

CMPT 413/713: Natural Language Processing

Sequence to Sequence Models (Seq2Seq)

Adapted from slides from Dangi Chen and Karthik Narasimhan (with some content from slides from Abigail See, Graham Neubig)

SFU Nat Lang Lab

Spring 2024 2024-02-05

Overview

- Sequence generation tasks
- Seq2Seq models Encoder/Decoder
- Decoding strategies
- Evaluating text generation
- Attention

Sequence Generation



Understanding what is said (encoding, parsing, feature extraction)

η γάτα κάθισε στο τραπέζι

Deciding what to say (decoding, generating)

Encoder-Decoder Model



Seq2Seq Tasks and Applications

Task/Application

- Input Output Machine Translation French English
 - Summarization
 - Dialogue Utterance
 - Parsing Sentence
- **Question Answering** Context + Question

Document

- Short Summary
 - Response
 - Parse tree (as sequence)
 - Answer

Cross-Modal Seq2Seq

Task/Application

Speech Recognition Spee

Image Captioning

Video Captioning

Vision-Language Navigation

Input	Output
ech Signal	Transcript
Image	Text
Video	Text
Text	Actions

Cross-modal sequence generation

Video captioning (video frames to text)



Embodied AI (text + frames to actions)



Seq2Seq Tasks and Applications

Task/Application	Input	Output
Machine Translation	French	English
Summarization	Document	Short Summary
Dialogue	Utterance	Response
Parsing	Sentence	Parse tree (as sequence)
Question Answering Co	ntext + Question	Answer

Sequence to sequence models

Neural Machine Translation

- to target
- Architecture: Encoder-Decoder
 - Two main components:
 - matrix
 - language (output)

A single neural network is used to translate from source

Encoder: Convert source sentence (input) into a vector/

Decoder: Convert encoding into a sentence in target

Sequence to Sequence learning (Seq2seq)



- Encode entire input sequence into a single vector (using an RNN)
- Decode one word at a time (again, using an RNN!)
- Beam search for better inference
- Learning is not trivial! (vanishing/exploding gradients)

(Sutskever et al., 2014)

Sentence: This cat is cute



This

cat

Encoder



Sentence: This cat is cute



Encoder



cute

Sentence: This cat is cute



Encoder





Encoder









• A conditioned language model





Seq2seq training

- Similar to training a language model!
- Minimize cross-entropy loss:



- Need a really big corpus





Russian: Машинный перевод - это круто!

English: Machine translation is cool!

Seq2seq training



Efficient Training: Batching

- Apply RNNs to batches of sequences
- Present data as 3D tensor of $(T \times B \times F)$

Padded sequences

1	1	1	1
1	0	0	C
1	1	1	1
1	1	1	O

• Use mask matrix to aid with computations that ignore padded zeros

0 0 0 0 1 1 0 0

Lengths



Batching

- Sorting (partially) can help to create more efficient mini-batches
- However, the input is less randomized Unsorted





Decoding strategies

- How can we use our model (decoder) to generate sentences?

Generation

• Sampling: Try to generate a *random* sentence according the the probability distribution

• Argmax: Try to generate the *best* sentence, the sentence with the *highest* probability

Decoding Strategies

- Ancestral sampling
- Greedy decoding
- Exhaustive search
- Beam search

27

Ancestral Sampling

- Randomly sample words one by one
- Provides diverse output (high variance)

Until $\tilde{x}_t = \langle eos \rangle$ The $x_0 Y$ z_0 $Y = h_7$



(figure credit: Luong, Cho, and Manning)

Greedy decoding



Compute argmax at every step of decoder to generate word

What's wrong?

Exhaustive search?

- Find arg max $P(y_1, ..., y_T | x_1, ..., x_n)$ $y_1, ..., y_T$
- Requires computing all possible sequences
 - $O(V^T)$ complexity!
 - Too expensive

Recall: Beam search (a middle ground)

- Key idea: At every step, keep track of the k most probable partial translations (hypotheses)
- Score of each hypothesis = log probability

$$\sum_{t=1}^{j} \log P(y_t | y_1, \dots, y_{t-1}, x_1, \dots, x_n)$$

- Not guaranteed to be optimal
- More efficient than exhaustive search

Beam size = k = 2. Blue numbers = score



Beam decoding

$$e(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{LM}(y_i | y_1, \dots, y_{i-1}, x)$$

Beam size = k = 2. Blue numbers = score



Beam decoding

$$e(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{LM}(y_i | y_1, \dots, y_{i-1}, x)$$



Beam decoding



Backtrack



- Different hypotheses may produce $\langle eos \rangle$ (end) token at different time steps
 - When a hypothesis produces $\langle eos \rangle$, stop expanding it and place it aside
- Continue beam search until:
 - All k hypotheses produce $\langle eos \rangle$ OR
 - Hit max decoding limit T
- Select top hypotheses using the *normalized* likelihood score

$$\frac{1}{T}\sum_{t=1}^{T}\log P(y)$$

Otherwise shorter hypotheses have higher scores

Beam decoding

- $y_t | y_1, \dots, y_{t-1}, x_1, \dots, x_n)$

Evaluating text generation

Evaluating translation quality

- Two main criteria:
 - content of $w^{(s)}$

To Vinay it like Py Vinay debugs mem Vinay likes Python

Different translations of A Vinay le gusta Python

• Adequacy: Translation $w^{(t)}$ should adequately reflect the linguistic

• Fluency: Translation $w^{(t)}$ should be fluent text in the target language

	Adequate?	Fluent?
thon	yes	no
iory leaks	no	yes
1	yes	yes

Evaluation metrics

- Manual evaluation is most accurate, but expensive
- Automated evaluation metrics:
 - Compare system hypothesis with reference translations
 - BiLingual Evaluation Understudy (BLEU) (Papineni et al., 2002)
 - Modified n-gram precision

 $p_n =$

number of *n*-grams appearing in both reference and hypothesis translations number of *n*-grams appearing in the hypothesis translation



Example

Reference: Vinay likes programming in Python

Hypothesis/Candidate

Vinay likes Python

To Vinay it like Python

BLEU

BLEU-N = $\exp \frac{1}{N} \sum_{n=1}^{N} \log p_n$ n-gram precision geometric mean over several values of n (up to N=4)

p_1	p 2	BLEU-2
3/3	1/2	0.7071
2/5	0	???

https://www.aclweb.org/anthology/P02-1040.pdf

Two modifications:

- Each n-gram in reference can be used at most once

should not get a unigram precision of 1 ($p_1 = 2/5$) clipped count Precision-based metrics favor short translations brevity penalty 8.0 and 9.0 2.0 Solution: Multiply score with a brevity penalty for translations output/reference length shorter than reference, $BP = e^{1-r/h}$ r = reference length, h = hypothesis length

BLEU

n-gram precision BLEU-N = exp $\frac{1}{N} \sum_{n=1}^{N} \log p_n$ geometric mean over several values of n (up to N=4)

• To avoid log 0, all precisions are smoothed Various smoothing techniques add 1 to numerator/denominator

• Ex. Hypothesis: to to to to to vs Reference: to be or not to be





• Correlates somewhat well with human judgements



Human Judgments

BLEU

(G. Doddington, NIST)

BLEU scores

https://www.nltk.org/_modules/nltk/translate/bleu_score.html

Sample BLEU scores for various system outputs $BP = e^{1-r/2}$							
Translation $p_1 p_2 p_3 p_4$ BI						BP	BLEU
Reference	Vinay likes programming in Python						
Sys1	To Vinay it like to program Python	$\frac{2}{7}$	0	0	0	1	.21
Sys2	Vinay likes Python	$\frac{3}{3}$	$\frac{1}{2}$	0	0	.51	.33
Sys3	Vinay likes programming in his pajamas	$\frac{4}{6}$	$\frac{3}{5}$	$\frac{2}{4}$	$\frac{1}{3}$	1	.76

	Sample BLEU scores for various system outputs				$BP = e^{1 - r/h}$				
ength		Translation	p_1	p_2	p_3	p_4	BP	BLEU	_
5	Reference	Vinay likes programming in Python							-
7	Sys1	To Vinay it like to program Python	$\frac{2}{7}$	0	0	0	1	.21	
3	Sys2	Vinay likes Python	$\frac{3}{3}$	$\frac{1}{2}$	0	0	.51	.33	
6	Sys3	Vinay likes programming in his pajamas	$\frac{4}{6}$	$\frac{3}{5}$	$\frac{2}{4}$	$\frac{1}{3}$	1	.76	

Example from: https://github.com/jacobeisenstein/gt-nlp-class/tree/master/notes

- Alternatives have been proposed:
 - METEOR: weighted F-measure
 - Translation Error Rate (TER): Edit distance between hypothesis and reference

Issues?

- Number is not that meaningful (BLEU will be higher for some language than others)
- Does not account for different word choices (synonyms)
- Does not account for morphology
- Does not penalize omitting important words





BLEU useful despite issues

- easy to compute
- automated

- consistent

)ensity



Minor note about <UNK>

Make sure you compare against the original reference (Don't have <UNK>s in your reference)

Re-evaluating Automatic Metrics for Image Captioning [Kilickaya et al, EACL 2017]

Sequence to sequence models with attention

Issues with vanilla seq2seq



- A single encoding vector, h^{enc}, needs to capture all the information about source sentence
- Longer sequences can lead to vanishing gradients
- Overfitting

Issues with vanilla seq2seq



- A single encoding vector, h^{enc} , needs to capture all the information about source sentence
- Longer sequences can lead to vanishing gradients
- Overfitting

Attention

- The neural MT equivalent of alignment models
- Key idea: At each time step during decoding, focus on a particular part of source sentence
 - This depends on the decoder's current hidden state (i.e. notion of what you are trying to decode)
 - Usually implemented as a probability distribution over the hidden states of the encoder (h_i^{enc})







Take softmax to turn the scores into a probability distribution







Use the attention distribution to take a **weighted sum** of the **encoder hidden** states.

The attention output mostly contains information from the hidden states that received high attention.









Computing attention



• Encoder hidden states: $h_1^{enc}, \ldots, h_n^{enc}$

• Decoder hidden state at time t: h_t^{dec}

First, get attention scores for this time step (we will see what g is soon!): $e^{t} = [g(h_1^{enc}, h_t^{dec}), \dots, g(h_n^{enc}, h_t^{dec})]$

Obtain the attention distribution using softmax:

 $\alpha^{t} = \operatorname{softmax} (e^{t}) \in \mathbb{R}^{n}$

Compute weighted sum of encoder hidden states:

$$a_t = \sum_{i=1}^n \alpha_i^t h_i^{enc} \in \mathbb{R}^h$$

Finally, concatenate with decoder state and pass on to output layer:

Types of attention

Dot-product attention: $g(h_i, z) = z^T h_i \in \mathbb{R}$ more efficient (matrix 2. Bilinear / multiplicative attention: multiplication)

Additive attention (essentially MLP): 3.

• Assume encoder hidden states h_1, h_2, \ldots, h_n and decoder hidden state z

Simplest (no extra parameters) requires z and h_i to be same size

 $g(h_i, z) = z^T W h_i \in \mathbb{R}$, where W is a weight matrix

More flexible than dot-product (W is trainable)

 $g(h_i, z) = v^T \tanh(W_1h_i + W_2z) \in \mathbb{R}$

where W_1, W_2 are weight matrices and v is a weight vector

Perform better for larger dimensions





Attention can be applied to other modalities





• Agent experience



Attention on other modalities

proposals

Object

Image Credit: Peter Anderson

$$C = \{\boldsymbol{h}_1, \dots, \boldsymbol{h}_5\}$$

or $C = \{\boldsymbol{v}_1, \dots, \boldsymbol{v}_6\}$

Image captioning example



Xu et al. ICML 2015

Soft vs Hard Attention

- Soft: Each attention candidate is weighted by α_i $\widehat{\boldsymbol{v}} = \sum_{i=1}^{k} \alpha_i \, \boldsymbol{v}_i$
 - Easy to train (smooth and differentiable) • But can be expensive over large input
- Hard: Use α_i as a sample probability to pick one attention candidate as input to subsequent layers Trainable with REINFORCE approaches (Xu et al. ICML) 2015), or Gumbel-Softmax (Jang et al. ICLR 2017)



Xu et al. ICML 2015

Global vs Local Attention

- Global: attention over the entire input
- Local: attention over a window (or subset) of the input



e input (or subset) of the inpu[.]



Local: subset of source states.

Luong et al, 2015

Self-Attention

• Attention (correlation) with different parts of itself



https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html

• Transformers: modules with scaled dot-product self-attention

The	The
animal	animal
didn't	didn't
cross	cross
the	the
street	street
because	because
it	it
was	was
too	too
wide	wide

Transformers: self-attention



• More recent models (e.g. Transformer,

- Vaswani et al., 2017) have replaced
- RNNs entirely with attention
- mechanisms
- Theoretically limiting (since recurrence
 - can help handle arbitrarily long
 - sequences)
- Huge gains in practical performance

Issues with vanilla seq2seq



- A single encoding vector, h^{enc}, needs to capture all the information about source sentence
- Longer sequences can lead to vanishing gradients
- Overfitting

Exposure bias

- Discrepancy in model input between training and generation time
- During training, model inputs are gold context tokens

$$\mathcal{L}_{MLE} = -\sum_{t=1}^{T} \log P(y_t^* | \{y_{< t}^*\})$$

• At generation time, inputs are previouslydecoded tokens

$$\mathcal{L}_{dec} = -\sum_{t=1}^{I} \log P(\hat{y}_t | \{ \hat{y}_{< t} \})$$

Student forcing: use predicted tokens during training Scheduled sampling: use decoded token with some probability p, increase p over time







(figure credit: Bengio et al, 2015)

Regularization

- Weight Decay
- Dropout
- Ensembling

Regularization: Dropout

- Form of regularization for RNNs (and any NN in general)
- Idea: "Handicap" NN by removing hidden units stochastically
 - set each hidden unit in a layer to 0 with probability p
 during training (p = 0.5 usually works well)
 - scale outputs by 1/(1-p)
 - hidden units forced to learn more general patterns
- Test time: Use all activations (no need to rescale)



(a) Standard Neural Net



(b) After applying dropout.

Dropout and attention improves translation

System

Winning WMT'14 system – phrase-based + larg

Existing NMT systems

RNNsearch (Jean et al., 2015)

RNNsearch + unk replace (Jean et al., 2015)

RNNsearch + unk replace + large vocab + ensen

Our NMT systems

Base

Base + reverse

Base + reverse + dropout

- Base + reverse + dropout + global attention (local
- Base + reverse + dropout + global attention (local
- Base + reverse + dropout + local-p attention (gen Base + reverse + dropout + local-p attention (gen

Ensemble 8 models + unk replace

WMT'14 English to German Results

	Ppl	BLEU
ge LM (Buck et al., 2014)		20.7
		16.5
		19.0
nble 8 models (Jean et al., 2015)		21.6
	10.6	11.3
	9.9	12.6 (+1.3)
	8.1	14.0 (+1.4)
ation)	7.3	16.8 (+2.8)
ation) + feed input	6.4	18.1 (+ <i>1</i> .3)
neral) + feed input	50	19.0 (+0.9)
neral) + feed input + unk replace	5.9	20.9 (+1.9)
		23.0 (+2.1)

(Luong et al, 2015)



Other challenges with NMT

- Out of vocabulary (OOV)
- Low-resource languages
- Long-term context
- Common sense knowledge (e.g. hot dog, paper jam)
- Fairness and bias
- Interpretability

Out of vocabulary (OOV)

- Subword-modeling
 - Character level GRU
 - Byte-pair encoding

Fully Character-Level Neural Machine Translation without Explicit Segmentation

Jason Lee, Kyunghyun Cho, Thomas Hoffmann. 2017. Encoder as below; decoder is a char-level GRU



(Lee et al, 2017)

Copy mechanism



- Probability of generating from vocabulary or copying from input
- Probability of copying specific word (similar to attention)

(See et al, 2017)

