

# Report for Week 10

## I. Summary

1. I built textual information feature set and compared it with other Feature Sets
2. I performed Statistical Testing for our current results
3. I prepared and updated the repository containing all the work so far

## II. Actions performed

### 1. Textual information

As discussed, I built a new feature set based on the textual information of the users' discourses. Thus, in order to obtain the new features, I detect the top 100 most frequent words (excepting stop words) and consider this to be my vocabulary.

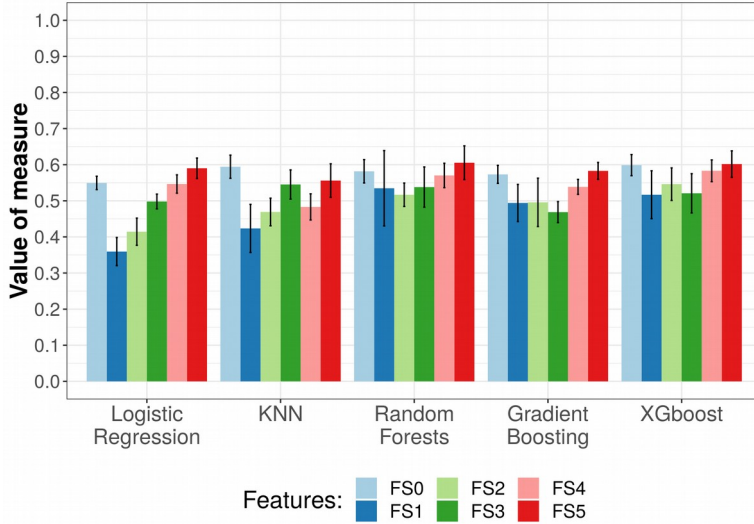
Then, for every time period, I aggregate the users' replies and build the discourse of each user (the aggregation of replies). Now I build a DTM, a Document Term Matrix, having as documents the users discourses and as terms the vocabulary formed earlier.

The weighting method applied is TF-IDF. Thus, now every user has 100 features corresponding to the tf-idf value of each vocabulary word present in his aggregated discourse. I also add the current political stance as a feature. The predicted / learned variable is the next political stance.

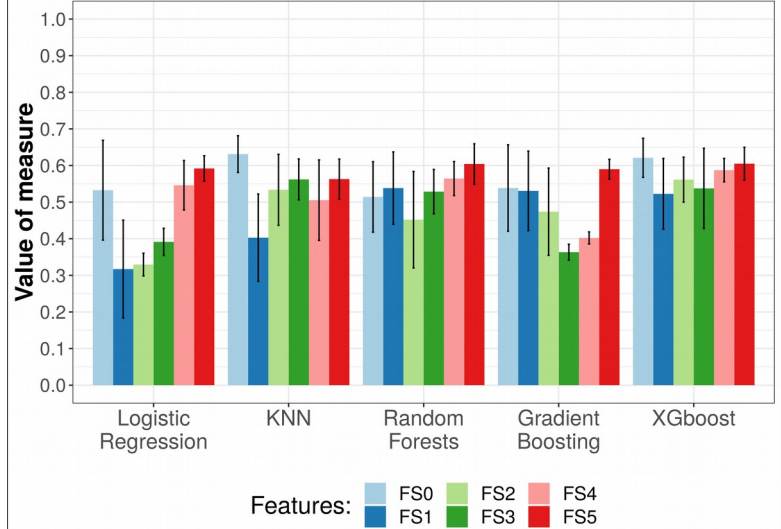
In the graphs below, the shortcuts mean:

Shortcut	Description
FS0	Textual Features
FS1	User activity (#posts, #initial posts, #replies per post)
FS2	User activity per group (same as FS1, just splitted per groups, ie #replies from A, B, N)
FS3	Structure of diffusions
FS4	Conversation Structure (FS1 + FS2 + FS3)
FS5	All = conversation structure + text (FS0 + FS4)

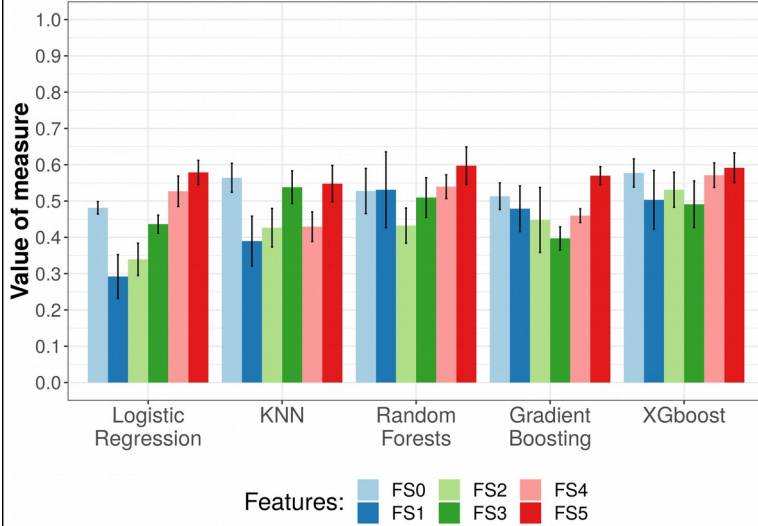
**Prediction performance: Recall**



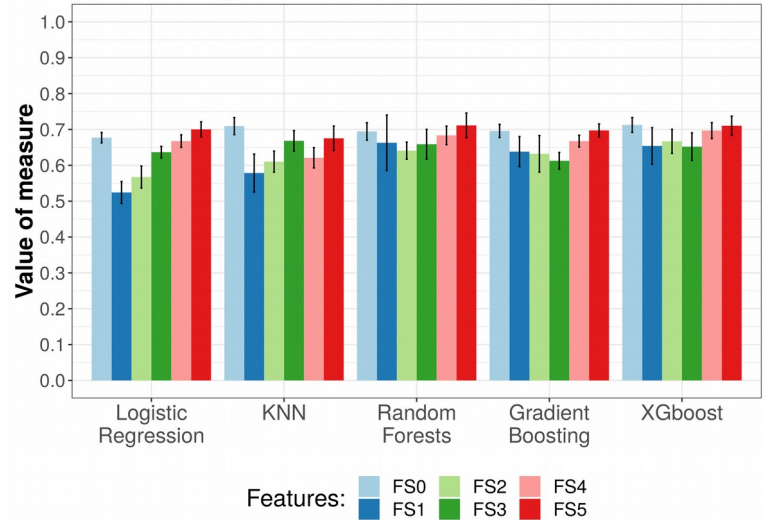
**Prediction performance: Precision**



**Prediction performance: F1**

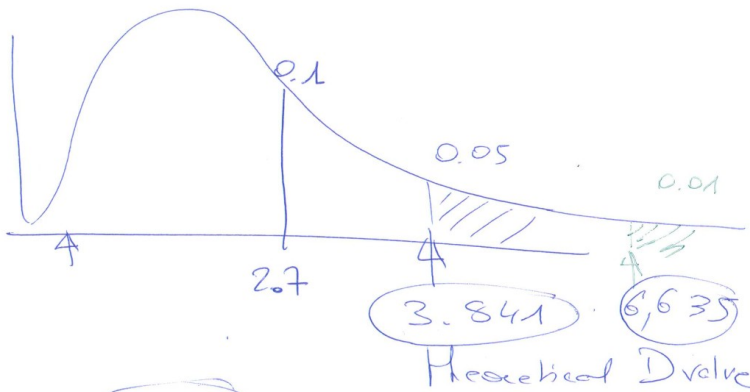


**Prediction performance: Accuracy**



## 2. Statistical Testing

I performed Chi Squared Test to determine if there is significant difference between the features sets when we add textual information. In the left figure, there is the probability density of the Chi Squared distribution, used to determine if the differences are significant or not.



if  $D_{obs} < D_{the}$   
 $\Rightarrow H_0$  can not be refused  $\Rightarrow$  Independent  
 else  $D_{obs} > D_{the}$   $\Rightarrow H_0$  is refused  $\Rightarrow H_1$  is accepted  $\Rightarrow$  significant difference

XGBOOST			
$F_0$ (text)	66.77%	} $\chi^2 = 4.4925$ $p\text{-value} = 0.03405$	
$F_4$ (all without text)	63.2%		
$F_0$ (text)	66.77%	} $\chi^2 = 1.4674$ $p\text{-value} = 0.2257$	
$F_5$ (all)	64.34%		
$F_0$ (text)	66.77%	} $\chi^2 = 22.739$ $p\text{-value} = 1.856 \cdot 10^{-6}$	
$F_3$ (diffusion structure)	58.49%		
$F_0$ (text)	66.77%	} $\chi^2 = 13.092$ $p\text{-value} = 0.0002965$	
$F_2$ (diffusion received replies)	59.77%		

The comparisons are made for XGBOOST between textual information features and the other feature sets. We can see that every time text info is used, there is a significant difference (thus a plus) and the accuracy is always better when we have text features. The only case when we don't have a significant improvement is at FS5 but that set already contains FS0.

## 3. Repository

I updated and completed the repository of this project:

<https://github.com/andreimardale/InformationDiffusion.git>

There I have the Final Documents, the Python scripts used for getting the reddit data and training the classifiers and ALL partial reports throughout this internship. In the README, there is a clean description of every code file and a link to the Data Sets.

Link to the Overleaf Report:

<https://www.overleaf.com/4784255338pvmxffwqtpsz>