# Natural Language Processing

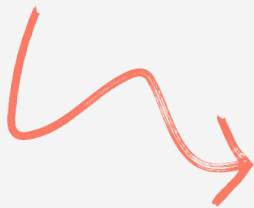Data Science Unit 4

# Before we start…

➔ Make sure you are comfortable

➔ Have water and maybe a strong coffee handy

➔ If you need a break… take it!

➔ If you need a stretch – please go ahead!

➔ Please mute yourselves if you are not talking

➔ Have your video on at all times

…and let's get started!

# In this session we will...

→ **Define** Natural Language Processing

→ **Explore** CountVectorizer

→ **Investigate** texts with Sentiment Analysis

# NLP

# NLP

## Analysis

NLP techniques provide tools to allow us to understand and analyse large amounts of text. For example:

➔ Analyse the positivity/negativity of comments on different websites.

➔ Extract keywords from meeting notes and visualise how meeting topics change over time.

## Vectorising for machine learning

When building a machine learning model, we typically must transform our data into numeric features. This process of transforming non-numeric data such as natural language into numeric features is called vectorisation. For example:

➔ Understanding related words. Using stemming, NLP lets us know that "swim", "swims", and "swimming" all refer to the same base word. This allows us to reduce the number of features used in our model.

➔ Identifying important and unique words. Using TF-IDF (term frequency-inverse document frequency), we can identify which words are most likely to be meaningful in a document.

# What is it?

Using computers to process (analyse, understand, generate) natural human languages.

Making sense of human knowledge stored as unstructured text.

Building probabilistic models using data about a language.

# Examples

# High Level

**Chatbots:**
Understand natural language from the user and return intelligent responses.

➔ Api.ai

**Text simplification:**
Preserve the meaning of text, but simplify the grammar and vocabulary.

➔ Rewordify

➔ Simple English Wikipedia

**Information retrieval:**
Find relevant results and similar results.

➔ Google

**Predictive text input:**
Faster or easier typing.

➔ Phrase completion application

➔ A much better application

**Information extraction:**
Structured information from unstructured documents.

➔ Events from Gmail

**Sentiment analysis:**
Attitude of speaker.

➔ Hater News

**Machine translation:**
One language to another.

➔ Google Translate

**Speech recognition and generation:**
Speech-to-text, text-to-speech.

➔ Google's Web Speech API demo

➔ Vocalware Text-to-Speech demo

# Low Level

**Tokenization:**
Breaking text into tokens
(words, sentences, n-grams)

**TF-IDF:**
word importance

**Spelling correction:**
"New Yrok City"

**Language detection:**
"translate this page"

**Stop-word removal:**
a/an/the

**Part-of-speech tagging:**
noun/verb/adjective

**Word sense
disambiguation:**
"buy a mouse"

**Machine learning:**
specialized models that
work well with text

**Stemming and
lemmatization:**
root word

**Named entity
recognition:**
person/organization
/location

**Segmentation:**
"New York City subway"

# Challenges

**Ambiguity**:
➔ Hospitals Are Sued by 7 Foot Doctors
➔ Juvenile Court to Try Shooting Defendant
➔ Local High School Dropouts Cut in Half

**Non-standard English:**
text messages

**Idioms:**
"throw in the towel"

**Newly coined words:**
"retweet"

**Tricky entity names:**
"Where is A Bug's Life playing?"

**World knowledge:**
"Mary and Sue are sisters", "Mary and Sue are mothers"
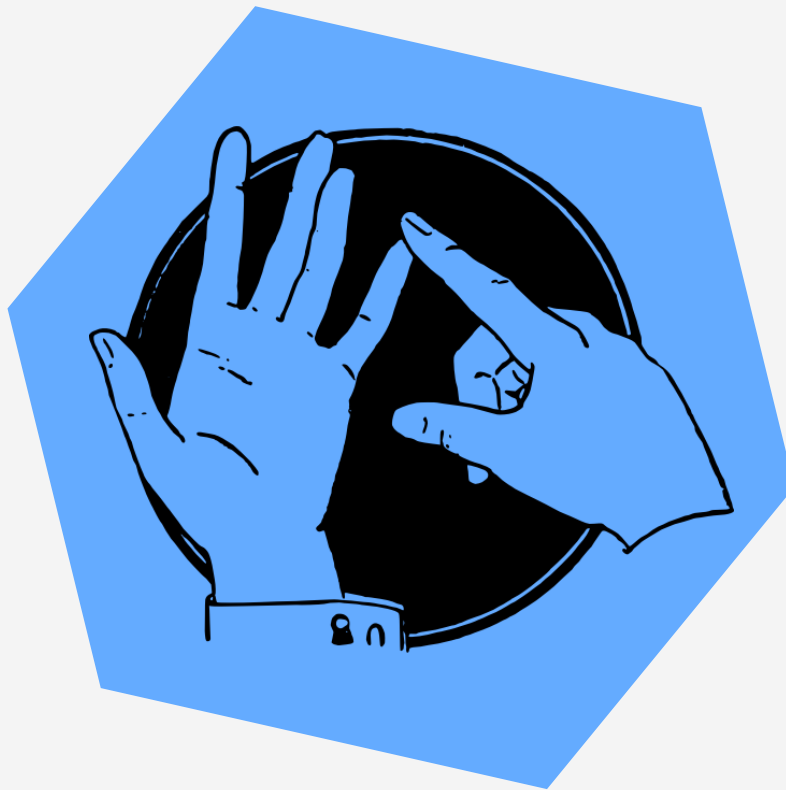
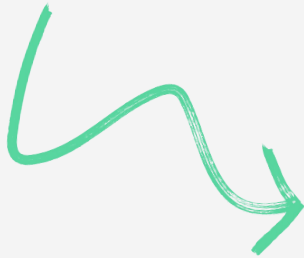# Text Classification

# Text Classification

# Text Classification

# Text Classification

# Text Classification

# Text Classification

```python
from sklearn.feature_extraction.text import CountVectorizer
# Use CountVectorizer to create document-term matrices from X_train and X_test.
vect = CountVectorizer()
X_train_dtm = vect.fit_transform(X_train)
X_test_dtm = vect.transform(X_test)
```

# Text Classification

multiverse

Let's Practice

# NGrams

# NGrams

```
my cat is awesome
Unigrams (1-grams): 'my', 'cat', 'is', 'awesome'
Bigrams (2-grams): 'my cat', 'cat is', 'is awesome'
Trigrams (3-grams): 'my cat is', 'cat is awesome'
4-grams: 'my cat is awesome'
```

# NGrams

```python
# Include 1-grams and 2-grams.
vect = CountVectorizer(ngram_range=(1, 2))
X_train_dtm = vect.fit_transform(X_train)
X_train_dtm.shape
```
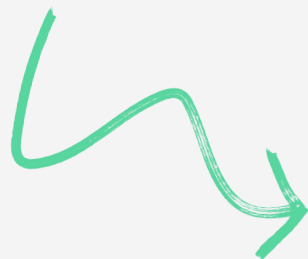
multiverse

Let's Practice

# Stop Words and Other Options

# Stop Words

# Stop Words

```python
# Remove English stop words.
vect = CountVectorizer(stop_words='english')
```

```
vect = CountVectorizer(max_df=0.7)
```

```
vect = CountVectorizer(max_features=100)
```

multiverse

Let's Practice

# Stemming

# Lemmatization

| Lemmatization | Stemming |
|---|---|
| shouted – shout | badly – bad |
| best – good | computing – comput |
| better – good | computed – comput |
| good – good | wipes – wip |
| wiping – wipe | wiped – wip |
| hidden – hide | wiping – wip |

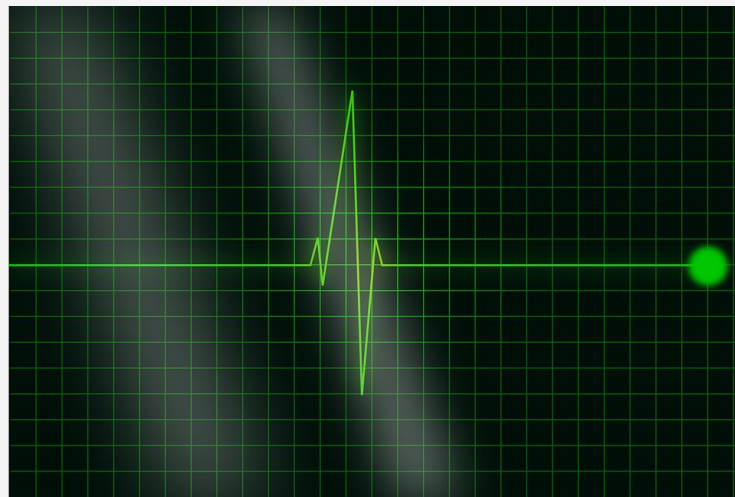multiverse

Let's Practice

# TF-IDF

```python
# TfidfVectorizer
vect = TfidfVectorizer()
pd.DataFrame(vect.fit_transform(simple_train).toarray(), columns=vect.get_feature_names())
```

|   | cab | call | me | please | tonight | you |
|---|-----|------|-----|--------|---------|-----|
| 0 | 0.000000 | 0.385372 | 0.000000 | 0.000000 | 0.652491 | 0.652491 |
| 1 | 0.720333 | 0.425441 | 0.547832 | 0.000000 | 0.000000 | 0.000000 |
| 2 | 0.000000 | 0.266075 | 0.342620 | 0.901008 | 0.000000 | 0.000000 |

# Let's Practice

# I Appreciate
# the Sentiment

# I Appreciate the Sentiment