

# Natural Language Understanding, Generation, and Machine Translation

## Lecture 6: Modelling Data and Words

---

Alexandra Birch

24 January 2025 (week 2)

School of Informatics  
University of Edinburgh  
a.birch@ed.ac.uk

Based on slides by Rico Sennrich

# Refresher

---

# Input Representation

## how do we represent input?

- 1-hot encoding
  - lookup of word embedding for input
  - probability distribution over vocabulary for output
- large vocabularies
  - increase network size
  - decrease training and decoding speed
- typical network vocabulary size: 10 000–100 000 symbols

vocabulary		representation of "cat"	
		1-hot vector	embedding
0	the	0	0.1
1	cat	1	0.3
2	is	0	0.7
.	.	.	0.5
1024	mat	0	

NLU and NLG are open-vocabulary problems

- many training corpora contain millions of word types
- productive word formation processes (compounding; derivation) allow formation and understanding of unseen words
- names, numbers are morphologically simple, but open word classes

Research: we can download a clean corpus eg. Hansard, Europarl, Penn Treebank

Real life: nothing like this

- What if we need data that is not available publically? medical, financial, conversational
- What if it is for a low-resource language and none available?
- What if data we have is very noisy?
- What if there are very long dependencies over many sentences?

# Modelling Data

---

# Language Identification

Many datasets are crawls from the internet: CommonCrawl (250B pages) or the Internet Archive (850B pages)

- First step is Language Identification: LID
- What languages are these?
  - ndiyahamba (I'm going)
  - This is lekker bru (This is lovely)
- Generally solved task for high-resource languages
- Challenge with low resource, closely related language, code-switched or noisy data
- Solution: Fast lightweight classifiers
- Fasttext - Pre-trained models for 157 different languages

Real data comes unsegmented into sentences.

- Sentence Segmentation:  
! ? Mostly unambiguous but "." is very ambiguous (eg. the U.N.)
- Many scripts do not have end of sentence marker
- Speech is often not easy to segment into complete sentences
- Clean sentences - normally what models are trained on can struggle with shorter/longer sequences
- Solution: Language specific rules



What is a word?

- Lookup in dictionary - but morphology makes this harder
- Thing between spaces - what about language without spaces or Finnish?
- Punctuation? Contractions? “that’s” → “that” “s”
- Solution: For languages with spaces use spaces + punctuation + rules
- For Chinese etc. large dictionaries, punctuation + rules

Critical for input to neural network - what is the input?

What sequence?

- Document, sentence, window, turn or utterance in a conversation

Sequence of what?

- Words, tokenized words, word stems, morphemes

Very long sequences are harder to model.

Vocabulary size needs to be limited as it has a huge effect on model size and efficiency.

## **Modelling words - open vocabulary models**

---

## Non-Solution: Ignore Rare Words

- replace out-of-vocabulary words with UNK
- a vocabulary of 50 000 words covers 95% of text
- this gets you 95% of the way...
  - ... if you only care about automatic metrics

## Non-Solution: Ignore Rare Words

- replace out-of-vocabulary words with UNK
- a vocabulary of 50 000 words covers 95% of text
- this gets you 95% of the way...  
... if you only care about automatic metrics

### why 95% is not enough

rare outcomes have high self-information

source      Mr **Gallagher** has offered a ray of hope.

reference    Herr **Gallagher** hat einen hoffnungsstrahl ausgesandt .

## Solution 1: Back-off Models

## back-off models [Jean et al., 2015, Luong et al., 2015]

- replace rare words with UNK at training time
- when system produces UNK, align UNK to source word, and translate this with back-off method

source	Das <b>Raumklima</b> ist sehr angenehm.	
reference	The <b>indoor temperature</b> is very pleasant.	
[Bahdanau et al., 2015]	The <b>UNK</b> is very nice.	✗
[Jean et al., 2015]	The <b>temperature</b> is very nice.	✗

## limitations

- compounds: hard to model 1-to-many relationships
- morphology: hard to predict inflection with back-off dictionary
- names: if alphabets differ, we need transliteration
- alignment: attention model unreliable

# Subwords for NMT: Motivation

## Subwords units could be meaningful useful for translation

- compounding and other productive morphological processes
  - they charge a carry-on bag fee.
  - sie erheben eine Hand|gepäck|gebühr.
- names
  - Edinburgh(English)
  - Edimburgo(Spanish)
- Morphological variation: slightly exaggerated eg. Turkish
  - OSMANLILAŞTIRAMAYABİLECEKLERİİMİZDENMİŞSİNİZ
  - OSMAN-LI-LAŞ-TIR-AMA-YABİL-ECEK-LER-İİMİZ-DEN-MİS-ŞİNİZ
- technical terms, numbers, etc.:
  - 10-12-2020.
  - December 10 2020.

## segmentation algorithms: wishlist

- **open-vocabulary NMT**: encode *all* words through small vocabulary
- encoding generalizes to unseen words
- small text size
- good translation quality

## our experiments [Sennrich et al., 2016]

- after preliminary experiments, we propose:
  - character n-grams (with shortlist of unsegmented words)
  - segmentation via *byte pair encoding* (BPE)



# Byte pair encoding for word segmentation

## bottom-up character merging

- starting point: character-level representation  
→ computationally expensive
- compress representation based on information theory  
→ byte pair encoding [Gage, 1994]
- repeatedly replace most frequent symbol pair ('A','B') with 'AB'
- hyperparameter: when to stop  
→ controls vocabulary size

word	freq	vocabulary: l o w e r n s t i d
'l o w'	5	
'l o w e r'	2	
'n e w e s t'	6	
'w i d e s t'	3	

# Byte pair encoding for word segmentation

## bottom-up character merging

- starting point: character-level representation  
→ computationally expensive
- compress representation based on information theory  
→ byte pair encoding [Gage, 1994]
- repeatedly replace most frequent symbol pair ('A','B') with 'AB'
- hyperparameter: when to stop  
→ controls vocabulary size

word	freq	
'l o w'	5	vocabulary: l o w e r n s t i d <b>e s</b>
'l o w e r'	2	
'n e w <b>e s</b> t'	6	
'w i d <b>e s</b> t '	3	

# Byte pair encoding for word segmentation

## bottom-up character merging

- starting point: character-level representation  
→ computationally expensive
- compress representation based on information theory  
→ byte pair encoding [Gage, 1994]
- repeatedly replace most frequent symbol pair ('A','B') with 'AB'
- hyperparameter: when to stop  
→ controls vocabulary size

word	freq	
'l o w'	5	vocabulary: l o w e r n s t i d e s e s t
'l o w e r'	2	
'n e w e s t'	6	
'w i d e s t '	3	

# Byte pair encoding for word segmentation

## bottom-up character merging

- starting point: character-level representation  
→ computationally expensive
- compress representation based on information theory  
→ byte pair encoding [Gage, 1994]
- repeatedly replace most frequent symbol pair ('A','B') with 'AB'
- hyperparameter: when to stop  
→ controls vocabulary size

word	freq	
'lo w'	5	vocabulary: l o w e r n s t i d e s e s t l o
'lo w e r'	2	
'n e w <b>est</b> '	6	
'w i d <b>est</b> '	3	

# Byte pair encoding for word segmentation

## why BPE?

- open-vocabulary:  
operations learned on training set can be applied to unknown words
- compression of frequent character sequences improves efficiency  
→ trade-off between text length and vocabulary size

'l o w e s t'

e s → es

es t → est

l o → lo

# Byte pair encoding for word segmentation

## why BPE?

- open-vocabulary:  
operations learned on training set can be applied to unknown words
- compression of frequent character sequences improves efficiency  
→ trade-off between text length and vocabulary size

'l o w **es** t'

<b>e s</b>	→	<b>es</b>
es t	→	est
l o	→	lo

# Byte pair encoding for word segmentation

## why BPE?

- open-vocabulary:  
operations learned on training set can be applied to unknown words
- compression of frequent character sequences improves efficiency  
→ trade-off between text length and vocabulary size

'l o w **est**'

e s → es

**es t** → **est**

l o → lo

# Byte pair encoding for word segmentation

## why BPE?

- open-vocabulary:  
operations learned on training set can be applied to unknown words
- compression of frequent character sequences improves efficiency  
→ trade-off between text length and vocabulary size

'lo w est'

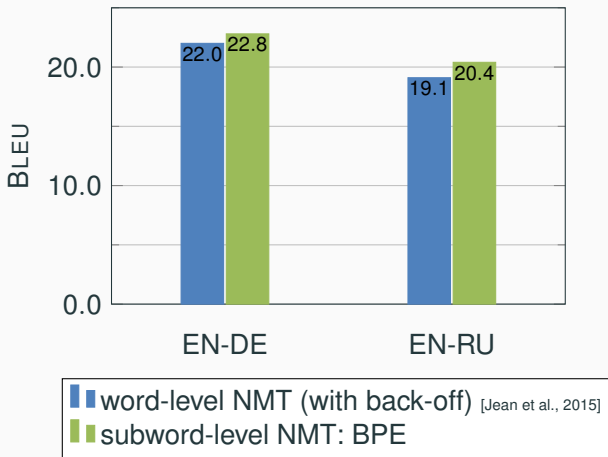
e s → es

es t → est

**l o** → **lo**



## Subword NMT: Translation Quality



# Subword Models: BPE-Dropout

u-n-r-e-l-a-t-e-d  
u-n re-l-a-t-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-related  
unrelated

(a)

u-n-r-e-l-a-t-e\_d  
u-n re-l-a-t-e\_d  
u-n re\_l-at-e\_d  
un re-l-at-e-d  
un re-l-at-ed  
un re-lat-ed  
un relat\_ed

u-n-r-e-l-a-t-e-d  
u\_n re\_l-a-t-e-d  
u\_n re-l-at-e-d  
u\_n re-l-ate\_d  
u\_n rel-ate-d  
u\_n relate\_d

u-n\_r\_e-l-a-t-e-d  
u-n-r\_e-l-at-e-d  
u-n-r\_e-l\_at\_ed  
un-r-e-l-at-ed  
un re-l\_at-ed  
un re-l-ated  
un rel\_ated

(b)

BPE

BPE dropout

From [Provilkov et al., 2020]

- Hyphen - possible merge
- merges performed - in green
- merges dropped - in red

# Subword Models: BPE-Dropout

- BPE-Dropout: Simple and effective Subword Regularizations

[Provilkov et al., 2020]

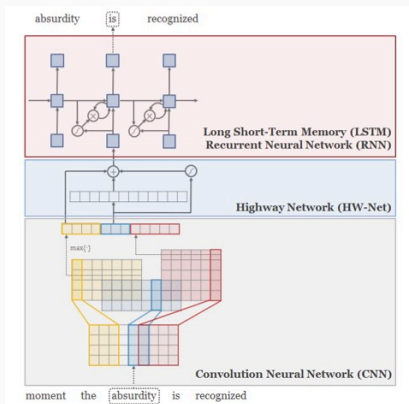
- Adding stochastic noise to increase model robustness
- BPE: most frequent words are intact in vocabulary, learns how to compose with infrequent words
- If we sometimes forget to merge, we will learn how words compose, and better transliteration
- forget 1 in 10 times for most scripts, 6/10 in CKJ scripts
- Consistently give 1+ BLEU scores across language pairs - widely used

# Character-level Models

- advantages:
  - (mostly) open-vocabulary
  - no heuristic or language-specific segmentation
  - neural network can conceivably learn from raw character sequences
- drawbacks:
  - increasing sequence length slows training/decoding (reported x2–x8 increase in training time)
- open questions
  - on which level should we represent meaning?
  - on which level should attention operate?

# Character Aware Neural Language Model [Kim et al., 2016]

- goal: vocabulary over character set
- Convolution over characters, highway network over words, and LSTM layers



# Character Aware Neural Language Model [Kim et al., 2016]

(Based on cosine similarity)

	In Vocabulary				
	<i>while</i>	<i>his</i>	<i>you</i>	<i>richard</i>	<i>trading</i>
Word Embedding	<i>although</i>	<i>your</i>	<i>conservatives</i>	<i>jonathan</i>	<i>advertised</i>
	<i>letting</i>	<i>her</i>	<i>we</i>	<i>robert</i>	<i>advertising</i>
	<i>though</i>	<i>my</i>	<i>guys</i>	<i>neil</i>	<i>turnover</i>
	<i>minute</i>	<i>their</i>	<i>i</i>	<i>nancy</i>	<i>turnover</i>
<b>Characters</b> (before highway)	<i>chile</i>	<i>this</i>	<i>your</i>	<i>hard</i>	<i>heading</i>
	<i>whole</i>	<i>hhs</i>	<i>young</i>	<i>rich</i>	<i>training</i>
	<i>meanwhile</i>	<i>is</i>	<i>four</i>	<i>richer</i>	<i>reading</i>
	<i>white</i>	<i>has</i>	<i>youth</i>	<i>richter</i>	<i>leading</i>
<b>Characters</b> (after highway)	<i>meanwhile</i>	<i>hhs</i>	<i>we</i>	<i>eduard</i>	<i>trade</i>
	<i>whole</i>	<i>this</i>	<i>your</i>	<i>gerard</i>	<i>training</i>
	<i>though</i>	<i>their</i>	<i>doug</i>	<i>edward</i>	<i>traded</i>
	<i>nevertheless</i>	<i>your</i>	<i>i</i>	<i>carl</i>	<i>trader</i>

# Beyond Character-level

- Massively multilingual settings character-level models can result in a very large vocabulary. eg. Unicode 1,112,064 codepoints
- Byte level:
  - better robustness to noise but longer training time ByT5: Towards a token-free future with pre-trained byte-to-byte models [Xue et al., 2021]
  - Claim: token free - but really use fixed Unicode tokenisation which is not linguistically motivated
  - Potentially unfair: Unicode characters beyond ASCII are much longer byte sequences - more expensive to model
- Pixel level:
  - similarities that human readers might pick up on eg. to generalise to rare Chinese characters
  - Makes translation significantly more robust to induced noise (including unicode errors) Robust Open-Vocabulary Translation from Visual Text

# Conclusion

- Understand how your data was preprocessed
- Important to model it correctly
- BPE and BPE-dropout is widely used
- There is no perfect method of handling tokenization.
- Opposing goals:
  - Decompose maximally for simple and robust processing
  - Desire to be computationally efficient in a way that is fair across languages
- Still not learning entities jointly with the rest of the model: separate preprocessing step
- How well these methods generalise from character strings to higher level of representation still to be fully studied





Bahdanau, D., Cho, K., and Bengio, Y. (2015).

**Neural Machine Translation by Jointly Learning to Align and Translate.**

In

Proceedings of the International Conference on Learning Representations



Gage, P. (1994).

**A New Algorithm for Data Compression.**

C Users J., 12(2):23–38.



Jean, S., Cho, K., Memisevic, R., and Bengio, Y. (2015).

**On Using Very Large Target Vocabulary for Neural Machine Translation.**

In

Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics, pages 1–10, Beijing, China. Association for Computational Linguistics.



Kim, Y., Jernite, Y., Sontag, D., and Rush, A. (2016).

**Character-aware neural language models.**

In Proceedings of the AAAI conference on artificial intelligence, volume 30.



Luong, T., Sutskever, I., Le, Q., Vinyals, O., and Zaremba, W. (2015).

**Addressing the Rare Word Problem in Neural Machine Translation.**

In

Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics, pages 11–19, Beijing, China. Association for Computational Linguistics.



Provilkov, I., Emelianenko, D., and Voita, E. (2020).

**Bpe-dropout: Simple and effective subword regularization.**

In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 1882–1892.



Salesky, E., Etter, D., and Post, M. (2021).

**Robust open-vocabulary translation from visual text representations.**

In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 7235–7252, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.



Sennrich, R., Haddow, B., and Birch, A. (2016).

**Neural Machine Translation of Rare Words with Subword Units.**

In

Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics,  
pages 1715–1725, Berlin, Germany.



Xue, L., Barua, A., Constant, N., Al-Rfou, R., Narang, S., Kale, M., Roberts, A., and Raffel, C. (2021).

**Byt5: Towards a token-free future with pre-trained byte-to-byte models.**