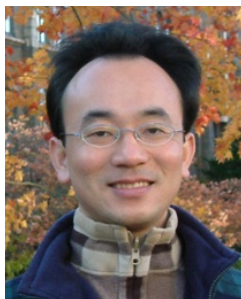# 文本表達特徵工程
# (Feature Engineering for Text Representation)

Time: 2020/05/29 (Fri) (9:10 -12:00)
Place: 國立臺北護理健康大學 (台北市明德路365號) G210
Host: 祝國忠 院長 (健康科技學院院長)

**Min-Yuh Day**
**戴敏育**
**Associate Professor**
**副教授**
**Dept. of Information Management, Tamkang University**
**淡江大學 資訊管理學系**

http://mail. tku.edu.tw/myday/
2020-05-29

# Topics

1. 自然語言處理核心技術與文字探勘
   (Core Technologies of Natural Language Processing and Text Mining)

2. 人工智慧文本分析基礎與應用
   **(Artificial Intelligence for Text Analytics: Foundations and Applications)**

3. 文本表達特徵工程
   **(Feature Engineering for Text Representation)**

4. 語意分析和命名實體識別
   (Semantic Analysis and Named Entity Recognition; NER)

5. 深度學習和通用句子嵌入模型
   (Deep Learning and Universal Sentence-Embedding Models)
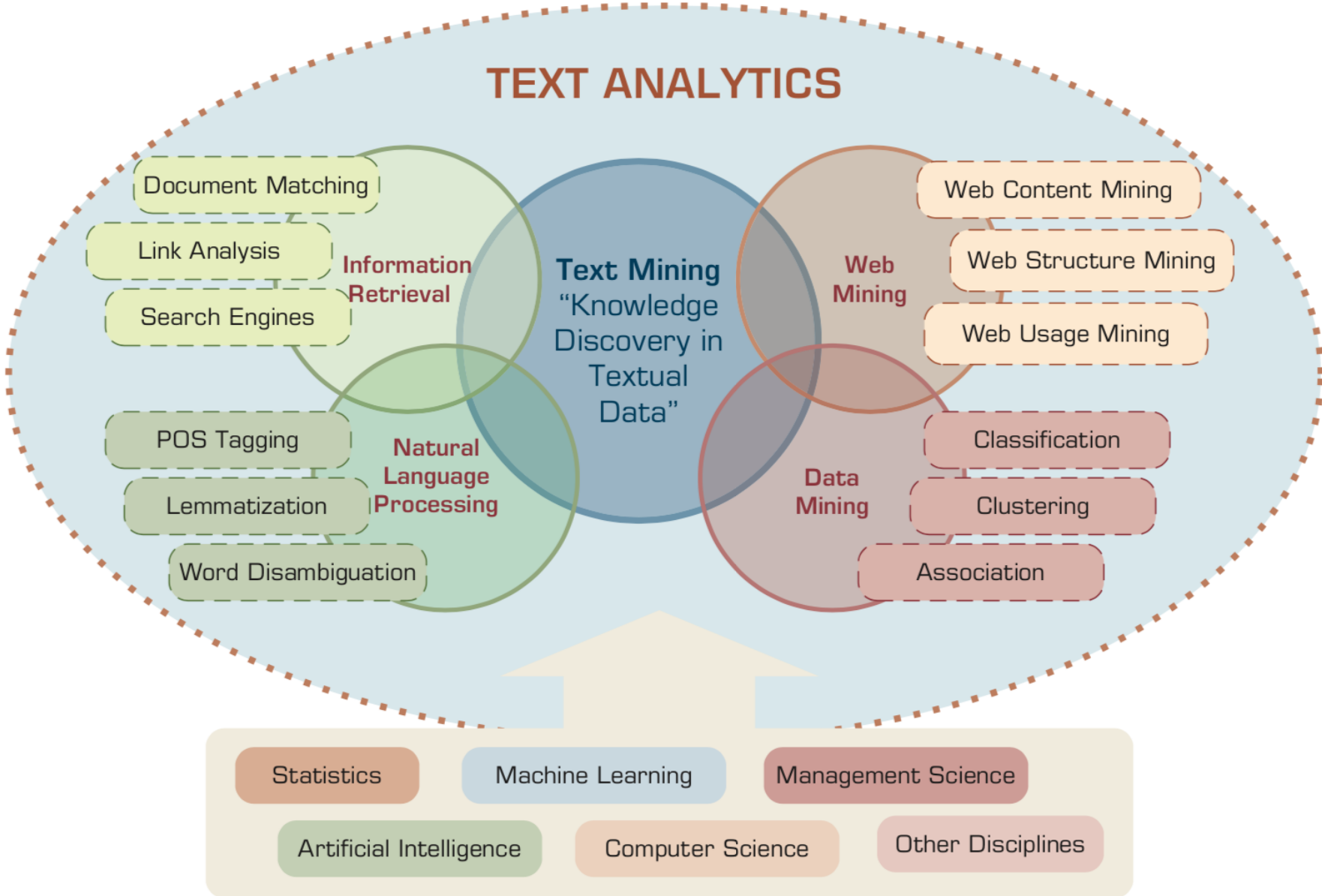
6. 問答系統與對話系統
   (Question Answering and Dialogue Systems)

# Outline

- Traditional Feature Engineering for Text Data
  - Bag of Words Model
  - Bag of N-Grams Model
  - TF-IDF Model
- Advanced Word Embeddings with Deep Learning
  - Word2Vec Model
  - Robust Word2Vec Models with Gensim
  - GloVe Model
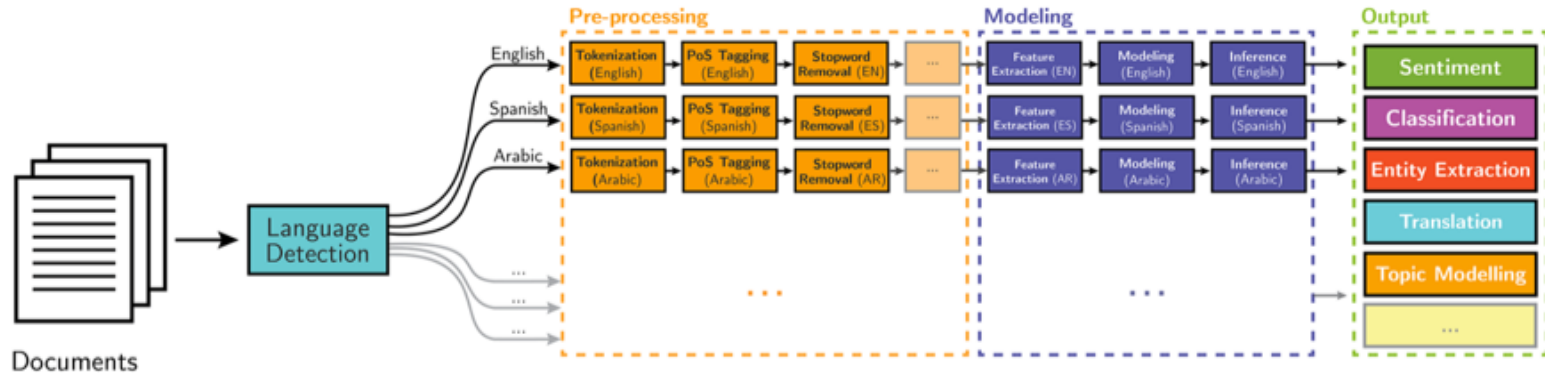  - FastText Model

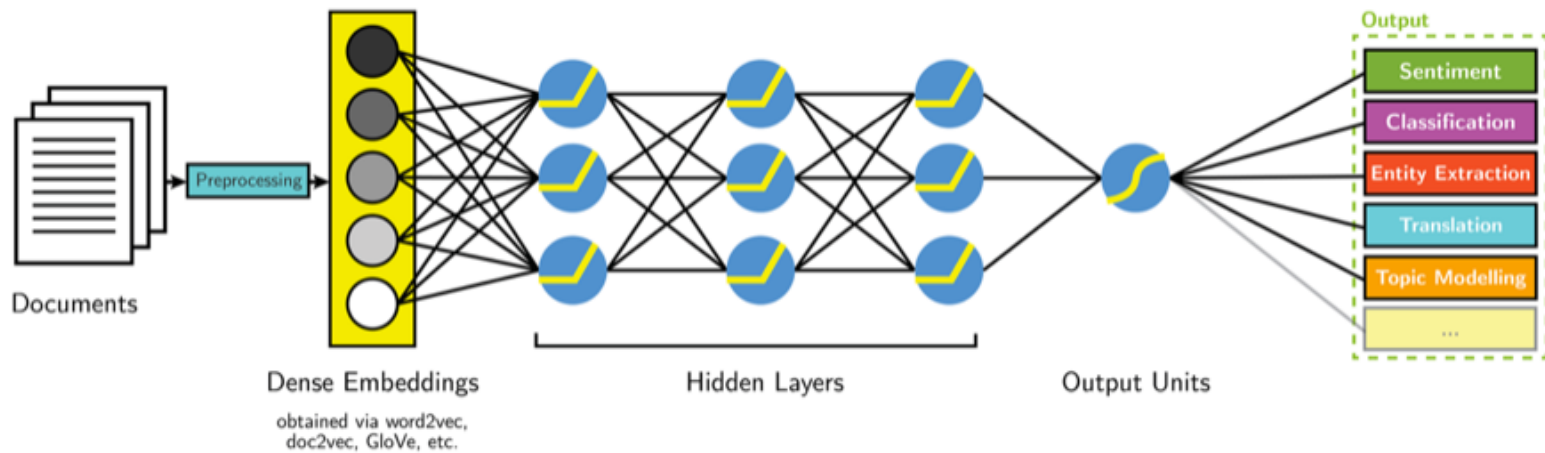# Feature Engineering for Text Representation

# Text Analytics and Text Mining



## TEXT ANALYTICS

**Information Retrieval**
- Document Matching
- Link Analysis
- Search Engines

**Natural Language Processing**
- POS Tagging
- Lemmatization
- Word Disambiguation

**Text Mining** "Knowledge Discovery in Textual Data"

**Web Mining**
- Web Content Mining
- Web Structure Mining
- Web Usage Mining

**Data Mining**
- Classification
- Clustering
- Association

- Statistics
- Machine Learning
- Management Science
- Artificial Intelligence
- Computer Science
- Other Disciplines
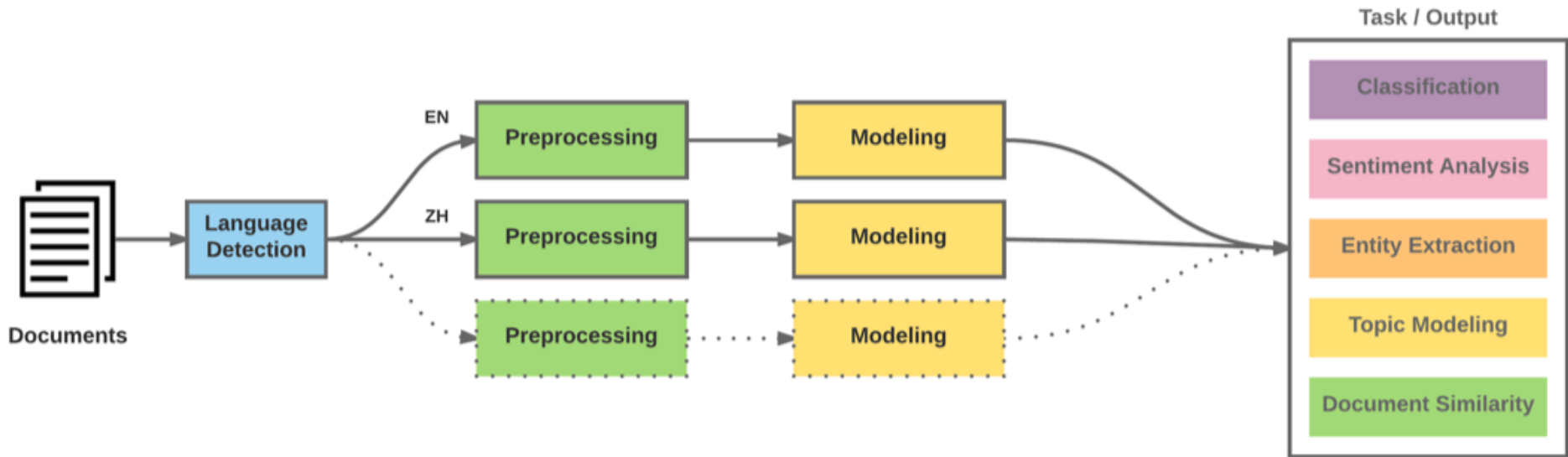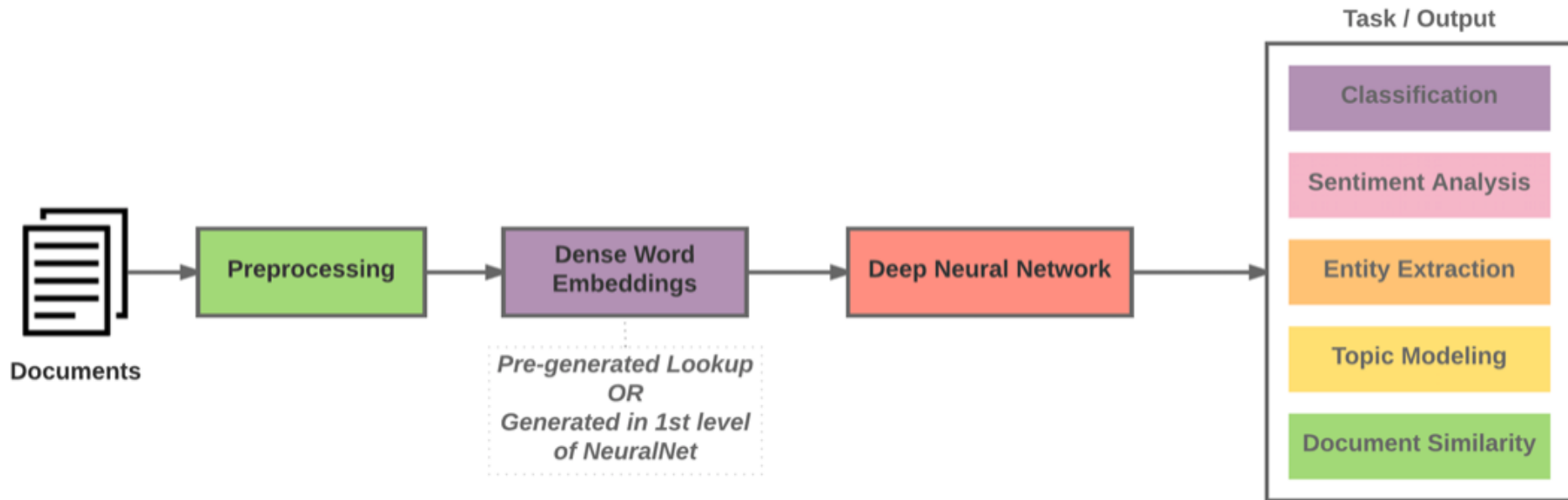
# NLP

# Modern NLP Pipeline

# Modern NLP Pipeline

# Deep Learning NLP

# Overview of
# Text Vectorization Methods

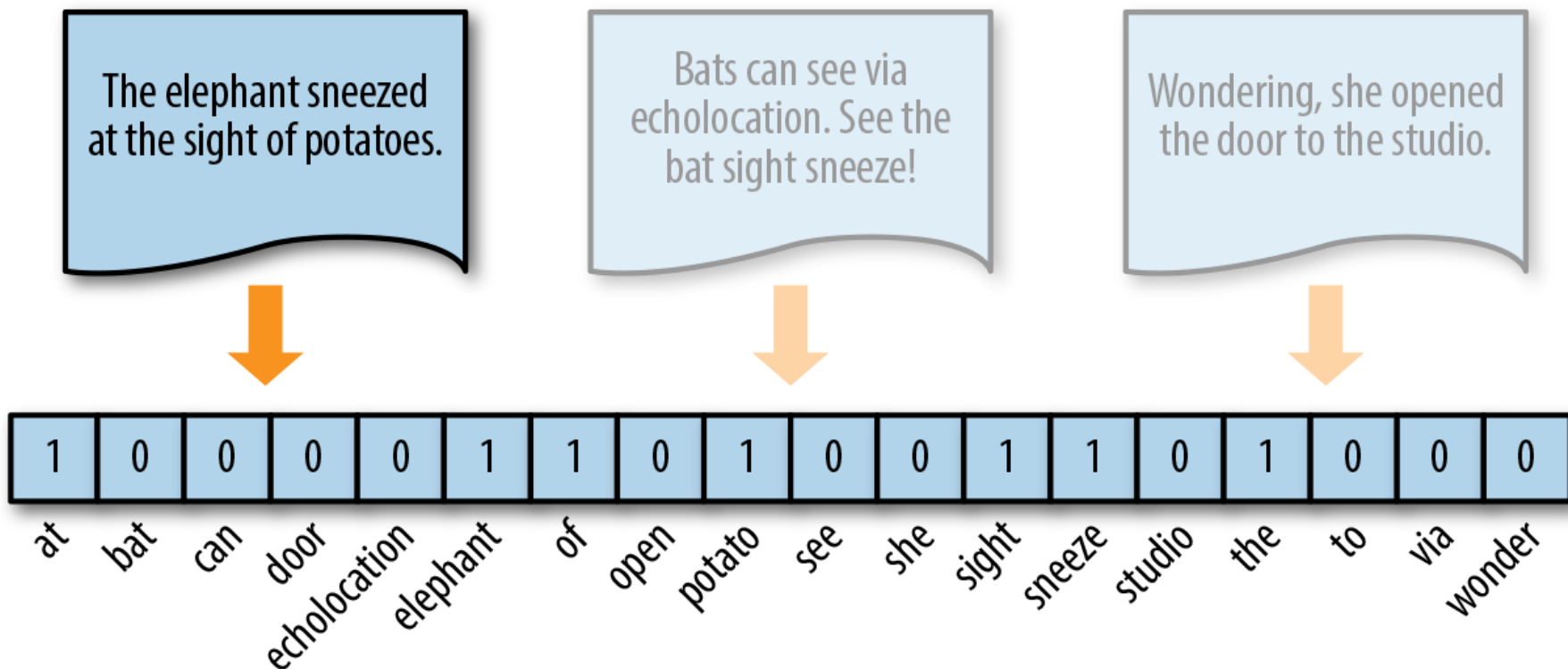| Vectorization Method | Function | Good For | Considerations |
|---|---|---|---|
| Frequency | Counts term frequencies | Bayesian models | Most frequent words not always most informative |
| One-Hot Encoding | Binarizes term occurrence (0, 1) | Neural networks | All words equidistant, so normalization extra important |
| TF–IDF | Normalizes term frequencies across documents | General purpose | Moderately frequent terms may not be representative of document topics |
| Distributed Representations | Context-based, continuous term similarity encoding | Modeling more complex relationships | Performance intensive; difficult to scale without additional tools (e.g., Tensorflow) |

# Encoding Documents as Vectors



Source: Benjamin Bengfort, Rebecca Bilbro, and Tony Ojeda (2018), Applied Text Analysis with Python, O'Reilly Media.
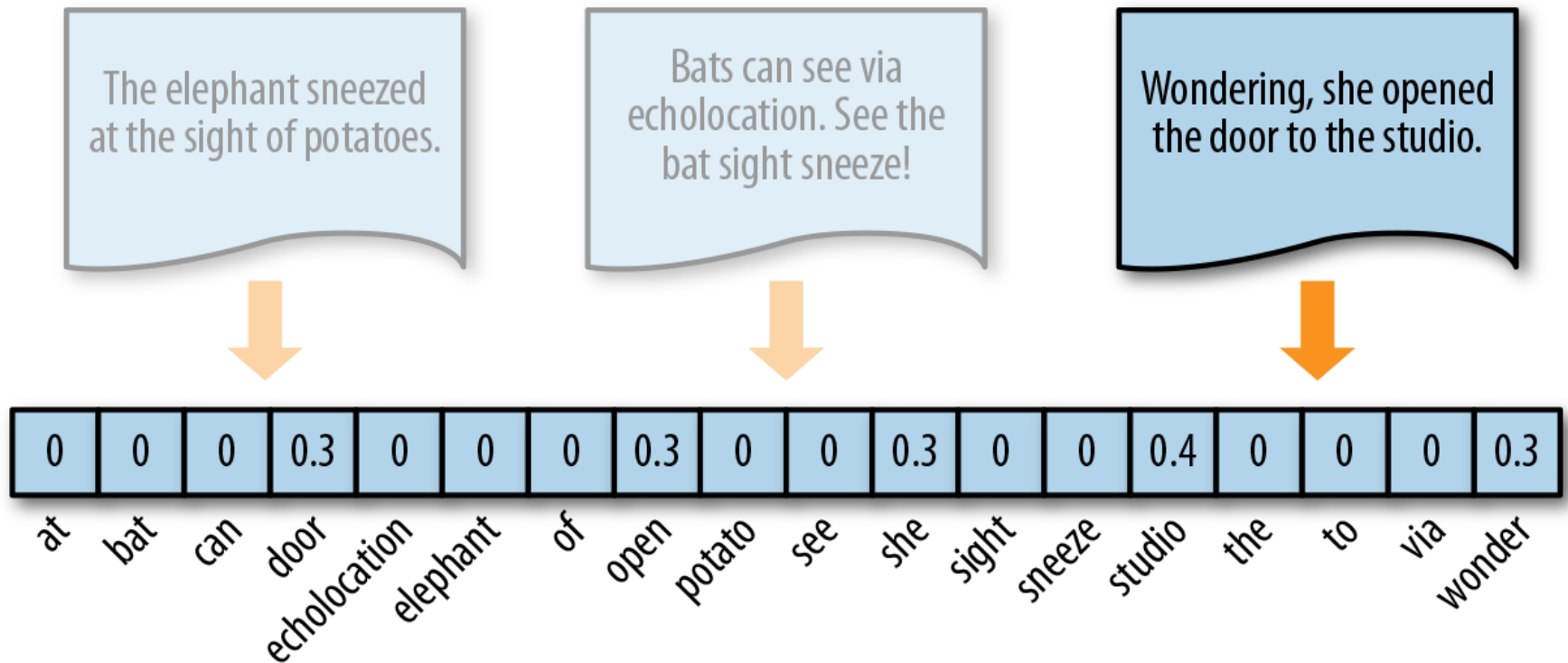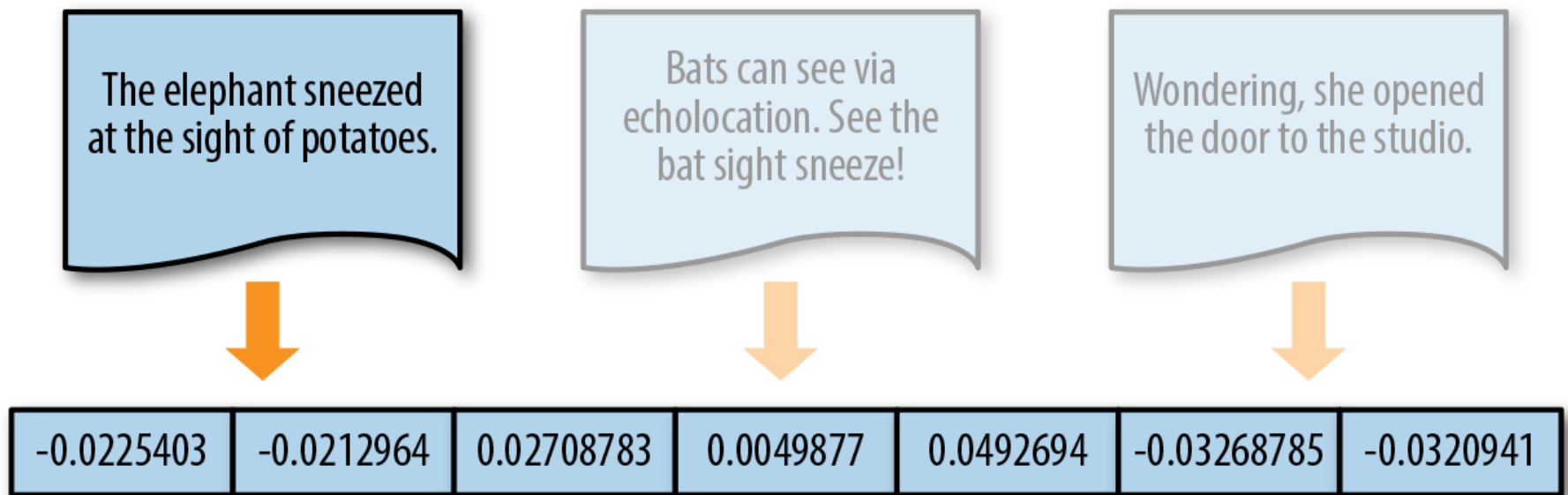
11

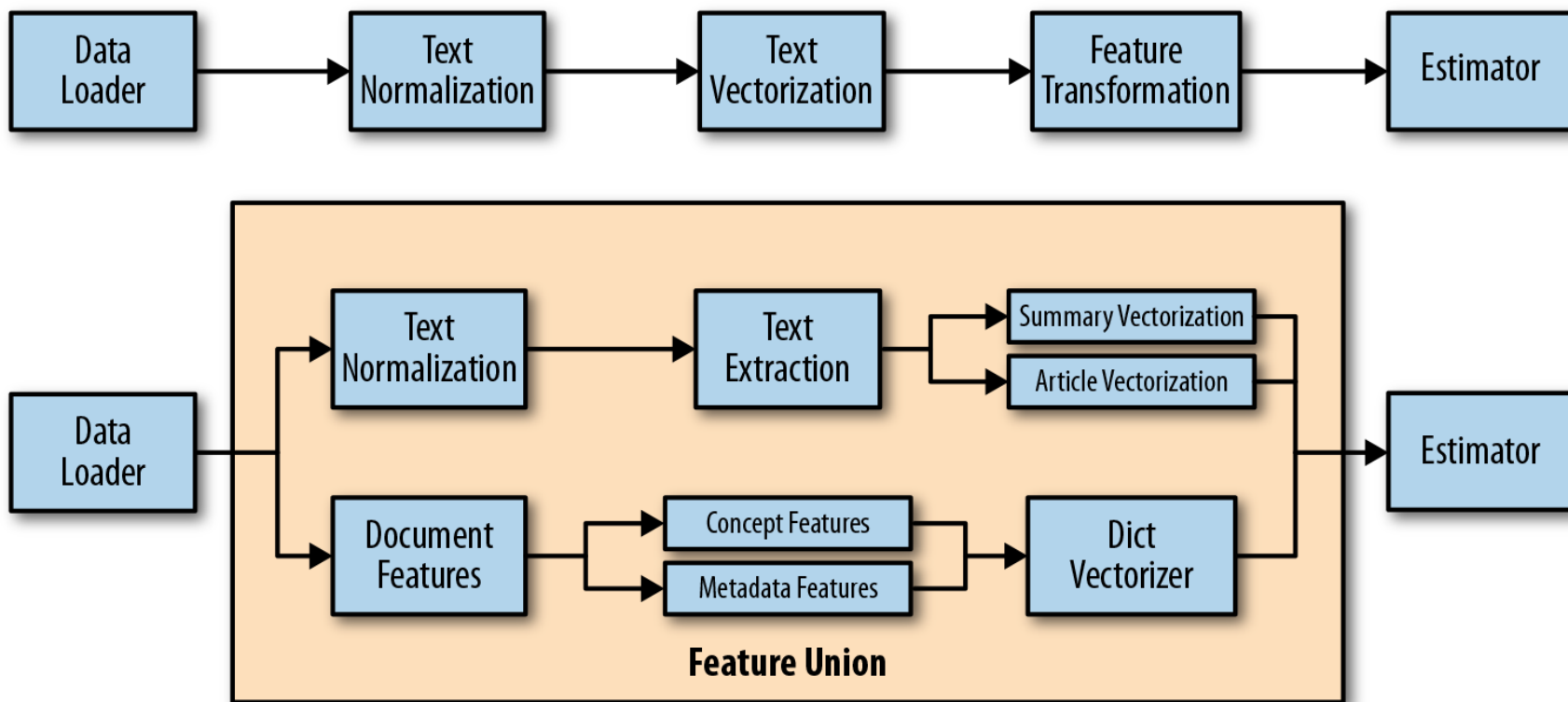# Token Frequency as Vector Encoding

# One-hot Encoding
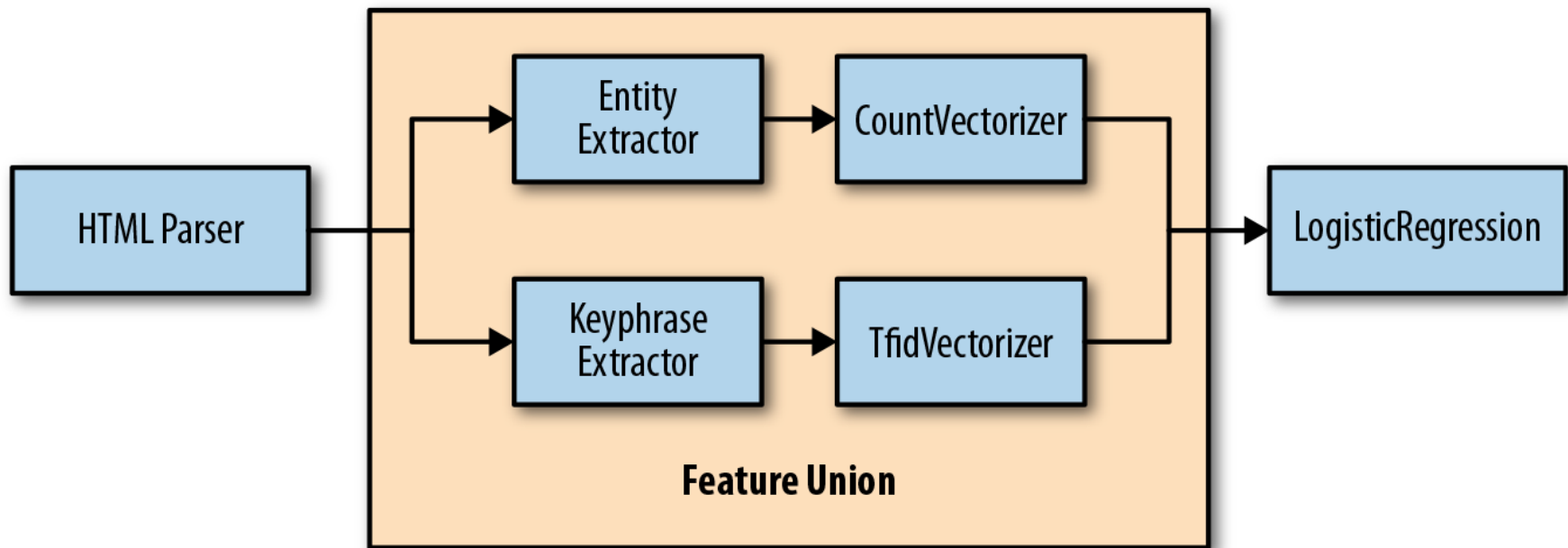
# TF-IDF Encoding

# Distributed Representation

# Pipelines for Text Vectorization and Feature Extraction

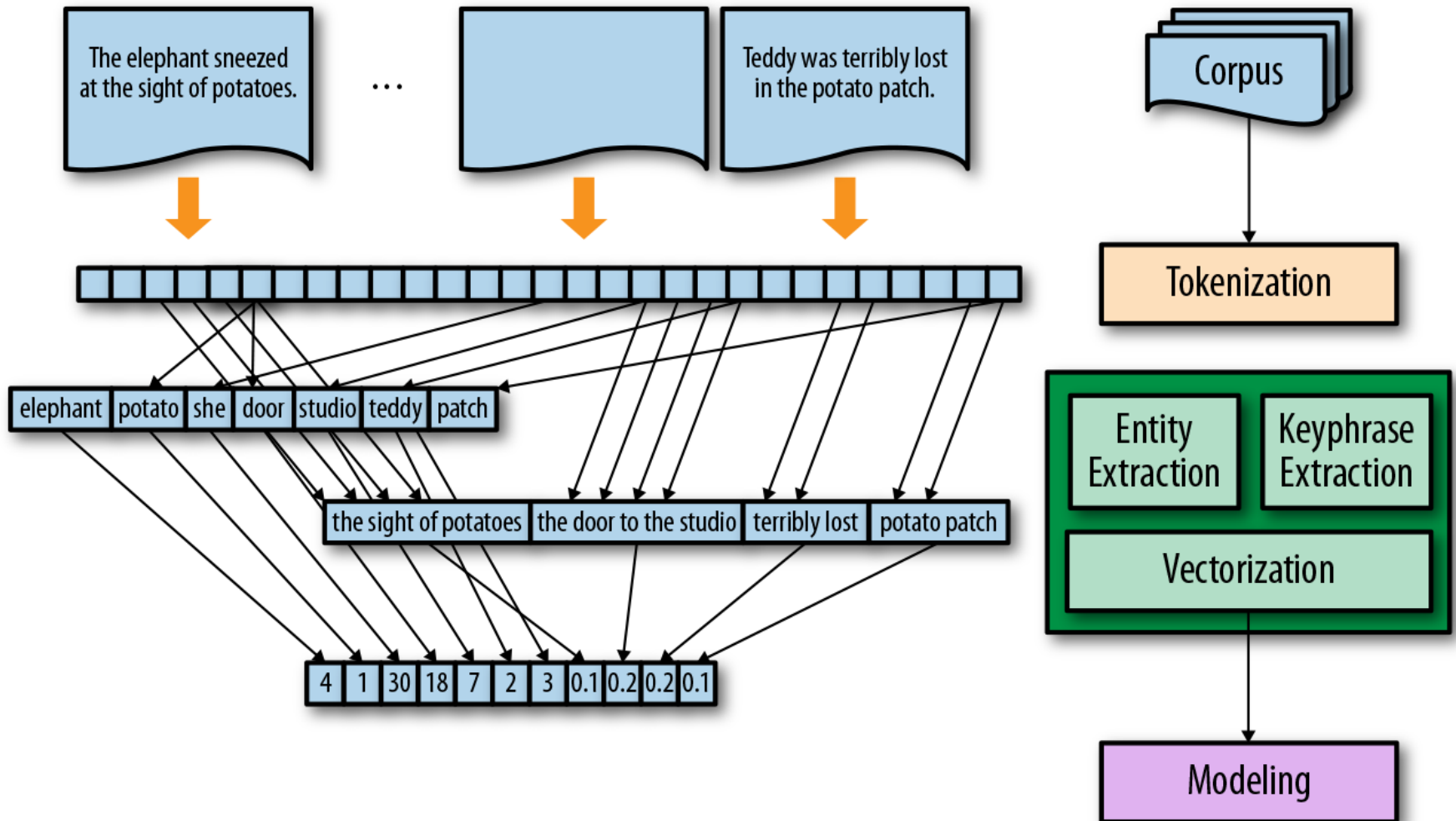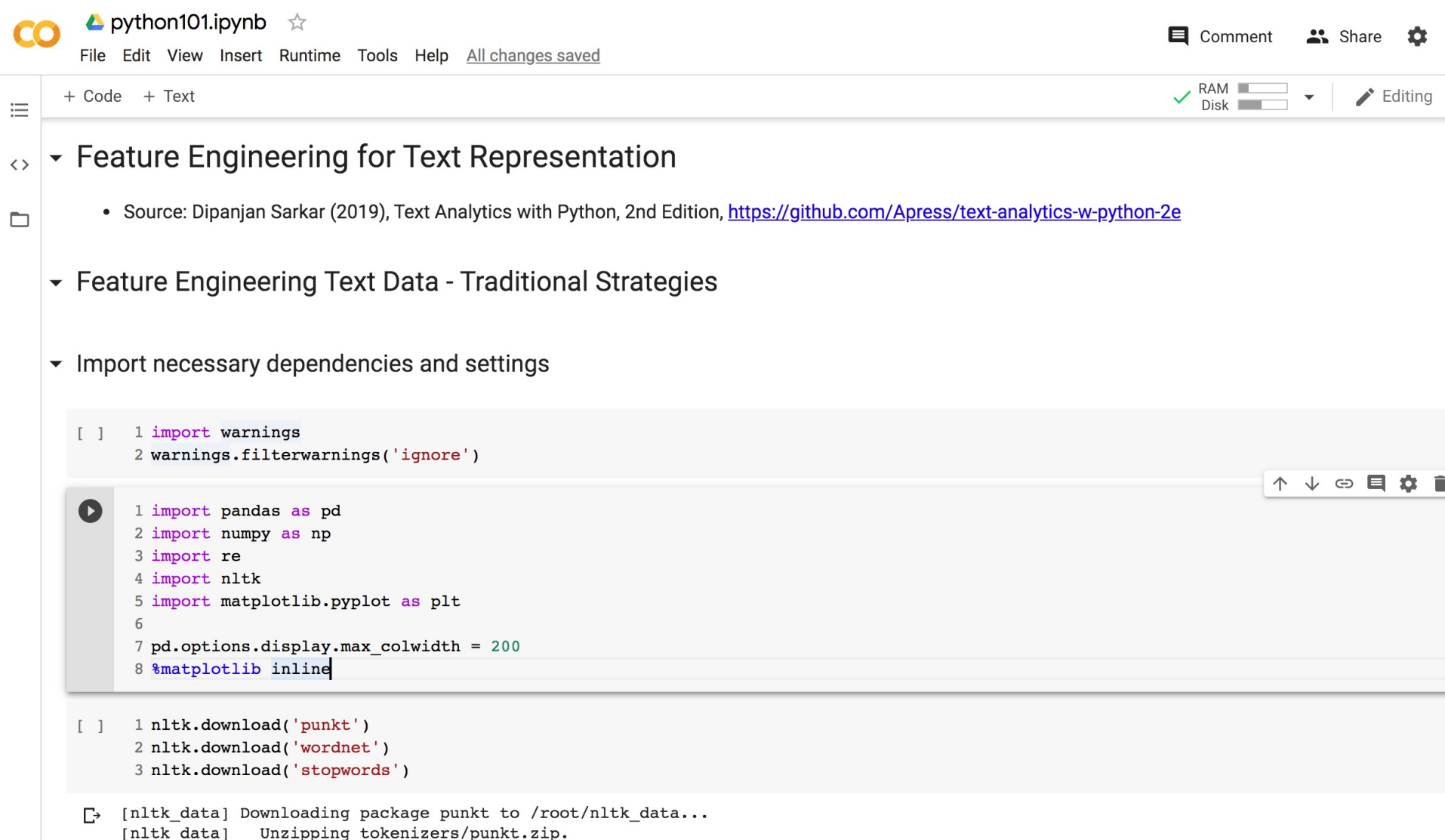# Feature Unions for Branching Vectorization

# Feature Extraction and Union

# Python in Google Colab (Python101)

https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT

python101.ipynb ☆

File   Edit   View   Insert   Runtime   Tools   Help   All changes saved

+ Code   + Text

RAM
Disk

✏ Editing

## Feature Engineering for Text Representation

- Source: Dipanjan Sarkar (2019), Text Analytics with Python, 2nd Edition, https://github.com/Apress/text-analytics-w-python-2e

## Feature Engineering Text Data - Traditional Strategies

## Import necessary dependencies and settings

```python
1 import warnings
2 warnings.filterwarnings('ignore')
```

```python
1 import pandas as pd
2 import numpy as np
3 import re
4 import nltk
5 import matplotlib.pyplot as plt
6
7 pd.options.display.max_colwidth = 200
8 %matplotlib inline
```

```python
1 nltk.download('punkt')
2 nltk.download('wordnet')
3 nltk.download('stopwords')
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt.zip.
```

https://tinyurl.com/imtkupython101

19

```python
corpus = ['The sky is blue and beautiful.',
'Love this blue and beautiful sky!',
'The quick brown fox jumps over the lazy dog.',
"A king's breakfast has sausages, ham, bacon, eggs, toast and
beans",
'I love green eggs, ham, sausages and bacon!',
'The brown fox is quick and the blue dog is lazy!',
'The sky is very blue and the sky is very beautiful today',
'The dog is lazy but the brown fox is quick!'
]
labels = ['weather', 'weather', 'animals', 'food', 'food',
'animals', 'weather', 'animals']

corpus = np.array(corpus)
corpus_df = pd.DataFrame({'Document': corpus,
'Category': labels})
corpus_df = corpus_df[['Document', 'Category']]
corpus_df
```

https://tinyurl.com/imtkupython101

```
corpus = np.array(corpus)
corpus_df = pd.DataFrame({'Document': corpus,
'Category': labels})
corpus_df = corpus_df[['Document', 'Category']]
corpus_df
```

| | Document | Category |
|---|---|---|
| 0 | The sky is blue and beautiful. | weather |
| 1 | Love this blue and beautiful sky! | weather |
| 2 | The quick brown fox jumps over the lazy dog. | animals |
| 3 | A king's breakfast has sausages, ham, bacon, eggs, toast and beans | food |
| 4 | I love green eggs, ham, sausages and bacon! | food |
| 5 | The brown fox is quick and the blue dog is lazy! | animals |
| 6 | The sky is very blue and the sky is very beautiful today | weather |
| 7 | The dog is lazy but the brown fox is quick! | animals |

https://tinyurl.com/imtkupython101

```python
wpt = nltk.WordPunctTokenizer()
stop_words = nltk.corpus.stopwords.words('english')

def normalize_document(doc):
# lower case and remove special characters\whitespaces
doc = re.sub(r'[^a-zA-Z\s]', '', doc, re.I|re.A)
doc = doc.lower()
doc = doc.strip()
# tokenize document
tokens = wpt.tokenize(doc)
# filter stopwords out of document
filtered_tokens = [token for token in tokens if token not in stop_words]
# re-create document from filtered tokens
doc = ' '.join(filtered_tokens)
return doc

normalize_corpus = np.vectorize(normalize_document)
norm_corpus = normalize_corpus(corpus)
norm_corpus
```

https://tinyurl.com/imtkupython101

```python
from sklearn.feature_extraction.text import CountVectorizer
# get bag of words features in sparse format
cv = CountVectorizer(min_df=0., max_df=1.)
cv_matrix = cv.fit_transform(norm_corpus)
cv_matrix
```

```python
# view non-zero feature positions in the sparse matrix
print(cv_matrix)
```

```python
# view dense representation
# warning might give a memory error if data is too big
cv_matrix = cv_matrix.toarray()
cv_matrix
```

```
array([
[0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0],
[0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0],
[0, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0],
[1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 1, 0],
[1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0],
[0, 0, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0],
[0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 2, 0, 1],
[0, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0]])
```

https://tinyurl.com/imtkupython101

```
# get all unique words in the corpus
vocab = cv.get_feature_names()
# show document feature vectors
pd.DataFrame(cv_matrix, columns=vocab)
```

```
1 # get all unique words in the corpus
2 vocab = cv.get_feature_names()
3 # show document feature vectors
4 pd.DataFrame(cv_matrix, columns=vocab)
```

| | bacon | beans | beautiful | blue | breakfast | brown | dog | eggs | fox | green | ham | jumps | kings | lazy | love | quick | sausages | sky | toast | today |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| 3 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 |
| 4 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 1 |
| 7 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |

https://tinyurl.com/imtkupython101

```python
# you can set the n-gram range to 1,2 to get unigrams as well as bigrams
bv = CountVectorizer(ngram_range=(2,2))
bv_matrix = bv.fit_transform(norm_corpus)

bv_matrix = bv_matrix.toarray()
vocab = bv.get_feature_names()
pd.DataFrame(bv_matrix, columns=vocab)
```

| | bacon eggs | beautiful sky | beautiful today | blue beautiful | blue dog | blue sky | breakfast sausages | brown fox | dog lazy | eggs ham | eggs toast | fox jumps | fox quick | green eggs | ham bacon | ham sausages | jumps lazy | kings breakfast |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 |
| 3 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| 5 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |

```python
1 from sklearn.feature_extraction.text import TfidfTransformer
2
3 tt = TfidfTransformer(norm='l2', use_idf=True, smooth_idf=True)
4 tt_matrix = tt.fit_transform(cv_matrix)
5
6 tt_matrix = tt_matrix.toarray()
7 vocab = cv.get_feature_names()
8 pd.DataFrame(np.round(tt_matrix, 2), columns=vocab)
```

| | bacon | beans | beautiful | blue | breakfast | brown | dog | eggs | fox | green | ham | jumps | kings | lazy | love | quick | sausages | sky | toast | today |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.00 | 0.00 | 0.60 | 0.53 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.60 | 0.00 | 0.0 |
| 1 | 0.00 | 0.00 | 0.49 | 0.43 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.57 | 0.00 | 0.00 | 0.49 | 0.00 | 0.0 |
| 2 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.38 | 0.38 | 0.00 | 0.38 | 0.00 | 0.00 | 0.53 | 0.00 | 0.38 | 0.00 | 0.38 | 0.00 | 0.00 | 0.00 | 0.0 |
| 3 | 0.32 | 0.38 | 0.00 | 0.00 | 0.38 | 0.00 | 0.00 | 0.32 | 0.00 | 0.00 | 0.32 | 0.00 | 0.38 | 0.00 | 0.00 | 0.00 | 0.32 | 0.00 | 0.38 | 0.0 |
| 4 | 0.39 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.39 | 0.00 | 0.47 | 0.39 | 0.00 | 0.00 | 0.00 | 0.39 | 0.00 | 0.39 | 0.00 | 0.00 | 0.0 |
| 5 | 0.00 | 0.00 | 0.00 | 0.37 | 0.00 | 0.42 | 0.42 | 0.00 | 0.42 | 0.00 | 0.00 | 0.00 | 0.00 | 0.42 | 0.00 | 0.42 | 0.00 | 0.00 | 0.00 | 0.0 |
| 6 | 0.00 | 0.00 | 0.36 | 0.32 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.72 | 0.00 | 0.5 |
| 7 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.45 | 0.45 | 0.00 | 0.45 | 0.00 | 0.00 | 0.00 | 0.00 | 0.45 | 0.00 | 0.45 | 0.00 | 0.00 | 0.00 | 0.0 |

https://tinyurl.com/imtkupython101

```
1  from sklearn.feature_extraction.text import TfidfVectorizer
2
3  tv = TfidfVectorizer(min_df=0., max_df=1., norm='l2',
4                       use_idf=True, smooth_idf=True)
5  tv_matrix = tv.fit_transform(norm_corpus)
6  tv_matrix = tv_matrix.toarray()
7
8  vocab = tv.get_feature_names()
9  pd.DataFrame(np.round(tv_matrix, 2), columns=vocab)
```

| | bacon | beans | beautiful | blue | breakfast | brown | dog | eggs | fox | green | ham | jumps | kings | lazy | love | quick | sausages | sky | toast | today |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.00 | 0.00 | 0.60 | 0.53 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.60 | 0.00 | 0.0 |
| 1 | 0.00 | 0.00 | 0.49 | 0.43 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.57 | 0.00 | 0.00 | 0.49 | 0.00 | 0.0 |
| 2 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.38 | 0.38 | 0.00 | 0.38 | 0.00 | 0.00 | 0.53 | 0.00 | 0.38 | 0.00 | 0.38 | 0.00 | 0.00 | 0.00 | 0.0 |
| 3 | 0.32 | 0.38 | 0.00 | 0.00 | 0.38 | 0.00 | 0.00 | 0.32 | 0.00 | 0.00 | 0.32 | 0.00 | 0.38 | 0.00 | 0.00 | 0.00 | 0.32 | 0.00 | 0.38 | 0.0 |
| 4 | 0.39 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.39 | 0.00 | 0.47 | 0.39 | 0.00 | 0.00 | 0.00 | 0.39 | 0.00 | 0.39 | 0.00 | 0.00 | 0.0 |
| 5 | 0.00 | 0.00 | 0.00 | 0.37 | 0.00 | 0.42 | 0.42 | 0.00 | 0.42 | 0.00 | 0.00 | 0.00 | 0.00 | 0.42 | 0.00 | 0.42 | 0.00 | 0.00 | 0.00 | 0.0 |
| 6 | 0.00 | 0.00 | 0.36 | 0.32 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.72 | 0.00 | 0.5 |
| 7 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.45 | 0.45 | 0.00 | 0.45 | 0.00 | 0.00 | 0.00 | 0.00 | 0.45 | 0.00 | 0.45 | 0.00 | 0.00 | 0.00 | 0.0 |

https://tinyurl.com/imtkupython101

```
1 from scipy.cluster.hierarchy import dendrogram, linkage
2
3 Z = linkage(similarity_matrix, 'ward')
4 pd.DataFrame(Z, columns=['Document Cluster 1', 'Document Cluster 2',
5                          'Distance', 'Cluster Size'], dtype='object')
```

| | Document Cluster 1 | Document Cluster 2 | Distance | Cluster Size |
|---|---|---|---|---|
| **0** | 2 | 7 | 0.253098 | 2 |
| **1** | 0 | 6 | 0.308539 | 2 |
| **2** | 5 | 8 | 0.386952 | 3 |
| **3** | 1 | 9 | 0.489845 | 3 |
| **4** | 3 | 4 | 0.732945 | 2 |
| **5** | 11 | 12 | 2.69565 | 5 |
| **6** | 10 | 13 | 3.45108 | 8 |

https://tinyurl.com/imtkupython101

```
1 plt.figure(figsize=(8, 3))
2 plt.title('Hierarchical Clustering Dendrogram')
3 plt.xlabel('Data point')
4 plt.ylabel('Distance')
5 dendrogram(Z)
6 plt.axhline(y=1.0, c='k', ls='--', lw=0.5)
```

`<matplotlib.lines.Line2D at 0x7ff7b5d793c8>`



Hierarchical Clustering Dendrogram

```
1 from scipy.cluster.hierarchy import fcluster
2 max_dist = 1.0
3
4 cluster_labels = fcluster(Z, max_dist, criterion='distance')
5 cluster_labels = pd.DataFrame(cluster_labels, columns=['ClusterLabel'])
6 pd.concat([corpus_df, cluster_labels], axis=1)
```

|   | Document | Category | ClusterLabel |
|---|---|---|---|
| 0 | The sky is blue and beautiful. | weather | 2 |
| 1 | Love this blue and beautiful sky! | weather | 2 |
| 2 | The quick brown fox jumps over the lazy dog. | animals | 1 |
| 3 | A king's breakfast has sausages, ham, bacon, eggs, toast and beans | food | 3 |
| 4 | I love green eggs, ham, sausages and bacon! | food | 3 |
| 5 | The brown fox is quick and the blue dog is lazy! | animals | 1 |
| 6 | The sky is very blue and the sky is very beautiful today | weather | 2 |
| 7 | The dog is lazy but the brown fox is quick! | animals | 1 |

https://tinyurl.com/imtkupython101

```python
1 from sklearn.decomposition import LatentDirichletAllocation
2 lda = LatentDirichletAllocation(n_components=3, max_iter=10000, random_state=0)
3 #lda = LatentDirichletAllocation(n_topics=3, max_iter=10000, random_state=0)
4 dt_matrix = lda.fit_transform(cv_matrix)
5 features = pd.DataFrame(dt_matrix, columns=['T1', 'T2', 'T3'])
6 features
```

|   | T1 | T2 | T3 |
|---|----|----|----|
| 0 | 0.832191 | 0.083480 | 0.084329 |
| 1 | 0.863554 | 0.069100 | 0.067346 |
| 2 | 0.047794 | 0.047776 | 0.904430 |
| 3 | 0.037243 | 0.925559 | 0.037198 |
| 4 | 0.049121 | 0.903076 | 0.047802 |
| 5 | 0.054902 | 0.047778 | 0.897321 |
| 6 | 0.888287 | 0.055697 | 0.056016 |
| 7 | 0.055704 | 0.055689 | 0.888607 |

https://tinyurl.com/imtkupython101

```
1 tt_matrix = lda.components_
2 for topic_weights in tt_matrix:
3     topic = [(token, weight) for token, weight in zip(vocab, topic_weights)]
4     topic = sorted(topic, key=lambda x: -x[1])
5     topic = [item for item in topic if item[1] > 0.6]
6     print(topic)
7     print()
```

[('sky', 4.332439442470133), ('blue', 3.373774254787669), ('beautiful', 3.3323650509884386), ('today', 1.3325579855138987), ('love

[('bacon', 2.33269586574902), ('eggs', 2.33269586574902), ('ham', 2.33269586574902), ('sausages', 2.33269586574902), ('love', 1.33

[('brown', 3.3323473548404405), ('dog', 3.3323473548404405), ('fox', 3.3323473548404405), ('lazy', 3.3323473548404405), ('quick',

```
1 from sklearn.cluster import KMeans
2
3 km = KMeans(n_clusters=3, random_state=0)
4 km.fit_transform(features)
5 cluster_labels = km.labels_
6 cluster_labels = pd.DataFrame(cluster_labels, columns=['ClusterLabel'])
7 pd.concat([corpus_df, cluster_labels], axis=1)
```

| | Document | Category | ClusterLabel |
|---|---|---|---|
| **0** | The sky is blue and beautiful. | weather | 1 |
| **1** | Love this blue and beautiful sky! | weather | 1 |
| **2** | The quick brown fox jumps over the lazy dog. | animals | 2 |
| **3** | A king's breakfast has sausages, ham, bacon, eggs, toast and beans | food | 0 |
| **4** | I love green eggs, ham, sausages and bacon! | food | 0 |
| **5** | The brown fox is quick and the blue dog is lazy! | animals | 2 |
| **6** | The sky is very blue and the sky is very beautiful today | weather | 1 |
| **7** | The dog is lazy but the brown fox is quick! | animals | 2 |

https://tinyurl.com/imtkupython101

```python
from gensim.models import word2vec

# tokenize sentences in corpus
wpt = nltk.WordPunctTokenizer()
tokenized_corpus = [wpt.tokenize(document) for document in norm_bible]

# Set values for various parameters
feature_size = 100 # Word vector dimensionality
window_context = 30 # Context window size
min_word_count = 1 # Minimum word count
sample = 1e-3 # Downsample setting for frequent words

w2v_model = word2vec.Word2Vec(tokenized_corpus, size=feature_size,
window=window_context, min_count=min_word_count,
sample=sample, iter=50)

# view similar words based on gensim's model
similar_words = {search_term: [item[0] for item in
w2v_model.wv.most_similar([search_term], topn=5)]
for search_term in ['god', 'jesus', 'noah', 'egypt', 'john', 'gospel',
'moses','famine']}
similar_words
```

https://tinyurl.com/imtkupython101

```
w2v_model.wv.most_similar([search_term], topn=5)]
```

```
{'egypt': ['egyptians', 'pharaoh', 'bondage', 'flowing',
'rod'], 'famine': ['pestilence', 'peril', 'deaths',
'morever', 'sword'], 'god': ['lord', 'worldly', 'soberly',
'reasonable', 'unto'], 'gospel': ['christ', 'faith',
'repentance', 'sufferings', 'afflictions'], 'jesus':
['peter', 'messias', 'immediately', 'apostles',
'synagogue'], 'john': ['james', 'baptist', 'devine',
'peter', 'simon'], 'moses': ['congregation', 'elisheba',
'naashon', 'joshua', 'children'], 'noah': ['shem',
'japheth', 'ham', 'noe', 'hoglah']}
```

https://tinyurl.com/imtkupython101

```python
from sklearn.decomposition import PCA

pca = PCA(n_components=2, random_state=0)
pcs = pca.fit_transform(w2v_feature_array)
labels = ap.labels_
categories = list(corpus_df['Category'])
plt.figure(figsize=(8, 6))

for i in range(len(labels)):
label = labels[i]
color = 'orange' if label == 0 else 'blue' if label == 1
else 'green'
annotation_label = categories[i]
x, y = pcs[i]
plt.scatter(x, y, c=color, edgecolors='k')
plt.annotate(annotation_label, xy=(x+1e-4, y+1e-3),
xytext=(0, 0), textcoords='offset points')
```

# BERT:
# Pre-training of Deep Bidirectional Transformers for Language Understanding

**BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding**

**Jacob Devlin**   **Ming-Wei Chang**   **Kenton Lee**   **Kristina Toutanova**

Google AI Language

{jacobdevlin,mingweichang,kentonl,kristout}@google.com

# BERT

## Bidirectional Encoder Representations from Transformers



## Pre-training model architectures

**BERT** uses a bidirectional Transformer.
**OpenAI GPT** uses a left-to-right Transformer.
**ELMo** uses the concatenation of independently trained left-to-right and right- to-left LSTM to generate features for downstream tasks.
Among three, only BERT representations are jointly conditioned on both left and right context in all layers.

# BERT input representation

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

The input embeddings is the sum of the token embeddings, the segmentation embeddings and the position embeddings.

# BERT Sequence-level tasks



(a) Sentence Pair Classification Tasks: MNLI, QQP, QNLI, STS-B, MRPC, RTE, SWAG

(b) Single Sentence Classification Tasks: SST-2, CoLA

43

# BERT Token-level tasks



(c) Question Answering Tasks: SQuAD v1.1

(d) Single Sentence Tagging Tasks: CoNLL-2003 NER

44

# General Language Understanding Evaluation (GLUE) benchmark

# GLUE Test results

| System | MNLI-(m/mm) | QQP | QNLI | SST-2 | CoLA | STS-B | MRPC | RTE | Average |
|---|---|---|---|---|---|---|---|---|---|
| | 392k | 363k | 108k | 67k | 8.5k | 5.7k | 3.5k | 2.5k | - |
| Pre-OpenAI SOTA | 80.6/80.1 | 66.1 | 82.3 | 93.2 | 35.0 | 81.0 | 86.0 | 61.7 | 74.0 |
| BiLSTM+ELMo+Attn | 76.4/76.1 | 64.8 | 79.9 | 90.4 | 36.0 | 73.3 | 84.9 | 56.8 | 71.0 |
| OpenAI GPT | 82.1/81.4 | 70.3 | 88.1 | 91.3 | 45.4 | 80.0 | 82.3 | 56.0 | 75.2 |
| $BERT_{BASE}$ | 84.6/83.4 | 71.2 | 90.1 | 93.5 | 52.1 | 85.8 | 88.9 | 66.4 | 79.6 |
| $BERT_{LARGE}$ | **86.7/85.9** | **72.1** | **91.1** | **94.9** | **60.5** | **86.5** | **89.3** | **70.1** | **81.9** |

**MNLI**: Multi-Genre Natural Language Inference
**QQP**: Quora Question Pairs
**QNLI**: Question Natural Language Inference
**SST-2**: The Stanford Sentiment Treebank
**CoLA**: The Corpus of Linguistic Acceptability
**STS-B**:The Semantic Textual Similarity Benchmark
**MRPC**: Microsoft Research Paraphrase Corpus
**RTE**: Recognizing Textual Entailment

# Facebook Research FastText

Pre-trained word vectors
Word2Vec
wiki.zh.vec (861MB)
332647 word
300 vec

Pre-trained word vectors for 90 languages,
trained on Wikipedia using fastText.

These vectors in dimension 300 were obtained using
the skip-gram model with default parameters.

https://github.com/facebookresearch/fastText/blob/master/pretrained-vectors.md

# Facebook Research FastText
# Word2Vec: wiki.zh.vec
## (861MB) (332647 word 300 vec)

https://github.com/facebookresearch/fastText/blob/master/pretrained-vectors.md

# Word Embeddings in LSTM RNN



Time Expanded LSTM Network

LSTM Internal States

Word Embeddings

Input Question: Is this person dancing ?

Fixed length question vector encoded by the LSTM

# Transformer (Attention is All You Need)
## (Vaswani et al., 2017)

# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

## BERT (Bidirectional Encoder Representations from Transformers)

## Overall pre-training and fine-tuning procedures for BERT

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).
"Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

## BERT (Bidirectional Encoder Representations from Transformers)

### BERT input representation

# BERT, OpenAI GPT, ELMo

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).
"Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

# Fine-tuning BERT on Different Tasks



(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC, RTE, SWAG

(b) Single Sentence Classification Tasks:
SST-2, CoLA

(c) Question Answering Tasks:
SQuAD v1.1

(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).
"Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

53

# Pre-trained Language Model (PLM)



**Semi-supervised Sequence Learning**
**context2Vec**
**Pre-trained seq2seq**

ULMFiT — ELMo — Transformer — BERT — Bidirectional LM — GPT

Multi-lingual → MultiFiT

Cross-lingual

Multi-task

+ Generation

XLM UDify

MT-DNN

Knowledge distillation

MT-DNN$_{KD}$

MASS UniLM

Span prediction
Remove NSP

Longer time
Remove NSP
More data

SpanBERT

RoBERTa

Permutation LM
Transformer-XL
More data

XLNet

+Knowledge Graph

ERNIE (Tsinghua)

Neural entity linker

KnowBert

Cross-modal

VideoBERT
CBT
ViLBERT
VisualBERT
B2T2
Unicoder-VL
LXMERT
VL-BERT
UNITER

Whole Word Masking

ERNIE (Baidu)
BERT-wwm

Larger model
More data → GPT-2 — Defense → Grover

By Xiaozhi Wang & Zhengyan Zhang @THUNLP

# Turing Natural Language Generation (T-NLG)



T-NLG
17b

MegatronLM
8.3b

GPT-2
1.5b

BERT-Large
340m

RoBERTa
355m

DistilBERT
66m

2018                    2019                    2020

# Transformers

## State-of-the-art Natural Language Processing for TensorFlow 2.0 and PyTorch

- Transformers
  - pytorch-transformers
  - pytorch-pretrained-bert
- provides state-of-the-art general-purpose architectures
  - (BERT, GPT-2, RoBERTa, XLM, DistilBert, XLNet, CTRL...)
  - for Natural Language Understanding (NLU) and
    Natural Language Generation (NLG)
    with over 32+ pretrained models
    in 100+ languages
    and deep interoperability between
    TensorFlow 2.0 and
    PyTorch.

# Transfer Learning
# in Natural Language Processing

Source: Sebastian Ruder, Matthew E. Peters, Swabha Swayamdipta, and Thomas Wolf (2019), "Transfer learning in natural language processing." In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Tutorials, pp. 15-18.

# NLP Benchmark Datasets

| Task | Dataset | Link |
|---|---|---|
| Machine Translation | WMT 2014 EN-DE<br>WMT 2014 EN-FR | http://www-lium.univ-lemans.fr/~schwenk/cslm_joint_paper/ |
| Text Summarization | CNN/DM<br>Newsroom<br>DUC<br>Gigaword | https://cs.nyu.edu/~kcho/DMQA/<br>https://summari.es/<br>https://www-nlpir.nist.gov/projects/duc/data.html<br>https://catalog.ldc.upenn.edu/LDC2012T21 |
| Reading Comprehension<br>Question Answering<br>Question Generation | ARC<br>CliCR<br>CNN/DM<br>NewsQA<br>RACE<br>SQuAD<br>Story Cloze Test<br>NarativeQA<br>Quasar<br>SearchQA | http://data.allenai.org/arc/<br>http://aclweb.org/anthology/N18-1140<br>https://cs.nyu.edu/~kcho/DMQA/<br>https://datasets.maluuba.com/NewsQA<br>http://www.qizhexie.com/data/RACE_leaderboard<br>https://rajpurkar.github.io/SQuAD-explorer/<br>http://aclweb.org/anthology/W17-0906.pdf<br>https://github.com/deepmind/narrativeqa<br>https://github.com/bdhingra/quasar<br>https://github.com/nyu-dl/SearchQA |
| Semantic Parsing | AMR parsing<br>ATIS (SQL Parsing)<br>WikiSQL (SQL Parsing) | https://amr.isi.edu/index.html<br>https://github.com/jkkummerfeld/text2sql-data/tree/master/data<br>https://github.com/salesforce/WikiSQL |
| Sentiment Analysis | IMDB Reviews<br>SST<br>Yelp Reviews<br>Subjectivity Dataset | http://ai.stanford.edu/~amaas/data/sentiment/<br>https://nlp.stanford.edu/sentiment/index.html<br>https://www.yelp.com/dataset/challenge<br>http://www.cs.cornell.edu/people/pabo/movie-review-data/ |
| Text Classification | AG News<br>DBpedia<br>TREC<br>20 NewsGroup | http://www.di.unipi.it/~gulli/AG_corpus_of_news_articles.html<br>https://wiki.dbpedia.org/Datasets<br>https://trec.nist.gov/data.html<br>http://qwone.com/~jason/20Newsgroups/ |
| Natural Language Inference | SNLI Corpus<br>MultiNLI<br>SciTail | https://nlp.stanford.edu/projects/snli/<br>https://www.nyu.edu/projects/bowman/multinli/<br>http://data.allenai.org/scitail/ |
| Semantic Role Labeling | Proposition Bank<br>OneNotes | http://propbank.github.io/<br>https://catalog.ldc.upenn.edu/LDC2013T19 |

# Aurélien Géron (2019),
# Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems, 2nd Edition
# O'Reilly Media, 2019



https://github.com/ageron/handson-ml2

# Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow

**Notebooks**
1. The Machine Learning landscape
2. End-to-end Machine Learning project
3. Classification
4. Training Models
5. Support Vector Machines
6. Decision Trees
7. Ensemble Learning and Random Forests
8. Dimensionality Reduction
9. Unsupervised Learning Techniques
10. Artificial Neural Nets with Keras
11. Training Deep Neural Networks
12. Custom Models and Training with TensorFlow
13. Loading and Preprocessing Data
14. Deep Computer Vision Using Convolutional Neural Networks
15. Processing Sequences Using RNNs and CNNs
16. Natural Language Processing with RNNs and Attention
17. Representation Learning Using Autoencoders
18. Reinforcement Learning
19. Training and Deploying TensorFlow Models at Scale

https://github.com/ageron/handson-ml2

# Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow

# Sequences using RNNs and CNNs

# TensorFlow

# TensorFlow

## is an

## Open Source
## Software Library

### for

# Machine Intelligence

https://www.tensorflow.org/

# TensorFlow 2.0

```python
import tensorflow as tf
mnist = tf.keras.datasets.mnist

(x_train, y_train),(x_test, y_test) = mnist.load_data()
x_train, x_test = x_train / 255.0, x_test / 255.0

model = tf.keras.models.Sequential([
  tf.keras.layers.Flatten(input_shape=(28, 28)),
  tf.keras.layers.Dense(128, activation='relu'),
  tf.keras.layers.Dropout(0.2),
  tf.keras.layers.Dense(10, activation='softmax')
])

model.compile(optimizer='adam',
              loss='sparse_categorical_crossentropy',
              metrics=['accuracy'])

model.fit(x_train, y_train, epochs=5)
model.evaluate(x_test, y_test)
```

https://www.tensorflow.org/overview/

# TensorFlow Playground

# Tensor

- **3**
  - # a rank 0 tensor; this is a **scalar** with shape []
- **[1. ,2., 3.]**
  - # a rank 1 tensor; this is a **vector** with shape [3]
- **[[1., 2., 3.], [4., 5., 6.]]**
  - # a rank 2 tensor; a **matrix** with shape [2, 3]
- **[[[1., 2., 3.]], [[7., 8., 9.]]]**
  - # a rank 3 **tensor** with shape [2, 1, 3]

https://www.tensorflow.org/

**Scalar**      80

**Vector**      $[50 \ 60 \ 70]$

**Matrix**      $\begin{bmatrix} 50 & 60 & 70 \\ 55 & 65 & 75 \end{bmatrix}$

**Tensor**      $\begin{bmatrix} [50 \ 60 \ 70] & [70 \ 80 \ 90] \\ [55 \ 65 \ 75] & [75 \ 85 \ 95] \end{bmatrix}$

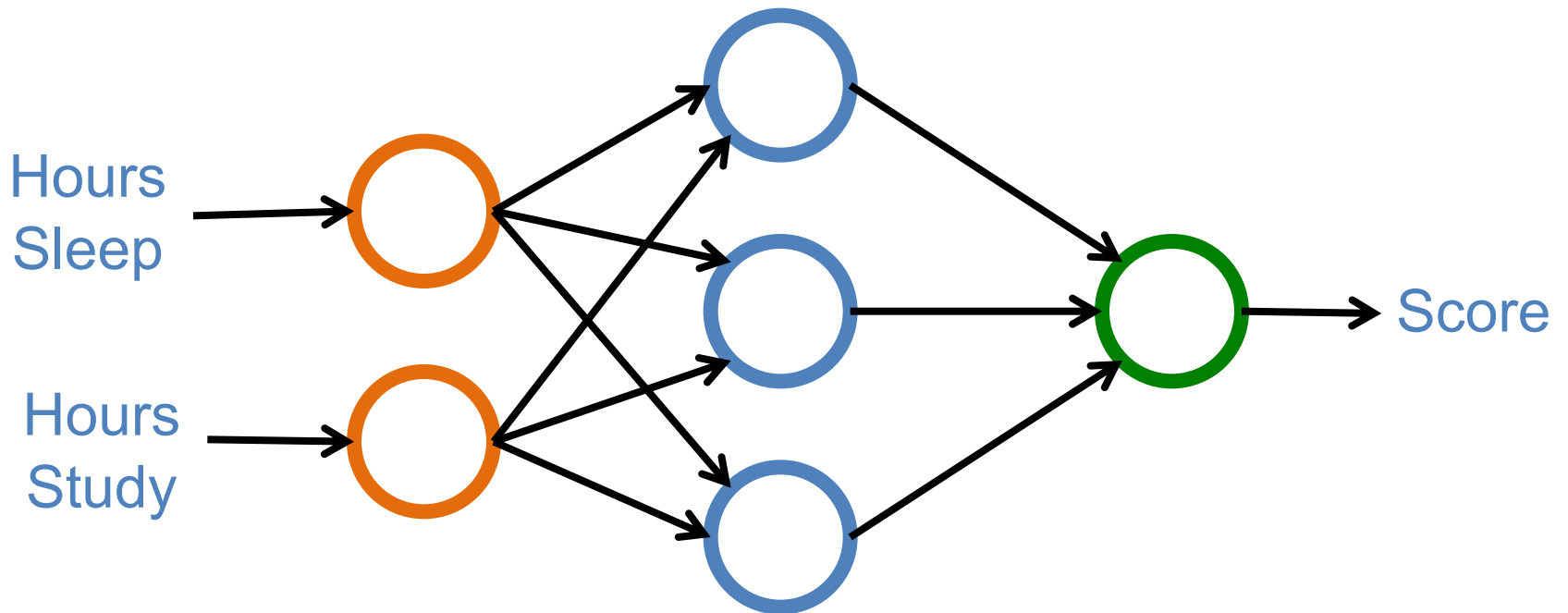# Deep Learning and Neural Networks

# Deep Learning and Neural Networks

**Input Layer (X)**    Hidden Layer (H)    **Output Layer (Y)**

# Deep Learning and Neural Networks
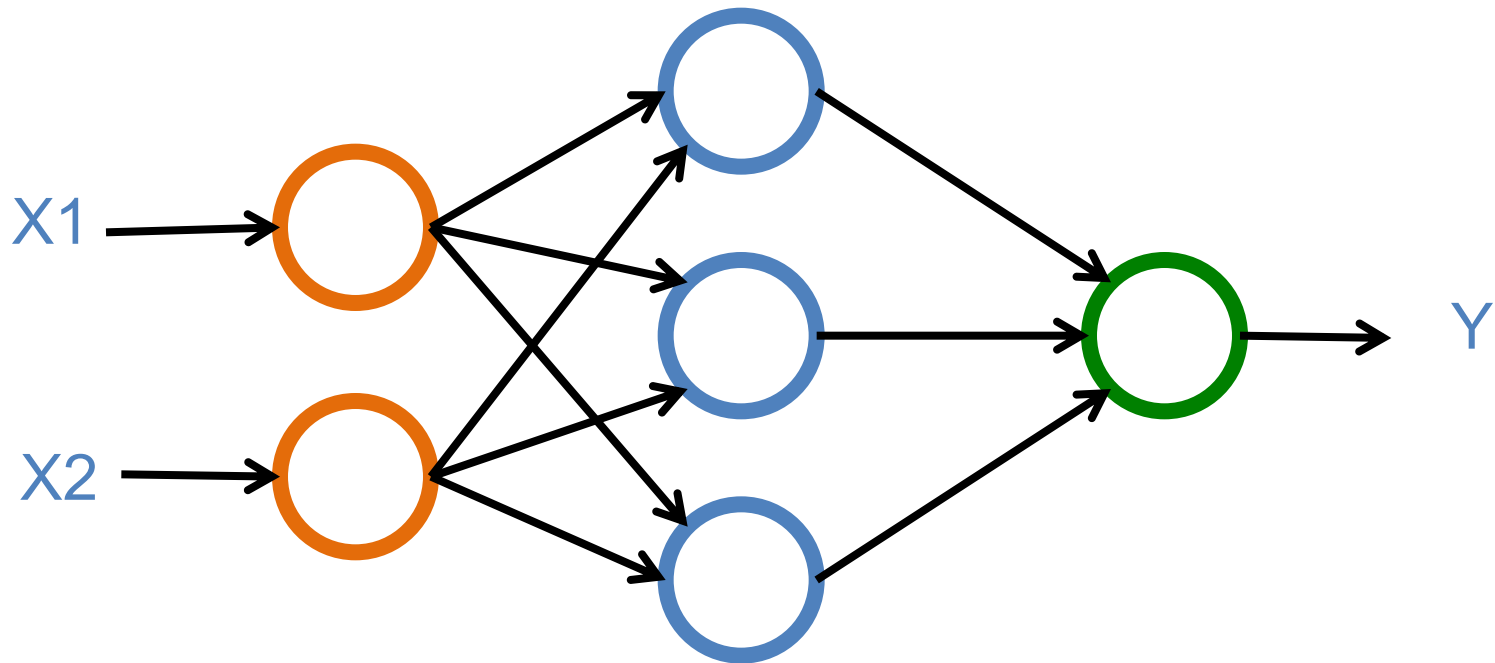
Input Layer
(X)

Hidden Layer
(H)

Output Layer
(Y)

# Deep Learning and Neural Networks

**Input Layer (X)**

Hidden Layers (H)

**Output Layer (Y)**

Deep Neural Networks
Deep Learning

# Convolutional Neural Networks

## (CNN or Deep Convolutional Neural Networks, DCNN)



LeCun, Yann, et al. "Gradient-based learning applied to document recognition." Proceedings of the IEEE 86.11 (1998): 2278-2324.

Source: http://www.asimovinstitute.org/neural-network-zoo/

# Recurrent Neural Networks (RNN)



Elman, Jeffrey L. "Finding structure in time." Cognitive science 14.2 (1990): 179-211

# Long / Short Term Memory (LSTM)

# Gated Recurrent Units (GRU)

Chung, Junyoung, et al. "Empirical evaluation of gated recurrent neural networks on sequence modeling." arXiv preprint arXiv:1412.3555 (2014).
Source: http://www.asimovinstitute.org/neural-network-zoo/

# Generative Adversarial Networks (GAN)



Goodfellow, Ian, et al. "Generative adversarial nets." Advances in Neural Information Processing Systems. 2014.

# Support Vector Machines (SVM)



Cortes, Corinna, and Vladimir Vapnik. "Support-vector networks." Machine learning 20.3 (1995): 273-297.

# Neural Networks

**Input Layer (X)** — Hidden Layer (H) — **Output Layer (Y)**

# Neural Networks

**Input Layer (X)**     Hidden Layer (H)     **Output Layer (Y)**

# The Neuron

# Neuron and Synapse

# The Neuron

$$y = F\left(\sum_i w_i x_i\right)$$



$x_1$    $w_1$

$x_2$   $w_2$

$\ldots$

$x_n$   $w_n$

$y$

$$F(x) = \max(0, x)$$

84

$$y = max ( 0, \text{-}0.21 * x_1 + 0.3 * x_2 + 0.7 * x_3 )$$

Weights

Inputs

$x_1$

$x_2$

$x_3$

-0.21

0.3

0.7

$y$

# Neural Networks

# Neural Networks

**Input Layer (X)**     Hidden Layer (H)     **Output Layer (Y)**

# Neural Networks

**Input Layer (X)**     Hidden Layers (H)     **Output Layer (Y)**

Deep Neural Networks
Deep Learning

# Neural Networks

**Input Layer (X)**  Hidden Layer (H)  **Output Layer (Y)**

Neuron

Synapse   Synapse

Neuron

X1

X2

Y

# Neural Networks



Input Layer (X)  Hidden Layer (H)  Output Layer (Y)

Hours Sleep, Hours Study → Score

# Neural Networks

Input Layer
(X)

Hidden Layer
(H)

Output Layer
(Y)

# Neural Networks

**Input Layer (X)**     Hidden Layer (H)     **Output Layer (Y)**



X1

X2

Y

| X | | Y |
|---|---|---|
| **Hours Sleep** | **Hours Study** | **Score** |
| 3 | 5 | 75 |
| 5 | 1 | 82 |
| 10 | 2 | 93 |
| 8 | 3 | ? |

|  | X | | Y |
| --- | --- | --- | --- |
|  | **Hours Sleep** | **Hours Study** | **Score** |
| **Training** | 3 | 5 | 75 |
| | 5 | 1 | 82 |
| | 10 | 2 | 93 |
| **Testing** | 8 | 3 | ? |

$$Y = WX + b$$

$$Y = W X + b$$

Output input

Weights bias

Trained

$$\mathbf{W}\,\mathbf{X} + \mathbf{b} = \mathbf{Y}$$

| 2.0 |
| 1.0 |
| 0.1 |

| 0.7 |
| 0.2 |
| 0.1 |

Scores ⟶ Probabilities

# SoftMAX

$$\mathbf{W} \, \mathbf{X} + \mathbf{b} = \mathbf{Y}$$

$$\begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix}$$

$$S(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

$$\begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$$

**Logits**          Scores ──────────────▶ Probabilities

$$S(y_i) = \frac{e^{y_i}}{\sum_j e^{y_i}} = \frac{e^{2.0}}{e^{2.0} + e^{1.0} + e^{0.1}} = \frac{2.7182^{2.0}}{2.7182^{2.0} + 2.7182^{1.0} + 2.7182^{0.1}} = 0.7$$

$$S(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}} = \frac{e^{1.0}}{e^{2.0} + e^{1.0} + e^{0.1}} = \frac{2.7182^{1.0}}{2.7182^{2.0} + 2.7182^{1.0} + 2.7182^{0.1}} = 0.2$$

$$S(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}} = \frac{e^{0.1}}{e^{2.0} + e^{1.0} + e^{0.1}} = \frac{2.7182^{0.1}}{2.7182^{2.0} + 2.7182^{1.0} + 2.7182^{0.1}} = 0.1$$

**W X + b = Y**

$$\begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix} \quad S(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}} \quad \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$$

**Logits**          Scores ⟶ Probabilities

# Training a Network

# =

# Minimize the Cost Function

# Training a Network

# =

# Minimize the Cost Function
# Minimize the Loss Function

**Error** = **Predict Y** - **Actual Y**
**Error : Cost : Loss**

# **Error** = **Predict Y** - **Actual Y**
## **Error : Cost : Loss**

**Error = Predict Y - Actual Y**
**Error : Cost : Loss**

# Activation Functions

# Activation Functions

## Sigmoid

## TanH

## ReLU
(Rectified Linear Unit)



**[0, 1]**

**[-1, 1]**

$f(x) = max(0, x)$

# Activation Functions

## Sigmoid

$$f(x) = \frac{1}{1 + e^{-x}}$$

## TanH

$$\tanh(x) = \frac{2}{1 + e^{-2x}} - 1$$

## ReLU

$$f(x) = \begin{cases} 0 & \text{for} \quad x < 0 \\ x & \text{for} \quad x \geq 0 \end{cases}$$

# Loss Function

# Binary Classification: 2 Class

# Activation Function: Sigmoid

# Loss Function: Binary Cross-Entropy

# Multiple Classification: 10 Class

# Activation Function: SoftMAX

# Loss Function: Categorical Cross-Entropy

# Dropout

Dropout: a simple way to prevent neural networks from overfitting



(a) Standard Neural Net

(b) After applying dropout.

Source: Srivastava, Nitish, Geoffrey E. Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. "Dropout: a simple way to prevent neural networks from overfitting." *Journal of machine learning research* 15, no. 1 (2014): 1929-1958.

# Learning Algorithm

While not done:

Pick a random training example "(input, label)"

Run neural network on "input"

Adjust weights on edges to make output closer to "label"

$$y = max\ (\ 0,\ \textcolor{red}{-0.21}\ *\ x_1\ +\ \textcolor{red}{0.3}\ *\ x_2\ +\ \textcolor{red}{0.7}\ *\ x_3\ )$$

Weights

Inputs

$x_1$

$x_2$

$x_3$

-0.21

0.3

0.7

$y$

# Next time:

$$y = max(0, -0.23 * x_1 + 0.31 * x_2 + 0.65 * x_3)$$

$$y = max(0, -0.21 * x_1 + 0.3 * x_2 + 0.7 * x_3)$$

Weights



Inputs

$x_1$

$x_2$

$x_3$

-0.23
-0.21

0.31
0.3

0.65
0.7

$y$

# Optimizer:
## Stochastic Gradient Descent (SGD)

*This shows a function of 2 variables: real neural nets are functions of hundreds of millions of variables!*

# Neural Network and Deep Learning

# Gradient Descent
## how neural networks learn

# Backpropagation



Source: 3Blue1Brown (2017), What is backpropagation really doing? | Chapter 3, deep learning,
https://www.youtube.com/watch?v=Ilg3gGewQ5U

119

# From image to text



Vision
Deep CNN

Language
Generating RNN

A group of people shopping at an outdoor market.

There are many vegetables at the fruit stand.

A woman is throwing a **frisbee** in a park.

A **dog** is standing on a hardwood floor.

A **stop** sign is on a road with a mountain in the background

A little **girl** sitting on a bed with a teddy bear.

A group of **people** sitting on a boat in the water.

A giraffe standing in a forest with **trees** in the background.

# From image to text

**Image: deep convolution neural network (CNN)**
**Text: recurrent neural network (RNN)**



A group of **people** sitting on a boat in the water.

# Convolutional Neural Networks (CNN)

# Convolutional Neural Networks (CNN)



Architecture of LeNet-5 (7 Layers)
(LeCun et al., 1998)

# Convolutional Neural Networks (CNN)

- Convolution

- Pooling

- Fully Connection (FC) (Flattening)

# CNN Architecture



Input      Conv      Pool      Conv      Pool      FC    FC    Softmax

# CNN Convolution Layer

**Convolution** is a **mathematical operation** to **merge two sets** of information

### 3x3 convolution

| | | | | |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

| | | |
|---|---|---|
| 1 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 1 |

Input

Filter / Kernel

# CNN Convolution Layer
# Input x Filter --> Feature Map

receptive field: 3x3



Input x Filter

Feature Map

# CNN **Convolution** Layer
# Input x Filter --> Feature Map

receptive field: 3x3



Input x Filter

Feature Map

# CNN **Convolution** Layer

| | | | | |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

Input

| | | |
|---|---|---|
| 1 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 1 |

Filter / Kernel

| | | | | |
|---|---|---|---|---|
| 1x1 | 1x0 | 1x1 | 0 | 0 |
| 0x0 | 1x1 | 1x0 | 1 | 0 |
| 0x1 | 0x0 | 1x1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

| | | |
|---|---|---|
| 4 | | |
| | | |
| | | |

Example convolution operation shown in 2D using a 3x3 filter

Source: Arden Dertat (2017), Applied Deep Learning - Part 4: Convolutional Neural Networks,
https://towardsdatascience.com/applied-deep-learning-part-4-convolutional-neural-networks-584bc134c1e2

# CNN Convolution Layer

10 different filters  10 feature maps of size 32x32x1



5x5x3

1x1x1

32x32x1

32

32

3

32

32

32

10

final output of the convolution layer:
a volume of size 32x32x10

# CNN **Convolution** Layer
## Sliding operation at 4 locations

# CNN **Convolution** Layer

two feature maps

# CNN **Convolution** Layer

**Stride** specifies how much
we move the convolution filter at each step



Stride 1

Feature Map

# CNN **Convolution** Layer

**Stride** specifies how much
we move the convolution filter at each step



Stride 2

Feature Map

# CNN **Convolution** Layer

## **Stride** 1 with **Padding**



Stride 1 with Padding

Feature Map

# CNN Pooling Layer

## Max Pooling



max pool with 2x2 window and stride 2

# CNN Pooling Layer



32

pooling →

16

16

32

10

10

# CNN Architecture
## 4 convolution + pooling layers, followed by 2 fully connected layers



Input    Conv + Maxpool    Conv + Maxpool    Conv + Maxpool    Conv + Maxpool    FC    FC    Output

# CNN Architecture
## 4 convolution + pooling layers, followed by 2 fully connected layers

https://gist.github.com/ardendertat/0fc5515057c47e7386fe04e9334504e3

```python
model = Sequential()
model.add(Conv2D(32, (3, 3), activation='relu', padding='same', name='conv_1',
                 input_shape=(150, 150, 3)))
model.add(MaxPooling2D((2, 2), name='maxpool_1'))
model.add(Conv2D(64, (3, 3), activation='relu', padding='same', name='conv_2'))
model.add(MaxPooling2D((2, 2), name='maxpool_2'))
model.add(Conv2D(128, (3, 3), activation='relu', padding='same', name='conv_3'))
model.add(MaxPooling2D((2, 2), name='maxpool_3'))
model.add(Conv2D(128, (3, 3), activation='relu', padding='same', name='conv_4'))
model.add(MaxPooling2D((2, 2), name='maxpool_4'))
model.add(Flatten())
model.add(Dropout(0.5))
model.add(Dense(512, activation='relu', name='dense_1'))
model.add(Dense(128, activation='relu', name='dense_2'))
model.add(Dense(1, activation='sigmoid', name='output'))
```

# Dropout



No Dropout                                        With Dropout

# Model Performance



Train Loss: 0.054, Val Loss: 1.345

Starts Overfitting

Train Accuracy: 0.981, Val Accuracy: 0.732

141

# The activations of an example ConvNet architecture.

# ConvNets

32x32x3 CIFAR-10 image

first Convolutional layer



32

32

3

# ConvNets



$x_0$

$w_0$

synapse

axon from a neuron

$w_0 x_0$

dendrite

cell body

$w_1 x_1$

$\sum_i w_i x_i + b$  $f$

$f\left(\sum_i w_i x_i + b\right)$

output axon

$w_2 x_2$

activation function

# Convolution Demo

**Input Volume (+pad 1) (7x7x3)**

x[:,:,0]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 0 | 2 | 1 | 0 |
| 0 | 2 | 2 | 2 | 1 | 1 | 0 |
| 0 | 2 | 2 | 2 | 0 | 1 | 0 |
| 0 | 2 | 2 | 1 | 2 | 1 | 0 |
| 0 | 2 | 1 | 2 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,1]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 0 | 2 | 2 | 1 | 2 | 0 |
| 0 | 1 | 2 | 0 | 0 | 2 | 0 |
| 0 | 0 | 1 | 2 | 1 | 0 | 0 |
| 0 | 2 | 2 | 2 | 2 | 0 | 0 |
| 0 | 2 | 2 | 2 | 0 | 2 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,2]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 0 | 0 | 2 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 1 | 1 | 0 |
| 0 | 2 | 2 | 2 | 1 | 2 | 0 |
| 0 | 1 | 2 | 0 | 0 | 2 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

**Filter W0 (3x3x3)**

w0[:,:,0]

| -1 | -1 | 0 |
|----|----|---|
| 1  | 1  | 1 |
| -1 | 0  | 1 |

w0[:,:,1]

| 0 | 0 | 1 |
|---|---|---|
| 0 | 1 | 0 |
| 0 | 0 | 1 |

w0[:,:,2]

| -1 | -1 | 0  |
|----|----|----|
| 1  | 0  | -1 |
| -1 | 0  | -1 |

**Bias b0 (1x1x1)**

b0[:,:,0]

| 1 |
|---|

**Filter W1 (3x3x3)**

w1[:,:,0]

| 1 | -1 | 0 |
|---|----|---|
| 0 | 1  | 1 |
| 0 | -1 | 1 |

w1[:,:,1]

| -1 | 1  | 0 |
|----|----|---|
| -1 | -1 | 1 |
| 0  | 0  | 0 |

w1[:,:,2]

| 1 | 0 | -1 |
|---|---|----|
| 0 | 0 | -1 |
| 1 | 0 | 1  |

**Bias b1 (1x1x1)**

b1[:,:,0]

| 0 |
|---|

**Output Volume (3x3x2)**

o[:,:,0]

| 6 | 3  | 6  |
|---|----|----|
| 7 | -1 | -2 |
| 2 | 3  | -2 |

o[:,:,1]

| 7  | -1 | -3 |
|----|----|----|
| 4  | 3  | 2  |
| -1 | 0  | -1 |

toggle movement

http://cs231n.github.io/convolutional-networks/

145

# ConvNets

input volume of size [224x224x64]
is pooled with **filter** size 2, **stride** 2
into output volume of size [112x112x64]



224x224x64

pool →

112x112x64

224

224

downsampling →

112

112

# ConvNets
# max pooling

Single depth slice

x

| 1 | 1 | 2 | 4 |
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

y

max pool with 2x2 filters
and stride 2

| 6 | 8 |
| 3 | 4 |

# Convolutional Neural Networks (CNN) (LeNet)

# Recurrent Neural Networks (RNN)

# Recurrent Neural Networks (RNN)

# Recurrent Neural Networks (RNN)

# Recurrent Neural Networks (RNN) Sentiment Analysis

**output**

**hidden** $h_{t-2}$ → $h_{t-1}$ → $h_t$ → $h_{t+1}$ → $h_{t+2}$ → **y**

**Input** $X_{t-2}$ $X_{t-1}$ $X_t$ $X_{t+1}$ $X_{t+2}$

This        movie        is        very        good

# Recurrent Neural Networks (RNN) Sentiment Analysis

**output**

**hidden** $h_{t-2}$ → $h_{t-1}$ → $h_t$ → $h_{t+1}$ → $h_{t+2}$

**Input** $X_{t-2}$ $X_{t-1}$ $X_t$ $X_{t+1}$ $X_{t+2}$

This      movie      is      very      boring

**y**

# Recurrent Neural Network (RNN)

# RNN

# RNN long-term dependencies



I grew up in France… I speak fluent French.

# Vanishing Gradient
# Exploding Gradient

# Recurrent Neural Networks (RNN)



output $\mathbf{y}_{t-2}$ $\mathbf{y}_{t-1}$ $\mathbf{y}_t$ $\mathbf{y}_{t+1}$ $\mathbf{y}_{t+2}$

V V V V V

hidden $\mathbf{h}_{t-2}$ $\mathbf{h}_{t-1}$ $\mathbf{h}_t$ $\mathbf{h}_{t+1}$ $\mathbf{h}_{t+2}$

W W W W

U U U U

Input $\mathbf{X}_{t-2}$ $\mathbf{X}_{t-1}$ $\mathbf{X}_t$ $\mathbf{X}_{t+1}$ $\mathbf{X}_{t+2}$

# RNN
## Vanishing Gradient problem
## Exploding Gradient problem

**Error**



output $y_{t-2}$ $y_{t-1}$ $y_t$ $y_{t+1}$ $y_{t+2}$

V V V V

W W W W

hidden $h_{t-2}$ $h_{t-1}$ $h_t$ $h_{t+1}$ $h_{t+2}$

U U U U

Input $X_{t-2}$ $X_{t-1}$ $X_t$ $X_{t+1}$ $X_{t+2}$

**if |W| < 1 (Vanishing)**
**if |W| > 1 (Exploding)**

159

# RNN
# Vanishing Gradient problem



**Error**

output — $y_{t-2}$  $y_{t-1}$  $y_t$  $y_{t+1}$  $y_{t+2}$

V   V   V   V

0.9   0.9   0.9   0.9

hidden — $h_{t-2}$  $h_{t-1}$  $h_t$  $h_{t+1}$  $h_{t+2}$

U   U   U   U

Input — $X_{t-2}$  $X_{t-1}$  $X_t$  $X_{t+1}$  $X_{t+2}$

**W = 0.9 < 1 (Vanishing)**

# RNN
# Exploding Gradient problem

**Error**



output $y_{t-2}$ $y_{t-1}$ $y_t$ $y_{t+1}$ $y_{t+2}$

V V V V

1.1 1.1 1.1 1.1

hidden $h_{t-2}$ $h_{t-1}$ $h_t$ $h_{t+1}$ $h_{t+2}$

U U U U

Input $X_{t-2}$ $X_{t-1}$ $X_t$ $X_{t+1}$ $X_{t+2}$

**W = 1.1 > 1 (Exploding)**

# RNN LSTM

**RNN**



**LSTM**

162

# Long Short Term Memory (LSTM)

# Long Short Term Memory (LSTM)

# Gated Recurrent Unit (GRU)

# Gated Recurrent Unit (GRU)

# LSTM

# LSTM vs GRU



## LSTM

## GRU

i, f and o are the input, forget and output gates, respectively.
c and c˜ denote the memory cell and the new memory cell content.

r and z are the reset and update gates, and h and h˜ are the activation and the candidate activation.

Source: Chung, Junyoung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. "Empirical evaluation of gated recurrent neural networks on sequence modeling." *arXiv preprint arXiv:1412.3555* (2014).

# Long Short Term Memory (LSTM)

# Long Short Term Memory (LSTM)

# LSTM
# Memory state (C)

# LSTM
# *forget gate (f)*



$$f_t = \sigma \left( W_f \cdot [h_{t-1}, x_t] \ + \ b_f \right)$$

# LSTM
# input gate (i)



$$i_t = \sigma\left(W_i \cdot [h_{t-1}, x_t] + b_i\right)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

# LSTM
# Memory state (C)



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

# *LSTM*
# *output gate (o)*



$$o_t = \sigma \left( W_o \left[ h_{t-1}, x_t \right] + b_o \right)$$

$$h_t = o_t * \tanh \left( C_t \right)$$

# LSTM
# *forget (f), input (i), output (o) gates*



$$f_t = \sigma \left( W_f \cdot [\boldsymbol{C_{t-1}}, h_{t-1}, x_t] + b_f \right)$$

$$i_t = \sigma \left( W_i \cdot [\boldsymbol{C_{t-1}}, h_{t-1}, x_t] + b_i \right)$$

$$o_t = \sigma \left( W_o \cdot [\boldsymbol{C_t}, h_{t-1}, x_t] + b_o \right)$$

# Gated Recurrent Unit (GRU)
## *update (z), reset (r) gates*



$$z_t = \sigma \left( W_z \cdot [h_{t-1}, x_t] \right)$$

$$r_t = \sigma \left( W_r \cdot [h_{t-1}, x_t] \right)$$

$$\tilde{h}_t = \tanh \left( W \cdot [r_t * h_{t-1}, x_t] \right)$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

# LSTM Recurrent Neural Network



| one to one | one to many | many to one | many to many | many to many |
| --- | --- | --- | --- | --- |
| Traditional Neural Network | Music Generation | Sentiment Classification | Name Entity Recognition | Machine Translation |

# The Sequence to Sequence model (seq2seq)

# Sequence to Sequence (Seq2Seq)

# Transformer (Attention is All You Need)
## (Vaswani et al., 2017)
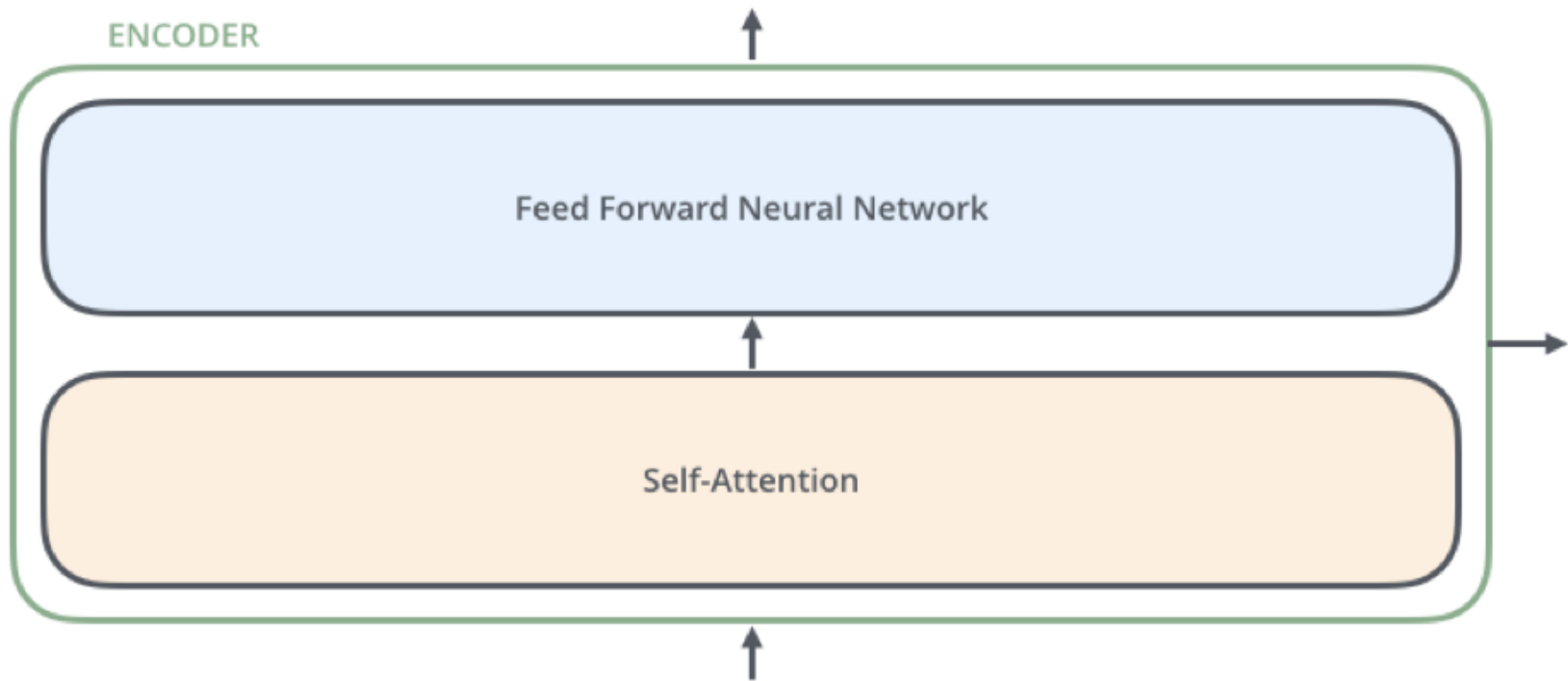
# Transformer

# Transformer
# Encoder Decoder

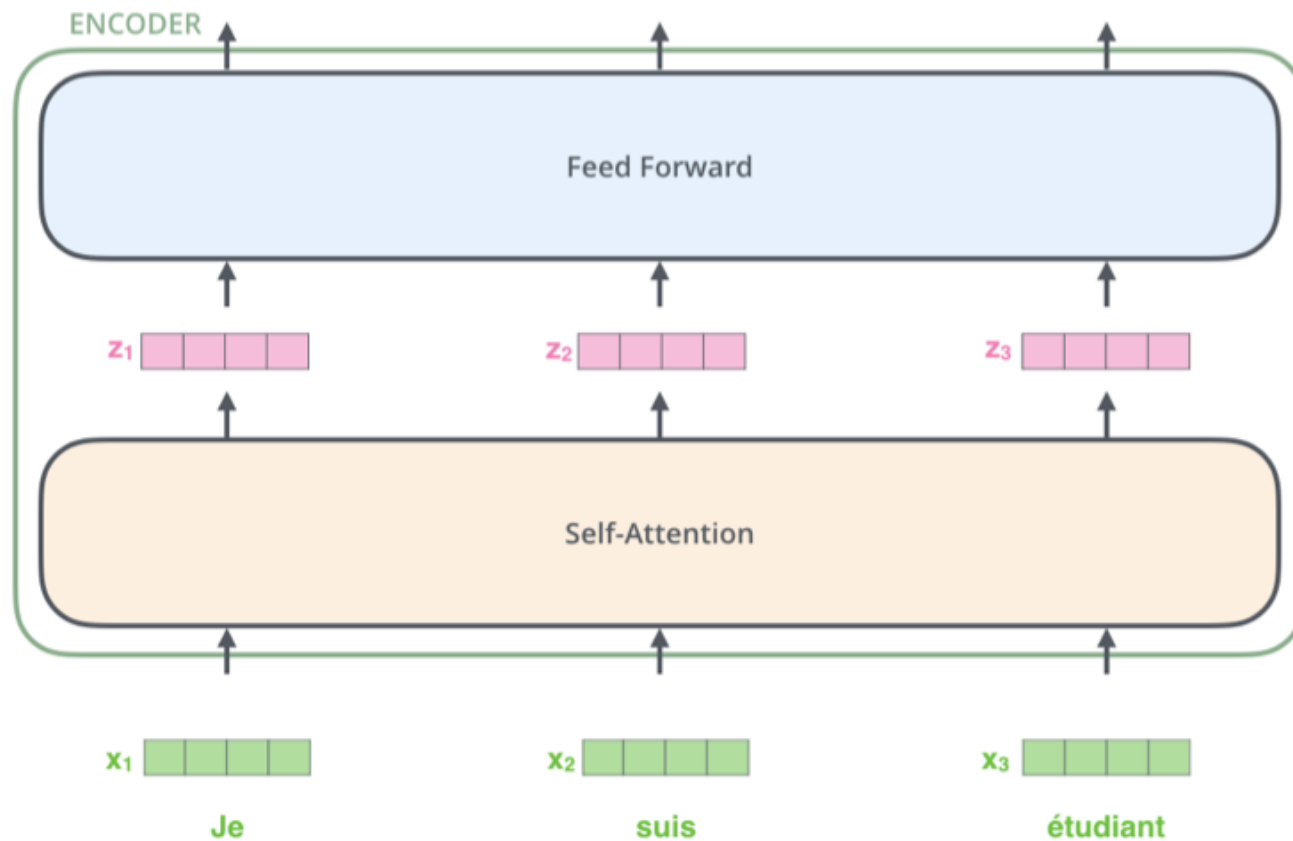# Transformer Encoder Decoder Stack

# Transformer
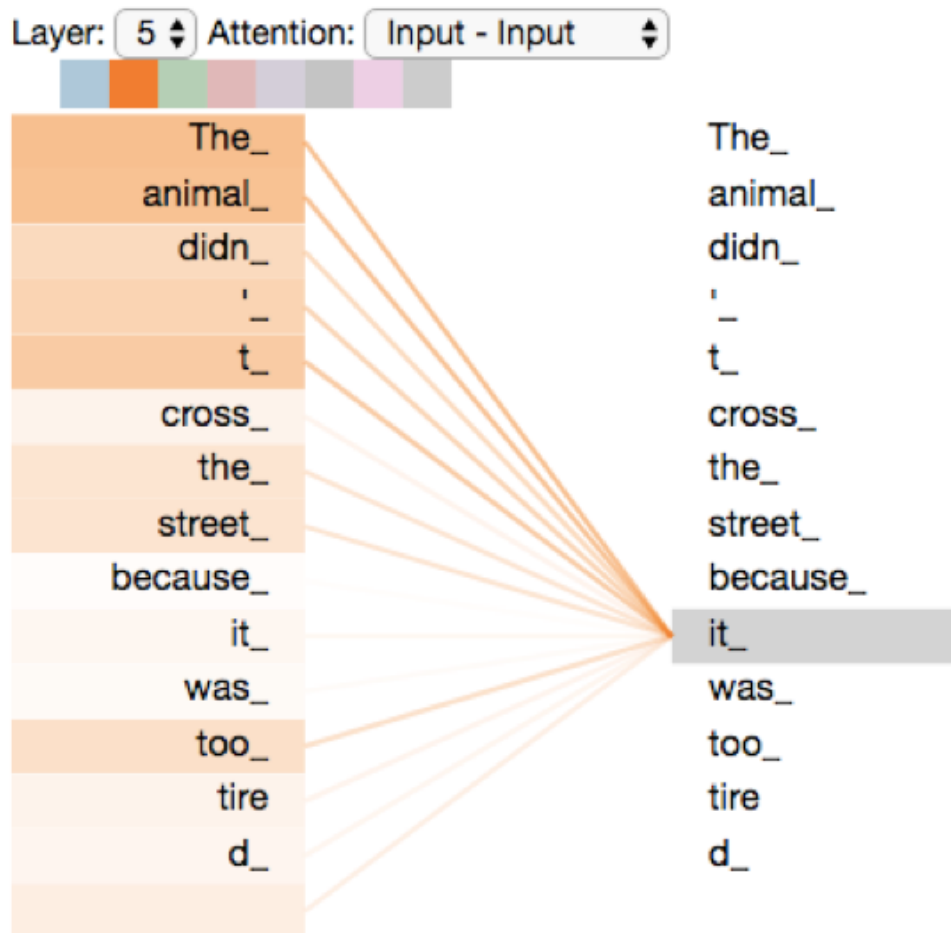# Encoder Self-Attention

# Transformer Decoder
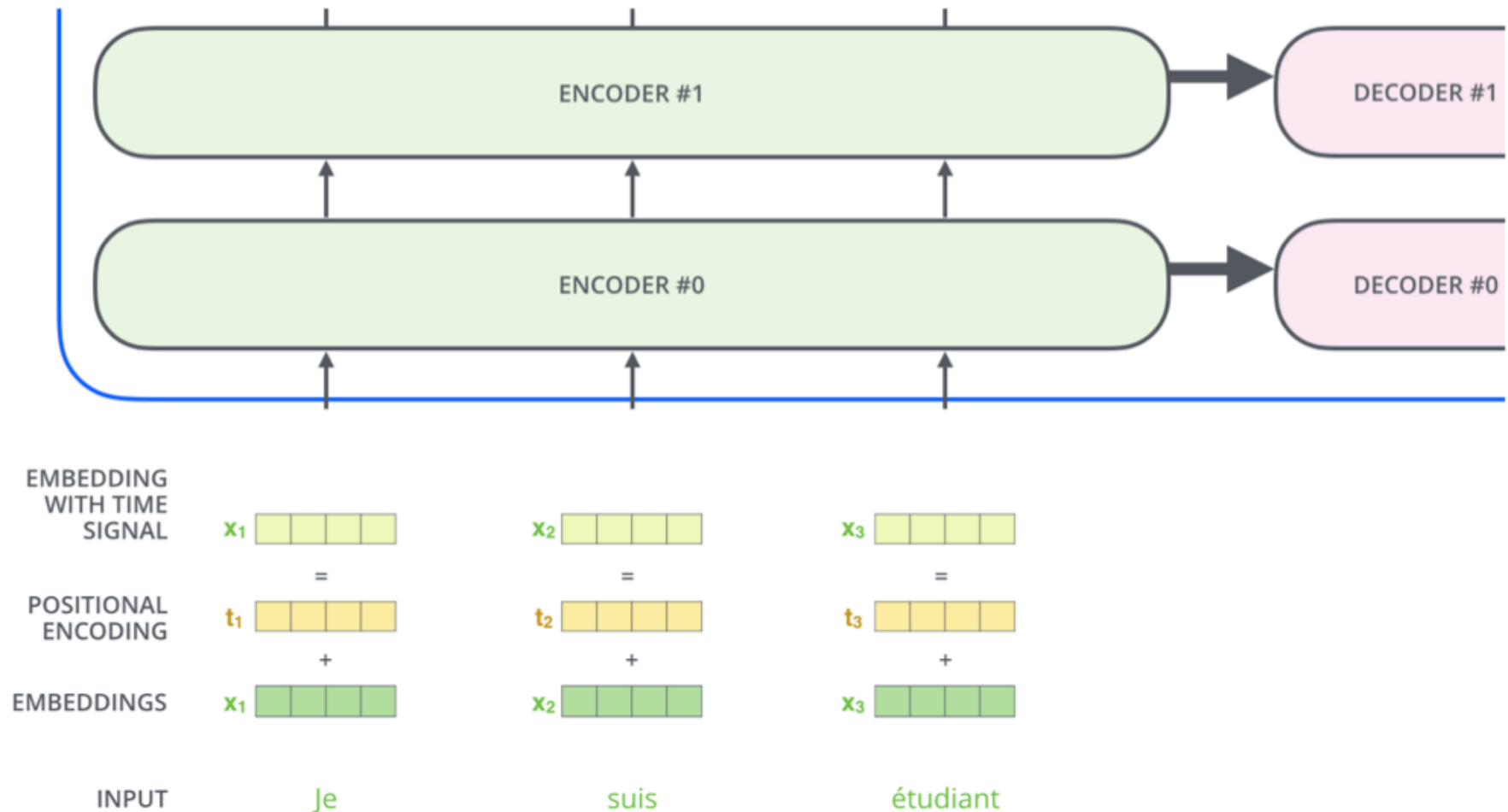
# Transformer
# Encoder with Tensors
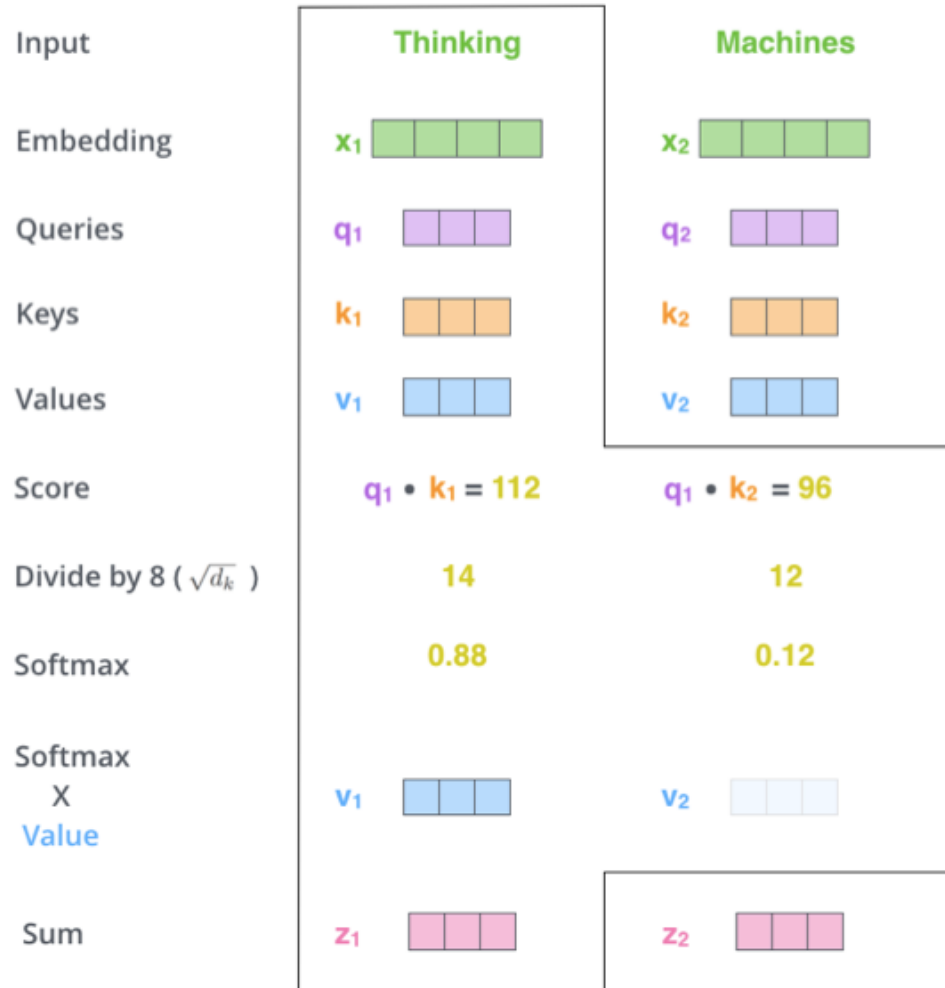# Word Embeddings

# Transformer
# Self-Attention Visualization

# Transformer
# Positional Encoding Vectors
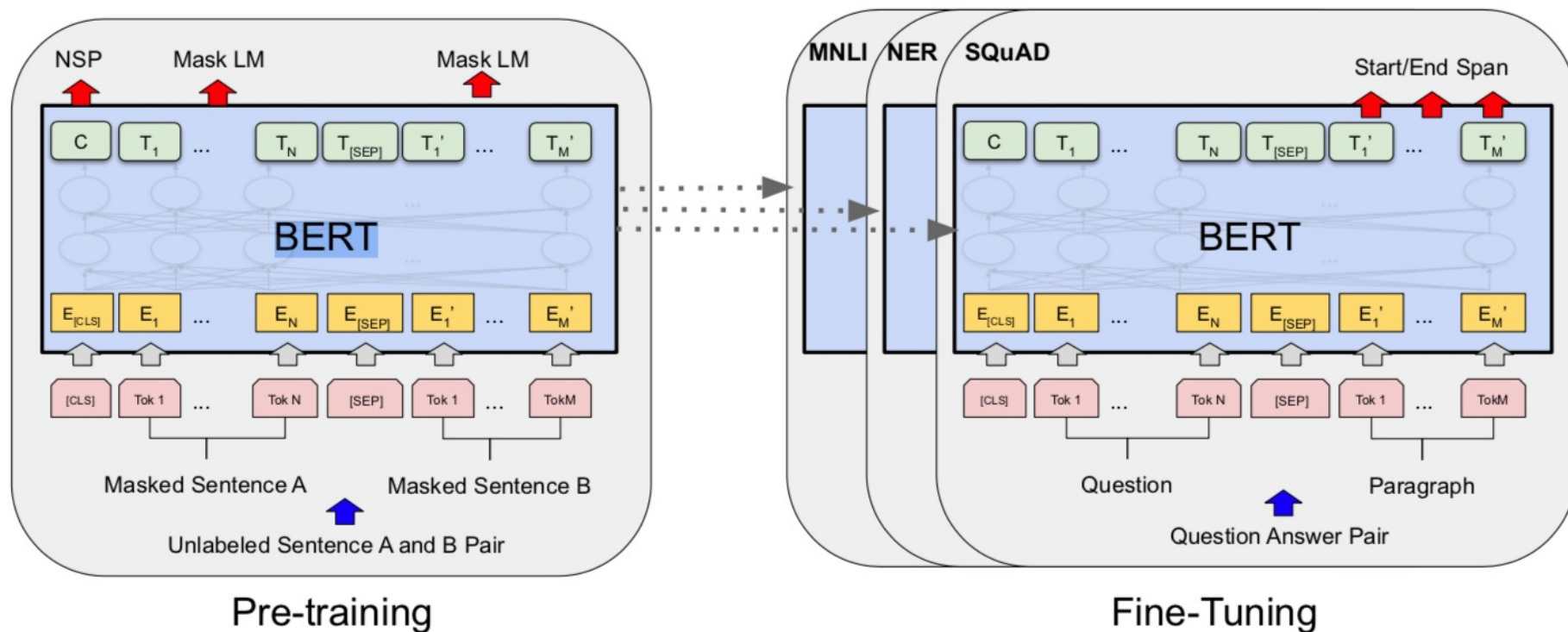
# Transformer
# Self-Attention Softmax Output

# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

## BERT
### (Bidirectional Encoder Representations from Transformers)

## Overall pre-training and fine-tuning procedures for BERT



Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).
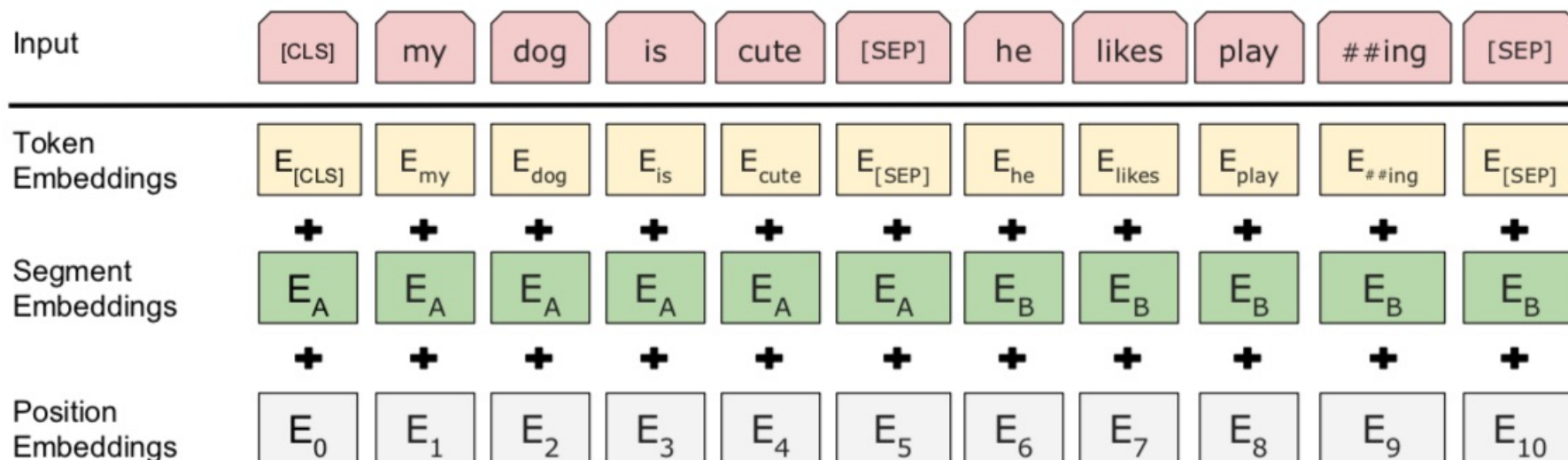"Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

## BERT
## (Bidirectional Encoder Representations from Transformers)
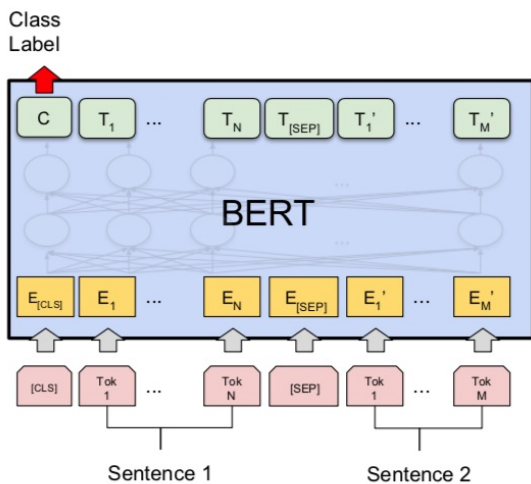
## BERT input representation

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).
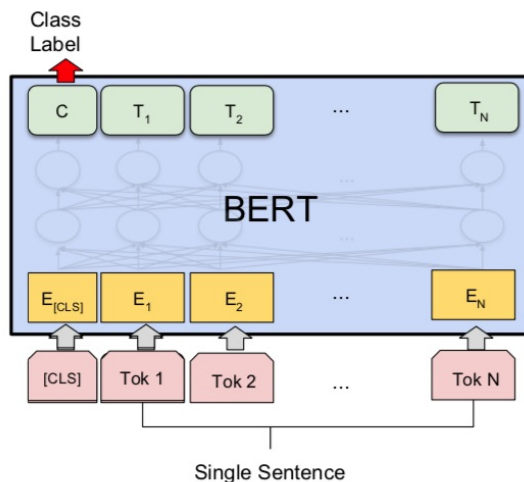"Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.
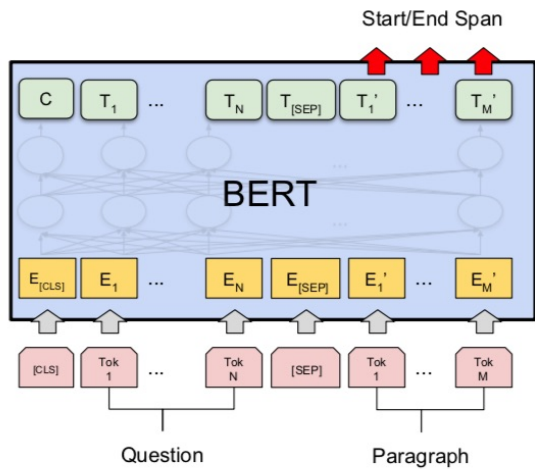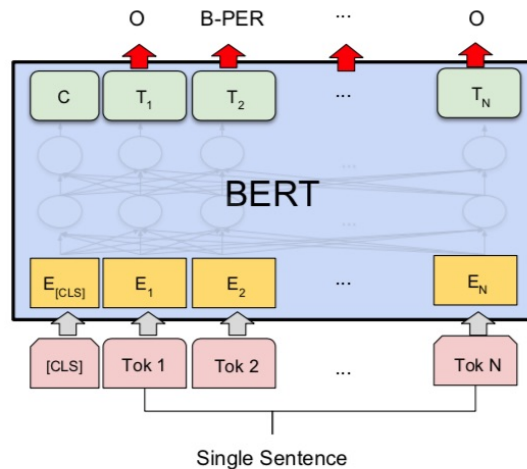
# Fine-tuning BERT on Different Tasks



(a) Sentence Pair Classification Tasks: MNLI, QQP, QNLI, STS-B, MRPC, RTE, SWAG

(b) Single Sentence Classification Tasks: SST-2, CoLA

(c) Question Answering Tasks: SQuAD v1.1

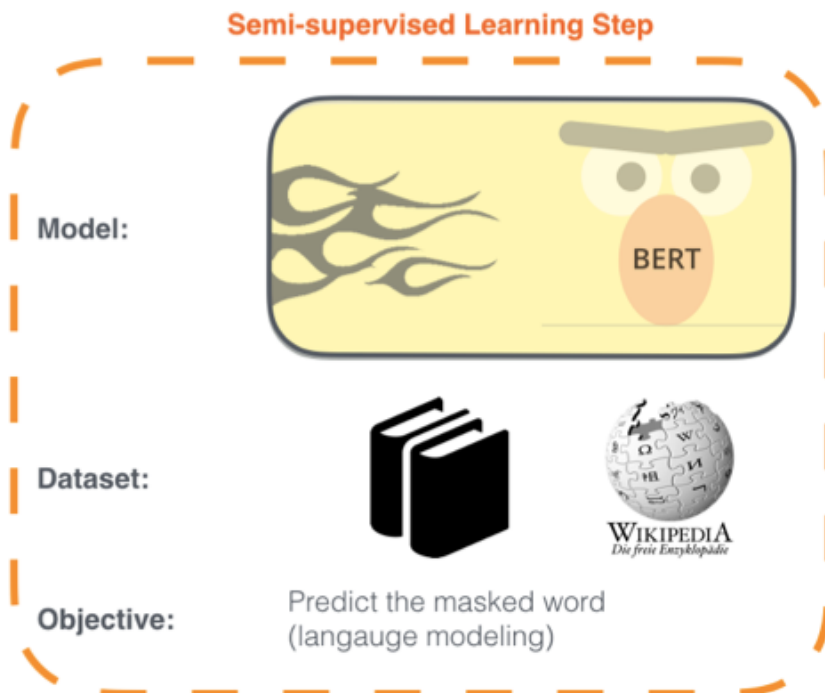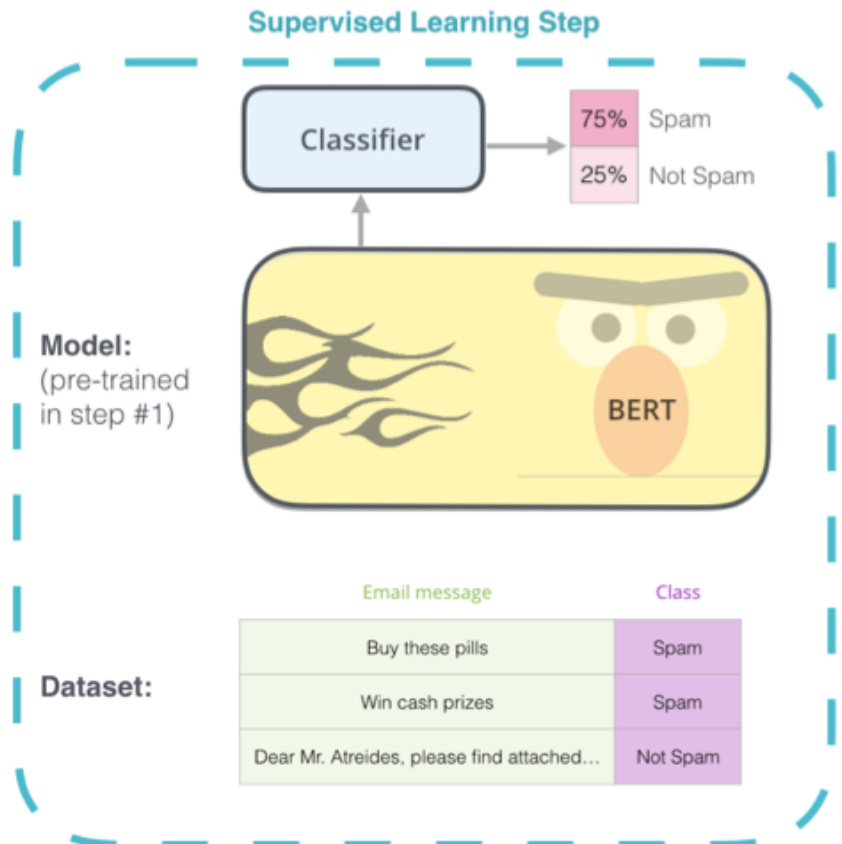(d) Single Sentence Tagging Tasks: CoNLL-2003 NER

193

# Illustrated BERT



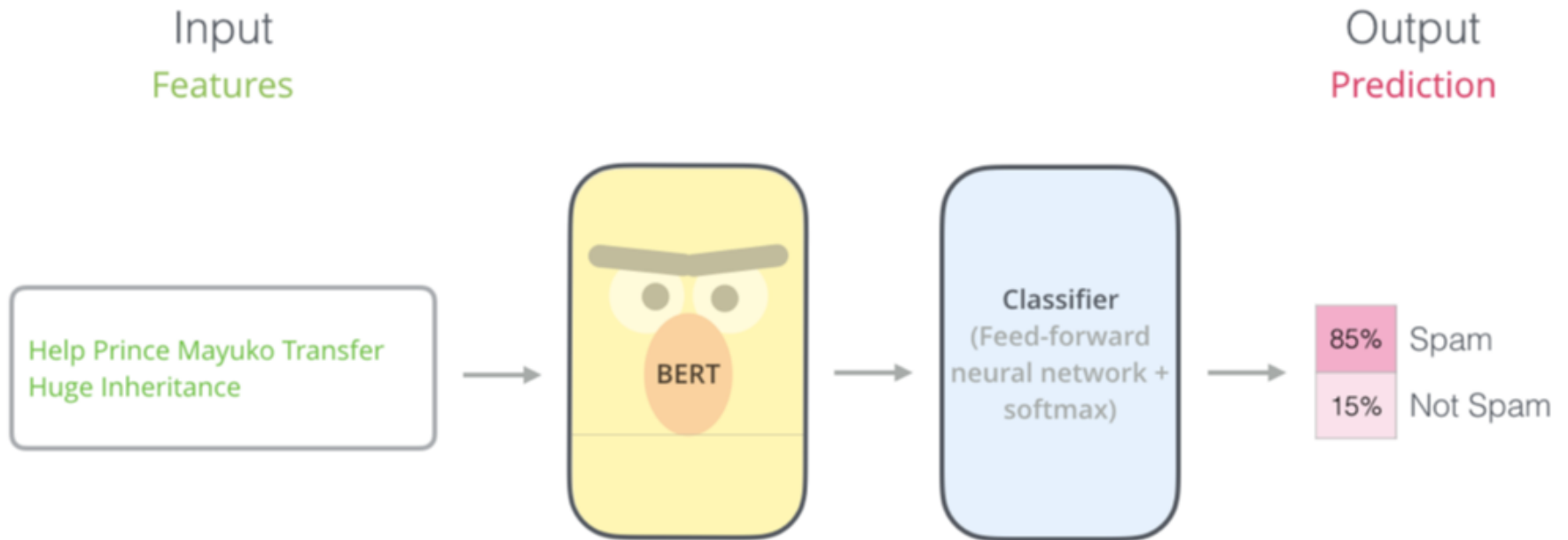1 - Semi-supervised training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.

**Semi-supervised Learning Step**

Model:

Dataset:

Objective: Predict the masked word (langauge modeling)

2 - Supervised training on a specific task with a labeled dataset.

**Supervised Learning Step**

Classifier → 75% Spam / 25% Not Spam

Model: (pre-trained in step #1)

Dataset:

| Email message | Class |
|---|---|
| Buy these pills | Spam |
| Win cash prizes | Spam |
| Dear Mr. Atreides, please find attached… | Not Spam |

# BERT Classification Input Output

# BERT Encoder Input

# BERT Classifier

# Summary

- Traditional Feature Engineering for Text Data
  - Bag of Words Model
  - Bag of N-Grams Model
  - TF-IDF Model
- Advanced Word Embeddings with Deep Learning
  - Word2Vec Model
  - Robust Word2Vec Models with Gensim
  - GloVe Model
  - FastText Model

# References

- Dipanjan Sarkar (2019), Text Analytics with Python: A Practitioner's Guide to Natural Language Processing, Second Edition. APress. https://github.com/Apress/text-analytics-w-python-2e
- Benjamin Bengfort, Rebecca Bilbro, and Tony Ojeda (2018), Applied Text Analysis with Python, O'Reilly Media. https://www.oreilly.com/library/view/applied-text-analysis/9781491963036/
- Aurélien Géron (2019), Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems, 2nd Edition, O'Reilly Media, 2019, https://github.com/ageron/handson-ml2
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018). "Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.
- Deep Learning Basics: Neural Networks Demystified, https://www.youtube.com/playlist?list=PLiaHhY2iBX9hdHaRr6b7XevZtgZRa1PoU
- Deep Learning SIMPLIFIED, https://www.youtube.com/playlist?list=PLjJh1vlSEYgvGod9wWiydumYl8hOXixNu
- 3Blue1Brown (2017), But what *is* a Neural Network? | Chapter 1, deep learning, https://www.youtube.com/watch?v=aircAruvnKk
- 3Blue1Brown (2017), Gradient descent, how neural networks learn | Chapter 2, deep learning, https://www.youtube.com/watch?v=IHZwWFHWa-w
- 3Blue1Brown (2017), What is backpropagation really doing? | Chapter 3, deep learning, https://www.youtube.com/watch?v=Ilg3gGewQ5U
- Jay Alammar (2019), The Illustrated Transformer, http://jalammar.github.io/illustrated-transformer/
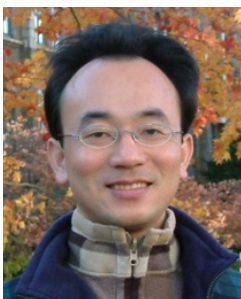- Jay Alammar (2019), The Illustrated BERT, ELMo, and co. (How NLP Cracked Transfer Learning), http://jalammar.github.io/illustrated-bert/
- Min-Yuh Day (2020), Python 101, https://tinyurl.com/imtkupython101