# Capstone Project
## Topic Modeling on News Articles

# Content

- **Problem Statement**
- **Data Summary**
- **Data Preprocessing**
- **Models**
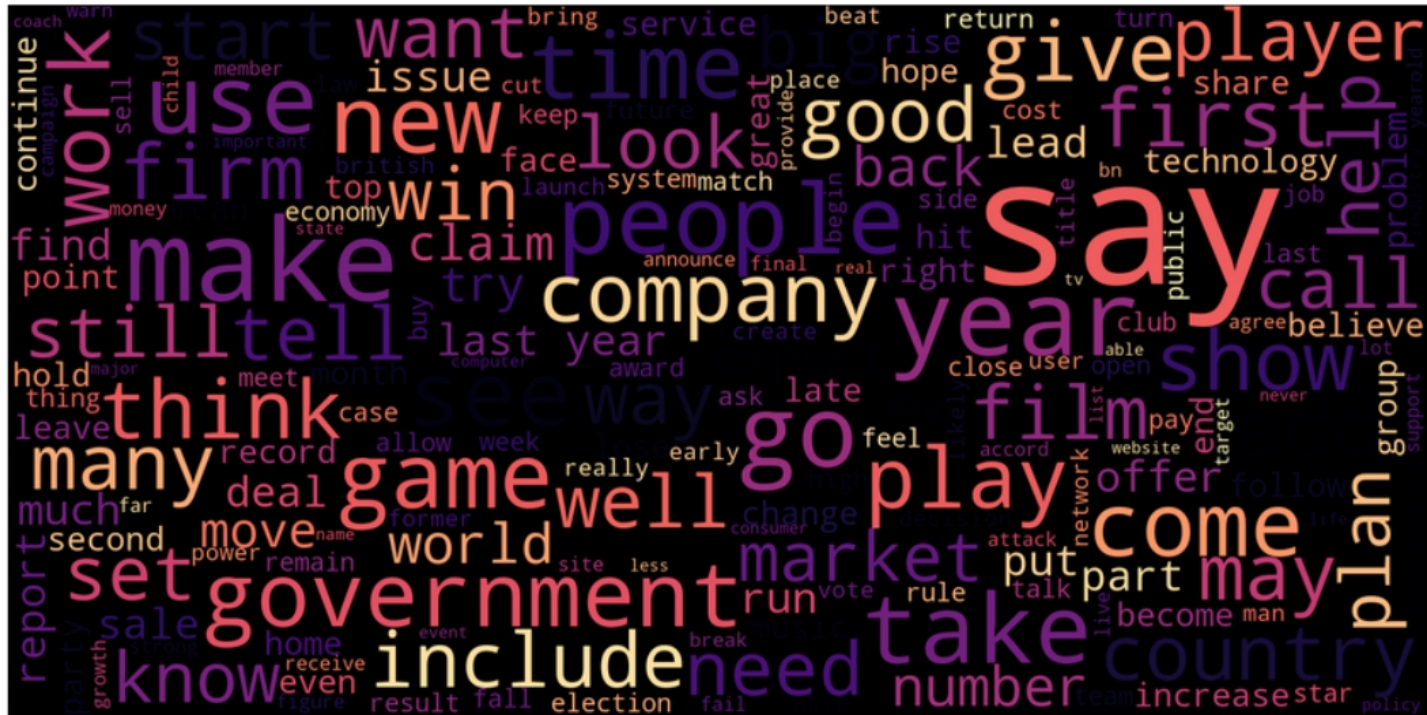- **Challenges**
- **Conclusions**
- **Q&A**

# Problem Statement

- **Identify major themes/topics across a collection of BBC news articles using Latent Dirichlet Allocation (LDA) and Latent Semantic Analysis (LSA).**
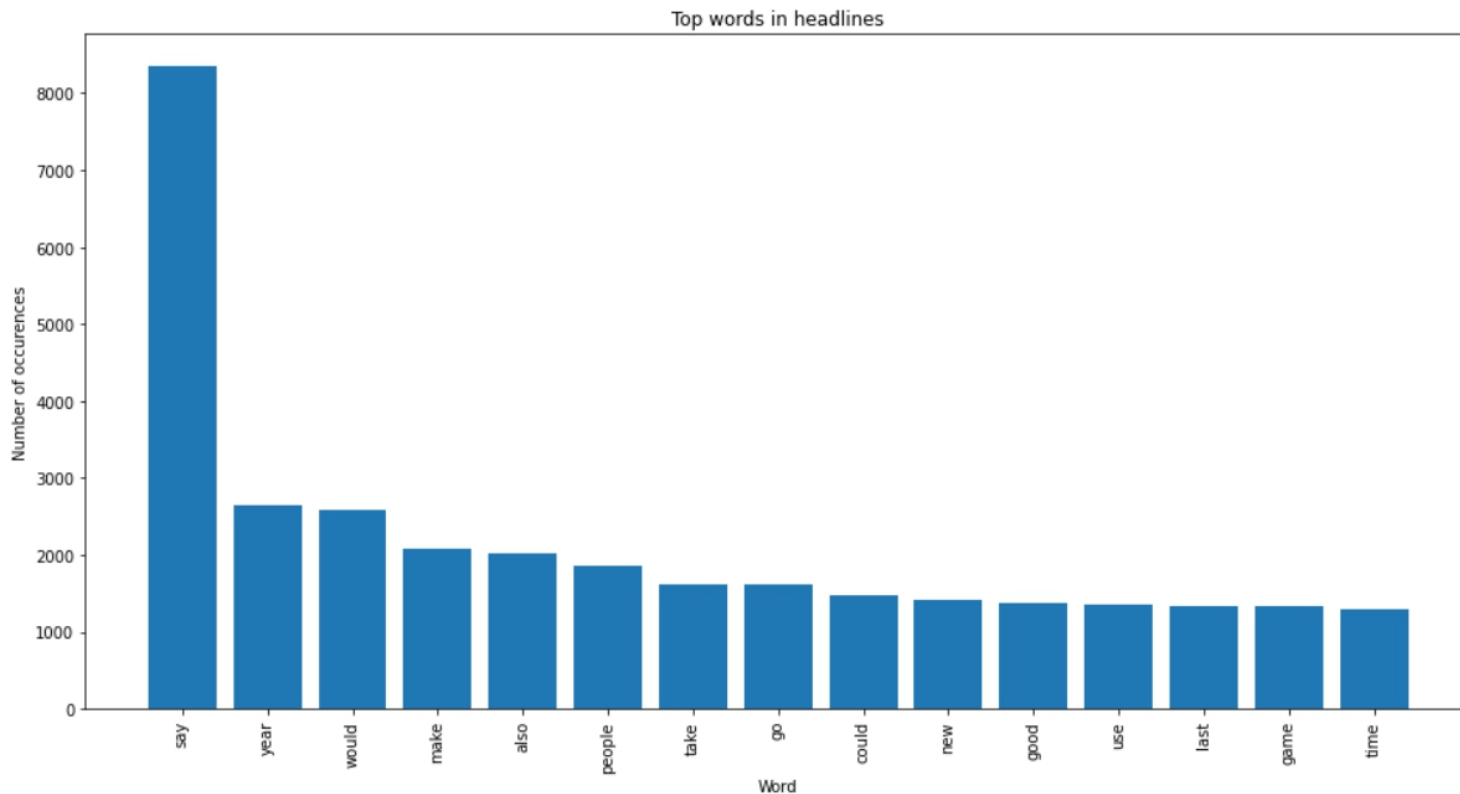
# Data Preprocessing

- **Remove duplicate values**
- **Remove Html tags**
- **Remove URLs**
- **Remove punctuation**
- **Remove numbers**
- **Remove small length words**
- **Remove stop words**
- **Lemmatization**

# Word Cloud

# Frequent words



Top words in headlines

# Countplot of News Types

# Length of Documents



Length of Documents

# Word Count of Documents

# Word Cloud for Business

# Word Cloud for Entertainment

# Word Cloud for Politics

# Word Cloud for Sport

# Word Cloud for Tech

# Latent Dirichilet Allocation



Intertopic Distance Map (via multidimensional scaling)

Top-30 Most Salient Terms[1]

Marginal topic distribution
2%
5%
10%

Overall term frequency
Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)

# Similar Words

**Word : Film**

| Good |
|---|
| Award |
| Star |
| Play |
| Actor |

**Word:Price**

| Rate |
|---|
| Rise |
| High |
| Profit |
| Fall |

**Word:Stock**

| Gross |
|---|
| Soar |
| Copy |
| Late |
| Monthly |

AI

# Latent Semantic Analysis

# Challenges

- **Text Preprocessing**
- **Limited visualization techniques**

# Conclusion

- **In Latent Dirichilet Allocation(LDA) with TF-IDF vectorizer we find best clustering for our news article dataset .**

- **Using of genism we can find a similar words .**

- **As a future we can implement topic modeling using different techniques like neural network.**

Q & A

# Thank You