

Text Summarization part 3

NLG Intro

applications: Document Summarization, E-mail Summarization

Meeting Summarization

Write story, narratives, poetry

NLG: Any task involving text production for human consumption requires natural language generation

how to generate language?

at each time step t , our model computes a vector of scores for each token in our vocabulary, S_{t+1}

$$S = f(\{y_t\}) \quad f(\cdot) \text{ is model}$$

Then, compute a probability distribution P over these scores (usually with a softmax function):

$$P(y_t = w | \{y_t\}) = \frac{\exp(S_w)}{\sum_{w' \in V} \exp(S_{w'})}$$

decoding algorithm defines a function to select a token from this distribution:

$$\hat{y}_t = g(P(y_t | \{y_t\})) \quad g(\cdot) \text{ is decoding algorithm}$$

exhaustive search Decoding

want to find a (length T) translation y that maximizes

$$\begin{aligned} P(y|x) &= P(y_1|x)P(y_2|y_1, x)P(y_3|y_1, y_2, x)\dots P(y_T|y_1, \dots, y_{T-1}, x) \\ &= \prod_{t=1}^T P(y_t|y_1, \dots, y_{t-1}, x) \end{aligned}$$

could try compute all possible sequences y

This means that on each step t of the decoder, we're tracking V possible partial translations, where V is vocab size

Complexity is far too expensive: $O(V^T)$

NLG Decoding strategy

Greedy decoding

Select the highest probability token in $P(y_t|y_{\leq t})$

$$\hat{y}_t = \underset{w \in V}{\operatorname{argmax}} P(w|y_{\leq t})$$

在每一步找出当前这一步的最优解，然后把这个解当成正确的解，代入下一步进行计算

Stop criterion:

到 $\langle \text{END} \rangle$ token 停止 decode

实际 code: np.argmax(conditional_probability)

Greedy search 的问题:

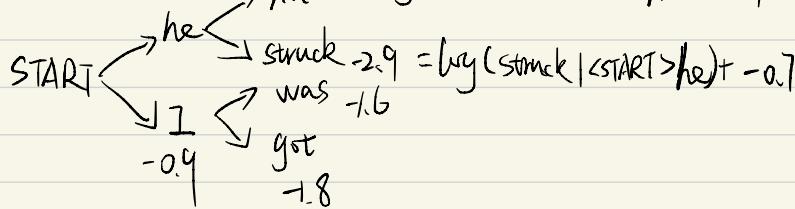
① 局部最优解不是全局的最优解

② 一错全错

Beam decoding:

从 top-1 到 top-N 个 (usually N: 5-10), 每一步看 top-N 个最后 scores 并

ex: Beam size = k = 2 Blue numbers = score (y_1, \dots, y_k) = $\sum_{i=1}^k \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



For each of the k hypotheses, find top k next words and calculate scores.

Stop Criterion

In beam search decoding, different hypotheses may produce <END> tokens on different time steps

when a hypothesis produces, hypothesis is complete

beam search 是在生成 softmax 后进行的, 也就是 inference 部分

what's the effect of changing beam size k ?

small k has similar problem to greedy decoding ($k=1$)

Ungrammatical, unnatural, nonsensical, incorrect

Larger k means consider more hypotheses

increasing k reduces some problems above

Larger k is more compute expensive

But increasing K can introduce other problems:

• for NMT, increasing K too much decreases BLEU score

• large- k beam search produces too short translations (even with score normalization!)

In open-ended tasks like chit-chat dialogue, large k can make output more generic

Effect of beam size in chitchat dialogue

low beam size: More on-topic but nonsensical;

High beam size: Converges to safe, 'correct' response, but it's generic and less relevant

Penalize longer sequences

$$S(Y) = \sum_{t=1}^T \log P(Y_t | \{Y\}_{t-1})$$

↓
shorter sequences will score better!

Solution: normalization $s(Y) = \frac{1}{|Y|} \sum_{t=1}^{|Y|} (\log P(Y_t | \{y\}_{ct}))$

Solution: Normalize by taken length relative to reference sequence

$$s(Y) = \frac{1}{|Y|} \sum_{t=1}^{|Y|} (\log P(Y_t | \{y\}_{ct})) \quad lpc(Y) = \frac{(s+|Y|)^{\alpha}}{(s+1)^{\alpha}}$$

α 在 paper 中推荐 0.5 - 0.7

Repetitive Problem 重复性问题

重复的单词会使 loss 一直下降

how to reduce repetition:

① Don't repeat n-grams (Hacky, but works!)

② Minimize additional loss term for minimizing hidden state similarity
(LSTM_s)

$$\hat{g}_t = g(\log P(Y_t | \{y\}_{ct})) - s(h_t, h_{t-m}) \quad \text{希望惩罚相近的 hidden state}$$

③ 2020 年提出: Unlikelihood objective

Penalize generation of already-seen tokens

Sequence-level Unlikelihood Training

Given a set of undesired tokens \mathcal{U} , lower their likelihood in context

$$\mathbb{Y}^t_{MLE} = -\log P(Y_t^* | \{y\}_{ct}^*) \quad L_{UL} = -\sum_{y \in \mathcal{U}} \log (-P(y | \{y\}_{ct}^*))$$

$\mathcal{U} = \{y\}_{ct}^*$ Unlikelihood objective lowers the probability of certain tokens

$$L_{\text{MLE}}^t = L_{\text{MLE}} + \alpha L_{\text{VL}}^t$$

④ $\xrightarrow{\text{F}}$ Softmax Avoid likelihood issues by factorizing the softmax

$$P(y_t | w_n | \{y\}_{t-1}) = \frac{e^{v_{nh}}}{\sum_{m=1}^M e^{v_{mh}}}$$



$$P(y_t = w_n | \{y\}_{t-1}) = \left(\frac{e^{v_{1h}}}{\sum_{c=1}^C e^{v_{ch}}} \right) \left(\frac{e^{v_{nh}}}{\sum_{m=1}^M e^{v_{mh}}} \right)$$

实际上是应用了MF的思想

sampling Strategy

random Sampling

$y = \text{sample a token from the distribution of tokens}$

$$\hat{y}_t \sim P(y_t | w | \{y\}_{t-1})$$

Too random : Temperature Scaling

Too much randomness: distribution has too much entropy

Solution: Make the distribution more "peaky" with temperature scaling

$$P(x_i | x_{1:t-1}) = \frac{\exp(u_i/t)}{\sum_j \exp(u_j/t)}$$

Raise the temperature $t > 1$: P becomes more uniform

More diverse output (probability is spread around vocab)

Lower the temperature $t < 1$: P becomes more spiky

Less diverse output (probability is concentrated on top words)

Top-k sampling:

Random sample taken from top k highest probability tokens in P .

only sample from the top k tokens in the distribution

$$\hat{y}_t \sim p^*(y_t = w | \{y\}_{\leq t})$$

Question of top-k Sampling:

概率很低的单词被 Top-k Sampling 完全抛弃

Flat Distribution: 很容易舍弃与 top-k 概率相近的单词

Peaked distribution: 可能有的单词概率很低但是因为在 top-k 里被留下来了。

Top-p (Nucleus) Sampling

p : percent

$$\sum_{x \in V(p)} p(x|x_1, \dots, x_i) \geq p$$

本质上 Top-p Sampling & Top-k Sampling 都是从 truncated Vocabulary distribution 中 sample token，区别在于置信区间的选取

通常 top-k 和 top-p Sampling 我们都会一起用，然后做 trade-off

transformer:

先做 top-k 再做 top-p

Research: paperswithcode.com

ResearchGate.net

cluebenchmark.com 中文数据集

Fiduciate Metrics:

① ROUGE-N

N: N-gram

measures the number of matching 'n-grams' between our model-generated text and a 'reference'

参考答案

Precision = $\frac{\# \text{ of } n\text{-grams found in model and reference}}{\# \text{ of } n\text{-grams in model}}$

ex: Model: "the hello a cat dog fox jumps"

reference text: "the fox jumps" → ['the', 'fox', 'jumps']



$$\frac{\text{Count}_{\text{match}}(\text{gram}_n)}{\text{Count}(\text{gram})} = \frac{3}{7} = 43\% \text{ precision}$$

ex: Rouge-1 Rouge-2 Rouge-2 都是通过 n-gram

$$\text{F1-score: } 2 * \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

② Rouge-L

measures the longest common subsequence (LCS) between our model output and reference

ex:

Model → "the hello a cat dog fox jumps"

common subsequence length=2

Common Subsequence

length=1

reference text: "the fox jumps"

longest Common Subsequence length=2

Precision:

$$\frac{LCS(\text{gram}_n)}{\text{Count}(\text{gram}_n)} = \frac{2}{7} = 29\% \text{ precision}$$

Recall: $\frac{LCS(\text{gram}_n)}{\text{Count}(\text{gram}_n)} = \frac{2}{3} = 66\%$

"reference text"

F1 Score: $2 \neq \frac{\text{Precision} + \text{recall}}{\text{precision} + \text{recall}} = 0.6$

ROUGE-S

S for skip-Gram

ex: model: "the brown fox jumps" → Precision

"the fox" → $\frac{\text{match}(\text{gram}_2)}{\text{Count}(\text{gram}_2)} = \frac{1}{3} = 33\%$
reference text

Recall: "the brown fox jumps" → $\frac{\text{match}(\text{gram}_2)}{\text{Count}(\text{gram}_2)} = \frac{1}{1} = 100\%$

"the fox" →
reference text

ROUGE 的缺点 (Cons):

无法去覆盖语义的表达

SSHAgent: 通过给多个表达式（可能都是一个意思），在计算时分别对它们做计算，也就是生成好几个分值，然后再对这些 Range 的分值做 average，这样会一定程度的缓解

| php install range