# MOVIES

Analyzing subtitles

Samer Ghosn

# MOVIE SCRIPTS

- Analyzing the text in a movie and classify it accordingly.

- Movies scripts contain words that reveal their genre.

- Analyzing the frequency of these words in the subtitle files.

- If the words "Space", "Ship" appear frequently then the movie is "Sci-fi"

Let's Analyze some movie scripts

help let jarvis three head put mean gotta great captain
place iron father made dad old course lot long wanted guys
i'd hold hulk whoa playing kill stay bad big around even happy thought
tony saw happened listen peter wanna believe maybe say
people power coming find world look love else thought
trying thank kid job ever leave best always guy work thanks nothing
first suit name parker every home kind hell next
night talk things understand feel harry anything two uh grunting
stark seen done
never please keep told move ready god mr still may life better
sure call god remember nice actually
stop thing give
sorry spider
grunts

# More Examples

# STARWARS MOVIES

- may the **Force** be with you.

- The **Empire** is still out there!

-  Commander **Skywalker**

- Our plan, **captain**?

- **Luke Skywalker**, **Jedi knight**.

# MARVEL MOVIES

- Can we hold them? - They're the Avengers!!

- I'm Captain America.

- who is also known as Iron Man

- the strongest substance in the universe..

- Xandarian outposts throughout the galaxy

# LORD OF THE RINGS MOVIES

- What did you tell him about __Frodo__ and the __Ring__?

- An __Elf__ will go __underground__ where a __Dwarf__ dare not?

- But we swore to serve the __master__ of the __precious__.

- They think we have the __Ring__.

# HORROR MOVIES

- I'm not <u>afraid</u> of you!

- <u>Run</u>, <u>run</u>, <u>run</u>.

- Someone <u>save</u> me.

- <u>Please</u>, I need your <u>help</u>.

- You seem <u>sad</u>.

MORE ABOUT THE PROCESS

Analyzing subtitles

# LEARNING

- Several English subtitles are downloaded from the internet (public use).

- Subtitles from several movie genres are gathered in several folders.

- Cleaning was done on 2 phases, phase 1 is removing punctuation, and unneeded characters. Phase 2 is removing stopwords.

- Flatten subtitle file and label it. Each subtitle on a line and add genre next to it.

- Process each genre separately and get an idea on most used words.

# CLEANING DATA

- Example of a subtitle before cleaning:

```
1367
02:11:09,600 --> 02:11:12,355
BILBO: Gandalf?
GANDALF: Bilbo Baggins.

1368
02:11:12,355 --> 02:11:16,038
BILBO: My dear Gandalf! Ha, ha!
GANDALF: It's good to see you.
```

- Remove Numbers, columns and dots.

- Remove Stopwords. Stopwords are most commonly used words (such as "the", "a", "an", "in").

- Write a python program to clean data, remove new lines and label it.

# PROCESSING DATA

- Code is written in python.

- NLTK library is used (Natural Language Toolkit)

```python
df = pd.read_csv(file_name, sep="|", header=None,encoding="ISO-8859-1")
df.columns = ['subtitle','category']
df_list.append(df)
```

```python
df = pd.concat(df_list, axis=0, ignore_index=True) #axis = 0 concatenate row wise
```

- TweetTokenizer is used instead of word_tokenize to split data.

- FreqDist is used to get the frequency of the words.

```python
frequency_dist = nltk.FreqDist(reviews_text)
sorted(frequency_dist,key=frequency_dist.__getitem__, reverse=True)
```
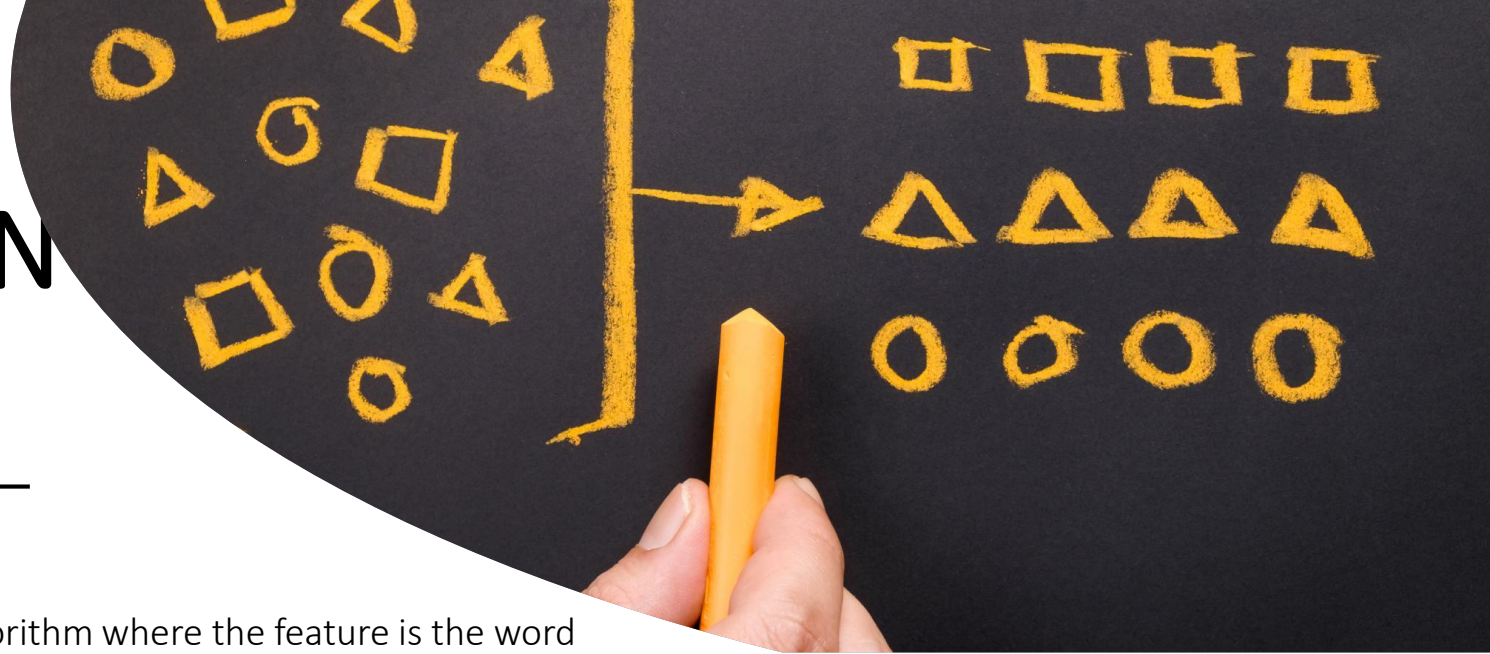
# PROCESSING DATA

- wordcloud

```python
from wordcloud import WordCloud
import matplotlib.pyplot as plt
text = df.subtitle.values
wordcloud = WordCloud(
    font_path="verdana",
    width=6400,
    height=3200,
    max_words=120,
    background_color="white",
    stopwords=stop_words
).generate_from_frequencies(frequency_dist)
#plt.figure(figsize=(20,10))
plt.imshow(wordcloud)
plt.axis("off")
plt.tight_layout(pad=0)
plt.savefig(directory + '\\movie_analysis.png')
plt.show()
plt.close()
```

# DATA CLASSIFICATION



- 2 Approaches to classify data:
  - Consider all columns as features and apply Knn algorithm where the feature is the word and the value is number of occurrences.
  - Using Deep learning to classify the movies. Tensor flow keras library
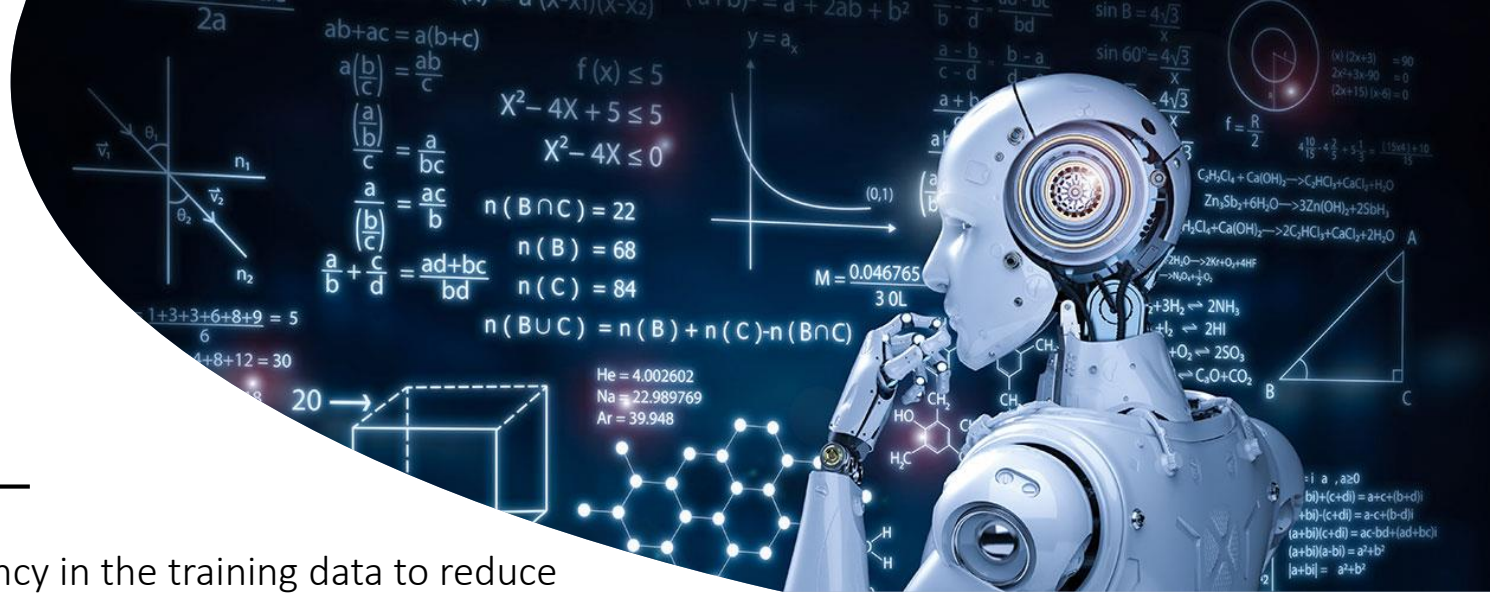
```python
from keras.layers import Dense

# Keras layers can be called on TensorFlow tensors:
x = Dense(128, activation='relu')(img)  # fully-connected layer with 128 units and ReLU activation
x = Dense(128, activation='relu')(x)
preds = Dense(10, activation='softmax')(x)  # output layer with 10 units and a softmax activation
```
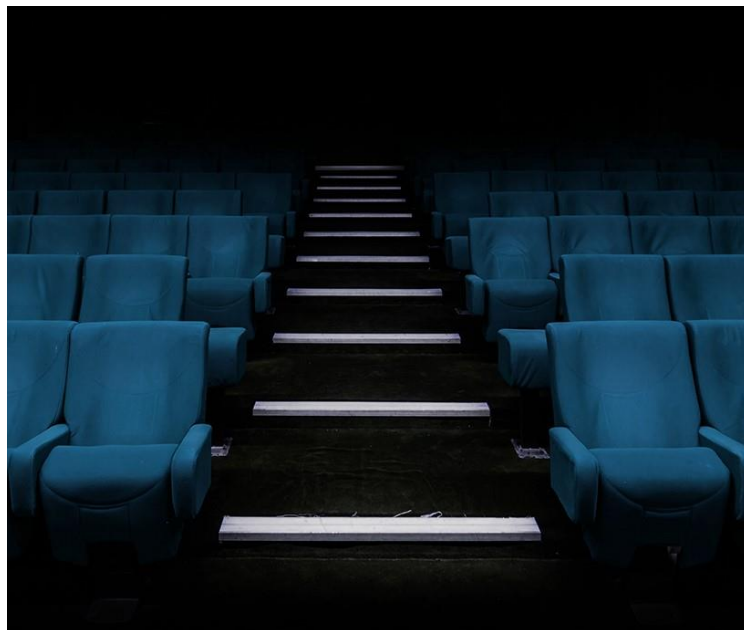
```python
# Run training loop
with sess.as_default():
    for i in range(100):
        batch = mnist_data.train.next_batch(50)
        train_step.run(feed_dict={img: batch[0],
                                  labels: batch[1]})
```

- https://blog.keras.io/keras-as-a-simplified-interface-to-tensorflow-tutorial.html

# IMPROVEMENTS



- Remove words that appear at the same frequency in the training data to reduce features.

- Apply Dimensionality reduction on the dataset, for example:
  - "Happy", "Joy"
  - "Run", "Sprint", "Rush", "Running", "rushing", "rushed"
  - "afraid", "fear", "scared"

- Ignore words that appear frequently in some movies, for example:
  The name "Rachel" in a horror movie has no importance. However, the name "Stark" as tony stark is important toward classifying a movie as a Marvel.

- Some action verbs are important not to ignore, like "run", "go", "move". If they appear frequently, this means that it's an action movie.

- Capture special characters ♪

- Maybe, Convert features to numbers.

# END OR QUESTIONS