

# GENERAL OVERVIEW AND REGULAR EXPRESSIONS

---

*Introduction to NLP  
Session 1*

2023

*David Buchaca Prats*

# NLP APPLICATIONS

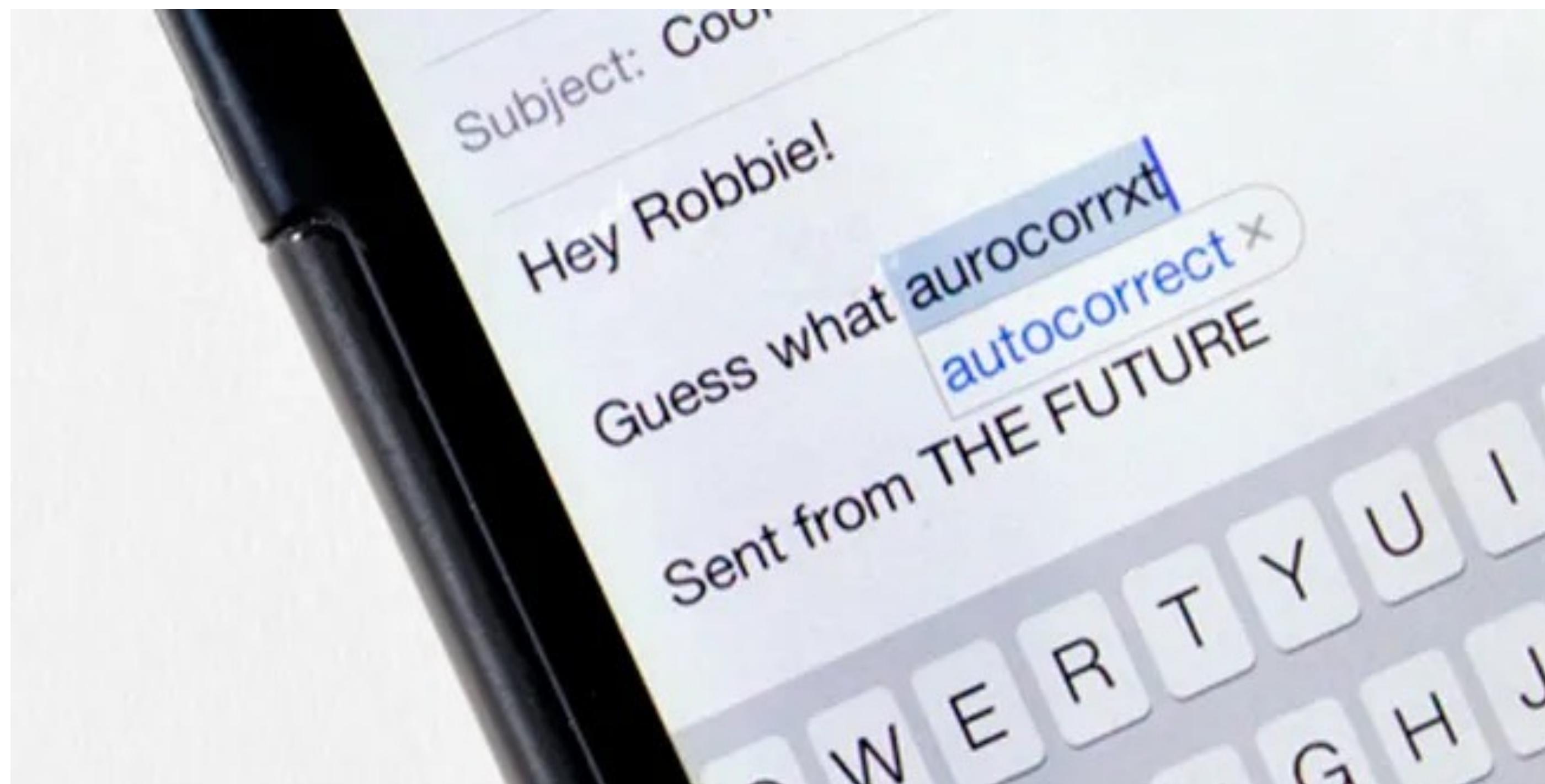
---

- Spellchecking and Autocorrection engines
- Sentiment analysis
- Search engines
- Question answering online services (Quora, Reddit, StackOverflow)
- Question answering systems
- Dialogue systems or "chatbots" (Microsoft just integrated ChatGPT in BING!)
- Machine translation
- Text Summarization
- Speech Recognition

# APPLICATIONS: SPELLCHECKING AND AUTOCORRECTION

---

- Spellchecking and autocorrection are key components of digital experiences
- Such models are mostly based on detecting very low level probability words, finding possible candidates, scoring the candidates in the context, and choosing the best one.



# APPLICATIONS: SENTIMENT ANALYSIS

---

- Sentiment analysis is the problem of detecting the "sentiment" or "class label" in a piece of text. Usually it refers to mapping text to a set of possible sentiments. For example: Positive, Negative, Neutral.
- Many applications might focus on other sentiments or classes. A popular use case consist on finding potentially offensive, racist, or insulting material.
- This is heavily used in
  - Survey responses
  - Customer service
  - Social media

# APPLICATIONS: SEARCH ENGINES

---

- Search is one of the most important applications of the internet and most of it is based on using text as input.
- There are applications where search takes other modalities, such as image to image, voice to image, image to video search ...
- Many companies integrate search as a key experience or component of the company:
  - General search engines: Google, Bing, Yahoo, duckduckgo
  - General second marked platforms: Ebay, Wallapop, Amazon
  - Specialized niche search platforms: Vlex, Habitaclia, fotocasa

# APPLICATIONS: QUESTION ANSWERING SYSTEMS EMBEDDED IN SEARCH

The screenshot shows a Microsoft Bing search results page. The search query is "what is bigger an elephant or a dog". The results are categorized under "ALL" (which is selected), along with other options like "IMAGES", "VIDEOS", "MAPS", "NEWS", "CHAT", and "MORE". Below the search bar, it says "About 144.000.000 results" and there is a date filter dropdown.

A large, prominent result box is titled "Elephant". It contains the following text:

A. An **elephant** is the most heavier than a dog.  
An elephant is ----- [heavy] than a dog. A. An ...  
✓ grammarquiz.net/24148/

Below this box, there is a "Was this helpful?" section with a thumbs up and thumbs down icon. To the right of the main search area, there is a vertical sidebar with letters V, E, T, I, A, W, O, A, W.

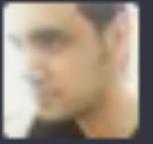
**People also ask**

**Is an elephant heavier than a dog?**  
An elephant is ----- [heavy] than a dog. A. An elephant is the most heavier than a dog. B. An elephant is more heavier than a dog. C. An elephant is the heaviest a dog.  
An elephant is ----- [hea...  
grammarquiz.net/24148/

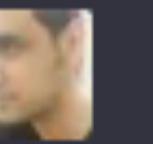
**How big are elephants compared to giraffes?**  
Answer Wiki. In terms of weight, elephants are by far much bigger than giraffes, weighing around 5–7 tons (for abnormally large males, 12 tons).  
Are giraffes or elephants bigger? - ...  
www.quora.com/Are-giraffes-or-eleph...

# APPLICATIONS: GENERAL QUESTION ANSWERING SYSTEMS FROM LANGUAGE MODELS

---

 are Catalans frugal ? 

 It is not accurate to make a blanket statement about the frugality of all people from Catalonia, as attitudes towards money and spending can vary greatly among individuals. However, it is often said that people from Catalonia have a reputation for being thrifty and careful with their money.  

 are Spaniards frugal ?

 It is not accurate to make a blanket statement about the frugality of all Spaniards, as attitudes towards money and spending can vary greatly among individuals and across regions. In general, the Spanish culture values family, friends, and socializing and it is common to spend money in these areas. However, some Spaniards may be more frugal and careful with their money.  

# APPLICATIONS: QUESTION ANSWERING ONLINE SERVICES (DUPLICATED QUESTIONS)

stack overflow Products How do I concatenate two lists in Python?

---

3244 votes **Q** How do I concatenate two lists in Python?  
✓ 31 answers  
4.0m views

How do I concatenate two lists in Python? Example: listone = [1, 2, 3] listtwo = [4, 5, 6] Expected outcome: >>> joinedlist [1, 2, 3, 4, 5, 6] ...  
python list concatenation

y2k 65k asked Nov 12, 2009 at 7:04

---

-3 votes **Q** how to concatenate two lists objects and join string objects in list [duplicate]  
3 answers  
36 views

How do I concatenate two lists in Python? example: input: x=['saeed'] y=['gharif'] output expected: z=['saeedgharif'] ...  
python

roini das 1 asked Nov 1, 2018 at 12:24

---

0 votes **Q** How to concatenate two list? [duplicate]  
✓ 1 answer  
47 views

How do I concatenate two lists in Python with the same size? ...  
python python-3.x

Danilo Dutra 25 asked Nov 4, 2019 at 10:05

---

300 votes **A** How do I concatenate two lists in Python?  
How do I concatenate two lists in Python? As of 3.9, these are the most popular stdlib methods for concatenating two (or more) lists in Python. Version Restrictions In-Place? ... But sum per-...  
python list concatenation

cs95 362k answered Jun 1, 2019 at 15:20

---

3 votes **Q** How do I concatenate two lists while taking elements 1 by 1 [duplicate]  
✓ 5 answers  
108 views

How do I concatenate two lists while taking elements 1 by 1 in Python? Example: listone = [1, 2, 3] listtwo = [4, 5, 6] Expected outcome: >>> joinedlist [1, 4, 2, 5, 3, 6] ...  
python list

user11597414 asked Apr 15, 2020 at 12:08

---

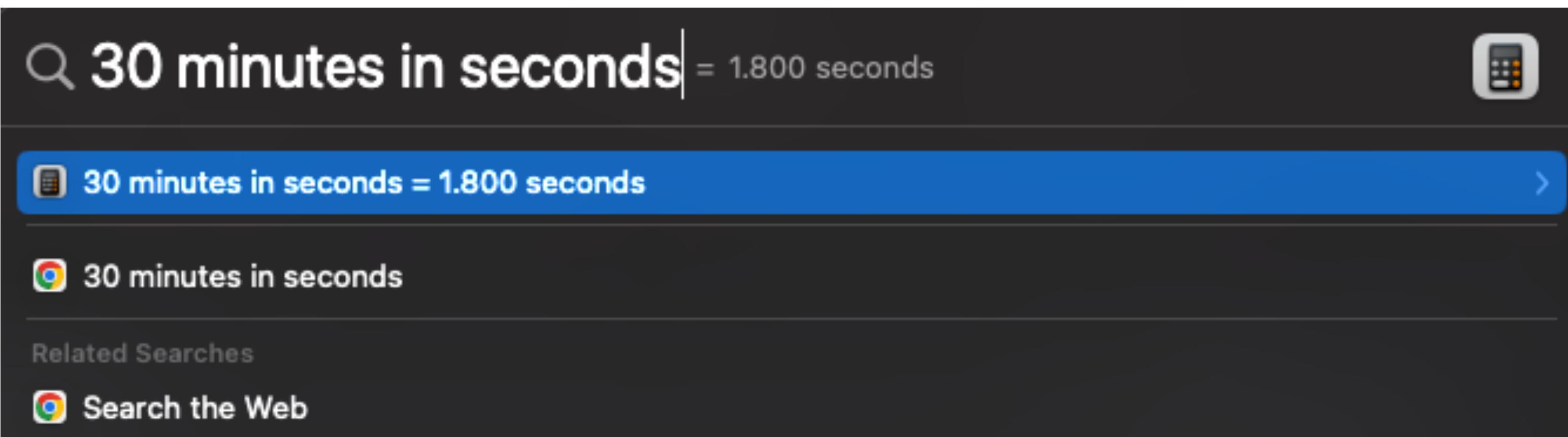
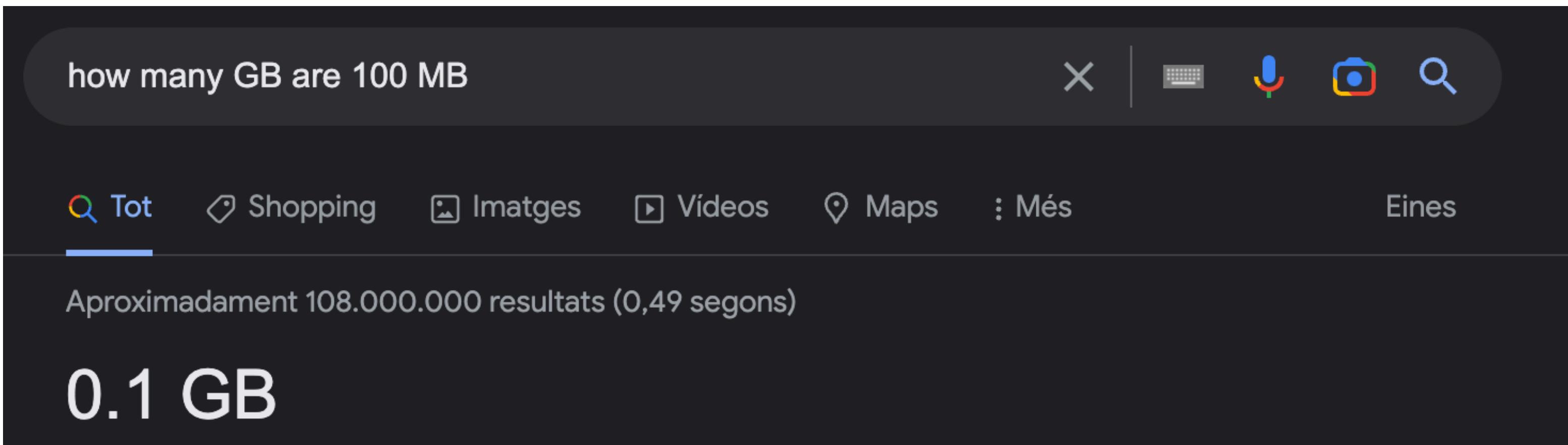
0 votes **Q** Merge two lists in reverse chronological order using regular expression python  
1 answer  
59 views

I am trying to merge two lists in Python in reverse chronological order using regular expression. ... I'm a little lost, the only thing I can do to merge them without errors so far is concatenate the...  
python regex python-re

tyrantGorilla 3 asked Feb 6 at 0:45

# APPLICATIONS: QUESTION ANSWERING SYSTEMS EMBEDDED IN SEARCH

- Some question answering systems have "specific purpose routines" build in.



# APPLICATIONS: MACHINE TRANSLATION

Català ↔ Anglès

Menjo truita de patates sense ceba ja que normalment no em senta bé després de menjar.

I eat potato omelette without onion as I usually don't feel well after eating.

Català ↔ Anglès

En cap cap

In no head



Català ↔ Anglès

En cap cap cap

No no no no



Anglès ↔ Espanyol

time flies like an arrow

El tiempo vuela como una flecha

Català ↔ Anglès

En cap cap cap el que cap en aquest cap

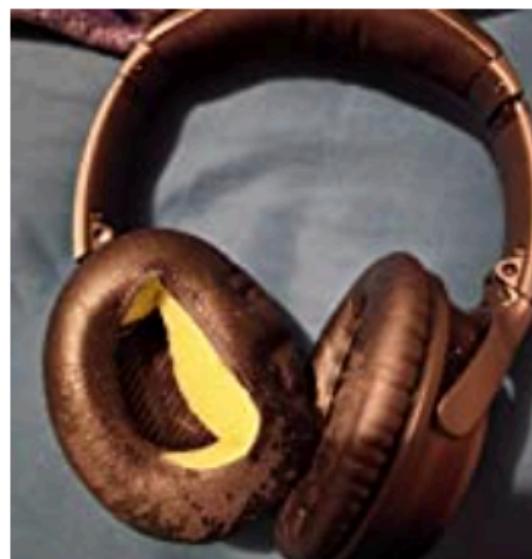
No no no no what no no no no



# APPLICATIONS: TEXT SUMMARYZATION / TOPIC RETRIEVAL

- Search based applications can retrieve relevant information or properties to help users navigate to interesting/relevant content.
- Information from reviews

## Reviews with images



[See all customer images](#)

## Read reviews that mention

sound quality   noise cancelling   battery life   google assistant  
highly recommend   every penny   worth the money   worth every  
headphones ever   much better   noise cancelation   bluetooth

## Read reviews that mention

sound quality   noise cancelling   battery life   google assistant  
highly recommend   every penny   worth the money   headphones ever  
worth every   much better   noise cancelation   ever owned

## Read reviews that mention

sound quality   noise cancelling   battery life   google assistant  
highly recommend   every penny   worth the money   headphones ever  
worth every   much better   noise cancelation   ever owned

# APPLICATIONS WE WILL SEE DURING THE COURSE

---

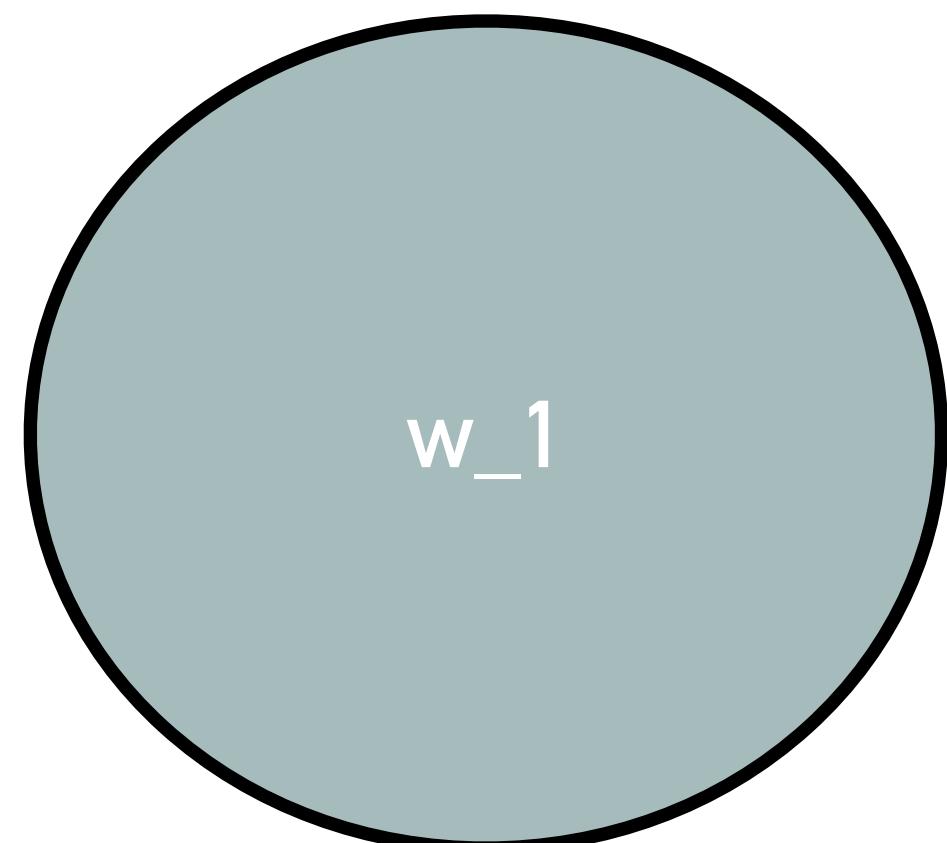
- Work efficiently with strings
- Perform preprocessing of text: Basic cleaning, tokenization and vocabulary learning
- Extract features from text (and for words)
- Perform Document classification
- Perform text cleaning/correction
- Search efficiently raw text data (Approximate nearest neighbour search)
- Word level tagging using structured prediction
- Language Models and applications of Language Models

# INTERESTING DETAILS FROM A PRACTICAL POINT OF VIEW

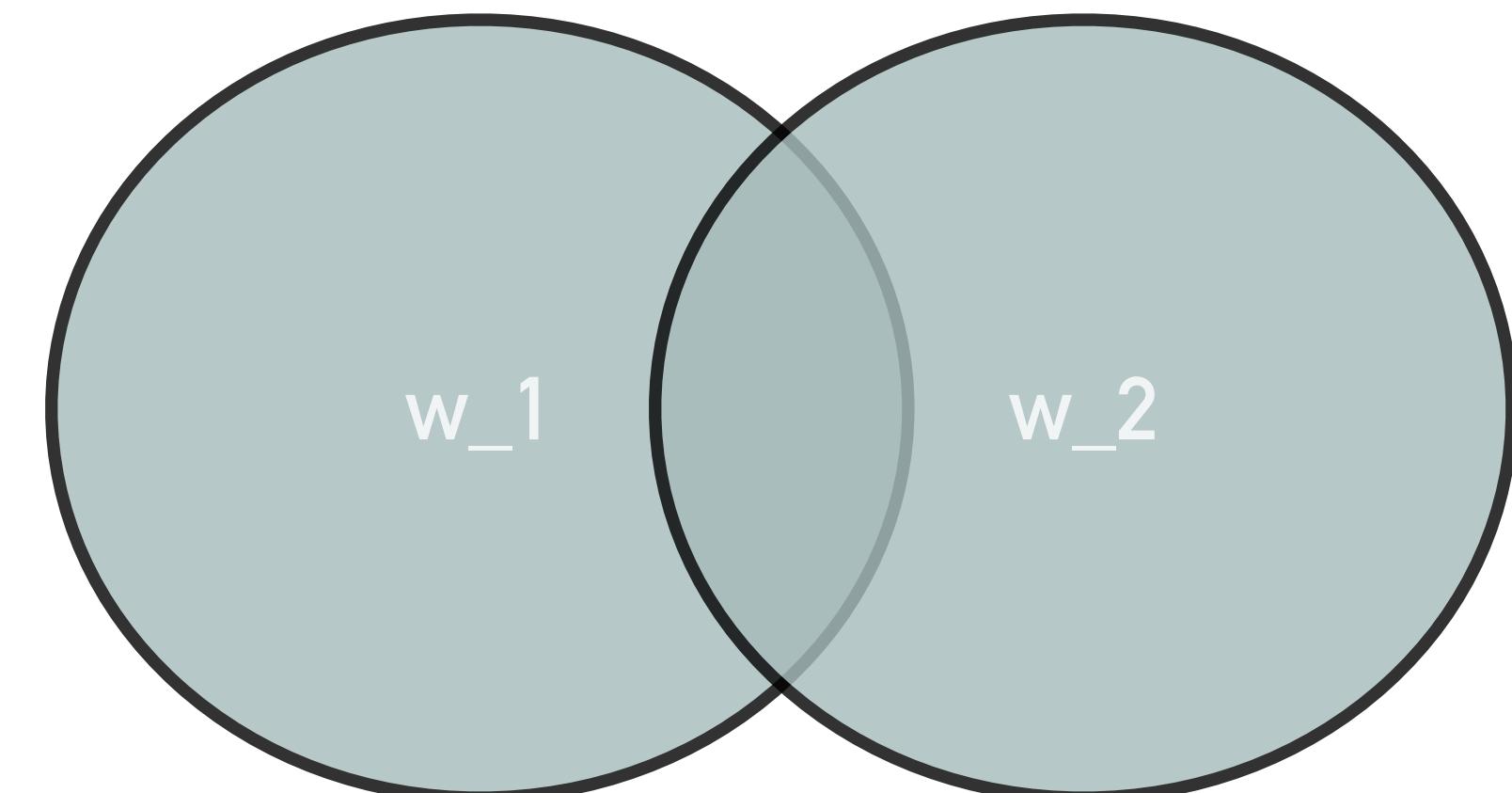
---

- The world is full of web services based on text data.
  - Google, Ebay, Amazon, Aliexpress, Wallapop
  - Most services already use NLP (or could be heavily improved using NLP)
- Let's look at a typical problem with key-word matching techniques for search

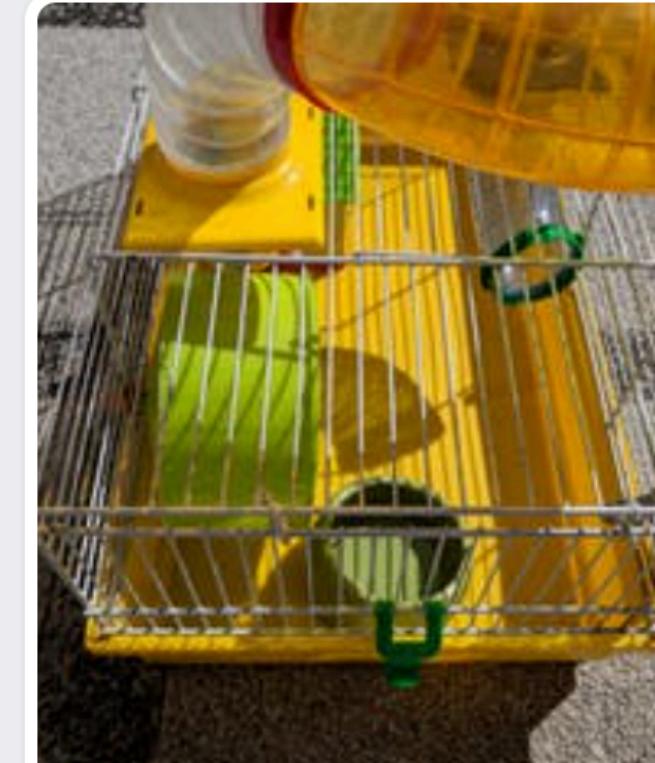
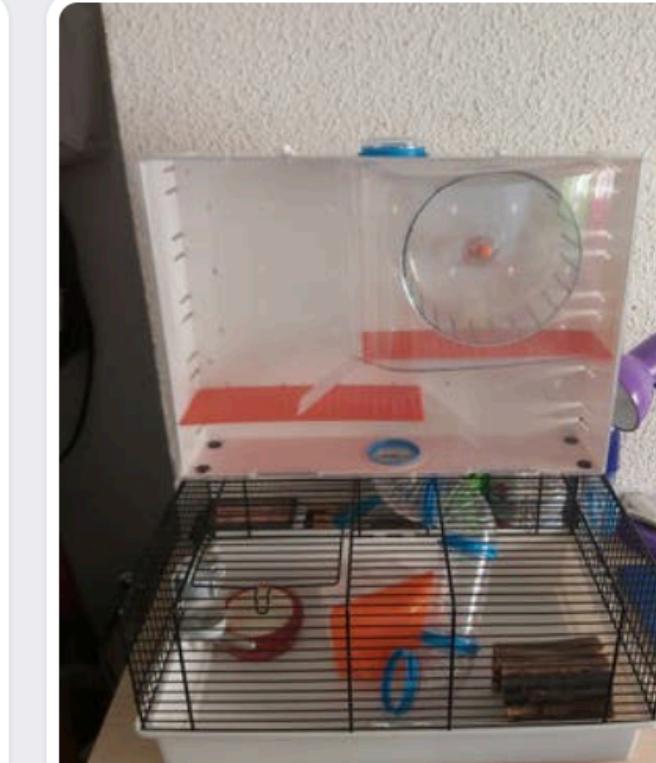
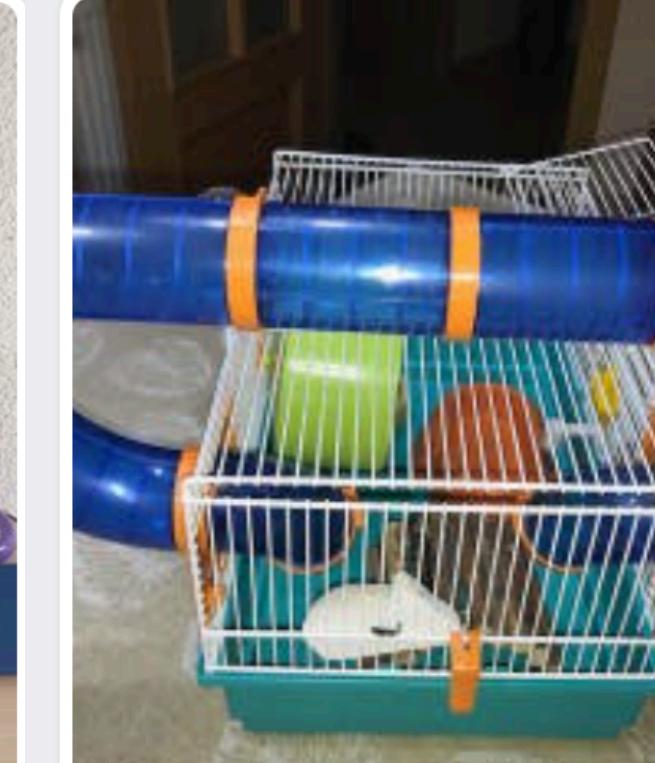
*Words have very distinct meaning*



*Words have similar meaning*

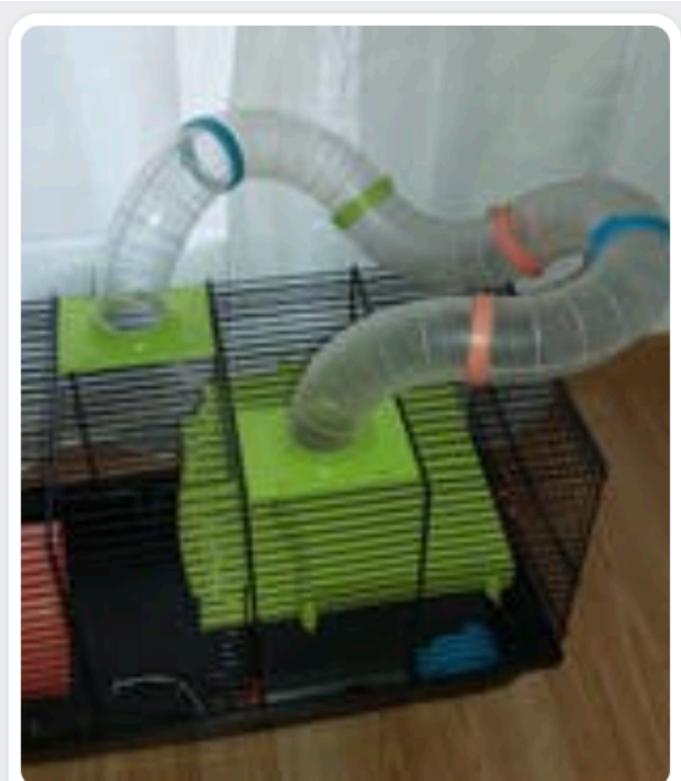


# LOOKING “HAMSTER” IN WALLAPOD

						
<b>10 €</b> La pandilla hamster  Juego de cooperación de 4 a 8 años.	<b>15 €</b> Jaula para hámster con...  ¡La nueva casita de tu amiguito peludo! El bebedero es ideal para que tu hámster apague su sed...	<b>20 €</b> Jaula grande de hámst...  Jaula grande con comedero, bebedero, casita, rueda y tubos para tu hámster. Si necesitas otro tamaño de...	<b>35 €</b> jaula  JAULA Hámster con accesorios	<b>30 €</b> Jaula hamster  Jaula de hamster con juego de tubos, rueda, comedero, bebedero y casita.	<b>25 €</b> Jaula para hámster  Jaula para hámster con poco uso, menos de 1 mes y está completamente limpia. Viene con los...	<b>11 €</b> Jaula Hámster  Jaula marca Hábitall hcon tobogán, comederos....y de regalo una bola de actividad para el Hámster,...

# LOOKING “JAMSTER” IN WALLAPOD

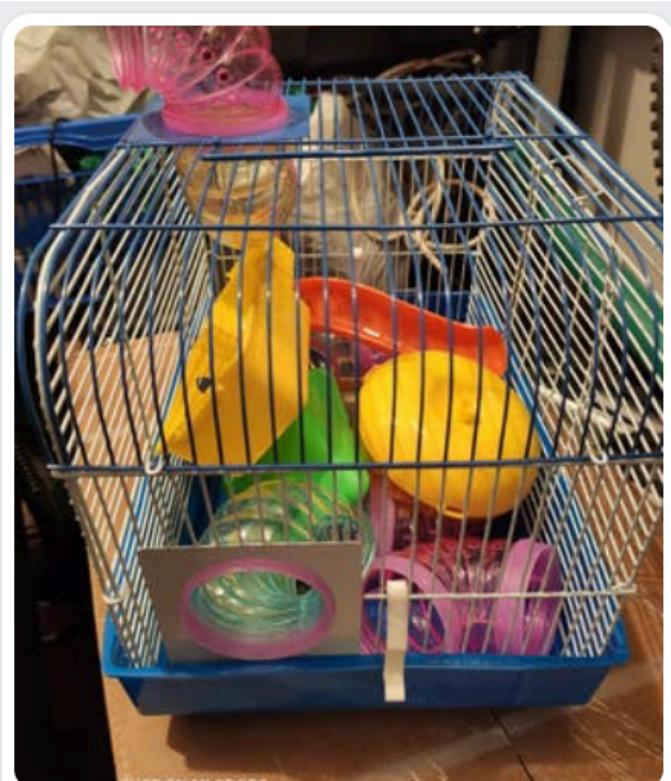
---



**25 €**

jaula jamster.

jaula jamster, contiene  
bebedero.



**20 €**

Jaula jamster

Vendo jaula de jamster por  
habernos abandonado. No  
tiene más de un mes de uso.  
Tiene tubos, rueda de...



**5 €**

Juego para Jamster

Juego para Jamster



**8 €**

jaula para jamster

esta perfecta, sin apenas  
uso. La usaba sólo cuando  
me iba de viaje, en casa  
usaba una más grande. La...



**3 €**

Jaula Jamster

Nueva, sin uso. Medidas  
21x29x20



**15 €**

Jaula de Jamster comp...

Jaula de Jamster completa y  
en perfecto estado.



**40 €**

jaula Hamster y mucho...

jaula de hámster Cunipic  
119C medidas 32.5X35.5,  
choza, comedero, noria  
para correr, altillo con...

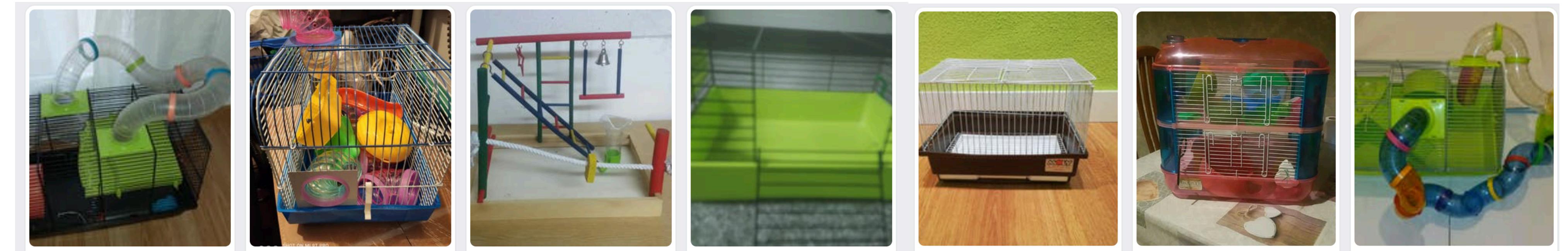
# HAMSTER VS JAMSTER

---

*Hamster*



*Jamster*



# INTERESTING PROBLEMS FROM A THEORETICAL POINT OF VIEW

---

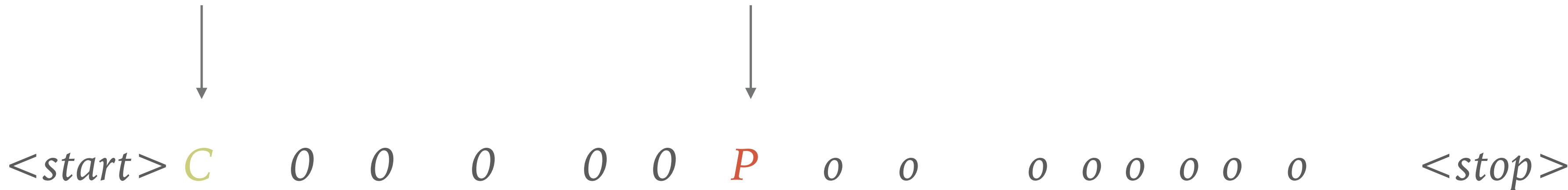
- We will see a different paradigm of Machine Learning with techniques that have to deal with sequential structured data
- In particular we will deal with independence assumptions in our models that can provide huge speedup benefits with good tradeoff in quality.
- Example: Let us consider  $P(X=x, Y=y)$  where  $x$  and  $y$  are sequences.
  - Assume  $x$  and  $y$  have length 10
  - Assume  $Y$  can have 12 possible values
  - How can you compute  $P(X=x)$  ?
  - $12^2 = 61.917.364.224$  sums

# PROBLEMS THAT WE WILL FACE

---

- Detecting similarity between documents/sentences
- Document Classification
- Tagging words as Named Entities (Person, Location, Company ...)

*<start> Ebay does not know what an ifone is (neither do I but I can guess) <stop>*



- Detecting a particular language from a given input string
- Predicting the next word/character within a stream of words/characters
- Predicting the probability of a given string

# EVALUATION (SCALE 0-10)

---

- 70 % Deliverables
  - Two projects have to be delivered. These projects consist on working on a “challenge”.
  - You will be given a task (dataset and metric) and you will have some freedom to explore different strategies to solve the task.
- 30 % Final Exam
  - A final exam will evaluate all the material in the course.
  - Even if you fail the exam you can pass the course with your deliverables.
- In the event that Exam + Deliverables < 5 there is a re-evaluation exam that you can apply if and only if you fail the course.

# REGULAR EXPRESSIONS

---

# REGULAR EXPRESSIONS

---

- Generate a regular expression using `p = re.compile(r'our_regular_expression')`
- `p` will be a `re.Pattern` object that we can apply to any string.
- `p.findall(s)` returns a list containing all substrings within `s` that satisfy our regular expression.
- `p.finditer(s)` returns a generator. Each element of the generator is a `SRE_Match` which contains
  - `span(x,y)` indicates the start position `x` and end position `y` (of `s`).
  - `match='...'` indicates the string that satisfies the regular expression.

# QUANTIFIERS

---

- Quantifiers are operators that are applied to the preceding symbol.
- '\*' previous symbol appears 0 or more matches.
- '+' previous symbol appears 1 or more matches.
- '?' previous symbol appears at most one (0 or 1).
- '{k}' previous symbol repeated k exact matches.
- '{min,max}' previous symbol appears any number between min and max.
- For example '{2,8}' would match numbers between 2 and 8

# QUANTIFIERS

---

- "r'a.\*'" the '\*' refers to '.' making this expression get triggered when 'a' is followed by any character.
- "r'\d{4}'" the '{4}' refers to '\d' making this regular expression get triggered with 4 consecutive digits.
- Example:

```
aux = "The girl who loved the black cat ended up with a catwomen costume."
```

- p = re.compile(r'cat\s')  
p.findall(aux)
- ['cat ']
- p = re.compile(r'cat\s\*')  
p.findall(aux)
- ['cat ', 'cat ']

# SPECIAL REGEX CHARACTERS

---

- ' .' Matches any character except new line.
- '\d' Matches any digit (0-9). Equivalent to '[0-9]'.
- '\D' Matches any NON digit. Equivalent to '[^0-9]'.
- '\w' Matches any "word character". Equivalent to '[a-zA-Z0-9\_]'
  - Equivalent to '[^a-zA-Z0-9\_]'
- '\s' Matches any whitespace (space, tab, newline).
  - Equivalent to '[\t\n\r\f\v]'
- '\S' Matches any NON whitespace (space, tab, newline).
  - Equivalent to '[^ \t\n\r\f\v]'

# REGEX EXAMPLE

---

- Let us consider the following example

```
aux = "The girl who loved the black cat ended up with a catwomen costume."
```

- ```
p = re.compile(r'cat')
```

```
p.findall(aux)
```
- ```
['cat', 'cat']
```
- ```
p = re.compile(r'cat\s')
```

```
p.findall(aux)
```
- ```
['cat ']
```
- ```
p = re.compile(r'cat.*')
```

```
p.findall(aux)
```
- ```
['cat ended up with a catwomen costume. ']
```

# REGEX EXAMPLE

---

- Let us consider the following example, we want to match any word with cat prefix

```
aux = "The girl who loved the black cat ended up with a catwomen costume."
```

```
➤ p = re.compile(r'cat\w+')
p.findall(aux)
```

```
➤ ['catwomen']
```

```
➤ p = re.compile(r'cat\w*')
p.findall(aux)
```

```
➤ ['cat', 'catwomen']
```

# WE CAN USE REGEX FOR TOKENIZATION

---

```
aux = 'The girl who loved the black cat, bought a catwomen costume'
```

- p = aux.split(' ')
  
- ['The', 'girl', 'who', 'loved',  
 'the', 'black', 'cat,', 'bought',  
 'a', 'catwomen', 'costume']

In [ ]:

# WE CAN USE REGEX FOR TOKENIZATION

---

```
aux = 'The girl who loved the black cat, bought a catwomen costume'
```

- p = re.compile(r'\w+')  
p.findall(aux)
  
- ['The', 'girl', 'who', 'loved',  
'the', 'black', 'cat', 'bought',  
'a', 'catwomen', 'costume']

In [ ]:

# ANCHERS

---

- '\b' Matches any word boundary.
- '\B' Matches any NON word boundary.
- '^' Matches beginning of a string
- '[^a-e]' negates the character set. That is matches any character outside 'a-e'.
- '\$' Matches a position that is end of a string.

# CHARACTER SETS

---

- '[]' allows us to specify sets of symbols.
- '[ab3-]' matches any character 'a' or 'b' or '3' or '-'.
- '[a-g]' matches any character 'a' to 'g' such as 'b', 'c', 'd', ..., 'g'
- '[A-G]' matches any character 'A' to 'G' such as 'B', 'C', 'D', ..., 'G'
- '[a-zA-G]' matches any lowercase character and any uppercase character from 'A' to 'G'.
- '[d1-d2]' matches any digit between 'd1' and 'd2'. For example '[0-5]' matches '0', '1', '2', '3', '4', '5'

# EXAMPLES

---

- $abc^*$  matches a string that has ab followed by zero or more c
- $abc^+$  matches a string that has ab followed by one or more c
- $abc?$  matches a string that has ab followed by zero or one c
- $abc\{2\}$  matches a string that has ab followed by 2 c
- $abc\{2,\}$  matches a string that has ab followed by 2 or more c
- $abc\{2,5\}$  matches a string that has ab followed by 2 up to 5 c
- $a(bc)^*$  matches a string that has a followed by zero or more copies of the sequence bc
- $a(bc)\{2,5\}$  matches a string that has a followed by 2 up to 5 copies of the sequence bc

# GROUPS

---

- Groups are created between parenthesis
- '(abc)?' Checks if the group 'abc' appears or not.