

(Generative) Language Modeling

Frank Ferraro – ferraro@umbc.edu

CMSC 473/673

Goal of Language Modeling

$$p_{\theta}([...text...])$$

Learn a **probabilistic model** of **text**

Accomplished through observing **text** and updating **model parameters** to make **text** more likely

Two Perspectives: Prediction vs. Generation

Prediction

Given observed word tokens $w_1 \dots w_{N-1}$, create a classifier p to predict the next word w_N

$$p(w_N = v | w_1 \dots w_{N-1})$$

Generation

Two Perspectives: Prediction vs. Generation

Prediction

Given observed word tokens $w_1 \dots w_{N-1}$, create a classifier p to predict the next word w_N

$$p(w_N = v | w_1 \dots w_{N-1}), \text{ e.g.,}$$
$$p(w_N = \text{meowed} | \text{The, fluffy, cat})$$

Generation

Two Perspectives: Prediction vs. Generation

Prediction

Given observed word tokens $w_1 \dots w_{N-1}$, create a classifier p to predict the next word w_N

$$p(w_N = v | w_1 \dots w_{N-1}), \text{ e.g.,} \\ p(w_N = \text{meowed} | \text{The, fluffy, cat})$$

Generation

Develop a probabilistic model p to *explain/score* the word sequence $w_1 \dots w_N$

$$p(w_1 \dots w_N), \text{ e.g.,} \\ p(\text{The, fluffy, cat, meowed})$$

Design Question 1: What Part of Language Do We Estimate?

$$p_{\theta}([...text...])$$

Is *[...text.]* a

- Full document?
- Sequence of sentences?
- Sequence of words?
- Sequence of characters?

A: It's task-dependent!

Design Question 2: How do we estimate robustly?

$$p_{\theta}([\dots \textit{typo-text} \dots])$$

What if *[...text.]* has a typo?

Design Question 3: How do we generalize?

$$p_{\theta}([\dotsynonymous\text{-}text\dots])$$

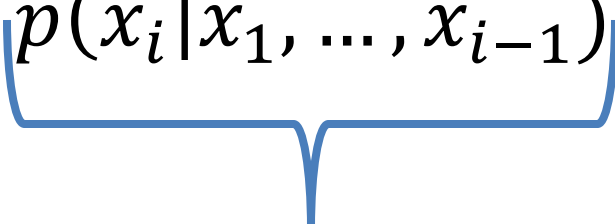
What if *[...text..]* has a word (or character or...) we've never seen before?

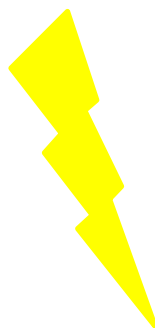
Key Idea: Probability Chain Rule

$$p(x_1, x_2, \dots, x_S) =$$

$$p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2) \cdots p(x_S | x_1, \dots, x_{S-1})$$

Key Idea: Probability Chain Rule

$$p(x_1, x_2, \dots, x_S) = p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2) \cdots p(x_S | x_1, \dots, x_{S-1}) = \prod_i^S p(x_i | x_1, \dots, x_{i-1})$$




Language modeling is about how to estimate each of these factors in {great, good, sufficient, ...} ways

Problem: Develop a Probabilistic Email Classifier

Input: an email (all text)

Output (Google categories):

Primary, Social, Forums, Spam

$$\operatorname{argmax}_y p(\text{label } Y = y \mid \text{email } X)$$

Approach #1: Discriminatively trained

Approach #2: Using Bayes rule

Classify Using Bayes Rule

$$p(\text{label } Y \mid \text{email } X) \propto p(X \mid Y) * p(Y)$$

Classify Using Bayes Rule

$$p(\text{label } Y \mid \text{email } X) \propto p(X \mid Y) * p(Y)$$

Q: Why is $p(Y \mid X)$ what we want to model?

Classify Using Bayes Rule

$$p(\text{label } Y \mid \text{email } X) \propto p(X \mid Y) * p(Y)$$



$$p(\text{Primary} \mid \text{Won't you please donate?}) \propto p(\text{Won't you please donate?} \mid \text{Primary}) p(\text{Primary})$$

A Closer Look at $p(\text{Primary})$

This is the **prior probability** of each *class*

Answers the question: without knowing anything specific about a document, how likely is each class?

A Closer Look at $p(\text{Primary})$

This is the **prior probability** of each *class*

Answers the question: without knowing anything specific about a document, how likely is each class?

Q: What's an easy way to estimate it?

A Closer Look at $p(\text{Won't you please donate?} \mid \text{Primary})$

This is a *class specific* language model

A Closer Look at $p(\text{Won't you please donate?} \mid \text{Primary})$

This is a *class specific* language model

$p(\text{Won't you please donate?} \mid \text{Primary})$ is different from

$p(\text{Won't you please donate?} \mid \text{Social})$ is different from

$p(\text{Won't you please donate?} \mid \text{Forums}) \dots$

A Closer Look at $p(\text{Won't you please donate?} \mid \text{Primary})$

This is a *class specific* language model

To learn $p(\text{Won't you please donate?} \mid \text{Class})$:

For each class **Class**:

Get a bunch of **Class** documents D_{Class}

Learn a new language model p_{Class} on just D_{Class}

Language Models & Smoothing

- Maximum likelihood (MLE): simple counting
- Other count-based models

- Laplace smoothing, add- λ
- Interpolation models
- Discounted backoff
- Interpolated (modified) Kneser-Ney
- Good-Turing
- Witten-Bell

- Maxent n-gram models

- Neural n-gram models

- Recurrent/autoregressive NNs

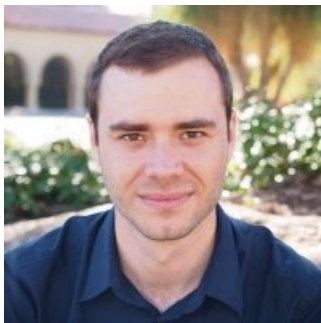
Easy to implement

Advanced/
out of
scope

Featureful LMs

Feedforward LMs

Super modern



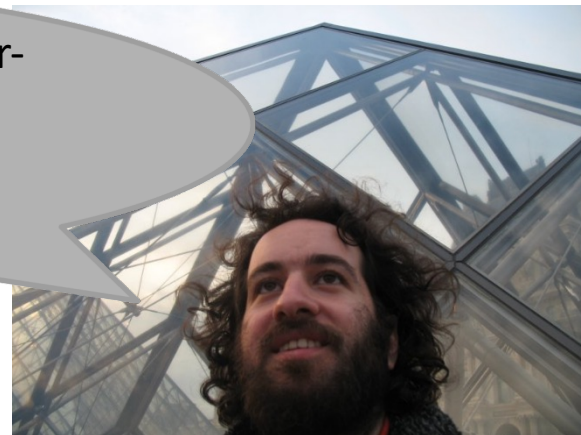
“The Unreasonable Effectiveness of
Recurrent Neural Networks”

<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

“The Unreasonable Effectiveness of Character-
level Language Models”

(and why RNNs are still cool)

<http://nbviewer.jupyter.org/gist/yoavg/d76121dfde2618422139>



Language Models & Smoothing

- **Maximum likelihood (MLE): simple counting**

- Other count-based models

- Laplace smoothing, add- λ

- Interpolation models

- Discounted backoff

- Interpolated (modified) Kneser-Ney

- Good-Turing

- Witten-Bell

Easy to implement

Advanced/
out of scope

- Maxent n-gram models

- Neural n-gram models

- Recurrent/autoregressive NNs

Featureful LMs

Feedforward LMs

Super modern

N-Grams

Maintaining an entire inventory over sentences
could be too much to ask

Store “smaller” pieces?

p(Colorless green ideas sleep furiously)

N-Grams

Maintaining an entire *joint* inventory over sentences could be too much to ask

Store “smaller” pieces?

$$p(\text{Colorless green ideas sleep furiously}) = p(\text{Colorless}) *$$

N-Grams

Maintaining an entire *joint* inventory over sentences could be too much to ask

Store “smaller” pieces?

$$\begin{aligned} p(\text{Colorless green ideas sleep furiously}) = \\ p(\text{Colorless}) * \\ p(\text{green} \mid \text{Colorless}) * \end{aligned}$$

N-Grams

Maintaining an entire *joint* inventory over sentences could be too much to ask

Store “smaller” pieces?

$$\begin{aligned} p(\text{Colorless green ideas sleep furiously}) = & \\ & p(\text{Colorless}) * \\ & p(\text{green} \mid \text{Colorless}) * \\ & p(\text{ideas} \mid \text{Colorless green}) * \\ & p(\text{sleep} \mid \text{Colorless green ideas}) * \\ & p(\text{furiously} \mid \text{Colorless green ideas sleep}) \end{aligned}$$

N-Grams

$p(\text{furiously} \mid \text{Colorless green ideas sleep})$

How much does “Colorless” influence the choice of “furiously?”

N-Grams

$p(\text{furiously} \mid \text{Colorless green ideas sleep})$

How much does “Colorless” influence the choice of “furiously?”

Remove history and contextual info

N-Grams

$p(\text{furiously} \mid \text{Colorless green ideas sleep})$

How much does “Colorless” influence the choice of “furiously?”

Remove history and contextual info

$p(\text{furiously} \mid \text{Colorless green ideas sleep}) \approx$
 $p(\text{furiously} \mid \text{Colorless-green-ideas sleep})$

N-Grams

p(furiously | Colorless green ideas sleep)

How much does “Colorless” influence the choice of “furiously?”

Remove history and contextual info

$p(\text{furiously} \mid \text{Colorless green ideas sleep}) \approx p(\text{furiously} \mid \text{ideas sleep})$

N-Grams

$$\begin{aligned} p(\text{Colorless green ideas sleep furiously}) = & \\ & p(\text{Colorless}) * \\ & p(\text{green} \mid \text{Colorless}) * \\ & p(\text{ideas} \mid \text{Colorless green}) * \\ & p(\text{sleep} \mid \text{Colorless green ideas}) * \\ & p(\text{furiously} \mid \text{Colorless green ideas sleep}) \end{aligned}$$

N-Grams

$$\begin{aligned} p(\text{Colorless green ideas sleep furiously}) = & \\ & p(\text{Colorless}) * \\ & p(\text{green} \mid \text{Colorless}) * \\ & p(\text{ideas} \mid \text{Colorless green}) * \\ & p(\text{sleep} \mid \text{Colorless green ideas}) * \\ & p(\text{furiously} \mid \text{Colorless green ideas sleep}) \end{aligned}$$

Trigrams

$$\begin{aligned} p(\text{Colorless green ideas sleep furiously}) = & \\ & p(\text{Colorless}) * \\ & p(\text{green} \mid \text{Colorless}) * \\ & p(\text{ideas} \mid \text{Colorless green}) * \\ & p(\text{sleep} \mid \text{green ideas}) * \\ & p(\text{furiously} \mid \text{ideas sleep}) \end{aligned}$$

Trigrams

$$\begin{aligned} p(\text{Colorless green ideas sleep furiously}) = & \\ & p(\text{Colorless}) * \\ & p(\text{green} \mid \text{Colorless}) * \\ & p(\text{ideas} \mid \text{Colorless green}) * \\ & p(\text{sleep} \mid \text{green ideas}) * \\ & p(\text{furiously} \mid \text{ideas sleep}) \end{aligned}$$

Trigrams

$$\begin{aligned} p(\text{Colorless green ideas sleep furiously}) = & \\ & p(\text{Colorless} \mid \langle \text{BOS} \rangle \langle \text{BOS} \rangle) * \\ & p(\text{green} \mid \langle \text{BOS} \rangle \text{Colorless}) * \\ & p(\text{ideas} \mid \text{Colorless green}) * \\ & p(\text{sleep} \mid \text{green ideas}) * \\ & p(\text{furiously} \mid \text{ideas sleep}) \end{aligned}$$

Consistent notation: Pad the left with <BOS> (beginning of sentence) symbols

Trigrams

$$\begin{aligned} p(\text{Colorless green ideas sleep furiously}) = & \\ & p(\text{Colorless} \mid \langle \text{BOS} \rangle \langle \text{BOS} \rangle) * \\ & p(\text{green} \mid \langle \text{BOS} \rangle \text{Colorless}) * \\ & p(\text{ideas} \mid \text{Colorless green}) * \\ & \quad p(\text{sleep} \mid \text{green ideas}) * \\ & p(\text{furiously} \mid \text{ideas sleep}) * \\ & p(\langle \text{EOS} \rangle \mid \text{sleep furiously}) \end{aligned}$$

Consistent notation: Pad the left with $\langle \text{BOS} \rangle$ (beginning of sentence) symbols

Fully proper distribution: Pad the right with a single $\langle \text{EOS} \rangle$ symbol

N-Gram Terminology

n	Commonly called	History Size (Markov order)	Example
1	unigram	0	$p(\text{furiously})$

N-Gram Terminology

n	Commonly called	History Size (Markov order)	Example
1	unigram	0	$p(\text{furiously})$
2	bigram	1	$p(\text{furiously} \mid \text{sleep})$

N-Gram Terminology

n	Commonly called	History Size (Markov order)	Example
1	unigram	0	$p(\text{furiously})$
2	bigram	1	$p(\text{furiously} \mid \text{sleep})$
3	trigram (3-gram)	2	$p(\text{furiously} \mid \text{ideas sleep})$

N-Gram Terminology

n	Commonly called	History Size (Markov order)	Example
1	unigram	0	$p(\text{furiously})$
2	bigram	1	$p(\text{furiously} \mid \text{sleep})$
3	trigram (3-gram)	2	$p(\text{furiously} \mid \text{ideas sleep})$
4	4-gram	3	$p(\text{furiously} \mid \text{green ideas sleep})$
n	n-gram	n-1	$p(w_i \mid w_{i-n+1} \dots w_{i-1})$

N-Gram Probability

$$p(w_1, w_2, w_3, \dots, w_S) =$$

$$\prod_{i=1}^S p(w_i | w_{i-N+1}, \dots, w_{i-1})$$

Count-Based N-Grams (Unigrams)

$$p(\text{item}) \propto \textit{count}(\text{item})$$

Count-Based N-Grams (Unigrams)

$$p(z) \propto \textit{count}(z)$$

Count-Based N-Grams (Unigrams)

$$\begin{aligned} p(\mathbf{z}) &\propto \text{count}(\mathbf{z}) \\ &= \frac{\text{count}(\mathbf{z})}{\sum_v \text{count}(\mathbf{v})} \end{aligned}$$

Diagram illustrating the formula for the probability of a unigram \mathbf{z} based on counts:

- The variable \mathbf{z} in the numerator is labeled "word type".
- The variable \mathbf{v} in the denominator is labeled "word type".

Count-Based N-Grams (Unigrams)

$$\begin{array}{c} \text{word type} \quad \quad \quad \text{word type} \\ \downarrow \quad \quad \quad \downarrow \\ p(\mathbf{z}) \propto \text{count}(\mathbf{z}) \\ = \frac{\text{count}(\mathbf{z})}{W} \\ \quad \quad \quad \uparrow \\ \quad \quad \quad \text{number of tokens observed} \end{array}$$

Count-Based N-Grams (Unigrams)

The film got a great opening and the film went on to become a hit .

Word (Type) z	Raw Count $\text{count}(z)$	Normalization	Probability $p(z)$
The	1		
film	2		
got	1		
a	2		
great	1		
opening	1		
and	1		
the	1		
went	1		
on	1		
to	1		
become	1		
hit	1		
.	1		

Count-Based N-Grams (Unigrams)

The film got a great opening and the film went on to become a hit .

Word (Type) z	Raw Count $\text{count}(z)$	Normalization	Probability $p(z)$
The	1	16	
film	2		
got	1		
a	2		
great	1		
opening	1		
and	1		
the	1		
went	1		
on	1		
to	1		
become	1		
hit	1		
.	1		

Count-Based N-Grams (Unigrams)

The film got a great opening and the film went on to become a hit .

Word (Type) z	Raw Count $\text{count}(z)$	Normalization	Probability $p(z)$
The	1	16	1/16
film	2		1/8
got	1		1/16
a	2		1/8
great	1		1/16
opening	1		1/16
and	1		1/16
the	1		1/16
went	1		1/16
on	1		1/16
to	1		1/16
become	1		1/16
hit	1		1/16
.	1		1/16

Count-Based N-Grams (Trigrams)

order matters in
conditioning



order matters in
count



$$p(z|x, y) \propto \textit{count}(x, y, z)$$

Count of the
sequence of items

“x y z”

Count-Based N-Grams (Trigrams)

order matters in
conditioning



order matters in
count



$$p(z|x, y) \propto \textit{count}(x, y, z)$$

$\textit{count}(x, y, z) \neq \textit{count}(x, z, y) \neq \textit{count}(y, x, z) \neq \dots$

Count-Based N-Grams (Trigrams)

$$\begin{aligned} p(z|x, y) &\propto \textit{count}(x, y, z) \\ &= \frac{\textit{count}(x, y, z)}{\sum_v \textit{count}(x, y, v)} \end{aligned}$$

Count-Based N-Grams (Trigrams)

The film got a great opening and the film went on to become a hit .

Context: x y	Word (Type): z	Raw Count	Normalization	Probability $p(z x y)$
The film	The	0	1	0/1
The film	film	0		0/1
The film	got	1		1/1
The film	went	0		0/1
...				
a great	great	0	1	0/1
a great	opening	1		1/1
a great	and	0		0/1
a great	the	0		0/1
...				

Count-Based N-Grams (Lowercased Trigrams)

the film got a great opening and the film went on to become a hit .

Context: x y	Word (Type): z	Raw Count	Normalization	Probability: $p(z x y)$
the film	the	0	2	0/2
the film	film	0		0/2
the film	got	1		1/2
the film	went	1		1/2
...				
a great	great	0	1	0/1
a great	opening	1		1/1
a great	and	0		0/1
a great	the	0		0/1
...				

Implementation: EOS Padding

Create an end of sentence (“chunk”) token
<EOS>

Don't estimate $p(\text{<BOS>} \mid \text{<EOS>})$

Training & Evaluation:

1. Identify “chunks” that are relevant (sentences, paragraphs, documents)
2. Append the <EOS> token to the end of the chunk
3. Train or evaluate LM as normal

Implementation: Memory Issues

Let V = vocab size, W = number of *observed* n-grams

Often, $W \ll V$

Dense count representation: $O(V^n)$, but many entries will be zero

Sparse count representation: $O(W)$

Sometimes selective precomputation is helpful (e.g., normalizers)

Implementation: Unknown words

Create an unknown word token <UNK>

Training:

1. Create a fixed lexicon L of size V
2. Change any word not in L to <UNK>
3. Train LM as normal

Evaluation:

Use UNK probabilities for any word not in training

A Closer Look at Count-based $p(\text{Won't you please donate?} \mid \text{Primary})$

This is a *class specific* language model

To learn $p(\text{Won't you please donate?} \mid \text{Class})$:

For each class **Class**:

Get a bunch of **Class** documents D_{Class}

Learn a new language model p_{Class} on just D_{Class}

Two Ways to Learn **Class**-specific Count-based Language Models

1. Create different count
table(s) $c_{\text{Class}}(\dots)$ for each
Class

e.g., record separate trigram
counts for **Primary** vs. **Social**
vs. **Forums** vs. **Spam**

Two Ways to Learn **Class**-specific Count-based Language Models

1. Create different count table(s)

$c_{\text{Class}}(\dots)$ for each **Class**

e.g., record separate trigram counts
for **Primary** vs. **Social** vs. **Forums** vs.
Spam

OR

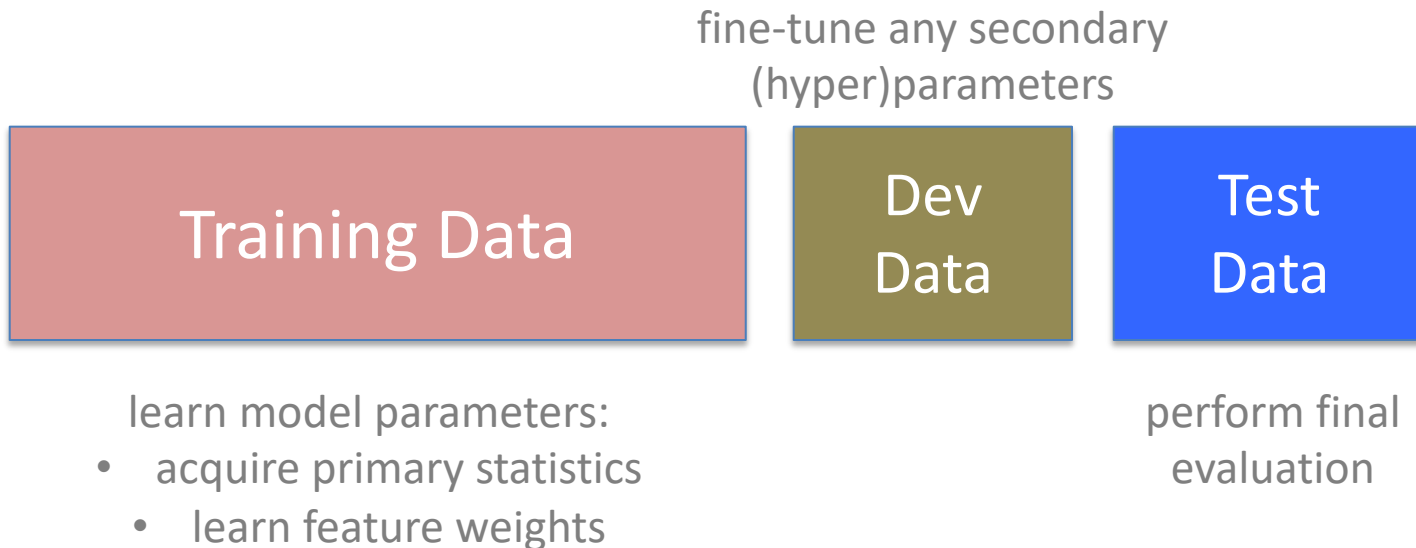
2. Add a dimension to your existing
tables $c(\text{Class}, \dots)$

e.g., record how often each trigram
occurs within **Primary** vs. **Social** vs.
Forums vs. **Spam** documents

Evaluating Language Models

What is “correct?”

What is working “well?”



DO NOT TUNE ON THE TEST DATA

Evaluating Language Models

What is “correct?”

What is working “well?”

Extrinsic: Evaluate LM in downstream task

Test an MT, ASR, etc. system and see which LM does better

Propagate & conflate errors

$$p(\text{label } Y \mid \text{doc } X) \propto p(X \mid Y) * p(Y)$$

Evaluating Language Models

What is “correct?”

What is working “well?”

Extrinsic: Evaluate LM in downstream task

Test an MT, ASR, etc. system and see which LM does better

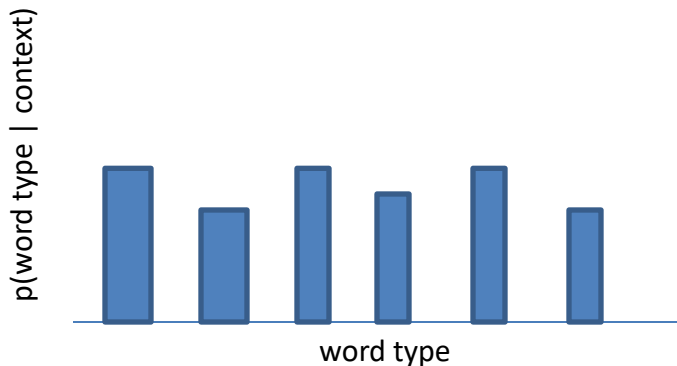
Propagate & conflate errors

Intrinsic: Treat LM as its own downstream task

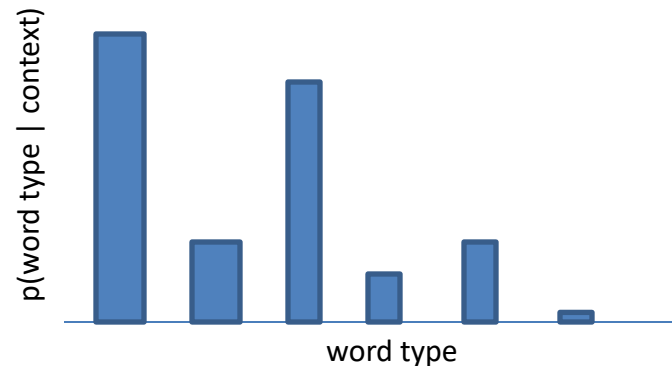
Use perplexity (from information theory)

Perplexity: Average “Surprisal”

Lower is better : lower perplexity → less surprised



Less certain →
More surprised →
Higher perplexity



More certain →
Less surprised →
Lower perplexity

Perplexity


Lower is better : lower perplexity → less surprised

$$\text{perplexity} = \exp(\text{avg xent})$$

Perplexity

Lower is better : lower perplexity → less surprised

*e.g., n-gram history
(n-1 items)*

$$\text{perplexity} = \exp\left(\frac{-1}{M} \sum_{i=1}^M \log p(w_i | h_i)\right)$$


Perplexity

Lower is better : lower perplexity → less surprised

$$\text{perplexity} = \exp\left(\frac{-1}{M} \sum_{i=1}^M \log p(w_i | h_i)\right)$$

≥ 0, ≤ 1: higher

Perplexity

Lower is better : lower perplexity → less surprised

$$\text{perplexity} = \exp\left(\frac{-1}{M} \sum_{i=1}^M \underbrace{\log p(w_i | h_i)}_{\substack{\geq 0, \leq 1: \text{higher}}} \right)$$

≤ 0 : higher

Perplexity

Lower is better : lower perplexity → less surprised

$$\text{perplexity} = \exp\left(\frac{-1}{M} \sum_{i=1}^M \log p(w_i | h_i)\right)$$

Annotations for the formula:

- For the sum $\sum_{i=1}^M \log p(w_i | h_i)$: ≤ 0 : higher
- For the sum $\sum_{i=1}^M \log p(w_i | h_i)$: $\geq 0, \leq 1$: higher
- For the entire expression $\exp\left(\frac{-1}{M} \sum_{i=1}^M \log p(w_i | h_i)\right)$: ≤ 0 , higher

Perplexity

Lower is better : lower perplexity → less surprised

$$\text{perplexity} = \exp\left(\frac{-1}{M} \sum_{i=1}^M \log p(w_i | h_i)\right)$$

Annotations for the equation above:

- Annotation for the exponent: ≥ 0 , lower is better
- Annotation for the sum: ≤ 0 : higher
- Annotation for the log term: $\geq 0, \leq 1$: higher
- Annotation for the overall expression: ≤ 0 , higher

Perplexity

Lower is better : lower perplexity → less surprised

$$\text{perplexity} = \exp\left(\frac{-1}{M} \sum_{i=1}^M \log p(w_i | h_i)\right)$$

The diagram includes several annotations with blue brackets and arrows:

- A bracket above the entire expression is labeled ≥ 0 , lower is better.
- A bracket above the sum $\sum_{i=1}^M \log p(w_i | h_i)$ is labeled ≤ 0 : higher.
- A bracket below the sum $\sum_{i=1}^M \log p(w_i | h_i)$ is labeled $\geq 0, \leq 1$: higher.
- A bracket below the entire expression is labeled ≤ 0 , higher.
- A bracket below the entire expression is labeled ≥ 0 , lower.

Perplexity

Lower is better : lower perplexity → less surprised

base must be the same

$$\text{perplexity} = \exp\left(\frac{-1}{M} \sum_{i=1}^M \log p(w_i | h_i)\right)$$

Annotations:

- ≥ 0 , lower is better (applies to the outer \exp function)
- ≤ 0 : higher (applies to the inner \log function)
- $\geq 0, \leq 1$: higher (applies to the probability $p(w_i | h_i)$)
- ≤ 0 , higher (applies to the sum $\sum_{i=1}^M \log p(w_i | h_i)$)
- ≥ 0 , lower (applies to the entire expression)

Perplexity

Lower is better : lower perplexity → less surprised

$$\text{perplexity} = \exp\left(\frac{-1}{M} \sum_{i=1}^M \log p(w_i | h_i)\right)$$

$$= \sqrt[M]{\underbrace{\prod_{i=1}^M \frac{1}{p(w_i | h_i)}}_{\text{weighted geometric average}}}$$

weighted
geometric
average

How to Compute Average Perplexity

- If you have a list of the probabilities for each observed n-gram “token:”

```
numpy.exp(-numpy.mean(numpy.log(probs_per_trigram_token)))
```

- If you have a list of observed n-gram “types” t and counts c, and log-prob. function lp:

```
numpy.exp(-numpy.mean(c*lp(t) for (t, c) in ngram_types.items()))
```

- If you’re computing a cross-entropy loss function (e.g., in Pytorch):

```
loss_fn = torch.nn.CrossEntropyLoss(reduction='mean')  
torch.exp(loss_fn(input_data))
```

What are the tri-grams for “The film , a hit !”

Trigrams	MLE $p(\text{trigram})$
<BOS> <BOS> The	1
<BOS> The film	1
The film ,	0
film , a	0
, a hit	0
a hit !	0
hit ! <EOS>	0
Perplexity	???

What are the tri-grams for “The film , a hit !”

Trigrams	MLE $p(\text{trigram})$
<BOS> <BOS> The	1
<BOS> The film	1
The film ,	0
film , a	0
, a hit	0
a hit !	0
hit ! <EOS>	0
Perplexity	Infinity

What are the tri-grams for “The film , a hit !”

Trigrams	MLE $p(\text{trigram})$	Smoothed $p(\text{trigram})$
<BOS> <BOS> The	1	2/17
<BOS> The film	1	2/17
The film ,	0	1/17
film , a	0	1/16
, a hit	0	1/16
a hit !	0	1/17
hit ! <EOS>	0	1/16
Perplexity	Infinity	???

What are the tri-grams for “The film , a hit !”

Trigrams	MLE $p(\text{trigram})$	Smoothed $p(\text{trigram})$
<BOS> <BOS> The	1	2/17
<BOS> The film	1	2/17
The film ,	0	1/17
film , a	0	1/16
, a hit	0	1/16
a hit !	0	1/17
hit ! <EOS>	0	1/16
Perplexity	Infinity	13.59

0s Are Not Your (Language Model's) Friend

$$p(\text{item}) \propto \text{count}(\text{item}) = 0 \rightarrow \\ p(\text{item}) = 0$$

0 probability \rightarrow item is *impossible*

0s annihilate: $x*y*z*0 = 0$

Language is creative:

new words keep appearing

existing words could appear in known contexts

How much do you trust your data?

Language Models & Smoothing

~~• Maximum likelihood (MLE): simple counting~~

• Other count-based models

– Laplace smoothing, add- λ

– Interpolation models

– Discounted backoff

– Interpolated (modified) Kneser-Ney

– Good-Turing

– Witten-Bell

• Maxent n-gram models

• Neural n-gram models

• Recurrent/autoregressive NNs

Easy to
implement

Advanced/
out of
scope

Featureful LMs

Feedforward LMs

Super modern

Add- λ estimation

Other names: Laplace
smoothing, Lidstone
smoothing

Pretend we saw each
word λ more times
than we did

$$p(z) \propto \textit{count}(z) + \lambda$$

Add λ to all the
counts

Add- λ estimation

Other names: Laplace
smoothing, Lidstone
smoothing

Pretend we saw each
word λ more times
than we did

Add λ to all the
counts

$$p(\mathbf{z}) \propto \text{count}(\mathbf{z}) + \lambda$$
$$= \frac{\text{count}(\mathbf{z}) + \lambda}{\sum_v (\text{count}(v) + \lambda)}$$

Add- λ estimation

Other names: Laplace
smoothing, Lidstone
smoothing

Pretend we saw each
word λ more times
than we did

Add λ to all the
counts

$$p(\mathbf{z}) \propto \frac{\text{count}(\mathbf{z}) + \lambda}{W + V\lambda}$$

tokens # types

Add- λ N-Grams (Unigrams)

The film got a great opening and the film went on to become a hit .

Word (Type)	Raw Count	Norm	Prob.	Add- λ Count	Add- λ Norm.	Add- λ Prob.
The	1	16	1/16			
film	2		1/8			
got	1		1/16			
a	2		1/8			
great	1		1/16			
opening	1		1/16			
and	1		1/16			
the	1		1/16			
went	1		1/16			
on	1		1/16			
to	1		1/16			
become	1		1/16			
hit	1		1/16			
.	1	1/16				

Add-1 N-Grams (Unigrams)

The film got a great opening and the film went on to become a hit .

Word (Type)	Raw Count	Norm	Prob.	Add-1 Count	Add-1 Norm.	Add-1 Prob.
The	1	16	1/16	2		
film	2		1/8	3		
got	1		1/16	2		
a	2		1/8	3		
great	1		1/16	2		
opening	1		1/16	2		
and	1		1/16	2		
the	1		1/16	2		
went	1		1/16	2		
on	1		1/16	2		
to	1		1/16	2		
become	1		1/16	2		
hit	1		1/16	2		
.	1	1/16	2			

Add-1 N-Grams (Unigrams)

The film got a great opening and the film went on to become a hit .

Word (Type)	Raw Count	Norm	Prob.	Add-1 Count	Add-1 Norm.	Add-1 Prob.
The	1	16	1/16	2	$16 + 14 * 1 = 30$	
film	2		1/8	3		
got	1		1/16	2		
a	2		1/8	3		
great	1		1/16	2		
opening	1		1/16	2		
and	1		1/16	2		
the	1		1/16	2		
went	1		1/16	2		
on	1		1/16	2		
to	1		1/16	2		
become	1		1/16	2		
hit	1		1/16	2		
.	1		1/16	2		

Add-1 N-Grams (Unigrams)

The film got a great opening and the film went on to become a hit .

Word (Type)	Raw Count	Norm	Prob.	Add-1 Count	Add-1 Norm.	Add-1 Prob.
The	1	16	1/16	2	$16 + 14 * 1 = 30$	=1/15
film	2		1/8	3		=1/10
got	1		1/16	2		=1/15
a	2		1/8	3		=1/10
great	1		1/16	2		=1/15
opening	1		1/16	2		=1/15
and	1		1/16	2		=1/15
the	1		1/16	2		=1/15
went	1		1/16	2		=1/15
on	1		1/16	2		=1/15
to	1		1/16	2		=1/15
become	1		1/16	2		=1/15
hit	1		1/16	2		=1/15
.	1		1/16	2		=1/15

An Extended Trigram Example

The film got a great opening and the film went on to become a hit .

Q: With OOV, EOS, and BOS,
how many types (for
normalization)?

Context: x y	Word (Type): z	Raw Count	Add-1 count	Norm.	Probability $p(z x y)$	
The film	The	0				
The film	film	0				
The film	got	1				
The film	went	0				
...						
The film	OOV	0				
The film	EOS	0				
...						
a great	great	0				
a great	opening	1				
a great	and	0				
a great	the	0				
...						

An Extended Trigram Example

The film got a great opening and the film went on to become a hit .

Q: With OOV, EOS, and BOS, how many types (for normalization)?

A: 16
(why don't we count BOS?)

Context: x y	Word (Type): z	Raw Count	Add-1 count	Norm.	Probability $p(z x y)$	
The film	The	0				
The film	film	0				
The film	got	1				
The film	went	0				
...						
The film	OOV	0				
The film	EOS	0				
...						
a great	great	0				
a great	opening	1				
a great	and	0				
a great	the	0				
...						

An Extended Trigram Example

The film got a great opening and the film went on to become a hit .

Q: With OOV, EOS, and BOS, how many types (for normalization)?

A: 16
(why don't we count BOS?)

Context: x y	Word (Type): z	Raw Count	Add-1 count	Norm.	Probability $p(z x y)$
The film	The	0	1	17 (=1+16*1)	1/17
The film	film	0	1		1/17
The film	got	1	2		2/17
The film	went	0	1		1/17
...					...
The film	OOV	0	1		1/17
The film	EOS	0	1		1/17
...					
a great	great	0	1	17	1/17
a great	opening	1	2		2/17
a great	and	0	1		1/17
a great	the	0	1		1/17
...					

An Extended Trigram Example

The film got a great opening and the film went on to become a hit .

Context: x y	Word (Type): z	Raw Count	Add-1 count	Norm.	Probability $p(z x y)$	
The film	The	0	1	17 (=1+16*1)	1/17	
The film	film	0	1		1/17	
The film	got	1	2		2/17	
The film	went	0	1		1/17	
...					...	
The film	OOV	0	1		1/17	
The film	EOS	0	1		1/17	
...						
a great	great	0	1	17	1/17	
a great	opening	1	2		2/17	
a great	and	0	1		1/17	
a great	the	0	1		1/17	
...						

Q: What is the perplexity for the sentence "The film , a hit !"

What are the tri-grams for “The film , a hit !”

Trigrams	MLE $p(\text{trigram})$
<BOS> <BOS> The	1
<BOS> The film	1
The film ,	0
film , a	0
, a hit	0
a hit !	0
hit ! <EOS>	0

What are the tri-grams for “The film , a hit !”

Trigrams	MLE $p(\text{trigram})$
<BOS> <BOS> The	1
<BOS> The film	1
The film ,	0
film , a	0
, a hit	0
a hit !	0
hit ! <EOS>	0

What are the tri-grams for “The film , a hit !”

Trigrams	MLE $p(\text{trigram})$	UNK-ed trigrams
<BOS> <BOS> The	1	<BOS> <BOS> The
<BOS> The film	1	<BOS> The film
The film ,	0	The film <UNK>
film , a	0	film <UNK> a
, a hit	0	<UNK> a hit
a hit !	0	a hit <UNK>
hit ! <EOS>	0	hit <UNK> <EOS>

What are the tri-grams for “The film , a hit !”

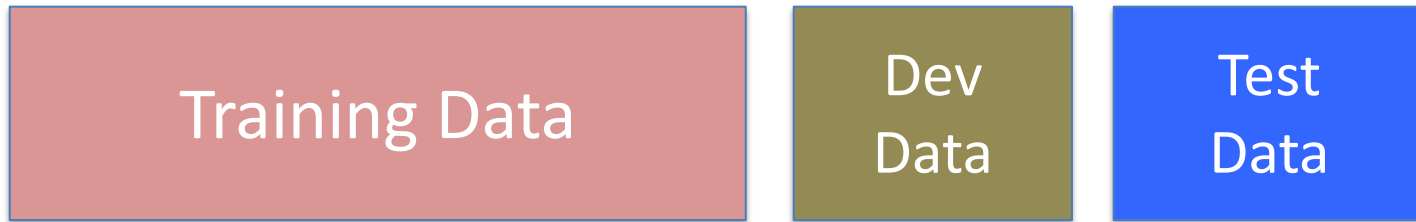
Trigrams	MLE $p(\text{trigram})$	UNK-ed trigrams	Smoothed $p(\text{trigram})$
<BOS> <BOS> The	1	<BOS> <BOS> The	2/17
<BOS> The film	1	<BOS> The film	2/17
The film ,	0	The film <UNK>	1/17
film , a	0	film <UNK> a	1/16
, a hit	0	<UNK> a hit	1/16
a hit !	0	a hit <UNK>	1/17
hit ! <EOS>	0	hit <UNK> <EOS>	1/16

What are the tri-grams for “The film , a hit !”

Trigrams	MLE $p(\text{trigram})$	UNK-ed trigrams	Smoothed $p(\text{trigram})$
<BOS> <BOS> The	1	<BOS> <BOS> The	2/17
<BOS> The film	1	<BOS> The film	2/17
The film ,	0	The film <UNK>	1/17
film , a	0	film <UNK> a	1/16
, a hit	0	<UNK> a hit	1/16
a hit !	0	a hit <UNK>	1/17
hit ! <EOS>	0	hit <UNK> <EOS>	1/16

Setting Hyperparameters

Use a **development** corpus



Choose λ s to maximize the probability of dev data:

- Fix the N-gram probabilities (on the training data)
- Then search for λ s that give largest probability to held-out set:



Advanced
topic

Other Kinds of Smoothing

- ~~Maximum likelihood (MLE):
simple counting~~
- Laplace smoothing, add- λ
- Interpolation models
- Discounted backoff
- Interpolated (modified)
Kneser-Ney
- Good-Turing
- Witten-Bell

Interpolated (modified) Kneser-Ney

Idea: How “productive” is a context?

How many different word *types* v
appear in a context x, y

Good-Turing

Partition words into classes of
occurrence

Smooth class statistics

Properties of classes are likely to
predict properties of other classes

Witten-Bell

Idea: Every observed type was at
some point novel

Give MLE prediction for novel *type*
occurring

Language Models & Smoothing

- ~~Maximum likelihood (MLE): simple counting~~

- ~~Other count based models~~

 - ~~– Laplace smoothing, add λ~~

 - ~~– Interpolation models~~

 - ~~– Discounted backoff~~

 - ~~– Interpolated (modified) Kneser-Ney~~

 - ~~– Good Turing~~

 - ~~– Witten-Bell~~

Easy to implement

Advanced/
out of
scope

- **Maxent n-gram models**

- Neural n-gram models

- Recurrent/autoregressive NNs

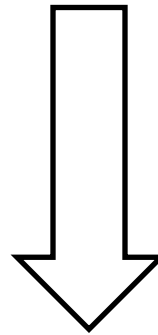
Featureful LMs

Feedforward LMs

Super modern

Maxent Models as Featureful n-gram Language Models

$$p(\text{Colorless green ideas sleep furiously} \mid \text{Label}) = p(\text{Colorless} \mid \text{Label}, \langle \text{BOS} \rangle) * \dots * p(\langle \text{EOS} \rangle \mid \text{Label}, \text{furiously})$$



Model each n-gram term with a maxent model

$$p(x_i \mid \mathbf{y}, x_{i-N+1:i-1}) = \text{maxent}(\mathbf{y}, x_{i-N+1:i-1}, x_i)$$

*generatively trained:
learn to model (class-specific) language*

Language Model with Maxent n-grams

$$p_n(\text{document} | y) = \prod_{i=1}^M \text{maxent}(y, \underbrace{x_{i-n+1:i-1}, x_i}_{\text{n-gram}})$$

label

n-gram

$$= \prod_{i=1}^M \frac{\exp(\theta_{x_i}^T f(y, x_{i-n+1:i-1}))}{\sum_{x'} \exp(\theta_{x'}^T f(y, x_{i-n+1:i-1}))}$$

Iterate through all possible output vocab types x' ---just like in count-based LMs

What Should These Features Do?

$$p(x_i | y, x_{i-N+1:i-1}) = \text{maxent}(y, x_{i-N+1:i-1}, x_i), \text{ e.g.,}$$

$$\begin{aligned} & p(\text{sleep} | y, \text{green}, \text{ideas}) = \\ & \text{maxent}(y, x_{i-2,i-1} = (\text{green}, \text{ideas}), x_i = \text{sleep}) \\ \propto & \exp(\theta_{x_i=\text{sleep}}^T f(y, x_{i-2,i-1} = (\text{green}, \text{ideas}))) \end{aligned}$$

(in-class discussion)

N-gram Language Models

given some context...

w_{i-3}

w_{i-2}

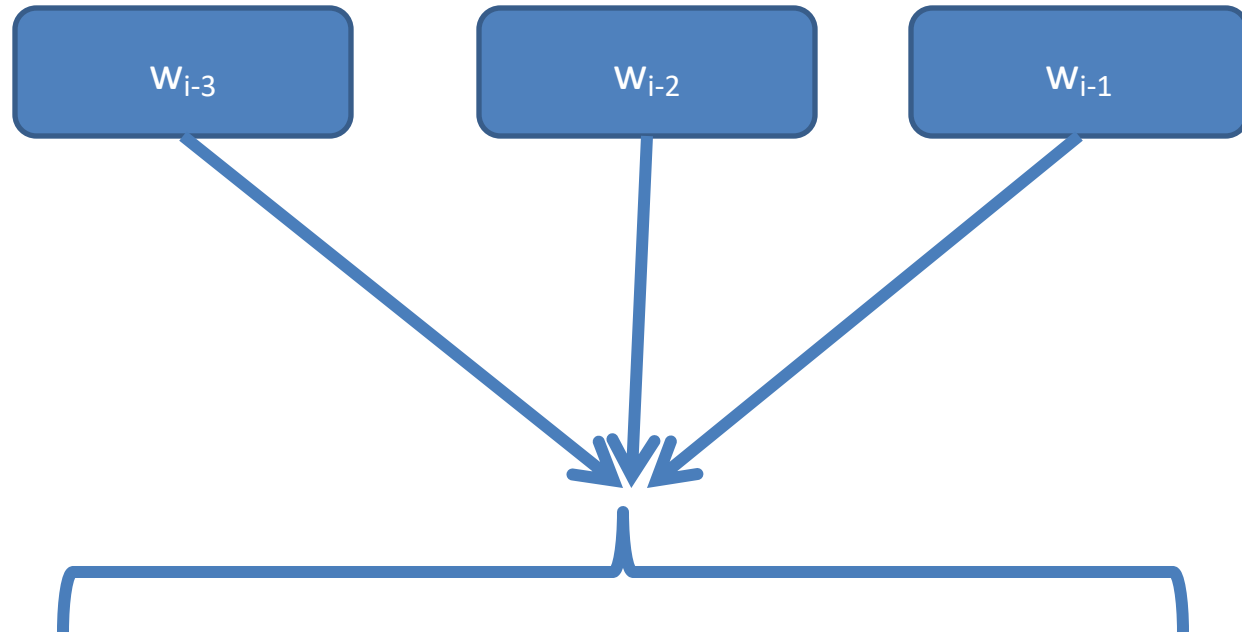
w_{i-1}

predict the next word

w_i

N-gram Language Models

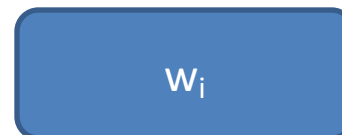
given some context...



compute beliefs about what is likely...

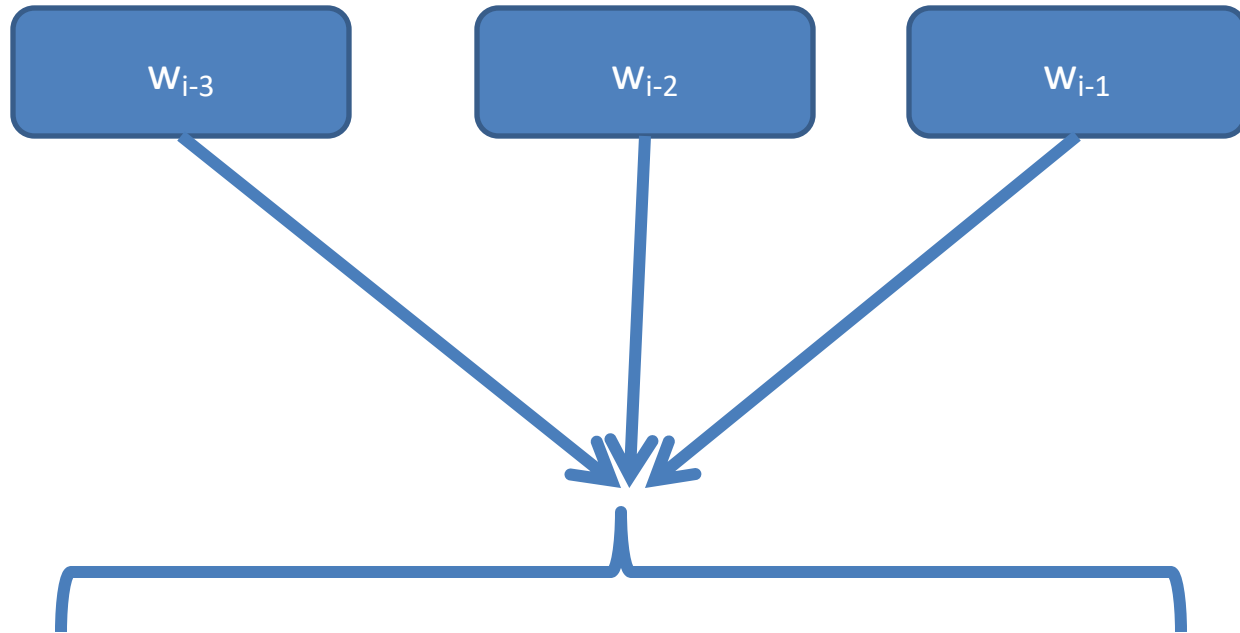
$$p(w_i | w_{i-3}, w_{i-2}, w_{i-1}) \propto \text{count}(w_{i-3}, w_{i-2}, w_{i-1}, w_i)$$

predict the next word



N-gram Language Models

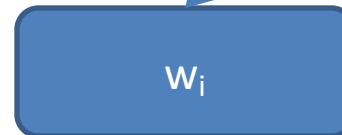
given some context...



compute beliefs about what is likely...

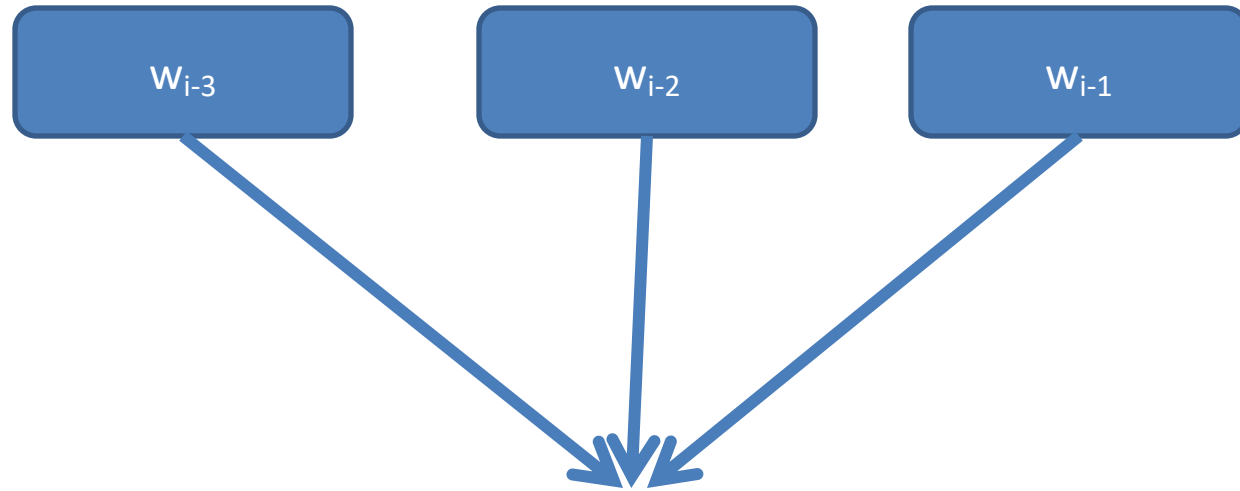
$$p(w_i | w_{i-3}, w_{i-2}, w_{i-1}) \propto \text{count}(w_{i-3}, w_{i-2}, w_{i-1}, w_i)$$

predict the next word



Maxent Language Models

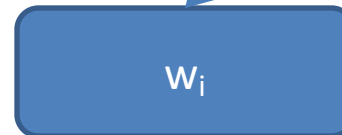
given some context...



compute beliefs about what is likely...

$$p(w_i | w_{i-3}, w_{i-2}, w_{i-1}) = \text{softmax}(\theta_{w_i} \cdot f(w_{i-3}, w_{i-2}, w_{i-1}))$$

predict the next word



A Closer Look at Maxent $p(\text{Won't you please donate?} \mid \text{Primary})$

This is a *class-based* language model, but incorporate the label into the features

To learn $p(\text{Won't you please donate?} \mid \text{Class})$:

Define features f that make use of the specific label **Class**

Unlike count-based models, you don't *need* "separate" models here

Language Models & Smoothing

~~• Maximum likelihood (MLE): simple counting~~

~~• Other count based models~~

~~— Laplace smoothing, add λ~~

~~— Interpolation models~~

~~— Discounted backoff~~

~~— Interpolated (modified) Kneser-Ney~~

~~— Good Turing~~

~~— Witten-Bell~~

~~• Maxent n-gram models~~

• **Neural n-gram models**

• Recurrent/autoregressive NNs

Easy to
implement

Advanced/
out of
scope

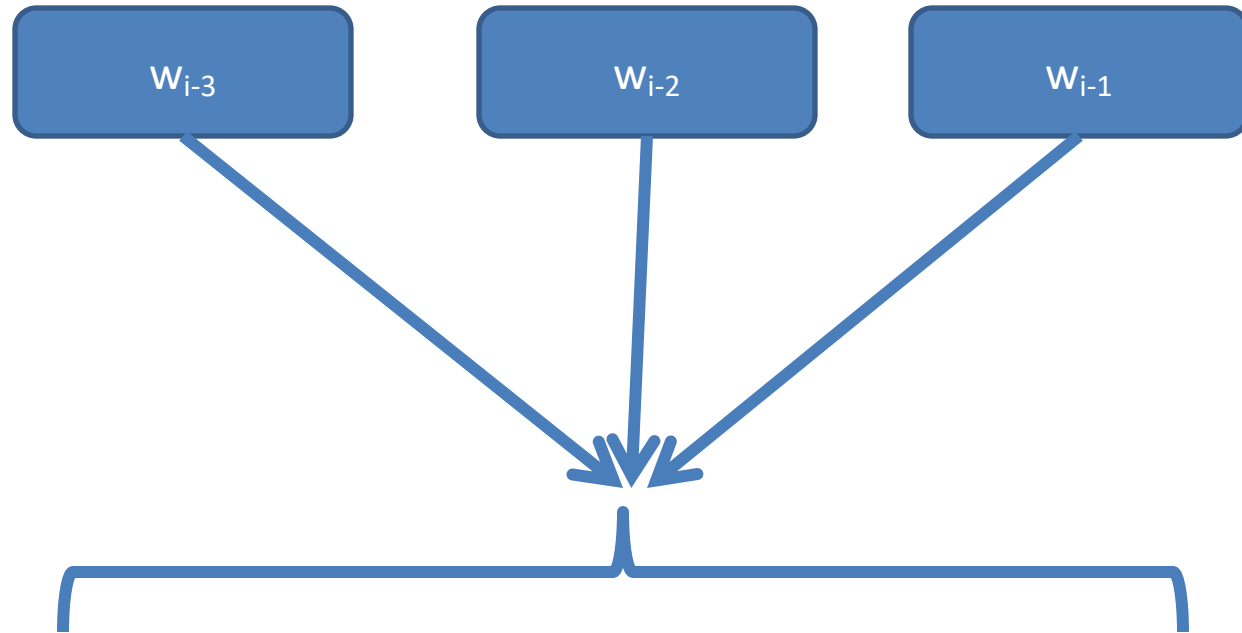
Featureful LMs

Feedforward LMs

Super modern

Maxent Language Models

given some context...

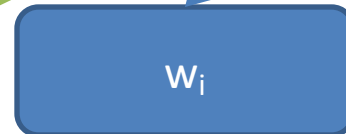


compute beliefs about what is likely...

$$p(w_i | w_{i-3}, w_{i-2}, w_{i-1}) = \text{softmax}(\theta_{w_i} \cdot f(w_{i-3}, w_{i-2}, w_{i-1}))$$

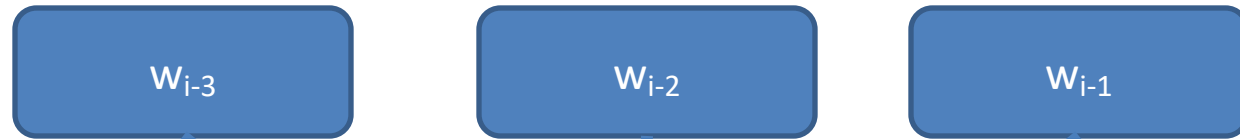
can we learn word-specific weights (by type)?

predict the next word



Neural Language Models

given some context...



can we *learn* the feature function(s) for *just* the context?

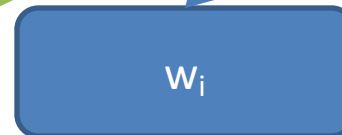
compute beliefs about what is likely...



$$p(w_i | w_{i-3}, w_{i-2}, w_{i-1}) = \text{softmax}(\theta_{w_i} \cdot f(w_{i-3}, w_{i-2}, w_{i-1}))$$

can we learn word-specific weights (by type)?

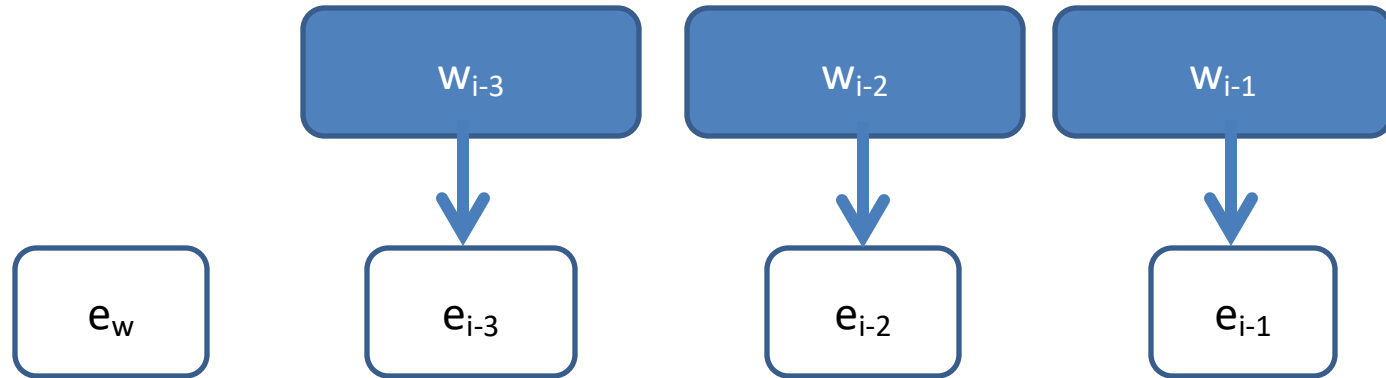
predict the next word



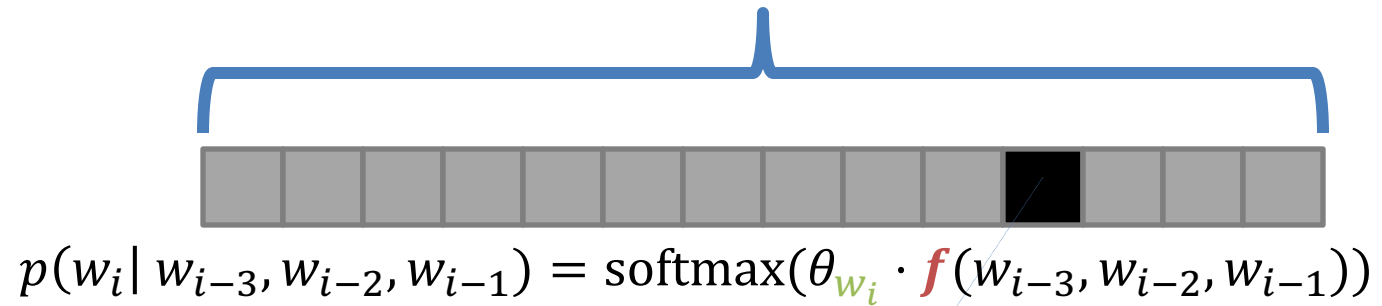
Neural Language Models

given some context...

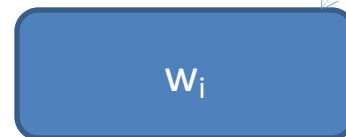
*create/use
“distributed
representations” ...*



*compute beliefs about
what is likely...*



predict the next word

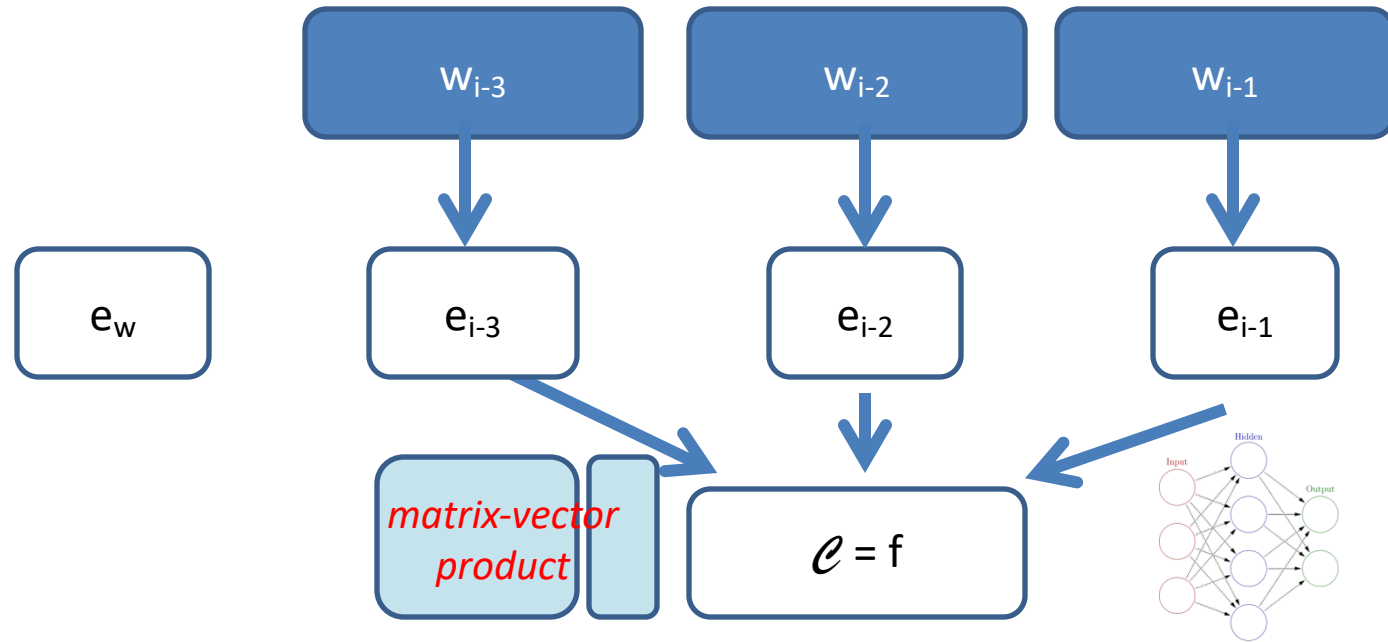


Neural Language Models

given some context...

create/use
"distributed
representations" ...

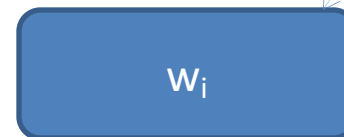
combine these
representations...



compute beliefs about
what is likely...

$$p(w_i | w_{i-3}, w_{i-2}, w_{i-1}) = \text{softmax}(\theta_{w_i} \cdot f(w_{i-3}, w_{i-2}, w_{i-1}))$$

predict the next word

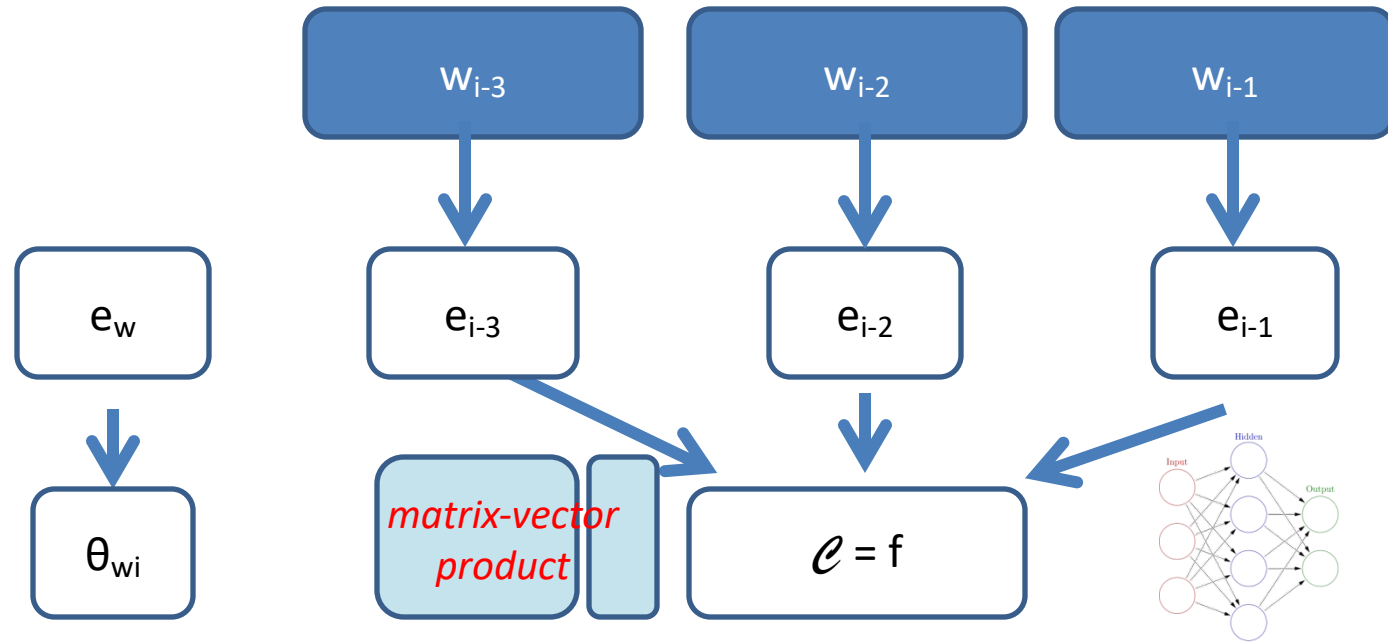


Neural Language Models

given some context...

create/use
“distributed
representations” ...

combine these
representations...



compute beliefs about
what is likely...

$$p(w_i | w_{i-3}, w_{i-2}, w_{i-1}) = \text{softmax}(\theta_{w_i} \cdot f(w_{i-3}, w_{i-2}, w_{i-1}))$$

predict the next word



Neural Language Models

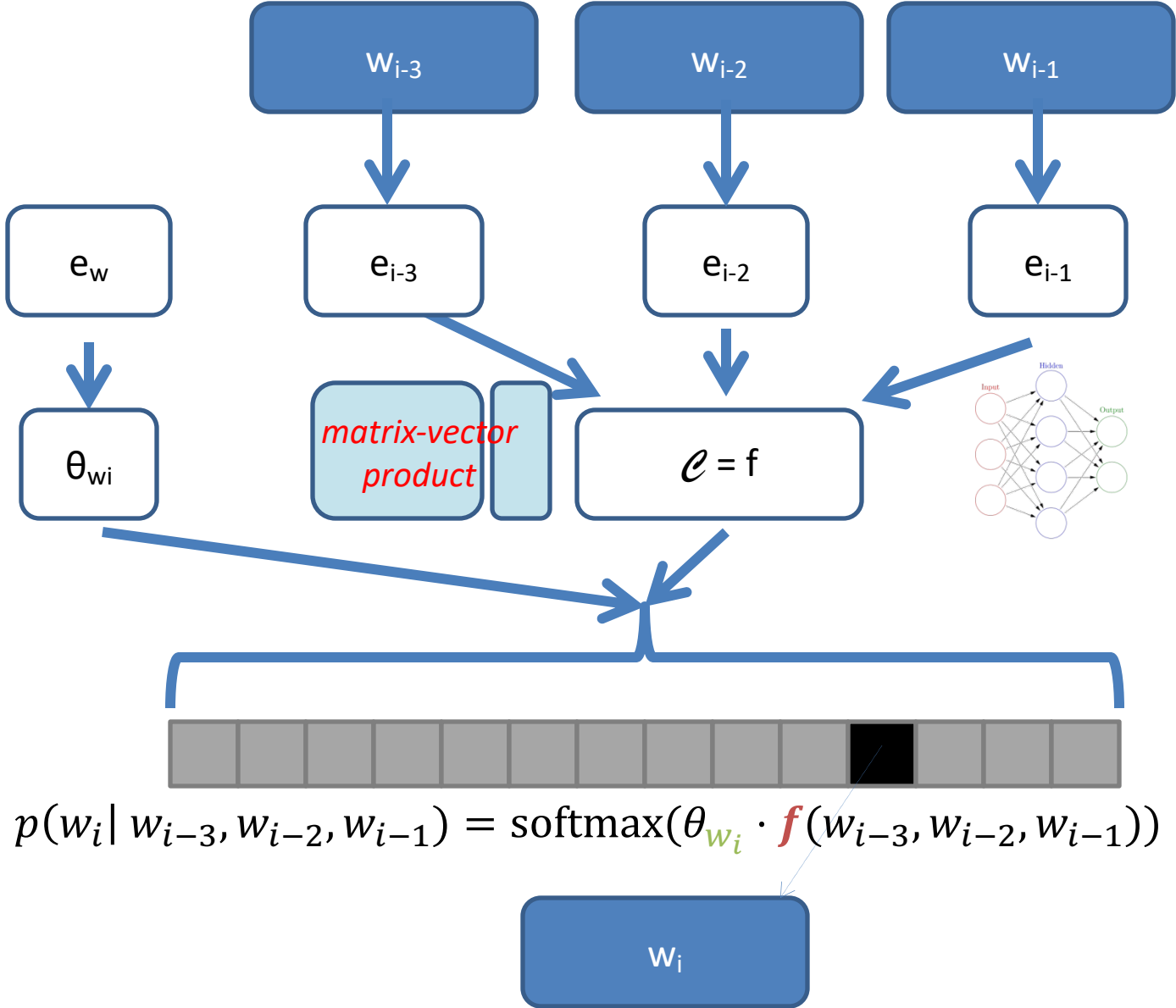
given some context...

create/use
"distributed
representations" ...

combine these
representations...

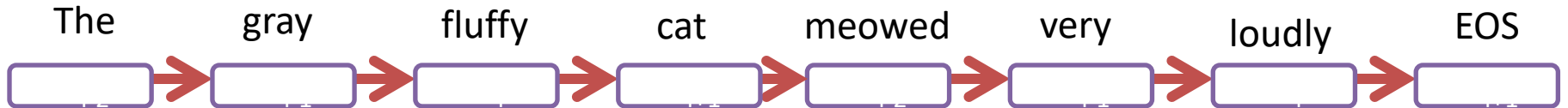
compute beliefs about
what is likely...

predict the next word



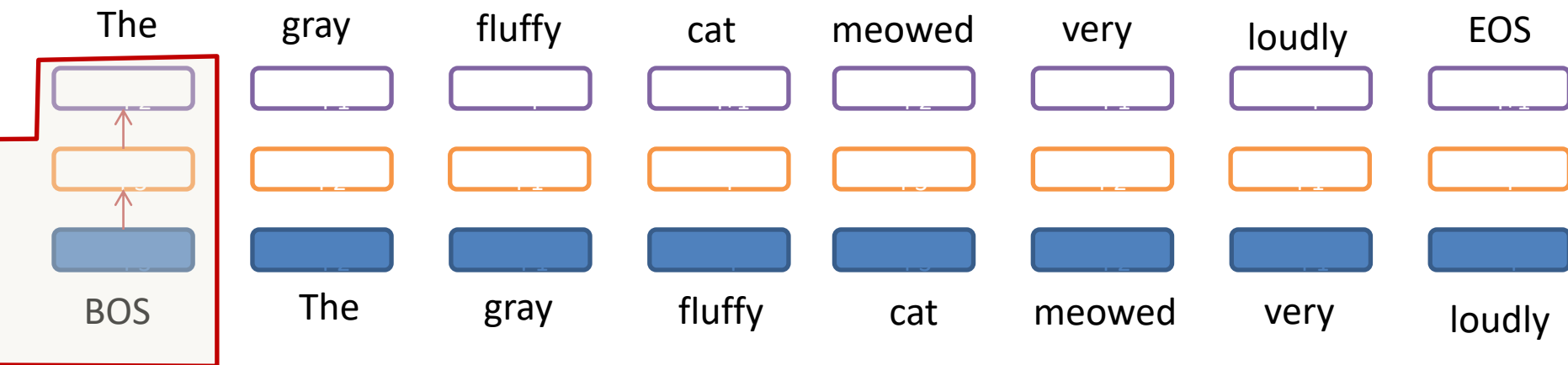
A Neural N-Gram Model

The gray fluffy cat meowed very loudly



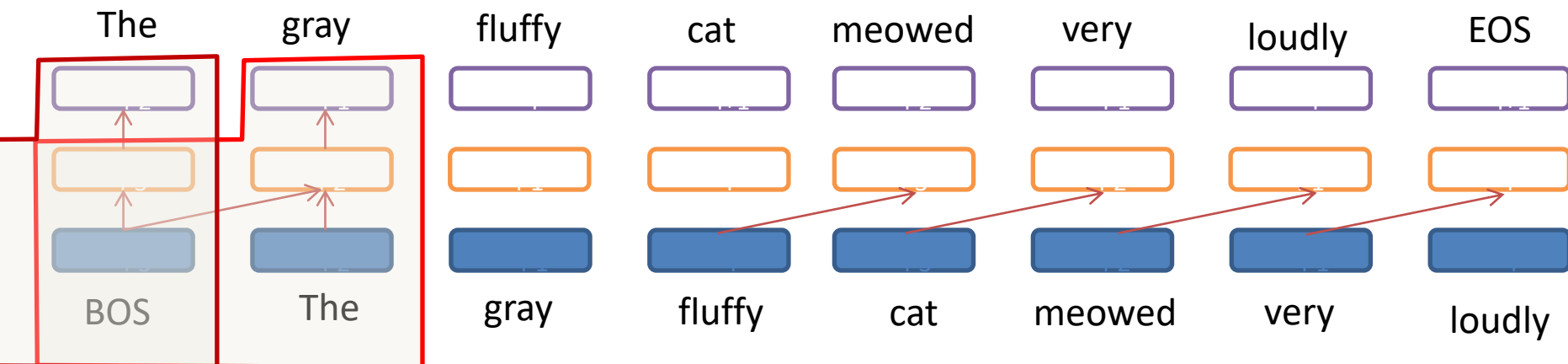
A Neural N-Gram Model (N=3)

The gray fluffy cat meowed very loudly



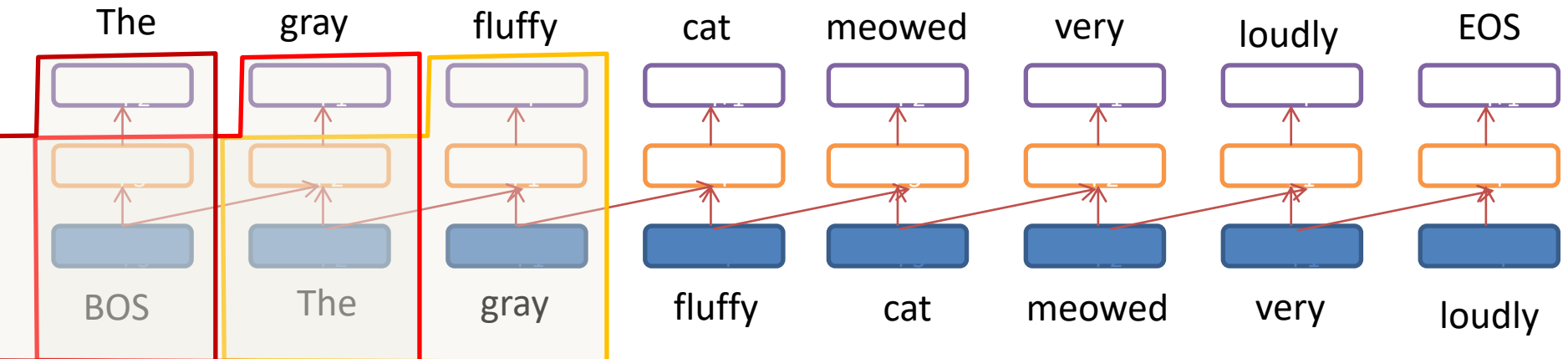
A Neural N-Gram Model (N=3)

The gray fluffy cat meowed very loudly



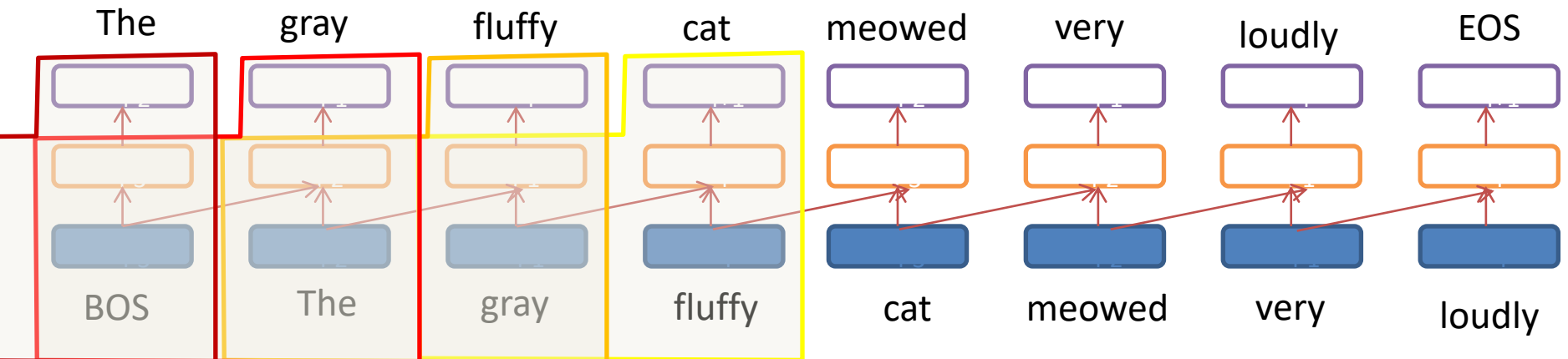
A Neural N-Gram Model (N=3)

The gray fluffy cat meowed very loudly



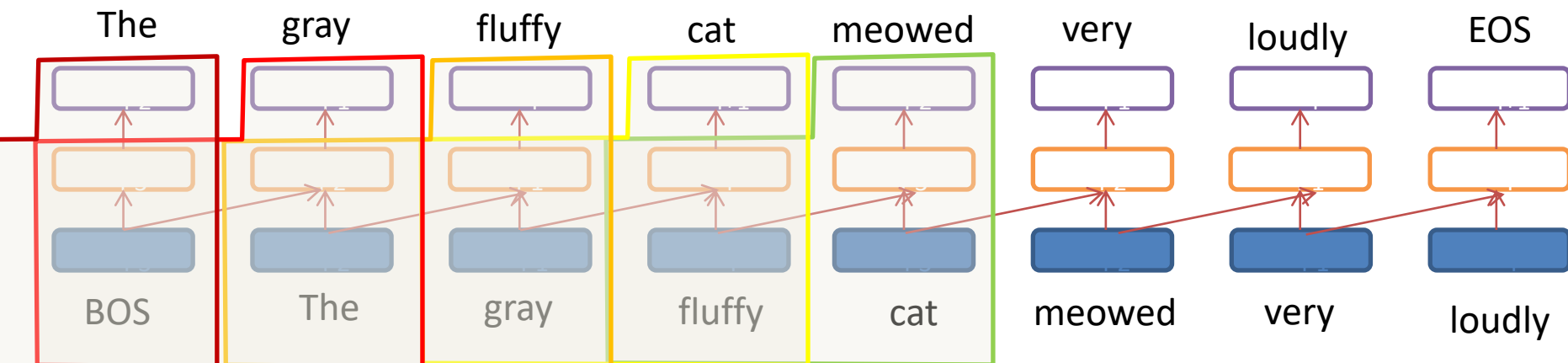
A Neural N-Gram Model (N=3)

The gray fluffy cat meowed very loudly



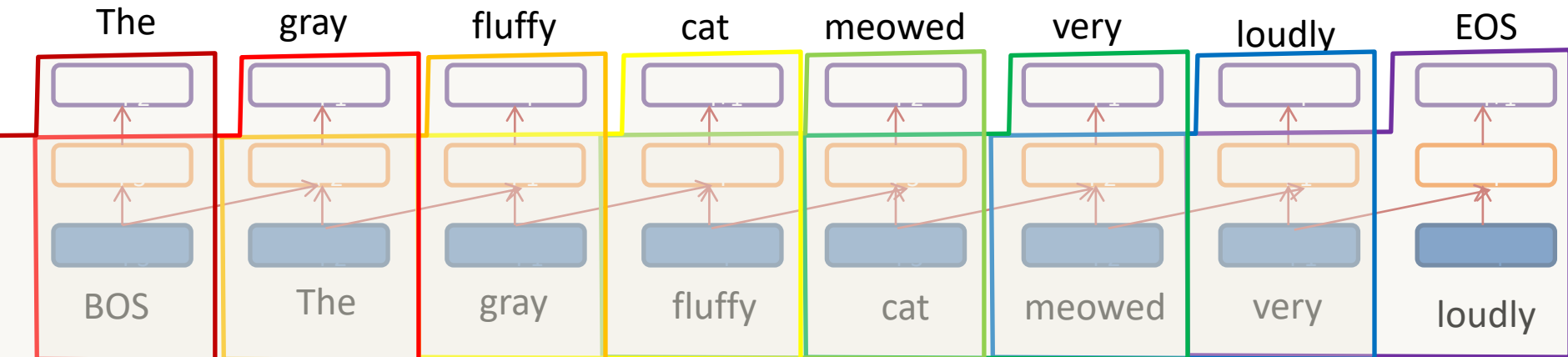
A Neural N-Gram Model (N=3)

The gray fluffy cat meowed very loudly



A Neural N-Gram Model (N=3)

The gray fluffy cat meowed very loudly



“A Neural Probabilistic Language Model,” Bengio et al. (2003)

Baselines

LM Name	N-gram	Params.	Test Ppl.
Interpolation	3	---	336
Kneser-Ney backoff	3	---	323
Kneser-Ney backoff	5	---	321
Class-based backoff	3	500 classes	312
Class-based backoff	5	500 classes	312

“A Neural Probabilistic Language Model,” Bengio et al. (2003)

Baselines

LM Name	N-gram	Params.	Test Ppl.
Interpolation	3	---	336
Kneser-Ney backoff	3	---	323
Kneser-Ney backoff	5	---	321
Class-based backoff	3	500 classes	312
Class-based backoff	5	500 classes	312

NPLM

N-gram	Word Vector Dim.	Hidden Dim.	Mix with non-neural LM	Ppl.
5	60	50	No	268
5	60	50	Yes	257
5	30	100	No	276
5	30	100	Yes	252

“A Neural Probabilistic Language Model,” Bengio et al. (2003)

Baselines

LM Name	N-gram	Params.	Test Ppl.
Interpolation	3	---	336
Kneser-Ney backoff	3	---	323
Kneser-Ney backoff	5	---	321
Class-based backoff	3	500 classes	312
Class-based backoff	5	500 classes	312

NPLM

N-gram	Word Vector Dim.	Hidden Dim.	Mix with non-neural LM	Ppl.
5	60	50	No	268
5	60	50	Yes	257
5	30	100	No	276
5	30	100	Yes	252

“we were not able to see signs of over-fitting (on the validation set), possibly because we ran only 5 epochs (over 3 weeks using 40 CPUs)” (Sect. 4.2)

A Closer Look at Neural $p(\text{Won't you please donate?} \mid \text{Primary})$

This is a *class-based* language model, but incorporate the label into the *embedding representation*

To learn $p(\text{Won't you please donate?} \mid \text{Class})$:

Define an embedding method that makes use of the specific label **Class**

Unlike count-based models, you don't *need* "separate" models here

Language Models & Smoothing

~~• Maximum likelihood (MLE): simple counting~~

~~• Other count based models~~

~~— Laplace smoothing, add λ~~

~~— Interpolation models~~

~~— Discounted backoff~~

~~— Interpolated (modified) Kneser-Ney~~

~~— Good Turing~~

~~— Witten-Bell~~

~~• Maxent n-gram models~~

~~• Neural n-gram models~~

• Recurrent/autoregressive NNs

Easy to implement

Advanced/
out of scope

Featureful LMs

Feedforward LMs

Super modern

Recurrent/Autoregressive LMs

- coming next class...