

N-GRAM LANGUAGE MODELS



Evaluating Language Models

Unigram vs. Bigram vs. Trigram vs. n-gram

Evaluating Language Models

The more n in n-gram, the better language model?
The more history, the better prediction of future?

Evaluating Language Models

Qualitative → Generate

```
Unigram: stream = []
while (w != '</s>'):
    w = unigrams_freq.select()
    stream.append(w)
```

Evaluating Language Models

Qualitative → Generate

```
Bigram: stream = []
while (w != '</s>'):
    w = bigrams_freq[stream[-1]].select()
    stream.append(w)
```

Evaluating Language Models

Qualitative → Generate

```
Trigram: stream = []
          while (w != '</s>'):
              w = trigrams_freq[stream[-2:]].select()
              stream.append(w)
```

1 gram	<p>–To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have</p> <p>–Hill he late speaks; or! a more to leg less first you enter</p>
2 gram	<p>–Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow.</p> <p>–What means, sir. I confess she? then all sorts, he is trim, captain.</p>
3 gram	<p>–Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done.</p> <p>–This shall forbid it should be branded, if renown made it empty.</p>
4 gram	<p>–King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in;</p> <p>–It cannot be but so.</p>

Figure 3.3 Eight sentences randomly generated from four n-grams computed from Shakespeare's works. All characters were mapped to lower-case and punctuation marks were treated as words. Output is hand-corrected for capitalization to improve readability.

1
gram

Months the my and issue of year foreign new exchange's september
were recession exchange new endorsed a acquire to six executives

2
gram

Last December through the way to preserve the Hudson corporation N.
B. E. C. Taylor would seem to complete the major central planners one
point five percent of U. S. E. has already old M. X. corporation of living
on information such as more frequently fishing to keep her

3
gram

They also point to ninety nine point six billion dollars from two hundred
four oh six three percent of the rates of interest stores as Mexico and
Brazil on market conditions

Figure 3.4 Three sentences randomly generated from three n-gram models computed from 40 million words of the *Wall Street Journal*, lower-casing all characters and treating punctuation as words. Output was then hand-corrected for capitalization to improve readability.

1
gram

–To him swallowed confess hear both. Which. Of save on trail for are ay device and
rote life have

–Hill he late speaks; or! a more to leg less first you enter

2
gram

–Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live
king. Follow.

–What means, sir. I confess she? then all sorts, he is trim, captain.

3
gram

–Fly, and will rid me these news of price. Therefore the sadness of parting, as they say,
'tis done.

–This shall forbid it should be branded, if renown made it empty.

4
gram

–King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A
great banquet serv'd in;

–It cannot be but so.

Figure 3.3 Eight sentences randomly generated from four n-grams computed from Shakespeare's works. All characters were mapped to lower-case and punctuation marks were treated as words. Output is hand-corrected for capitalization to improve readability.

1
gram

1
gram

Months the my and issue of year foreign new exchange's september
were recession exchange new endorsed a acquire to six executives

2
gram

2
gram

Last December through the way to preserve the Hudson corporation N.
B. E. C. Taylor would seem to complete the major central planners one
point five percent of U. S. E. has already old M. X. corporation of living
on information such as more frequently fishing to keep her

3
gram

3
gram

They also point to ninety nine point six billion dollars from two hundred
four oh six three percent of the rates of interest stores as Mexico and
Brazil on market conditions

Figure 3.4 Three sentences randomly generated from three n-gram models computed from
40 million words of the *Wall Street Journal*, lower-casing all characters and treating punctuation
as words. Output was then hand-corrected for capitalization to improve readability.

Cross Evaluating Language Models

Qualitative → Generate

Biased toward the corpus! Dialect, Genre, ...

Evaluating Language Models

Quantitative → Likelihood

Models

M^*

$$\begin{aligned} P(L | M^*) &= 1 \\ P(l | M^*) &=? \end{aligned}$$

M_2

$$\begin{aligned} P(L | M_2) &= 0.001 \\ P(l | M_2) &= 1 \end{aligned}$$

M_1

$$\begin{aligned} P(L | M_1) &= 0 \\ P(l | M_1) &= 0 \end{aligned}$$

Language
 L

l

Find M^*

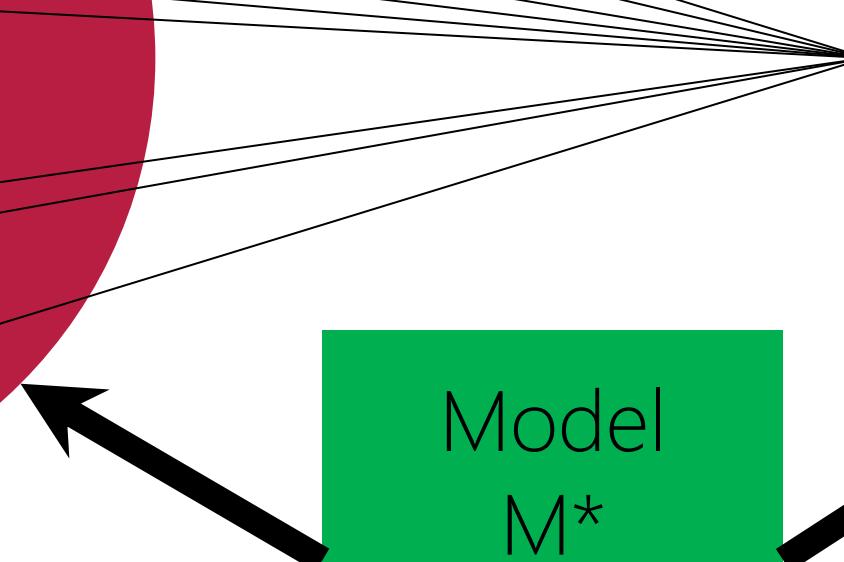
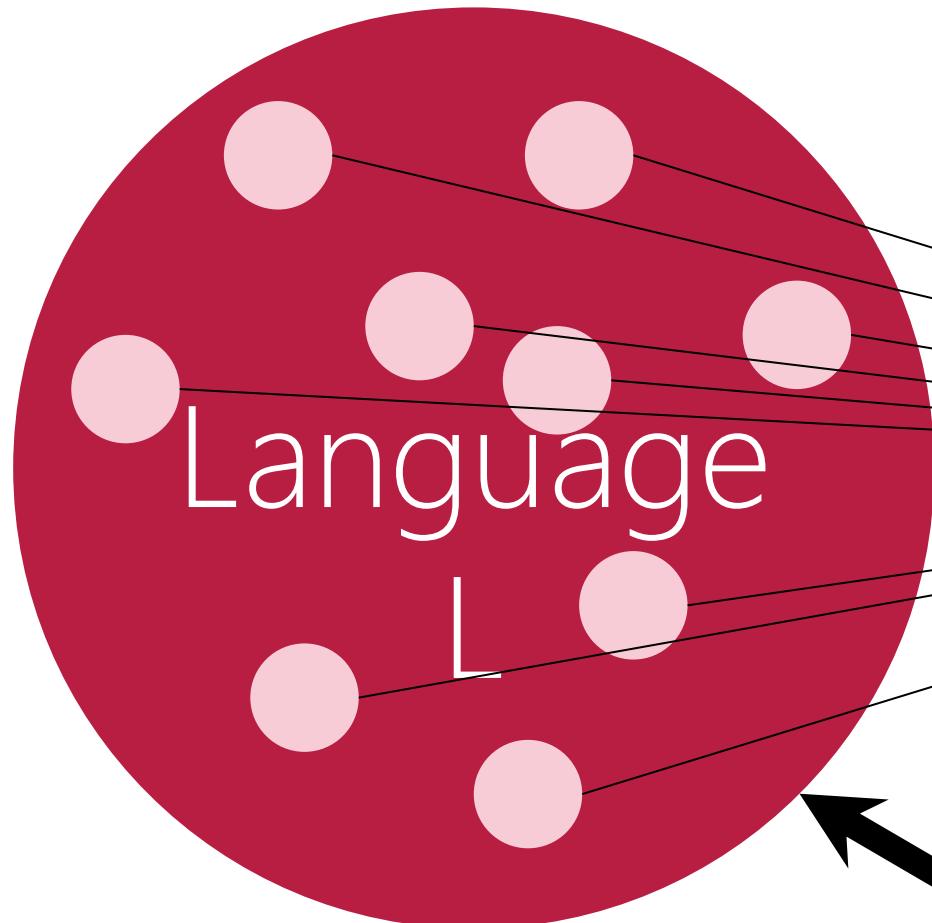
Golden Model

- The language space is available.
- Enumerate the model space to find M^* , assuming it exists!

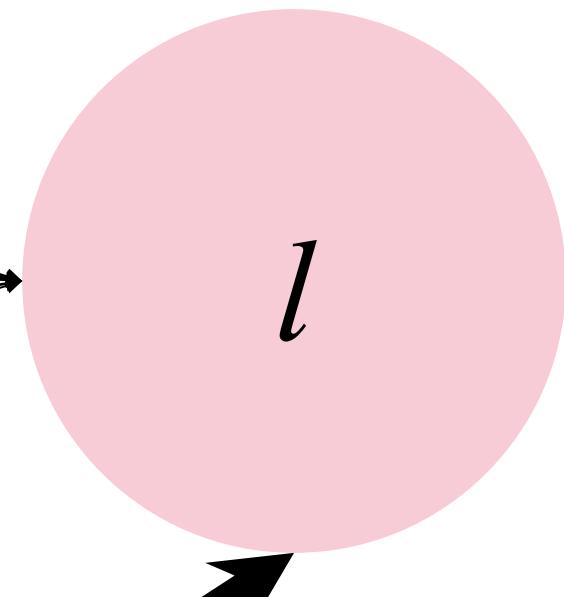
Find M^*

Golden Model

- ~~The language space is available.~~ → Randomly selected subset
- Enumerate the model space to find M^* , assuming it exists!



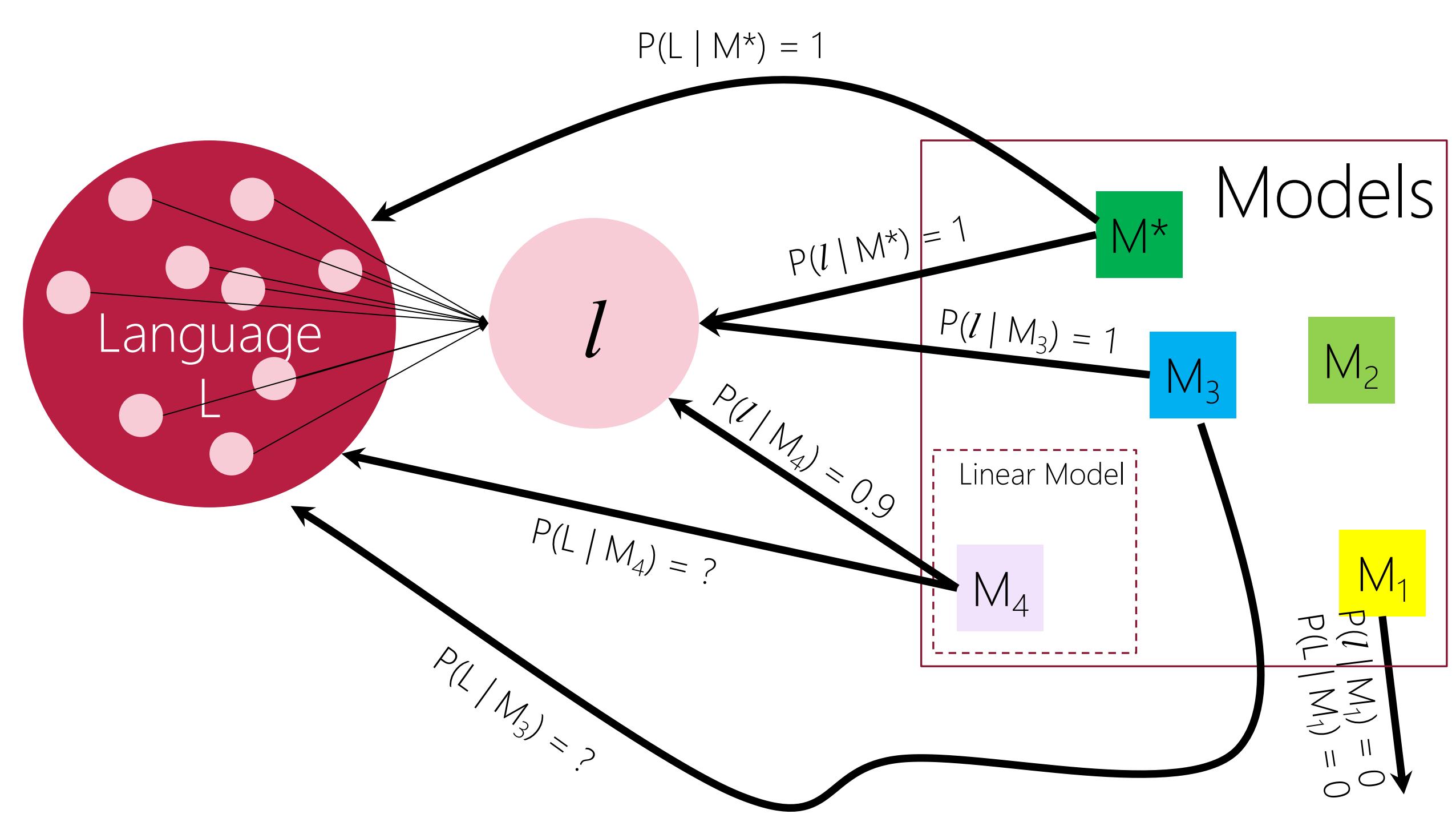
$$P(L | M^*) = 1$$
$$P(l | M^*) = 1$$



Find M^* M^{\wedge}

Golden Silver Model

- The language space is available. → Randomly selected subset
 - Enumerate the model space to find M^* , assuming it exists!
- Search subspace, e.g., just linear models!



$$M^* = \operatorname{argmax}_{M \in \text{Models}} P(l | M)$$

- The language space is available. → Randomly selected subset
- Enumerate the model space to find M^* , assuming it exists!
→ Search subspace, e.g., just linear models!

$$\hat{M} = \operatorname{argmax}_{M \in \text{Models}} P(l | M)$$

$P(l | M)$ is also called Likelihood $\rightarrow \mathcal{L}(l | M)$

Argmax is also called Maximum Likelihood Estimation \rightarrow MLE

Argmax sometimes is expressed in log \rightarrow Log Likelihood

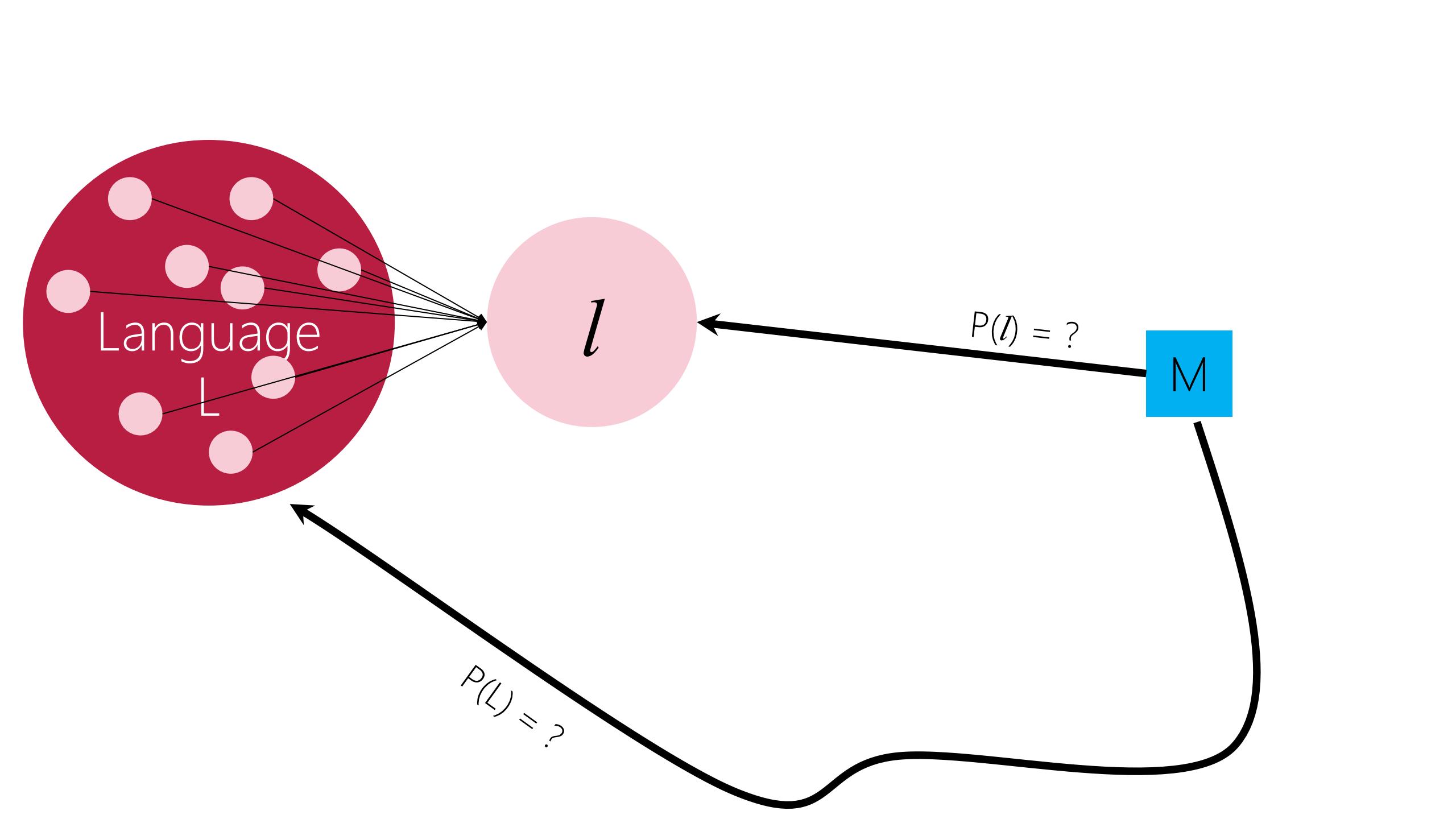
$$P(l \mid M) = \mathcal{L}(l \mid M) = \frac{P(l, M)}{P(M)}$$

$P(M)$ is necessary for models that has prior distribution over their parameters such as Gaussian models $N(\mu, \Sigma)$.

$$M^{\wedge} = \operatorname{argmax}_{M \in \text{Models}} P(l, M)$$

Assuming M is in effect, the likelihood of M is:

$$P(l, M) = P_M(l) = \mathcal{L}_M(l)$$



Likelihood of a Language Modeling

l : ['The', 'course', 'COMP8730', 'is', 'about', 'nlp', '.', 'The', 'instructor', "s", 'name', 'is', 'Hossein', '.', 'There', 'are', '13', 'students', 'in', 'the', 'class', '.']

M1 = |token|-gram model = 22-gram model $\rightarrow P(l) = \mathcal{L}_{\text{22-gram}}(l) = 1$

Likelihood of a Language Modeling

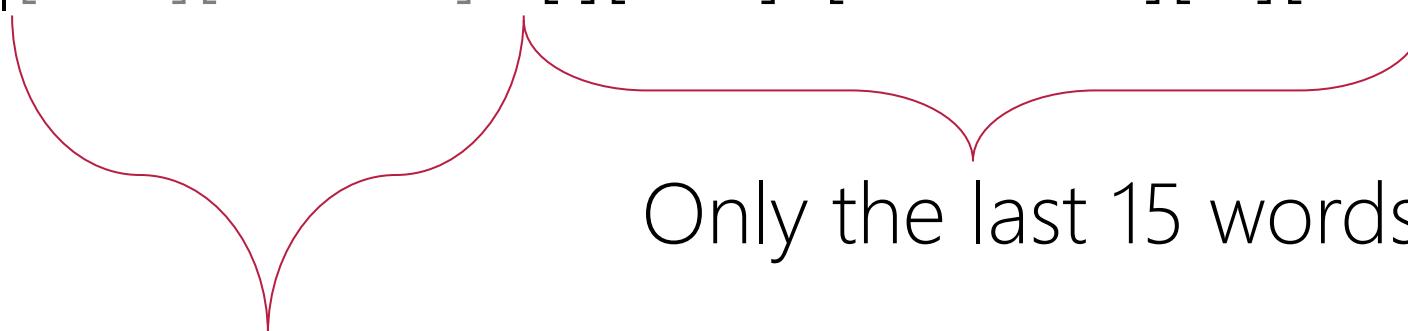
l : ['The', 'course', 'COMP8730', 'is', 'about', 'nlp', '.', 'The', 'instructor', "s", 'name', 'is', 'Hossein', '.', 'There', 'are', '13', 'students', 'in', 'the', 'class', '.']

M1 = |token|-gram model = 22-gram model $\rightarrow P(l) = \mathcal{L}_{22\text{-gram}}(l) = 1$

M2 = |vocab|-gram model = 16-gram model $\rightarrow P(l) = \mathcal{L}_{16\text{-gram}}(l)$

$P([\text{the}])P([\text{course}][\text{The}])P([\text{COMP8730}][\text{The}][\text{course}]) \dots$

$P([\text{class}][\text{The}][\text{course}]\dots[.][\text{The}]\dots[\text{students}][\text{in}][\text{the}])P(\dots)$



Only the last 15 words

These are not considered in 16-gram model

Likelihood of a Language Modeling

\mathbf{l} : ['The', 'course', 'COMP8730', 'is', 'about', 'nlp', '.', 'The', 'instructor', "s", 'name', 'is', 'Hossein', '.', 'There', 'are', '13', 'students', 'in', 'the', 'class', '.']

M1 = |token|-gram model = 22-gram model $\rightarrow P(\mathbf{l}) = \mathcal{L}_{\text{22-gram}}(\mathbf{l}) = 1$

M2 = |vocab|-gram model = 16-gram model $\rightarrow P(\mathbf{l}) = \mathcal{L}_{\text{16-gram}}(\mathbf{l}) = ?$

M3 = 3-gram model = $P(\mathbf{l}) = \mathcal{L}_{\text{3-gram}}(\mathbf{l})$

$P([\text{The}])P([\text{course}][\text{The}])P([\text{COMP8730}][\text{The}][\text{course}])P([\text{is}][\text{course}][\text{CO}]\text{MP8730})P([\text{about}][\text{COMP8730}][\text{is}]) \dots$



Only the last 2 words

Likelihood of a Language Modeling

\mathbf{l} : ['The', 'course', 'COMP8730', 'is', 'about', 'nlp', '.', 'The', 'instructor', "s", 'name', 'is', 'Hossein', '.', 'There', 'are', '13', 'students', 'in', 'the', 'class', '.']

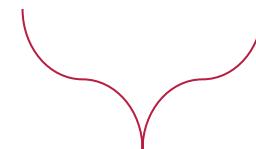
M1 = |token|-gram model = 22-gram model $\rightarrow P(\mathbf{l}) = \mathcal{L}_{\text{22-gram}}(\mathbf{l}) = 1$

M2 = |vocab|-gram model = 16-gram model $\rightarrow P(\mathbf{l}) = \mathcal{L}_{\text{16-gram}}(\mathbf{l}) = ?$

M3 = 3-gram model = $P(\mathbf{l}) = \mathcal{L}_{\text{3-gram}}(\mathbf{l}) = ?$

M4 = 2-gram model = $P(\mathbf{l}) = \mathcal{L}_{\text{2-gram}}(\mathbf{l})$

$P([\text{The}])P([\text{course}])P([\text{The}])P([\text{COMP8730}])P([\text{course}])P([\text{is}])P([\text{COMP8730}])P([\text{about}])P([\text{is}])P([\text{nlp}])P([\text{about}])\dots$



Only the last word \rightarrow Markovian

Likelihood of a Language Modeling

\mathbf{l} : ['The', 'course', 'COMP8730', 'is', 'about', 'nlp', '.', 'The', 'instructor', "s", 'name', 'is', 'Hossein', '.', 'There', 'are', '13', 'students', 'in', 'the', 'class', '.']

M1 = |token|-gram model = 22-gram model $\rightarrow P(\mathbf{l}) = \mathcal{L}_{\text{22-gram}}(\mathbf{l}) = 1$

M2 = |vocab|-gram model = 16-gram model $\rightarrow P(\mathbf{l}) = \mathcal{L}_{\text{16-gram}}(\mathbf{l}) = ?$

M3 = 3-gram model = $P(\mathbf{l}) = \mathcal{L}_{\text{3-gram}}(\mathbf{l}) = ?$

M4 = 2-gram model = $P(\mathbf{l}) = \mathcal{L}_{\text{2-gram}}(\mathbf{l}) = ?$

M5 = 1-gram model = $P(\mathbf{l}) = \mathcal{L}_{\text{1-gram}}(\mathbf{l})$

$P([\text{The}])P([\text{course}])P([\text{COMP8730}])P([\text{is}])P([\text{about}])P([\text{nlp}])\dots$

No history!

Likelihood of a Language Modeling

l : ['The', 'course', 'COMP8730', 'is', 'about', 'nlp', '.', 'The', 'instructor', "s", 'name', 'is', 'Hossein', '.', 'There', 'are', '13', 'students', 'in', 'the', 'class', '.']

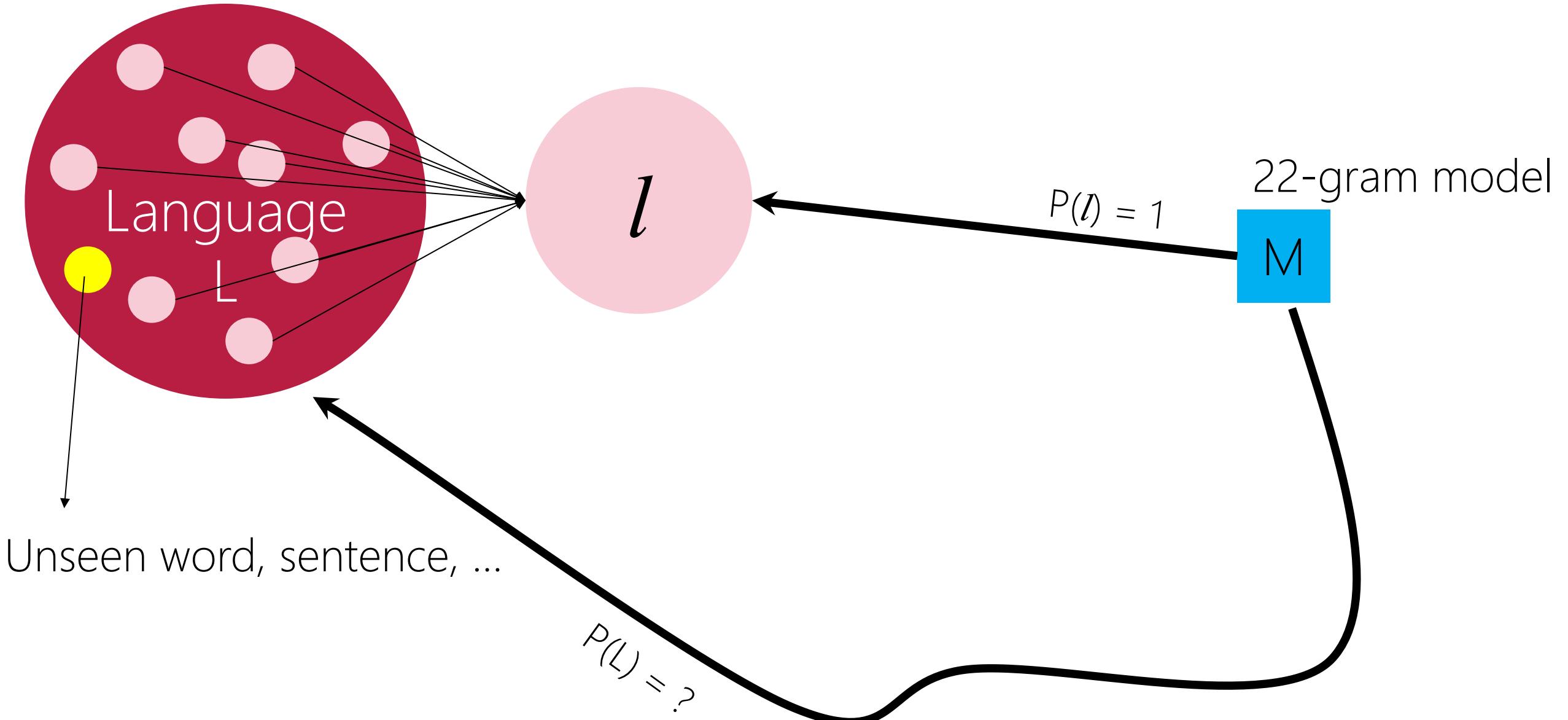
M1 = |token|-gram model = 22-gram model $\rightarrow P(l) = \mathcal{L}_{\text{22-gram}}(l) = 1$

M2 = |vocab|-gram model = 16-gram model $\rightarrow P(l) = \mathcal{L}_{\text{16-gram}}(l) = ?$

M3 = 3-gram model = $P(l) = \mathcal{L}_{\text{3-gram}}(l) = ?$

M4 = 2-gram model = $P(l) = \mathcal{L}_{\text{2-gram}}(l) = ?$

M5 = 1-gram model = $P(l) = \mathcal{L}_{\text{1-gram}}(l) = ?$



Likelihood of a Language Modeling

l : ['The', 'course', 'COMP8730', 'is', 'about', 'nlp', '.', 'The', 'instructor', "s", 'name', 'is', 'Hossein', '.', 'There', 'are', '13', 'students', 'in', 'the', 'class', '.']

l' : ['Hossein', 'is', 'the', 'name', '.']

M1 = 22-gram model $\rightarrow P(l') = \mathcal{L}_{22\text{-gram}}(l') = ?$

M2 = 16-gram model $\rightarrow P(l') = \mathcal{L}_{16\text{-gram}}(l') = ?$

M3 = 3-gram model = $P(l') = \mathcal{L}_{3\text{-gram}}(l') = ?$

M4 = 2-gram model = $P(l') = \mathcal{L}_{2\text{-gram}}(l') = ?$

M5 = 1-gram model = $P(l') = \mathcal{L}_{1\text{-gram}}(l') = ?$

Likelihood of a Language Modeling

l : ['The', 'course', 'COMP8730', 'is', 'about', 'nlp', '.', 'The', 'instructor', "s", 'name', 'is', 'Hossein', '.', 'There', 'are', '13', 'students', 'in', 'the', 'class', '.']

l' : ['Hossein', 'is', 'the', 'name', '.']

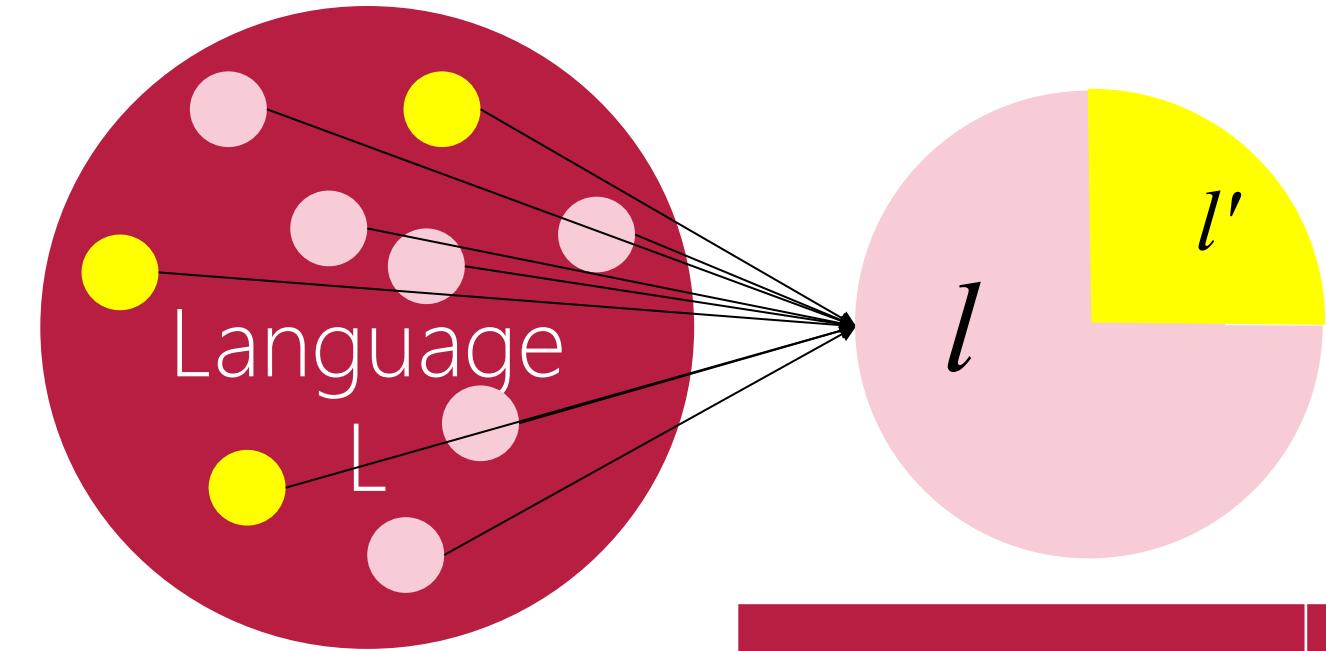
M1 = 22-gram model $\rightarrow P(l') = \mathcal{L}_{22\text{-gram}}(l') = 0$

M2 = 16-gram model $\rightarrow P(l') = \mathcal{L}_{16\text{-gram}}(l') = 0$

M3 = 3-gram model = $P(l') = \mathcal{L}_{3\text{-gram}}(l')$ = 0

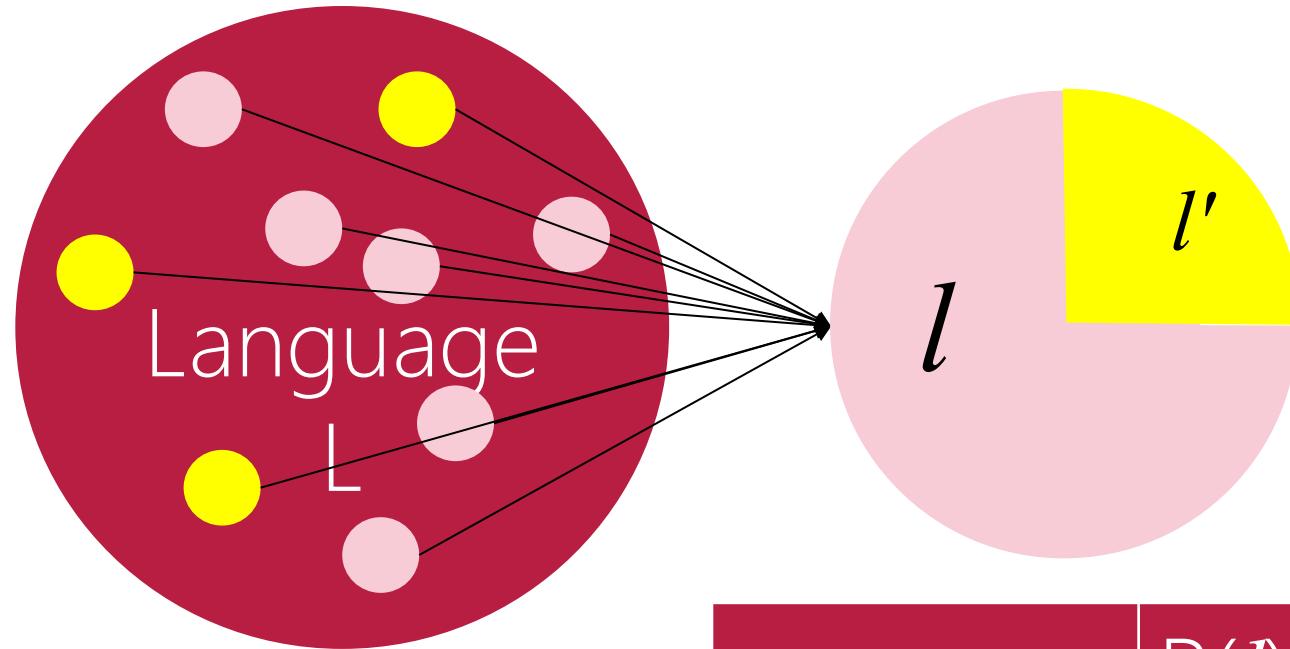
M4 = 2-gram model = $P(l') = \mathcal{L}_{2\text{-gram}}(l')$ = 0

M5 = 1-gram model = $P(l') = \mathcal{L}_{1\text{-gram}}(l')$ = nonzero!



	$P(l) = \mathcal{L}_M(l)$	$P(l') = \mathcal{L}_M(l')$
M1	High	High
M2	High	Low
M3	Low	High
M4	Low	Low

P is trained on \mathcal{I} and does not know anything about \mathcal{I}'



	$P(\mathcal{I}) = \mathcal{L}_M(\mathcal{I})$	$P(\mathcal{I}') = \mathcal{L}_M(\mathcal{I}')$	
M1	High	High	
M2	High	Low	Overfit
M3	Low	High	Not Possible
M4	Low	Low	Underfit

Evaluating Language Models

Quantitative → Perplexity

How perplexed (confused) a language model is to generate **seen** or **unseen** stream of tokens! lower perplexity, the better!

$$\text{Perplexity} = \text{PP}(l') = P(l')^{-1/|l'|}$$

$$= |l'| \sqrt{\frac{1}{P(l')}}$$

P is trained on l and does not know anything about l'

Higher $P(l')$, lower perplexity, the better!

$$\text{Perplexity} = \text{PP}(l') = \sqrt{\frac{1}{P(l')}}$$

$$\text{Unigram approx.: } \sqrt{\frac{1}{P(l')}} = \sqrt{\frac{1}{\prod_{k=1}^{|l'|} P(w_k)}} ; w_i \in l'$$

$$\text{Bigram approx.: } \sqrt{\frac{1}{P(l')}} = \sqrt{\frac{1}{\prod_{k=1}^{|l'|} P(w_i | w_{i-1})}} ; w_{i-1}w_i \in l'$$

$$\text{Trigram approx.: } \sqrt{\frac{1}{P(l')}} = \sqrt{\frac{1}{\prod_{k=1}^{|l'|} P(w_i | w_{i-2}w_{i-1})}} ; w_{i-2}w_{i-1}w_i \in l'$$

Corpus: Wall Street Journal

Size: 38 million words

Vocab (Types): 19,979

	Unigram	Bigram	Trigram
Perplexity	962	170	109

Test Set: 1.5 million words

Zeros!
Unseen

Likelihood of a Language Modeling

l : ['The', 'course', 'COMP8730', 'is', 'about', 'nlp', '.', 'The', 'instructor', "s", 'name', 'is', 'Hossein', '.', 'There', 'are', '13', 'students', 'in', 'the', 'class', '.']

l' : ['Hossein', 'is', 'the', 'name', 'of', 'a', 'person', '.']

M1 = 22-gram model $\rightarrow P(l') = \mathcal{L}_{22\text{-gram}}(l') = 0$

M2 = 16-gram model $\rightarrow P(l') = \mathcal{L}_{16\text{-gram}}(l') = 0$

M3 = 3-gram model = $P(l') = \mathcal{L}_{3\text{-gram}}(l') = 0$

M4 = 2-gram model = $P(l') = \mathcal{L}_{2\text{-gram}}(l') = 0$

M5 = 1-gram model = $P(l') = \mathcal{L}_{1\text{-gram}}(l') = 0$ (Why?)

Zeros!

Unseen

Not all unigrams are available in training set! E.g., [of]
Not all bigrams are available in training set! E.g., [Hossein][is]
Not all trigrams are available in training set! ...

Zeros!

- 1) Vocabulary + <UNK>
- 2) Train Vocabulary + Learn Unseen Tokens (Subwords)
- 3) Smoothing

<UNK>

learn the stats of unseen tokens

- 1) Pick a dictionary D
- 2) From $w \in l$ such that $w \notin D$, (OOV) replace it with <UNK>
- 3) Train model
- 4) At test, from $w \in l'$, if $w \notin l$ (unseen), replace it with <UNK>

Likelihood of a Language Modeling

D : ['The', 'course', 'is', 'about', 'instructor', 'name', 'There', 'are', '13', 'students', 'in', 'class']

l : ['The', 'course', '<UNK>', 'is', 'about', '<UNK>', '<UNK>', 'The', 'instructor', '<UNK>', 'name', 'is', '<UNK>', '<UNK>', 'There', 'are', '13', 'students', 'in', 'the', 'class', '.']

l' : ['<UNK>', 'is', 'the', 'name', '<UNK>', '<UNK>', '<UNK>', '<UNK>']

M1 = 22-gram model $\rightarrow P(l') = \mathcal{L}_{22\text{-gram}}(l') = 0$

M2 = 16-gram model $\rightarrow P(l') = \mathcal{L}_{16\text{-gram}}(l') = 0$

M3 = 3-gram model = $P(l') = \mathcal{L}_{3\text{-gram}}(l') = 0$

M4 = 2-gram model = $P(l') = \mathcal{L}_{2\text{-gram}}(l') = 0$

M5 = 1-gram model = $P(l') = \mathcal{L}_{1\text{-gram}}(l') = \text{Nonzero! (Why?)}$

<UNK>

learn the stats of unseen tokens

1) Pick a dictionary $D \rightarrow \text{small}$

Gives higher probability (lower perplexity) at test

All the model should generate is stream of <UNK>s!

Likelihood of a Language Modeling

D : ['The', 'is']

l : ['The', '<UNK>', '<UNK>', 'is', '<UNK>', '<UNK>', '<UNK>', 'The', '<UNK>', '<UNK>', '<UNK>', '<UNK>', '<UNK>', '<UNK>', '<UNK>', '<UNK>', '<UNK>', 'the', '<UNK>', '<UNK>']

l' : ['<UNK>', 'is', '<UNK>', 'the', '<UNK>', '<UNK>', '<UNK>', '<UNK>']

M1 = 22-gram model $\rightarrow P(l') = \mathcal{L}_{22\text{-gram}}(l') = 0$

M2 = 16-gram model $\rightarrow P(l') = \mathcal{L}_{16\text{-gram}}(l') = 0$

M3 = 3-gram model = $P(l') = \mathcal{L}_{3\text{-gram}}(l') = 0$

M4 = 2-gram model = $P(l') = \mathcal{L}_{2\text{-gram}}(l') = \text{Nonzero!}$

M5 = 1-gram model = $P(l') = \mathcal{L}_{1\text{-gram}}(l') = \text{Nonzero!}$

Subwords

learn the unseen tokens from seen tokens

Please refer to: Text Normalization → Learn to Tokenize

Smoothing

- 1) Add-1 (Laplace) or Add- k , $\rightarrow k=\{1,2,\dots\}$
- 2) Backoff
- 3) Interpolation
- 4) ...

Add-k

Add k unit to all counts, so zero entries become k. Add-1 is also called Laplace.

$$\text{Unigram model: } P(w_i) = \frac{\#w_i + k}{|\text{tokens}| + k \times |\text{vocabs}|}$$

$$\text{Bigram model: } P(w_i|w_{i-1}) = \frac{\#(w_{i-1}w_i) + k}{\#(w_{i-1}) + k \times ?}$$

$$\text{Trigram model: } P(w_i|w_{i-2}w_{i-1}) = \frac{\#(w_{i-2}w_{i-1}w_i) + k}{\#(w_{i-2}w_{i-1}) + k \times ?}$$

Backoff

If n-gram does have not seen, try (n-1)-gram.

Trigram model: $P(w_i|w_{i-2}w_{i-1}) = \frac{\#(w_{i-2}w_{i-1}w_i)}{\#(w_{i-2}w_{i-1})} = 0$



Bigram model: $P(w_i|w_{i-1}) = \frac{\#(w_{i-1}w_i)}{\#(w_{i-1})} = 0$



Unigram model: $P(w_i) = \frac{\#w_i}{|\text{tokens}|}$

Interpolation

$P(n\text{-gram})$ is linear interpolation of all $(n-i)$ -grams: $i=\{1,2,\dots, n-1\}$.

$$\text{Trigram model: } P(w_i|w_{i-2}w_{i-1}) = \lambda_1 \frac{\#(w_{i-2}w_{i-1}w_i)}{\#(w_{i-2}w_{i-1})} +$$

$$\text{Bigram model: } P(w_i|w_{i-1}) = \lambda_2 \frac{\#(w_{i-1}w_i)}{\#(w_{i-1})} +$$

$$\text{Unigram model: } P(w_i) = \lambda_3 \frac{\#w_i}{|\text{tokens}|} \quad \sum \lambda_i = 1$$

Kneser-Ney Smoothing

Kneser, R. and Ney, H. (1995). Improved back-off for M-gram language modeling. In ICASSP-95, Vol. 1, 181–184.

Chen, S. F. and Goodman, J. (1999). An empirical study of smoothing techniques for language modeling. Computer Speech and Language, 13, 359–394.