

Reading Tea Leaves: How Humans Interpret Topic Models

Paper by Jonathan Chang, Jordan Boyd-Graber, Sean Gerrish, Chong Wang, David M. Blei

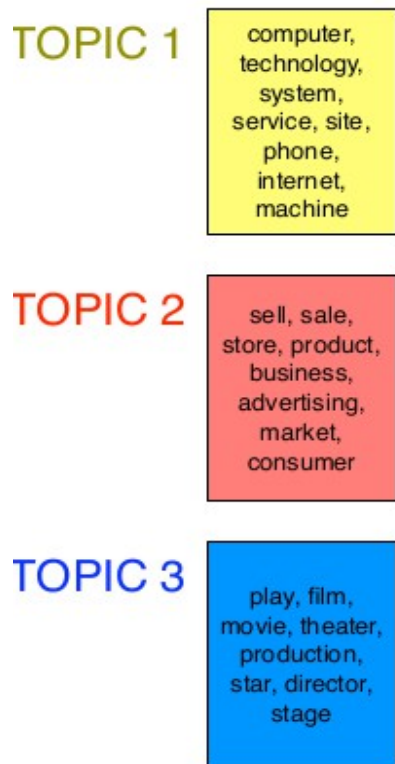
Presentation by M. Falis

Roadmap

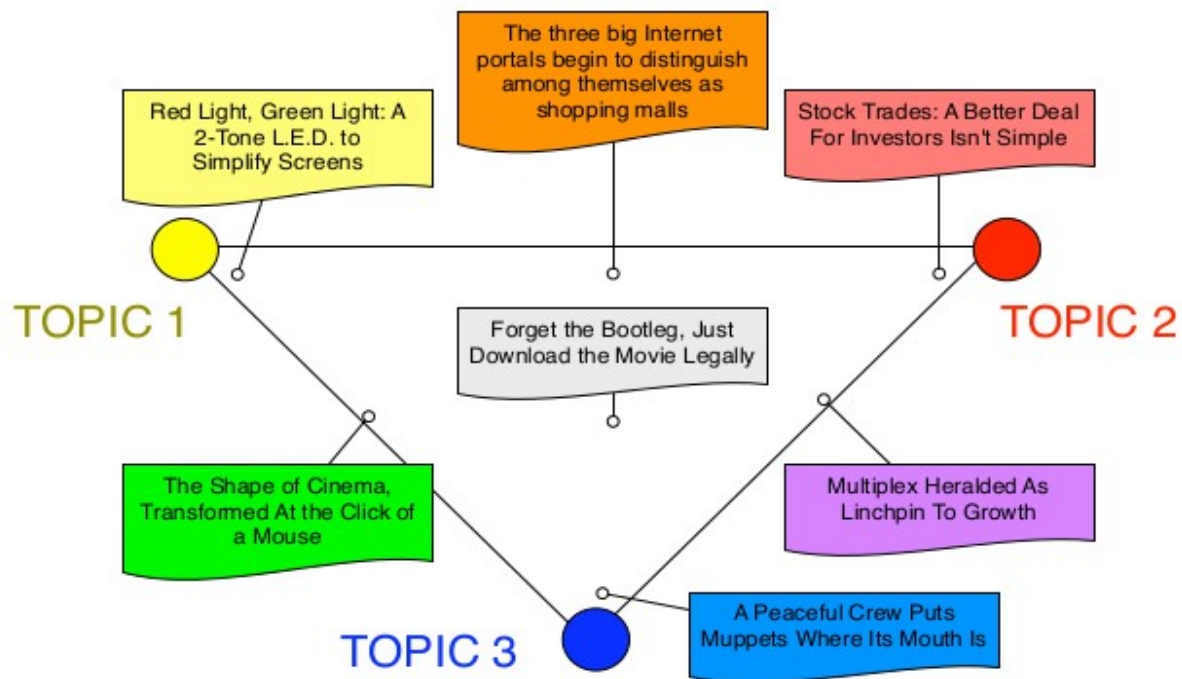
- Intro to topic modelling
- Models
- LDA
- Measuring Performance
- MP and TLO
- Results
- Summary

Intro to Topic Modelling

- Motivation
- Topic = Distribution over words
- Document = Mixture of topics
- Problem definition



(a) Topics

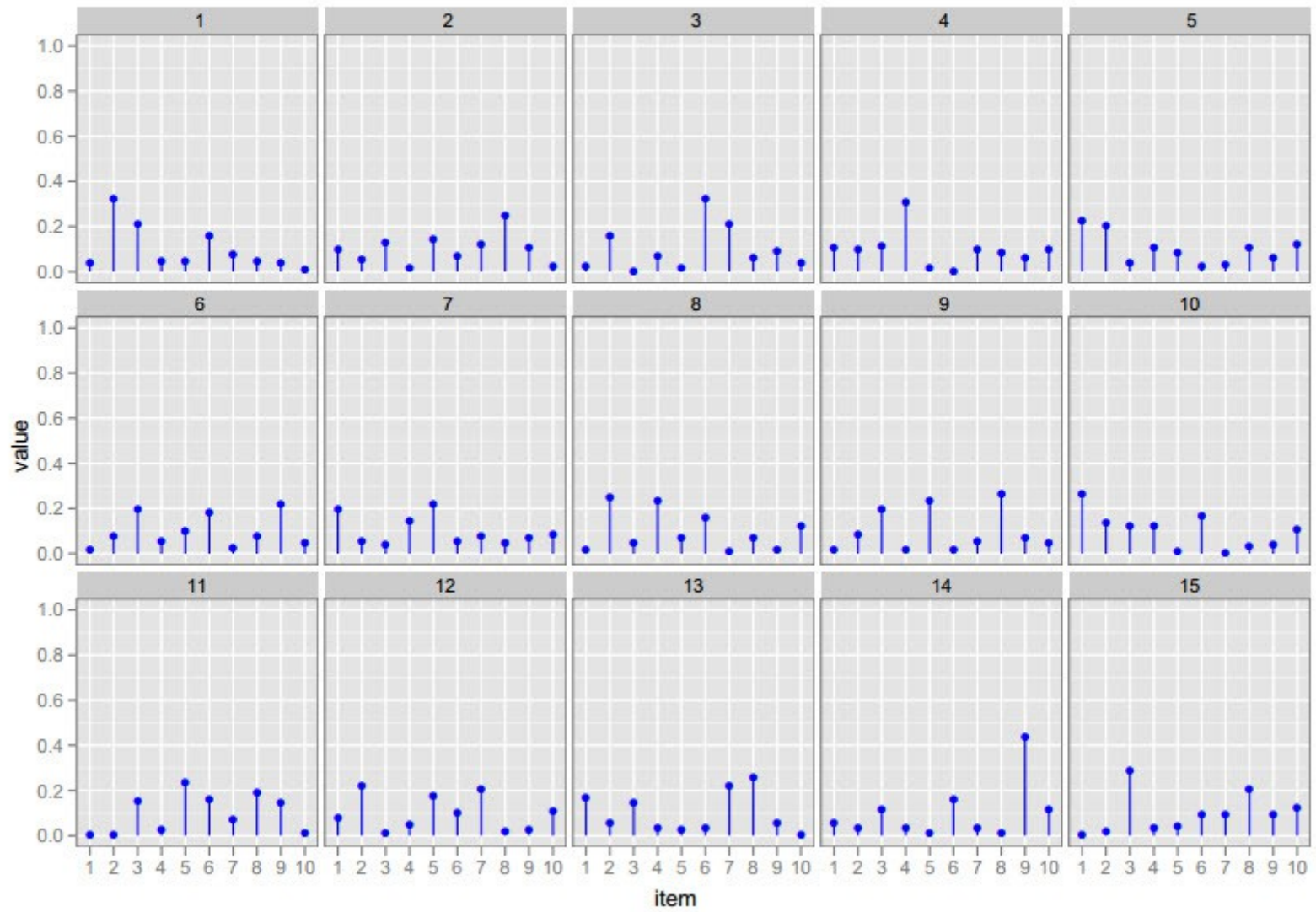


(b) Document Assignments to Topics

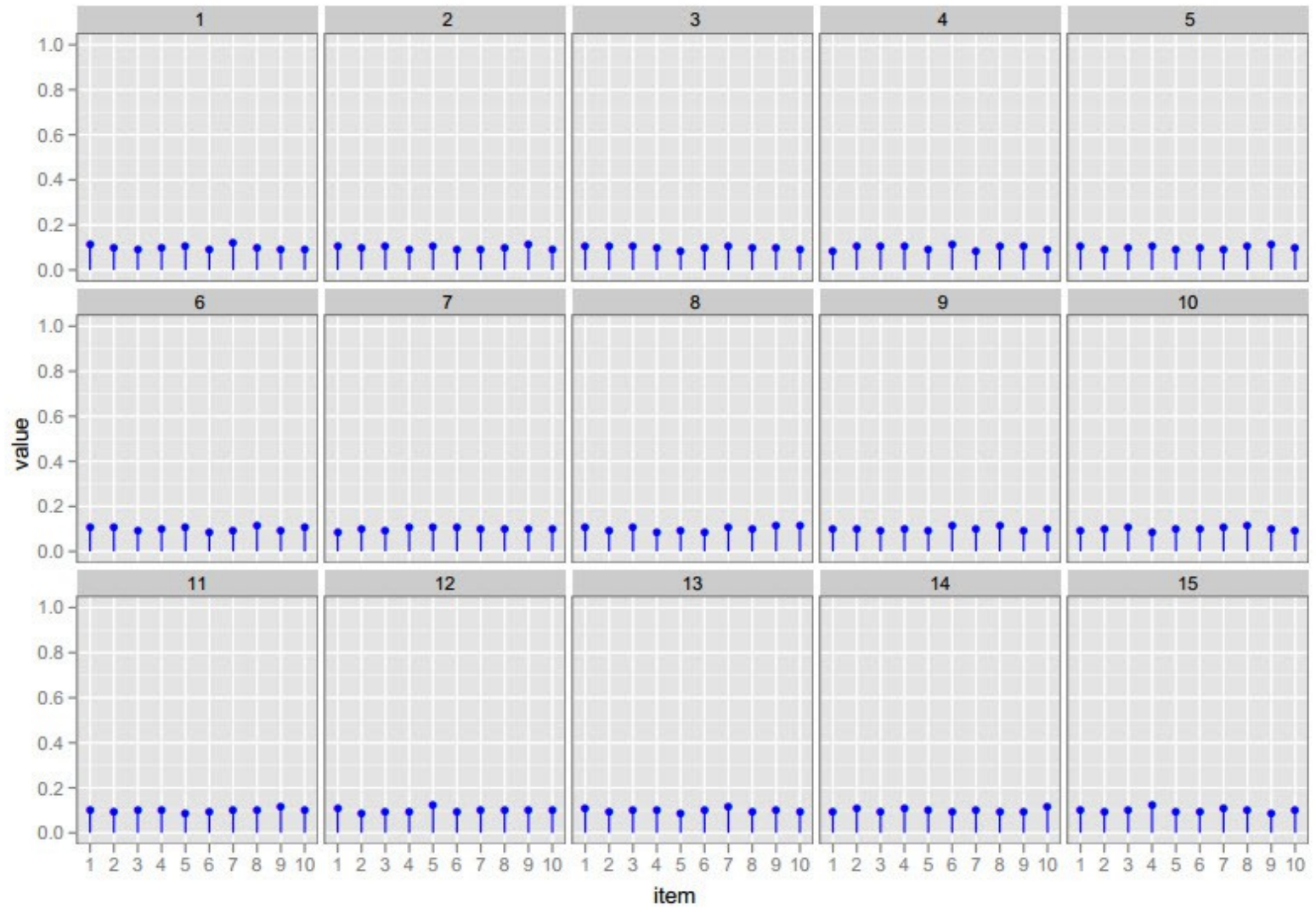
Models

- pLSI
- LDA
- CTM

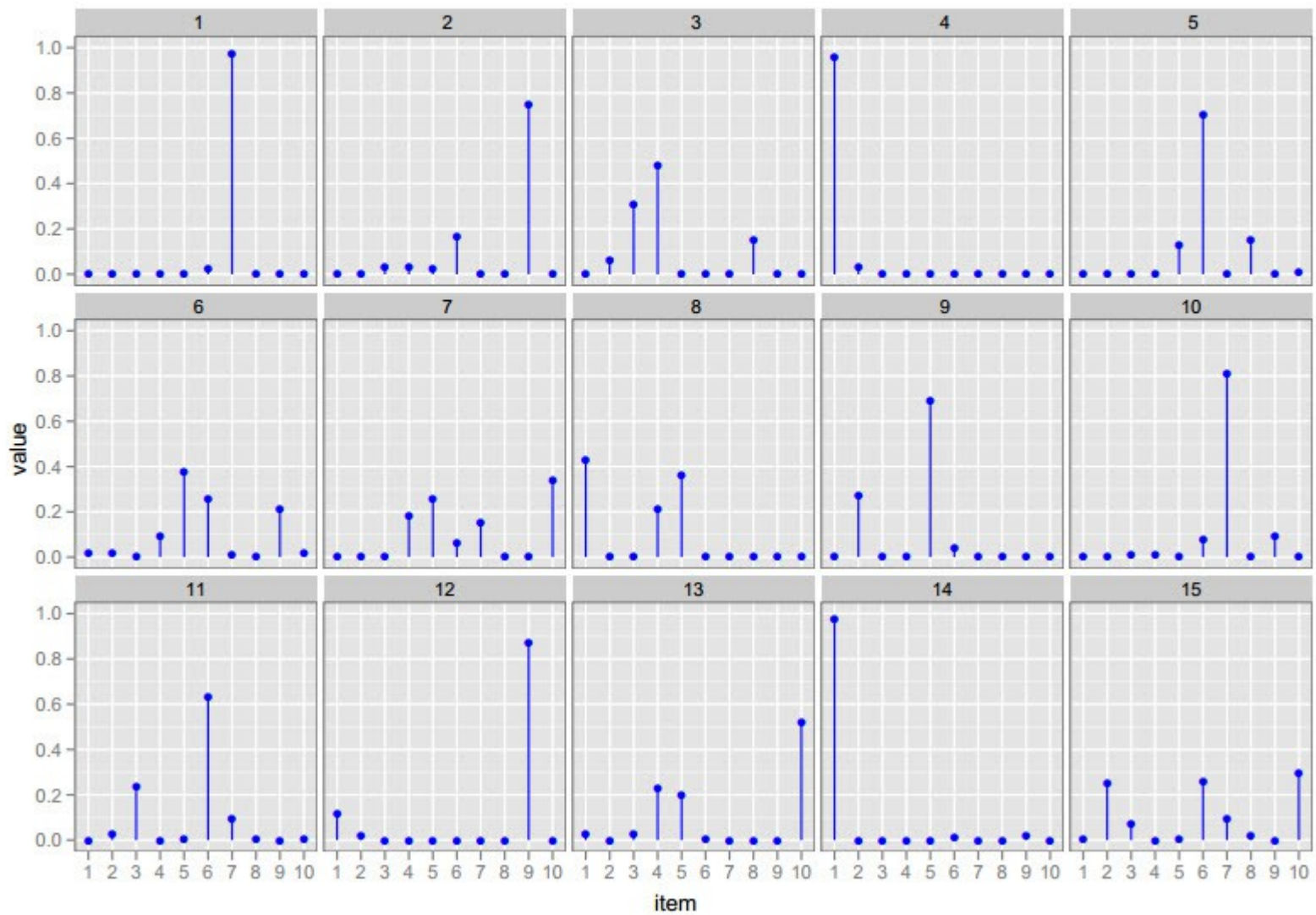
$$\alpha = 1$$



$\alpha = 100$



$\alpha = 0.1$



LDA

(fitting the model)

- Initial assignment: go through each document in the corpus and to each word assign a random topic from the set of topics T .
- For each document go through each word and for each topic calculate:
 - (a) the probability of topic given the document
 - (b) and the probability of the word given the topic
- reassign each word an new topic, chose the topic with probability $(a)*(b)$
- repeat until convergence

Measuring Performance

- External Tasks
- Held-out likelihood
- Humans?

$$\text{perplexity}(\text{test set } w) = \exp\left\{-\frac{\mathcal{L}(w)}{\text{count of tokens}}\right\}$$

CORPUS	TOPICS	LDA	CTM	PLSI
NEW YORK TIMES	50	-7.3214 / 784.38	-7.3335 / 788.58	-7.3384 / 796.43
	100	-7.2761 / 778.24	-7.2647 / 762.16	-7.2834 / 785.05
	150	-7.2477 / 777.32	-7.2467 / 755.55	-7.2382 / 770.36
WIKIPEDIA	50	-7.5257 / 961.86	-7.5332 / 936.58	-7.5378 / 975.88
	100	-7.4629 / 935.53	-7.4385 / 880.30	-7.4748 / 951.78
	150	-7.4266 / 929.76	-7.3872 / 852.46	-7.4355 / 945.29

MP and TLO

- Word Intrusion and Model Precision
- Topic Intrusion and Topic Log Odds

Word Intrusion

1 / 10	floppy	alphabet	computer	processor	memory	disk
2 / 10	molecule	education	study	university	school	student
3 / 10	linguistics	actor	film	comedy	director	movie
4 / 10	islands	island	bird	coast	portuguese	mainland

Topic Intrusion

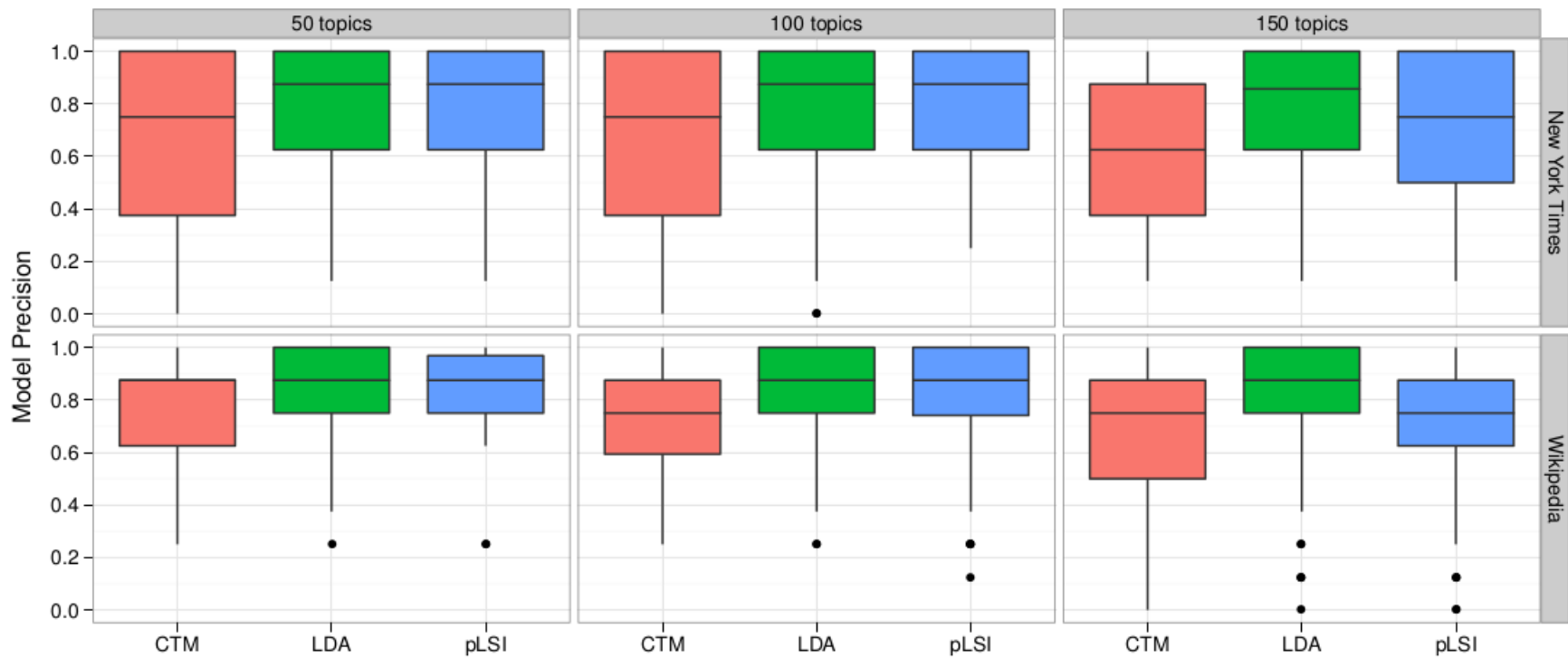
6 / 10	DOUGLAS_HOFSTADTER							
	Douglas Richard Hofstadter (born February 15, 1945 in New York, New York) is an American academic whose research focuses on consciousness, thinking and creativity. He is best known for " , first published in Show entire excerpt							
student	school	study	education	research	university	science	learn	
human	life	scientific	science	scientist	experiment	work	idea	
play	role	good	actor	star	career	show	performance	
write	work	book	publish	life	friend	influence	father	

$$\text{MP}_k^m = \sum_s \mathbf{1}(i_{k,s}^m = \omega_k^m) / S.$$

$$\text{TLO}_d^m = (\sum_s \log \hat{\theta}_{d,j_{d,*}^m}^m - \log \hat{\theta}_{d,j_{d,s}^m}^m) / S.$$

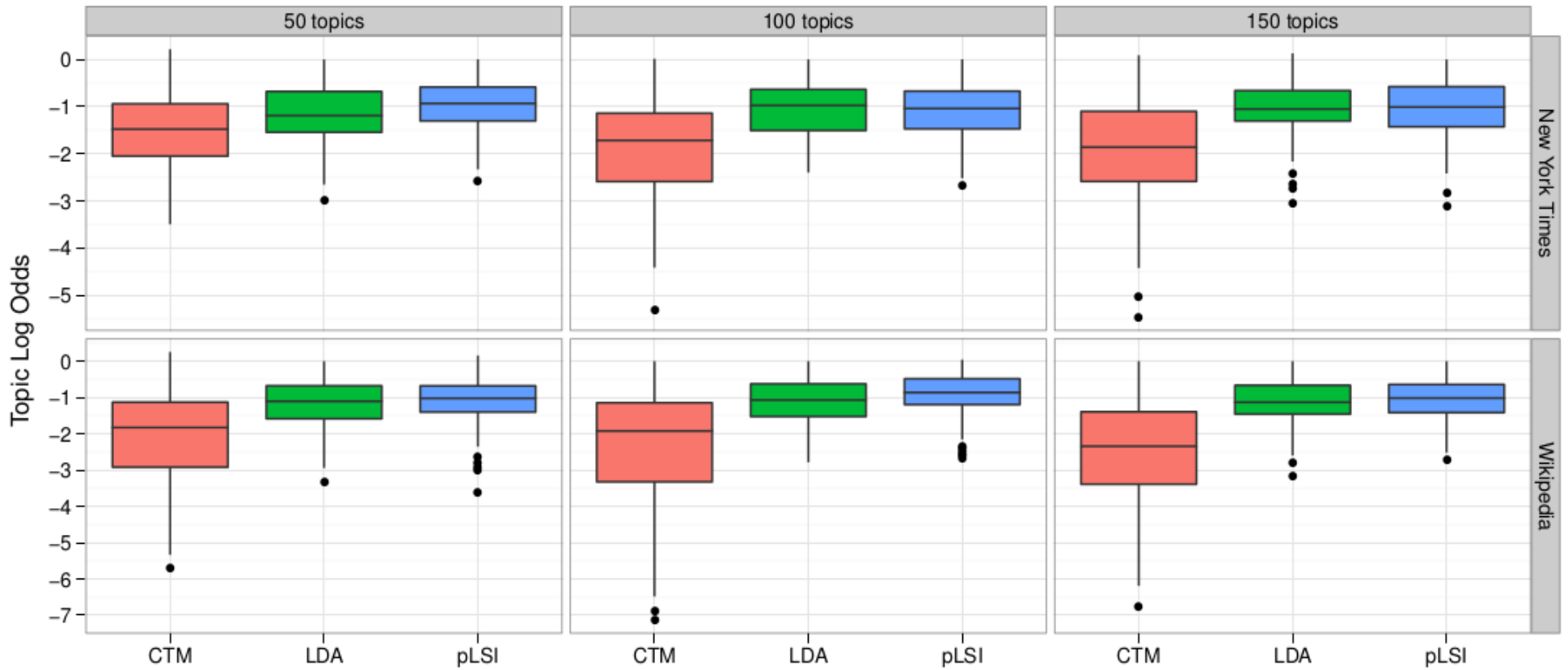
Results

- Model Precision

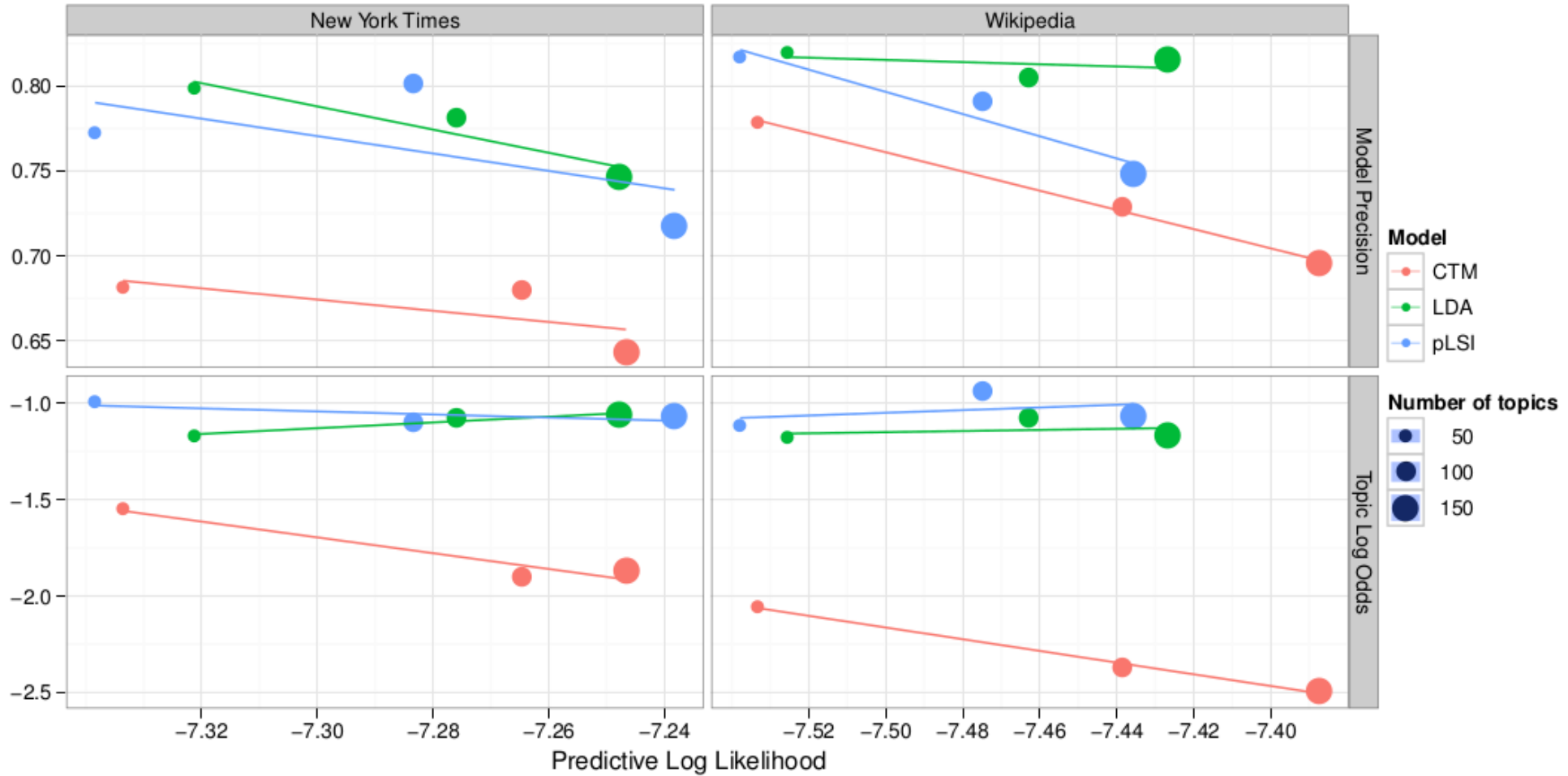


Results

- TLO



Results



Summary

- The use of metrics evaluating real world task performance
- Lower perplexity (higher likelihood) is not necessarily correlated to better coherence of topics

References

- <http://www.umiacs.umd.edu/~jbg/docs/nips2009-rtl.pdf>
- http://www.inf.ed.ac.uk/teaching/courses/nlu/lectures/nlu_l03-lda.pdf