Language Models



IT Systems Engineering | Universität Potsdam

Dr. Mariana Neves (adapted from the original slides of Prof. Philipp Koehn)

November 28th, 2016

• Language models answer the question:

How likely is a string of English words good English?

- 3

Image: A math a math

• Help with reordering

 $p_{\text{LM}}(ext{the house is small}) > p_{ ext{LM}}(ext{small the is house})$

• Help with word choice

 $p_{\text{LM}}(\text{I am going home}) > p_{\text{LM}}(\text{I am going house})$

- Given: a string of English words $W = w_1, w_2, w_3, ..., w_n$
- Question: what is p(W)?

• We collect large amount of text and count how often *W* occurs to estimate p(W)

- 31

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

- Sparse data: Many good English sentences will not have been seen before
- Decomposing p(W) using the chain rule:

 $p(w_1, w_2, w_3, ..., w_n) = p(w_1) p(w_2|w_1) p(w_3|w_1, w_2) ... p(w_n|w_1, ..., w_{n-1})$

(not much gained yet, $p(w_n|w_1, w_2, ..., w_{n-1})$ is equally sparse)

▲□▶ ▲□▶ ▲□▶ ▲□▶ = ののの

• Markov assumption:

- only previous history matters
- limited memory: only last k words are included in history (older words less relevant)
- \rightarrow *k*th order Markov model

・ 何 ト ・ ヨ ト ・ ヨ ト

• For instance 2-gram language model:

 $p(w_1, w_2, w_3, ..., w_n) \simeq p(w_1) p(w_2|w_1) p(w_3|w_2)...p(w_n|w_{n-1})$

• What is conditioned on, here w_{i-1} is called the **history**

- 31

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

• More training data allows for longer histories (higher **kth**).

- Most commonly, trigram (3-grams) models are used.
- But bigrams (2-grams), unigrams (single words) or any other order of n-grams is possible.

Maximum likelihood estimation

$$p(w_2|w_1) = rac{\mathsf{count}(w_1, w_2)}{\mathsf{count}(w_1)}$$

- Collect counts over a large text corpus
- Millions to billions of words are easy to get (trillions of English words available on the web)

Example: 3-Gram

• Counts for trigrams and estimated word probabilities

the gre	<mark>en</mark> (tot	al: 1748)	the ree	d (tota	l: 225)	the blu	le (tot	tal: 54)
word	с.	prob.	word	c.	prob.	word	c.	prob.
paper	801	0.458	cross	123	0.547	box	16	0.296
group	640	0.367	tape	31	0.138		6	0.111
light	110	0.063	army	9	0.040	flag	6	0.111
party	27	0.015	card	7	0.031	,	3	0.056
ecu	21	0.012	,	5	0.022	angel	3	0.056

• 225 trigrams in the Europarl corpus start with the red

- 123 of them end with cross
- \rightarrow maximum likelihood probability is $\frac{123}{225} = 0.547$.

- 4 同 6 4 日 6 4 日 6

How good is the LM?

- A good model assigns a text of real English W a high probability
- This can be also measured with cross entropy:

$$H(W) = -\frac{1}{n} \log p(W_1^n) \\ -\frac{1}{n} \sum_{i=1}^n \log p(w_i | w_1, w_2, ..., w_{i-1})$$

• Or, perplexity

 $perplexity(W) = 2^{H(W)}$

riana		

Example: trigrams

prediction	$oldsymbol{ ho}_{ ext{LM}}$	$-\log_2 p_{\text{LM}}$
$p_{\text{LM}}(i <\mathrm{s}>)$	0.109	3.197
$p_{\text{LM}}(\text{would} <\text{s}>\text{i})$	0.144	2.791
$p_{\text{\tiny LM}}(\text{like} \text{i would})$	0.489	1.031
p_{LM} (to would like)	0.905	0.144
${m ho}_{\scriptscriptstyle m LM}({ m commend} { m like to})$	0.002	8.794
$p_{\scriptscriptstyle \mathrm{LM}}(\mathrm{the} \mathrm{to}\;\mathrm{commend})$	0.472	1.084
p_{LM} (rapporteur commend the)	0.147	2.763
$\rho_{\text{\tiny LM}}(ext{on} ext{the rapporteur})$	0.056	4.150
p_{LM} (his rapporteur on)	0.194	2.367
$p_{\text{LM}}(\text{work} \text{on his})$	0.089	3.498
$p_{\text{\tiny LM}}(. ext{his work})$	0.290	1.785
$p_{\scriptscriptstyle \mathrm{LM}}(<\!\!/\mathrm{s}\!\!>\!\! \mathrm{work} .)$	0.99999	0.000014
	average	2.634

I would like to commend the rapporteur on his work.

-

3

< □ > < ---->

Comparison 1-4-Gram

word	unigram	bigram	trigram	4-gram
i	6.684	3.197	3.197	3.197
would	8.342	2.884	2.791	2.791
like	9.129	2.026	1.031	1.290
to	5.081	0.402	0.144	0.113
commend	15.487	12.335	8.794	8.633
the	3.885	1.402	1.084	0.880
rapporteur	10.840	7.319	2.763	2.350
on	6.765	4.140	4.150	1.862
his	10.678	7.316	2.367	1.978
work	9.993	4.816	3.498	2.394
	4.896	3.020	1.785	1.510
	4.828	0.005	0.000	0.000
average	8.051	4.072	2.634	2.251
perplexity	265.136	16.817	6.206	4.758

- 2

<ロ> (日) (日) (日) (日) (日)

- We have seen $i \ like \ to \ in \ our \ corpus$
- \bullet We have never seen $i\ like\ to\ smooth\ in\ our\ corpus$
- $\rightarrow p(\text{smooth}|\text{i like to}) = 0$

• Any sentence that includes i like to smooth will be assigned probability 0

• For all possible n-grams, add the count of one.

$$p = \frac{c+1}{n+v}$$

- c = count of n-gram in corpus
- *n* = count of history
- v = vocabulary size (total number of possible n-grams)

- But there are many more unseen n-grams than seen n-grams
- Example: Europarl 2-bigrams:
 - 86,700 distinct words
 - $86,700^2 = 7,516,890,000$ possible bigrams
 - but only about 30,000,000 bigrams in corpus

• Add $\alpha < 1$ to each count

$$p = \frac{c + \alpha}{n + \alpha v}$$

- What is a good value for α ?
- Could be optimized on held-out set

Example: 2-Grams in Europarl

Count	Adjusted count		Test count
С	$(c+1)\frac{n}{n+v^2}$	$(c+\alpha)\frac{n}{n+\alpha v^2}$	t _c
0	0.00378	0.00016	0.00016
1	0.00755	0.95725	0.46235
2	0.01133	1.91433	1.39946
3	0.01511	2.87141	2.34307
4	0.01888	3.82850	3.35202
5	0.02266	4.78558	4.35234
6	0.02644	5.74266	5.33762
8	0.03399	7.65683	7.15074
10	0.04155	9.57100	9.11927
20	0.07931	19.14183	18.95948

- Add- α smoothing with $\alpha = 0.00017$
- t_c are average counts of n-grams in test set that occurred c times in corpus

Mariana Neves

- Estimate true counts in held-out data
 - split corpus in two halves: training and held-out
 - counts in training $C_t(w_1, ..., w_n)$
 - number of n-grams with training count r: N_r
 - total times n-grams of training count r seen in held-out data: T_r

Example: Deleted estimation (bigrams)

Count	Count of counts	Counts in held-out	Exp. count
r	N _r	T _r	$E(r) = T_r/N_r$
0	7,515,623,434	938,504	0.00012
1	753,777	353,383	0.46900
2	170,913	239,736	1.40322
3	78,614	189,686	2.41381
4	46,769	157,485	3.36860
5	31,413	134,653	4.28820
6	22,520	122,079	5.42301
8	13,586	99,668	7.33892
10	9,106	85,666	9.41129
20	2,797	53,262	19.04992

・ロン ・四 ・ ・ ヨン ・ ヨン

- 2

- We can adjust the real counts to these expected counts
 - better estimates for both seen and unseen events
- Both halves can be switched and results combined

$$r_{del} = rac{T_r^1 + T_r^2}{N_r^1 + N_r^2} \ ext{where} \ r = count(w_1, ..., w_n)$$

Mariana Neves

• Adjust actual counts r to expected counts r^* with formula

$$r^* = (r+1)\frac{N_{r+1}}{N_r}$$

• N_r number of n-grams that occur exactly r times in corpus

	Neves	

Good-Turing for 2-Grams in Europarl

Count	Count of counts	Adjusted count	Test count
r	N _r	r*	t
0	7,514,941,065	0.00015	0.00016
1	1,132,844	0.46539	0.46235
2	263,611	1.40679	1.39946
3	123,615	2.38767	2.34307
4	73,788	3.33753	3.35202
5	49,254	4.36967	4.35234
6	35,869	5.32928	5.33762
8	21,693	7.43798	7.15074
10	14,880	9.31304	9.11927
20	4,546	19.54487	18.95948

adjusted count fairly accurate when compared against the test count

Mariana Neves

- In given corpus, we may never observe
 - Scottish beer drinkers
 - Scottish beer eaters
- Both have count 0

 \rightarrow our smoothing methods will assign them the same probability

3

18 A.

Back-Off

- Better: backoff to bigrams:
 - beer drinkers
 - $\bullet\,$ beer eaters

3

<ロ> (日) (日) (日) (日) (日)

- Higher and lower order n-gram models have different strengths and weaknesses
 - high-order n-grams are sensitive to more context, but have sparse counts
 - low-order n-grams consider only very limited context, but have robust counts

4 AR & 4 E & 4 E &

• Combine them

$$p_{I}(w_{3}|w_{1}, w_{2}) = \lambda_{1} p_{1}(w_{3}) \\ + \lambda_{2} p_{2}(w_{3}|w_{2}) \\ + \lambda_{3} p_{3}(w_{3}|w_{1}, w_{2})$$

 $\forall \lambda_n : 0 \le \lambda_n \le 1 \\ \sum_n \lambda_n = 1$

- We can trust some histories $w_{i-n+1}, ..., w_{i-1}$ more than others
- Condition interpolation weights on history: $\lambda_{w_{i-n+1},...,w_{i-1}}$
- Recursive definition of interpolation

$$p'_n(w_i|w_{i-n+1},...,w_{i-1}) = \lambda_{w_{i-n+1},...,w_{i-1}} p_n(w_i|w_{i-n+1},...,w_{i-1}) + + (1 - \lambda_{w_{i-n+1},...,w_{i-1}}) p'_{n-1}(w_i|w_{i-n+2},...,w_{i-1})$$

Back-Off

Trust the highest order language model that contains n-gram

$$p_n^{BO}(w_i|w_{i-n+1},...,w_{i-1}) = \begin{cases} d_n(w_{i-n+1},...,w_{i-1}) \ p_n(w_i|w_{i-n+1},...,w_{i-1}) \\ \text{if count}_n(w_{i-n+1},...,w_i) > 0 \\ \\ \alpha_n(w_{i-n+1},...,w_{i-1}) \ p_{n-1}^{BO}(w_i|w_{i-n+2},...,w_{i-1}) \\ \\ \text{otherwise} \end{cases}$$

- Requires
 - adjusted prediction model $\alpha_n(w_i|w_{i-n+1},...,w_{i-1})$
 - discounting function $d_n(w_1, ..., w_{n-1})$

3

Diversity of Predicted Words

• Consider the bigram histories spite and constant

- both occur 993 times in Europarl corpus
- only 9 different words follow spite almost always followed by of (979 times), due to expression in spite of
- 415 different words follow constant most frequent: and (42 times), concern (27 times), pressure (26 times), but huge tail of singletons: 268 different words
- More likely to see new bigram that starts with constant than spite
- Witten-Bell smoothing considers diversity of predicted words

Witten-Bell Smoothing

- Recursive interpolation method
- Number of possible extensions of a history $w_1, ..., w_{n-1}$ in training data

$$N_{1+}(w_1,...,w_{n-1},\bullet) = |\{w_n : c(w_1,...,w_{n-1},w_n) > 0\}|$$

Lambda parameters

$$1 - \lambda_{w_1,...,w_{n-1}} = \frac{N_{1+}(w_1,...,w_{n-1},\bullet)}{N_{1+}(w_1,...,w_{n-1},\bullet) + \sum_{w_n} c(w_1,...,w_{n-1},w_n)}$$

Witten-Bell Smoothing: Examples

Let us apply this to our two examples:

$$1 - \lambda_{spite} = \frac{N_{1+}(\text{spite}, \bullet)}{N_{1+}(\text{spite}, \bullet) + \sum_{w_n} c(\text{spite}, w_n)}$$
$$= \frac{9}{9 + 993} = 0.00898$$

$$1 - \lambda_{constant} = \frac{N_{1+}(constant, \bullet)}{N_{1+}(constant, \bullet) + \sum_{w_n} c(constant, w_n)}$$
$$= \frac{415}{415 + 993} = 0.29474$$

-

3

- Consider the word York
 - fairly frequent word in Europarl corpus, occurs 477 times
 - $\bullet\,$ as frequent as foods, indicates and providers
 - ightarrow in unigram language model: a respectable probability
- However, it almost always directly follows New (473 times)
- Recall: unigram model only used, if the bigram model inconclusive
 - York unlikely second word in unseen bigram
 - in back-off unigram model, York should have low probability

- Kneser-Ney smoothing takes diversity of histories into account
- Count of histories for a word

 $N_{1+}(\bullet w) = |\{w_i : c(w_i, w) > 0\}|$

• Recall: maximum likelihood estimation of unigram language model

$$p_{ML}(w) = \frac{c(w)}{\sum_i c(w_i)}$$

• In Kneser-Ney smoothing, replace raw counts with count of histories

$$p_{KN}(w) = \frac{N_{1+}(\bullet w)}{\sum_{w_i} N_{1+}(\bullet w_i)}$$

Evaluation of smoothing methods:

Perplexity for language models trained on the Europarl corpus

Smoothing method	bigram	trigram	4-gram
Good-Turing	96.2	62.9	59.9
Witten-Bell	97.1	63.8	60.4
Modified Kneser-Ney	95.4	61.6	58.6

3

A B A A B A

- ∢ 🗇 እ

 Millions to billions of words are easy to get (trillions of English words available on the web)

• But: huge language models do not fit into RAM

Number of Unique N-Grams

Number of unique n-grams in Europarl corpus 29,501,088 tokens (words and punctuation)

Order	Unique n-grams	Singletons
unigram	86,700	33,447 (38.6%)
bigram	1,948,935	1,132,844 (58.1%)
trigram	8,092,798	6,022,286 (74.4%)
4-gram	15,303,847	13,081,621 (85.5%)
5-gram	19,882,175	18,324,577 (92.2%)

 \rightarrow remove singletons of higher order n-grams

- Language models too large to build
- What needs to be stored in RAM?
 - maximum likelihood estimation

$$p(w_n|w_1,...,w_{n-1}) = \frac{\operatorname{count}(w_1,...,w_n)}{\operatorname{count}(w_1,...,w_{n-1})}$$

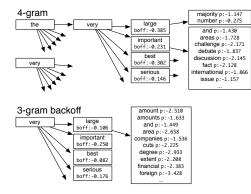
• can be done separately for each history $w_1, ..., w_{n-1}$

- ∢ ศ⊒ ▶

- Keep data on disk
 - extract all n-grams into files on-disk
 - sort by history on disk
 - only keep n-grams with shared history in RAM
- Smoothing techniques may require additional statistics

- Need to store probabilities for
 - the very large majority
 - the very large number
- Both share history the very large
- $\rightarrow\,$ no need to store history twice
- \rightarrow Trie

Efficient Data Structures



2-gram backoff

large	accept p:-3.791
boff:-0.470	acceptable p:-3.778
	accession p:-3.762
	accidents p:-3.806
	accountancy p:-3.416
	accumulated p:-3.885
	accumulation p:-3.895
	action p:-3.510
	additional p: -3.334
	administration p:-3.729

1-gram backoff

aa-afns p:-6.154
aachen p:-5.734
aaiun p:-6.154
aalborg p:-6.154
aarhus p:-5.734
aaron p:-6.154
aartsen p:-6.154
ab p:-5.734
abacha p:-5.156
aback p:-5.876

Mariana Neves				

3

・ロト ・聞ト ・ヨト ・ヨト

Reducing Vocabulary Size

- For instance: each number is treated as a separate token
- Replace them with a number token NUM
 - but: we want our language model to prefer

 $\rho_{\rm LM}(\text{I pay 950.00 in May 2007}) > \rho_{\rm LM}(\text{I pay 2007 in May 950.00})$

• not possible with number token

 $p_{\text{LM}}(\text{I pay NUM in May NUM}) = p_{\text{LM}}(\text{I pay NUM in May NUM})$

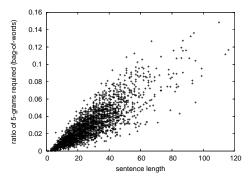
• Replace each digit (with unique symbol, e.g., @ or 5), retain some distinctions

 $ho_{\rm LM}({
m I}~{
m pay}~555.55~{
m in}~{
m May}~5555) >
ho_{\rm LM}({
m I}~{
m pay}~5555~{
m in}~{
m May}~555.55)$

▲□▶ ▲□▶ ▲□▶ ▲□▶ = ののの

Filtering Irrelevant N-Grams

- We use language model in decoding
 - we only produce English words in translation options
 - filter language model down to n-grams containing only those words
- Ratio of 5-grams needed to all 5-grams (by sentence length):



Summary

- Language models: *How likely is a string of English words good English?*
- N-gram models (Markov assumption)
- Perplexity
- Count smoothing
 - $\bullet\,$ add-one, add- α
 - deleted estimation
 - Good Turing
- Interpolation and backoff
 - Good Turing
 - Witten-Bell
 - Kneser-Ney
- Managing the size of the model

• Statistical Machine Translation, Philipp Koehn (chapter 7).

3

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >