### Foundations for Natural Language Processing Evaluation for generation, tokenization

Ivan Titov (with graphics/materials from Elena Voita)



## Plan for today

Last time:

- Talked about generation and discussed inference algorithms
- Introduce vanilla encoder-decoder algorithms

Today, we will

- discuss how to evaluate text generation systems
- discuss tokenization and specifically subword tokenization

### Recap: Generating text

To generate text using a language model, you could just *sample* tokens from the probability distribution predicted by a model



## Recap: Sequence-to-Sequence modeling



### Recap: Encoder-decoder framework



#### Recap: simplest RNN-based Model:



### Inference (aka decoding)

$$y' = \arg \max_{y} p(y|x) = \arg \max_{y} \prod_{t=1}^{n} p(y_t|y_{< t}, x)$$
 How to find the argmax?

The simplest idea – greedy decoding, at each step, pick the most likely token, but note:

$$\arg \max_{y} \prod_{t=1}^{n} p(y_t | y_{< t}, x) \neq \prod_{t=1}^{n} \arg \max_{y_t} p(y_t | y_{< t}, x)$$

We can also do sampling (e.g., nucleus), but this is not argmax, how can approximate the <u>global</u> argmax better: beam-search

#### **Beam search**

Maintaining top hypotheses as you go



#### Start with the begin of sentence token or with an empty sequence

(video, not visible in pdf)

#### Beam search

#### Maintaining top hypotheses as you go



All hypotheses are complete - generation ended

## Why not sampling?

Actually, we can also sample in machine translation too (as with language modelling)

The risk is that a sample translation can deviate from the source sentence in meaning (i.e. hallucinate)

As discussed last time, sampling may be preferable to beam-search for other seq2seq task, where there is more freedom in the choice of generations (e.g., generate a story given keywords) Evaluating text generation models

#### How to evaluate text generation?

Consider French to English machine translation

Source sentence: Le chat est assis sur le tapis

Human translation into English: The cat is on the carpet

#### How to evaluate text generation?

Consider French to English machine translation

Source sentence: Le chat est assis sur le tapis

Human translation into English: The cat is on the carpet

How can we design a metric which would score MTI > MT2?

MTI: The cat is seated on the mat

MT2: The chat is assassinated on the tape

(We are looking into *automatic extrinsic evaluation*)

Idea: count overlapping ngrams - BLEU

Typically, we need more than I human (aka reference) translation per sentence to have reliable evaluation.

Let's focus on unigrams (individual tokens) for now

MT: The the the the the the a

Reference I: The cat is on the mat

Reference 2: There is a cat on the mat

(ignore capitalization for evaluation, i.e. treat 'The' and 'the' as the same word)

Idea: count overlapping ngrams - BLEU

MT: <u>The the the the the the the</u> a Reference I: <u>The</u> cat is on <u>the</u> mat Reference 2: There is a cat on the mat

'the' appears 7 times

'the' appears 2 times

'the' appears 1 time

Modified unigram precision: 2 / 7

Idea: count overlapping ngrams - BLEU

MT: The the the the the the the  $\underline{a}$  'a' appears 1 time Reference 1: The cat is on the mat 'a' appears 0 times Reference 2: There is  $\underline{a}$  cat on the mat 'a' appears 1 time Modified unigram precision: (2 + 1) / (7 + 1) = 3 / 8Aggregate over all unigrams in the MT ('candidate')

### **BLEU** metric

Actual BLEU is considerably more complicated, as needs to

- aggregate over the entire test set
- aggregate over ngrams of different order (unigrams, bigrams, ...)
- penalize short translation (remember from parsing: *precision* favors models producing short outputs)

There are other ngram overlap metrics which can be more suitable for other text generation problems (e.g., ROUGE for summarization)

Ngram overlap metrics - weaknesses

- do not account for lexical paraphrases (e.g., substituting words with their synonyms)
- even more problematic for long text generation (e.g., document machine translation)
- unreliable for tasks with less restricted outputs (e.g., generate "a scary novel about Edinburgh")
- do not sufficiently penalize hallucinations

#### What can we do if they are so unreliable?

••

- Human evaluation (expensive, hard to relate to results of older experiments)
- Neural model-based metrics (e.g., BERTScore, GPTScore)
- Specialized metrics (e.g., FActScore for hallucinations)



We considered tokenization of sentences into 'words' (whatever we mean by a 'word')

Word-level

- fixed vocabulary
- can process only a fixed number of words

Subword-level

- open vocabulary
- rare and unknown tokens are encoded as sequences of subword units

Instead of 'unrelated', we get two tokens 'un@@' 'related'

Tokenization

why subword tokenization?

Word-level

- fixed vocabulary
- can process only a fixed number of words

Subword-level

- open vocabulary
- rare and unknown tokens are encoded as sequences of subword units

Instead of 'unrelated', we get two tokens 'un@@' 'related'

Tokenization

why subword tokenization?

- reduces sparsity
- memory requirement
- results in a speeds-up (recall: softmax involves summation over all token types, few token typs -> faster computation)
- <u>may</u> provide an appropriate bias for compositionality, e.g., tokens for "unrelated" and "related" <u>may</u> be shared (e.g., include related)

Word-level

- fixed vocabulary
- can process only a fixed number of words

Subword-level

- open vocabulary
- rare and unknown tokens are encoded as sequences of subword units

Instead of 'unrelated', we get two tokens 'un@@' 'related'

Tokenization

#### Especially crucial for morphologically-rich languages

Standard segmentation algorithms rely on character ngram frequency, not on morphology (e.g., Byte-Pair Encoding)

Used in virtually any modern neural model

## Byte pair encoding algorithm

Byte Pair Encoding (BPE) algorithm (Gage, 1994) was applied to subword segmentation in machine translation by Sennrich, Haddow, and Birch (ACL 2015)

BPE, two phases

- Training: learn "BPE rules", i.e., which pairs of symbols to merge;
- Inference: apply learned rules to segment a text.

Start with individual characters as tokens

Repeat:

- count pairs of tokens: how many times each pair occurs together in the training data;
- find the most frequent pair of tokens;
- merge this pair add a merge to the merge table, and the new token to the vocabulary.

In practice, the algorithm first counts how many times each word appeared in the data.

=> Using this information, one can count pairs of adjacent tokens more easily.

Note the tokens do not cross word boundary - everything happens within words.

Initial vocabulary:	word	count	Current merge table:
characters	cat	4	(empty)
ļ	mat	5	(empty)
Split each word	mats	2	
into characters	mate	3	
	ate	3	
	eat	2	

Count pairs of tokens:	word	count	Current merge table:
<b>a + †</b> : 20	cat	4	(empty)
<b>m + a</b> : 10	m a t	5	
<b>† + e</b> : 7	mats	2	
<b>c + a</b> : 4	mate	3	
<b>† + s</b> : 2	ate	3	
<b>e + a</b> : 2	e a t	2	

Words in the data:



Current merge table:

a t -> at

	word	count	Current merge table:
	c a t	4	a t -> at
	m a t	5	
Apply the merge	mats	2	
to the data	m a t e	3	
	ate	3	
	e a t	2	



Count pairs of new tokens:	word	count	Current merge table:
<b>m + at</b> : 10	c at	4	a t -> at
<b>at + e</b> : 7	m at	5	
<b>c + at</b> : 4	m at s	2	
<b>at + s</b> : 2	m at e	3	
<b>e + at</b> : 2	at e	3	
	e at	2	

Words in the data:

Count pairs of new tokens:	word	count
<b>m + at</b> : 10	c at	4
at + e: 7	m at	5
c + at: 4 Add the most	m at s	2
at + s: 2 frequent pair to	m at e	3
<b>e + at</b> : 2	at e	3
	e at	2

Current merge table:
a t -> at
m at -> mat

	word	count	Current merge table:
	c at	4	a t -> at
	m at	5	m at -> mat
Apply the merge	m at s	2	
to the data	m at e	3	
	at e	3	
	e at	2	

	word	count	Current merge table:
	c at	4	a † -> at
	mat	5	m at -> mat
Apply the merge	mat s	2	
to the data	mat e	3	
	at e	3	
	e at	2	

c + at: 4  c at  4  a + -> at    at + e: 4  mat  5  m at -> mat    mat + e: 3  mat s  2	Count pairs of new tokens:	word	count	Current merge table:
at + e: 4  mat  5  m at -> mat    mat + e: 3  mat s  2    mat + s: 2  mat e  3    e + at: 2  at e  3    e at  2  2	<b>c + at</b> : 4	c at	4	a † -> at
mat + e: 3  mat s  2    mat + s: 2  mat e  3    e + at: 2  at e  3    e at  2	<b>at + e</b> : 4	mat	5	m at -> mat
mat + s: 2  mat e  3    e + at: 2  at e  3    e at  2	<b>mat + e</b> : 3	mat s	2	
e + at: 2 at e 3 e at 2	<b>mat + s</b> : 2	mat e	3	
e at 2	<b>e + at</b> : 2	at e	3	
		e at	2	

Count pairs of new tokens:	word	count
<b>c + at</b> : 4	c at	4
at + e: 4	mat	5
mat + e: 3 Add the most	mat s	2
mat + s: 2 frequent pair to	mat e	3
<b>e + at</b> : 2	at e	3
	e at	2

Current merge table:
a † -> at
m at -> mat
c at -> cat
--------------------------------
Apply the merge to the data

	word	count	Current merge table:
	cat	4	a t -> at
Apply the merge to the data	mat	5	m at -> mat
	mat s	2	c at -> cat
	mat e	3	
	at e	3	
	e at	2	

Count pairs of new tokens:	word	count	Current merge table:
<b>at + e</b> : 4	cat	4	a † -> at
<b>mat + e</b> : 3	mat	5	m at -> mat
<b>mat + s</b> : 2	mat s	2	c at -> cat
<b>e + at</b> : 2	mat e	3	
	at e	3	
	e at	2	





Count pairs of new tokens:	word	count	Current merge table:
<b>mat + e</b> : 3	cat	4	a † -> at
<b>mat + s</b> : 2	mat	5	m at -> mat
<b>e + at</b> : 2	mat s	2	c at -> cat
	mat e	3	at e -> ate
	ate	3	
	e at	2	

Count pairs of new tokens:	word	count	Current merge table:
<b>mat + e</b> : 3	cat	4	a t -> at
<b>mat + s</b> : 2	mat	5	m at -> mat
e + at: 2 Add the most	mat s	2	c at -> cat
frequent pair to the merge table	mat e	3	at e -> ate
	ate	3	mat e -> mate
	e at	2	

	word	count	Current merge table:
Apply the merge to the data	cat	4	a † -> at
	mat	5	m at -> mat
	mat s	2	c at -> cat
	► mat e	3	at e -> ate
	ate	3	mat e -> mate
	e at	2	

	word	count	Current merge table:
Apply the merge to the data	cat	4	a † -> at
	mat	5	m at -> mat
	mat s	2	c at -> cat
	► mate	3	at e -> ate
	ate	3	mat e -> mate
	e at	2	

	word	count	Current merge table:
Reached maximum	cat	4	a t -> at
number of merges	mat	5	m at -> mat
	mat s	2	c at -> cat
stop	mate	3	at e -> ate
	ate	3	mat e -> mate
	e at	2	



After training, we are left with a vocabulary and a merge tables (ordered set of merges)

The algorithm starts with segmenting a word into a sequence of characters. After that, it iteratively makes the following two steps until no merge it possible:

- among all possible merges at this step, find the highest merge in the table (highest -> more frequent);
- apply this merge.

hyphens are possible merges

u-n-r-e-l-a-t-e-d

hyphens are possible merges u-n-r-e-l-a-t-e-d the highest merge: merge with the highest priority (the highest in the merge table)

**u-n-r-e-l-a-t-e-d** ↓ merge

u-n-<mark>r-e</mark>-l-a-t-e-d u-n <u>re</u>-l-a-t-e-d

u-n-r-e-l-a-t-e-d u-n <u>re</u>-l-a-t-e-d

u-n-r-e-l-a-t-e-d u-n <u>re</u>-l-a-t-e-d

pick the highest merge

u-n-<mark>r-e</mark>-l-a-t-e-d u-n <u>re</u>-l-a-t-e-d

pick the highest merge

u-n-r-e-l-a-t-e-d u-n <u>re</u>-l-a-t-e-d ↓ merge

u-n-<mark>r-e</mark>-l-a-t-e-d u-n <u>re</u>-l-<u>a-</u>t-e-d u-n re-l-<u>at</u>-e-d

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-</u>e-d u-n re-l-<u>at</u>-e-d pick the highest merge

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-e-d</u> u-n re-l-<u>at</u>-e-d pick the highest merge

u-n-r-e-l-a-t-e-d u-n <u>re</u>-l-a-t-e-d u-n re-l-<u>at</u>-e-d ↓ merge

u-n-r-e-l-a-t-e-d u-n <u>re</u>-l-a-t-e-d u-n re-l-<u>at</u>-e-d ↓ merge

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-e-d</u> u-n re-l-<u>at-e-d</u> u-n re-l-at-<u>ed</u> pick the highest merge

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-e-d</u> u-n re-l-<u>at-e-d</u> u-n re-l-at-<u>ed</u> pick the highest merge

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-</u>e-d u-n re-l-<u>at</u>-e-d u-n re-l-at-<u>ed</u> <u>werge</u>

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-</u>e-d u-n re-l-<u>at-</u>e-d u-n re-l-at-<u>ed</u> <u>un</u> re-l-at-ed

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-</u>e-d u-n re-l-<u>at-</u>e-d <u>u-n re-l-at-ed</u> <u>un re-l-at-ed</u> pick the highest merge

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-e-d</u> u-n re-l-<u>at-e-d</u> u-n re-l-at-<u>ed</u> <u>un re-l-at-ed</u> pick the highest merge

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-</u>e-d u-n re-l-<u>at-</u>e-d <u>u-n re-l-at-ed</u> <u>un re-l-at-ed</u> <u>merge</u>

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-e-d</u> u-n re-l-<u>at-e-d</u> <u>u-n re-l-at-ed</u> <u>un re-l-at-ed</u> un re-l-<u>at-ed</u>

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-e-d</u> u-n re-l-<u>at-e-d</u> <u>u-n re-l-at-ed</u> <u>un re-l-at-ed</u> un re-l-<u>at-ed</u>

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-</u>e-d u-n re-l-<u>at-</u>e-d <u>u-n re-l-at-ed</u> <u>un re-l-at-ed</u> un <u>re-l-at-ed</u> un <u>re-l-ated</u> pick the highest merge

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-e-d</u> u-n re-l-<u>at-e-d</u> <u>u-n re-l-at-ed</u> <u>un re-l-at-ed</u> un <u>re-l-at-ed</u>
u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-e-d</u> u-n re-l-<u>at-e-d</u> <u>u-n re-l-at-ed</u> <u>un re-l-at-ed</u> un <u>re-l-at-ed</u>

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-e-d</u> u-n re-l-<u>at-e-d</u> u-n re-l-at-<u>ed</u> <u>un re-l-at-ed</u> un <u>re-l-ated</u> un <u>rel</u>-ated

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-</u>e-d u-n re-l-<u>at-</u>e-d u-n re-l-at-<u>ed</u> un re-l-at-ed un <u>re-l-ated</u> un <u>rel-ated</u> <u>pick the</u> highest merge

u-n-r-e-l-a-t-e-d u-n re-l-at-e-d u-n re-l-at-e-d u-n re-l-at-ed un re-l-at-ed un re-l-ated un rel-ated bighest merge

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-e-d</u> u-n re-l-<u>at-e-d</u> u-n re-l-at-<u>ed</u> un re-l-at-ed un re-l-<u>ated</u> un <u>rel-ated</u>

u-n-r-e-l-a-t-e-d u-n <u>re-l-a-t-e-d</u> u-n re-l-<u>at-e-d</u> u-n re-l-at-<u>ed</u> <u>un re-l-at-ed</u> un <u>re-l-ated</u> un <u>rel-ated</u> un <u>rel-ated</u> un <u>related</u>

u-n-r-e-l-a-t-e-d u-n re-l-at-e-d u-n re-l-at-e-d u-n re-l-at-ed un re-l-at-ed un re-l-ated un rel-ated un rel-ated un related No merges un@@ related

## BPE – drawbacks / discussion

- The induced segmentation does not align with linguistic morphological boundaries.
- Frequent words that share the same morphemes (e.g., roots) are segmented differently from less frequent words.
- The segmentation does not directly correspond to the information content of a token or the computational effort required in a neural model.

#### What can we do about it?

# Summary

- Done with inference algorithms (greedy, beam-search, sampling, temperature,...)
- Evaluating text generation (e.g., BLEU)
- Subword tokenization