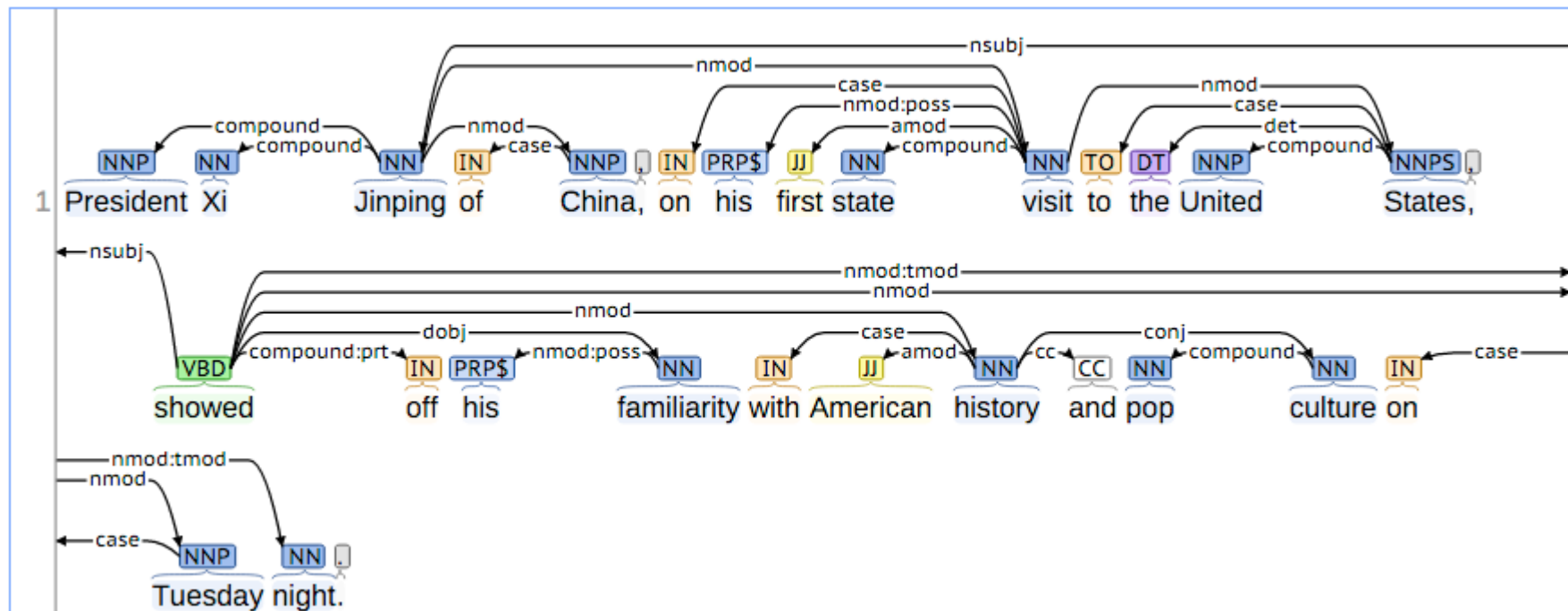
Every man asked some actress that he met about some play that she appeared in

Diagram illustrating a sentence structure tree for the sentence: "Every man asked some actress that he met about some play that she appeared in". The tree shows hierarchical groupings into NP (Noun Phrase) and VP (Verb Phrase) structures, with POS tags (DT, NN, VBD, IN, PRP, PP, NP, VP, PRT, RP) assigned to each word.



## Conceptos de NLP

### POS tagging y dependencias







## Conceptos de NLP

### POS tagging y dependencias en Spacy

Editable Code

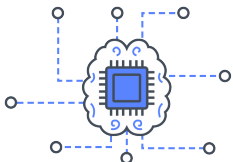
spaCy v3.0 - Python 3 - via Binder

```
import spacy

nlp = spacy.load("en_core_web_sm")
doc = nlp("Apple is looking at buying U.K. startup for $1 billion")

for token in doc:
    print(token.text, token.lemma_, token.pos_, token.tag_, token.dep_,
          token.shape_, token.is_alpha, token.is_stop)
```

RUN

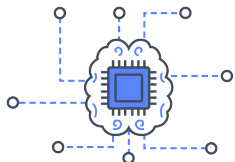




## Conceptos de NLP

### POS tagging y dependencias en Spacy

TEXT	LEMMA	POS	TAG	DEP	SHAPE	ALPHA	STOP
Apple	apple	PROPN	NNP	nsubj	Xxxxx	True	False
is	be	AUX	VBZ	aux	xx	True	True
looking	look	VERB	VBG	ROOT	xxxx	True	False
at	at	ADP	IN	prep	xx	True	True
buying	buy	VERB	VBG	pcomp	xxxx	True	False
U.K.	u.k.	PROPN	NNP	compound	X.X.	False	False
startup	startup	NOUN	NN	dobj	xxxx	True	False
for	for	ADP	IN	prep	xxx	True	True
\$	\$	SYM	\$	quantmod	\$	False	False
1	1	NUM	CD	compound	d	False	False
billion	billion	NUM	CD	pobj	xxxx	True	False



► Dependencia sintáctica



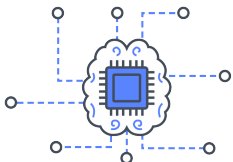
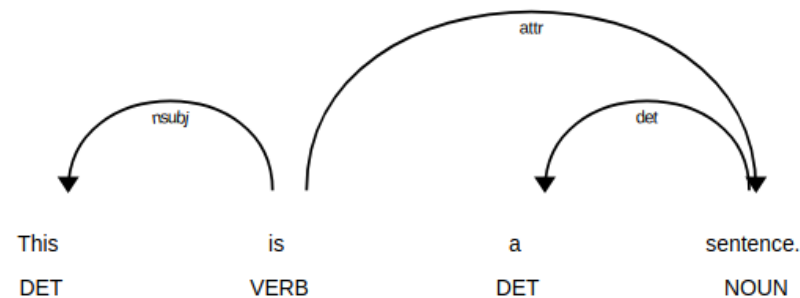
# Conceptos de NLP

## POS tagging y dependencias en Spacy

### DEPENDENCY EXAMPLE

```
import spacy
from spacy import displacy

nlp = spacy.load("en_core_web_sm")
doc = nlp("This is a sentence.")
displacy.serve(doc, style="dep")
```

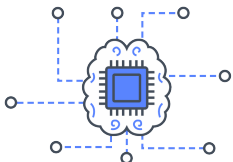




## Conceptos de NLP

### Dependencias ¿Cómo funciona?

<b>Clausal Argument Relations</b>	<b>Description</b>
NSUBJ	Nominal subject
DOBJ	Direct object
IOBJ	Indirect object
CCOMP	Clausal complement
XCOMP	Open clausal complement
<b>Nominal Modifier Relations</b>	<b>Description</b>
NMOD	Nominal modifier
AMOD	Adjectival modifier
NUMMOD	Numeric modifier
APPOS	Appositional modifier
DET	Determiner
CASE	Prepositions, postpositions and other case markers
<b>Other Notable Relations</b>	<b>Description</b>
CONJ	Conjunct
CC	Coordinating conjunction

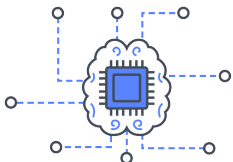
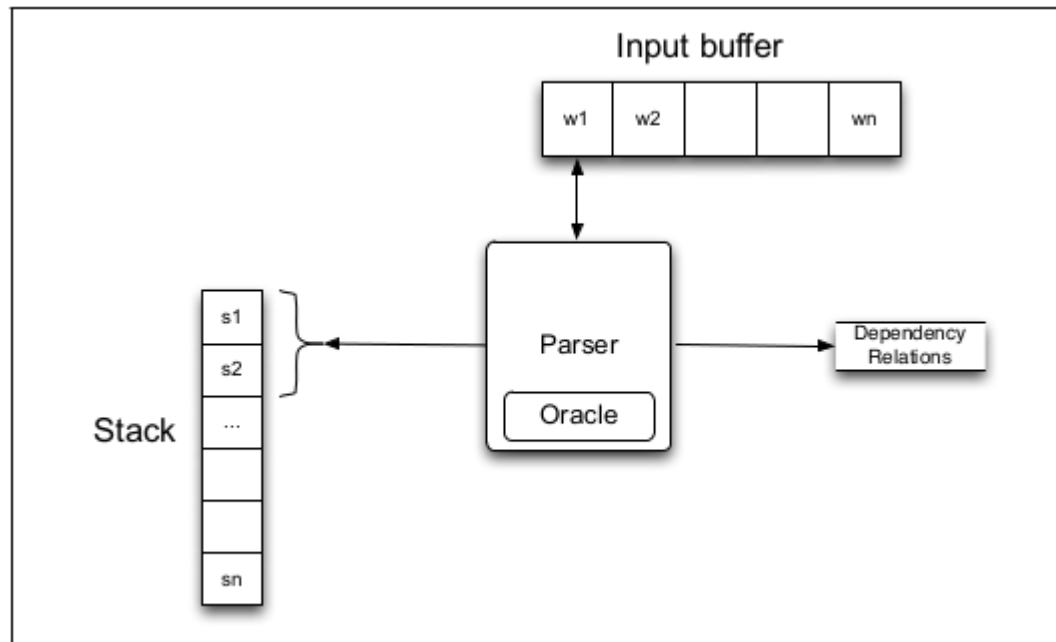




## Conceptos de NLP

### Dependencias ¿Cómo funciona?

Transition-based parser:



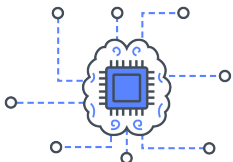


## Conceptos de NLP

### Dependencias ¿Cómo funciona?

Transition-based parser: Deducimos training instances desde un treebank. Podemos determinísticamente registrar las operaciones correctas del parser, creando un training dataset.

Step	Stack	Word List	Action	Relation Added
0	[root]	[book, me, the, morning, flight]	SHIFT	
1	[root, book]	[me, the, moming, flight]	SHIFT	
2	[root, book, me]	[the, morning, flight]	RIGHTARC	(book → me)
3	[root, book]	[the, morning, flight]	SHIFT	
4	[root, book, the]	[morning, flight]	SHIFT	
5	[root, book, the, morning]	[flight]	SHIFT	
6	[root, book, the, morning, flight]	[]	LEFTARC	(moming ← flight)
7	[root, book, the, flight]	[]	LEFTARC	(the ← flight)
8	[root, book, flight]	[]	RIGHTARC	(book → flight)
9	[root, book]	[]	RIGHTARC	(root → book)
10	[root]	[]	Done	





## Conceptos de NLP

### Dependencias ¿Cómo funciona?

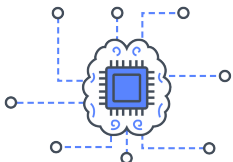
Training instances (intermedia):

Stack	Word buffer	Relations
[root, canceled, flights]	[to Houston]	(canceled → United) (flights → morning) (flights → the)

Actúa sobre el primer token del buffer

La transición correcta del parser es *shift*. Luego, creamos training instances de este tipo en el dataset:

$\langle s_1.w = flights, op = shift \rangle$	← Tokens del stack
$\langle s_2.w = canceled, op = shift \rangle$	
$\langle s_1.t = NNS, op = shift \rangle$	← POS-tags del stack
$\langle s_2.t = VBD, op = shift \rangle$	
$\langle b_1.w = to, op = shift \rangle$	← Token y tag del buffer
$\langle b_1.t = TO, op = shift \rangle$	
$\langle s_1.wt = flightsNNS, op = shift \rangle$	← Template compuesto





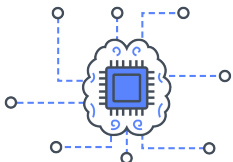
## Conceptos de NLP

### Dependencias ¿Cómo funciona?

Source	Feature templates		
One word	$s_1.w$	$s_1.t$	$s_1.wt$
	$s_2.w$	$s_2.t$	$s_2.wt$
	$b_1.w$	$b_1.w$	$b_0.wt$
Two word	$s_1.w \circ s_2.w$	$s_1.t \circ s_2.t$	$s_1.t \circ b_1.w$
	$s_1.t \circ s_2.wt$	$s_1.w \circ s_2.w \circ s_2.t$	$s_1.w \circ s_1.t \circ s_2.t$
	$s_1.w \circ s_1.t \circ s_2.t$	$s_1.w \circ s_1.t$	

El clasificador predice dep tag a nivel de palabra usando las características indicadas.

Típicamente se usan SVM, LR multinomial o ANN con softmax para esta tarea.







## Conceptos de NLP

### Dependencias ¿Cómo funciona?

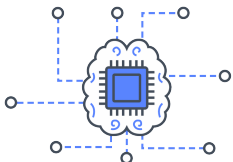
Source	Feature templates		
One word	$s_1.w$	$s_1.t$	$s_1.wt$
	$s_2.w$	$s_2.t$	$s_2.wt$
	$b_1.w$	$b_1.w$	$b_0.wt$
Two word	$s_1.w \circ s_2.w$	$s_1.t \circ s_2.t$	$s_1.t \circ b_1.w$
	$s_1.t \circ s_2.wt$	$s_1.w \circ s_2.w \circ s_2.t$	$s_1.w \circ s_1.t \circ s_2.t$
	$s_1.w \circ s_1.t \circ s_2.t$	$s_1.w \circ s_1.t$	

El clasificador predice dep tag a nivel de palabra usando las características indicadas.

Típicamente se usan SVM, LR multinomial o ANN con softmax para esta tarea.



Danqi Chen, [Christopher D. Manning](#): A Fast and Accurate Dependency Parser using Neural Networks. [EMNLP2014](#): 740-750

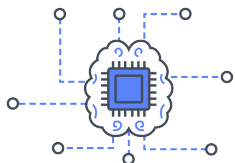
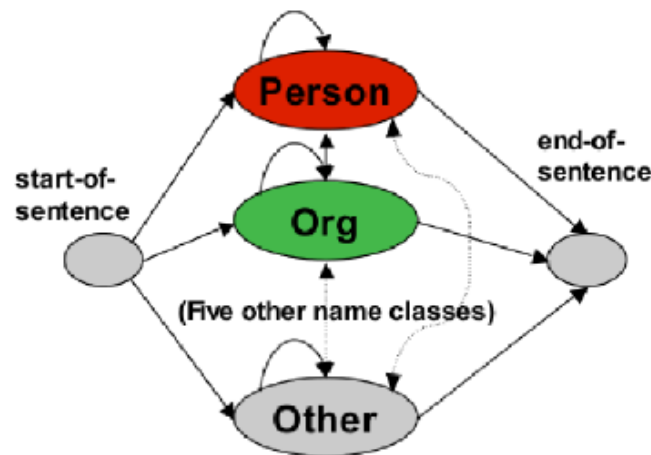




## Conceptos de NLP

### Named Entity Recognition

- ▶ Tarea: Identificar entidades en texto (personas, organizaciones, etc.)
- ▶ Separa el text en chunks, y para cada cual asocia una NE. Opera sobre texto tagged.
- ▶ NER types: organization, person, location, date, time, money, percent, facility (human made artifacts), gpe (geo-political ents).
- ▶ POS tagging puede ayudar, agregando *entity* como un estado mas.





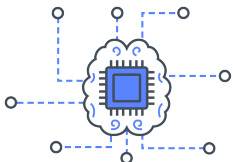
## Conceptos de NLP

### Named Entity Recognition

1 President Xi Jinping of China, on his first state visit to the United States, showed off his familiarity with American history and pop culture on Tuesday night.

Named Entity Recognition (NER) labels for the sentence above:

- Person: President Xi Jinping
- Loc: China
- ORDINAL: first
- Location: United States
- Misc: American
- Date: Tuesday
- Time: night

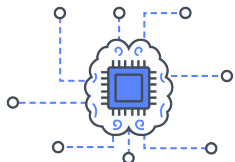
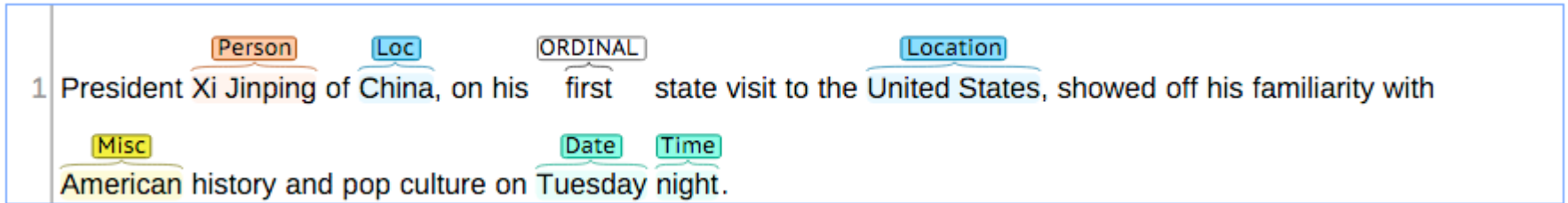




# Conceptos de NLP

## Named Entity Recognition

bi-grama





## Conceptos de NLP

### Named Entity Recognition

bi-grama



1 President Xi Jinping of China, on his first state visit to the United States, showed off his familiarity with American history and pop culture on Tuesday night.

Person Loc ORDINAL Location Misc Date Time

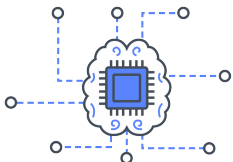
Named Entity Recognition puede ser muy desafiante:

Back in 2000 , People Magazine PUBLISHER highlighted Prince Williams' PERSON style who at the time was a little more fashion-conscious , even making fashion statements at times .

Now-a-days the prince mainly wears navy COLOR suits ITEM ( sometimes double-breasted DESIGN ) , light blue COLOR button-ups ITEM with classic LOOK pointed DESIGN collars PART , and burgundy COLOR ties ITEM .

But who knows what the future holds ...

Duchess Kate PERSON did wear an Alexander McQueen BRAND dress ITEM to the wedding OCCASION in the fall of 2017 SEASON .

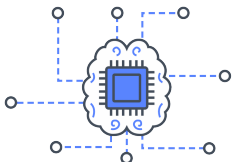




## Conceptos de NLP

### Named Entity Recognition en Spacy

```
nlp_md = es_core_news_md.load()  
article_text = '''La ONG Fundación del Río explicó este viernes ( ...  
doc = nlp_md(article_text)  
SVG(data = displacy.render(doc, style="ent"))
```



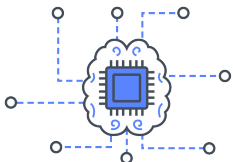


## Conceptos de NLP

### Named Entity Recognition en Spacy

```
nlp_md = es_core_news_md.load()
article_text = '''La ONG Fundación del Río explicó este viernes ( ...
doc = nlp_md(article_text)
SVG(data = displacy.render(doc, style="ent"))
```

La **ONG Fundación del Río** **ORG** explicó este viernes que la decisión de la **Organización de la ONU** **ORG** para la **Educación** **ORG**, la **Ciencia** **LOC** y la **Cultura** **LOC** ( **Unesco** **ORG** ) de declarar como geoparque el río **Coco** **LOC**, ubicado en el norte de **Nicaragua** **LOC**, obliga a las autoridades nicaragüenses a proteger su ecosistema, ya que se encuentra en el área más deforestada de la cuenca. **La Unesco** **LOC** está reconociendo la importancia del **río Coco** **LOC**, pero también está haciendo un llamado al **Gobierno** **LOC** a que actúe en la protección y la conservación de esos ecosistemas, dijo a **Efe** **PER** el presidente de la **Fundación del Río** **ORG**, **Amaru Ruiz** **PER**.

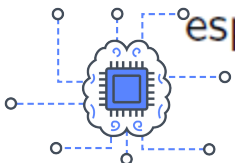




## Conceptos de NLP

### Collocations: n-gramas frecuentes

- ▶ El significado conjunto es más que la suma de las partes (**compositionality**)
  1. Armas de destrucción masiva
  2. Strong tea
  3. Libre de sodio
  4. Intel inside
  5. Fast food
  6. Nuclear war
- ▶ Detectar colocaciones mejora la representación del contenido.
- ▶ Cada colocación puede ser procesada como un término.
- ▶ Se pueden detectar analizando co ocurrencias, etiquetando el par como una colocación si su co ocurrencia es mucho mayor que la esperada (azar, equiprobable).



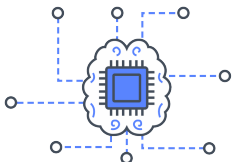




## Conceptos de NLP

*Collocations:* ¿Cómo se detectan en NLTK?

```
>>> import nltk
>>> from nltk.collocations import *
>>> bigram_measures = nltk.collocations.BigramAssocMeasures()
>>> trigram_measures = nltk.collocations.TrigramAssocMeasures()
>>> finder = BigramCollocationFinder.from_words(
...     nltk.corpus.genesis.words('english-web.txt'))
>>> finder.nbest(bigram_measures.pmi, 10) # doctest: +NORMALIZE_WHITESPACE
[(u'Allon', u'Bacuth'), (u'Ashteroth', u'Karnaim'), (u'Ben', u'Ammi'),
 (u'En', u'Mishpat'), (u'Jegar', u'Sahadutha'), (u'Salt', u'Sea'),
 (u'Whoever', u'sheds'), (u'appoint', u'overseers'), (u'aromatic', u'resin'),
 (u'cutting', u'instrument')]
```





## Conceptos de NLP

*Collocations*: ¿Cómo se detectan en NLTK?

```
>>> import nltk
>>> from nltk.collocations import *
>>> bigram_measures = nltk.collocations.BigramAssocMeasures()
>>> trigram_measures = nltk.collocations.TrigramAssocMeasures()
>>> finder = BigramCollocationFinder.from_words(
...     nltk.corpus.genesis.words('english-web.txt'))
>>> finder.nbest(bigram_measures.pmi, 10) # doctest: +NORMALIZE_WHITESPACE
[(u'Allon', u'Bacuth'), (u'Ashteroth', u'Karnaim'), (u'Ben', u'Ammi'),
 (u'En', u'Mishpat'), (u'Jegar', u'Sahadutha'), (u'Salt', u'Sea'),
 (u'Whoever', u'sheds'), (u'appoint', u'overseers'), (u'aromatic', u'resin'),
 (u'cutting', u'instrument')]
```

Podemos usar un umbral de frecuencia absoluta para encontrar *collocations* frecuentes y luego rankear usando PMI:

$$\text{pmi}(x; y) \equiv \log \frac{p(x, y)}{p(x)p(y)}$$

```
>>> finder.apply_freq_filter(3)
>>> finder.nbest(bigram_measures.pmi, 10) # doctest: +NORMALIZE_WHITESPACE
[(u'Beer', u'Lahai'), (u'Lahai', u'Roi'), (u'gray', u'hairs'),
 (u'Most', u'High'), (u'ewe', u'lambs'), (u'many', u'colors'),
 (u'burnt', u'offering'), (u'Paddan', u'Aram'), (u'east', u'wind'),
 (u'living', u'creature')]
```

