

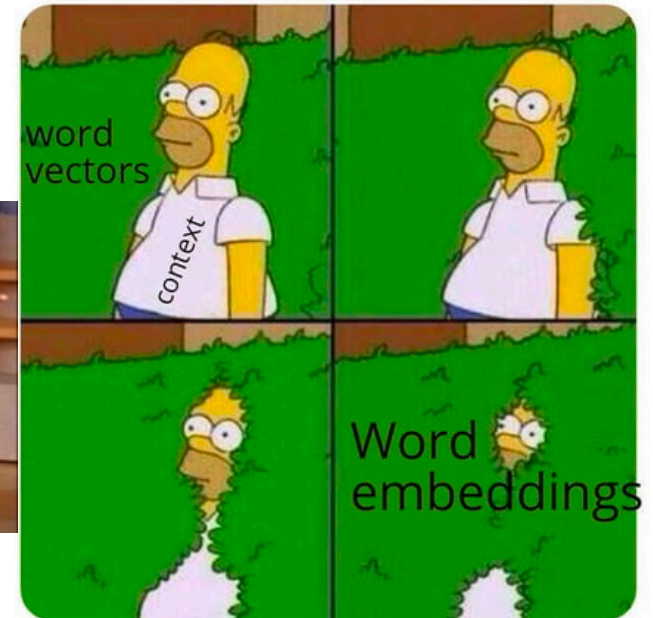
April 17, 2023 .

DSML : NLP module .

Word embeddings
in a nutshell

Introduction to NLP

Class starts
@ 9:05



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



Computer Vision / Deep Learning: Main ideas:

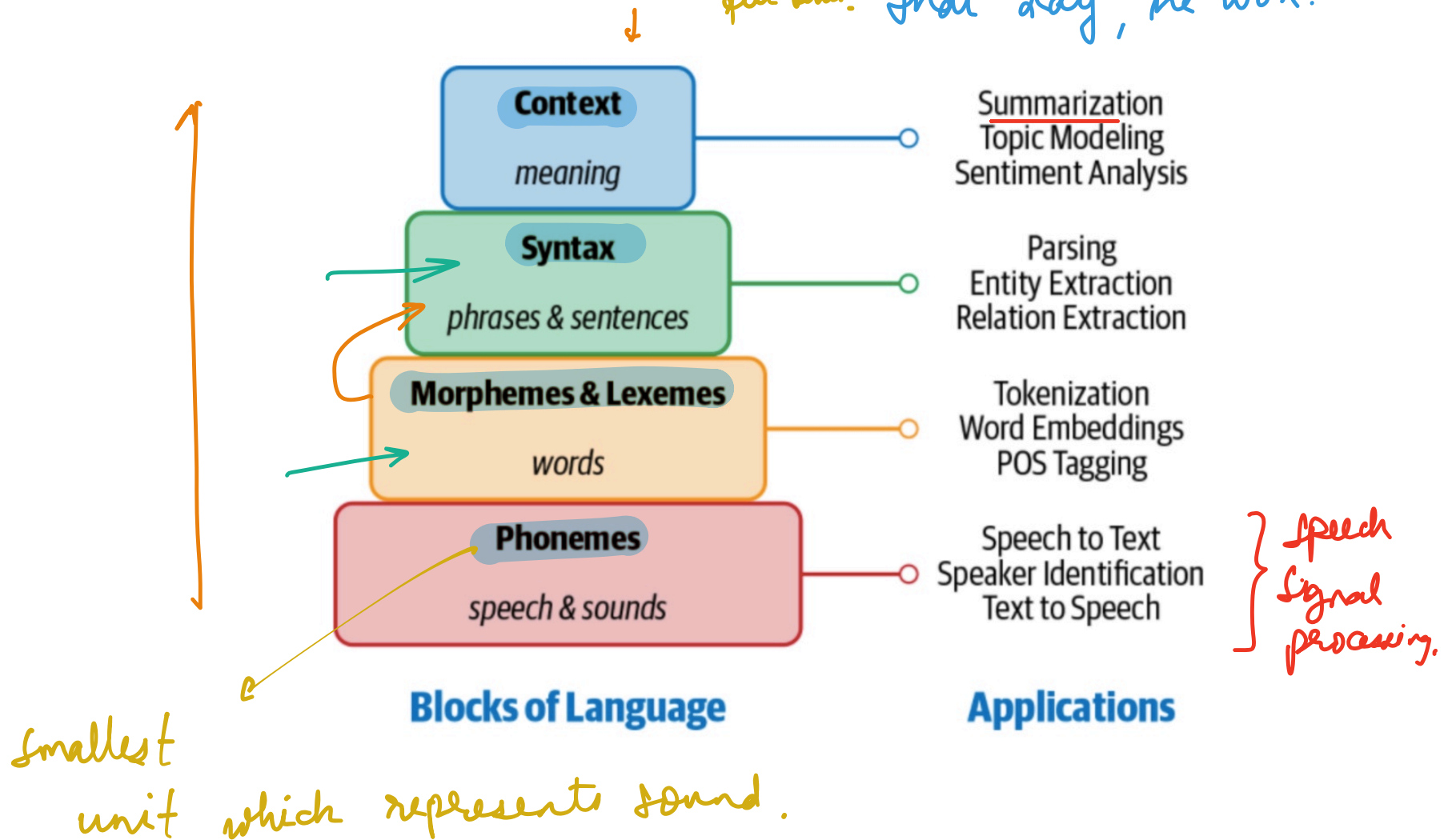
- 1] Low-level computational blocks: Conv, Pooling, Activation.
- 2] Transfer learning.
- 3] Base task: Image classification
- 4] Architectures: skip connections, Inception module, 1×1 convs. ROI, RPN etc.

* BERT, word2vec, RNN, LSTM, Transformers.

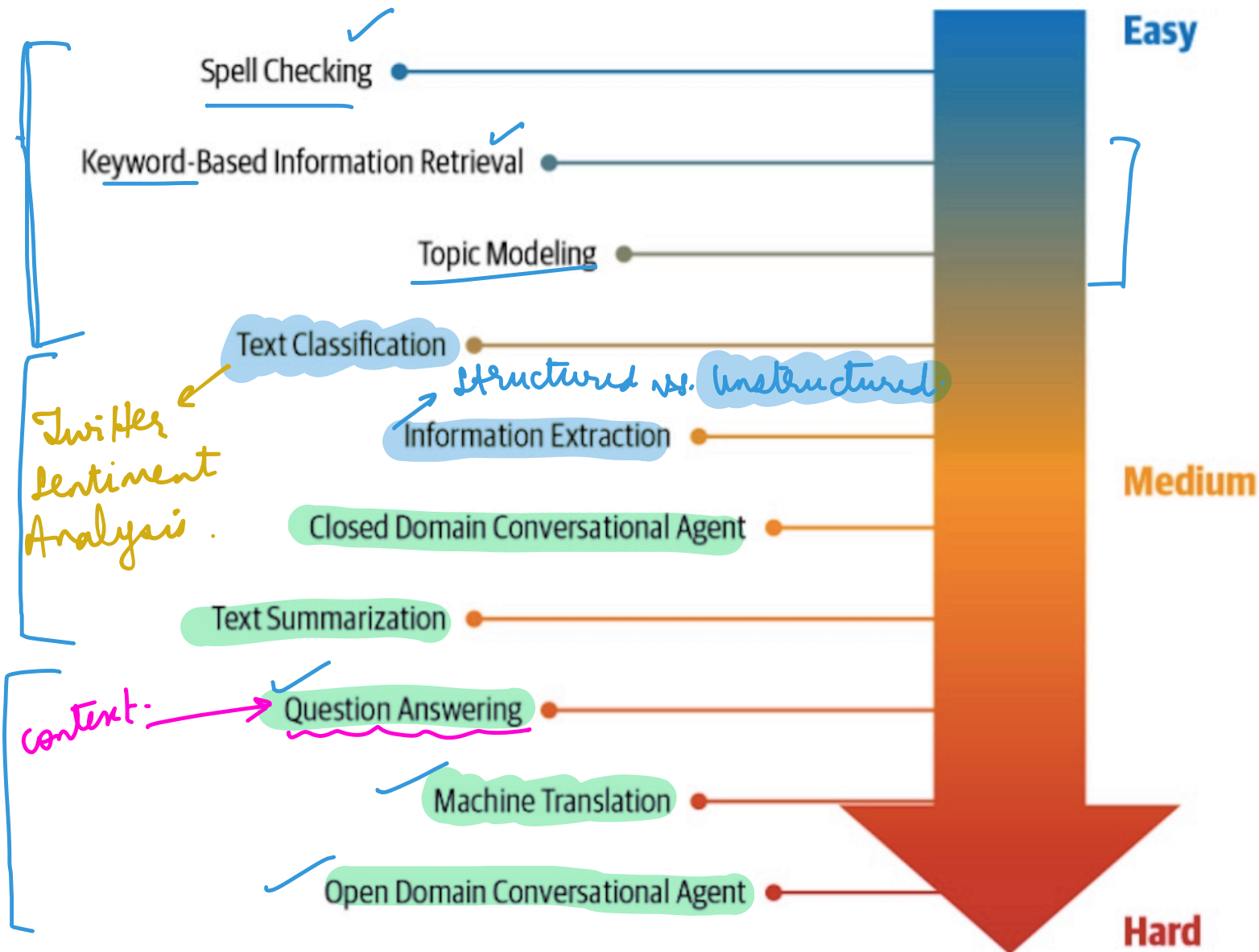
Natural language processing

→ Words, sounds.

He usually lost to make her
feel better. That day, he won.



NLP: hierarchy of difficulty.



* pre - suffix meaning before .

$\textcircled{x} \rightarrow \textcircled{y}, \textcircled{z} \rightarrow$ Suffixes meaning the same thing .
↑ ↑ ↑
morphemes .

$\textcircled{un}, \textcircled{in}, \textcircled{dis}, \textcircled{mis}$.
↑ ↑

Broad Approaches to NLP Problems;

1] Heuristics based Approach.

2] Rule-based approaches.

3] Optimization / learning.

→ Develop a Heuristics based approach to solve Twitter sentiment analysis.

Our first NLP problem:

Twitter dataset - Covid-19 tweets.

Objective: sentiment analysis on Covid-19 tweets.

↓
positive

↓
negative.

Regular Expressions.

↳ Computer Science: Theory of Computation / compiler theory.

* What is a regular expression?

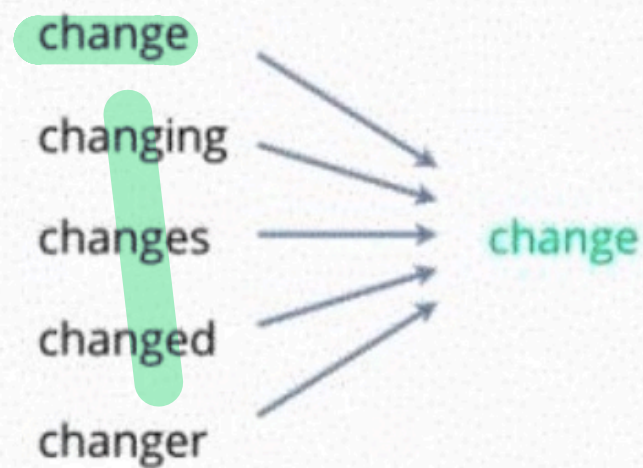
→ It is a set of pre-defined rules about characters or text, which help us "cluster" them into unnamed clusters.

* Application: Tokenization

↓
The process of extracting individual words from a tweet.

These methods help us further
reduce words to its root form.

Stemming vs Lemmatization



Converting tweets to feature vectors.

Feature vector, All tweets have variable lengths.

$\begin{bmatrix} 1, & \text{count_positive}, & \text{count_negative} \end{bmatrix}$

Every tweet vectorization: $\in \mathbb{R}^3$.

Eg: $\{$

$\begin{aligned} & ("eat", 0) \rightarrow \text{wp} \\ & ("eat", 1) \rightarrow \text{rw} \\ & ("hell", 0) \end{aligned}$

counts: $\begin{aligned} & 67, \\ & 63, \\ & 117, \\ & \vdots \end{aligned}$

Feature extraction

freqs: dictionary mapping from (word, class) to frequency

$$X_m = [1, \sum_w \textit{freqs}(w, 1), \sum_w \textit{freqs}(w, 0)]$$

↓ ↓ ↓ ↓

Features of Bias Sum Pos. Sum Neg.
tweet m Frequencies Frequencies

$$D = \{(\bar{x}_i, \bar{y}_i)\}_{i=1}^n$$

$$\bar{x}_i \in \mathbb{R}^d \quad d = 3$$

$$\bar{y}_i \in \mathbb{R}.$$

ANN, logistic / Linear Regression.

DT \rightarrow Decision trees.

NB \rightarrow Naive Bayes.

What could we have done better?

- * We could have gone for sparse (Bow) representation.
- * Use alternative representation (vectorizers, embeddings etc.)

* Normalizing our vector representation.

* Counts of synonym words. (No semantics captured in the representation).

* We have not used contextual information in any way.

Problem : Driver Assistance.

Driver drowsiness.

Road actions monitoring.

Trip monitoring.

object detection.

Pedestrians,
Emergency vehicles.

Free space