SUBJECT: ISM 6359
DATA MINING


TOPIC: FAKE NEWS DETECTION


SHEETAL MURALI

# BUSINESS REASON:
# FAKE NEWS DETECTION
# TOOL: RAPID MINER

➢ In today's world of Social Media, one of the major problem we face is misinformation.

➢ Millions of articles are being published every minute-the scope of human manual detection of a real or fake news is not feasible.

➢ The model created will test the unseen data, and accordingly detect fake articles and can be used and integrated with any system for future use.

# DATA STRUCTURE

- The data I have collected is from Kaggle.com. This data consists of the following information:

- title: title of the article – 21724 Unique Values (Sampled it down to 5000)

- news_url: URL of the article

- source domain: web domain where article was posted.

- tweet_num: number of retweets for this article.

- real: label column, where 1 is 'real' and 0 is 'fake'.

# TEXT PREPARATION

- Loading the data

- Set Role – Label "Real"

- Select Attributes- Remove irrelevant attribute "URL"

- Numerical to Binomial – Label

- Filter Examples – Missing values

- Sample – 5000 random values

- Normalize – number of tweets ranging from (from '0' to '29060')

- Nominal to Text

| Name | | Type | Missing | Statistics | | | Filter (4 / 4 attributes): | Search for Attributes |
|---|---|---|---|---|---|---|---|---|
| Label **real** | | Binominal | 0 | Negative<br>false | Positive<br>true | Values<br>true (3750), false (1225) | | |
| **title** | | Text | 0 | Least<br>ï»¿Lily [...] Feel' (0) | Most<br>Connecti [...] News (7) | Values<br>Connecti [...] ough News (7), A Comple [...] ations | | |
| **source_domain** | | Text | 0 | Least<br>zimbabwe-today.com (0) | Most<br>people.com (379) | Values<br>people.com (379), www.dailymail.co.uk (217), ...[ | | |
| **tweet_num** | | Real | 0 | Min<br>-0.380 | Max<br>19.164 | Average<br>-0 | | |

# PROCESS DOCUMENT FROM DATA

- Tokenization

- Transform Cases

- Filter Stopwords

- Filter Token by length

- Stemming-Porter
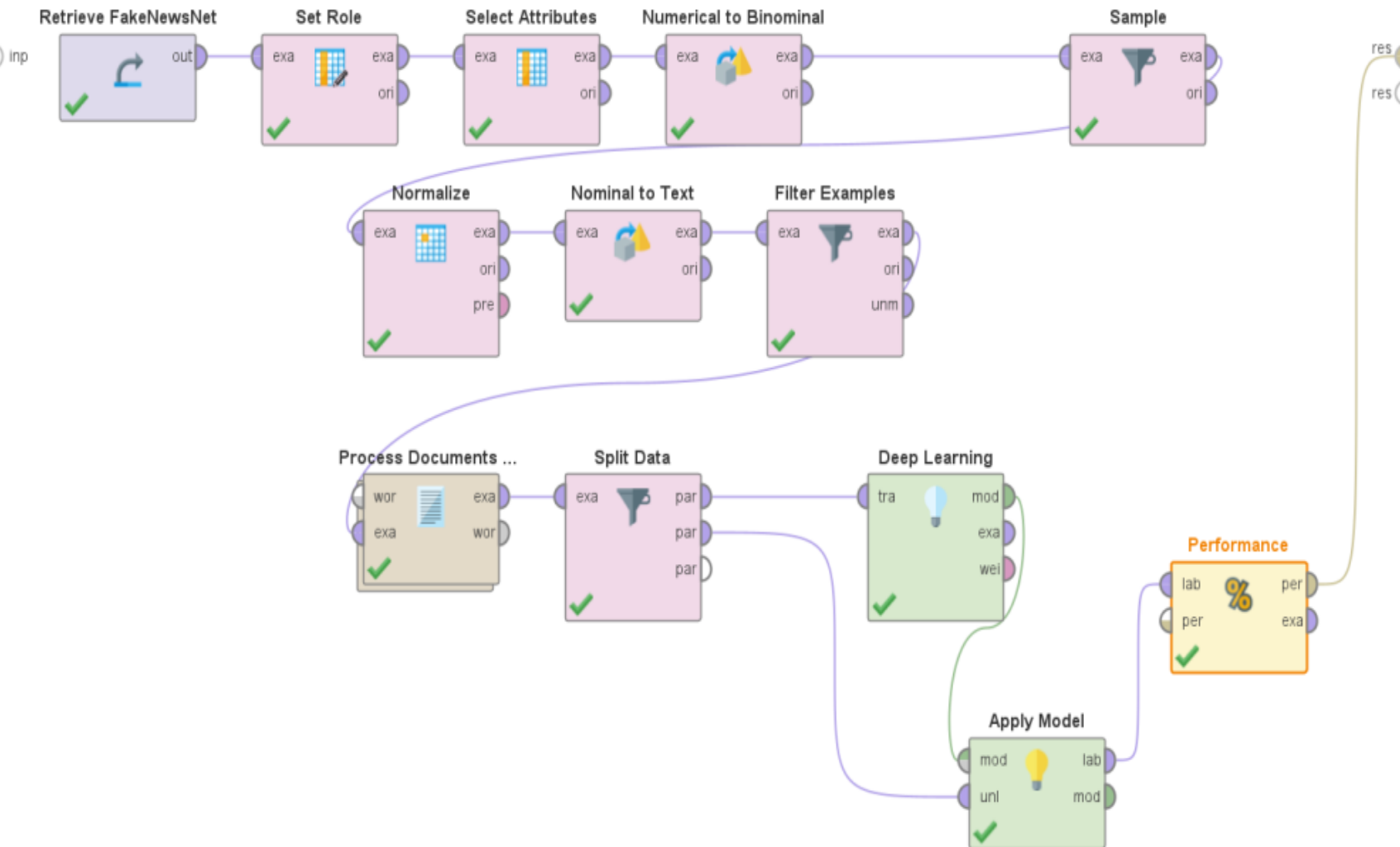
SPLIT DATA: 70% Training set

30% Test set

# ALGORITHMS USED

- Naïve Bayes

- Decision Tree

- Random forest

- Deep Learning

- SVM

# ANALYSIS OF THE DATA

| Operators | Attribute count |
|---|---|
| Tokenization | 12282 |
| Transform cases | 10079 |
| Filter Stopwords | 9804 |
| Filter tokens by length | 9171 |
| Stem | 7538 |

# Accuracy for different Algorithms

| Model | Split Data 70%-30% |
|---|---|
| Decision Tree | 78.19% |
| Random Forest | 75.38% |
| Sequential Minimal Optimization (SVM) | 81.51% |
| Naïve Bayes | 72.28% |
| Deep Learning | 81.57% |

# PERFORMANCE BASED ON WEIGHTS

| | |
|---|---|
| Binary Term Occurrences | 77.39% |
| Term Occurrences | 79.2% |
| Term-Frequency | 75.88% |
| TF-IDF | 81.57% |

## DATA SOURCE:

I found the dataset used to build this model on Kaggle.com

https://www.kaggle.com/datasets/algord/fake-news?select=FakeNewsNet.csv

# 3 W'S

1) **What went well?**

Topics taught in class could be easily related to the business problem I was solving. Video lectures helped me simultaneously work on RapidMiner as we kept learning new features.

2) **What did not go well?**

Initially I started the project with 20000+ documents, but text mining operator did run because of space constraints on the laptop due to which I had to reduce the data to 5000 examples by using random sampling

3) **What would you use Differently Next Time?**

Since the local host does not have enough memory to run Text Mining on a large data set, I would explore cloud services like AWS to scale my algorithms on a better computive and memory intensive host.