

# Perplexity on Reduced Corpora

— Analysis of Cutoff by Power Law

Hayato Kobayashi  
Yahoo Japan Corporation

# Cutoff

---

- ▶ Removing low-frequency words from a corpus
- ▶ Common practice to save computational costs in learning
- ▶ Language modeling
  - ▶ Needed even in a distributed environment, since the feature space of k-grams is quite large [Brants+ 2007]
- ▶ Topic modeling
  - ▶ Enough for roughly analyzing topics, since low-frequency words have a small impact on the statistics [Steyvers&Griffiths 2007]

# Question

---

- ▶ **How many low-frequency words can we remove while maintaining sufficient performance?**
  - ▶ More generally, how much can we reduce a corpus/model using a certain strategy?
- ▶ **Many experimental studies addressing the question**
  - ▶ [Stoleke 1998], [Buchsbaum+ 1998], [Goodman&Gao 2000], [Gao&Zhang 2002], [Ha+ 2006], [Hirsimaki 2007], [Church+ 2007]
  - ▶ Discussing trade-off relationships between the size of reduced corpus/model and its performance
- ▶ **No theoretical study!**

# This work

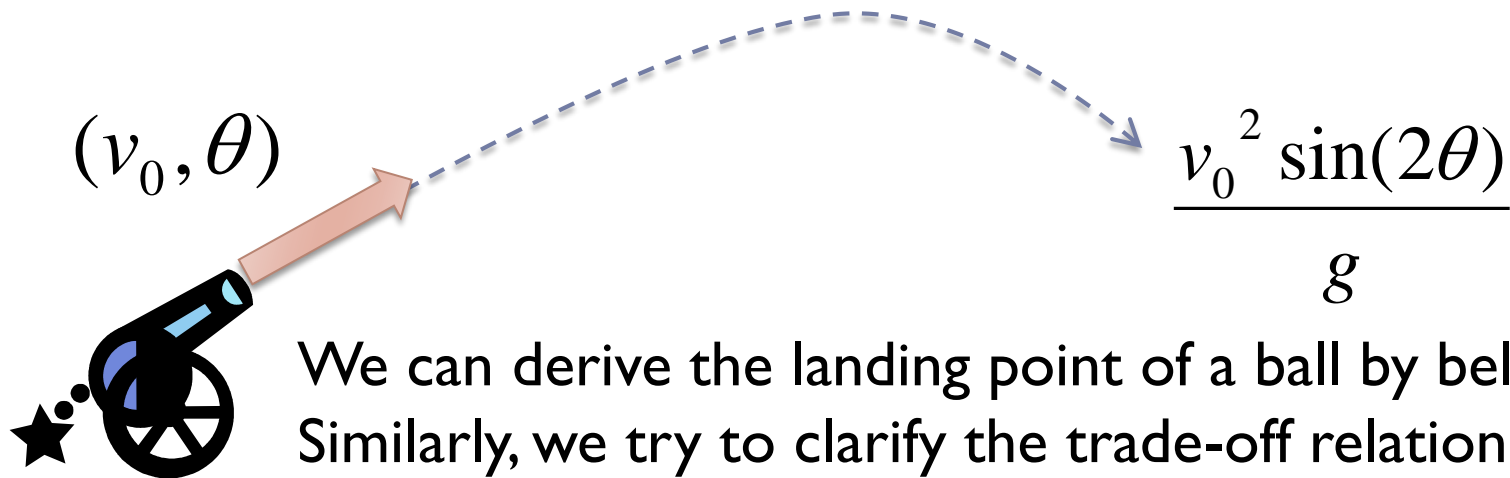
---

- ▶ First address the question from a theoretical standpoint
- ▶ Derive the trade-off formulae of the cutoff strategy for k-gram models and topic models
  - ▶ Perplexity vs. reduced vocabulary size
- ▶ Verify the correctness of our theory on synthetic corpora and examine the gap between theory and practice on several real corpora

# Approach

---

- ▶ Assume a corpus follows Zipf's law (power law)
  - ▶ Empirical rule representing a long-tail property in a corpus
- ▶ Essentially the same approach as in physics
  - ▶ Constructing a theory while believing experimentally observed results (e.g., gravity acceleration  $g$ )



We can derive the landing point of a ball by believing  $g$ . Similarly, we try to clarify the trade-off relationships by believing Zipf's law.

# Outline

---

- ▶ **Preliminaries**
  - ▶ Zipf's law
  - ▶ Perplexity (PP)
  - ▶ Cutoff and restoring
- ▶ PP of unigram models
- ▶ PP of k-gram models
- ▶ PP of topic models
- ▶ Conclusion

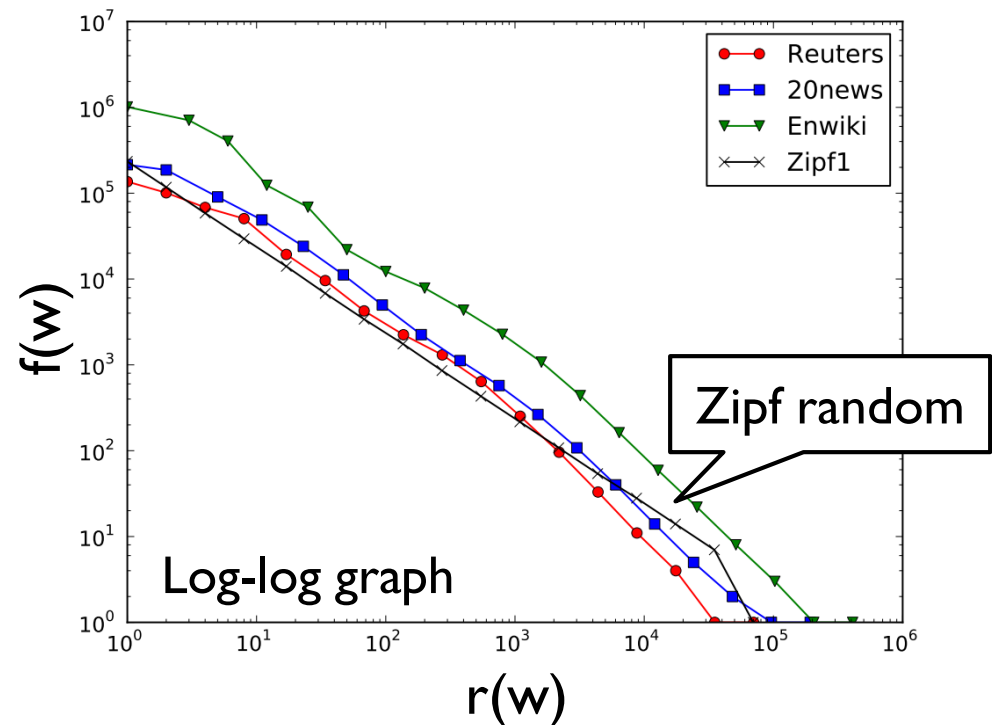
# Zipf's law

- ▶ Empirical rule discovered on real corpora [Zipf, 1935]
  - ▶ Word frequency  $f(w)$  is inversely proportional to its frequency ranking  $r(w)$

$$\underbrace{f(w)}_{\text{Frequency}} = \frac{\text{Max. frequency } C}{\underbrace{r(w)}_{\text{Frequency ranking}}}$$

(Linear on a log-log graph)

Real corpora roughly follow Zipf's law



# Perplexity (PP)

---

- ▶ Widely used evaluation measure of statistical models
  - ▶ Geometric mean of the inverse of the per-word likelihood on the held-out test corpus

$$PP := \left( \prod_{w \in \underline{\mathbf{w}}_{\tau}} \frac{1}{p(w)} \right)^{\frac{1}{N_{\tau}}}$$

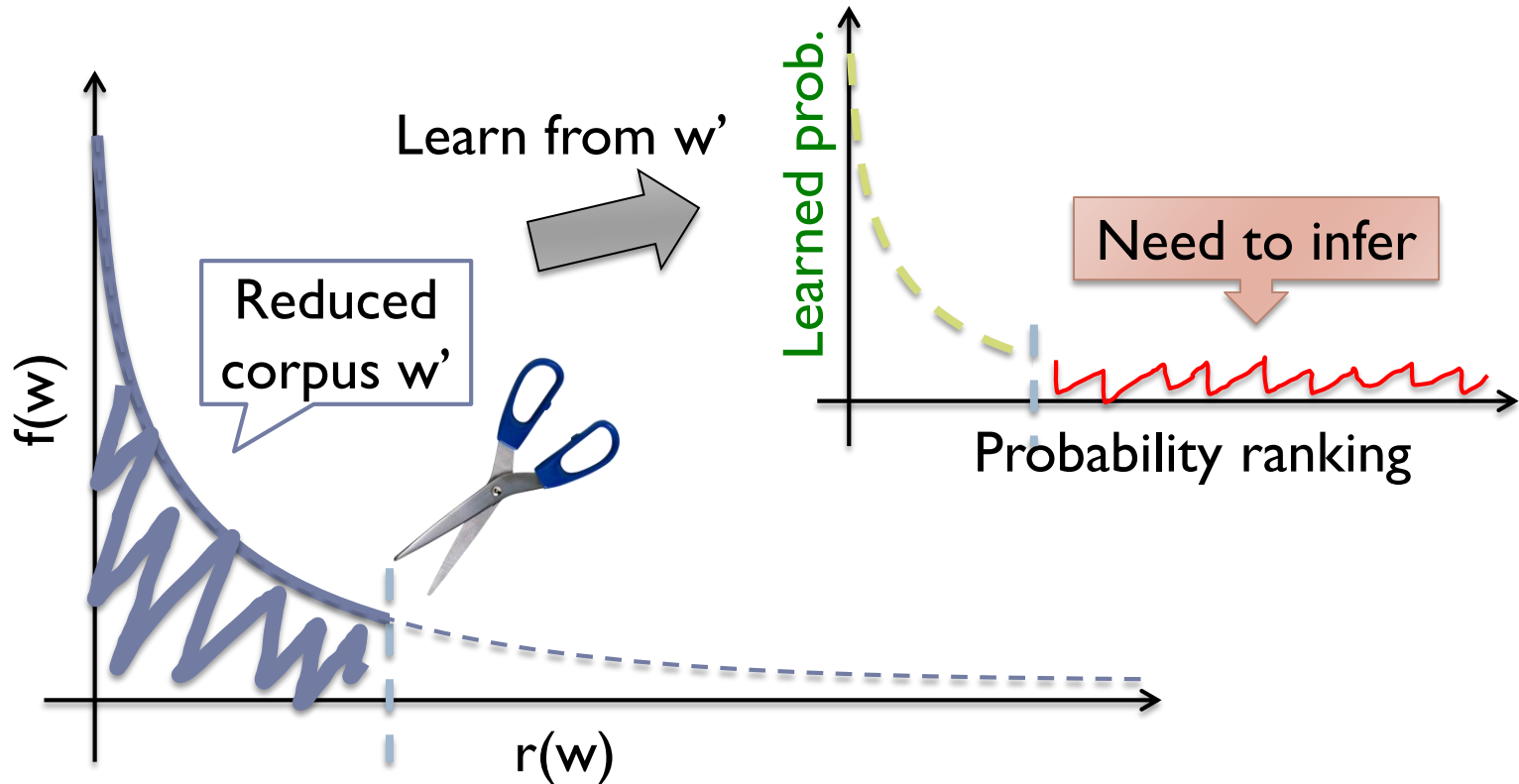
Test corpus      Corpus size

- ▶ PP means how many possibilities one has for estimating the next word
  - ▶ Lower perplexity means better generalization performance



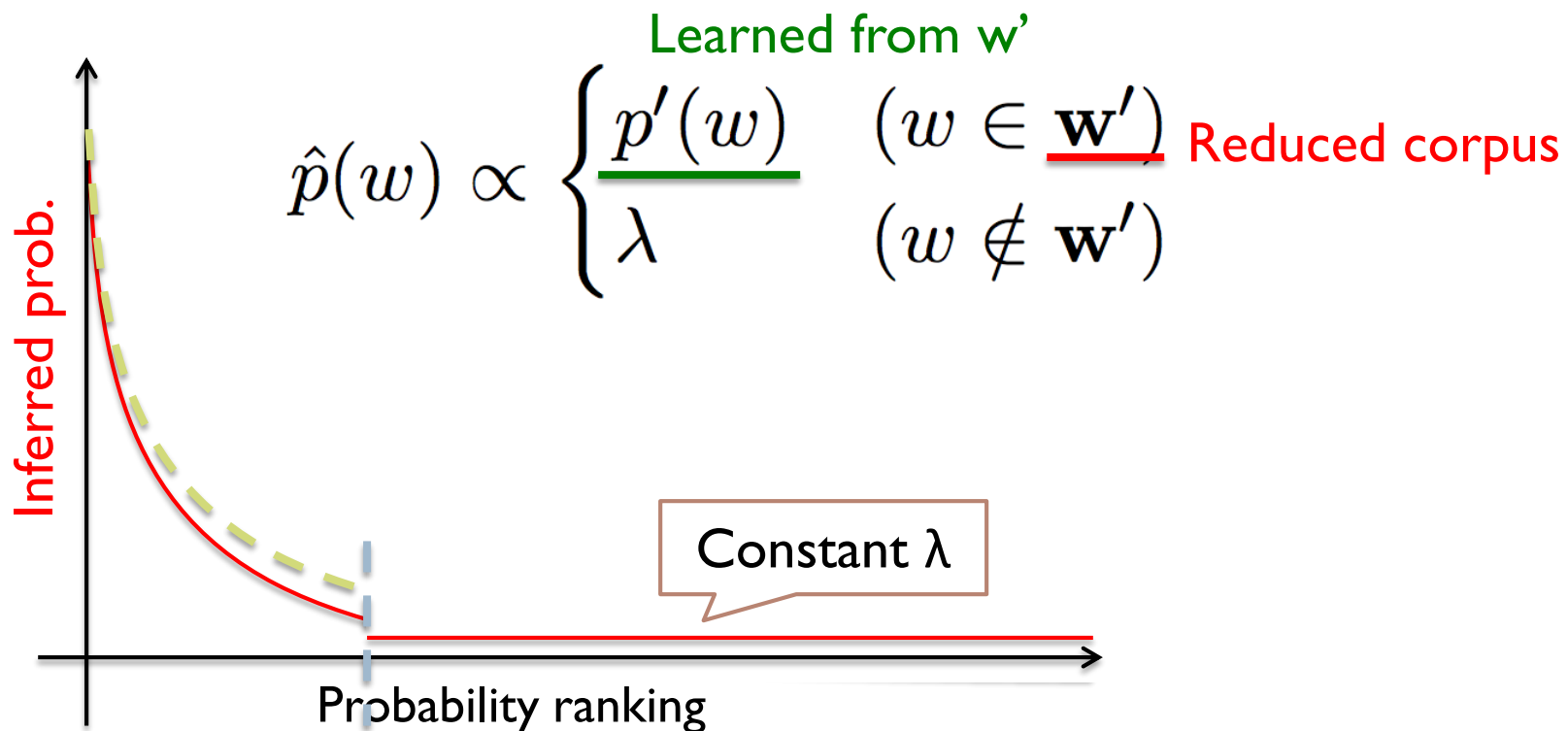
# Cutoff

- ▶ Removing low frequency words
  - ▶  $f(\text{remaining word}) \geq f(\text{removed word})$  holds



# Constant restoring

- ▶ Infer the prob. of the removed words as a constant
  - ▶ Approximate the result learned from the original corpus



# Outline

---

- ▶ Preliminaries
  - ▶ Zipf's law
  - ▶ Perplexity (PP)
  - ▶ Cutoff and restoring
- ▶ **PP of unigram models**
- ▶ PP of k-gram models
- ▶ PP of topic models
- ▶ Conclusion

# Perplexity of unigram models

---

- ▶ Predictive distribution of unigram models

$$p'(w') = \frac{f(w')}{\underline{N'}} \quad \text{Reduced corpus size}$$

- ▶ Optimal restoring constant

- ▶ Obtained by minimizing PP w.r.t. a constant  $\lambda$ , after substituting the restored probability  $\hat{p}(w)$  into PP

$$\lambda^* = \frac{\text{Corpus size} \quad \underline{N} - N'}{\text{Vocab. size} \quad \underline{(W - W')} N'}$$

# Theorem (PP of unigram models)

---

- ▶ For any reduced vocabulary size  $W'$ , the perplexity  $PP_1$  of the optimal restored distribution of a unigram model is calculated as

$$\hat{PP}_1(W') = H(W) \exp\left(\frac{B(W')}{H(W)}\right) \left(\frac{W - W'}{H(W) - H(W')}\right)^{1 - \frac{H(W')}{H(W)}}$$

$$\underline{H(X) := \sum_{x=1}^X \frac{1}{x}} \quad \text{Harmonic series}$$

$$\underline{B(X) := \sum_{x=1}^X \frac{\ln(x)}{x}} \quad \text{Bertrand series (special form)}$$

# Approximation of PP of unigrams

---

- ▶  $H(X)$  and  $B(X)$  can be approximated by definite integrals

$$H(X) \approx \ln X + \underline{\gamma} \quad \text{Euler-Mascheroni const.}$$

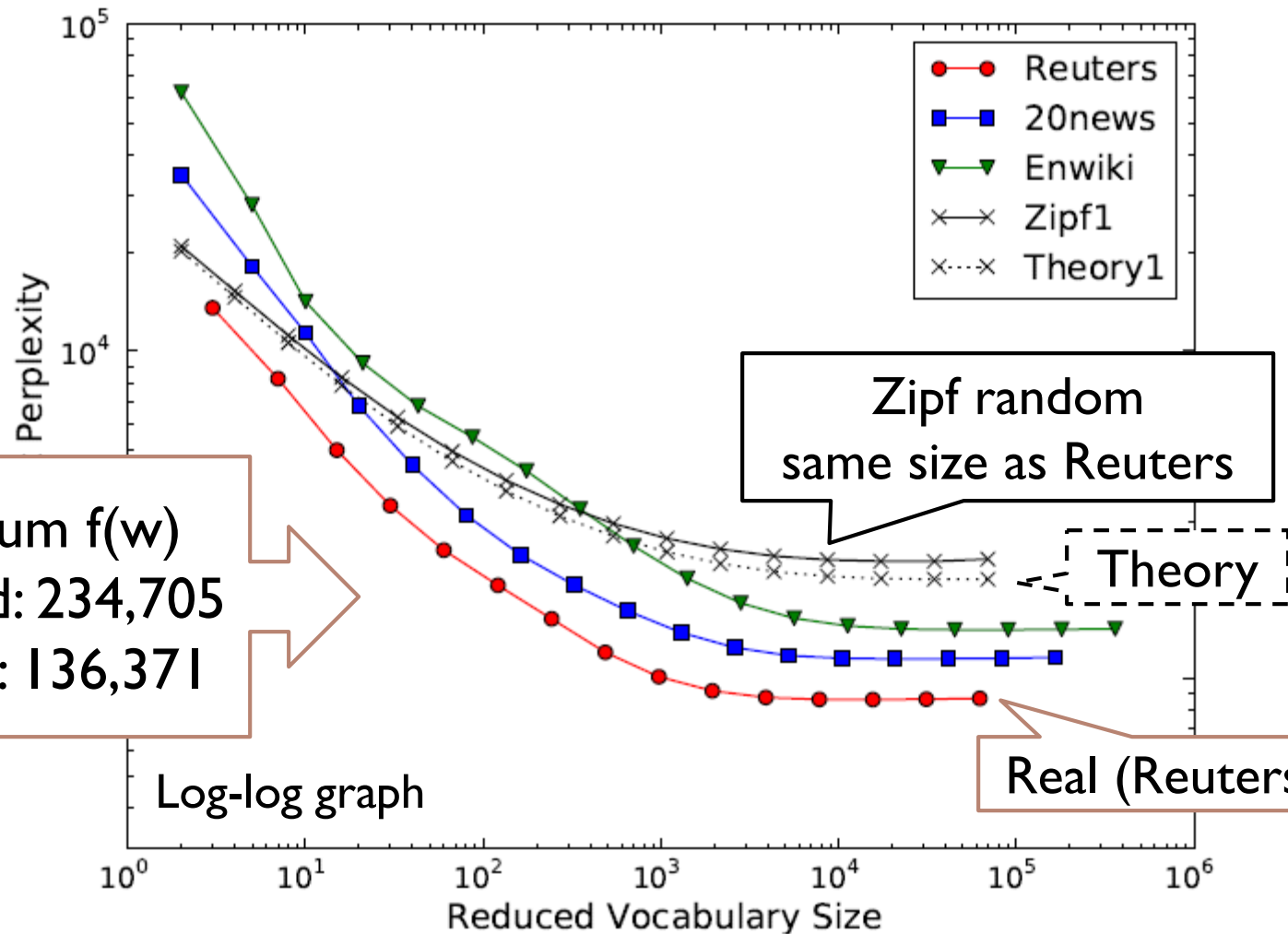
$$B(X) \approx \frac{1}{2} \ln^2 X$$

- ▶ Approximate formula  $\tilde{P}P_1(W')$  is obtained as

$$\tilde{P}P_1(W') = \sqrt{W} \ln W \exp \frac{(\ln \underline{W'} - \ln W)^2}{2 \ln W} \quad \text{Reduced vocab. size}$$

- ▶  $\tilde{P}P_1(W')$  is quasi polynomial (quadratic)
  - ▶ Behaves as a quadratic function on a log-log graph

# PP of unigrams vs. reduced vocab. size



Our theory is suited for inferring the growth rate of perplexity rather than the perplexity value itself

# Outline

---

- ▶ Preliminaries
  - ▶ Zipf's law
  - ▶ Perplexity (PP)
  - ▶ Cutoff and restoring
- ▶ PP of unigram models
- ▶ **PP of k-gram models**
- ▶ PP of topic models
- ▶ Conclusion



# Perplexity of k-gram models

---

- ▶ Simple model where k-grams are calculated from a random word sequence based on Zipf's law
- ▶ The model is “stupid”

- ▶ Bigram “is is” is quite frequent

$$p(\text{"is is"}) = p(\text{"is"})p(\text{"is"})$$

- ▶ Two bigrams “is a” and “a is” have the same frequency

$$p(\text{"is a"}) = p(\text{"is"})p(\text{"a"}) = p(\text{"a is"})$$

- ▶ Later experiment will uncover the fact that the model can roughly capture the behavior of real corpora

# Frequency of a k-gram

---

- ▶ Frequency  $f_k$  of a k-gram  $w_k$  is defined by

$$f_k(w_k) = \frac{C_k}{\underline{g_k(r_k(w_k))}} \text{ Decay function}$$

- ▶ Decay function  $g_2$  of bigrams is as follows

$$\begin{aligned}(g_2(i))_i &:= (g_2(1), g_2(2), g_2(3), \dots) \\ &= (1 \cdot 1, 1 \cdot 2, 2 \cdot 1, 1 \cdot 3, 3 \cdot 1, \dots) \\ &= (1, 2, 2, 3, 3, 4, 4, 4, 5, 5, 6, \dots)\end{aligned}$$

- ▶ Decay function  $g_k$  of k-grams is defined through its

inverse:

$$g_k^{-1}(\ell) := \sum_{n=1}^{\ell} d_k(n)$$

$$\underline{d_k(n) := \sum_{i_1 \cdot i_2 \cdots i_k = n} 1}$$

Piltz divisor function that represents # of divisors of n

# Exponent of k-gram distributions

---

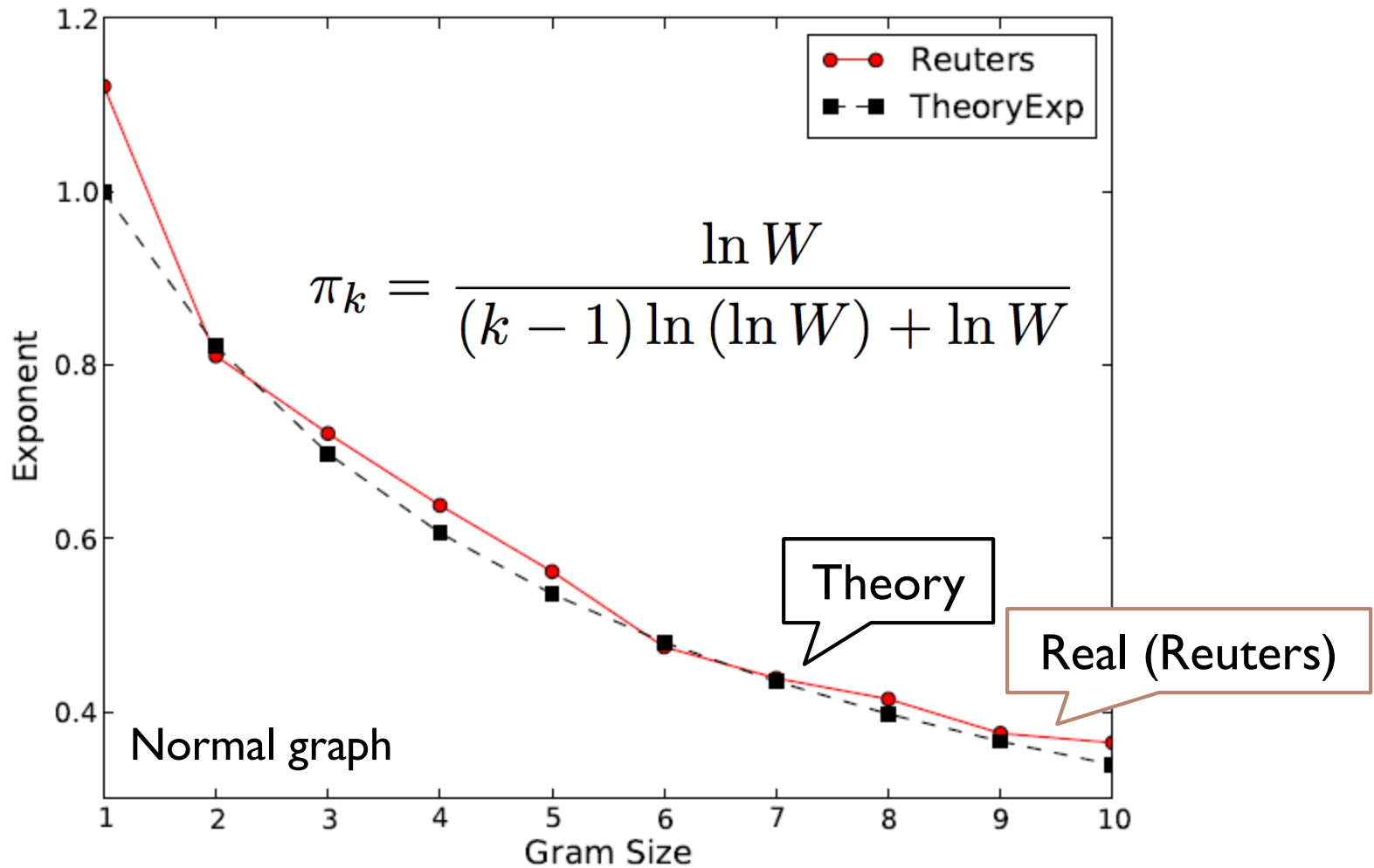
- ▶ Assume k-gram frequencies follow a power law
  - ▶ [Ha+ 2006] found k-gram frequencies roughly follow a power law, whose exponent  $\pi_k$  is smaller than 1 ( $k > 1$ )

$$f_k(w_k) \propto r_k(w_k)^{-\pi_k}$$

- ▶ Optimal exponent in our model based on the assumption
  - ▶ By minimizing the sum of squared errors between the inverse gradients  $g_k^{-1}(r)$  and  $r^{1/\pi_k}$  on a log-log graph

$$\pi_k = \frac{\ln W}{(k - 1) \ln (\ln W) + \ln W}$$

# Exponent of k-grams vs. gram size



## Corollary (PP of k-gram models)

- ▶ For any reduced vocabulary size  $W'$ , the perplexity of the optimal restored distribution of a k-gram model is calculated as

$$\hat{P}P_k(W') = H_{\pi_k}(W) \exp\left(\frac{B_{\pi_k}(W')}{H_{\pi_k}(W)}\right) \left(\frac{W - W'}{H_{\pi_k}(W) - H_{\pi_k}(W')}\right)^{1 - \frac{H_{\pi_k}(W')}{H_{\pi_k}(W)}}$$

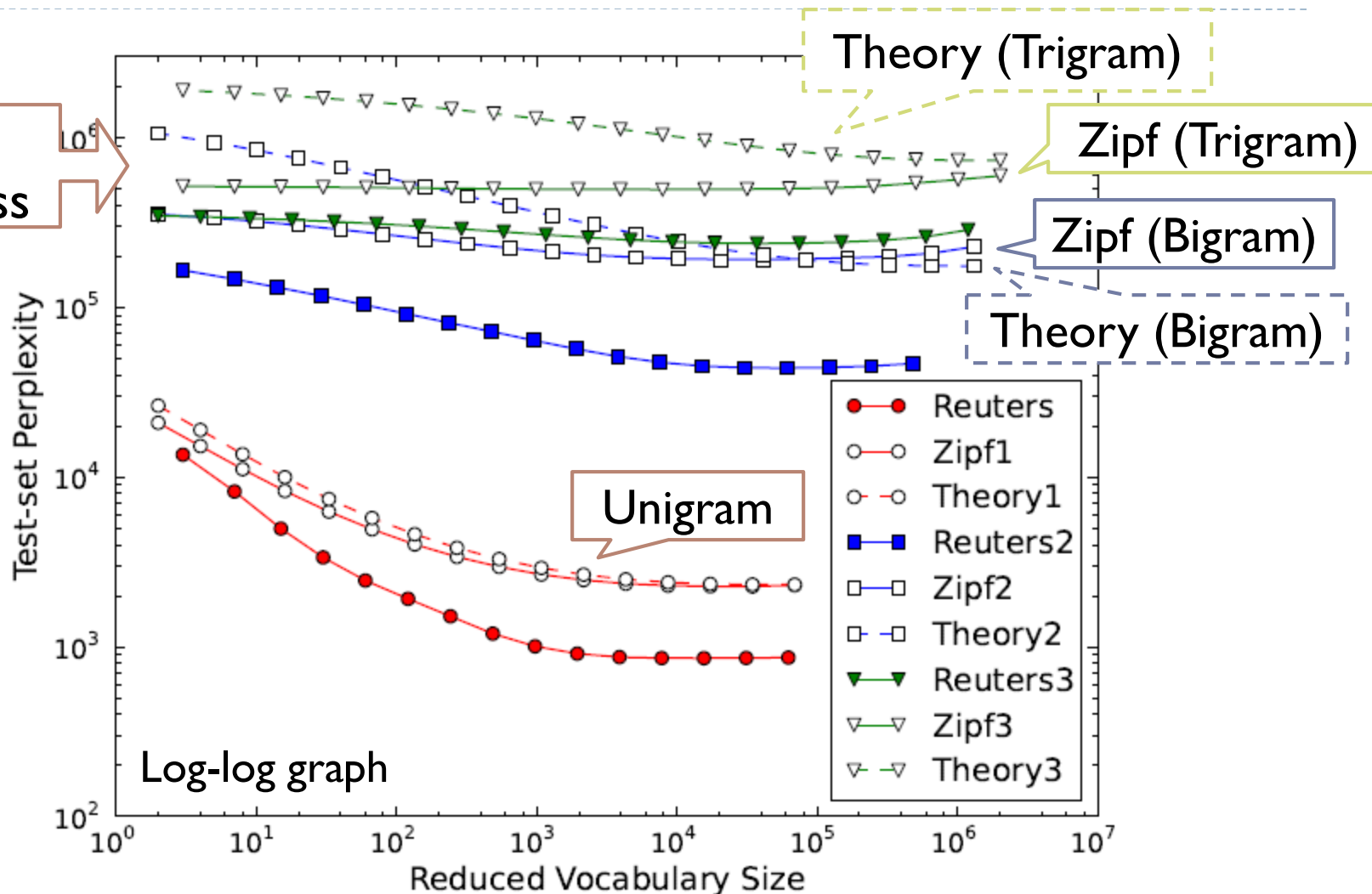
$$H_a(X) := \sum_{x=1}^X \frac{1}{x^a}$$

Hyper harmonic series

$$B_a(X) := \sum_{x=1}^X \frac{a \ln x}{x^a}$$

Bertrand series (another special form)

# PP of k-grams vs. reduced vocab. size



We need to make assumptions that include backoff and smoothing for higher order k-grams

# Additional properties by power-law

---

- ▶ Treat as a variant of the coupon collector's problem
  - ▶ How many trials are needed for collecting all coupons whose occurrence probabilities follow some stable distribution
  - ▶ There exists several works about power law distributions
- ▶ Corpus size for collecting all of the k-grams, according to [Boneh&Papanicolaou 1996]  $\frac{kW^k}{1-\pi_k}$ 
  - ▶ When  $\pi_k = 1$ ,  $W \ln^2 W$ , otherwise,  $\frac{kW^k}{1-\pi_k}$
- ▶ Lower and upper bound of the number of k-grams from the corpus size N and vocab. size W, according to [Atsonios+ 2011]

$$(\pi_k + 1) \left( 1 - e^{-\frac{(1-\pi_k)N}{W^k-1} - \ln \frac{W^k-1}{W^k}} \right) \leq \tilde{W}_k \leq \frac{\pi_k}{\pi_k - 1} \left( \frac{N}{H_{\pi_k}(W^k)} \right)^{\frac{1}{\pi_k}} - \frac{N}{(\pi_k - 1)H_{\pi_k}(W^k)} W^{1-\pi_k}$$

# Outline

---

- ▶ Preliminaries
  - ▶ Zipf's law
  - ▶ Perplexity (PP)
  - ▶ Cutoff and restoring
- ▶ PP of unigram models
- ▶ PP of k-gram models
- ▶ **PP of topic models**
- ▶ Conclusion



# Perplexity of topic models

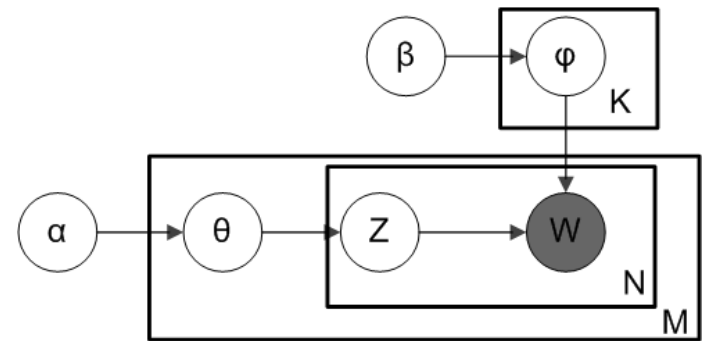
## ▶ Latent Dirichlet Allocation (LDA) [Blei+ 2003]

$$\theta_{d_i} \sim \text{Dirichlet}(\alpha)$$

$$z_i | \theta_{d_i} \sim \text{Multi}(\theta_{d_i})$$

$$\phi_{z_i} \sim \text{Dirichlet}(\beta)$$

$$w_i | z_i, \phi_{z_i} \sim \text{Multi}(\phi_{z_i}),$$



[Griffiths&Steyvers 2004]

## ▶ Learning with Gibbs sampling

▶ Obtain a “good” topic assignment  $z_i$  for each word  $w_i$

## ▶ Posterior distributions of two hidden parameters

$$\hat{\theta}_d(z) \propto n_z^{(d)} + \alpha \quad \text{Document-topic distribution}$$

Mixture rate of topic  $z$  in document  $d$

$$\hat{\phi}_z(w) \propto n_z^{(w)} + \beta \quad \text{Topic-word distribution}$$

Occurrence rate of word  $w$  in topic  $z$

# Rough assumptions of $\phi$ and $\theta$

---

- ▶ Assumption of  $\phi$

- ▶ Word distribution  $\phi_z$  of each topic  $z$  follows Zipf's law

It is natural, regarding each topic as a corpus

- ▶ Assumptions of  $\theta$  (two extreme cases)

- ▶ Case All: Each document evenly has all topics
- ▶ Case One: Each document only has one topic (uniform dist.)

The curve of actual perplexity is expected to be between their values

- ▶ Case All: PP of a topic model  $\approx$  PP of a unigram

- ▶ Marginal predictive distribution is independent of  $d$

$$\sum_{z=1}^T \underbrace{\hat{\theta}_d(z)}_{=1/T} \hat{\phi}_z(w) \propto \sum_{z=1}^T \frac{n_z^{(w)} + \beta}{T} \approx f(w)$$

# Theorem(PP of LDA models: Case One)

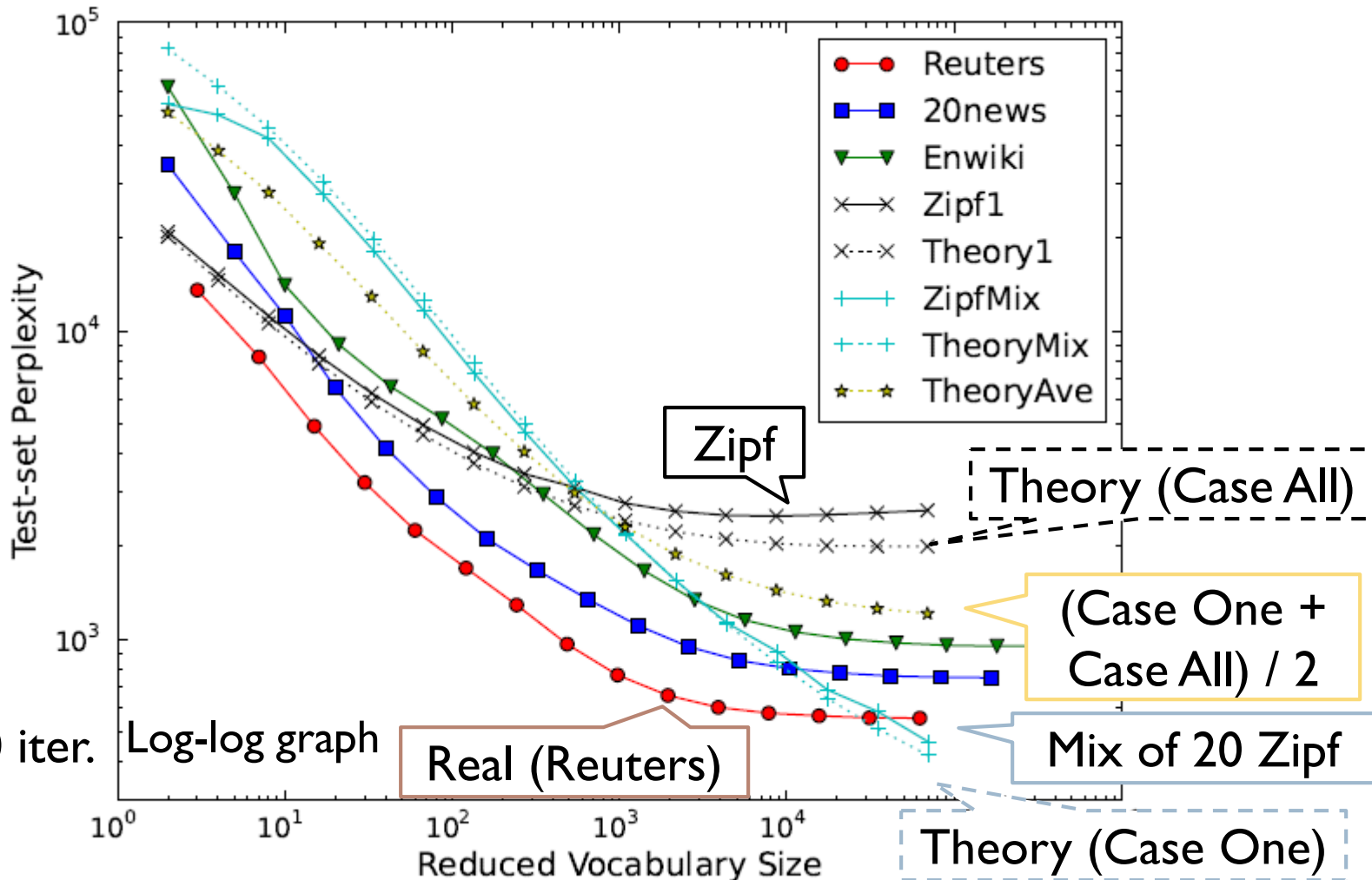
---

- ▶ For any reduced vocabulary size  $W'$ , the perplexity of the optimal restored distribution of a topic model in the Case One is calculated as

$$\hat{P}P_{Mix}(W') = H(W/T) \exp\left(\frac{B(W'/T)}{H(W/T)}\right) \left(\frac{W - W'}{H(W/T) - H(W'/T)}\right)^{1 - \frac{H(W'/T)}{H(W/T)}}$$

$T$  : # of topics in LDA

# PP of LDA models vs. reduced vocab. size



T=20  
CGS w/ 100 iter.  
 $\alpha=\beta=0.1$

Log-log graph

Real (Reuters)

Zipf

Theory (Case All)

(Case One + Case All) / 2

Mix of 20 Zipf


Theory (Case One)

# Time, memory, and PP of LDA learning

---

## ▶ Results of Reuters corpus

corpus	time	memory	perplexity
original	4m3.80s	71,548KB	500
(1/10)	3m55.70s	46,648KB	550
(1/20)	3m42.63s	34,024KB	611



- ▶ Memory usage of the (1/10)-corpus is only 60% of that of the original corpus
  - ▶ Helps in-memory computing for a larger corpus, although the computational time decreased a little

# Outline

---

- ▶ Preliminaries
  - ▶ Zipf's law
  - ▶ Perplexity (PP)
  - ▶ Cutoff and restoring
- ▶ PP of unigram models
- ▶ PP of k-gram models
- ▶ PP of topic models
- ▶ **Conclusion**

# Conclusion

---

- ▶ Trade-off formulae of the cutoff strategy for k-gram models and topic models based on Zipf'law
  - ▶ Perplexity vs. reduced vocabulary size
- ▶ Experiments on real corpora showed that the estimation of the perplexity growth rate is reasonable
- ▶ We can get the best cutoff parameter by maximizing the reduction rate ensuring an acceptable (relative) perplexity
- ▶ Possibility that we can theoretically derive empirical parameters, or “rules of thumb”, for different NLP problems

Can we derive other “rules of thumb” based on Zipf's law?

# Thank you

---