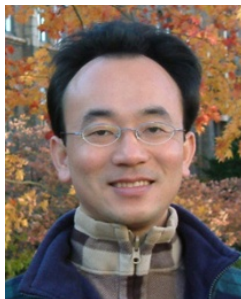# 人工智慧文本分析
# (AI for Text Analytics)

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

1091AITA05
MBA, IMTKU (M2455) (8418) (Fall 2020)
Thu 3, 4 (10:10-12:00) (B206)

**Min-Yuh Day**
**戴敏育**
**Associate Professor**
副教授
**Institute of Information Management**, **National Taipei University**
國立臺北大學 資訊管理研究所
https://web.ntpu.edu.tw/~myday

2020-10-22

# 課程大綱 (Syllabus)

週次 (Week)　日期 (Date)　內容 (Subject/Topics)

1 2020/09/17 人工智慧文本分析課程介紹
(Course Orientation on Artificial Intelligence for Text Analytics)

2 2020/09/24 文本分析的基礎：自然語言處理
(Foundations of Text Analytics: Natural Language Processing; NLP)

3 2020/10/01 中秋節 (Mid-Autumn Festival) 放假一天 (Day off)

4 2020/10/08 Python自然語言處理
(Python for Natural Language Processing)

5 2020/10/15 處理和理解文本
(Processing and Understanding Text)

6 2020/10/22 文本表達特徵工程
(Feature Engineering for Text Representation)

# 課程大綱 (Syllabus)

週次 (Week)　日期 (Date)　內容 (Subject/Topics)

7 2020/10/29 人工智慧文本分析個案研究 I
(Case Study on Artificial Intelligence for Text Analytics I)

8 2020/11/05 文本分類
(Text Classification)

9 2020/11/12 文本摘要和主題模型
(Text Summarization and Topic Models)

10 2020/11/19 期中報告 (Midterm Project Report)

11 2020/11/26 文本相似度和分群
(Text Similarity and Clustering)

12 2020/12/03 語意分析和命名實體識別
(Semantic Analysis and Named Entity Recognition; NER)

# 課程大綱 (Syllabus)

週次 (Week)　日期 (Date)　內容 (Subject/Topics)

13 2020/12/10 情感分析
(Sentiment Analysis)

14 2020/12/17 人工智慧文本分析個案研究 II
(Case Study on Artificial Intelligence for Text Analytics II)

15 2020/12/24 深度學習和通用句子嵌入模型
(Deep Learning and Universal Sentence-Embedding Models)

16 2020/12/31 問答系統與對話系統
(Question Answering and Dialogue Systems)

17 2021/01/07 期末報告 I (Final Project Presentation I)

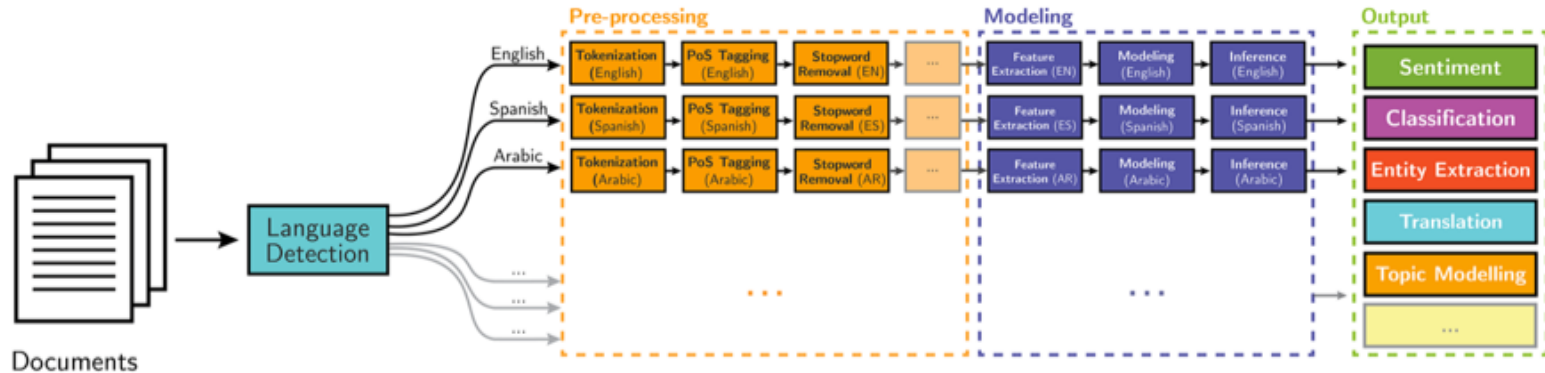18 2021/01/14 期末報告 II (Final Project Presentation II)

# 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
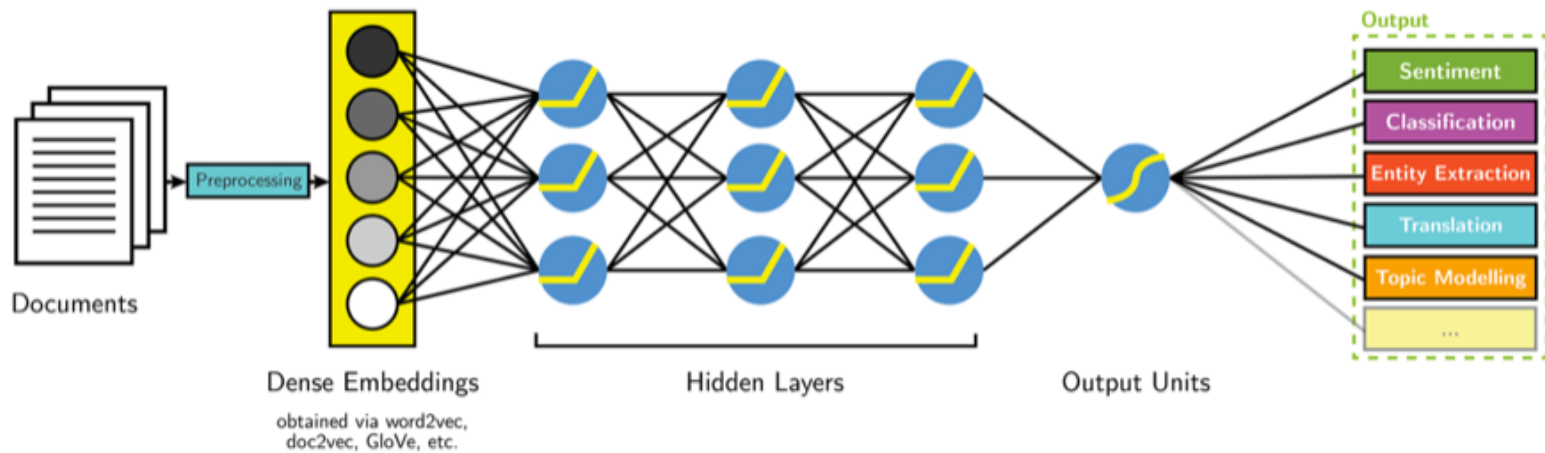  - GloVe Model
  - FastText Model
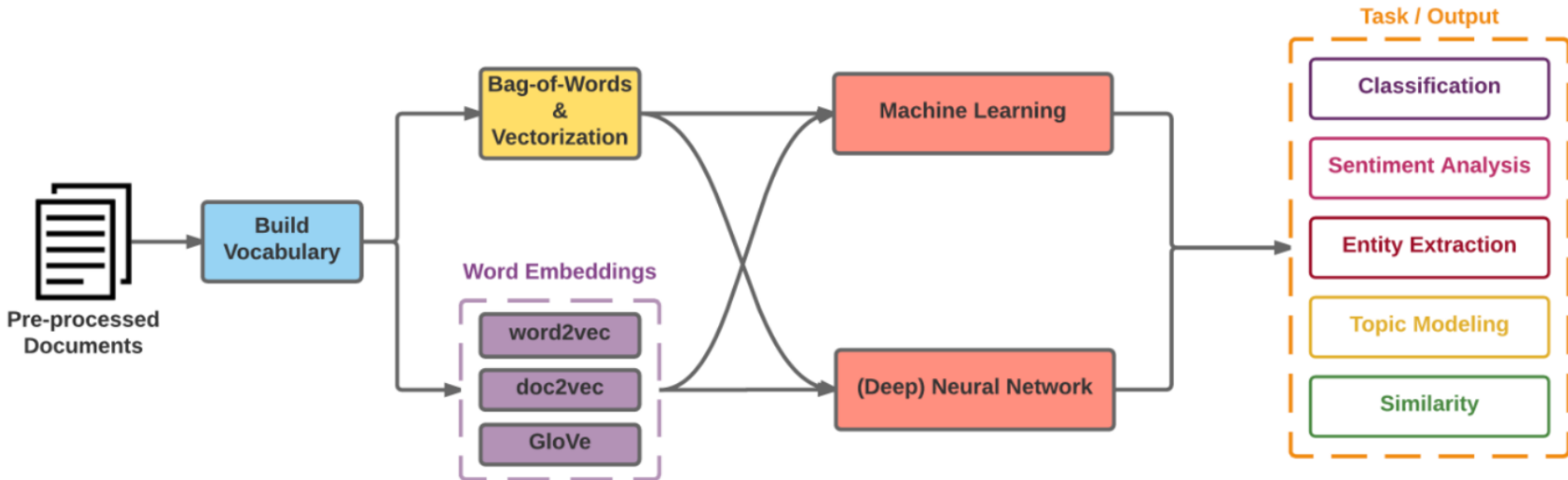
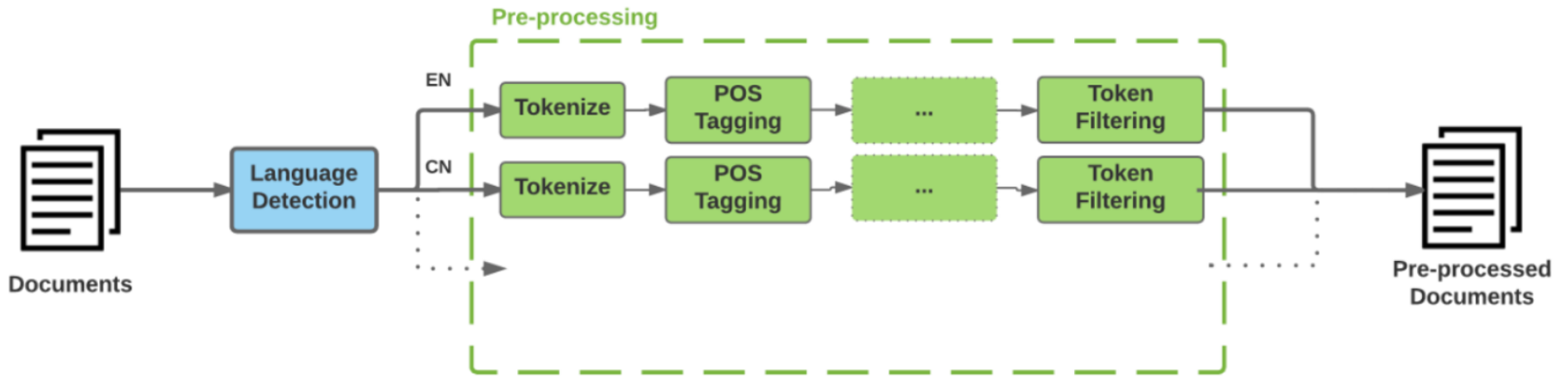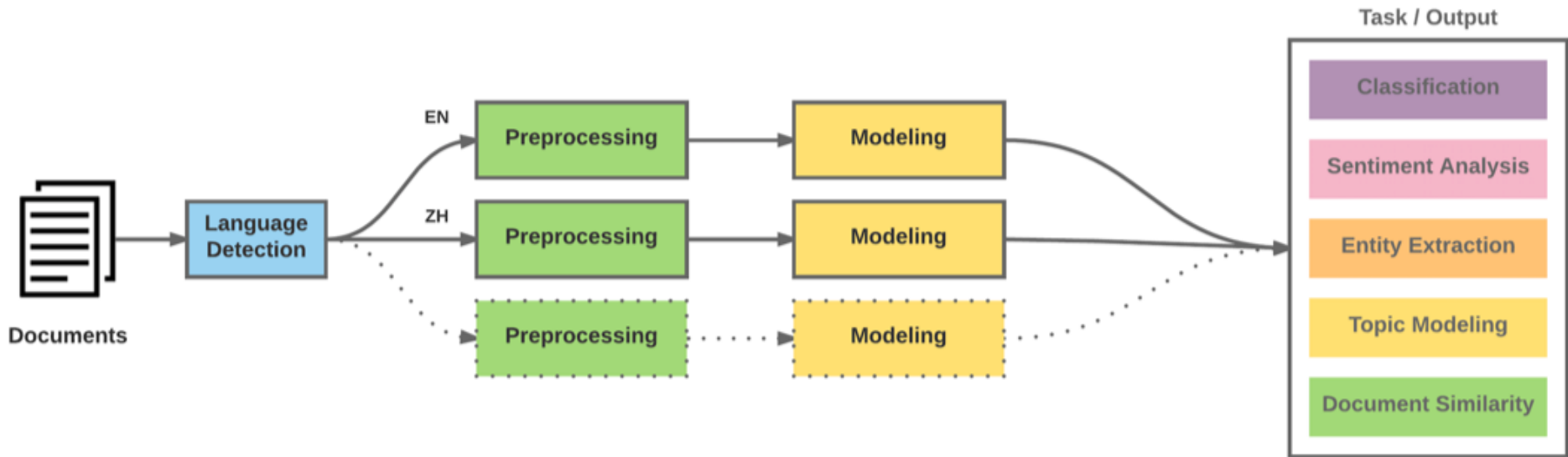# Feature Engineering for Text Representation

# NLP

# Modern NLP Pipeline

# Modern NLP Pipeline

# Deep Learning NLP

**Documents** → **Preprocessing** → **Dense Word Embeddings** → **Deep Neural Network** → **Task / Output**

*Pre-generated Lookup
OR
Generated in 1st level
of NeuralNet*

**Task / Output**
- Classification
- Sentiment Analysis
- Entity Extraction
- Topic Modeling
- Document Similarity

# Overview of
# Text Vectorization Methods

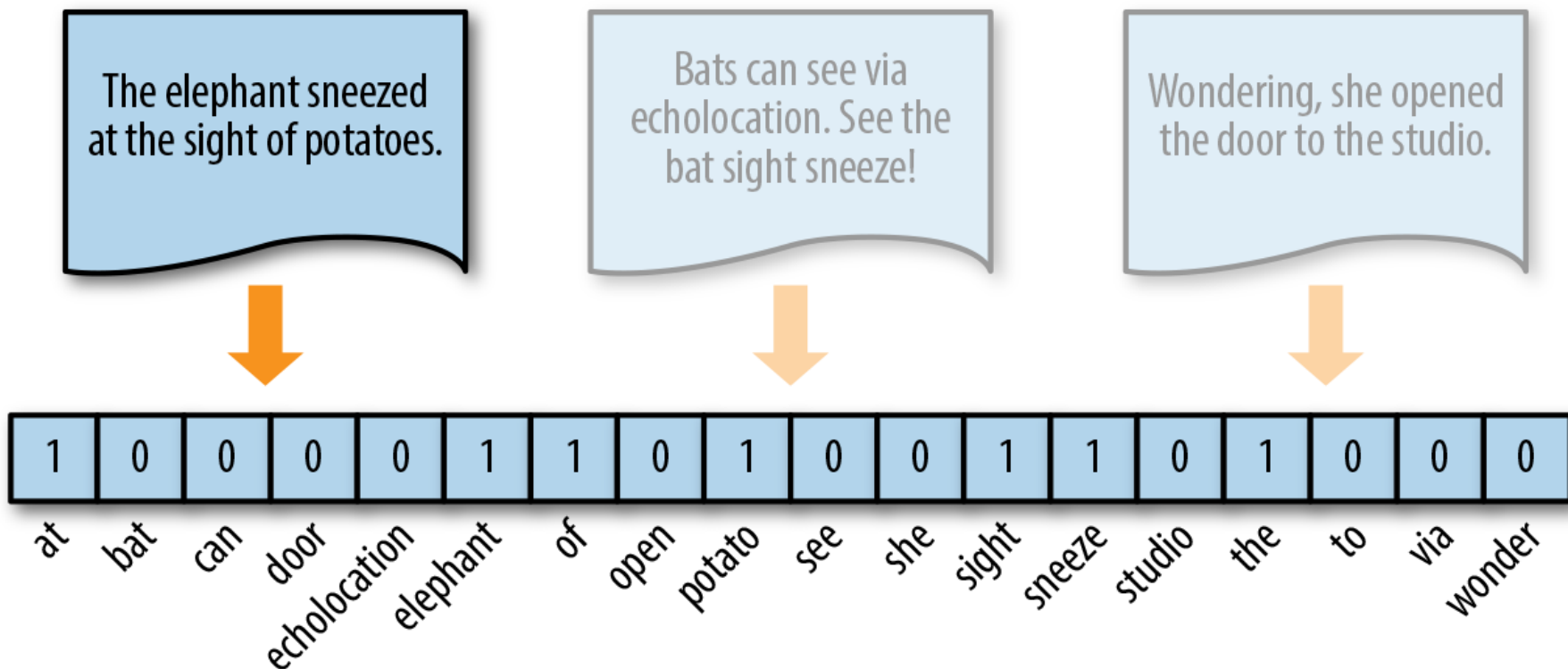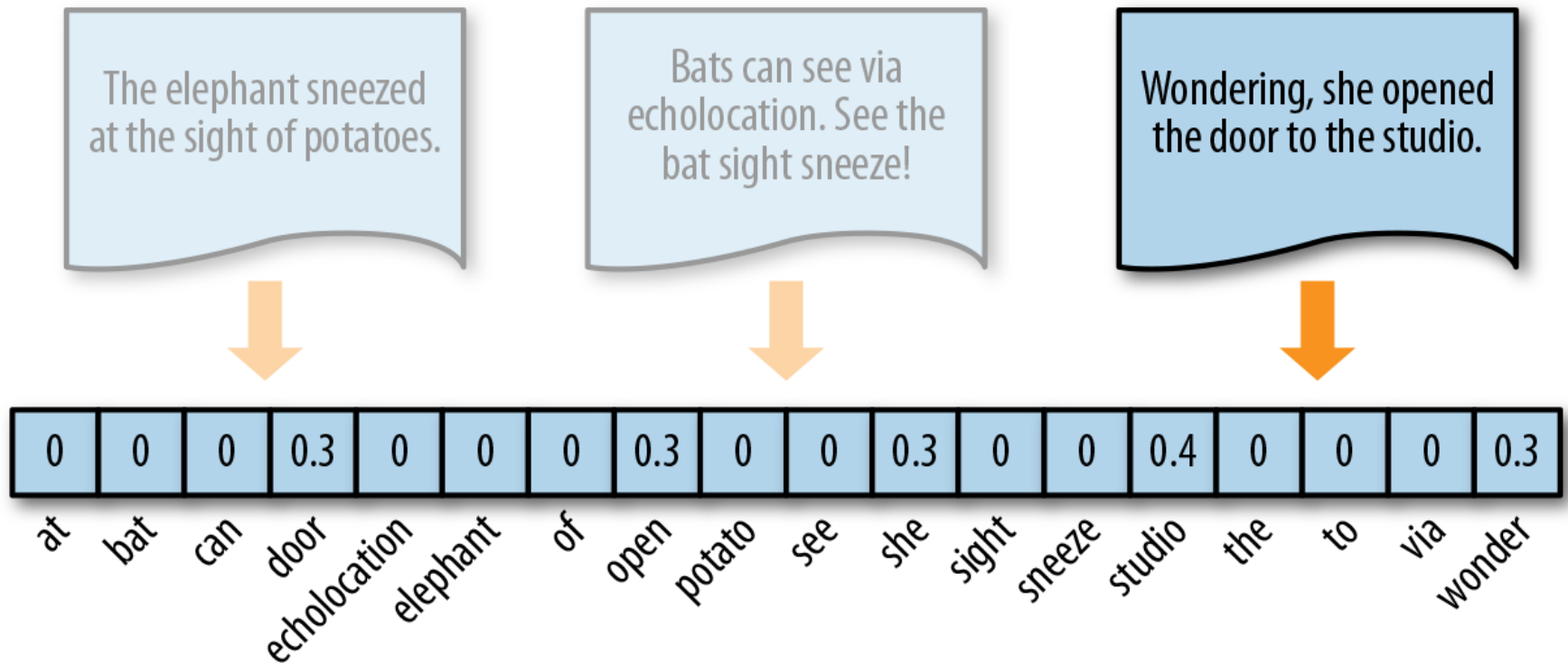| Vectorization Method | Function | Good For | Considerations |
|---|---|---|---|
| Frequency | Counts term frequencies | Bayesian models | Most frequent words not always most informative |
| One-Hot Encoding | Binarizes term occurrence (0, 1) | Neural networks | All words equidistant, so normalization 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
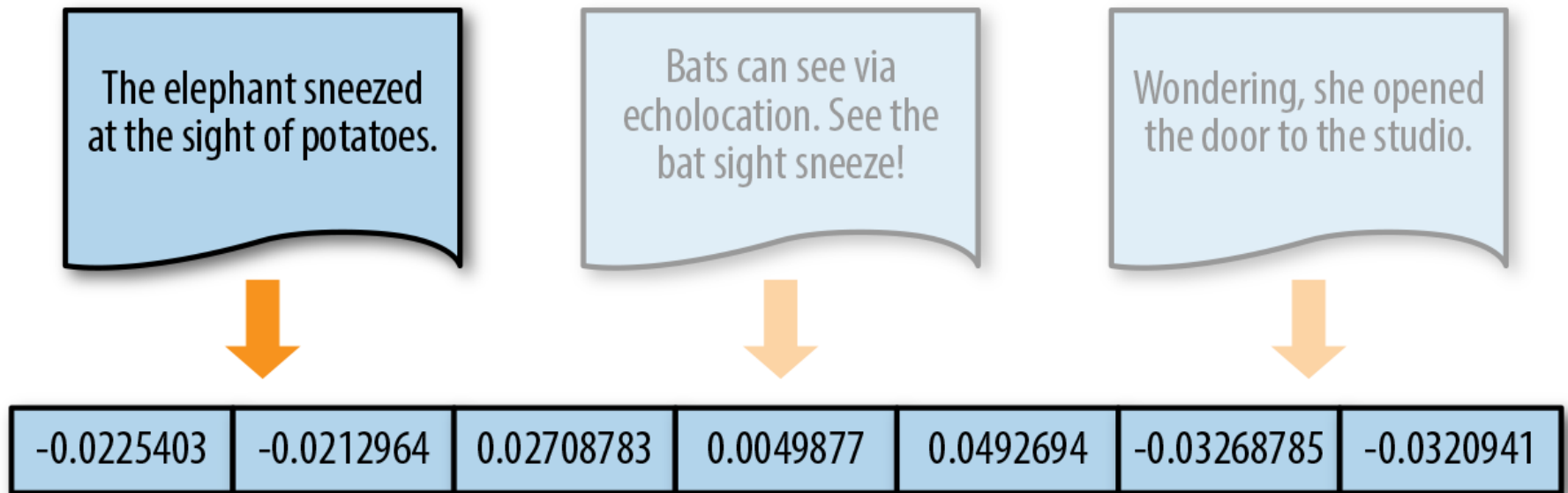
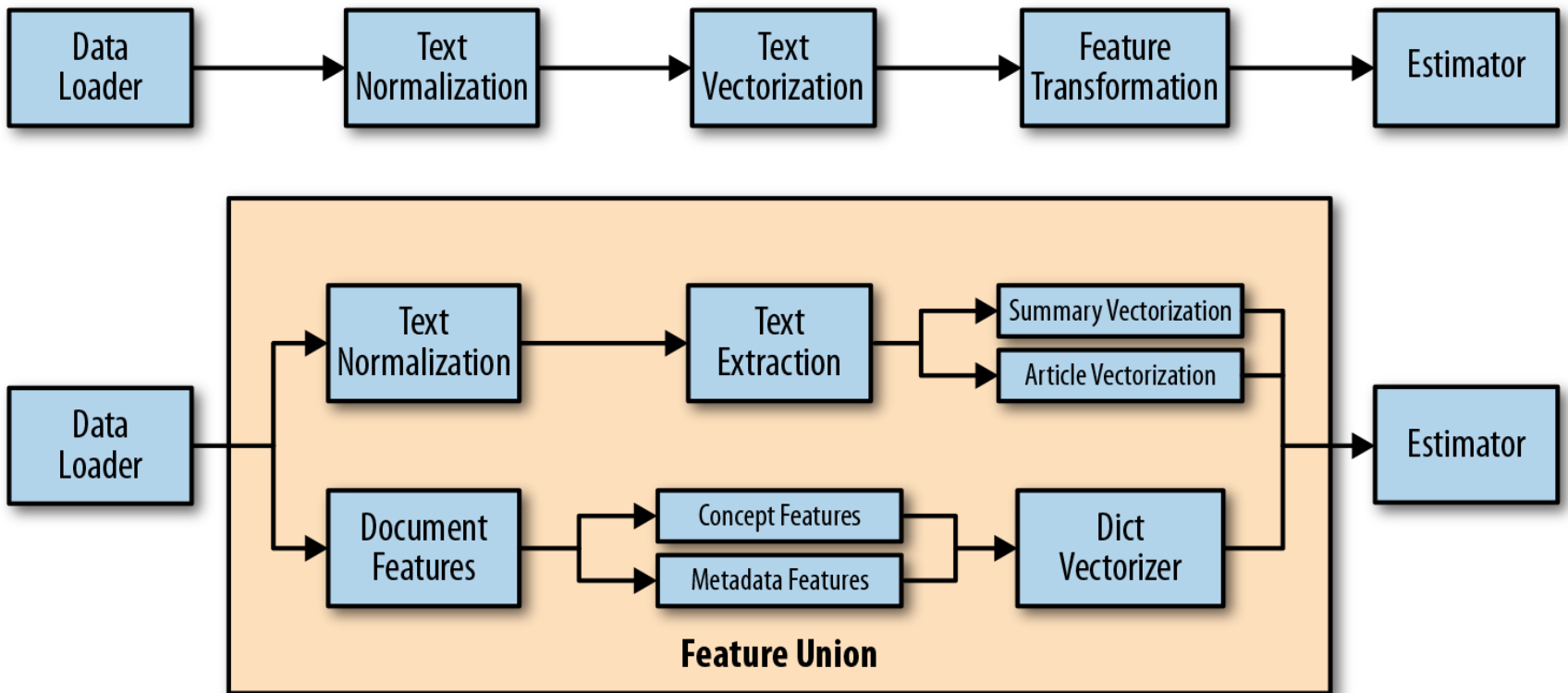# Token Frequency as Vector Encoding

# One-hot Encoding
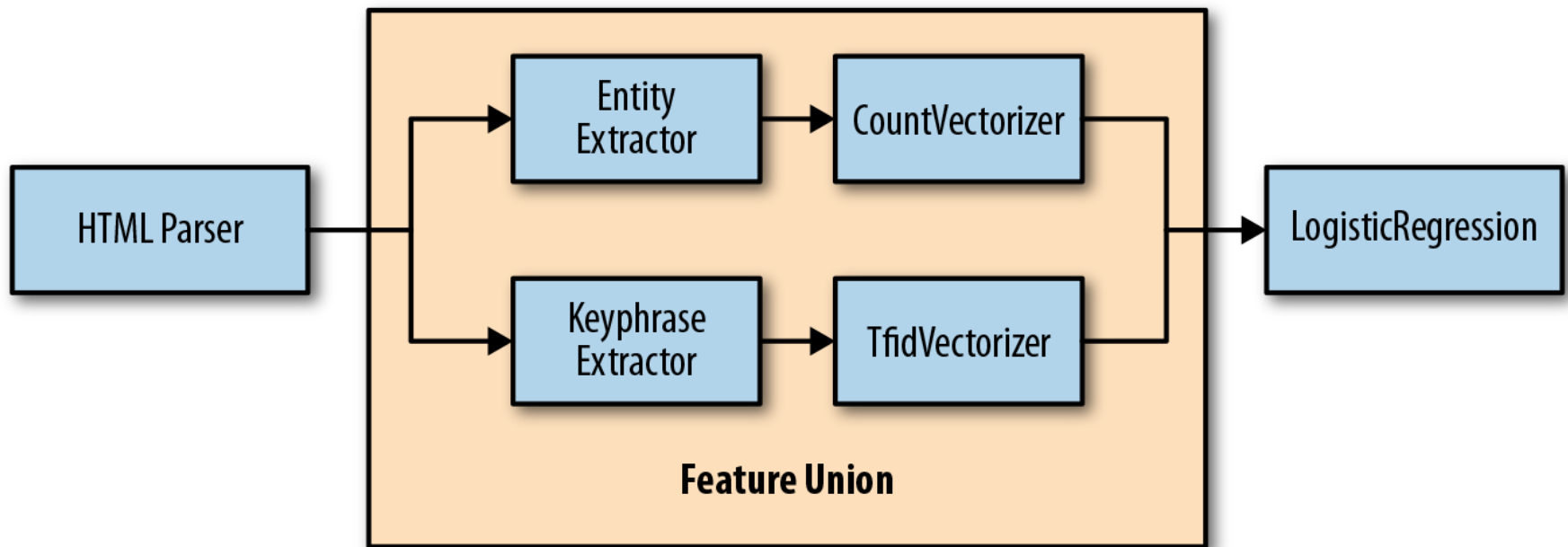
# TF-IDF Encoding

# 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

```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
```

```
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/aintpupython101

```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/aintpupython101

```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/aintpupython101

```python
# 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/aintpupython101

```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/aintpupython101

```
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 |

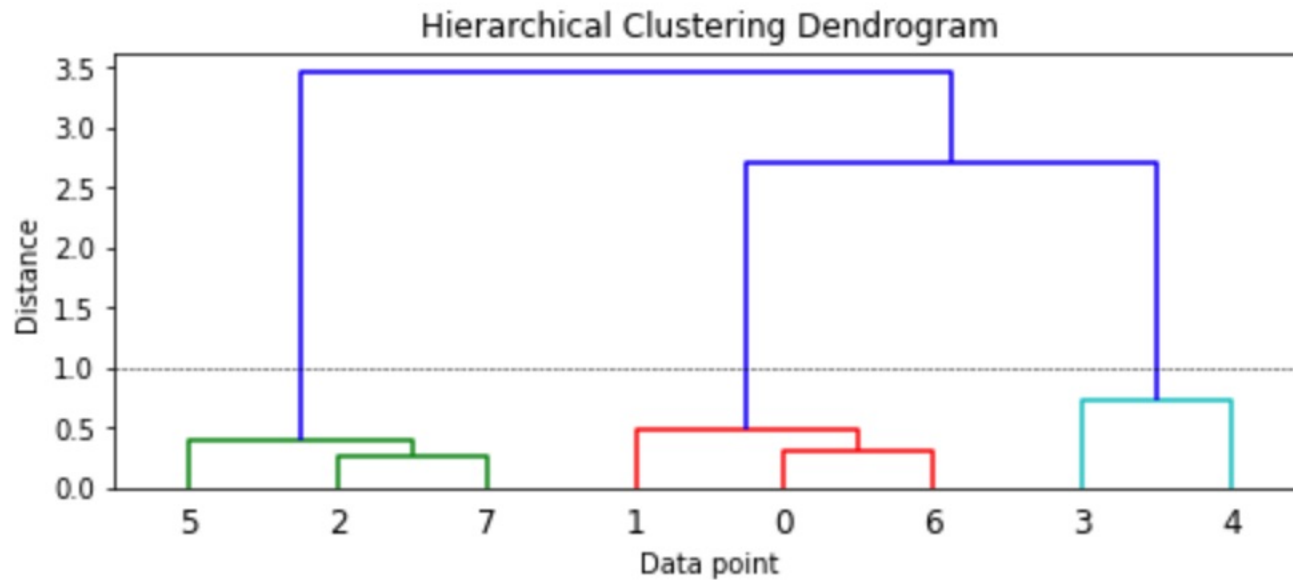https://tinyurl.com/aintpupython101

```
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 |

```
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>

```
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)
```

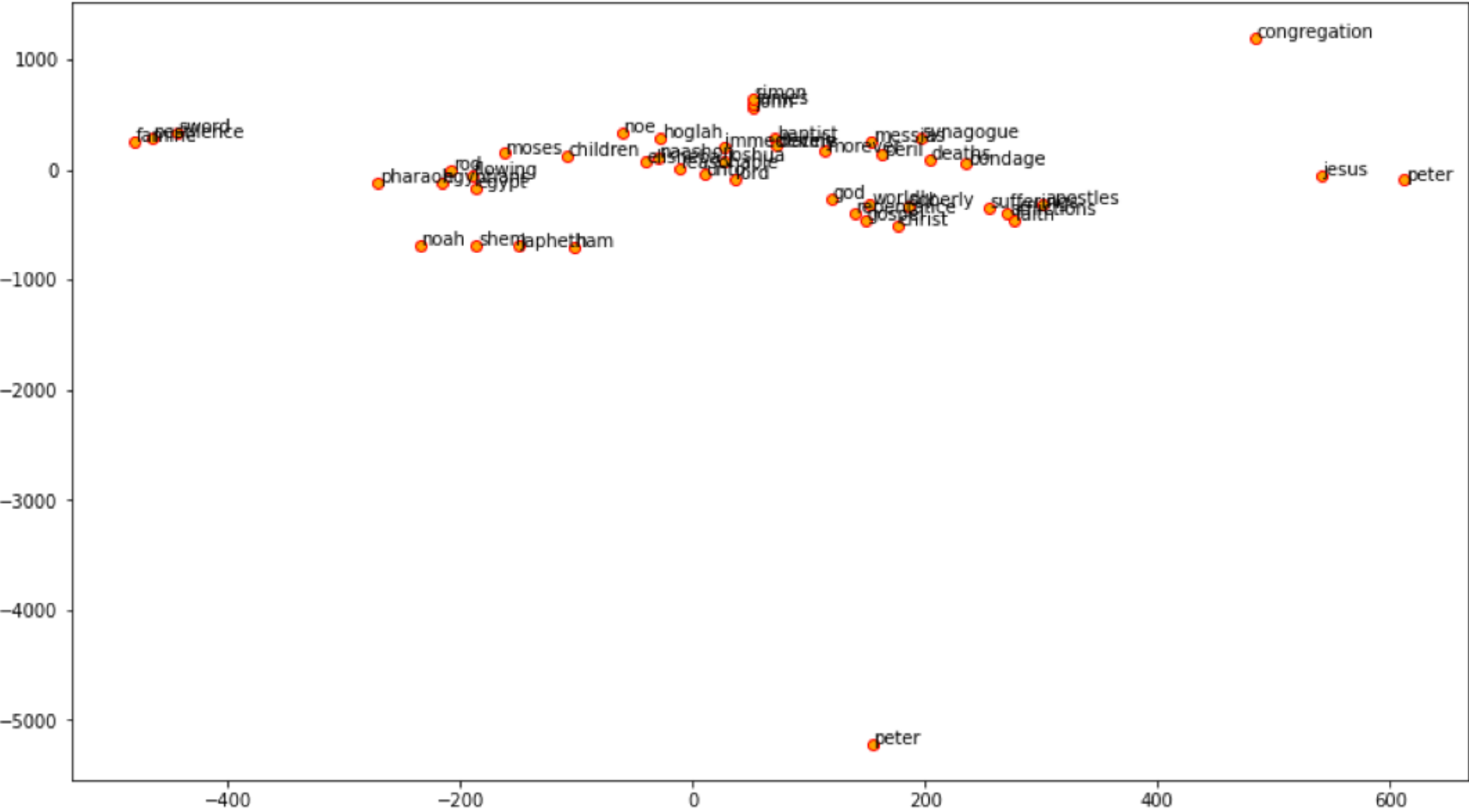|   | Document | Category | ClusterLabel |
|---|----------|----------|--------------|
| 0 | The sky is blue and beautiful. | weather | 2 |
| 1 | Love this blue and beautiful sky! | weather | 2 |
| 2 | The quick brown fox jumps over the lazy dog. | animals | 1 |
| 3 | A king's breakfast has sausages, 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 |

```
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/aintpupython101

```
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/aintpupython101

```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
```

```
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/aintpupython101

```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

41

# 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.

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 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

# BERT Token-level tasks



(c) Question Answering Tasks: SQuAD v1.1

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

45

# 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

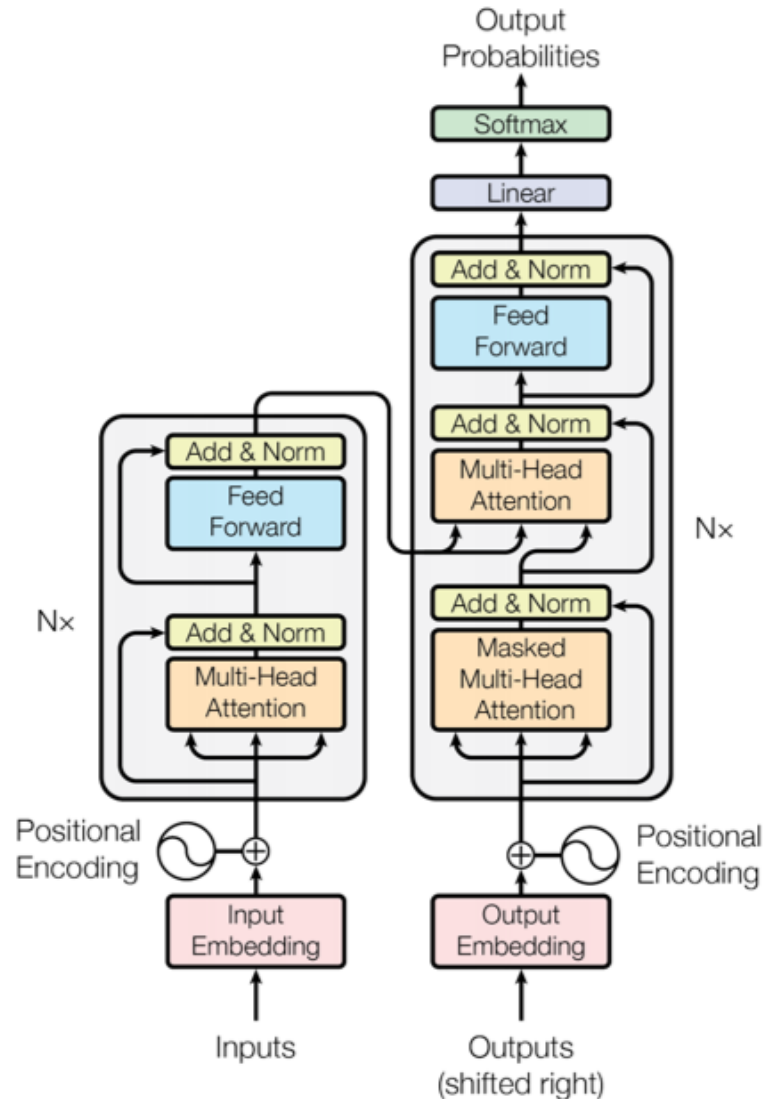Input Question: Is this person dancing ?

Fixed length question vector encoded by the LSTM

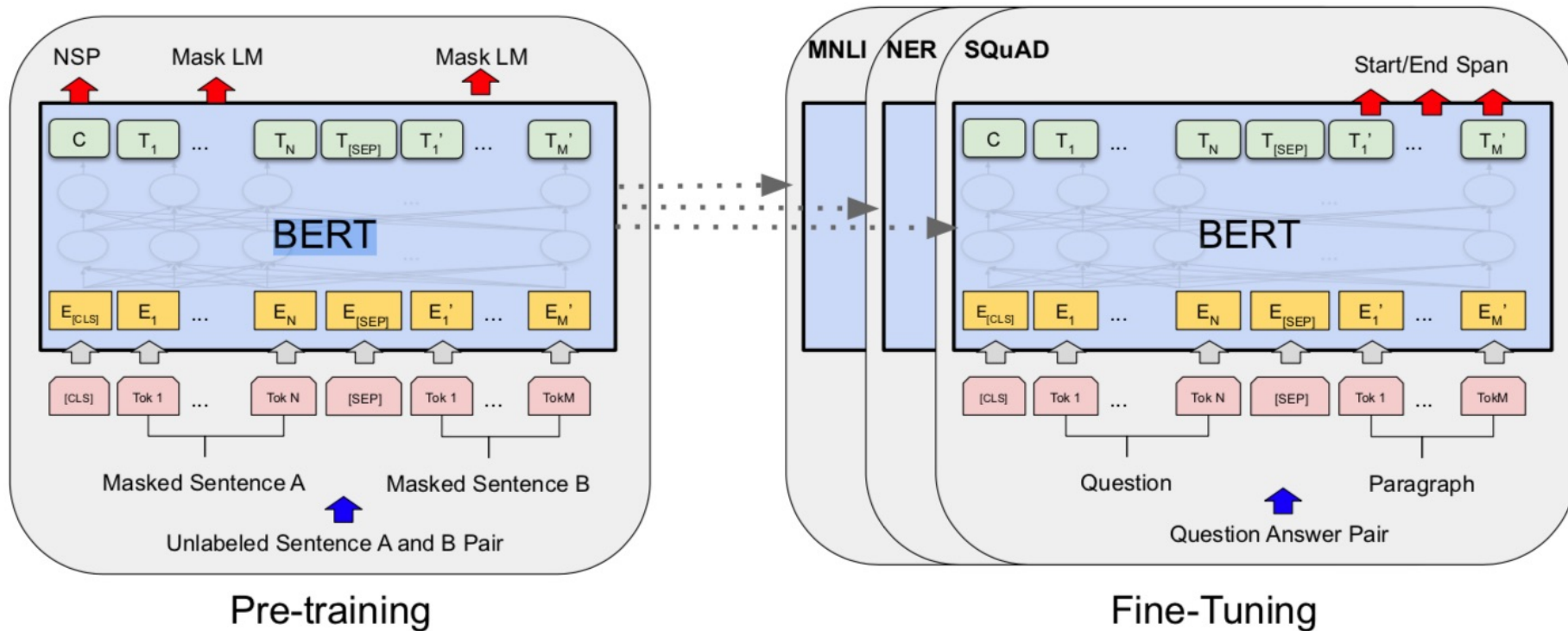# Transformer (Attention is All You Need)
## (Vaswani et al., 2017)

# 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

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

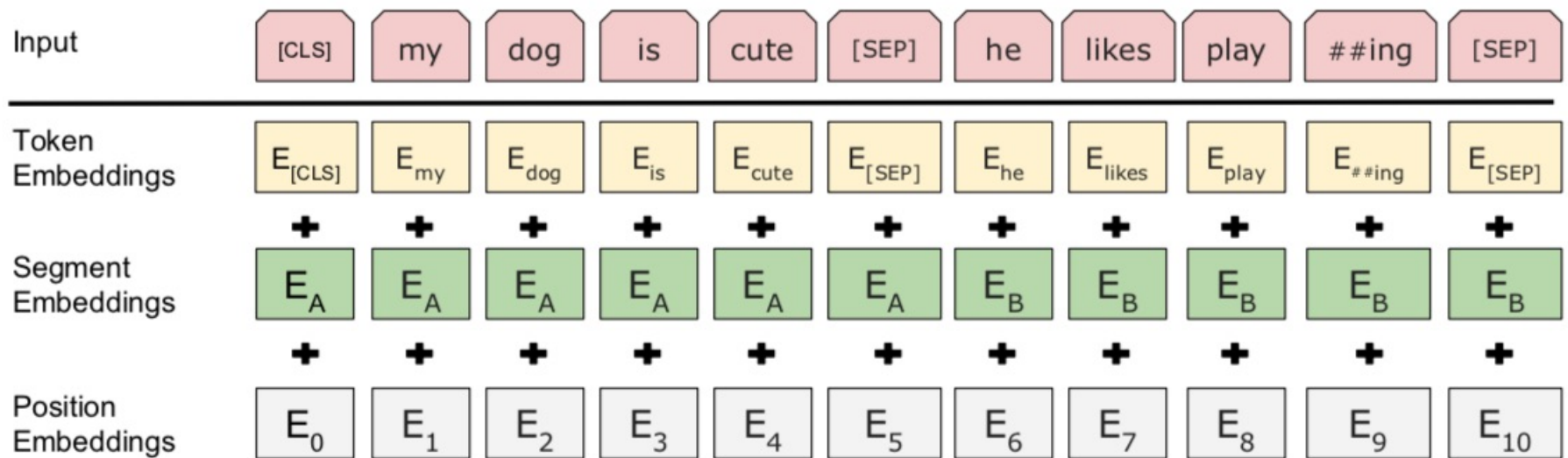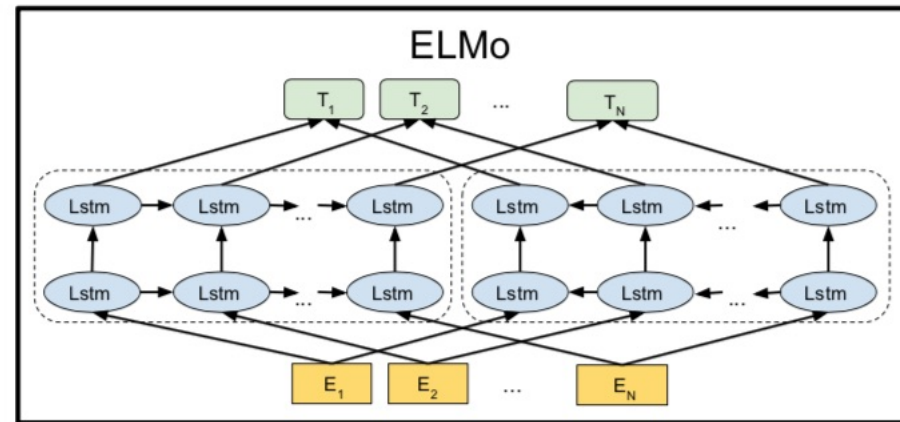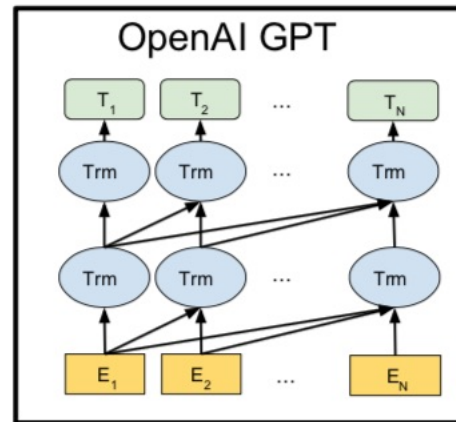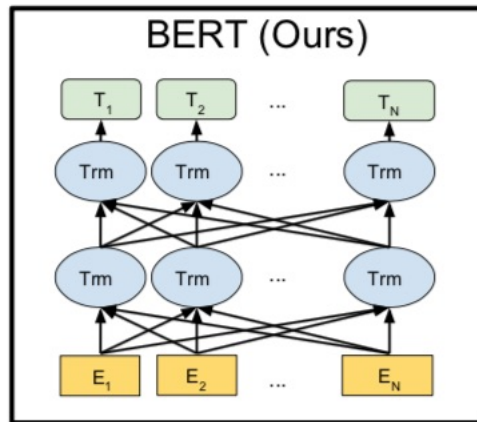## BERT (Bidirectional Encoder Representations from Transformers)

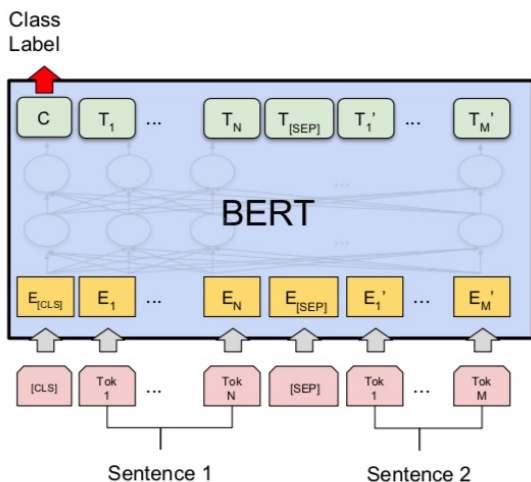## BERT input representation

# 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.

# Pre-trained Language Model (PLM)



Source: https://github.com/thunlp/PLMpapers

# Turing Natural Language Generation (T-NLG)

# 🤗 Transformers 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 |

# 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/