Oversmoothing in neural autoregressive modeling





Kyunghyun ChoWork done together with Ilia Kulikov and Maksim Eremeev



Language modeling

- Setup:
 - $Y = (y_1, y_2, \dots, y_L)$, where $y_l \in V$ with V a set of all unique tokens.
 - The length $L \in \{1,2,3,\ldots\}$ may vary from one sequence to another.
 - *V* includes a special symbol $\langle eos \rangle$, and $y_L = \langle eos \rangle$ for any *Y*.
- Goal:
 - Build a model (parametrized by θ) that computes a probability $p(Y; \theta)$.
 - A reasonable Y receives a high probaiblity

Language modeling is fascinating because language is a fascinating way to encapsulate ideas

- Relation extraction
 - <u>Barak Obama</u> was the 44th president of the United States.
- Common-sense reasoning
 - Jane performed an emergency surgery on John, because John was in a car accident.
- Question-answering
- Machine translation
 - **Iocated in New York City.**"

• Wikipedia is a free content, multilingual online encyclopedia written and maintained by a community of volunteers through a model of open collaboration, using a wiki-based editing system.

• "<u>뉴욕대학교는 뉴욕시에 위치한 사립대학교이다.</u>" is <u>in English</u> "New York University is a private university



Language modeling is fascinating because language is a fascinating way to encapsulate ideas

- All these problems can be solved by solving
 - *Y*_{missing}
- So, all we need is a language model.

arg max $\log p([y_{observed}, y_{missing}]; \theta)$

- (Closely related) properties of language modeling
 - We deal with variable-length sequences.

Each element within a sequence is from a finite set of unique tokens.

Autoregressive modeling

$p(Y; \theta)$

- One classifier is used over and over
 - It can handle sequences of varying length.
 - It is compositional (rather simple one, though.)

$$p(y_t | y_{1:t-1}; \theta) = \prod_{t=1}^{L} p(y_t | y_{1:t-1}; \theta)$$

- Autoregressive modeling $p(Y;\theta) = \prod_{t=1}^{L} p(y_t | y_{1:t-1};\theta)$
- A natural choice for language modeling since 1950

2. ENTROPY CALCULATION FROM THE STATISTICS OF ENGLISH

One method of calculating the entropy // is by a series of approximations 7? $F_1, F_2 > \cdots$, which successively take more and more of the statistics of the language into account and approach H as a limit. F_N may be called the N-gram entropy; it measures the amount of information or entropy due to statistics extending over N adjacent letters of text. F_N is given by¹

$$FN = -\sum_{i,j} p(b_i, j) \log_2 p_{b_i}(j) = -\sum_{i,j} p(b_i, j) \log_2 p(b_i, j) + \sum_i p(b_i) \log p(b_i)$$
(1)

in which: bi is a block of N-1 letters [(N-1)-gram]

j is an arbitrary letter following b_i

 $p(b_i, j)$ is the probability of the N-gram b_i, j

 $p_{b_i}(j)$ is the conditional probability of letter j after the block b_i ,

and is given by $p(b_i, j)/p(b_i)$.

The equation (1) can be interpreted as measuring the average uncertainty (conditional entropy) of the next letter /' when the preceding N-1 letters are known. As N is increased, F_N includes longer and longer range statistics and the entropy, H, is given by the limiting value of F_N as $N \to \infty$:

$$H = \lim_{N \to \infty} F_N \,. \tag{2}$$

The N-gram entropies F_N for small values of N can be calculated from standard tables of letter, digram and trigram frequencies.² If spaces and punctuation are ignored we have a twenty-six letter alphabet and F_6 may be taken (by definition) to be $\log_2 26$, or 4.7 bits per letter. F_1 involves letter frequencies and is given by

nnon, 1950]
$$F_1 = -\frac{X_1}{i-1} p(i) \log_2 p(i) = 4.14$$
 bits per letter. (3)



- Autoregressive modeling $p(Y;\theta) = \prod_{t=1}^{L} p(y_t | y_{1:t-1};\theta)$
- A natural choice for language modeling since 1950.
- Extended to a neural version in 2001.





- Autoregressive modeling $p(Y;\theta) = \prod p(y_t | y_{1:t-1};\theta)$ t=1
- A natural choice for language modeling since 1950.
- Extended to a neural version in 2001.
- Infinite-context version in 2011:
 - This is when things got so much more interesting.



[Mikolov et al., 2011]



- Autoregressive modeling $p(Y;\theta) = \prod p(y_t | y_{1:t-1};\theta)$ t=1
- A natural choice for language modeling since 1950.
- A neural version in 2001.
- Infinite-context version in 2011.
- Now: the war of "# parameters"

This is a bit ridiculous that we are spending so many scientists on this war, though...



10

Neural autoregressive language modeling

t - 1

- - Target $y_t \in V$ $\operatorname{Input}(y_1, \dots, y_{t-1}) \in \underbrace{V \times \cdots \times V}_{}$
 - Model parameters θ

Language modeling is reduced to text classification: $p(y_t | y_1, \dots, y_{t-1}; \theta)$

- Two natural ways to build this classifier $p(y_t | y_1, ..., y_{t-1}; \theta)$
 - A recurrent network: LSTM, GRU, etc.
 - **Pros:** constant memory/computation in the forward pass
 - Cons: difficult learning and inherent sequential processing
 - An attention network: transformers, etc.
 - Pros: highly parallel processing and easy learning
 - Cons: linear memory/computation in the forward pass

But, this is not what I want to talk about it today...



Training a language model

- Language modeling is reduced to text classification: $p(y_t | y_1, \dots, y_{t-1}; \theta)$
- For learning, we use cross-entropy:

$$FN = -\sum_{y > 1}^{2}$$

 \bullet

$$-\sum_{y \in V} p^*(y)$$

• Entropy: expected negative log-probability under the same distribution

 $\sum p(b_i, j) \log_2 p_{b_i}(j)$

Cross-entropy: expected negative log-probability under a different distribution

x)log $p(y | x; \theta)$

Training a language model almost exactly like training a classifier

- For training a language model, we use crossentropy which is the upperbound to entropy.
- In other words, training a language model is find a distribution that matches the inherent entropy of natural language.

because



$$-\sum_{y\in V}p^*(y|x)\log p(y|x; heta)\geq -\sum_{y\in V}p^*(y|x)\log p^*(y|x),$$

$$egin{aligned} y|x)\log p(y|x; heta) &= \sum_{y\in V} p^*(y|x)\lograc{p^*(y|x)}{p^*(y|x)}p(y|x; heta) \ &= \sum_{y\in V} p^*(y|x)\left(\log p^*(y|x) + \lograc{p(y|x; heta)}{p^*(y|x)}
ight) \ &= \sum_{y\in V} p^*(y|x)\log p^*(y|x) + \sum_{y\in V} p^*(y|x)\lograc{p(y|x)}{p^*(y|x)} \ &= \sum_{y\in V} p^*(y|x)\log p^*(y|x). \end{aligned}$$



Training a language model Almost exactly like training a classifier

• Overall, the loss function is

$$\mathsf{E}_{Y \sim D} \left[\sum_{t=1}^{|Y|} \log p(y_t | y_{< t}; \theta) \right] \approx \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{N} \log p(y_t^n | y_{< t}^n; \theta)$$

- Almost like training a classifier, but not quite.
- Examples $(y_{<t}^n, y_t^n)$ are highly correlated within each sequence Y^n .
 - Highly correlated, because they come from the same text.

But, it turned out no one really cares about this... ugh...

Taking a step back ...

- Setup:
 - $Y = (y_1, y_2, \dots, y_I)$, where $y_I \in V$ with V a set of all unique tokens.
 - The length $L \in \{1, 2, 3, ...\}$ may vary from one sequence to another.
 - V includes a special symbol $\langle eos \rangle$, and $y_L = \langle eos \rangle$ for any Y.
- Goal:
 - Build a model (parametrized by θ) that computes a probability $p(Y; \theta)$.
 - A reasonable Y receives a high probaiblity

If <eos> wasn't there ...

- A prefix $y_{1:t}$ is always more likely than the full sequence $y_{1:T}$: $p(y_1)p(y_2 | y_1) \cdots p(y_t | y_{< t}) \ge p(y_1)p(y_2 | y_1) \cdots p(y_t | y_{< t}) \cdots p(y_T | y_{< T})$
- because the probability is bounded from above by 1.
- Weird...

But, <eos> is here ...

- <eos> saves us from this weird behaviour.
 - Any valid sequence ends with <eos>.
- A prefix $y_{1:t}$ is always more likely than the full sequence $y_{1:T}$:
 - Instead of a prefix $y_{1:t}$, consider a premature sequence $y_{1:t} + \langle \cos \rangle$

 $p(y_1)p(y_2 | y_1) \cdots p(y_t | y_{< t}) p(\langle eos \rangle | y_{\le t})$

)?
$$p(y_1)p(y_2|y_1)\cdots p(y_t|y_{$$



But, <eos> is here ...

- A prefix $y_{1:t}$ is always more likely than the full sequence $y_{1:T}$:
 - Instead of a prefix $y_{1,t}$, consider a premature sequence $y_{1,t} + \langle eos \rangle$

• The probability of a premature seq. should be lower than the full sequence:

• In other words, as long as <eos> is *highly* unlikely after a prefix, the premature sequence can be less probable than the full sequence.

 $p(y_1)p(y_2|y_1)\cdots p(y_t|y_{<t})p(\langle eos \rangle | y_{<t})? p(y_1)p(y_2|y_1)\cdots p(y_t|y_{<t})\cdots p(y_T|y_{<T})$

 $p(y_1)p(y_2|y_1)\cdots p(y_t|y_{<t})p(\langle eos \rangle | y_{<t}) \le p(y_1)p(y_2|y_1)\cdots p(y_t|y_{<t})p(y_{t+1}|y_{<t+1})\cdots p(y_t|y_{<t})\cdots p(y_t|y_{<t})$



Oversmoothing rate

violated:

$$r(y) = \frac{1}{T} \sum_{t=1}^{T} I\left(p(\langle \cos \rangle | y_{< t}) > \prod_{t'=t}^{T} p(y_{t'} | y_{< t'}) \right)$$

- We want this oversmoothing rate to be small (pprox 0) in natural language.

• Given a sequence $y = (y_1, y_2, \dots, y_T)$, we check how often the inequality is

• It counts how many premature seq.'s are more probable than the full seq.

Oversmoothing rate is high!



Oversmoothing rate is high!

- It is not enough to simply lower $p(\langle e$
- $p(\langle eos \rangle | y_{< t})$ must be lowered below
 - i.e., we must take into account the sequence structure.

$$|y_{ when $y_t \neq \langle \cos \rangle$.$$

• i.e., it is not enough to treat it as a series of independent classification.

bw
$$\prod_{t'=t}^{T} p(y_{t'} | y_{< t'}).$$

Minimizing the oversmoothing rate

The oversmoothing rate is not easy to minimize

$$r(y) = \frac{1}{T} \sum_{t=1}^{T} I\left(p(\langle \cos \rangle \,|\, y_{< t}) > \prod_{t'=t}^{T} p(y_{t'} |\, y_{< t'}) \right)$$

• because it is an average of 0-1 losses (piece-wise constant.)



Oversmoothing loss

• A convex relaxation with a piece-wise linear loss and margin $m \ge 0$:

$$o(y) = \frac{1}{T} \sum_{t=1}^{T} \max\left(0, m + \log t\right)$$

$$\max\left(0, m + \log p(\langle \cos \rangle | y_{< t}) - \log \prod_{t'=t}^{T} p(y_t)\right)$$



Oversmoothing loss

• A convex relaxation with a piece-wise linear loss and margin $m \ge 0$:

$$o(y) = \frac{1}{T} \sum_{t=1}^{T} \max\left(0, m + \log p(\langle \cos \rangle | y_{< t}) - \sum_{t'=t}^{T} \log p(y_{t'} | y_{< t'})\right)$$

to show up: we treat <eos> specially within a sequence.

• Maximizing
$$\sum_{t'=t}^{T} \log p(y_{t'} | y_{ ensures earlier tokens but also emphasize the$$

Minimizing $\log p(\langle \cos \rangle | y_{< t})$ ensures <eos> is scored low when it is not supposed

s the model does not focus too much on the

tokens in the latter part of a sequence.

We can lower the oversmoothing rate. Dramatically!



What happens under the hood?

• The log-probability of an incorrectly place <eos> decreases dramatically.

$$o(y) = \frac{1}{T} \sum_{t=1}^{T} \max\left(0, m + \log p(\langle \cos \rangle | y_{< t}) - \sum_{t'=t}^{T} \log p(y_{t'})\right)$$



Figure 2: Log-probabilities of $\langle eos \rangle$ token within lengtht prefixes averaged across all positions per translation and then averaged across all translations.



What happens under the hood?

- The effect is however much more visible when we look at their rank.
- The rank is what matters when we consider approximate decoding.

Theorem 3.4 (Inconsistency of an incomplete decoding algorithm). There exists a consistent recurrent LM p_{θ} from which an incomplete decoding algorithm \mathcal{F} , that considers only up to (|V| - 1)most likely tokens according to $p_{\theta}(y_t | y_{\leq t}, C)$ at each step t, finds an infinite-length sequence Ywith probability 1, i.e., $q_{\mathcal{F}}(|Y| = \infty) = 1$.



Figure 4: Normalized rank of (eos) token within lengtht prefixes averaged across all positions per translation and then averaged across all translations.



But, it doesn't hurt nor improve modeling

- This implies that
 - MLE has ambiguity in modeling <eos>
 - Oversmoothing loss selects a solution that avoids the issue of oversmoothing



Figure 3: Perplexity measured on reference translations is stable as oversmoothing loss has more contribution.



Let's take a quick step back ...

Oversmoothing in the wild Koehn & Knowles in 2016

- Better search (larger beam) leads to worse translation.
- Worse translations tend to be dramatically shorter.







English-Romanian













Oversmoothing in the wild Stahlberg & Byrne in 2019

• For many source sentences, the most probable translations are *empty*!



Figure 1: BLEU over the percentage of search errors. Large beam sizes yield fewer search errors but the BLEU score suffers from a length ratio below 1.

Model	Beam-10		Exact
	BLEU	#Search err.	#Empty
LSTM*	28.6	58.4%	47.7%
SliceNet*	28.8	46.0%	41.2%
Transformer-Base	30.3	57.7%	51.8%
Transformer-Big*	31.7	32.1%	25.8%

Table 2: *: The recurrent LSTM, the convolutional SliceNet (Kaiser et al., 2017), and the Transformer-Big systems are strong baselines from a WMT'18 shared task submission (Stahlberg et al., 2018a).





Oversmoothing in the wild **Post-hoc fixes**

Length penalty during decoding [Cho et al., 2014; Wu et al., 2016; ...]

When we use the beam-search to find the k best translations, we do not use a usual log-probability but one normalized with respect to the length of the translation. This prevents the RNN decoder from favoring shorter translations, behavior which was observed earlier in, e.g., (Graves, 2013).

[Cho et al., 2014]

$$s(Y,X) = \log(P(Y|X))/lp(Y) + cp(X;Y)$$
$$lp(Y) = \frac{(5+|Y|)^{\alpha}}{(5+1)^{\alpha}}$$
$$cp(X;Y) = \beta * \sum_{i=1}^{|X|} \log(\min(\sum_{j=1}^{|Y|} p_{i,j}, 1.0)),$$

[Wu et al., 2016]

Oversmoothing in the wild **Post-hoc fixes**

Learned length penalty [Murray & Chiang, 2018]

Finally, some systems add a constant word re*ward* (He et al., 2016):

$$s'(e) = s(e) + \gamma m.$$

If $\gamma = 0$, this reduces to the baseline model. The advantage of this simple reward is that it can be computed on partial translations, making it easier to integrate into beam search.

If we approximate the expectation using the mode of the distribution, we get

$$rac{\partial L}{\partial \gamma} pprox - |e^*| + |\hat{e}|$$

where \hat{e} is the 1-best translation. Then the stochastic gradient descent update is just the familiar perceptron rule:

$$\gamma \leftarrow \gamma + \eta \left(|e^*| - |\hat{e}| \right),$$

Oversmoothing in the wild **Post-hoc fixes**

• Length penalty is a bandage and does not get to the heart of the problem.

However, none of these papers directly solves the mystery outlined in Section 1, i.e., why doesn't maximimum likelihood training already shift probability away from emtpy output candidates, toward appropriate-length ones?



Let's now come back to our method.

Shi et al. [2020]

Length degeneracy disappears

- Even without using any length penalty during beam search



Strong regularization of oversmoothing removes the issue of length degeneracy



Oversmoothing vs. translation quality

- With less approximate search (beam1000), we observe significant improvement in translation quality with lessened oversmoothing
 - Mostly from ruling out unreasonably short translations.



(a) IWSLT'17 DE \rightarrow EN

• With highly approximate search (beam5), no impact on translation quality.

(b) WMT'16 EN \rightarrow DE

(c) WMT'19 EN \rightarrow DE

Problem solved ...?

Obviously not ... :(



Oversmoothing **Positive news**

- Carefully determined/defined the issue of oversmoothing:
- Carefully designed a loss function to alleviate this issue:

$$o(y) = \frac{1}{T} \sum_{t=1}^{T} \max\left(0, m + \log p(\langle \cos \rangle | y_{< t}) - \log \prod_{t'=t}^{T} p(y_{t'} | y_{< t'})\right)$$

$r(y) = \frac{1}{T} \sum_{t=1}^{T} I\left(p(\langle \cos \rangle | y_{< t}) > \prod_{t'=t}^{T} p(y_{t'} | y_{< t'}) \right)$

Experimentally demonstrated the effectiveness of this oversmoothing loss.

Oversmoothing **Negative news**

- BUT, translation quality has not improved.
- AND, better translations are still found by more approximate search Beam search with a small beam
- Mystery continues ...
 - How is beam search finding those high-quality but not necessarily most probable translations from our model?

Fin.