



# Chapter 5:

## Backing-Off Language Models



## 5.1. Unknown words and unseen unigrams



# Goal of this section

- How to define a vocabulary for language tasks
- Notion of “out-of-vocabulary” words

# Example „Crime and Punishment“

- Training corpus: 11 000 tokens
- Test corpus:
  - 3000 tokens
  - 108 of them not seen in training corpus
- Usually training corpus defines vocabulary
- Unseen words in test corpus are hence not part of the vocabulary

OOV-words: “out-of-vocabulary” words  
(tokens)

# Example „Crime and Punishment“

- Training corpus: 11 000 tokens
- Test corpus:
  - 3000 tokens
  - 108 of them not seen in training corpus

Definition:  $\text{OOV - rate} = \frac{\text{\# unseen tokens in test corpus}}{\text{\# of tokens in test corpus}}$



# OOV-Rate on “Crime and Punishment”

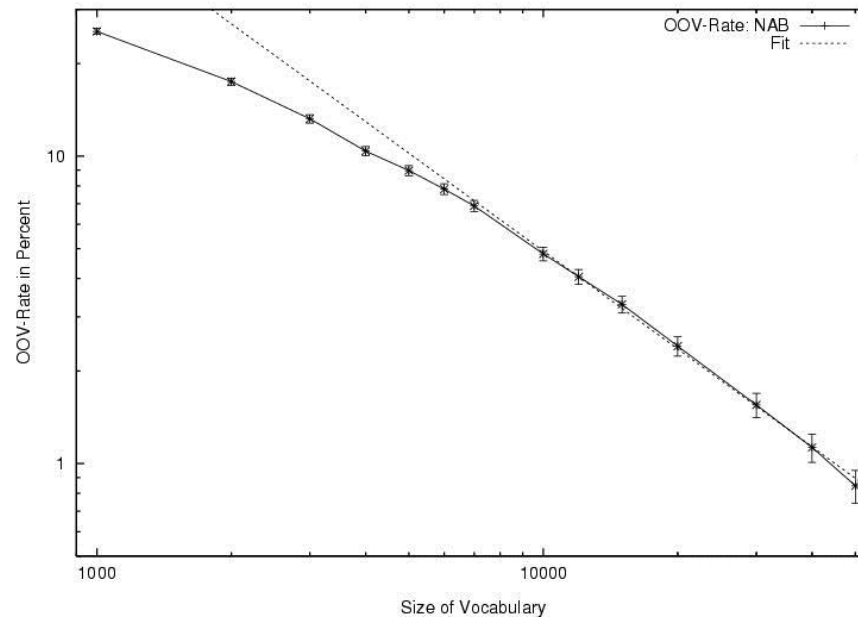
- $\text{OOV-Rate} = 108/3000 = 0.036 = 3.6\%$



# How to define a vocabulary

- Count the frequency of each word in a corpus
- Use the most frequent  $N$  words for the vocabulary

# OOV-Rate vs. Vocabulary



- OOV rate depends on vocabulary size
- Morphologically rich languages have higher OOV rates for the same vocabulary size
- Specifically the relation is a power law

$$\text{OOV - Rate} \propto \frac{1}{(\text{Size of Vocabulary})^\alpha}$$





# Corresponding Problem with sequences of words (M-Gram)

- Size of vocabulary:  $N$
- Number of sequences of length  $M$  with  $N$  words:

$$N^M$$

- Numerical example:  $N=64000$  and  $M=3$ 
  - $N^M = 3 \cdot 10^{14}$
  - No realistic possibility to observe all trigrams in a training corpus

⇒ need methods to account for unseen events



## 5.2 Maximum Likelihood Estimate of Probabilities



# Goal of this section

- Simple way to determine a language model  $P(w|h)$



# From Perplexity to Likelihood

$$PP = P(w_1 \dots w_N)^{-1/N}$$
$$= \exp \left( - \sum_{w,h} f(w,h) \log(P(w|h)) \right)$$

$$-\log(PP) = \sum_{w,h} f(w,h) \log(P(w|h)) = \frac{1}{N} \sum_{w,h} N(w,h) \log(P(w|h))$$

Define likelihood

$$F = \sum_{w,h} N(w,h) \log(P(w|h))$$

Minimize perplexity



Maximize likelihood



# Maximum-Likelihood Estimation of Probabilities

Maximize

$$F = \sum_{w,h} N(w,h) \log(P(w|h))$$

Using the normalization constraint

$$\sum_w P(w|h) = 1 \quad \text{for all } h$$

Modified likelihood using Lagrange multiplier

$$F' = \sum_{w,h} N(w,h) \log(P(w|h)) + \sum_h \lambda(h) \left[ 1 - \sum_w P(w|h) \right]$$

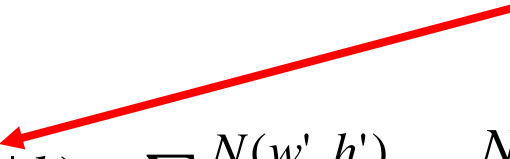


# Maximum-Likelihood Estimation of Probabilities (cont.)

$$F' = \sum_{w,h} N(w,h) \log(P(w|h)) + \sum_h \lambda(h) \left[ 1 - \sum_w P(w|h) \right]$$

$$\frac{\partial F'}{\partial P(w'|h')} = \frac{N(w',h')}{P(w'|h')} - \lambda(h') = 0 \quad \Rightarrow P(w'|h') = \frac{N(w',h')}{\lambda(h')}$$

Constraint

$$1 = \sum_w P(w|h) = \sum_{w'} \frac{N(w',h')}{\lambda(h')} = \frac{N(h')}{\lambda(h')} \quad \Rightarrow \lambda(h') = N(h')$$


$$P(w|h) = \frac{N(w,h)}{N(h)}$$

Maximum likelihood estimate results in relative frequencies



# Issues with Maximum Likelihood Estimate

- Unseen events
  - ↳ vanishing probability estimate
  - ↳ this sequence cannot be found in speech recognition
  - ↳ IR: this document does not cover all query terms, it cannot be relevant

Need other modeling scheme



## 5.3 Absolute discounting

(Kneser-Ney-Smoothing Part 1)





# Goal of this section

- Discuss smoothing techniques that avoid vanishing probabilities

# Backing-off language model

$$P(w | h) = \begin{cases} \frac{N(w, h)}{N(h)} + & \text{für } N(w, h) > 0 \\ 0 & \end{cases}$$

Backing-off weight  $\alpha(h)$

Backing-off distribution  $\beta(w|h)$  (normalized!)

Discounting parameter  $d$

# Backing-off language model (without animation)

$$P(w | h) = \begin{cases} \frac{N(w, h) - d}{N(h)} + \alpha(h) \beta(w | h) & \text{für } N(w, h) > 0 \\ \alpha(h) \beta(w | h) & \end{cases}$$

Backing-off weight  $\alpha(h)$

Backing-off distribution  $\beta(w|h)$  (normalized!)

Discounting parameter  $d$



# Calculating the backing-off weight $\alpha(h)$


Use normalization

$$1 = \sum_w P(w|h) = \sum_{w:N(w,h)>0} \frac{N(w,h)-d}{N(h)} + \sum_w \alpha(h) \beta(w|h)$$


$$= 1 - \frac{\sum_{w:N(w,h)>0} d}{N(h)} + \alpha(h)$$

Solve for  $\alpha(h)$

$$\alpha(h) = \frac{\sum_{w:N(w,h)>0} d}{N(h)} = \frac{d \sum_{w:N(w,h)>0} 1}{N(h)} = \frac{d R(h)}{N(h)} \quad \text{with } R(h) = \sum_{w:N(w,h)>0} 1$$



# Calculating the backing-off weight $\alpha(h)$



Backing-off weight

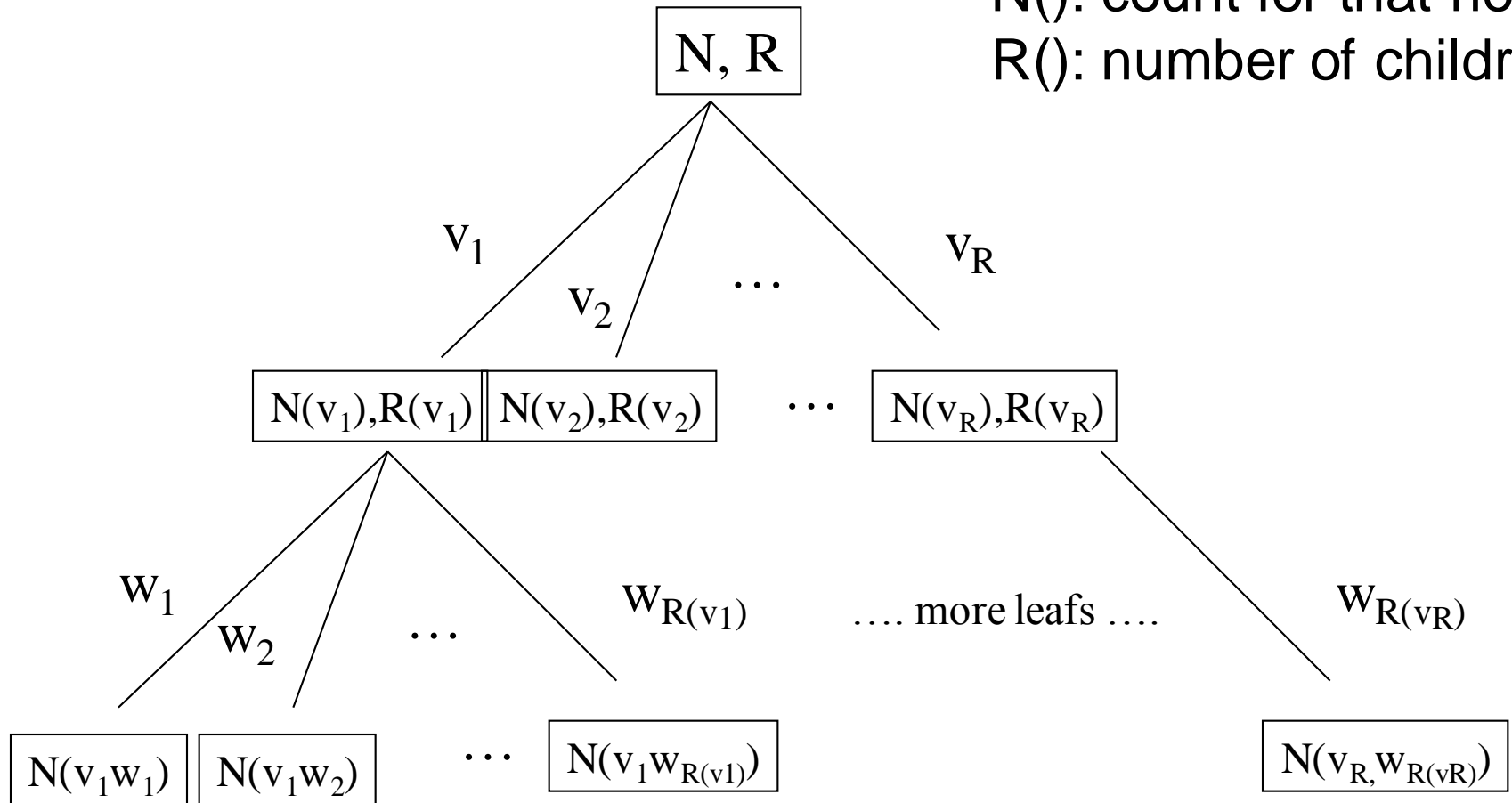
$$\alpha(h) = \frac{d R(h)}{N(h)}$$

$$\text{with } R(h) = \sum_{w: N(w, h) > 0} 1$$



# Storing the counts (Example: Bigram $v, w$ )

$N()$ : count for that node  
 $R()$ : number of children

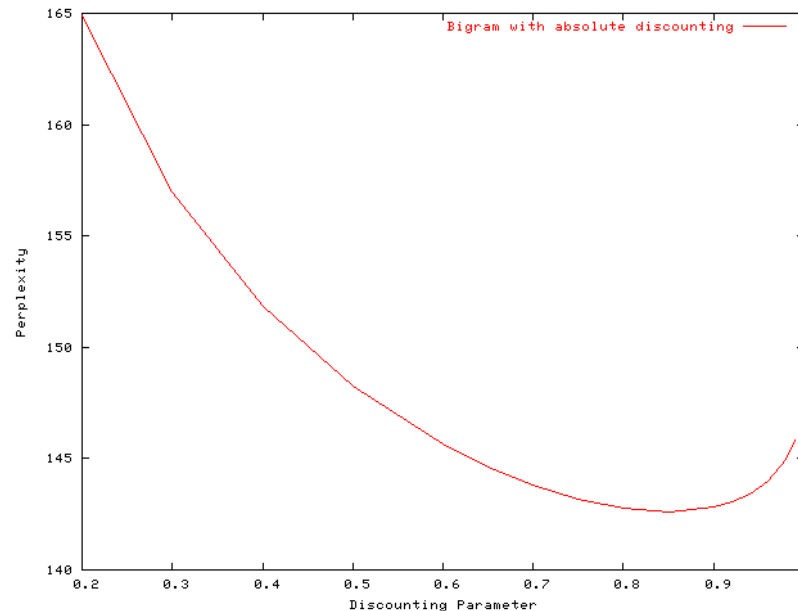




# Examples

- Example of backing-off code
- Example of language model

# Influence of Discounting Parameter



- Discounting parameter has a significant impact on perplexity
- Values of between 0.7 and 0.9 are typically good choices





## 5.4 Leaving-one-out estimate of the discounting parameters



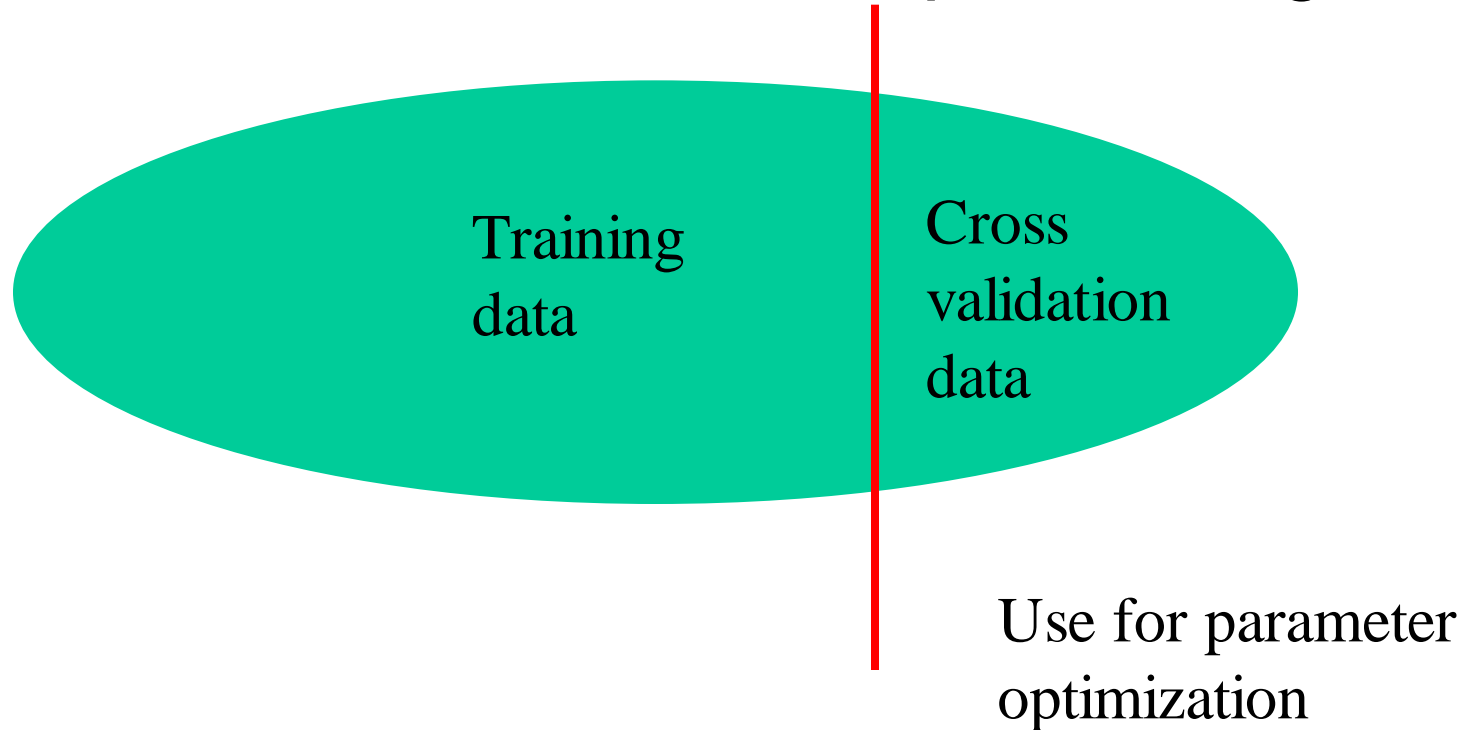
# Goal

- How to get a good estimate for the discounting parameter without tuning it explicitly?



# Closed form optimization solution for discounting parameter

- Idea of cross validation: split training data





# Variants of cross validation

- Suppose you want to train a unigram language model:
  - What is the smallest possible cross validation corpus?
  - How do you avoid fluctuations?



# Leaving-one-out

- Each unigram of the corpus can be the cross validation corpus
- Calculate average of all possible cross validation corpora

# Change of counts for leaving one out

- $W_1 \ W_2 \ W_3 \ W_4 \ W_5 \ ~~W_6~~ \ W_7 \ W_8 \ W_9 \ \dots \ W_N$

$$N(w_6) \quad \mapsto \quad N(w_6)-1$$

$$N \quad \mapsto \quad N-1$$



# Leaving-one-out likelihood

$$\begin{aligned} F_{LOO} &= \sum_{n=1}^N \log P_{LOO}(w_n) = \sum_{w \in W} N(w) \log P_{LOO}(w) \\ &= \sum_{w \in W: N(w) > 1} N(w) \log \left[ \frac{N(w) - 1 - d}{N - 1} + \frac{Rd}{N - 1} \beta \right] \\ &\quad + \sum_{w \in W: N(w) = 1} \log \left[ \frac{(R - 1)d}{N - 1} \beta \right] \end{aligned}$$

Calculate  $\frac{\partial F_{LOO}}{\partial d} = 0$



# Maximizing Leaving-one-out Likelihood

Calculate  $\frac{\partial F_{Loo}}{\partial d} = 0$

- Skip calculating the derivative
- Result

$$0 = \sum_{r=2}^{\infty} n_r r \frac{R\beta - 1}{r - 1 - d + R\beta d} + \frac{n_1}{d}$$

with

$n_1$  : number of words seen once (singeltons)

$n_r$  : number of words that have a count of  $r$  (count - of - counts)

$\mapsto$  no closed form solution to this euqation





# Iterative solution

Split off terms in  $n_1$  and  $n_2$

Solve for  $d$  and use as an iterative equation

$$d_{i+1} = \frac{n_1}{(1 - R\beta) \left( n_1 + 2n_2 + \sum_{r=3}^{\infty} \frac{n_r r (1 - d_i (1 - R\beta))}{r - 1 - d_i (1 - R\beta)} \right)}$$



# Approximate solution

Just use first iteration and initialize with  $d_0=0$

$$d_1 = \frac{n_1}{(1 - R\beta) \left( n_1 + 2n_2 + \sum_{r=3}^{\infty} \frac{n_r r}{r-1} \right)}$$



# Approximate solution

Usually  $R\beta < 1$

Ignore the sum (only a small change to discounting parameter)

$$d \approx \frac{n_1}{n_1 + 2n_2}$$

Works very well for all practical applications



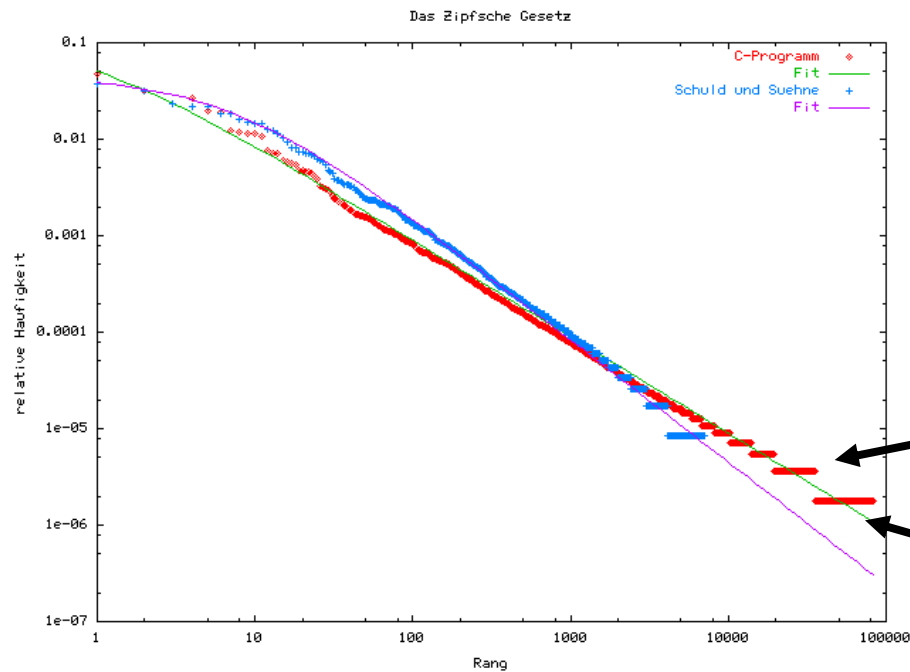
# Demonstrate Approximate Solution

- -> Test formula using switchboard bigram

# Remember Zipf's Law

Zipf (simplified):

$$N(r) = \frac{N_0}{r^\gamma}$$



$n_2$

$n_1$



# Estimate $n_1$ and $n_2$

- Length of the first step in the Zipf-Distribution

$$N(r) = \frac{N_0}{r^\gamma} \quad \Rightarrow \quad r = \left( \frac{N_0}{N} \right)^{\frac{1}{\gamma}}$$

$$\Rightarrow n_1 = \left( \frac{N_0}{1/2} \right)^{\frac{1}{\gamma}} - \left( \frac{N_0}{3/2} \right)^{\frac{1}{\gamma}}$$

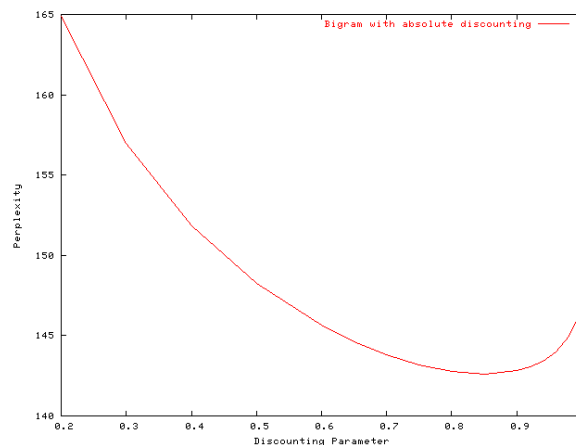
$$\text{and } n_2 = \left( \frac{N_0}{3/2} \right)^{\frac{1}{\gamma}} - \left( \frac{N_0}{5/2} \right)^{\frac{1}{\gamma}}$$



# Estimate Discounting Parameter

$$d \approx \frac{n_1}{n_1 + 2n_2} = \frac{1}{1 + 2 \frac{n_2}{n_1}}$$
$$= \frac{1}{1 + 2 \frac{((2/3)^{1/\gamma} - (2/5)^{1/\gamma})}{(2^{1/\gamma} - (2/3)^{1/\gamma})}}$$

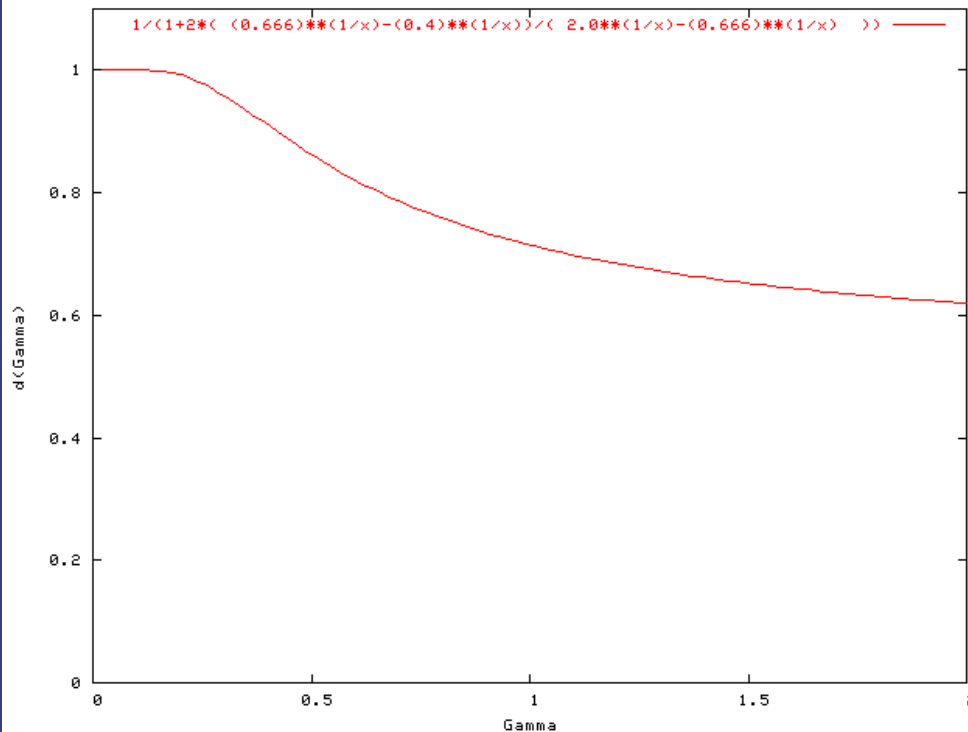
For  $\gamma = 1$  :  $d \approx \frac{5}{7} \approx 0.71$



Estimate only  
gives degradation  
as compared to  
optimal numerical  
solution



# Slope in Zipf vs. Discounting

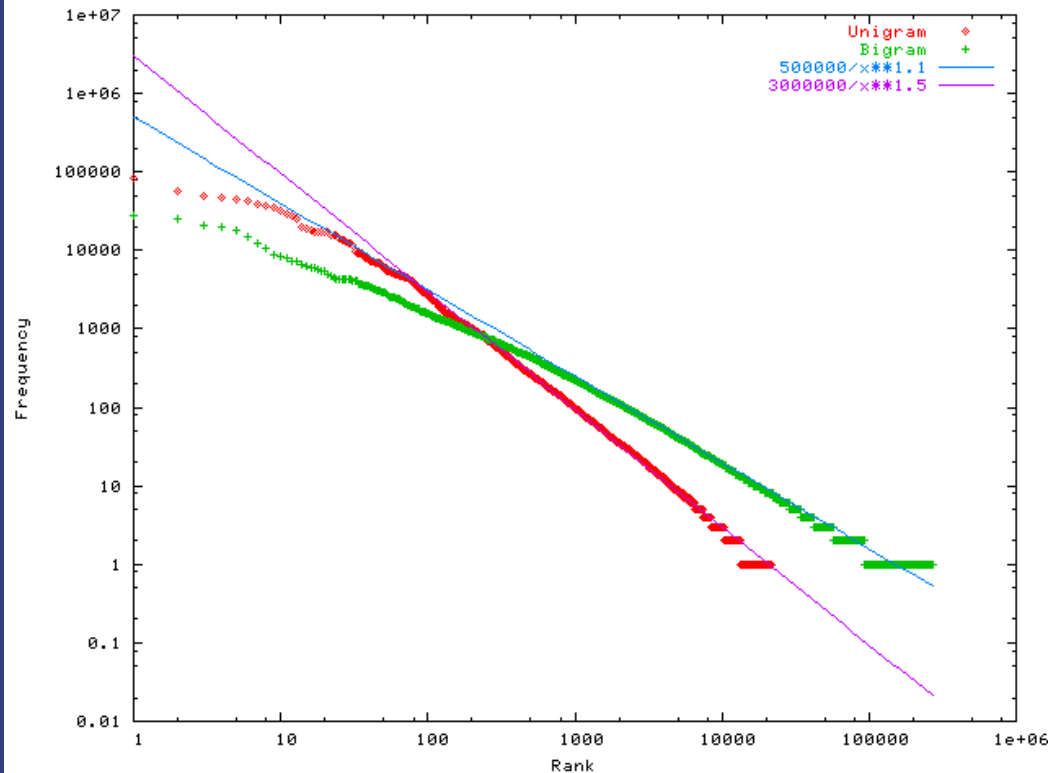


The larger the slope,  
the smaller the  
discounting  
parameter





# Zipf for Unigram and Bigram



Longer range  
models have larger  
discounting  
parameters



## 5.5 Other smoothing methods



# Goal

- Smoothing techniques that are simpler
- Smoothing techniques that are used for other reasons



# Floor discounting (“add epsilon”) aka additive smoothing; Lidstone smoothing, Laplace Smoothing, ...

Add a little bit to each count

$$P(w|h) = \frac{1}{Z} (N(w, h) + \varepsilon)$$

1/Z: Normalisation

$$1 = \sum_{w=1}^V P(w|h) = \frac{1}{Z} (N(h) + \varepsilon V) \quad \Rightarrow \quad Z = (N(h) + \varepsilon V)$$

$$P(w|h) = \frac{N(w, h) + \varepsilon}{N(h) + \varepsilon V}$$



# Linear Discounting

$$P(w|h) = (1 - \varepsilon) \frac{N(w, h)}{N(h)} + \varepsilon \beta(w|h)$$

$\varepsilon$ : determines amount of smoothing  
 $\beta(w|h)$ : "backing - off distribution"

Variant : 
$$P(w|h) = (1 - \varepsilon(h)) \frac{N(w, h)}{N(h)} + \varepsilon(h) \beta(w|h)$$

Popular value :  $\varepsilon(h) = \frac{n_1(h)}{N(h)}$  with  $n_1(h) = \sum_{w: N(w, h)=1} 1$

See example

# Good-Turing Discounting

Adjusted counts:

$$N^*(w, h) = \begin{cases} [N(w, h) + 1] \frac{n_{N(w, h) + 1}}{n_{N(w, h)}} & \text{für } N(w, h) \leq 5 \\ N(w, h) & \text{sonst} \end{cases}$$

Model:

$$P(w | h) = \frac{N^*(w, h)}{N(h)} + \lambda(h) \beta(w | h)$$

Backing-off weight from normalization:  $\lambda(h) = 1 - \frac{N^*(h)}{N(h)}$



## 5.6 Kneser-Ney Smoothing



# Kneser-Ney Smoothing

- Try to find a dedicated backing-off distribution
- Two different variants
  - “Marginal constraint” backing off
  - Singleton backing-off



# Idea

$$P(w|h) = \begin{cases} \frac{N(w,h)-d}{N(h)} + \alpha(h) \beta(w|\hat{h}) & \text{für } N(w,h) > 0 \\ \alpha(h) \beta(w|\hat{h}) & \end{cases}$$

With  $\hat{h}$  beeing the history shortened by one word

$$p(w|\hat{h}) = \sum_{h'} P(w|h)P(h'|\hat{h}) \quad (**)$$

with h' beeing the word removed from the history

⇒ look for backing off distribution such that (\*\*) is satisfied

# Solution

- Direct calculation
- Solution

$$\beta(w | \hat{h}) = \frac{N_+(w, \hat{h})}{\sum_w N_+(w, \hat{h})}$$

with

$$N_+(w, \hat{h}) = \sum_{h': N(w, \hat{h}, h') > 0} 1$$

Marginal  
constraint  
backing-off

# Alternative Determine Backing-Off Distribution using Leaving-One-Out

$$F = F_0 + \sum_{h' \hat{h} w: N(w \hat{h} h')=1} \log(\alpha(h) \beta(w | \hat{h})) \\ + \sum_{\hat{h}} \lambda(\hat{h}) \left( 1 - \sum_w \beta(w | \hat{h}) \right)$$

$F_0$       likelihood for events seen more than once

$\lambda(\hat{h})$     Lagrange multiplier

# Solution for Leaving-One-Out Approach to Backing-Off Distribution

- Direct calculation of first derivative
- Solution

$$\beta(w | \hat{h}) = \frac{N_1(w, \hat{h})}{\sum_w N_1(w, \hat{h})}$$

with

$$N_1(w, \hat{h}) = \sum_{h': N(w, \hat{h}, h')=1} 1$$



# Comparison Marginal-Constraint/Singleton Backing-Off



Backing-off Model	Bigram	Trigram	4-gram
Rel. Frequencies	160.1	146.4	147.8
Marginal Constraint	154.8	138.3	137.9
Singleton BO	156.2	142.6	141.9



## 5.7 Pruning of count trees

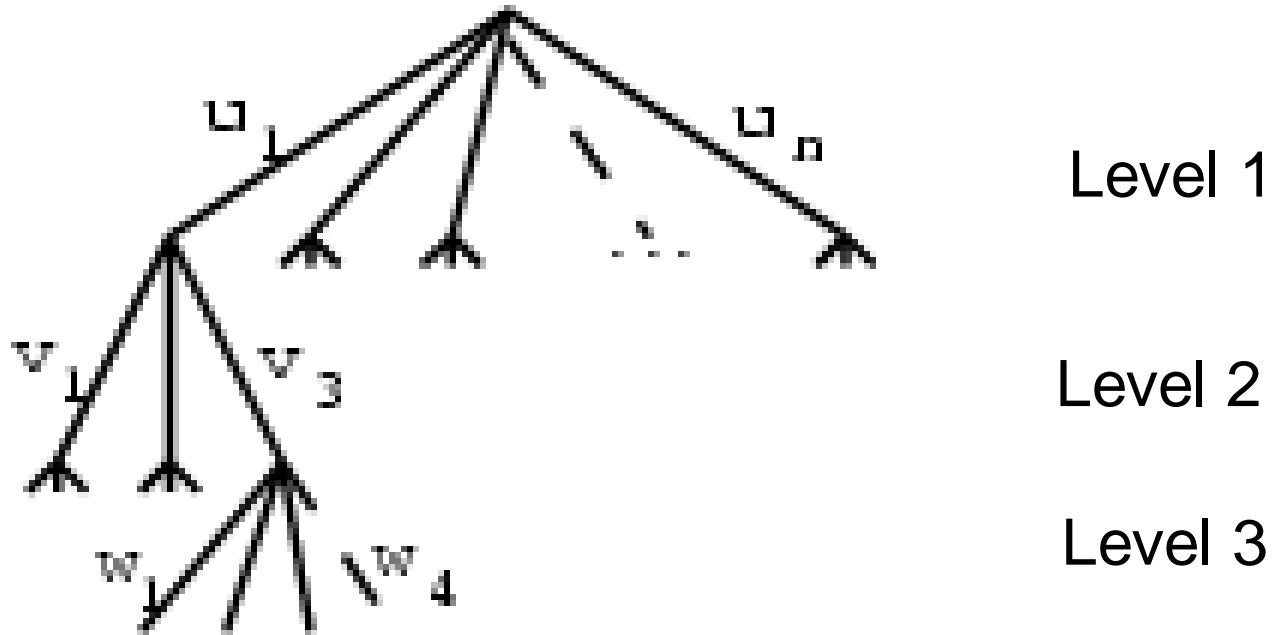


# Goal

- Method to decrease size of the tree



# Pruning Count-Trees



Simple criterion for pruning trees:  
remove all infrequent nodes





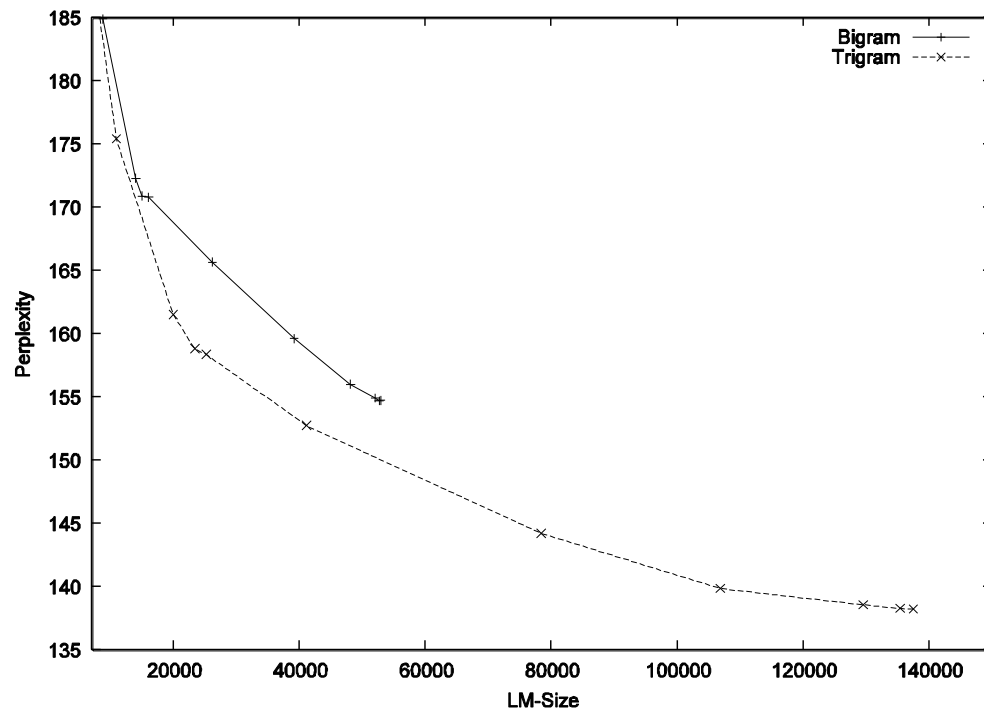
# Improved Criterion: Change in Likelihood

- Metric inspired by maximum likelihood objective
- In case a leaf is pruned, the backing off probability will be used instead

$$\begin{aligned}\Delta F_{Leaf} &= N(w, h) \log P(w | h) - N(wh) \log(\alpha(h) \beta(w | \hat{h})) \\ &= N(w, h) \log \left( \frac{P(w | h)}{\alpha(h) \beta(w | \hat{h})} \right)\end{aligned}$$

- This is a metric for the leaves of the tree
- Create measure for sub-trees from measures for leaves

# Effect of language model pruning



- Pruned Trigram better than pruned bigram
- For a given size pruning gives you the best LM
- Results generalize to longer N-grams



# Summary

- OOV-words
- Maximum likelihood estimator
- Backing-off language models
- Discounting parameter
- Kneser-Ney Smoothing
- Pruning of trees