

Statistical Structure in Natural Language

Tom Jackson

April 4th PACM seminar

Advisor: Bill Bialek



Some Linguistic Questions

- Why does language work so well?
 - Unlimited capacity and flexibility
 - Requires little conscious computational effort
- Why does language work at all?
 - Wittgenstein's dilemma
 - Rather accurate transmission of meaning from one mind to another
- How do we acquire language so easily?
 - Poverty of the stimulus
 - Everyone learns 'the same' language

Some Linguistic Answers



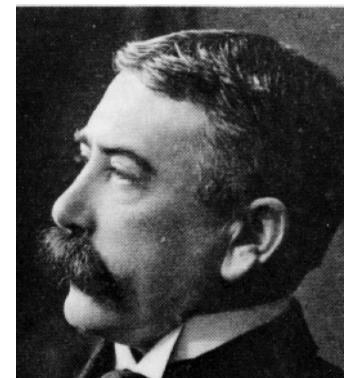
Noam Chomsky

- Language as ‘words and rules’
 - Semantic and syntactic components of language
- Hierarchical structure
 - Morphology (intra-word)
 - Syntax (intra-sentence)
 - Compositional rules and paragraph structure
- Chomsky: Grammar must be ‘universal’
- Nowak: evolution of syntactic communication
 - More efficient to express an increasing number of concepts using combinatorial rules.

Computational Linguistics

- Language performance vs. language competence
 - The ‘Chomskyan Revolution’
 - Linguistics as psychology
 - Corpus-based methods outside the mainstream
- Inherent drawbacks
 - No access to contextual meaning
 - Word meaning only defined through relationships to other words
 - Vive la différence!

Ferdinand de
Saussure



- Relevant for applications
- What are appropriate metrics to use?

Information Theory

- Entropy: information capacity

$$H(X) = - \sum_{\{x\}} p_x \log p_x$$



Claude Shannon

- Mutual information: Information X provides about Y and vice versa.

$$I(X;Y) = H(X) + H(Y) - H(X \cdot Y)$$

$$= \sum_{\{x,y\}} p_{xy} \log \frac{p_{xy}}{p_x p_y}$$

Naïve Entropy Estimation

$$\hat{H}(\{n_i\}) = - \sum_{i=1}^K \frac{n_i}{N} \log \frac{n_i}{N}$$

- Pro: Entirely straightforward
- Con: Useless for interesting case K<N
 - In the undersampled regime, just because we don't see an event doesn't mean it has zero probability.
 - Analogous to overfitting a curve to data
 - More appropriate to consider a "smooth" class of distributions
- Problem: How to apply Occam's principle in an unbiased way?

Bayesian Inference

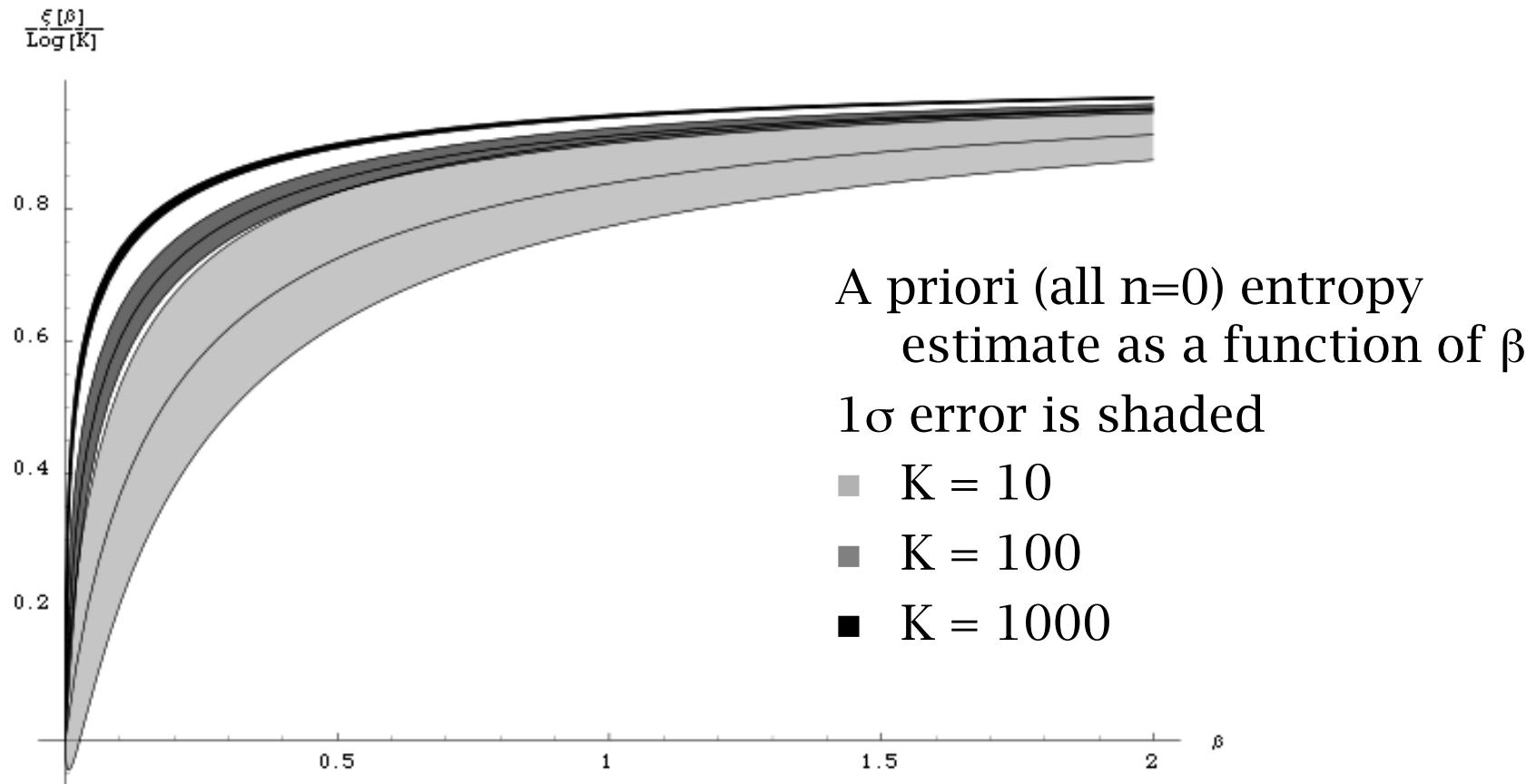
- Bayes' Theorem: $\Pr(\text{model} \mid \text{data}) = \frac{\Pr(\text{data} \mid \text{model})\Pr(\text{model})}{\Pr(\text{data})}$
- Here our prior is on the space of distributions...

$$\frac{1}{Z} \delta\left(1 - \sum_{i=1}^K q_i\right) \prod_{i=1}^K q_i^{\beta-1}$$

- Generalize uniform prior to the Dirichlet prior parameterized by β .

$$\langle q_i \rangle_\beta = \frac{n_i + \beta}{N + K\beta}$$

Bias in the Dirichelet Prior



- Prior highly biased for fixed values of β
- Obtain an unbiased prior by averaging over β

The estimator

$$\hat{S}^m = \frac{1}{Z(\{n_i\})} \int d\beta \frac{d\xi}{d\beta} \rho(\{n_i\}, \beta) \langle S^m \rangle_\beta$$

$$\rho(\{n_i, \beta\}) = \frac{\Gamma(K\beta)}{\Gamma(N + K\beta)} \prod_{i=1}^K \frac{\Gamma(n_i + \beta)}{\Gamma(\beta)}$$

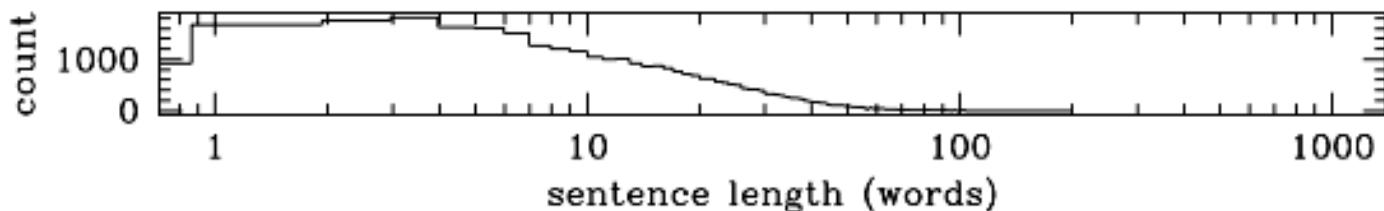
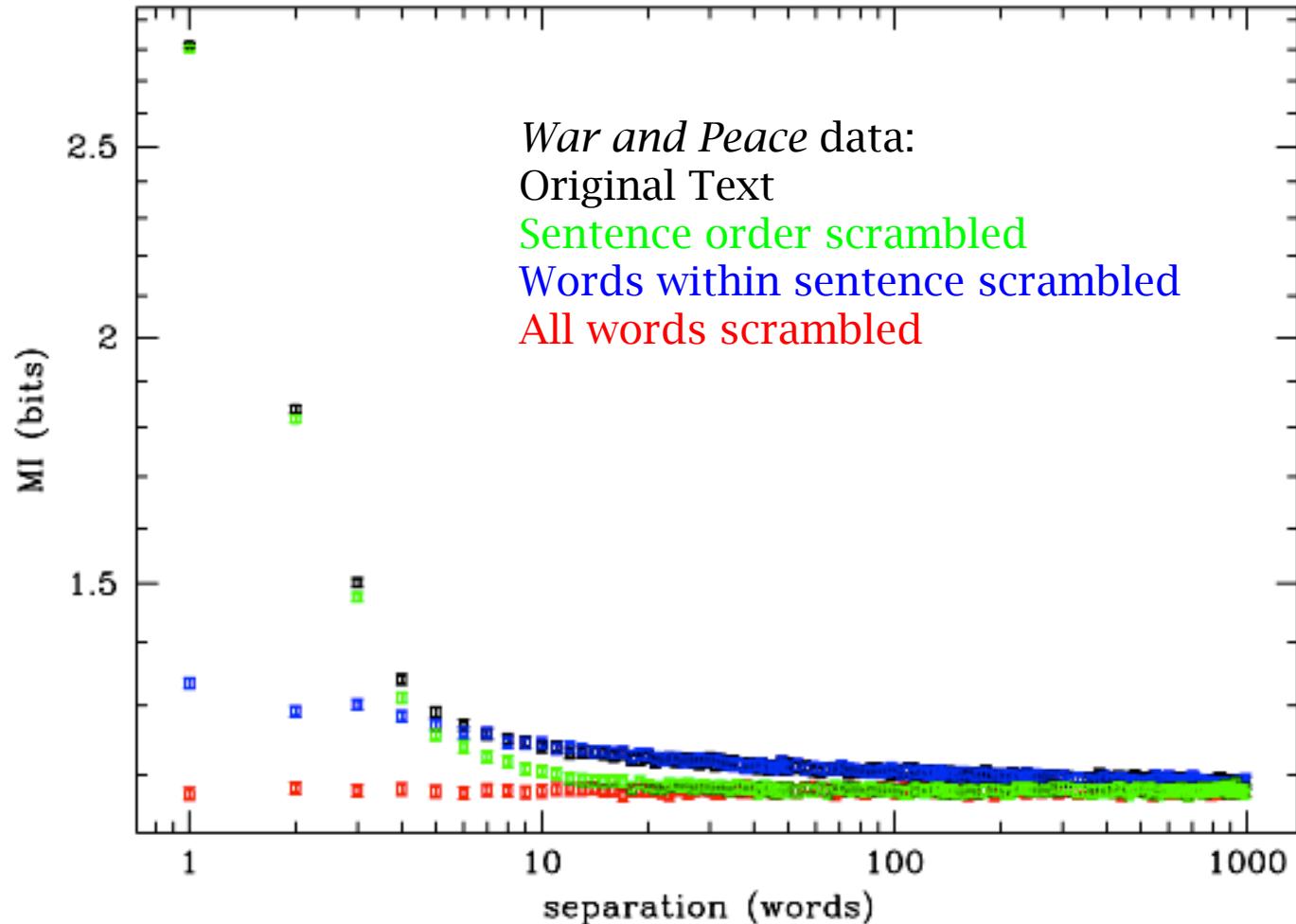
$$\langle S \rangle_\beta = - \sum_{i=1}^K \frac{n_i + \beta}{N + K\beta} (\Psi(n_i + \beta + 1) - \Psi(N + K\beta + 1))$$

- Computational complexity still linear in N

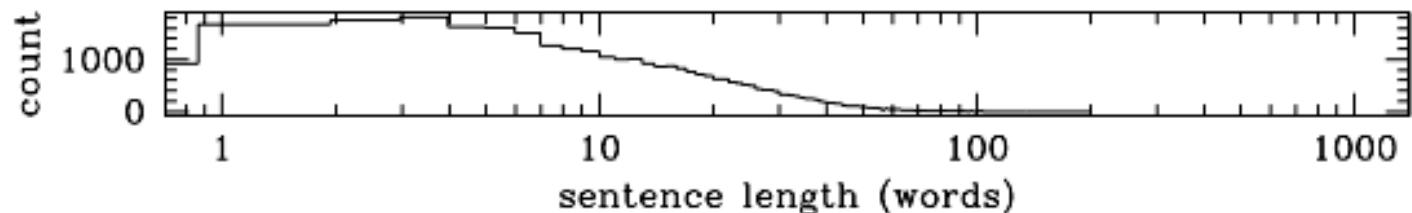
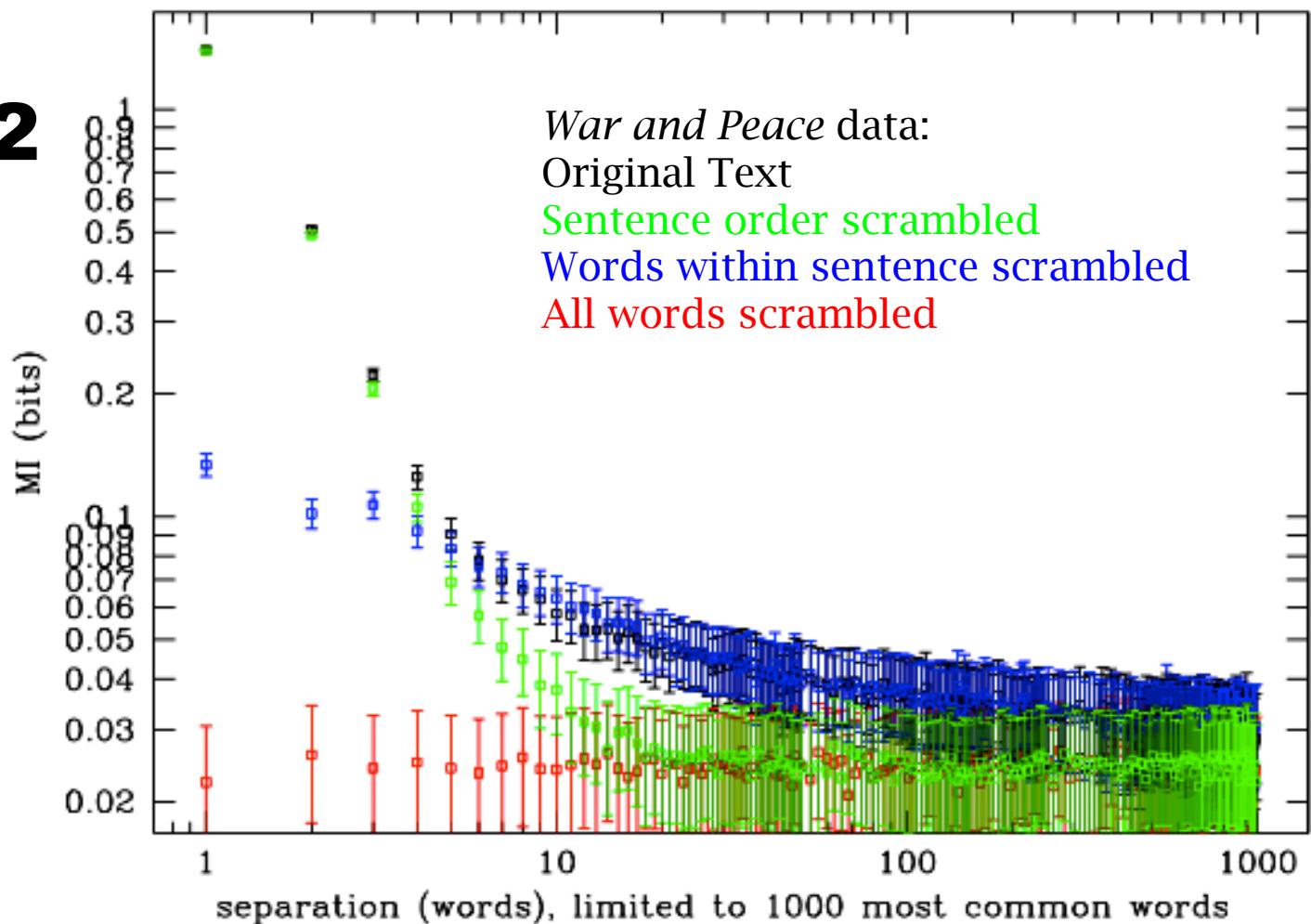
Getting to Work

- Summer '03 research project
- MI estimation algorithm implemented in ~5K lines of C code
- Texts used:
 - Project Gutenberg's *Moby Dick*
 - Project Gutenberg's *War and Peace*
 - K=18004, N = 567657
 - Reuters news articles corpus
 - K=40696, N = 2478116

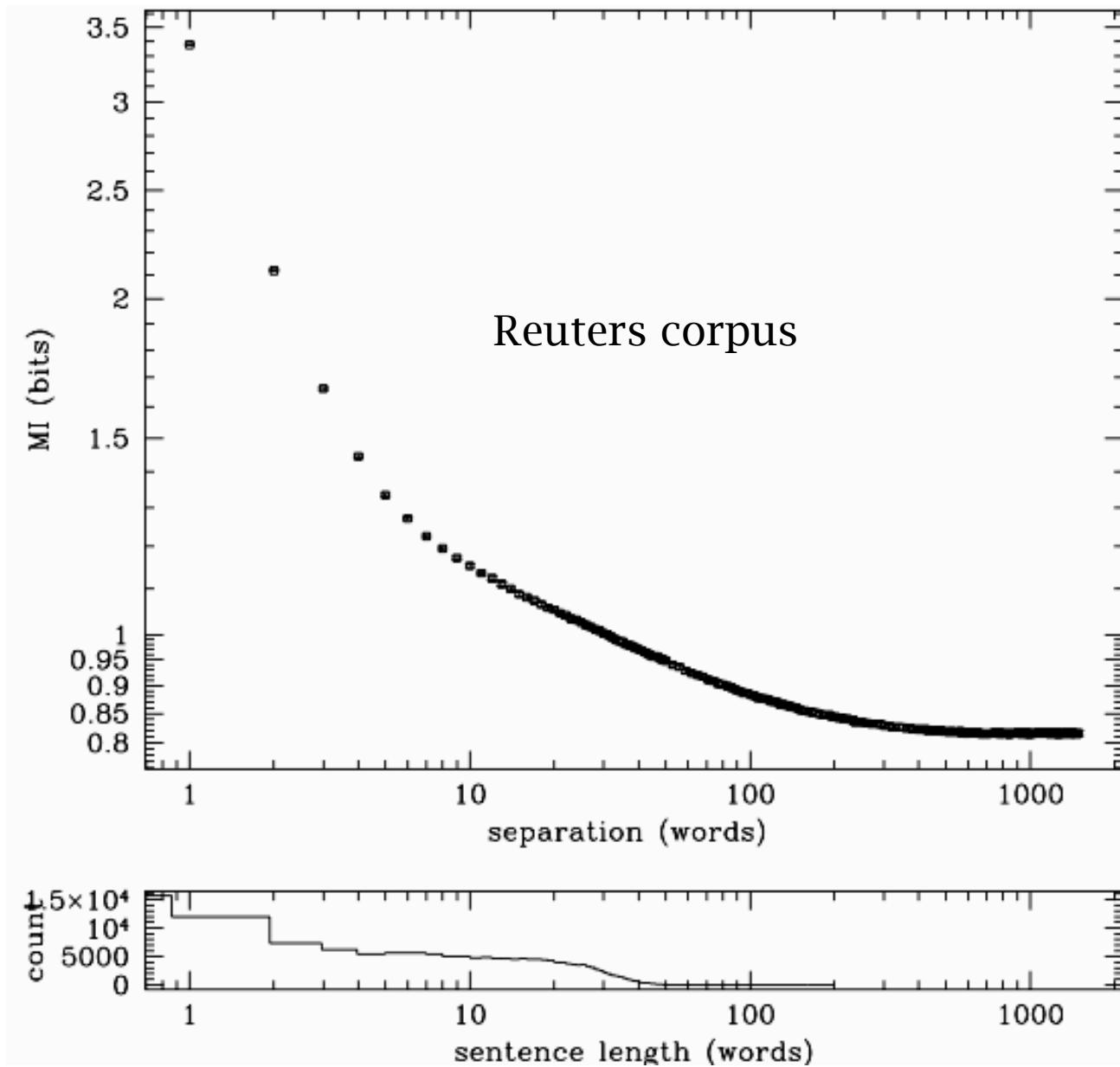
Results 1



Results 2



Results 3



Information Bottleneck

- Are we really observing a crossover between syntactic and semantic constraints?
 - What “carries” the information at a given separation?
- Information Bottleneck method: compress X into X' while maximizing $I(X';Y)$.
 - Two algorithms to determine $p(x'|x)$:
 - Simulated Annealing (top down clustering)
 - Agglomerative (bottom up clustering)
- Difficult to implement in this case
 - Start with $K \sim 10^5$; expect relevant clusters at $K' \sim 10$
 - Considerable degeneracy from *hapax legomena*

Conclusions

- Robust estimation of probability distribution functionals is possible even with undersampling
- Pairwise mutual information displays power law behavior, with breaks at semantically significant length scales

Possible Applications

- Semantic extraction
 - Automatic document summarizing
 - Search engines
- Data compression
 - Shannon compression theorem
 - Existing algorithms work at finite range only

For further reading

- Language
 - *Words and Rules: the ingredients of language*, by Steven Pinker
 - Nature 417:611
- Entropy estimation
 - Physical review E 52:6841
 - xxx.lanl.gov/abs/physics/0108025



End