

# Comparing ML models for single-label text classification

Ana Cardoso Cachopo  
Jan 2017

# Outline

- Problem statement, goals and success criteria
- Explain data and run statistical analysis
- Apply ML models to data
- Summary and future work

# Problem statement, goals and success criteria

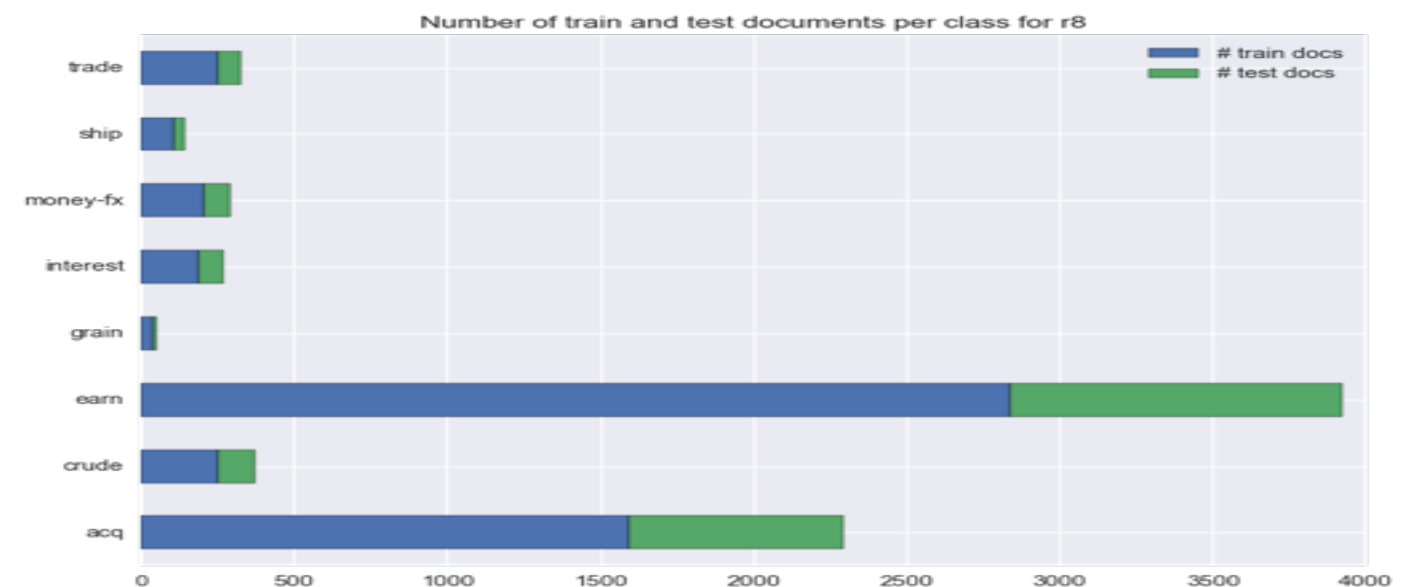
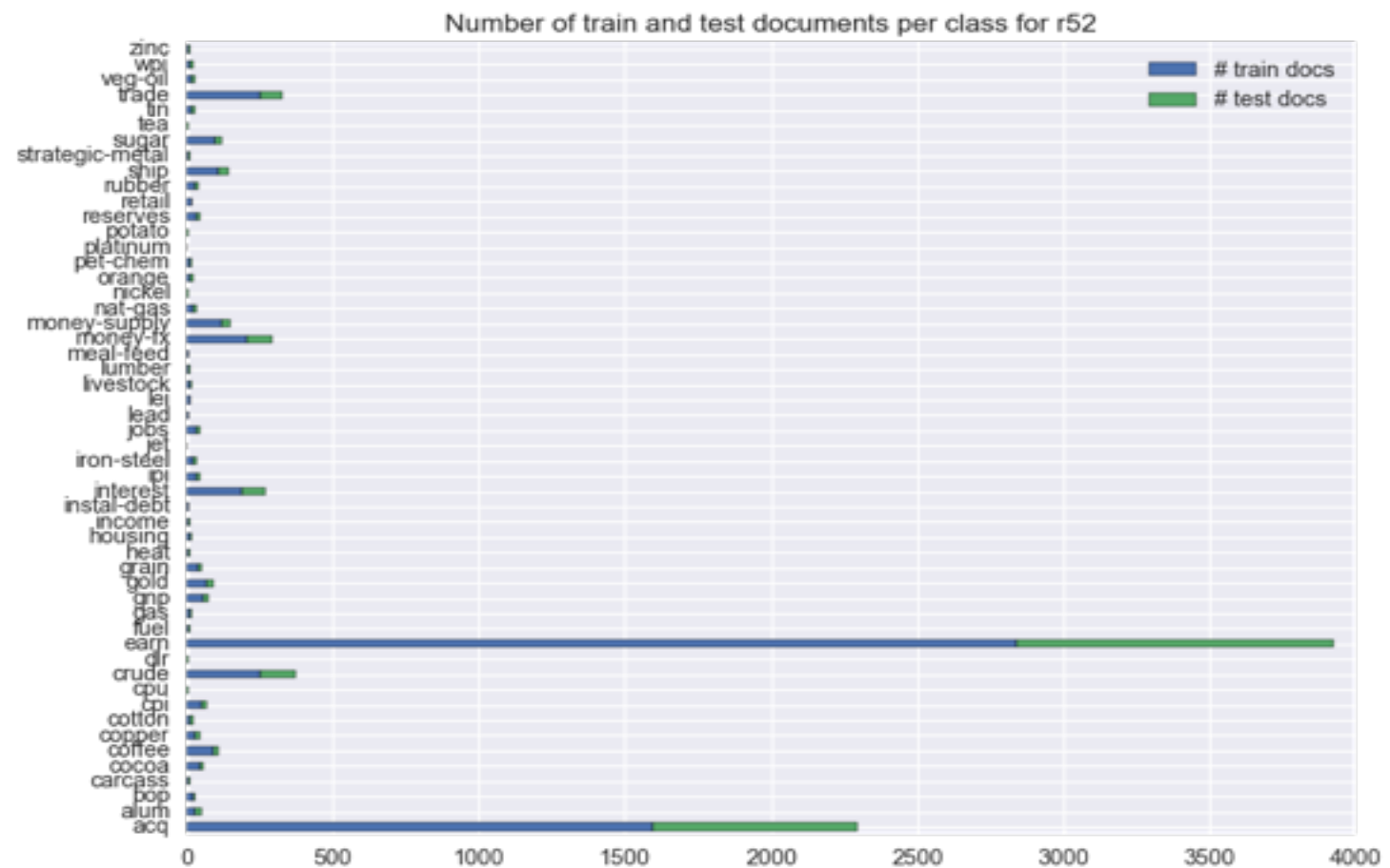
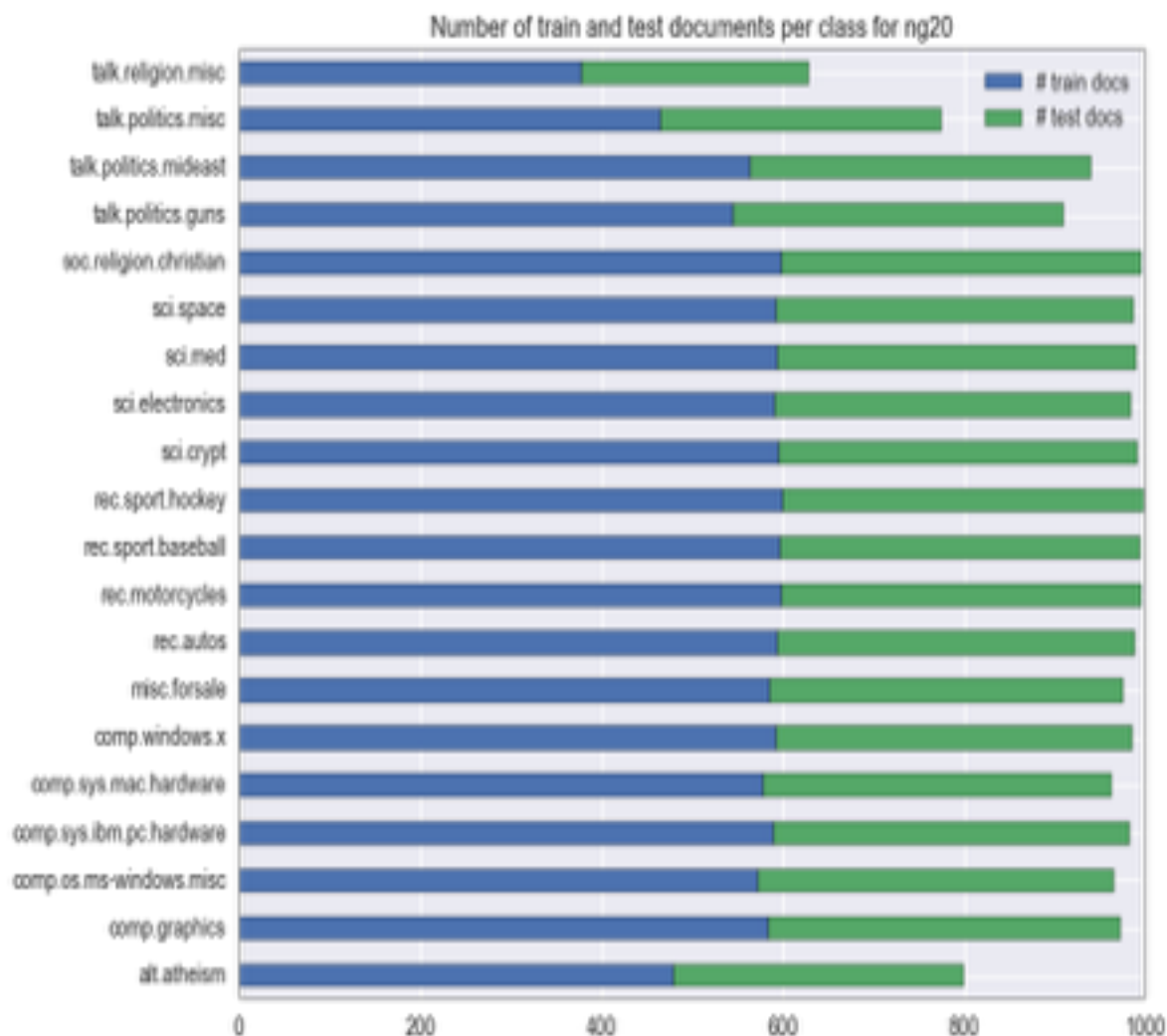
- How well do different Machine Learning models perform on single-label text classification tasks?
- Goals: reproduce part of my PhD work using state-of-the-art libraries in Python, assess how this area evolved in 10 years.
- Successful if: reproduce the initial "related work" from my thesis.

# Explain data and run statistical analysis

- Datasets: 20 Newsgroups, Reuters-21578
- Numbers of documents
- Numbers of features
- Word clouds

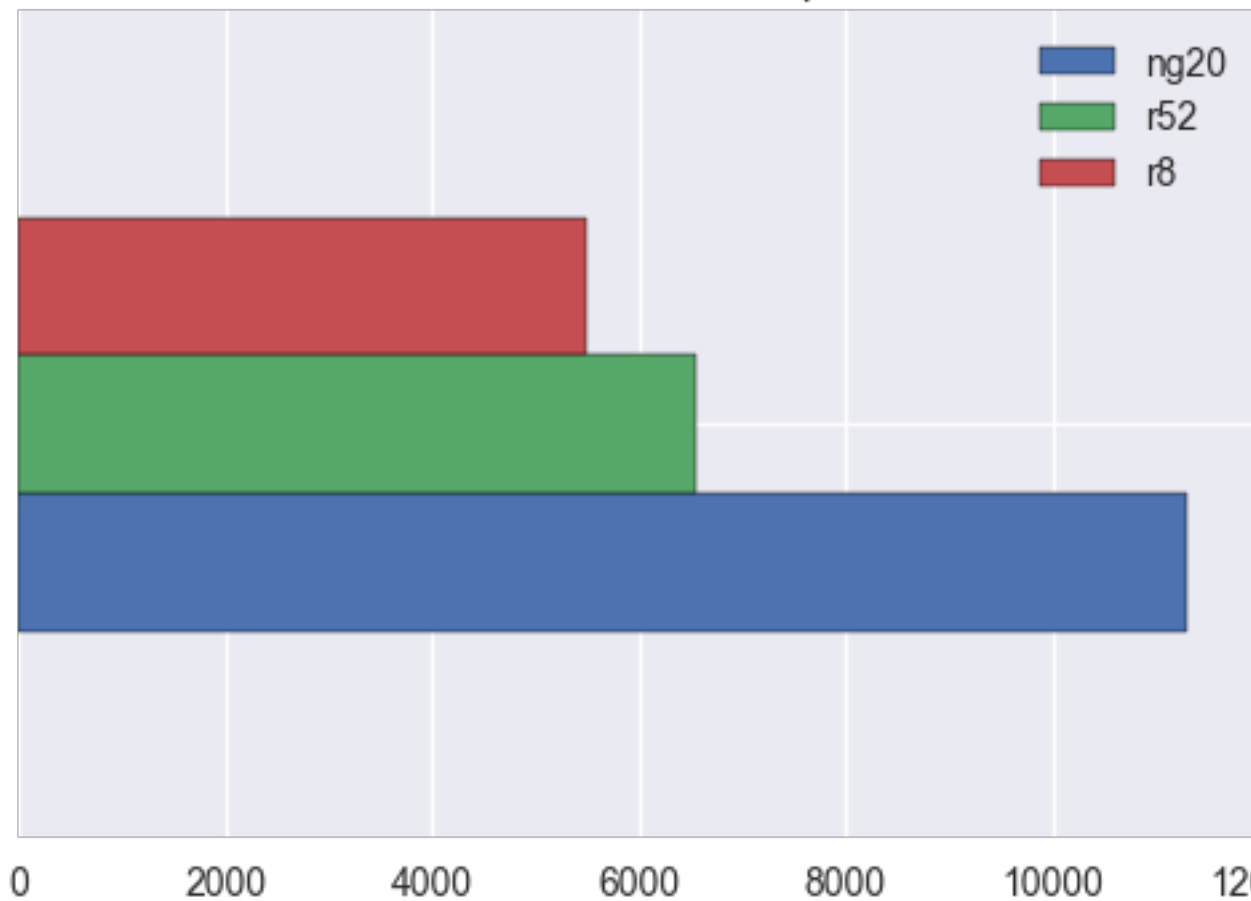
# Numbers of documents

	# train docs	# test docs	total # docs
<b>r8</b>	5485	2189	7674
<b>r52</b>	6532	2568	9100
<b>ng20</b>	11293	7528	18821

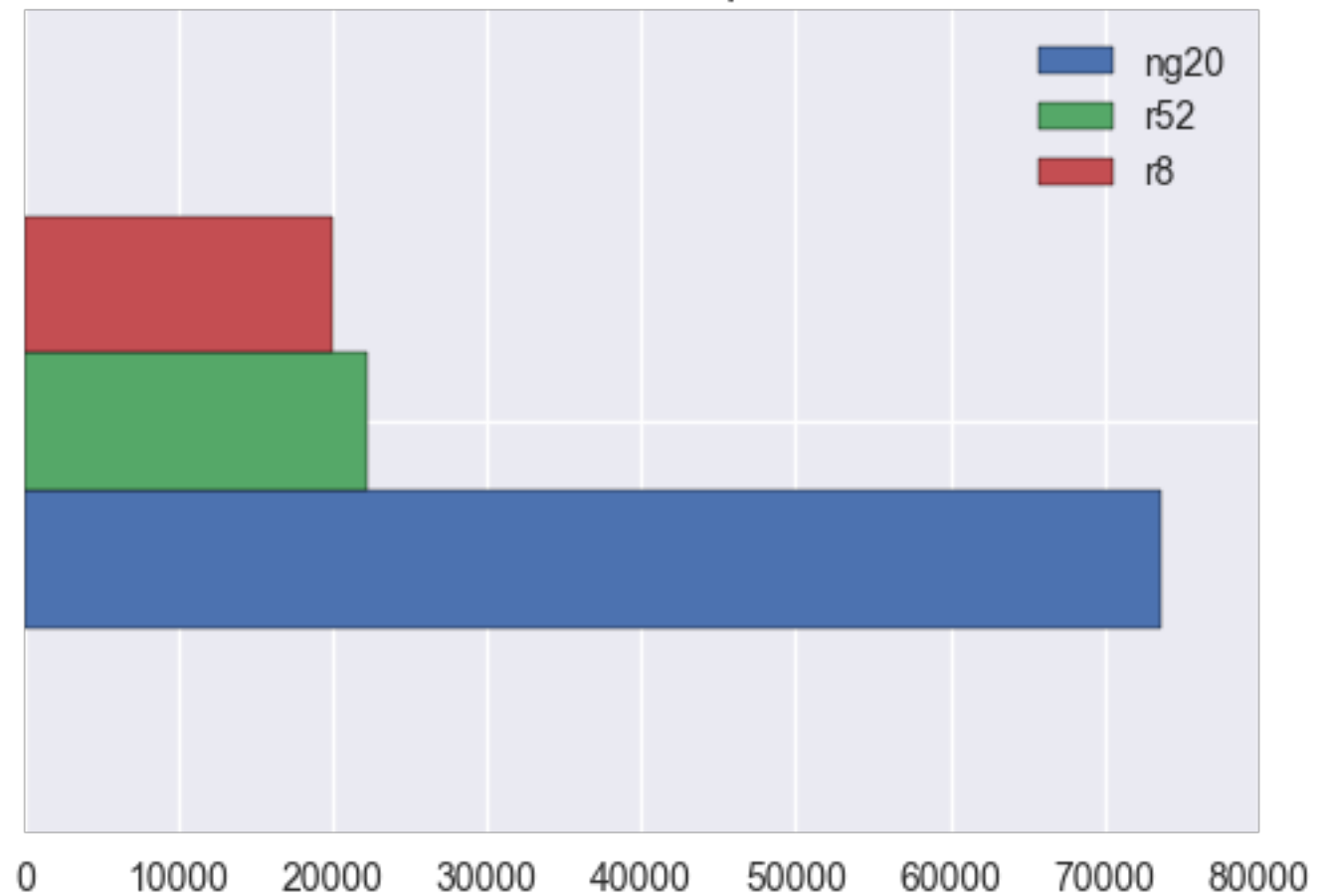


# Numbers of features

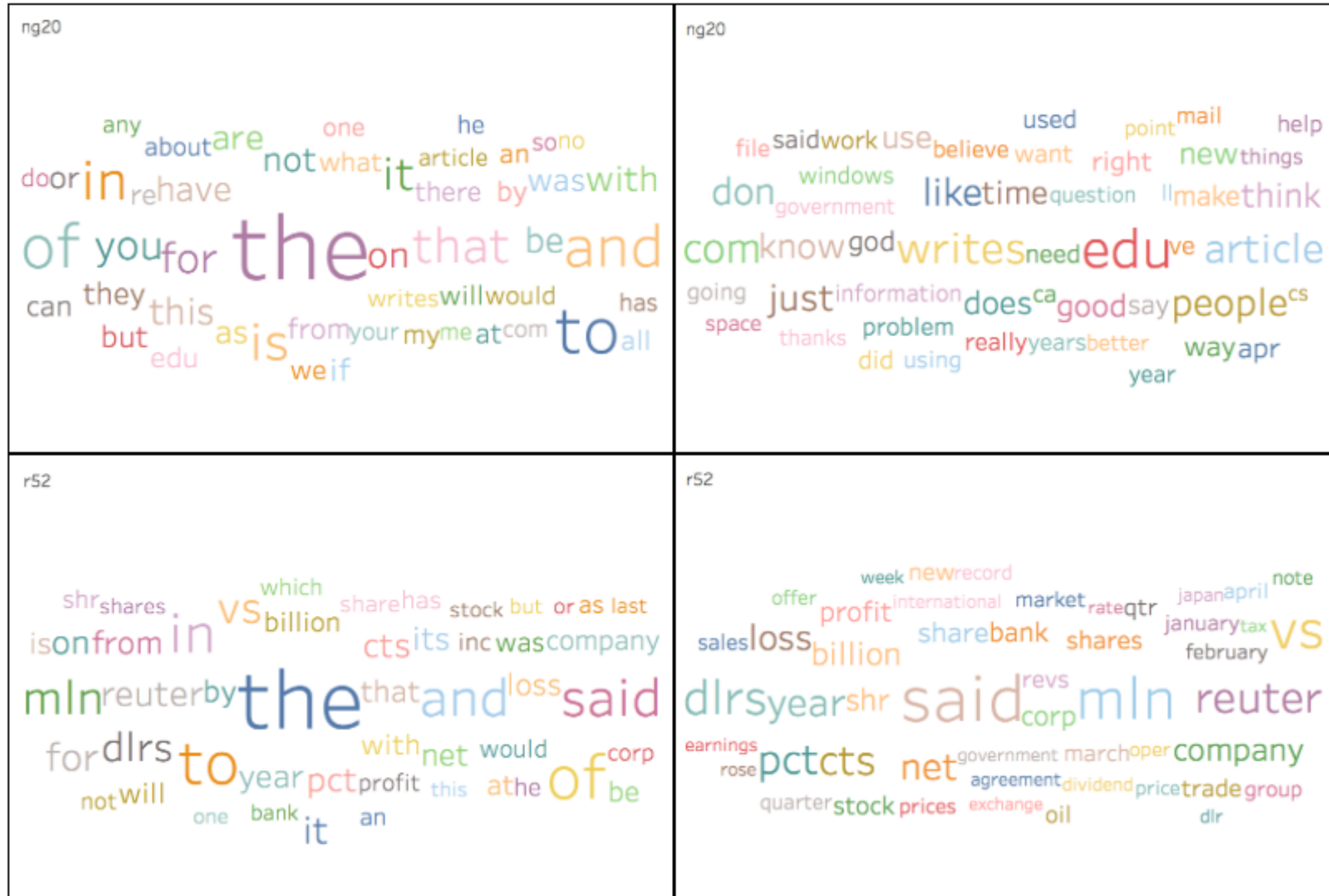
Number of train documents per dataset



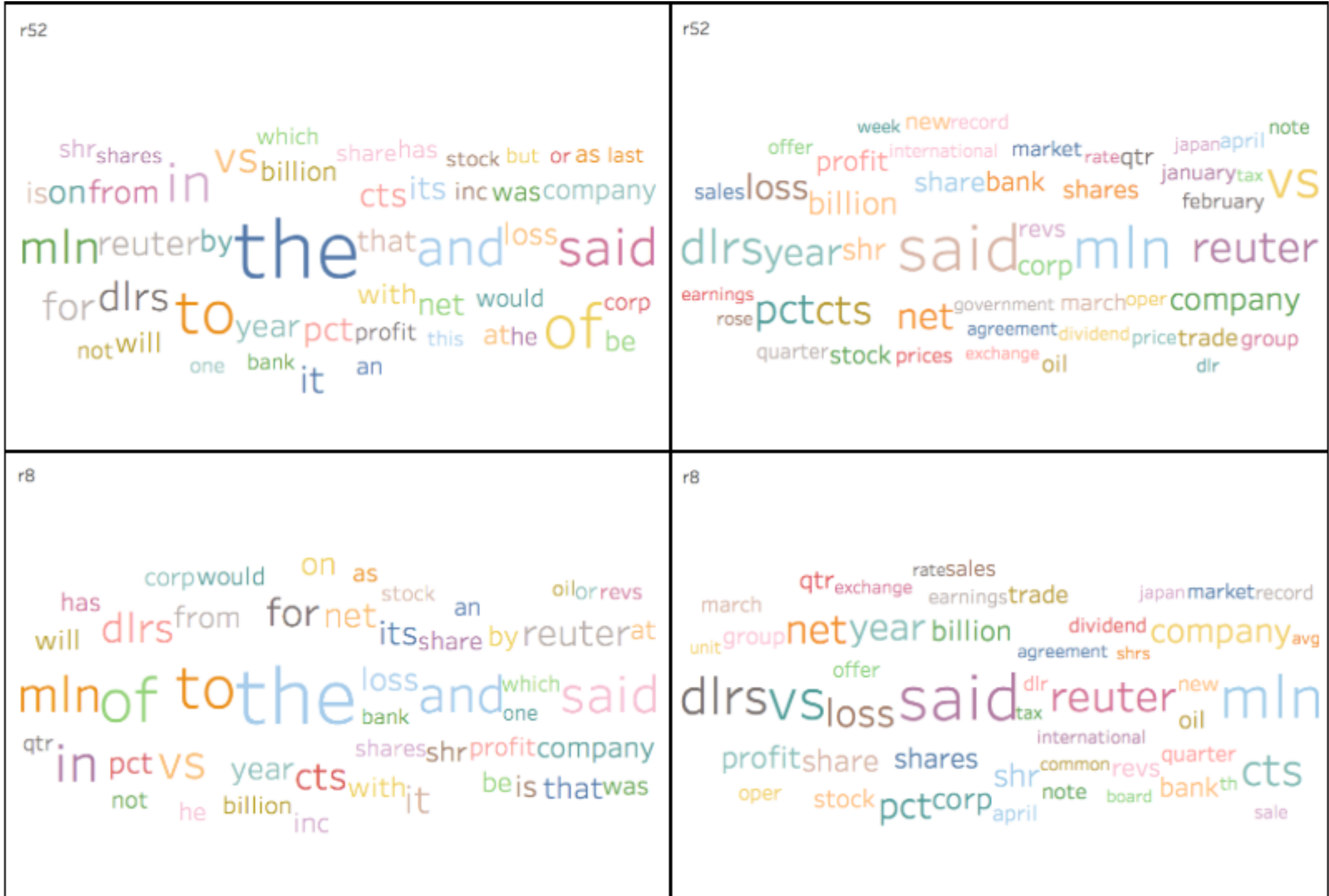
Number of features per dataset



# Word clouds



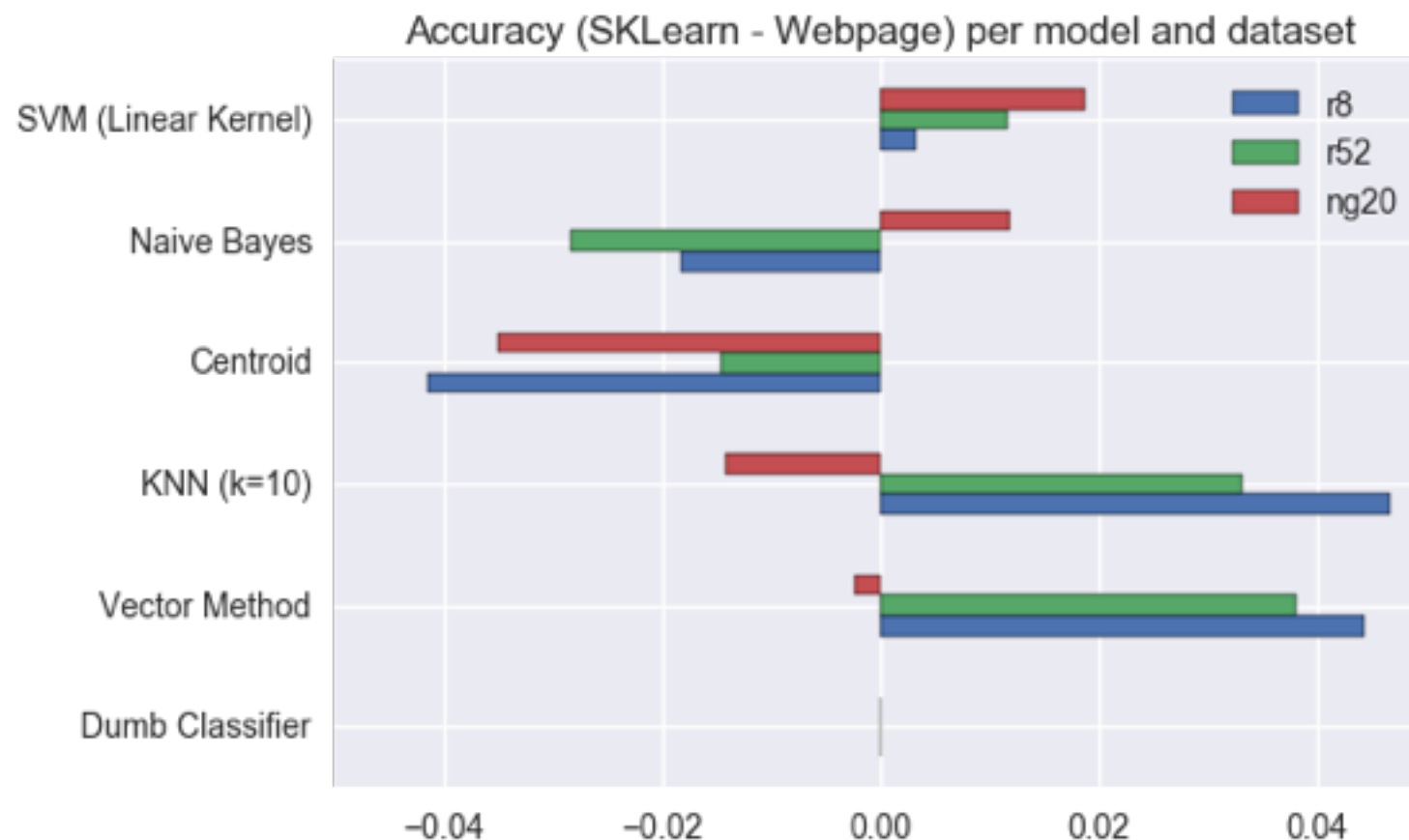
# Word clouds (cont.)





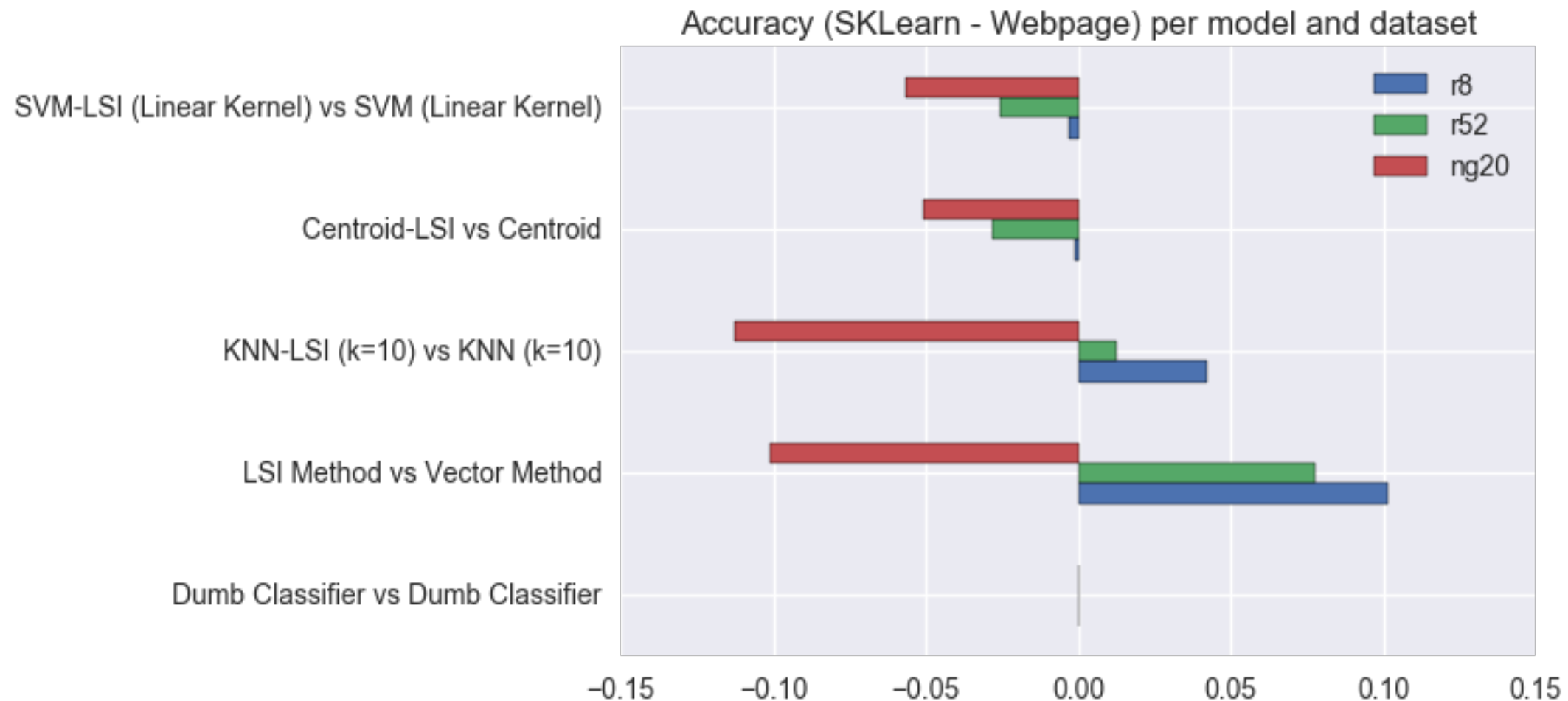
# Apply ML models (reproduce previous results)

Accuracy Values					
Classification Method	R8	R52	20Ng	Cade12	WebKb
Dumb classifier	0.4947	0.4217	0.0530	0.2083	0.3897
Vector Method	0.7889	0.7687	0.7240	0.4142	0.6447
kNN (k = 10)	0.8524	0.8322	0.7593	0.5120	0.7256
Centroid (Normalized Sum)	0.9356	0.8717	0.7885	0.5148	0.8266
Naive Bayes	0.9607	0.8692	0.8103	0.5727	0.8352
SVM (Linear Kernel)	0.9698	0.9377	0.8284	0.5284	0.8582



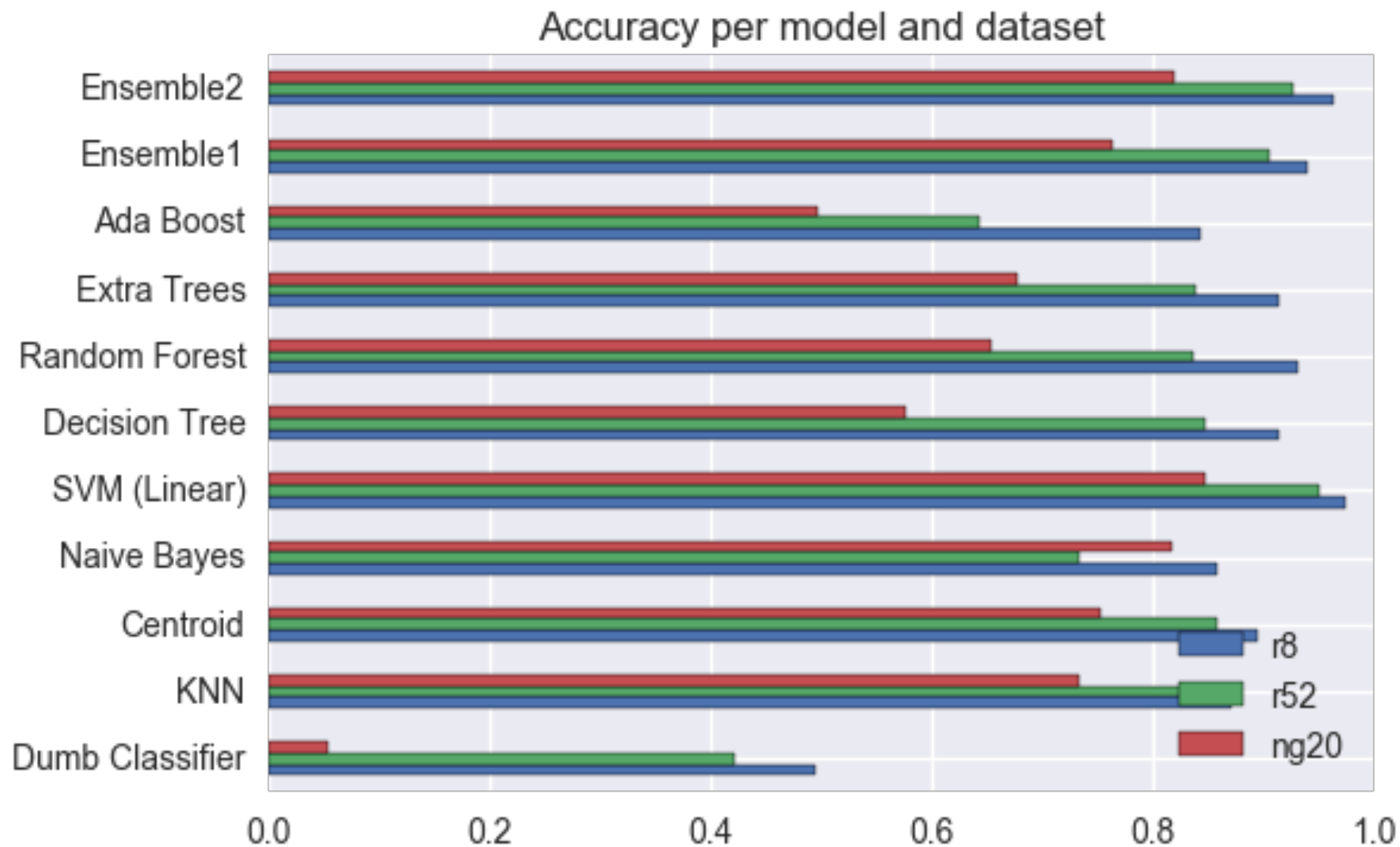
# Apply ML models

(use LSI/LSA/SVD to improve results)



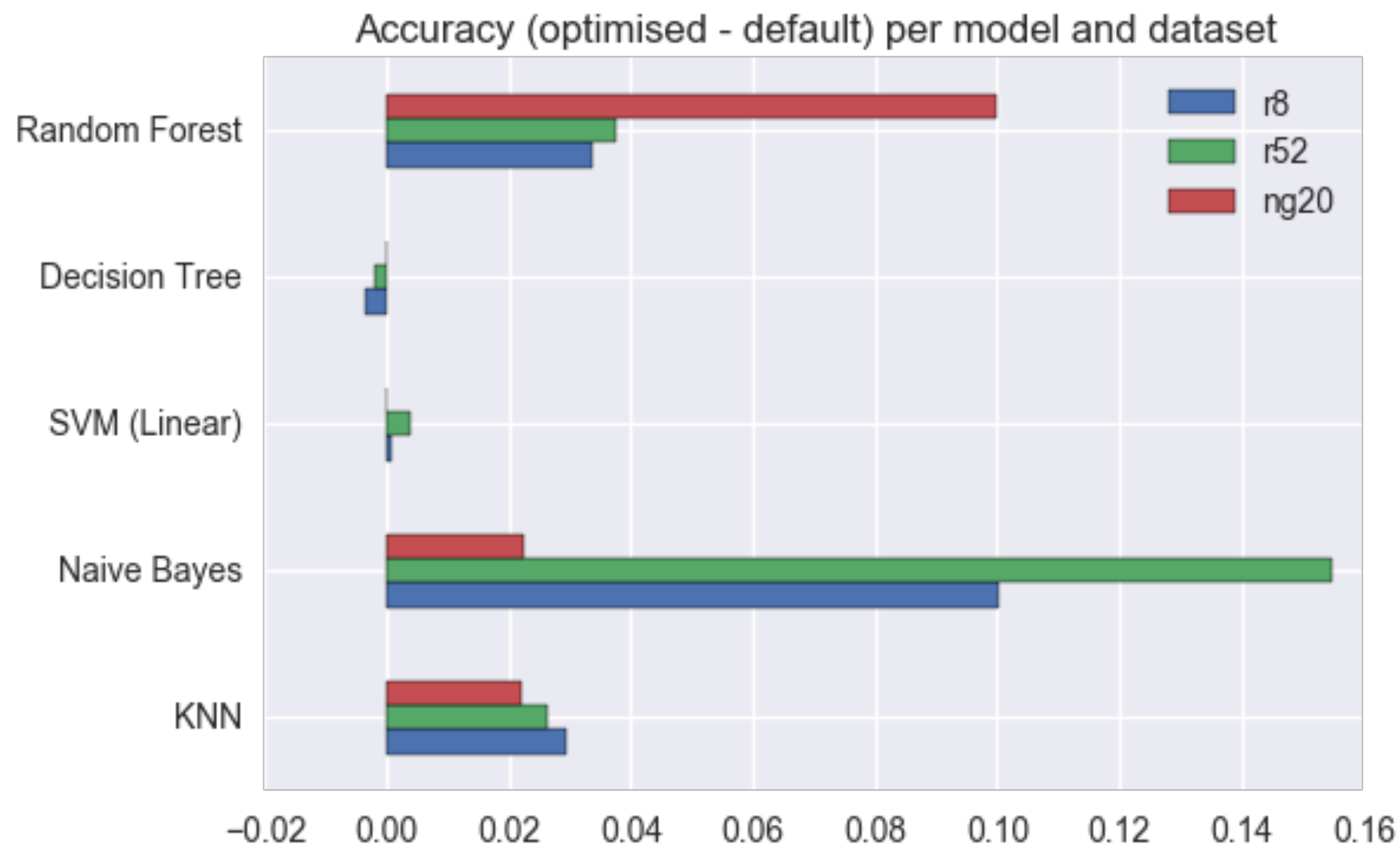
# Apply ML models

(use more RT/Boost/Ensemble)



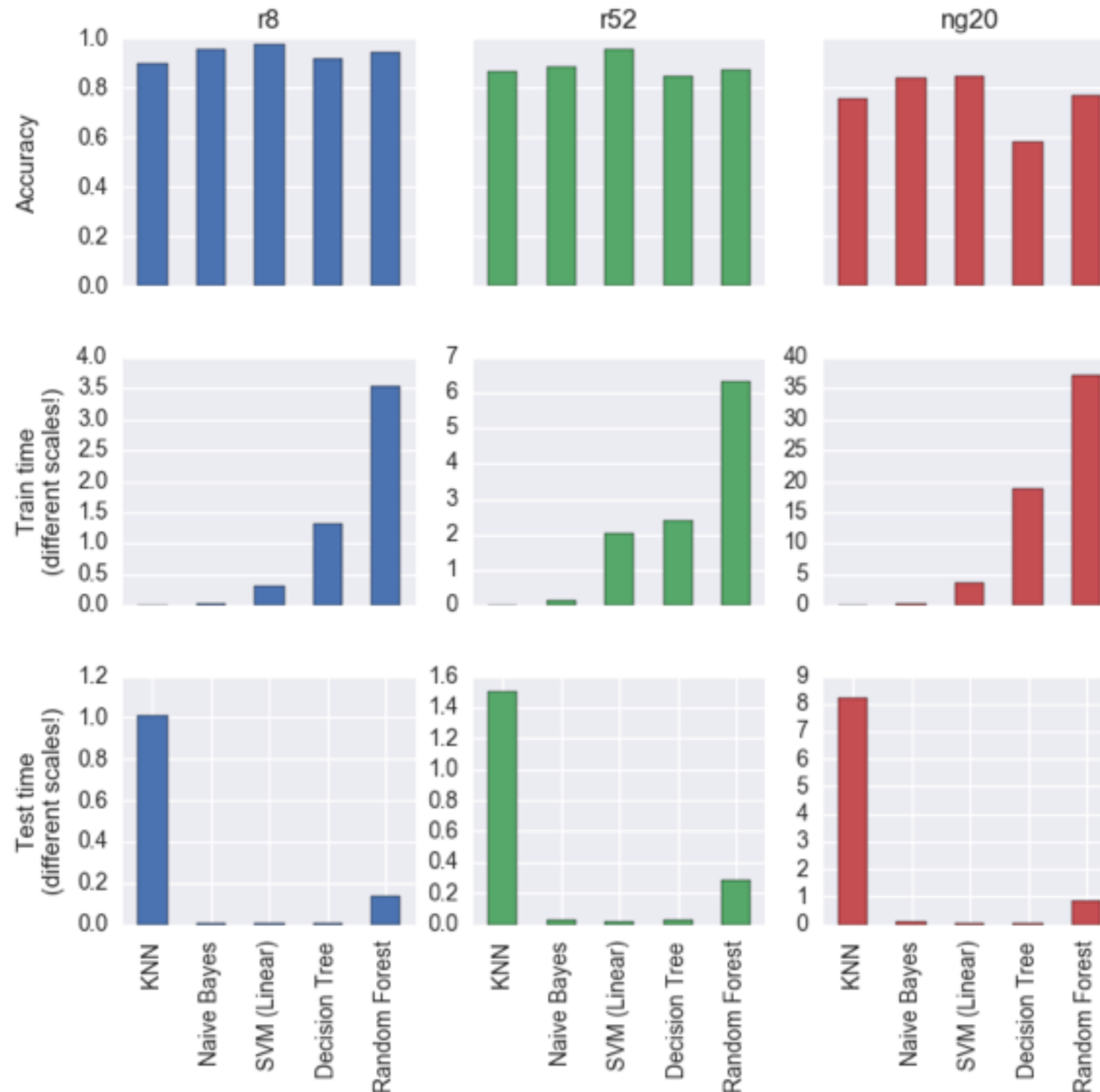
# Apply ML models

(optimise parameters — GridSearchCV)



# Apply ML models

(compare results and time)



# Summary

- **SVMs (still) rock!**
- Future work:
  - More ML models
  - More (and bigger) datasets
  - Different feature selection
  - Other success measures
  - AWS to parallelise tasks

# Thank you!

- Questions?