

Project Farseer

News Popularity Prediction

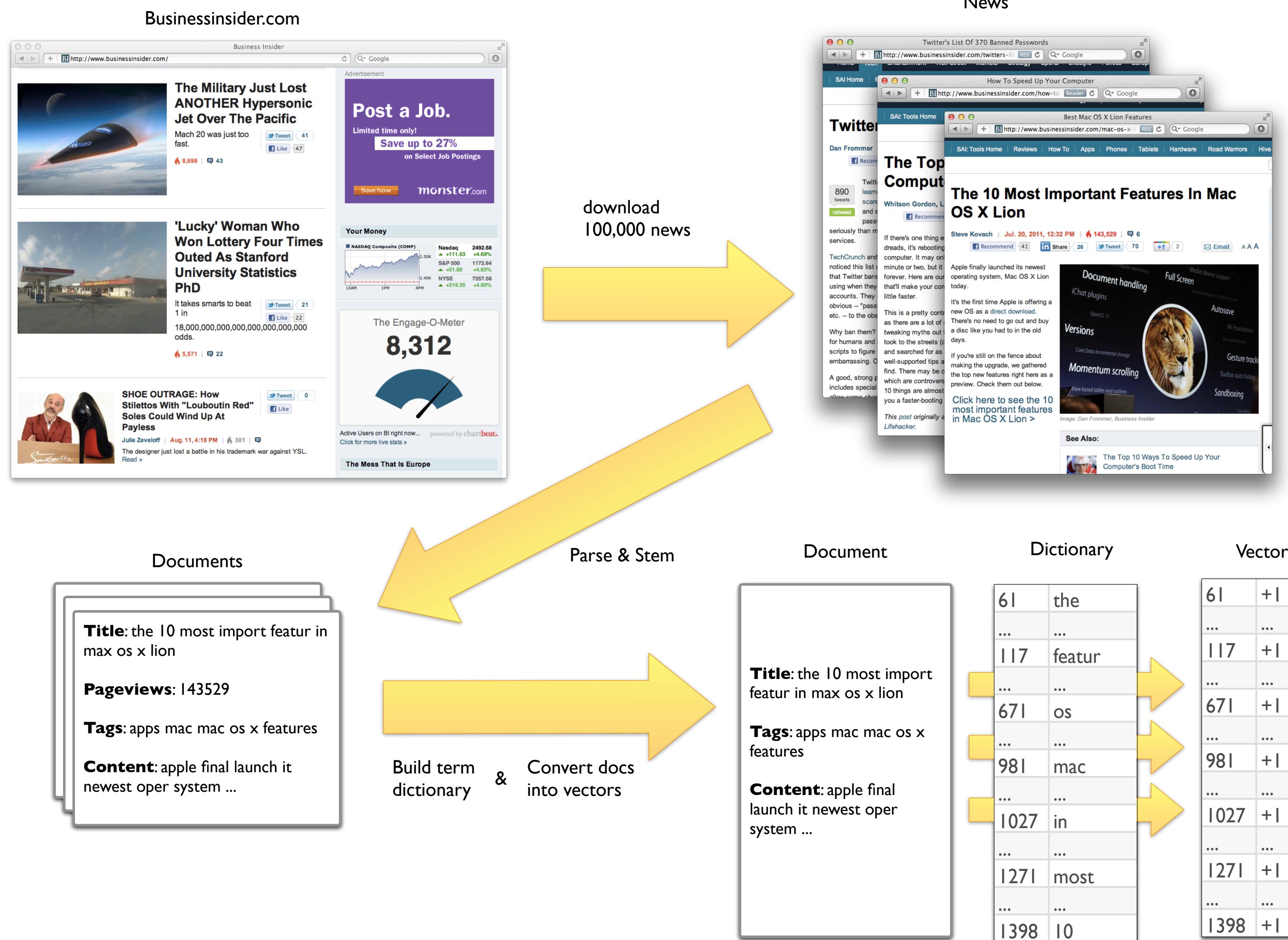
Vladimir Zaytsev
Pablo Paredes
Eduard Tuchfeld
UC Berkeley



I. Goal and Opportunity

- Predict the popularity (# page views) of a news article given its content
- Create an analysis tool to improve content creation.

2. Observations



Summary

- Preparation and selection of data is important.
- We can only predict lower bounds of page views.

Ideas for Future Work

- Use larger data sets (from various sources)
- Add new components (emotions, time, pictures)
- Define a multidimensional metric of popularity (pg. views, tweets, FB likes, others)
- Use other classification techniques (SVM, Random Trees, GBM, etc.)

3. Reduction

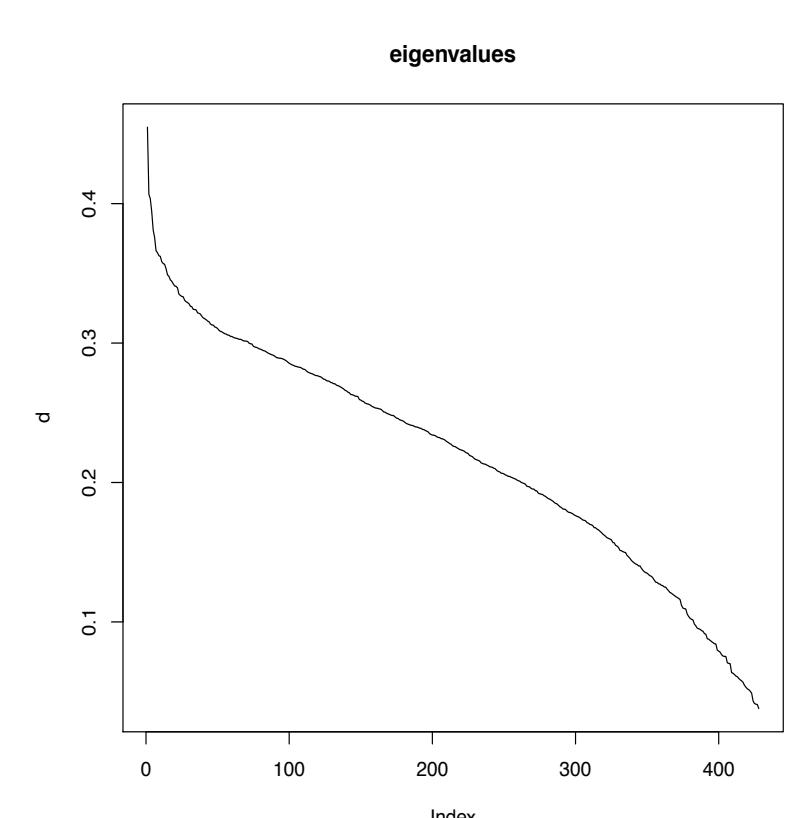
Calculate score - average number of page views

total page views	frequency	score	term
1800552	15	120037	valley
1626632	16	101665	silicon
1136660	12	94722	present
1306454	18	72581	woman
794235	11	72203	inspire
1405330	20	70267	tour
3438812	52	66131	2
1815864	29	62616	hot
1324035	22	60183	differ
1699716	31	54830	truth
1038872	22	47221	fact
1208126	26	46466	half
1172500	26	45096	men
656446	15	43763	epic
511907	13	39377	download
841427	22	38247	crush
978864	27	36254	pictur
526572	15	35105	shock
1023056	30	34102	os

Select the terms which are the most correlated with number of pageviews

total page views	frequency	score	term
a%			
TOP TERMS			valley silicon inspire 20 woman tour fact epic download picture shock os 14 21 phoe highest incredible future fake lion x comic beautiful most hollywood
BOTTOM TERMS			forecast tool dax trend outlook weak forex armagedon finance letterman monuments content fantacy hunt alert oil fuel cloud politics direct conan yanke radio tv summit
b%			

Positive influence on the number of page views
Negative influence on the number of page views

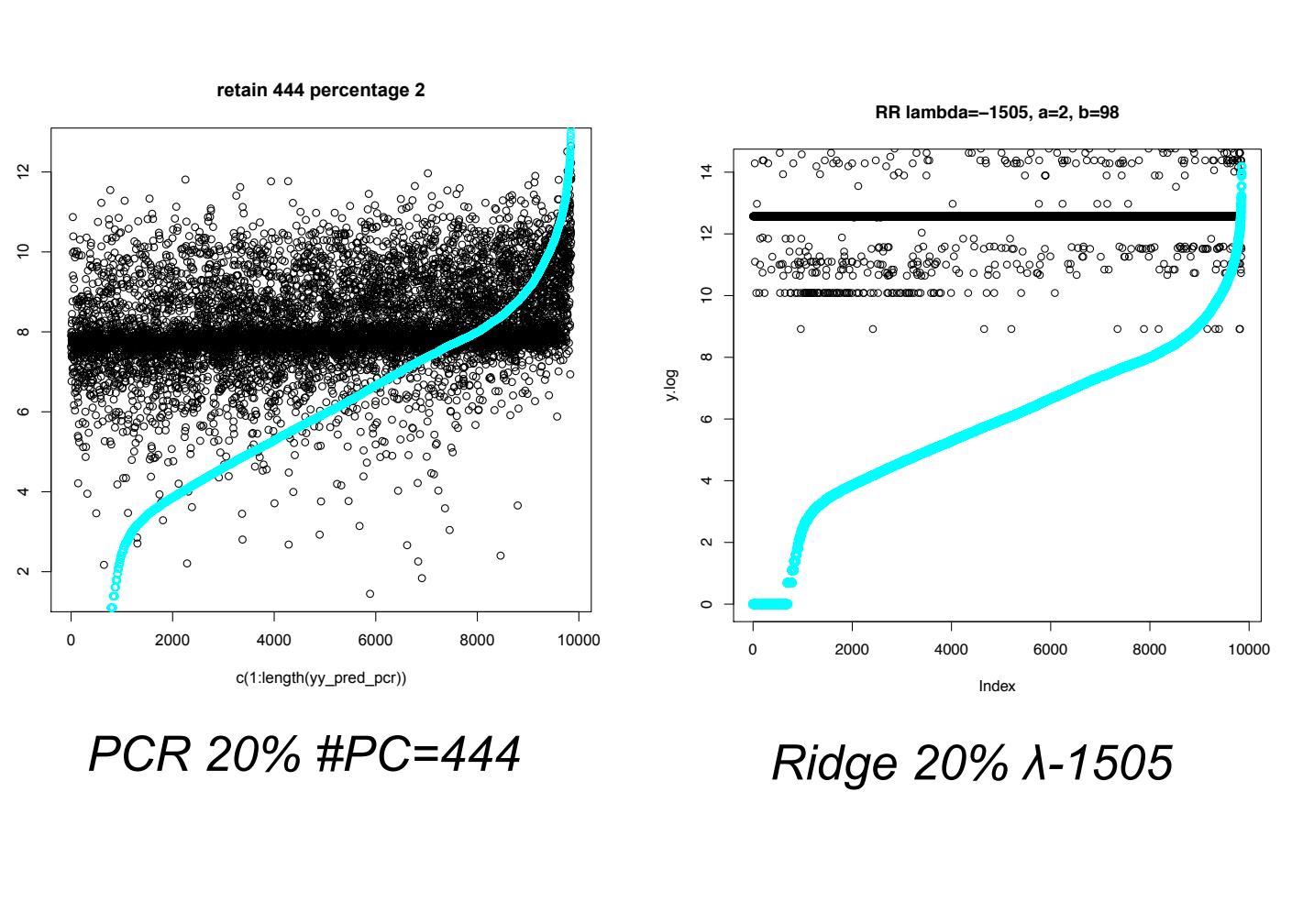
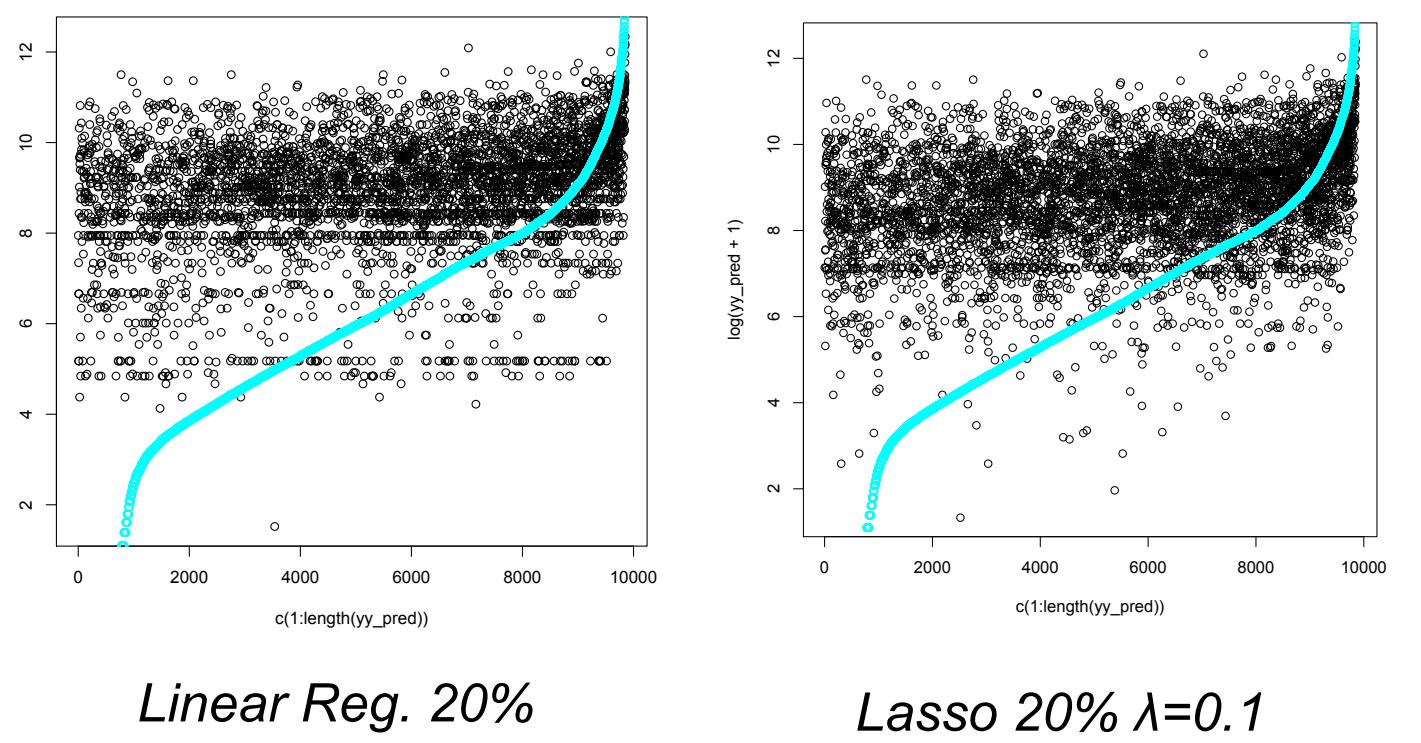
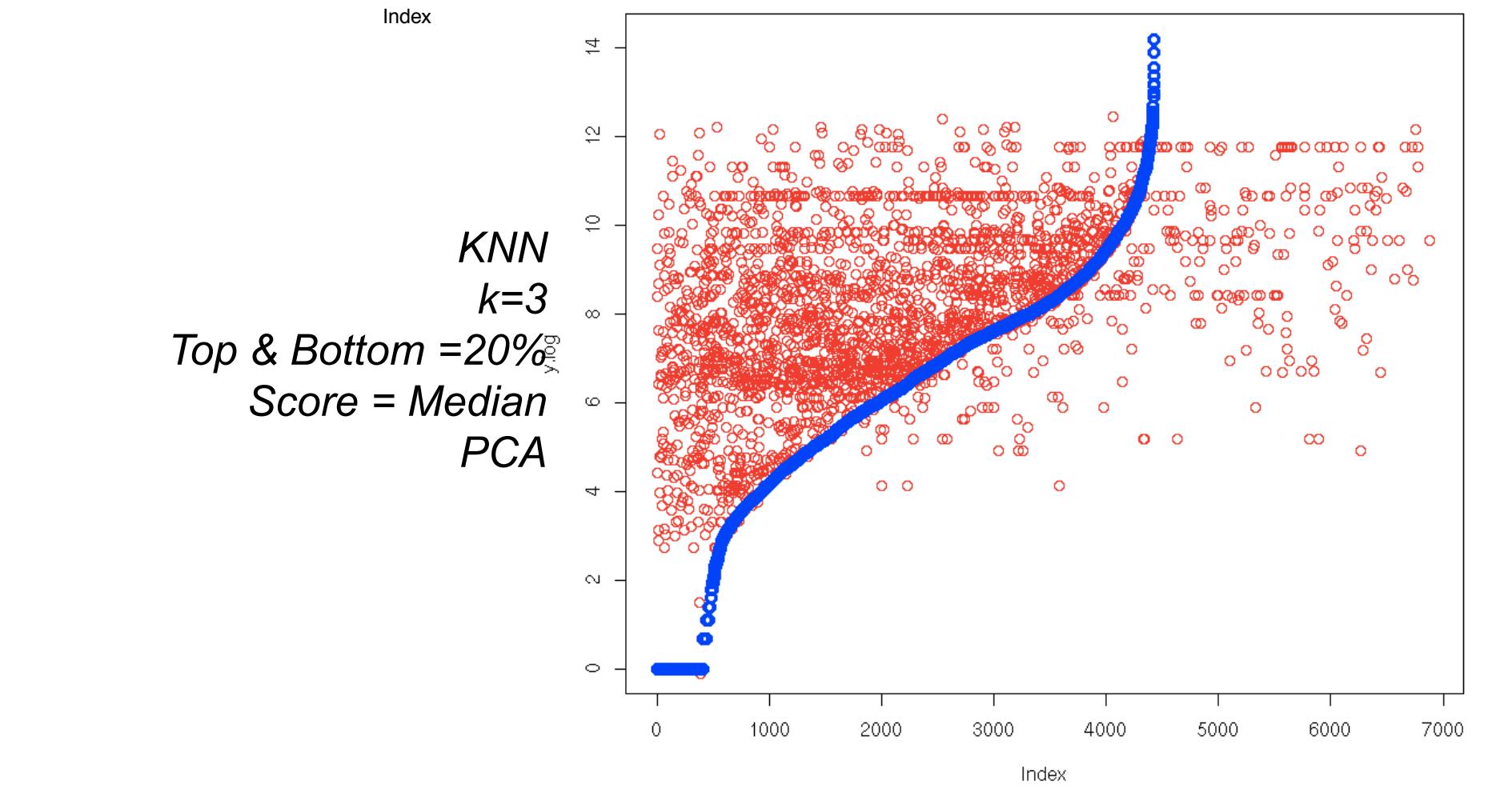
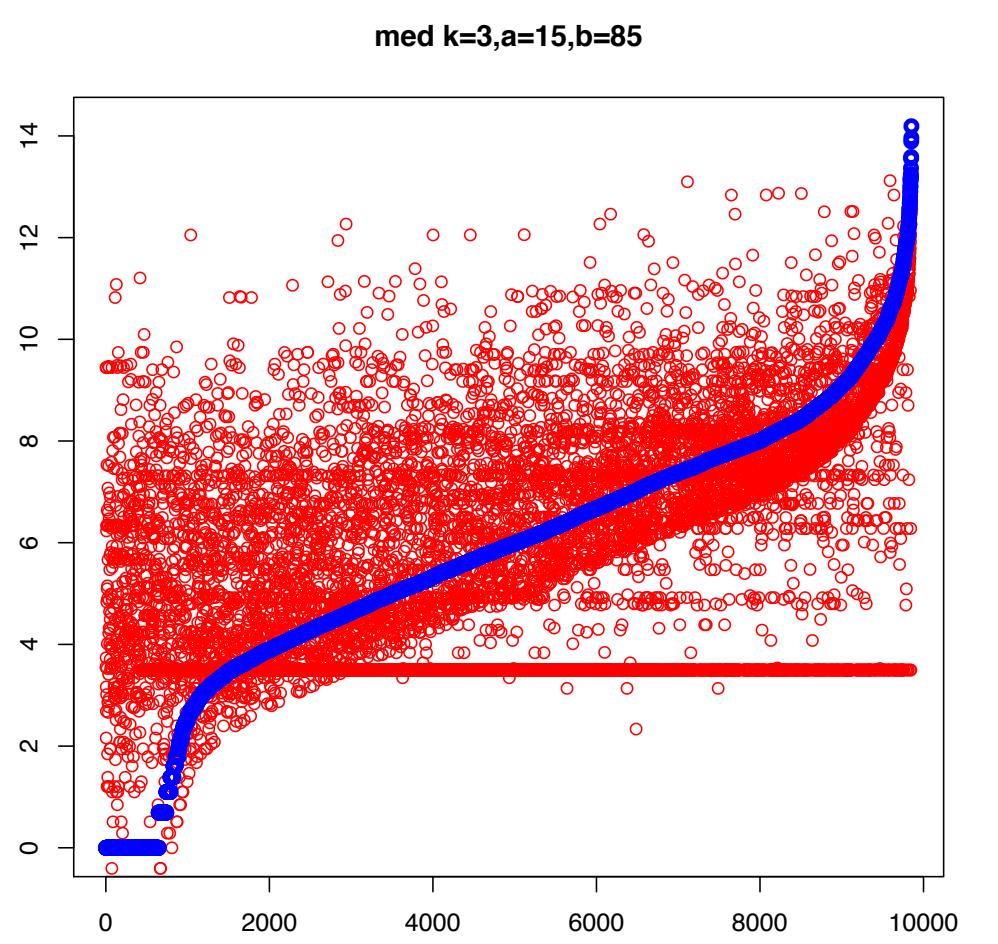


Apply PCA

Some Data Facts

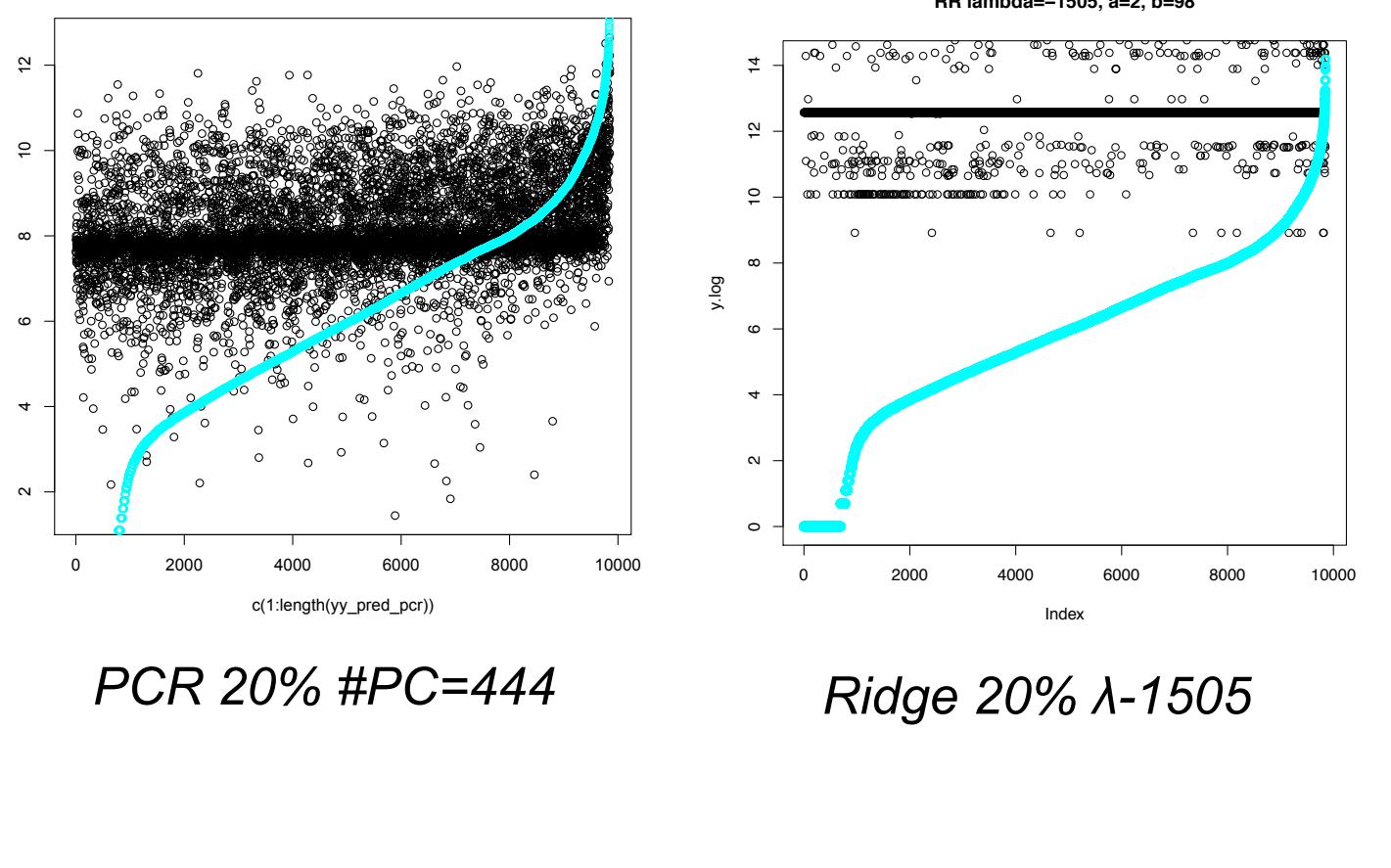
- Used different parameters a and b
- Used different datasets - with all the content of news or just only with titles
- With PCA or without PCA
- Tried to use medians to calculate scores
- Produced about 250 possible datasets with 100,000 documents and from 400 to 1500 terms

4. Results



RR lamda=1505, a=2, b=68

retain 444 percentage 2



Ridge 20% $\lambda=1505$