

# TIỀN XỬ LÝ DỮ LIỆU VĂN BẢN (EXPLORING AND PREPROCESSING TEXT DATA)

AI Academy Vietnam

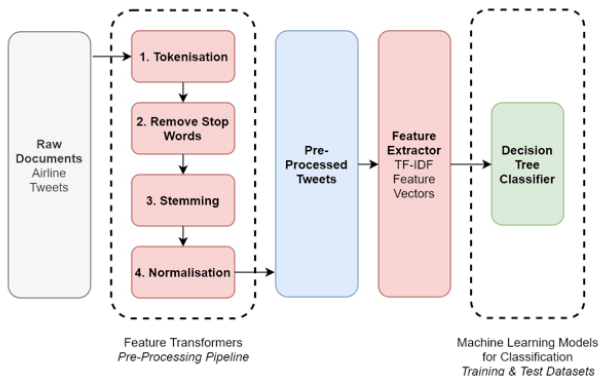
# Nội dung

## 1 Tiền xử lý dữ liệu

## 2 Các kỹ thuật tiền xử lý dữ liệu văn bản

- Chuyển chữ thường
- Loại bỏ dấu chấm (Removing Punctuation)
- Loại bỏ từ dừng (Removing stop words)
- Chuẩn hoá văn bản (Standardizing text)
- Chuẩn hoá chính tả (Correcting spelling)
- Tách từ (Tokenizing)
- Đưa về từ gốc (stemming and lemmatization)
- Khám phá dữ liệu văn bản
- Xây dựng chương trình tiền xử lý

# Tiền xử lý dữ liệu



<https://www.mlanalytics.in/how-does-text-preprocessing-in-nlp-work/>

- Tiền xử lý dữ liệu là bước quan trọng
- 70%-80% công sức là xử lý dữ liệu

# Tại sao cần tiền xử lý văn bản?



- Nhiều nguồn: web, HTML, documents....
- Chứa nhiều nhiễu (noise) và chưa được làm sạch (clean)

→ **Understandable format**

# Nội dung

## 1 Tiền xử lý dữ liệu

## 2 Các kỹ thuật tiền xử lý dữ liệu văn bản

- Chuyển chữ thường
- Loại bỏ dấu chấm (Removing Punctuation)
- Loại bỏ từ dừng (Removing stop words)
- Chuẩn hoá văn bản (Standardizing text)
- Chuẩn hoá chính tả (Correcting spelling)
- Tách từ (Tokenizing)
- Đưa về từ gốc (stemming and lemmatization)
- Khám phá dữ liệu văn bản
- Xây dựng chương trình tiền xử lý

# Chuyển chữ thường

Lệnh trừng phạt của Mỹ lên Huawei không chỉ tác động đến các công ty công nghệ Trung Quốc mà còn kéo theo nhiều hệ lụy tới ngành công nghiệp toàn cầu.



lệnh trừng phạt của mỹ lên huawei không chỉ tác động đến các công ty công nghệ trung quốc mà còn kéo theo nhiều hệ lụy tới ngành công nghiệp toàn cầu

Bài toán: tất cả dữ liệu văn bản cần được định dạng giống nhau, để chắc chắn rằng "NLP" và "nlp" là như nhau.

- **Đầu vào:** Một văn bản chứa nội dung
- **Đầu ra:** Dữ liệu văn bản đã được chuyển về chữ thường
- **Giải pháp:** Sử dụng hàm **lower()** trong python

# Loại bỏ dấu chấm (Removing Punctuation)

	A	B	C	D
1	Text with Punctuation			Remove Punctuation
2	"Apple"			Apple
3	(Pear). 5			Pear 5
4	{[Orange]}			Orange
5	Lemon;;; :::			Lemon
6	Lychee!			Lychee
7	<Blueberry>			Blueberry
8	Dash-test			Dashtest
9	TEST~!#\$%^&*()_+{} []";<>?.,			TEST



- Loại bỏ punctuation:
  - Quan trọng vì punctuation không mang thêm thông tin
  - Giảm kích thước dữ liệu, đồng thời tăng hiệu quả tính toán
- Bài toán:
  - Đầu vào:** Một văn bản chứa nội dung
  - Đầu ra:** Văn bản đã được loại bỏ các dấu chấm
  - Giải pháp:** Sử dụng biểu thức chính quy và hàm **replace()** trong python

## Loại bỏ từ dừng (Removing stop words)

Ngày cả khi trời mưa, trận đấu vẫn diễn ra



Removing stop words

Trời mưa, trận đấu diễn ra

Từ dừng: phổ biến, nhưng không mang nhiều ý nghĩa. Việc loại bỏ giúp: giảm kích thước dữ liệu, có thể cải thiện hiệu năng của mô hình.

- **Đầu vào:** Một văn bản
- **Đầu ra:** Văn bản đã được loại bỏ các từ dừng (stop words)
- **Giải pháp:**
  - Sử dụng thư viện **NLTK**
  - Xây dựng một danh sách các từ dừng, sau đó sử dụng nó để loại bỏ các từ dừng có trong văn bản. Ví dụ stop words cho tiếng Việt <sup>1</sup>

<sup>1</sup> <https://github.com/stopwords/vietnamese-stopwords>



# Chuẩn hoá văn bản (Standardizing text)

Raw	Normalized
2moro 2mrrw 2morrow 2mrw tomrw	tomorrow
b4	before
otw	on the way
:) :-) ;-)	smile

- **Đầu vào:** Một văn bản
- **Đầu ra:** Văn bản đã được chuẩn hoá
- **Giải pháp:** Tạo một từ điển để tìm kiếm các từ ngắn hoặc các từ viết tắt

# Sửa lỗi chính tả (Correcting spelling)

studing → studying

intresting → interesting

aquire → acquire

- Dữ liệu có chứa lỗi chính tả (đánh giá của người dùng, blogs, tweets...)
- Giảm số bản sao của các từ. Ví dụ: nếu không sửa thì "studing" và "studying" được coi là 2 từ khác nhau.
- **Bài toán:**
  - *Đầu vào:* Một văn bản chứa nội dung
  - *Đầu ra:* Văn bản đã được sửa lỗi chính tả
  - *Giải pháp:* Sử dụng thư viện **TextBlod**
- Ví dụ minh hoạ

# Tách từ (Tokenizing)

```
'hello e v e r y o n e           d o n t b u y t h i s p h o n e a t a l l f i r s t o
f a l l t h a t s a y s t h e p h o n e i n n e w   i t o o k i t t o t h e l a b a f t e r   m o n t h t h e
p h o n e i s d e a d d e a d   y o u c a n s a v e i t t h e y o p e n t h e p h o n e i n t h e l a b a n d s a y
s   t h e p h o n e i s r e n e w   a n d i t s c h e a p e s t c o m m o n e n t s i p a y e d   f o r o n
l y   m o n t h   n o w i n e e d t o b u y n e w o n e t h i s l q g   i s d e a d   n o t a b e s t t h i n g
p e o p l e   a r e s a y i n g t o m e d o n t b u y f r o m   a t a l l   i t s t r o u b l i n g   '
```

↓ Tokenization

```
['hello', 'e', 'v', 'e', 'r', 'y', 'o', 'n', 'e', 'dont', 'buy', 'this', 'phone', 'at', 'all', 'first', 'of', 'all', 'that', 'says', 'the', 'phone', 'in', 'new', 'i', 'took', 'it', 'to', 'the', 'lab', 'after', 'month', 'the', 'phone', 'is', 'dead', 'dead', 'you', 'can', 'save', 'it', 'they', 'open', 'the', 'phone', 'in', 'the', 'lab', 'and', 'says', 'the', 'phone', 'is', 'renew', 'and', 'its', 'cheapest', 'components', 'i', 'payed', 'for', 'only', 'month', 'now', 'i', 'need', 'to', 'buy', 'new', 'one', 'this', 'lq', 'g', 'is', 'dead', 'not', 'a', 'best', 'thing', 'people', 'are', 'saying', 'to', 'me', 'dont', 'buy', 'from', 'at', 'all', 'it', 's', 'troubling']
```

[captured from medium.com]

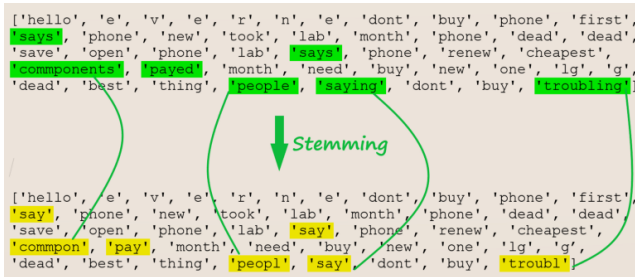
- **Đầu vào:** Một văn bản chứa nội dung
- **Đầu ra:** Các câu hoặc các từ được tách ra từ văn bản gốc
- **Giải pháp:**
  - Với tiếng Anh: Sử dụng thư viện NLTK, SpaCy, TextBlod
  - Với tiếng Việt: VnCoreNLP <sup>2</sup>, underthesea <sup>3</sup>, coccoc-tokenizer <sup>4</sup>

<sup>2</sup> <https://github.com/vncorenlp/VnCoreNLP>

<sup>3</sup> <https://github.com/undertheseanlp/underthesea>

<sup>4</sup> <https://github.com/coccoc/coccoc-tokenizer>

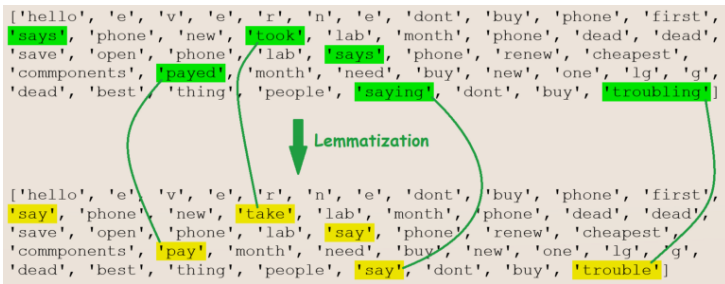
# Đưa về từ gốc (stemming)



[captured from medium.com]

- **Đầu vào:** Một một từ, như fishing, fishes
- **Đầu ra:**
  - Từ gốc của từ đó, như là fish (bỏ đi các phần tiền hoặc là hậu tố)
  - Hình thái học của từ
- **Giải pháp:**
  - Sử dụng thư viện: **NLTK**
  - Sử dụng thư viện: **TextBlod**

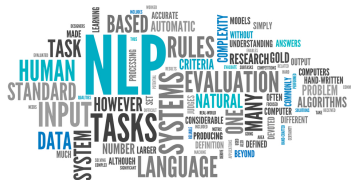
# Đưa về từ gốc (lemmatization)



[captured from medium.com]

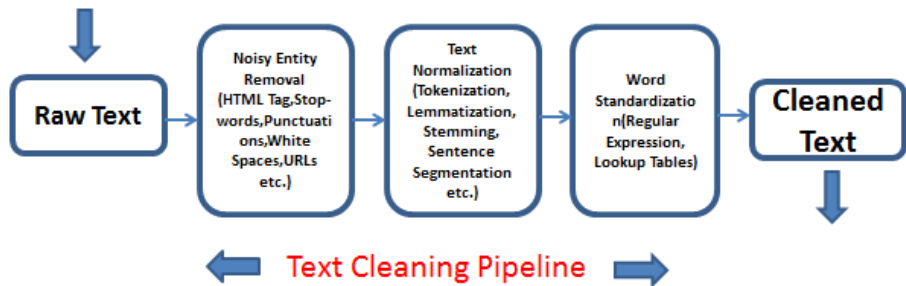
- **Đầu vào:** Một một từ, như good, best, better
- **Đầu ra:**
  - Từ gốc của từ đó, như là good.
  - Nghĩa của từ
- **Giải pháp:**
  - Sử dụng thư viện: **NLTK**
  - Sử dụng thư viện: **TextBlod**

# Phân tích và hiểu văn bản



- **Đầu vào:** Một văn bản
- **Đầu ra:**
  - Đếm số từ của một văn bản
  - Đếm tần số xuất hiện của các từ
  - Đếm từ với độ dài lớn hơn 3 và vẽ phân bố
  - Xây dựng đám mây từ (word cloud)
- **Giải pháp:**
  - Sử dụng thư viện: **NLTK**
  - Sử dụng thư viện: **TextBlod**

# Xây dựng một chương trình tiền xử lý dữ liệu văn bản



[captured from medium.com]

- **Đầu vào:** Một văn bản
- **Đầu ra:** Văn bản đã được tiền xử lý
- **Giải pháp:** Tạo một hàm tiền xử lý văn bản với tất cả các kỹ thuật tiền xử lý đã được trình bày ở trên.

# THANK YOU

## Q&A