# Machine Learning for NLP

Mouhcine MENDIL

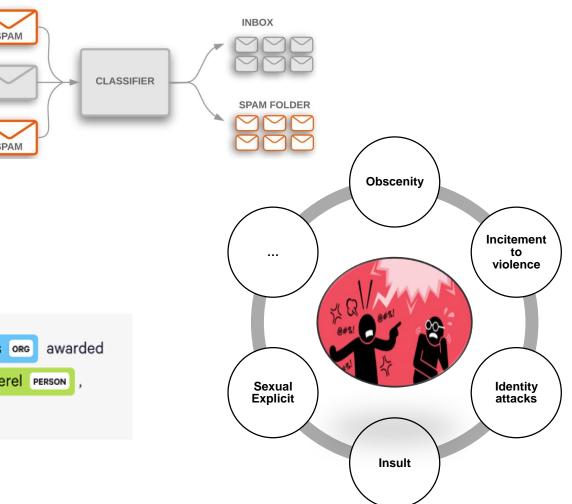# Some Tasks in NLP: Classification

**Assigning predefined labels or categories to text:**

- Hate speech detection in social media
- Spam detection
- Topic categorization
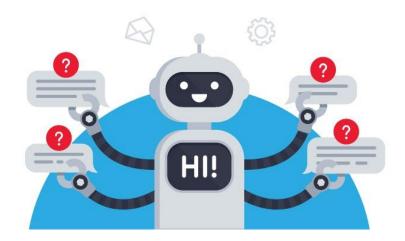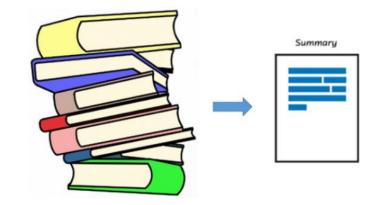- Named Entity Recognition
- ...

# Some Tasks in NLP: Text Generation

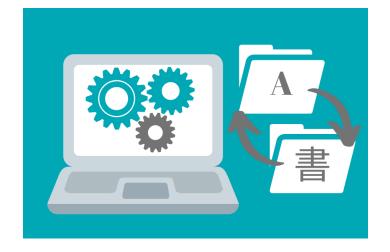**Producing coherent and contextually relevant text based on given input:**

- Machine Translation

- Text Summarization

- Paraphrasing

- Chatbots

- …

# Rules-based Systems

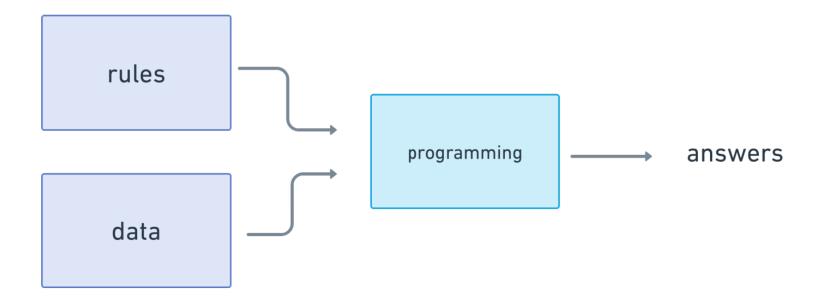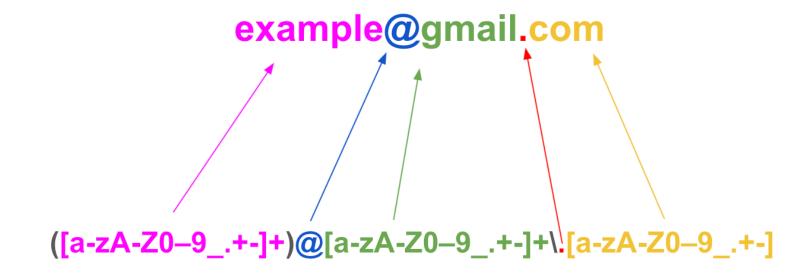- Early approaches (60's): symbolic methods, hand-written rules

# Rules-based Systems

- Regular expressions: language for specifying a matching pattern in text search

**example@gmail.com**

**([a-zA-Z0–9_.+-]+)@[a-zA-Z0–9_.+-]+\.[a-zA-Z0–9_.+-]**
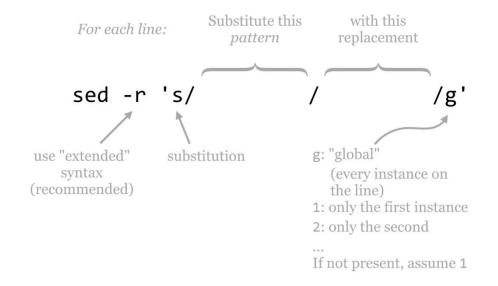
# Rules-based Systems

- Regular expressions: language for specifying a matching pattern in text search
- Useful for text normalization and substitution

```
In [3]: import re

        text = "i love python and PYTHON is great for many tasks"
        pattern = r'\bpython\b'
        repl = "Python"
        result = re.sub(pattern, repl, text, flags=re.IGNORECASE)

In [4]: print(result)

        i love Python and Python is great for many tasks
```

For each line:  Substitute this    with this
                pattern            replacement

```
sed -r 's/          /          /g'
```

use "extended"    substitution    g: "global"
syntax                            (every instance on
(recommended)                      the line)
                                  1: only the first instance
                                  2: only the second
                                  ...
                                  If not present, assume 1

**fit** | FRENCH INSTITUTES OF TECHNOLOGY

# Rules-based Systems

- ELIZA simulates a Rogerian psychologist



```
Welcome to
        EEEEEE  LL      IIII  ZZZZZZ  AAAAA
        EE      LL       II       ZZ  AA   AA
        EEEEE   LL       II      ZZZ  AAAAAAA
        EE      LL       II     ZZ    AA   AA
        EEEEEE  LLLLLL  IIII  ZZZZZZ  AA   AA

Eliza is a mock Rogerian psychotherapist.
The original program was described by Joseph Weizenbaum in 1966.
This implementation by Norbert Landsteiner 2005.


ELIZA: Is something troubling you ?
YOU:    Men are all alike.
ELIZA: What is the connection, do you suppose ?
YOU:    They're always bugging us about something or other.
ELIZA: Can you think of a specific example ?
YOU:    Well, my boyfriend made me come here.
ELIZA: Is it important to you that your boyfriend made you come here ?
YOU:    He says I'm depressed much of the time.
ELIZA: I am sorry to hear that you are depressed.
YOU:    It's true. I am unhappy.
ELIZA: Can you explain what made you unhappy ?
YOU:
```

https://web.njit.edu/~ronkowit/eliza.html

# Rules-based Systems

- ELIZA simulates a Rogerian psychologist
- Hard-coded rules based on pattern matching and substitution



https://web.njit.edu/~ronkowit/eliza.html

# Rules-based Systems

- ELIZA simulates a Rogerian psychologist
- Hard-coded rules based on pattern matching and substitution

Example of a condition and possible answers:

"?*x I want ?*y": [

    "What would it mean if you got ?y?",
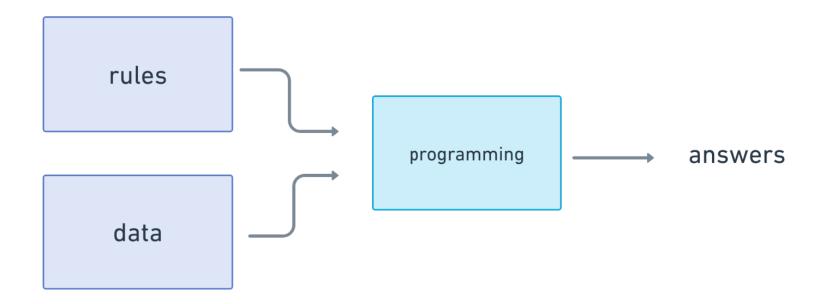
    "Why do you want ?y?",

    "Suppose you got ?y soon."

    …

    ]

```
Welcome to
                    EEEEEE  LL      IIII  ZZZZZZ   AAAAA
                    EE      LL       II       ZZ  AA   AA
                    EEEEE   LL       II      ZZZ  AAAAAAA
                    EE      LL       II     ZZ    AA   AA
                    EEEEEE  LLLLLL  IIII  ZZZZZZ  AA   AA

    Eliza is a mock Rogerian psychotherapist.
    The original program was described by Joseph Weizenbaum in 1966.
    This implementation by Norbert Landsteiner 2005.

ELIZA: Is something troubling you ?
YOU:   Men are all alike.
ELIZA: What is the connection, do you suppose ?
YOU:   They're always bugging us about something or other.
ELIZA: Can you think of a specific example ?
YOU:   Well, my boyfriend made me come here.
ELIZA: Is it important to you that your boyfriend made you come here ?
YOU:   He says I'm depressed much of the time.
ELIZA: I am sorry to hear that you are depressed.
YOU:   It's true. I am unhappy.
ELIZA: Can you explain what made you unhappy ?
YOU:
```
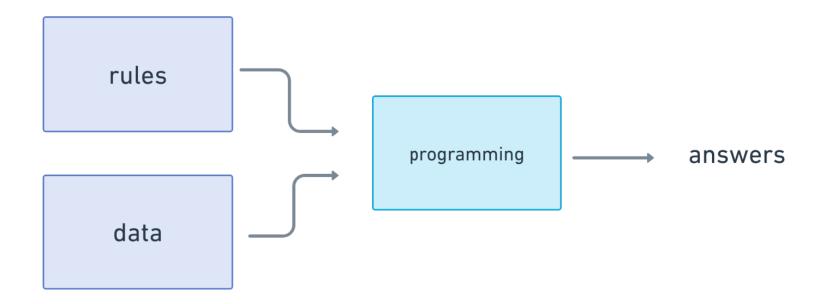
https://web.njit.edu/~ronkowit/eliza.html

FRENCH
INSTITUTES OF
TECHNOLOGY

# Rules-based Systems

- Early approaches (60's): symbolic methods, hand-written rules
- Advantages: based on expert knowledge, precise
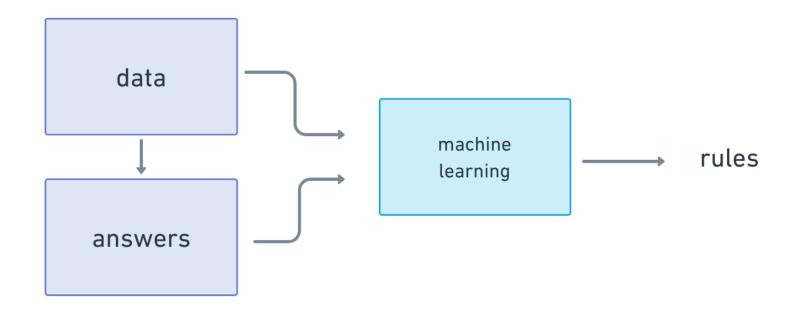
# Rules-based Systems

- Early approaches (60's): symbolic methods, hand-written rules
- Advantages: based on expert knowledge, precise
- Downsides: lack of coverage, expensive to build and maintain
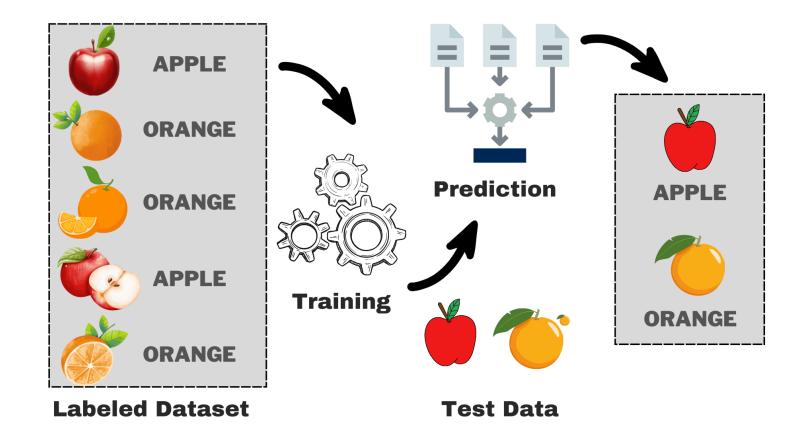
# Learning from Data (Supervised)

- Learn rules automatically: machine learning (90's), neural methods (2010's)

# Learning from Data (Supervised)

- Learn rules automatically: machine learning (90's), neural methods (2010's)

# Learning from Data (Supervised)

- Learn rules automatically: machine learning (90's), neural methods (2010's)
- Advantages: improved performance and generalization, fast

# Learning from Data (Supervised)
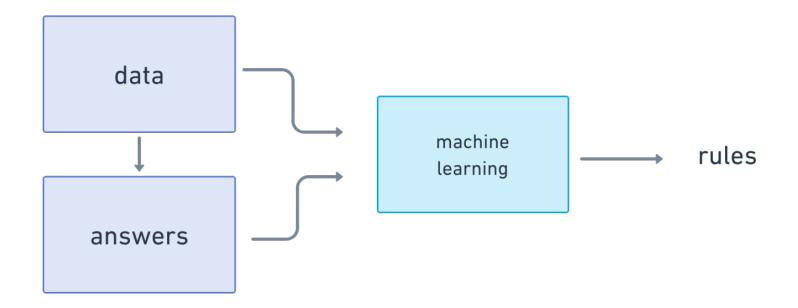
- Learn rules automatically: machine learning (90's), neural methods (2010's)
- Advantages: improved performance and generalization, fast
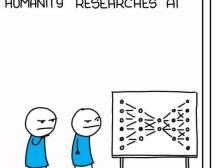- Downsides: complex, hard to interpret

# Questions so far ?

# ML Workflow (Supervised)



Raw Data

# ML Workflow (Supervised)

# ML Workflow (Supervised)

# ML Workflow (Supervised)

- How to split data into train/val/test subsets ?

# ML Workflow for NLP

- How to split data into train/val/test subsets ?

- Models need for numerical features. How to go from sentences/words/tokens to vector representation?

# ML Workflow for NLP

- How to split data into train/val/test subsets ?

- Models need for numerical features. How to go from sentences/words/tokens to vector representation?

- Are there suitable models for learning ?

FRENCH
INSTITUTES OF
TECHNOLOGY

- How to split data into train/val/test subsets ?

- Models need for numerical features. How to go from sentences/words/tokens to vector representation?

- Are there suitable models for learning ?

- Metrics for evaluation ?

# ML Models for NLP: Splitting the data

| Corpus | Document | Token | Vocabulary |
|--------|----------|-------|------------|
| Collection of documents | Collection of tokens | Collection of characters | Unique tokens |

Train          Validation          Test

FRENCH INSTITUTES OF TECHNOLOGY

# ML Models for NLP: Splitting the data

- Continuous text

- Random short sequences



Test

Validation

Training

Corpus

# ML for NLP: From language to Vectors

## Data Preparation

• Preprocessing & normlalization: lower case, special characters, stop words, stemming, tokenizing, …

## Data Preparation

• Preprocessing & normlalization: lower case, special characters, stop words, stemming, tokenizing, …

• Numerical feature vectors: ML methods requires numerical features

# ML for NLP: From language to Vectors

## Data Preparation

• Preprocessing & normlalization: lower case, special characters, stop words, stemming, tokenizing, …

• Numerical feature vectors: ML methods requires numerical features

# ML for NLP: From language to Vectors

## Data Preparation

• Preprocessing & normlalization: lower case, special characters, stop words, stemming, tokenizing, …

• Numerical feature vectors: ML methods requires numerical features



**Corpus**

**Document 1:** "Cyber threats pose a risk to organizations."

**Document 2:** "Inadequate cyber defenses expose organizations to cyber risks."

**Vocabulary** (after lemmatization and sorting by alphabetical order)
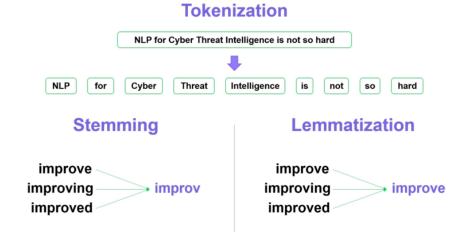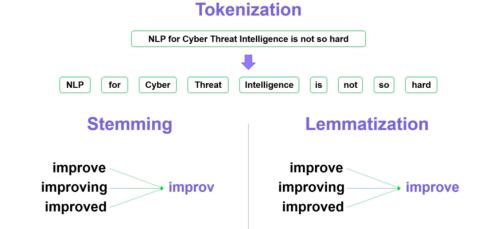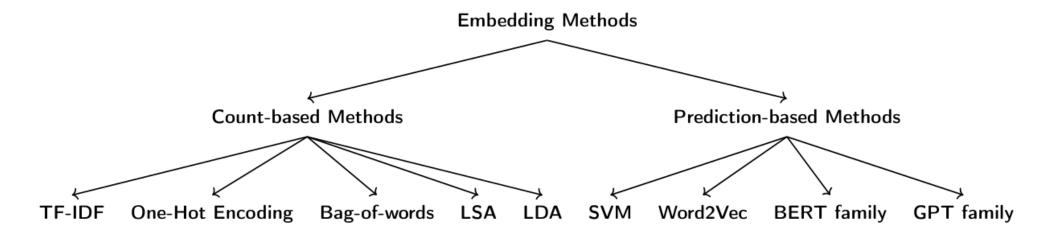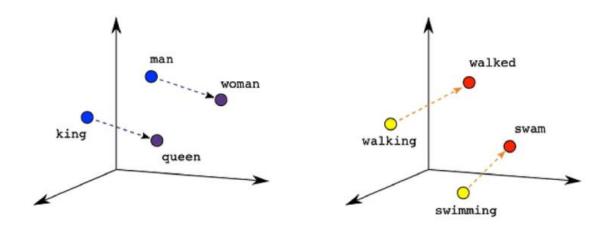
[cyber, defense, expose, inadequate, organization, pose, risk, threats]

**One-hot Encoding** of Document 1 (with the Vocabulary)

| Doc1\Voc | cyber | defense | expose | inadequate | organization | pose | risk | threats |
|---|---|---|---|---|---|---|---|---|
| Cyber | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| threats | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| pose | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| risk | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| organization | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |

**BoW** (with the Vocabulary)

| Doc\Voc | cyber | defense | expose | inadequate | organization | pose | risk | threats |
|---|---|---|---|---|---|---|---|---|
| Document 1 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| Document 2 | 2 | 1 | 1 | 1 | 1 | 0 | 1 | 1 |

**TF-IDF** of Document 2 (with the Corpus & Vocabulary)

| | cyber | defense | expose | inadequate | organization | pose | risk | threats |
|---|---|---|---|---|---|---|---|---|
| # docs w/ the token | 2 | 1 | 1 | 1 | 2 | 1 | 2 | 1 |
| IDF | $\log(\frac{2}{2})$ | $\log(\frac{2}{1})$ | $\log(\frac{2}{1})$ | $\log(\frac{2}{1})$ | $\log(\frac{2}{2})$ | $\log(\frac{2}{1})$ | $\log(\frac{2}{2})$ | $\log(\frac{2}{1})$ |
| TF(*,Doc2) | $\frac{2}{7}$ | $\frac{1}{7}$ | $\frac{1}{7}$ | $\frac{1}{7}$ | $\frac{1}{7}$ | $\frac{1}{7}$ | $\frac{1}{7}$ | $\frac{1}{7}$ |
| TF - IDF(*,Doc2) | 0 | 0.043 | 0.043 | 0.043 | 0 | 0.043 | 0 | 0 |

# ML Models for NLP: Classification

- *Input:*

  - a document d

  - a fixed set of classes $C = \{c_1, c_2, \dots, c_K\}$

  - A training set of n hand-labeled documents $(\text{d}_1, \text{c}_1), \dots, (d_n, c_n)$

- *Output:*

  - a learned classifier $\hat{f}: d \rightarrow c \in C$

- *Input:*

  - a document d
  - a fixed set of classes $C = \{c_1, c_2, \ldots, c_K\}$
  - A training set of n hand-labeled documents $(d_1, c_1), \ldots, (d_n, c_n)$

- *Output:*

  - a learned classifier $\hat{f}: d \rightarrow c \in C$
  - Actually $\hat{f}: d \rightarrow (\pi_1, \ldots, \pi_k) \in [0,1]^K$

# ML Models for NLP: Classification

- ## Any kind of classifier can be used:

  - Naïve Bayes

  - Logistic regression

  - Support-vector machines

  - Random Forest

  - k-Nearest Neighbors

  - …

# Metrics for Classification:

- Usual classification metrics can be used:



$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} = 0.75$$

$$\text{Recall} = \frac{TP}{TP + FN} = 0.25$$

$$\text{Precision} = \frac{TP}{TP + FP} = 0.33$$

$$\text{F1 Score} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} = 0.28$$

# ML Models for NLP: Language Models

# ML Models for NLP: Language Models

# ML Models for NLP: Language Models



$P(\text{table} \mid \text{The cat sat on the})$

$P(\text{chair} \mid \text{The cat sat on the})$

$P(\text{mat} \mid \text{The cat sat on the})$

...

$$P(\text{table} \mid \text{The cat sat on the}) = ?$$

# ML Models for NLP: Language Models



$P(\text{table} \mid \text{The cat sat on the})$

$P(\text{chair} \mid \text{The cat sat on the})$

$P(\text{mat} \mid \text{The cat sat on the})$

...

$$P(\text{table} \mid \text{The cat sat on the}) = ?$$

To be estimated from data

# ML Models for NLP: Language Models



$P(\text{table} \mid \text{The cat sat on the})$

$P(\text{chair} \mid \text{The cat sat on the})$

$P(\text{mat} \mid \text{The cat sat on the})$

...

$$P(\text{table} \mid \text{The cat sat on the}) = \frac{\#\ (\textit{The cat sat on the table})}{\#\ (\textit{The cat sat on the})}$$

The cat sat on the → table $\quad P(\text{table} \mid \text{The cat sat on the})$

The cat sat on the → chair $\quad P(\text{chair} \mid \text{The cat sat on the})$

The cat sat on the → mat $\quad P(\text{mat} \mid \text{The cat sat on the})$

...

Input

Prediction

$$P(\text{table} \mid \text{The cat sat on the}) = \frac{\#\ (\textit{The cat sat on the table})}{\#\ (\textit{The cat sat on the})}$$

**We'll hardly see enough data for estimating these**

# ML Models for NLP: Language Models



$P(\text{table} \mid \text{The cat sat on the})$

$P(\text{chair} \mid \text{The cat sat on the})$

$P(\text{mat} \mid \text{The cat sat on the})$

...

$$P(\text{table} \mid \text{The cat sat on the}) \approx \begin{cases} P(table) \ or \\ P(table \mid the) \ or \\ P(table \mid on \ the) \ or \\ P(table \mid sat \ on \ the) \ or \\ ... \\ P(w_i \mid w_1 w_2 ... w_{i-1}) \end{cases}$$

**Markov Assumption**

# ML Models for NLP: Language Models



$P(\text{table} \mid \text{The cat sat on the})$

$P(\text{chair} \mid \text{The cat sat on the})$

$P(\text{mat} \mid \text{The cat sat on the})$

...

$$P(\text{table} \mid \text{The cat sat on the}) \approx \begin{cases} P(table) \; or \\ P(table \mid the) \; or \\ P(table \mid on\; the) \; or \\ P(table \mid sat\; on\; the) \; or \\ ... \\ P(w_i \mid w_1 w_2 \ldots w_{i-1}) \end{cases}$$

**Markov Assumption**

**Increasing context and complexity**

# ML Models for NLP: Language Models

Training corpus

$$P(w_i \mid w_{i-1}) = \frac{\#\,(w_{i-1}, w_i)}{\#\,(w_{i-1})}$$

<s> I am a Student </s>
<s> Student I am </s>
<s> I do like Machine Learning </s>

$P(I \mid <s>) =$

$P(Student \mid <s>) =$

$P(am \mid I) =$

$P(</s> \mid Student) =$

$P(Student \mid am) =$

$P(do \mid I) =$

# ML Models for NLP: Language Models

Training corpus

$$P(w_i \mid w_{i-1}) = \frac{\#\,(w_{i-1}, w_i)}{\#\,(w_{i-1})}$$

<s> I am a Student </s>
<s> Student I am </s>
<s> I do like Machine Learning </s>

$$P(I \mid <s>) = \frac{2}{3}$$

$$P(Student \mid <s>) = \frac{1}{3}$$

$$P(am \mid I) = \frac{2}{3}$$

$$P(</s> \mid Student) = \frac{1}{2}$$

$$P(Student \mid am) = \frac{1}{2}$$

$$P(do \mid I) = \frac{1}{3}$$

FRENCH
INSTITUTES OF
TECHNOLOGY

# Metrics for text generation

- **Perplexity**: Measures how well the model predicts a (test) sequence (used for language models)

- **BLEU (Bilingual Evaluation Understudy)**: Measures n-gram overlap between generated and reference text (used for machine translation).

- **ROUGE (Recall-Oriented Understudy for Gisting Evaluation)**: Focuses on content overlap (used for summarization).

- …

**FRENCH INSTITUTES OF TECHNOLOGY**

# NLP Workflow

**01.**
**Data Collection**

Open Source database
Web scrapping
Crowdsourcing
Social Media data
…

**02.**
**Data Cleaning**

Unicode normalization
Regular Expression
OCR
Deduplication
…

**03.**
**Pre-Processing**

Tokenization
Lowercasing
StopWord removal
Stemming
…

**04.**
**Feature Eng.**

Bag of Words
TF-IDF
N-grams
Word Embedding
…

**Feedback loop**

**07.**
**Deployment**

Export
Input pipeline
Scaling
Monitoring
…

**06.**
**Evaluation**

F1-score
MAP
ROUGE
BLEU
…

**05.**
**Model Building**

Naive Bayes
SVM
RNN, LSTM
Transformers
…

FRENCH
INSTITUTES OF
TECHNOLOGY

When you penalize your Natural Language Generation model for large sentence lengths

Me think, why waste time say lot word, when few word do trick.