

A Utility Model of Authors in the Scientific Community

Yanchuan Sim, School of Computer Science, Carnegie Mellon University, ysim@cs.cmu.edu

Bryan R. Routledge, Tepper School of Business, Carnegie Mellon University, routledge@cmu.edu

Noah A. Smith, Computer Science & Engineering, University of Washington, nasmith@cs.washington.edu

Introduction

We introduce a probabilistic model of authoring scientific papers where:

- authors have individual preferences,
- writing a paper requires trading off among the preferences of authors as well as extrinsic rewards in the form of community response to their papers,
- preferences (of individuals and the community) and tradeoffs vary over time.

Our model lead to improved predictive accuracy of citations given texts and texts given authors. Further, our model's posterior suggests an interesting relationship between seniority and author choices.

Utility function of an author

$$U(\theta_d) = \kappa_a y_d - \frac{1}{2} \|\theta_d - (\eta_a + \epsilon_{d,a})\|_2^2$$

Document content Author preferences

Author trade-offs Response to document

We assume that author a is an optimizer: when writing document d , she seeks to increase response y_d while keeping the contents θ_d "close" to her preferences, η_a . κ_a is a trade-off between the extrinsic (citation-seeking) and intrinsic (preference-satisfying) objectives:

- If κ_a is large, a might be understood as a citation maximizing agent;
- If κ_a is small, a might appear to care more about writing certain papers than citations.

Assuming a linear model for y_d , the expected utility is

$$\mathbb{E}[U(\theta_d)] = \kappa_a \beta^\top \theta_d - \frac{1}{2} \|\theta_d - (\eta_a + \epsilon_{d,a})\|_2^2$$

Response largely driven by content

An author's decision will therefore be

$$\hat{\theta}_d = \arg \max_{\theta} \kappa_a \beta^\top \theta - \frac{1}{2} \|\theta - (\eta_a + \epsilon_{d,a})\|_2^2$$

Optimality implies that θ_d solves the first-order equations

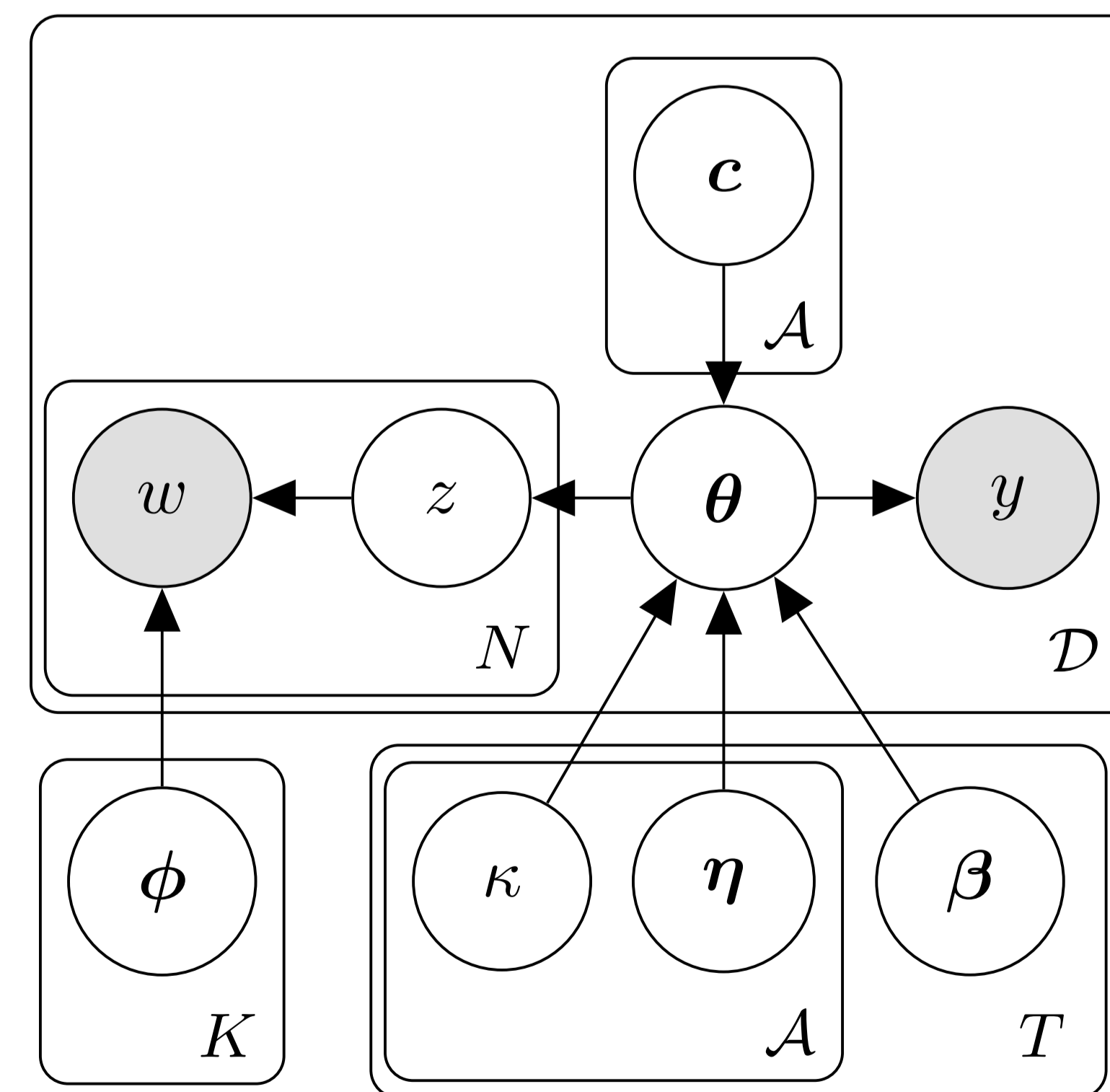
$$\kappa_a \beta_j - (\hat{\theta}_{d,j} - (\eta_{a,j} + \epsilon_{d,a,j})) = 0, \forall 1 \leq j \leq K$$

A probabilistic model incorporating utility

With multiple co-authors, the joint utility is simply the average of the group's utility, and the "cost" term for each author scaled by $c_{d,a}$ a fractional contribution term. We treat document content, θ_d , as a mixture of topics similar to latent Dirichlet allocation (LDA). Instead, θ_d is constrained as a result of the utility function

$$\theta_d \sim \mathcal{N} \left(\sum_{a \in \mathbf{a}_d} \kappa_a \beta + c_{d,a} \eta_a, \|\mathbf{c}_d\|_2^2 \mathbf{I} \right)$$

We extend our model to allow variations at different time-steps for β , η , and κ . To model the temporal dynamics, we place a multivariate Gaussian prior on the variables' differences between consecutive time steps.

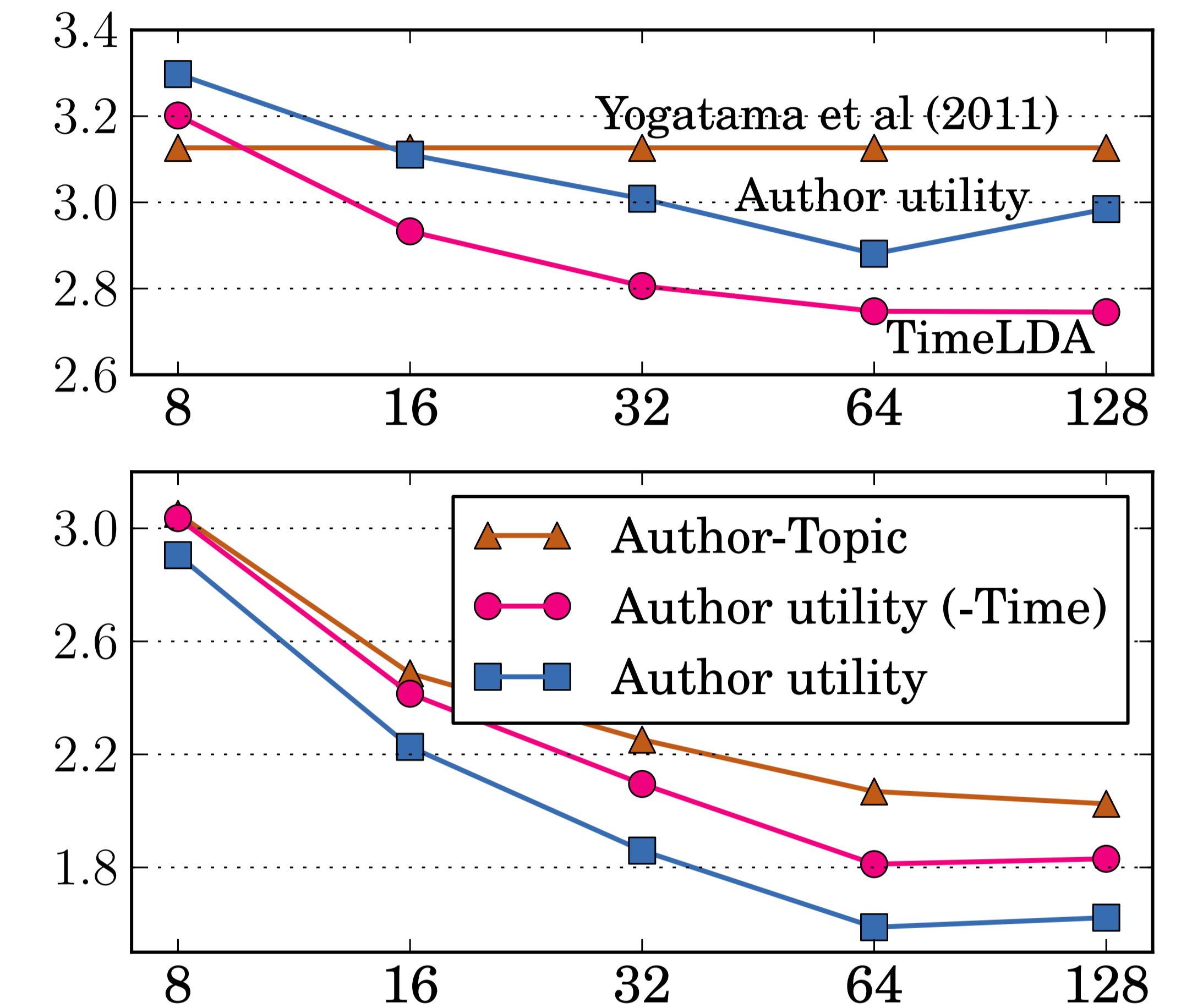


We adopt a Bayesian approach parameter estimation, using Monte Carlo EM to perform approximate inference on our model.

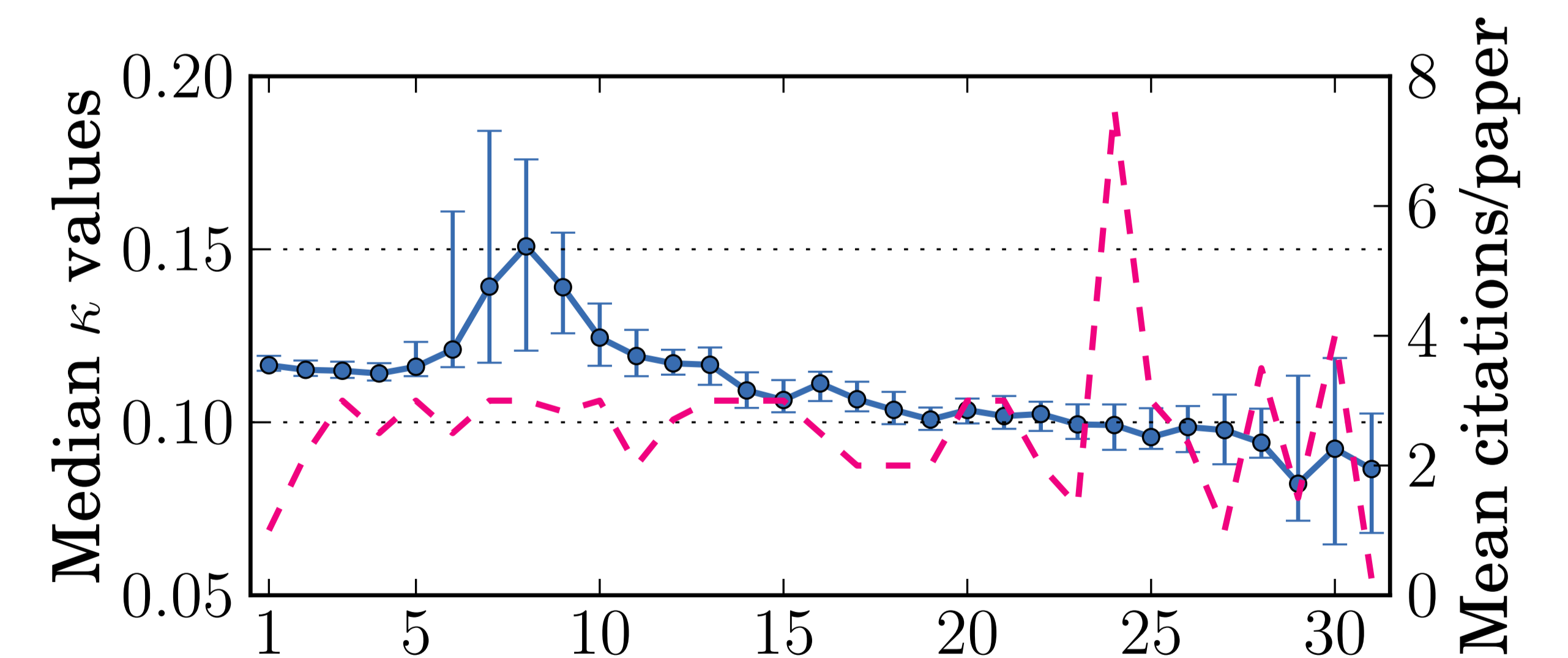
Data. We used a subset of the ACL anthology corpus, focusing on conference papers published between 1980 and 2010. The sub-corpus contains 5,498 documents authored by 5,575 scientists.

Experiments

Predicting citations and words. The mean absolute error (in citation counts, *top*) and held-out perplexity ($\times 10^3$, *bottom*) plotted against varying number of topics K .



Tradeoffs and Seniority. Plot of authors' median κ (blue, solid) and mean citation counts (magenta, dashed) against their academic age. The Spearman's rank correlation between κ of an author and her age is -0.870, with p -value $< 10^{-5}$.



References

- Michal Rosen-Zvi, Thomas Griffiths, Mark Steyvers, and Padhraic Smyth. 2004. The Author-Topic Model for Authors and Documents. In *Proc. of UAI*.
- Yanchuan Sim, Bryan Routledge, and Noah A. Smith. 2015. The Utility of Text: The Case of Amicus Briefs and the Supreme Court. In *Proc. of AACL*.
- Dani Yogatama, Michael Heilman, Brendan O'Connor, Chris Dyer, Bryan R. Routledge, and Noah A. Smith. 2011. Predicting a Scientific Community's Response to an Article. In *Proc. of EMNLP*.