

Lecture 14: Transformers and Self-Attention

Ashish Vaswani and Anna Huang

Learning Representation of Variable Length Data

RNNs

- LSTMs, GRUs and variants dominate recurrent models
- Natural fit for sentences and sequences of pixel

But

Sequential computation inhibits parallelization

No explicit modeling of long and short range dependencies

Attention \Rightarrow why not use attention for representation

Text generation

Classification & regression with self attention

[Parikh et al. 2016, Lin et al. 2016]

Self attention with RNNs [Long et al. 2016, Shao, Grews et al. 2017]

Recurrent attention [Sukhbaatar et al. 2015]

Attention is cheap

[FLOPS]

Self Attention $O(\text{length}^2 \cdot \text{dim})$ if $\text{length} \leq \text{dim}$

RNN(LSTM) $O(\text{length} \cdot \text{dim}^2)$

Convolution $O(\text{length} \cdot \text{dim}^2 \cdot \text{kernel-width})$

@ Multi-head attention:

this is to put attention in different positions
separately

Importance of Residuals

Residuals carry positional information to higher layer, among other information

Music Generation Finding similar motifs

Probabilistic Image Generation

Texture Synthesis with Self Similarity

Non-local NN (Wang 2015)

Selfattention

$$A(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Image Transformer

Combining Locality with Self Attention

Image Transformer (ICML 2018)

Music generation using Self Attention

Music Transformer (ICLR 2019)

→ score → performance sound
composer performance instrument → listener

$$\text{MultiAttention} = \text{softmax}\left(QK^T + \text{score}(QE_{\text{rel}}^T)\right)$$

Equivalence

Relative positions:

Translational Equivalence

- Fast Decoding sequence models ICML 2018 Kaiser

- Better understanding of Vector Quantized auto encoders Roy 2018

- Blockwise Parallel Decoding for Autoregressive model (NeurIPS, 2019) Stern

Lecture 15: Natural Language Generation (NLG)

Neural approaches to NLG

Section: LM and decoding algs: NLG

↳ sub component

- Machine Translation
- (Abstactive) summarization
- Dialogue (chit-chat and task-based)

- Creative writing: storytelling, poetry generation

- free-form QA (i.e. answer is generated)

- Image Captioning

Language Modeling: The task of predicting the next word, given the words so far

$$P(y_t | y_1, y_2, \dots, y_{t-1})$$

- A system that produces this probability distribution is called a Language Model

Conditional Language model:

$$P(y_t | y_1, y_2, \dots, y_{t-1} | x)$$

Decoding algo

- Greedy Decoding

- Beam Search

→ search algorithm which aims to find a high probability sequence (non-necessary the optimal sequence, though) by tracking multiple possible sequences at once

what's the effect of changing beam size k?

• small k has similar problems to greedy decoding

→ ungrammatical, unnatural, nonsensical, incorrect

• Larger k means you consider more hypotheses

→ increasing k reduces some of the problem

above

→ Larger k is more computationally expensive

→ Increasing k too much decreases BLEU

score. This is primarily because large-k

beam search produces too short translation

→ In open ended tasks, large-k can make it output
more generic

Sampling-based decoding

• Pure sampling

→ On each step t, randomly sample from

the prob dist P_t to obtain your next word

→ Like greedy decoding but sample based
instead of argmax

Top-n sampling:

- randomly sample from p_t -
 → on each step t , restricted to just the top- n most probable words.
 → increasing n to get more diverse/risky output
 → Decrease n to get more generic/safe output
~~effie both of them are more efficient than beam search~~

Softmax temperature:

On timestep t , the LM computes a prob dist p_t by applying the softmax function to a vector of scores $s \in \mathbb{R}^M$

$$p_t(w) = \frac{\exp(s_w)}{\sum_{w' \in V} \exp(s_{w'})}$$

We can apply a temperature hyperparameter τ to the softmax

$$p_t(w) = \frac{\exp(s_w/\tau)}{\sum_{w' \in V} \exp(s_{w'}/\tau)}$$

→ Raise the temperature T : P_f becomes more uniform
Thus more diverse output (probability is spread around vocab)

→ Lower the temperature T : P_f becomes more spiky
Thus less diverse output

Softmax temperature is not a decoding algo
we can apply this during test time, in
conjunction with a decoding algorithm (such
as beam search or sampling)

Section 2: NLP tasks and neural approaches for them:

Summarization

Biggaword \rightarrow news article \rightarrow headline

LSTMs \rightarrow paragraph \rightarrow sentence summary

NYT, CNN/Daily Mail \rightarrow news article \rightarrow (multi) sentence summary

With how few? full h/w to article \rightarrow summary sentences

Sent:

- Sentence simplification is a different but related task
 - rewrite the source text in a simpler way
- simple Wikipedia: standard Wikipedia sentence → simple version
- Newsela: news article → version for children

Two main strategies:

Extractive summarization

Select parts (typically sentences) of the original text to form a summary

• Easier

• Restrictive (no paraphrasing)

Abstractive Summarization

Generate new text using natural language techniques

• More difficult

• More flexible (more human)

Evaluation:

ROUGE (Recall)-Oriented Understudy for Existing Evaluation

Evaluation

Rouge-N

like BLEU, it's based on n-gram overlap.

• ROUGE has no brevity penalty

• ROUGE based on recall, while BLEU is based on precision

• Precision is more important for MT and recall

is more important for summarization (assuming you have a max length constraint)

• However, often a F1 version of ROUGE is reported (anyway)

• BLEU is reported as a single number, which is combination of the precisions for $n=1, 2, 3, 4$ n-grams

• ROUGE scores are reported separately for each n-grams

ROUGE-L \Rightarrow unigram overlap

ROUGE-D \Rightarrow bigram overlap

ROUGE-LCS \Rightarrow LCS overlap

Neural approaches for summarization (2015-present)

→ 2015: Rush et al. publish the first seq2seq

Summarization paper

→ Single-document abstractive summarization is

a translation task!

Thus we can apply standard seq2seq + attention NMT methods

since 2015, developments

- Making it easier to copy [prevent too much of it]
- Hierarchical/multi-level attention
- More global/high-level content selection
- Using RL to directly maximize ROUGE on discrete goals
- Resurrecting pre-neural ideas (graph algos) and

working them ~~out~~ into neural systems

Copy Mechanism:

- Seq2Seq + attention are good at writing fluent output but bad at copying over details correctly
- Copy mechanisms use attention to enable seq2seq system to easily copy words and phrases from the input to the output.
 - Allowing both copying and generating gives us hybrid extractive/abstactive approach

$$P(w) = P_{\text{gen}} \text{Prob}(a) + (1 - P_{\text{gen}}) \sum_{i: w_i = w} \alpha_i^t$$

- Big problems with copying mechanism

They copy too much

bad at overall content selection

Preneural summarization had separate stages for
content selection and surface realization

Bottom-up summarization:

- Content selection stage: Use a neural sequence-tagging model to tag words as include or don't include
 - Bottom-up attention stage: The seq2seq+attention system can't attend to words tagged don't-include (apply a mask)
- Simple but effective → better overall content selection
Less copying

Neural summarization via RL:

- Use RL to directly optimize ROUGE-L
- ML+RL gives best result

Dialogue

• Task-oriented dialogue

→ Assistive (customer service, recommendations)

→ Co-operative (two agents solve a task together)

→ Adversarial (two agents compete in a task through dialogue)

• Social dialogue

→ chit-chat dialogue

→ Therapy/mental wellbeing

Preneural dialogue systems used predefined

templates or retrieve an appropriate response from a corpus of responses.

After 2015: open ended freeform dialogue.

→ Seq2Seq-based dialogue

However it quickly became apparent that a naive application of seq2seq has serious deficiencies for (chit chat) dialogue

- Genericness/boring responses
- Irrelevant responses (not related to experience)
- Repetition (showing the same thing over and over again)
- Lack of context
- Lack of consistent persona

Irrelevant response problem:

Reason:

- because it's generic
- or because changing the subject to something unrelated

One solution: optimize for MMI (Maximum Mutual Information)

between input S and response T

$$\log \frac{p(s, T)}{p(s)p(T)}$$

$$\hat{T} = \arg \max \left\{ \log \frac{p(T|s)}{p(T)} \right\}$$

Genericness or boring response problem

→ Easy test-time fixes:

- Directly upweight rare words during beam search
- Use a sampling decoding algo rather than beam search

→ Conditioning fixes

- Condition the decoder on some additional content (e.g. sample some content words and attend them to them)
- Train a retrieve-and-refine model rather than a generate-from-scratch model

Repetition Problems

Simple Solution

→ Directly block repeating n-grams during beam search

More complex solutions

- Train a coverage mechanism - in seq2seq
this is an objective that prevents the attention mechanism from attending to the same words multiple times.

- Define a training objective to discourage repetition
- If this is a differentiable function of the generated output, then will need some techniques e.g. RL to train

NLG: Storytelling:

Most neural storytelling work uses some kind of prompt:

- Generate a story-like paragraph given an image
- Generate a story given a brief writing prompt
- Generate the next sentence of a story given the story so far (story continuation)

Question: How to get around the lack of parallel data?

Answer: Use a common sentence encoding space

→ skip-thought vectors are a type of general-purpose sentence embedding method

- The idea is similar to how we learn an embedding for a word by trying to predict the words around it.

• Using COCO (an image captioning dataset), learn a mapping from images to the skip-thought encodings of their captions

• Using the target style corpus, train an RNN-LM to decode a skip-thought vector to the original text

• Put the two together

Generating story from a writing prompt

In 2018, Fan et. al. released a new story generation dataset collected from Reddit's WritingPrompts subreddit.

- Each story has an associated brief writing prompt or title.
 - propose a convolution-based story model
- Gated Multi-head multiscale self attention
- The self-attention is important for capturing long-range context
 - The gates allow the attention mechanisms to be more selective
 - The different attention heads attend to different scales - this means there are different attention mechanisms dedicated to retrieving fine-grained information and coarse-grained information.

Model-fusion:

- Pretrain one seq2seq model, then train a second seq2seq model that has access to the hidden states of the first.
- The idea is that the first seq2seq model learns general LM and the second learns to condition on the prompt.
- The results are impressive

• Related to prompt

• Diverse; non-generic

• Stylistically dramatic

However,

- Mostly atmospheric / descriptive / scene-setting; less events / plot.
- When generating for longer, mostly stays on the same idea without moving forward to new ideas - coherence issues.

NLG Evaluation:

Word overlap based metrics (BLEU, ROUGE, METEOR, F1, etc.)

- We know they are not ideal for machine translation worse for summarization and even worse for dialogue and storytelling

What about perplexity?

- Captures how powerful your LM is, but doesn't tell you anything about generation (e.g. if your decoding algo is bad, perplexity is unaffected)

Word embedding based metrics?

- compare the similarity of the word embeddings not just the overlap of the words themselves.
- Captures semantics in a more flexible way
- still doesn't correlate well with human judgements for open-ended task

Automatic evaluations of NLG

- We can define more focused automatic capture particular aspects of generated text
- Fluency (compute probability wrt a well trained LM)
- Correct style prob wrt ~~to the~~ LM trained on target corpus
- Diversity (rare word usage, uniqueness of strings, n-grams)
- Relevance to input (semantic similarity measures)
- Simple things like lengths and repetitions
- Task specific metrics e.g. compression rate for summarization

Though these doesn't measure overall quality, they can help us track some important qualities that we care about.

Human evaluation

• regarded as gold standard but slow and expensive

• does human evaluation solve all your problems? No!

Conducting human evaluation effectively is very difficult

Humans are inconsistent

• can be illogical

• lose concentration

• misinterpret your question

• can't always explain why they feel the way they do.

Adversarial discrimination

- Test whether the NLG system can fool a discriminator which is trained to distinguish human text from artificially generated text

Section 4:

Thoughts on NLG research, current trends and the future:

- Incorporating discrete latent variables into NLG
- Alternatives to strict left-to-right generation
- Alternatives to maximum likelihood training with teacher forcing
- NLG is the wildest part remaining in NN
- Neural NLG community is rapidly expanding
- Biggest roadblock is still evaluation

8 things on NLG

- 1) The more open-ended the task the harder everything becomes
- 2) Aiming for specific improvement can also be more manageable than aiming to improve overall generation quality
- 3) If you're using a LM for NLG, improving the LM (i.e. perplexity) will most likely improve generation quality.
- 4) Look at your outputs a lot
- 5) You need an automatic metric, even if it's imperfect (probably several)
- 6) If we do human eval, make the questions as focused as possible
- 7) Reproducibility is a huge problem in today's NLP+DL, and a "huge" problem in NLG
• publicly release all your generated outputs.

Working with NLG is very frustrating.
But also very funny.

Non-autoregressive generation for NMT

In 2018 Gauvret et al. published Non-autoregressive generation for NMT. It generates the translation left-to-right, with each word depending on the ones before it.

- It generates the translation in parallel.
- It has obvious efficiency advantages but it is also intriguing from NLP point of view
- Transformer based architecture and decoder can run in parallel at text time