

Final project

Norbert Eke

David Rowswell

Aalisha Lakdawala

Intro: Data set

- Spam emails classification
- Objective: minimise the false positive classifications.
- Data set: 4601 emails with content frequencies
- Data source: UCI Machine Learning repository

Methodology & results

1	Hierarchical clustering (single linkage)	39.38 %
2	Hierarchical clustering (average linkage)	39.38 %
3	Hierarchical clustering (complete linkage)	39.38 %
4	knn Classification	18.6 %
5	K-means clustering	36.4 %
6	LDA cross validation	11.8 %
7	LDA	11.0 %
8	QDA	17.12671 %
9	QDA cross validation	16.90937 %
10	Logistic regression	[8.1; 10.5; 7.7] %
11	Artificial Neural Network	[7.7; 6.6; 7.0; 7.6; 6.4; 6.7] %
12	Principle Component Analysis	14.9098 %
13	Bagging	5.23 %
14	Random forrest	4.75 %
15	Classification Trees	9.90%
16	Multiple Linear Regression wo/ interactions	12.10 %
17	Multiple Linear Regression w/ interactions	22.00 %

Logistic regression

- Classification into groups instead of the continuous scale from 0-N
- Misclassification rate: 6.8%
- Results are too big to be shown here but are attached in the appendix

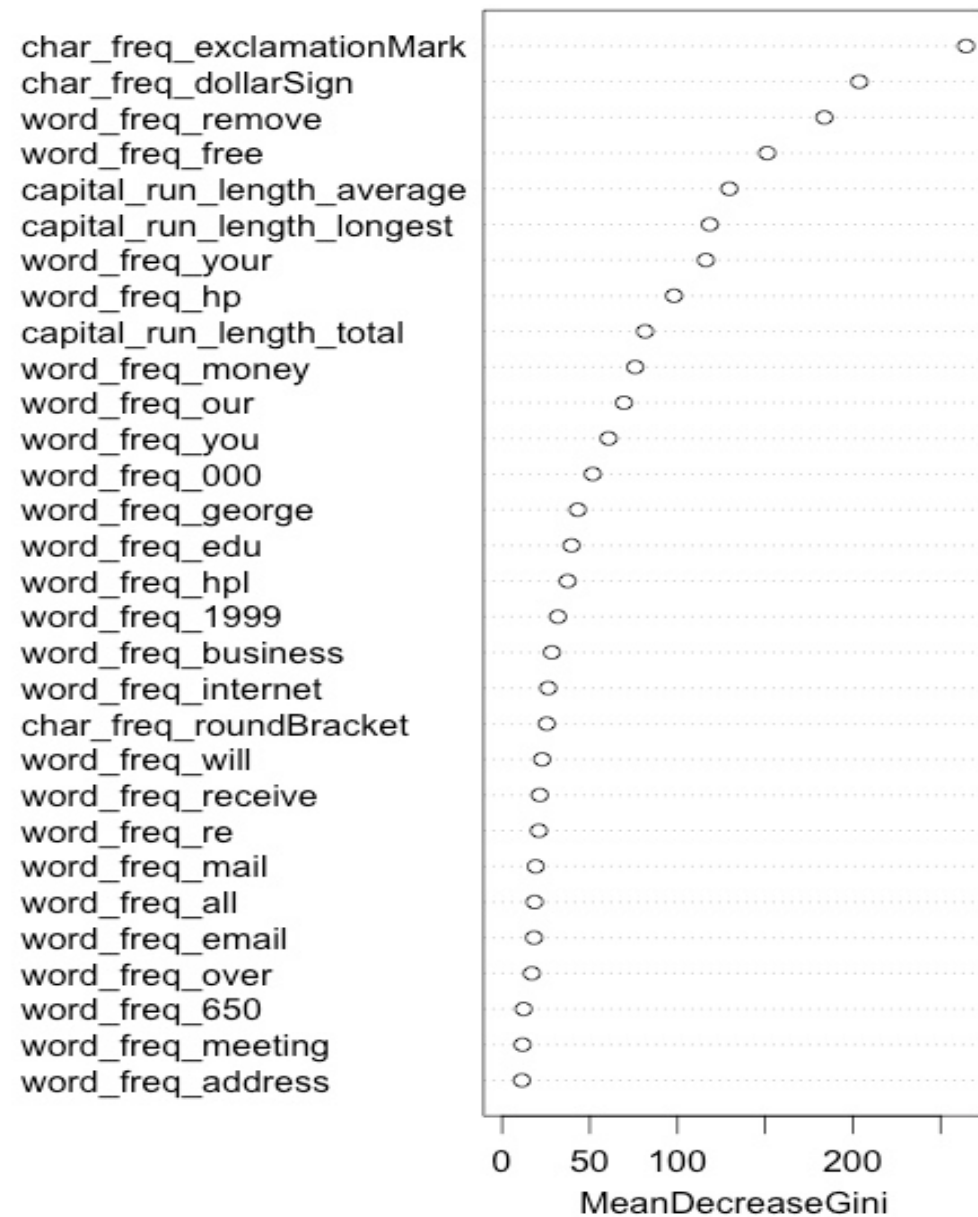
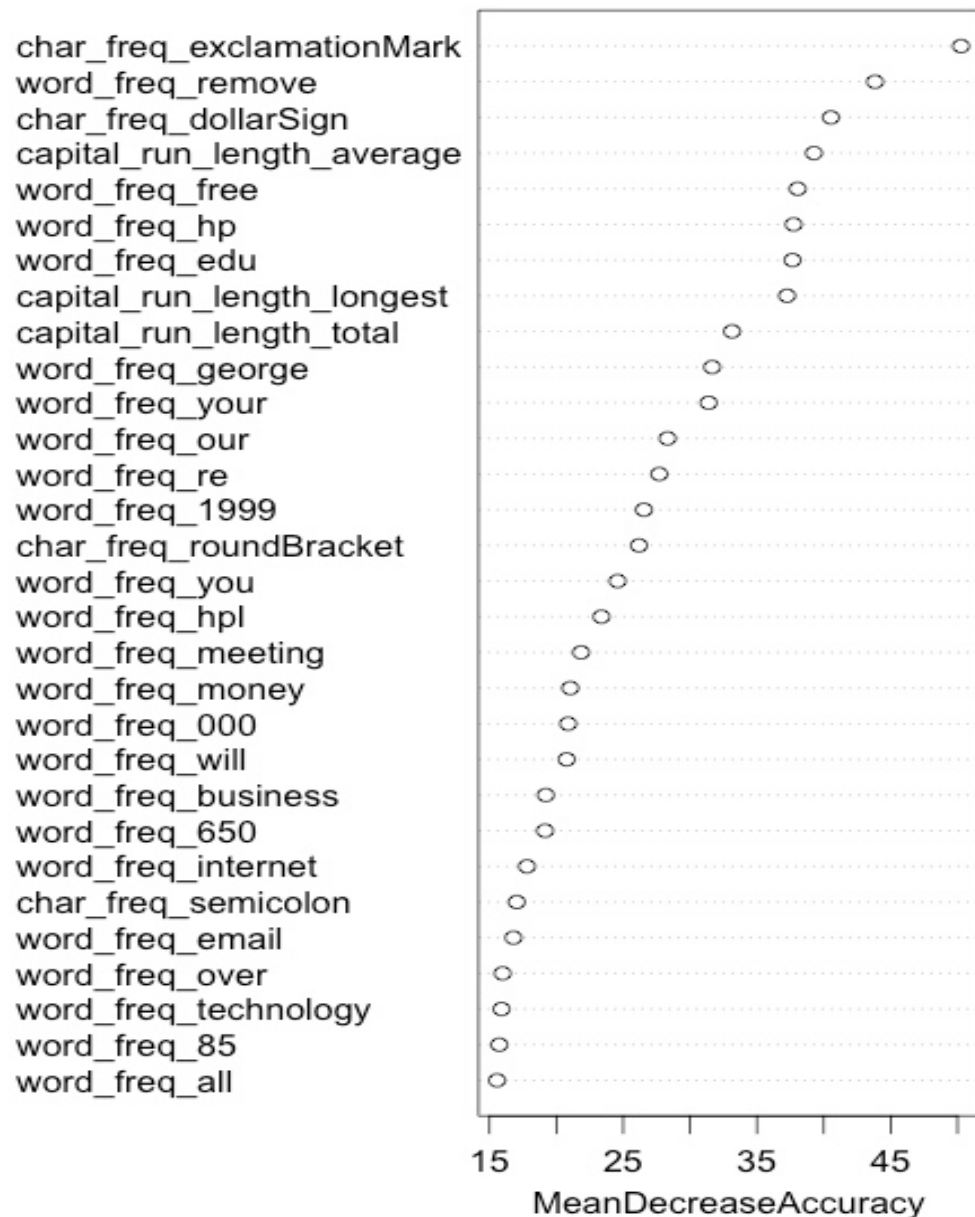
Bagging

- Mtry = 57
- Misclassification = 5.24%
- Since bagging is an average of models, it loses its interpretability but it does have good end result classification.

Random Forests

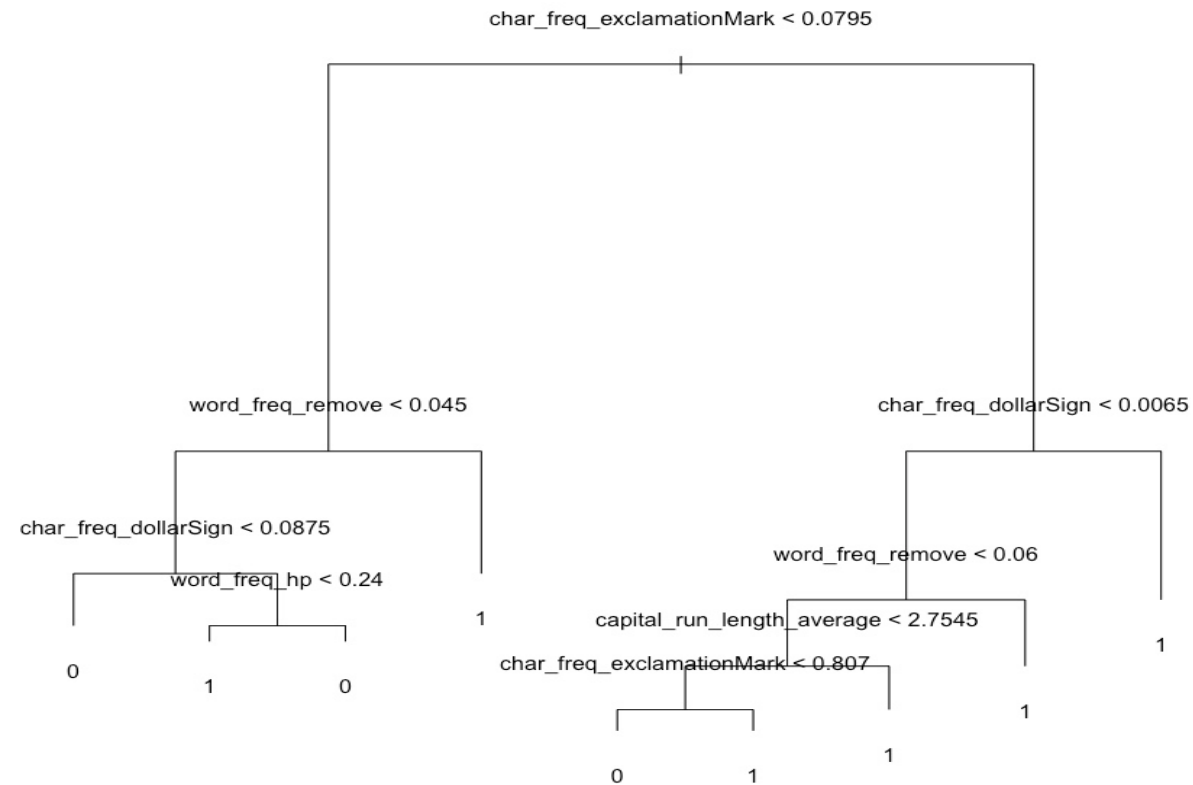
- Set mtry to default
- Decreases variance and increases stability
- Misclassification rate: 4.75%

trainRF



Classification Tree splits

Misclassification % = 9.90 %



0 = non-spam e-mails and 1 = spam emails

Principle Component Analysis

- First principle component loadings:

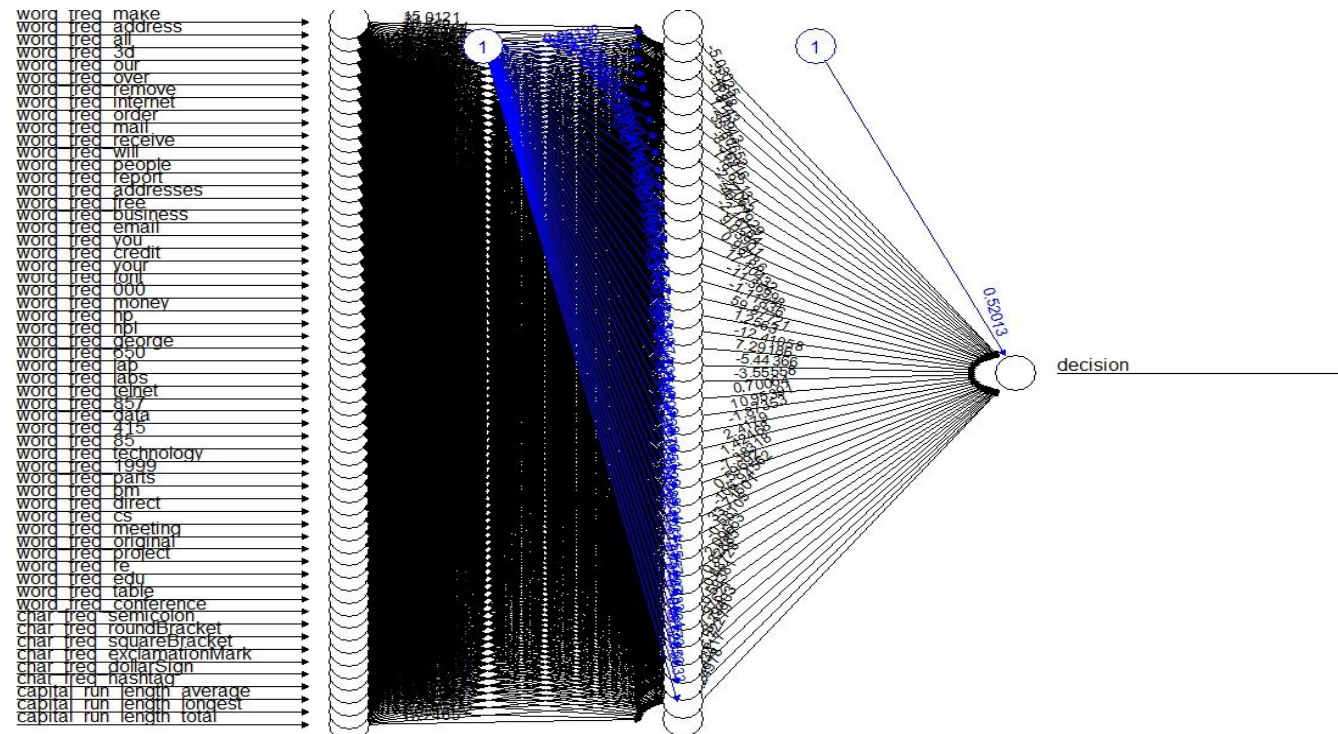
```
> pcaSpam <- prcomp(spam[,-58], scale.=TRUE)
> round(pcaSpam$rotation[,1], 2)

word_freq_make -0.04 word_freq_address -0.01 word_freq_all -0.05
word_freq_3d -0.01 word_freq_our -0.04 word_freq_over -0.05
word_freq_remove -0.05 word_freq_internet -0.03 word_freq_order -0.05
word_freq_mail -0.02 word_freq_receive -0.05 word_freq_will -0.02
word_freq_people -0.04 word_freq_report -0.02 word_freq_addresses -0.03
word_freq_free -0.04 word_freq_business -0.05 word_freq_email -0.02
word_freq_you -0.08 word_freq_credit -0.03 word_freq_your -0.08
word_freq_font -0.01 word_freq_000 -0.05 word_freq_money -0.04
word_freq_hp 0.21 word_freq_hpl 0.21 word_freq_george 0.04
word_freq_650 0.28 word_freq_lab 0.22 word_freq_labs 0.30
word_freq_telnet 0.31 word_freq_857 0.35 word_freq_data 0.01
word_freq_415 0.35 word_freq_85 0.27 word_freq_technology 0.32
word_freq_1999 0.05 word_freq_parts 0.00 word_freq_pm 0.04
word_freq_direct 0.32 word_freq_cs 0.01 word_freq_meeting 0.02
word_freq_original 0.07 word_freq_project 0.01 word_freq_re 0.01
word_freq_edu 0.00 word_freq_table 0.00 word_freq_conference 0.00
char_freq_semicolon 0.00 char_freq_roundBracket 0.02
char_freq_exclamationMark -0.04 char_freq_dollarSign 0.00
capital_run_length_average -0.02 capital_run_length_longest -0.03 capital_run_length_total -0.04
```

The word frequencies make sense, since the emails are Hewlett-Packard Internal-only Technical Reports.

Artificial Neural Networks

- `nnSpam <- neuralnet(formula, data=train, linear.output = FALSE, hidden=40)`
- `nnresults <- compute(nnSpam, test[, -58])`
- `table(testDecision, nnresults$net.result > 0.5)`



Neutral Network Result

Misclassification %	Hidden = ...	Misclassification %	Hidden = ...
7.7	5	6.7	c(12,11)
6.6	10	8.9	c(18,17)
7.0	20	8.7	c(20,10)
7.6	30	7.6	c(20,15)
7.0	35	7.5	c(20,18)
6.9	36	6.9	c(25,13)
6.4	37	7.6	c(25,15)
8.1	38	8.3	c(20,20)
6.7	39	8.5	c(25,20)
6.4	40	8.2	c(10,10,5)
7.1	45	8.1	c(10,10,10)
6.7	50	8.5	c(15,10,5)
6.8	c(7,7)	7.9	c(15,10,10)
7.3	c(10,10)	7.7	c(15,15,10)
7.4	c(15,10)	7.0	c(15,15,15)
6.7	c(15,15)	9.0	c(20,10,10)
7.8	c(13,13)	7.6	c(20,15,10)
7.4	c(11,11)	9.9	c(20,15,15)

Conclusion

- Bagging, random Forests, Neural networks & Logistic regression worked the best
- Misclassification error can be brought down to 4.75%