

Week 1 – Intro and TF-IDF

EGCO467 Natural Language and Speech Processing

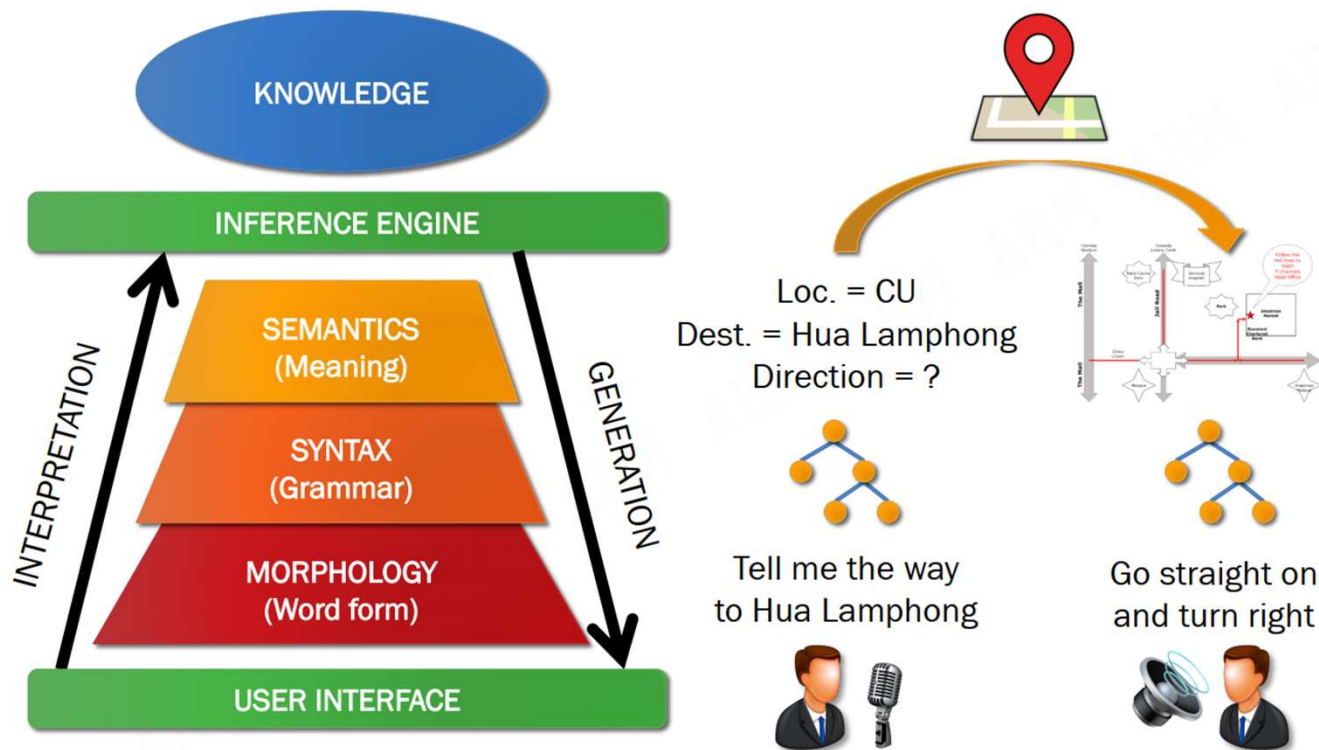
Topics covered (tentative)

- tf-idf
- word embedding
- neural networks / DL
- sentence classification
- language models
- NER/POS
- seq2seq models/attention models
- BERT and transformer models
- question answering models
- machine translation

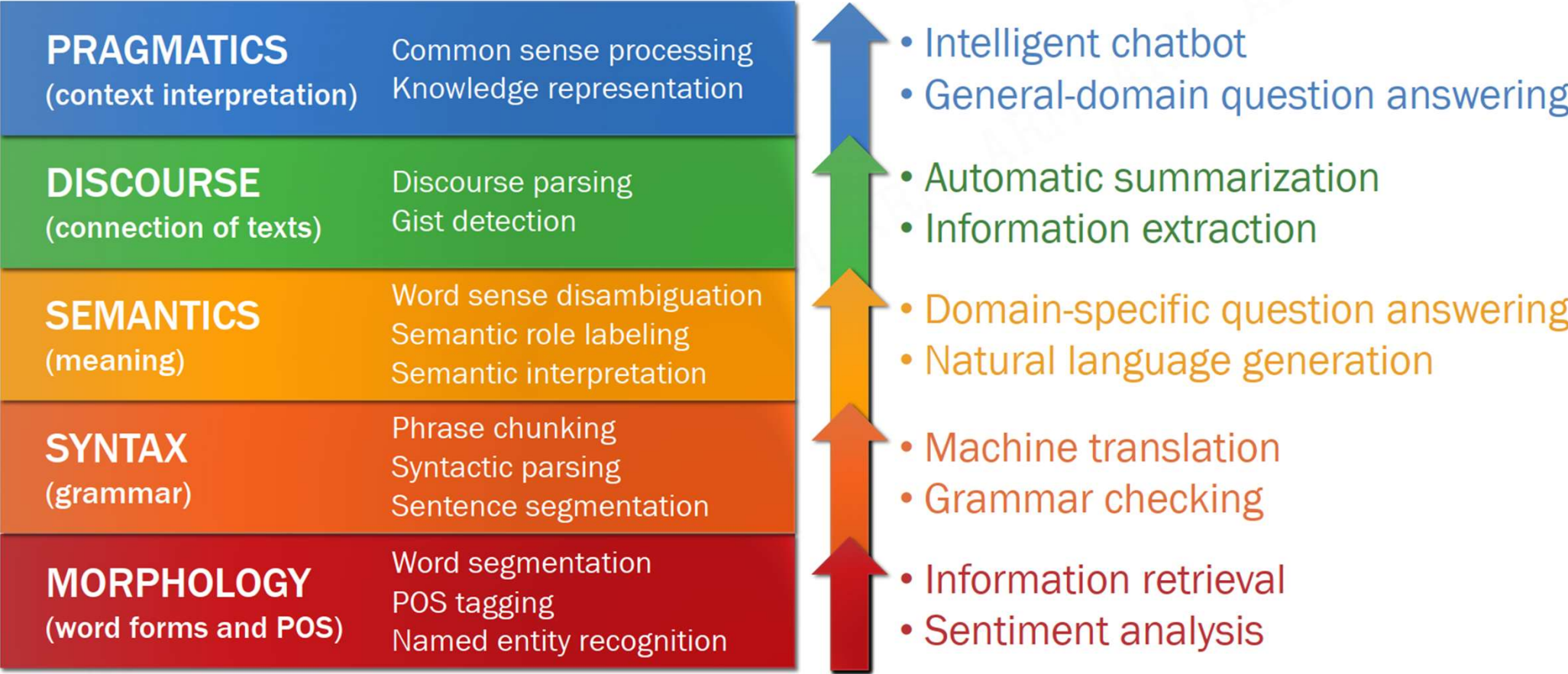
evaluation

- homework 30%
- project 70%

Natural Language Processing



Levels of NLP



Eras of NLP

$$\frac{}{\Delta, B \rightarrow B} \text{initial}^\dagger$$
$$\frac{\Delta \rightarrow B \quad \Delta \rightarrow C}{\Delta \rightarrow B \wedge C} \wedge R^\dagger$$
$$\frac{\Delta \rightarrow B}{\Delta \rightarrow B \vee C} \vee R$$
$$\frac{B, \Delta \rightarrow C}{\Delta \rightarrow B \supset C} \supset R^\dagger$$
$$\frac{A \supset (B \supset C), \Delta \rightarrow G}{(A \wedge B) \supset C, \Delta \rightarrow G} \supset L_2^\dagger$$
$$\frac{B \supset C, \Delta \rightarrow A \supset B \quad C, \Delta \rightarrow G}{(A \supset B) \supset C, \Delta \rightarrow G} \supset L_4$$

$$\frac{B, C, \Delta \rightarrow G}{B \wedge C, \Delta \rightarrow G} \wedge L^\dagger$$
$$\frac{B, \Delta \rightarrow G \quad C, \Delta \rightarrow G}{B \vee C, \Delta \rightarrow G} \vee L^\dagger$$
$$\frac{\Delta \rightarrow C}{\Delta \rightarrow B \vee C} \vee R$$
$$\frac{C, B, \Delta \rightarrow G}{B \supset C, B, \Delta \rightarrow G} \supset L_1$$
$$\frac{A \supset C, B \supset C, \Delta \rightarrow G}{(A \vee B) \supset C, \Delta \rightarrow G} \supset L_3^\dagger$$

[1]

Symbolic Logic (1949-1989)

- Handcrafted rules and constraints for logical deduction
- Computable knowledge representation
- Theoretical upper bound of language processing



[2]

Statistical Methods (1989-2009)

- Learning to generalize from a large dataset
- Explicit features required
- “Every time I fire a linguist, the performance of the speech recognizer goes up.” (F. Jelinek, 1985)



Deep Learning (2009-now)

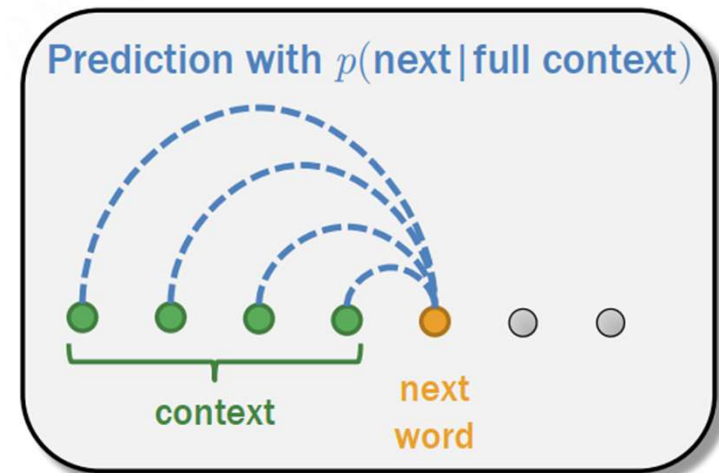
- Extracting features automatically from a very large dataset
- Yielding very high accuracy
- Requiring **powerful** hardware
- Lacking explainability because the models are rather treated as a blackbox

[1] <https://philoslife.wordpress.com/2017/05/22/course-note-symbolic-logic-uses/>
[2] <https://online.stanford.edu/courses/hrp259-introduction-probability-and-statistics-epidemiology>

Language Models

- Interpreted as a generative model
 1. Generate the first word w_1
 2. Keep generating the **next word** w_k based on the previous words (a.k.a. **context**) $w_1 \dots w_{k-1}$ until the whole sentence of length N is produced

$$P(w_1 \dots w_N) = p(w_1) \prod_{k=2}^N p(w_k | w_1 \dots w_{k-1})$$

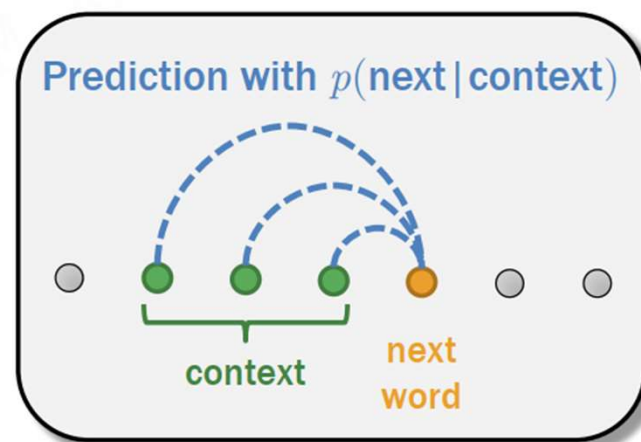


n-gram models

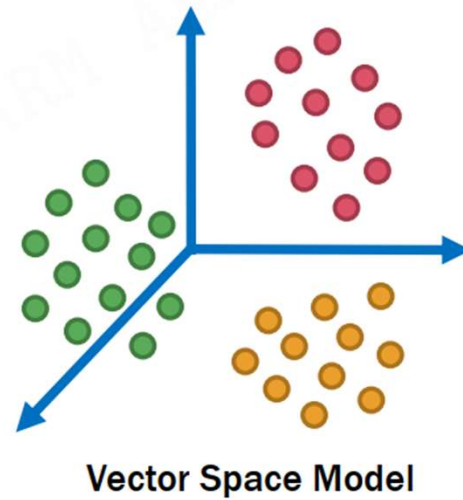
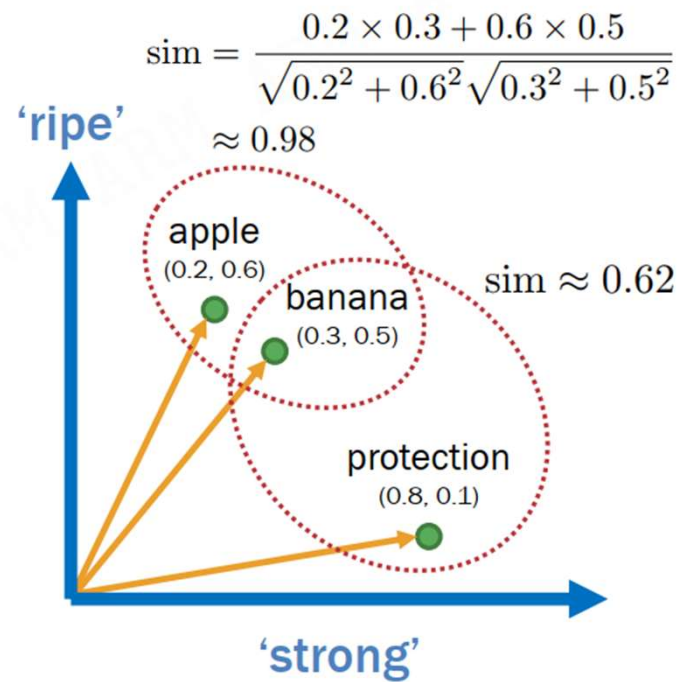
- Language models whose context is truncated to at most $n-1$ previous words

$$P(w_1 \dots w_N) = p(w_1) \prod_{k=2}^N p(\overset{\text{next word}}{w_k} | \overset{n-1 \text{ prev words}}{w_{k-n+1} \dots w_{k-1}})$$

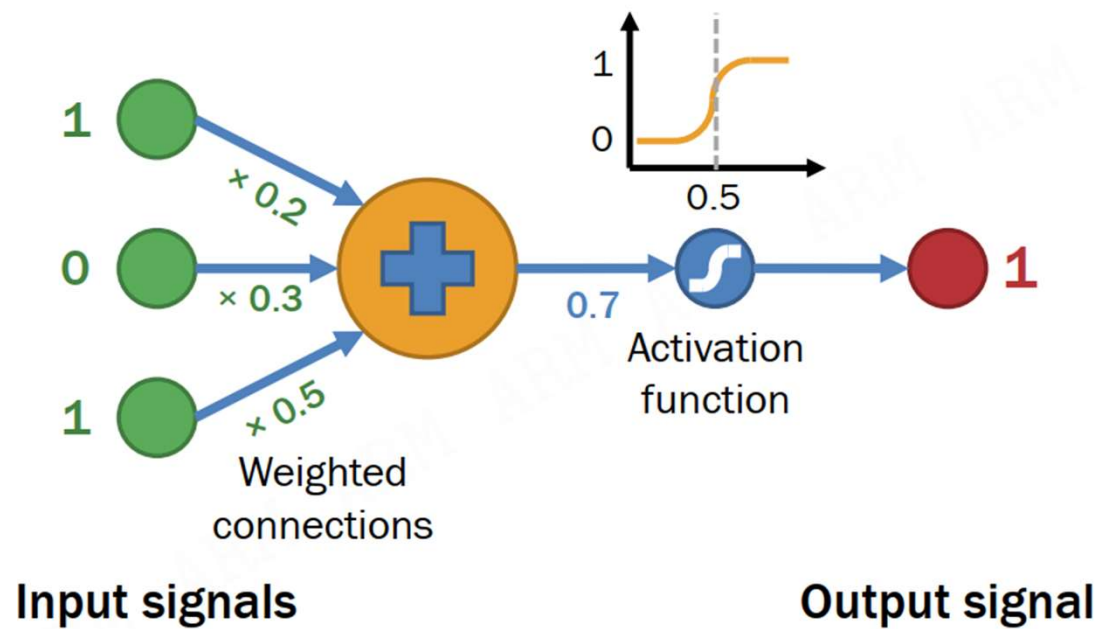
- Unigram** ($n=1$): $P(w_1 \dots w_N) = \prod_{k=1}^N p(w_k)$
- Bigram** ($n=2$): $P(w_1 \dots w_N) = p(w_1) \prod_{k=2}^N p(w_k | w_{k-1})$
- Trigram** ($n=3$): $P(w_1 \dots w_N) = p(w_1)p(w_2|w_1) \prod_{k=3}^N p(w_k | w_{k-2}, w_{k-1})$



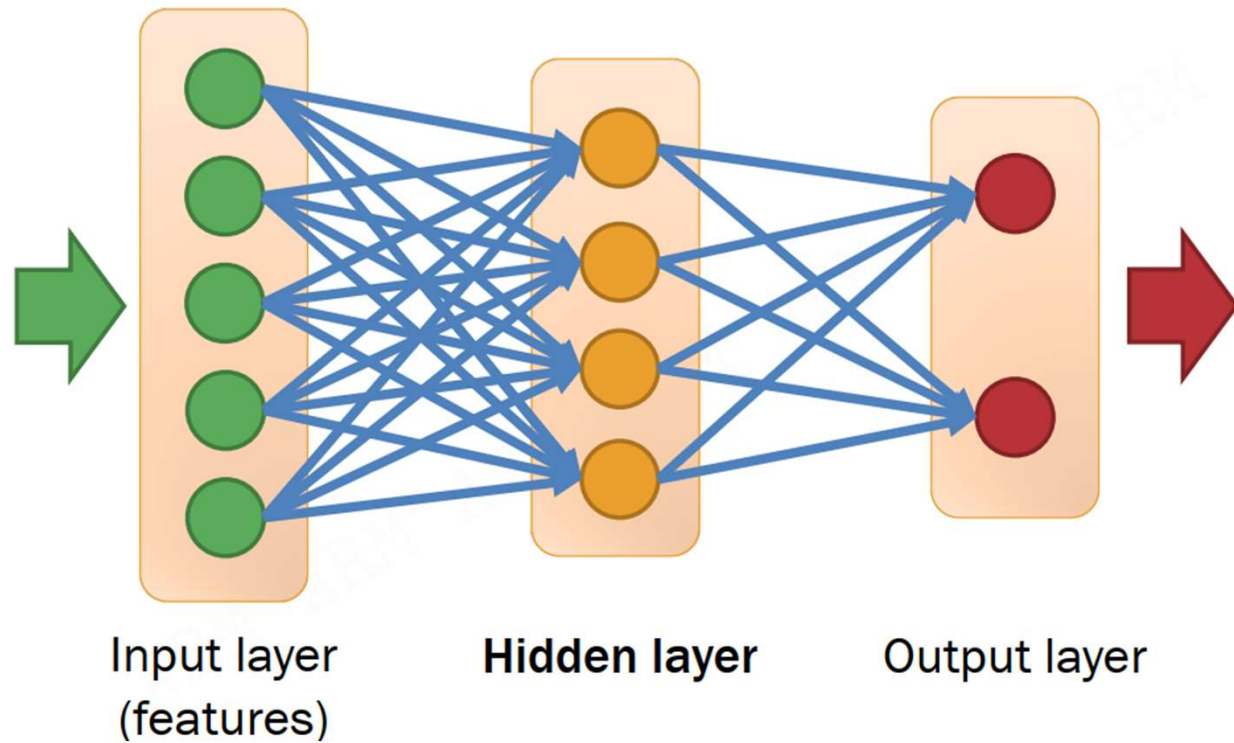
Representing Words



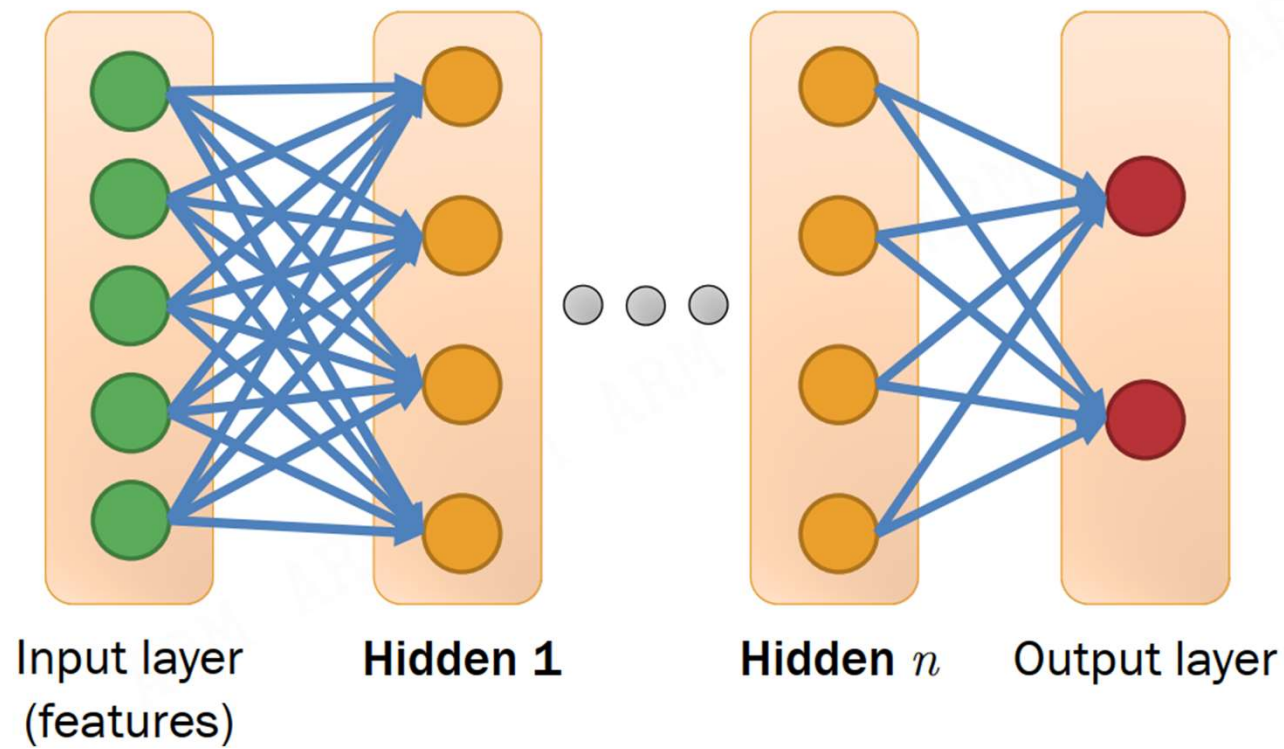
Perceptron



Multi-Layer Perceptron

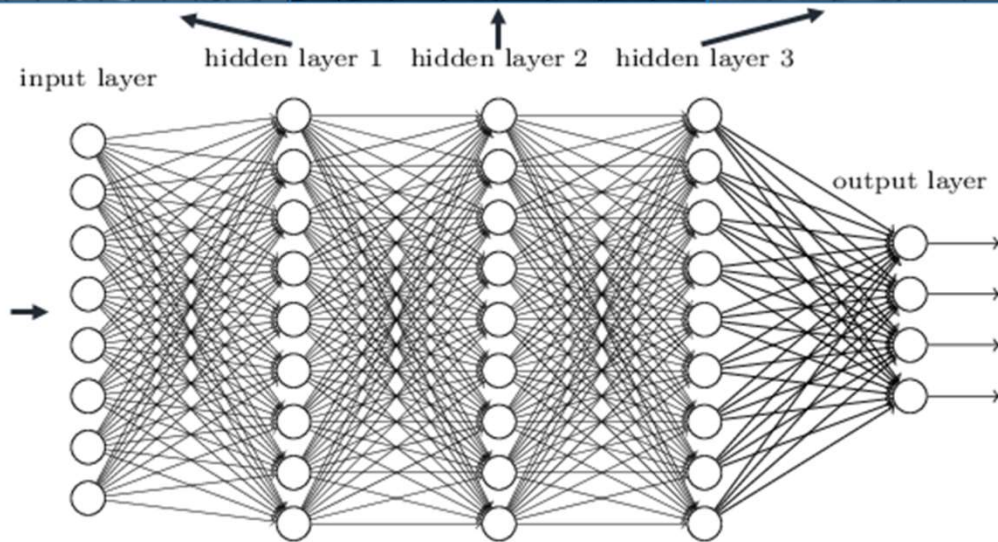


Deep Neural Network



Learning Representation

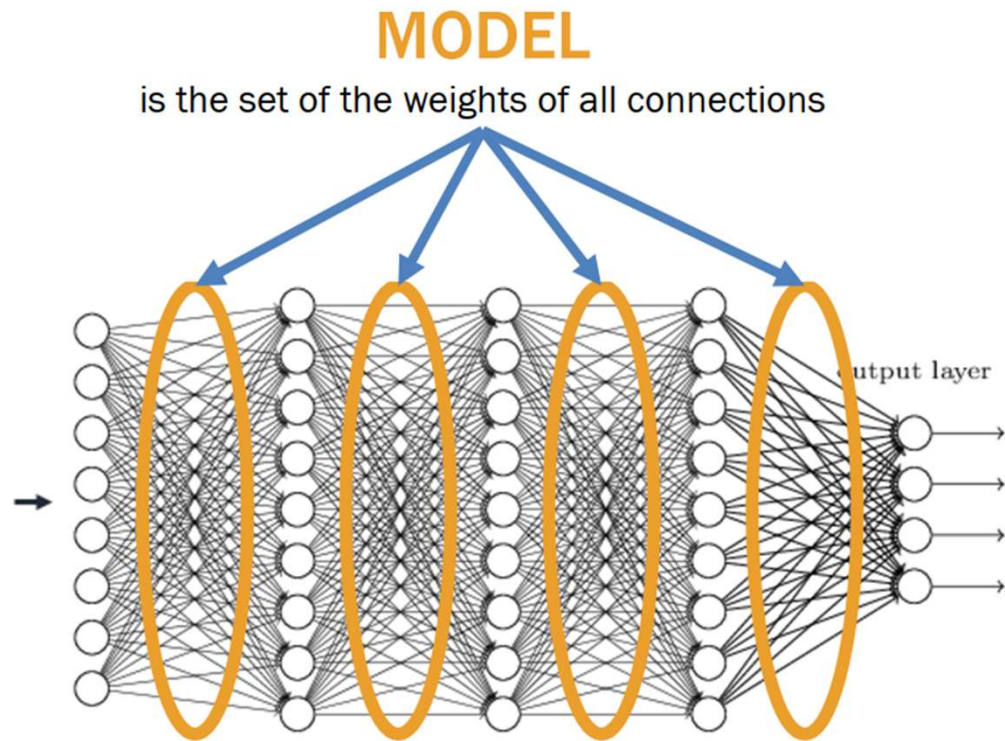
Deep neural networks learn hierarchical feature representations



<http://www.strong.io/blog/deep-neural-networks-go-to-the-movies>

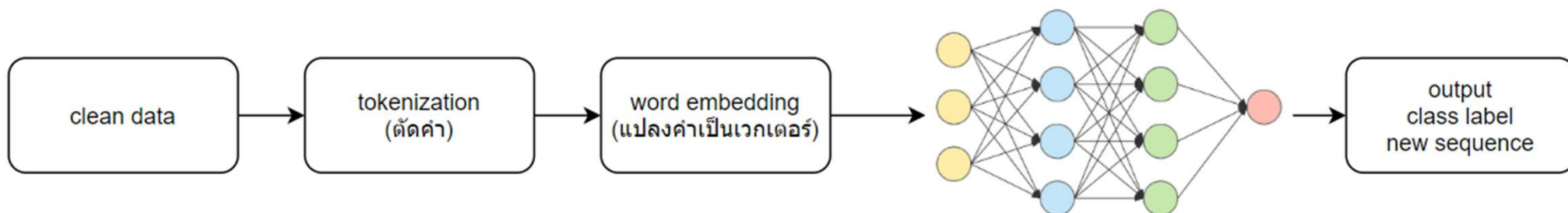
Learning Representation

Deep neural networks learn hierarchical feature representations



<http://www.strong.io/blog/deep-neural-networks-go-to-the-movies>

Workflow for NLP



Libraries

- Python
- Spacy <https://spacy.io/>
- Gensim <https://radimrehurek.com/gensim/>
- scikit-learn <https://scikit-learn.org/stable/index.html>
- Tensorflow/Pytorch
- Huggingface <https://huggingface.co/transformers/>

TF IDF

(Term frequency inverse document frequency)

TF-IDF

- TF = term frequency
- IDF = inverse document frequency
- TF-IDF = TF x IDF

Definition

- document = 1 piece of text d
- corpus = set of documents $D = \{d_1, d_2, \dots, d_N\}$
- term = a word in a document

Example TF

- d_1 = ฉันชอบกินข้าวมันไก่ร้านนี้ -> ฉัน ชอบ กิน ข้าว มัน ไก่ ร้าน นี้
- d_2 = ฉันไม่ชอบหนังเรื่องนี้เลย -> ฉัน ไม่ ชอบ หนัง เรื่อง นี้ เลย

term\document	tf(t,d1)	tf(t,d2)
ฉัน	1	1
ชอบ	1	1
กิน	1	0
ข้าว	1	0
มัน	1	0
ไก่	1	0
ร้าน	1	0
นี้	1	1
ไม่	0	1
หนัง	0	1
เรื่อง	0	1
เลย	0	1

TF has many forms

term\document	tf(t,d1)	tf(t,d2)
ฉัน	1	1
ชอบ	1	1
กิน	1	0
ข้าว	1	0
มัน	1	0
ไก่	1	0
ร้าน	1	0
นี้	1	1
ไม่	0	1
หนัง	0	1
เรื่อง	0	1
เลย	0	1

$$\text{tf}(t, d) = 0.5 + 0.5 \cdot \frac{f_{t,d}}{\max\{f_{t',d} : t' \in d\}}$$

Variants of term frequency (tf) weight

weighting scheme	tf weight
binary	0, 1
raw count	$f_{t,d}$
term frequency	$f_{t,d} / \sum_{t' \in d} f_{t',d}$
log normalization	$\log(1 + f_{t,d})$
double normalization 0.5	$0.5 + 0.5 \cdot \frac{f_{t,d}}{\max_{\{t' \in d\}} f_{t',d}}$
double normalization K	$K + (1 - K) \frac{f_{t,d}}{\max_{\{t' \in d\}} f_{t',d}}$

IDF

term\document	tf(t,d1)	tf(t,d2)	มีอยู่ในกี่ document
ฉัน	1	1	2
ชอบ	1	1	2
กิน	1	0	1
ข้าว	1	0	1
มัน	1	0	1
ไก่	1	0	1
ร้าน	1	0	1
นี้	1	1	2
ไม่	0	1	1
หนัง	0	1	1
เรื่อง	0	1	1
เลย	0	1	1

$$\text{idf}(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|}$$

term\document	tf(t,d1)	tf(t,d2)	มีอยู่ในกี่ document	idf(t,D)	tf-idf(t,d1,D)	tf-idf(t,d2,D)
ฉัน	1	1	2	0.0000	0.0000	0.0000
ชอบ	1	1	2	0.0000	0.0000	0.0000
กิน	1	0	1	0.3010	0.3010	0.0000
ข้าว	1	0	1	0.3010	0.3010	0.0000
มัน	1	0	1	0.3010	0.3010	0.0000
ไก่	1	0	1	0.3010	0.3010	0.0000
ร้าน	1	0	1	0.3010	0.3010	0.0000
นี้	1	1	2	0.0000	0.0000	0.0000
ไม่	0	1	1	0.3010	0.0000	0.3010
หนัง	0	1	1	0.3010	0.0000	0.3010
เรื่อง	0	1	1	0.3010	0.0000	0.3010
เลย	0	1	1	0.3010	0.0000	0.3010

$$\text{idf}(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|}$$

IDF also has many forms

Variants of inverse document frequency (idf) weight

weighting scheme	idf weight ($n_t = \{d \in D : t \in d\} $)
unary	1
inverse document frequency	$\log \frac{N}{n_t} = -\log \frac{n_t}{N}$
inverse document frequency smooth	$\log \left(\frac{N}{1 + n_t} \right) + 1$
inverse document frequency max	$\log \left(\frac{\max_{\{t' \in d\}} n_{t'}}{1 + n_t} \right)$
probabilistic inverse document frequency	$\log \frac{N - n_t}{n_t}$

$$\text{idf}(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|}$$

Concept of TF-IDF

- TF words that appear a lot in a document is indicative of the topic of that document. E.g. the word "interest rate" in an economic news article.
- IDF words that appear in every documents are not indicative of the topic. E.g. common words such as: "the", "a", "that", "is"
- Therefore, words with high TF-IDF values, are words that appear often in *some* documents only such as: cabinet (politics) gold (economic).

Exercise

calculate the TF-IDF table (by hand or using Excel)

- d1 = the man went out for a walk
- d2 = the children sat around the fire
- Use the *raw count* form for TF and the *inverse document frequency* form for IDF.

Solution

term	tf(t,d1)	tf(t,d2)	nt	idf	tf-idf(d1)	tf-idf(d2)
the	1	2	2	0	0	0
man	1	0	1	0.30103	0.30103	0
went	1	0	1	0.30103	0.30103	0
out	1	0	1	0.30103	0.30103	0
for	1	0	1	0.30103	0.30103	0
a	1	0	1	0.30103	0.30103	0
walk	1	0	1	0.30103	0.30103	0
children	0	1	1	0.30103	0	0.30103
sat	0	1	1	0.30103	0	0.30103
around	0	1	1	0.30103	0	0.30103
fire	0	1	1	0.30103	0	0.30103

Homework 1

- Write Python function `tf_idf(str1, str2)`
- Function takes two string. These are the two documents.
- Return the tf-idf table as Numpy array or Pandas DataFrame
- Submit your work as .ipynb file