Diving into the Training for the Sequence-to-sequence Model Yang Feng Institute of Computing Technology Chinese Academy of Sciences 2019.10.12





The Sequence-to-sequence Model

- Encoder-Decoder Framework
 - Encoder: encode the source into representations
 - Decoder: decode the representation into the target



Main Models



RNNSearch





Training

Teacher Forcing





Problems



Inference



Problems



Inference





Ground Truth:How should I dealwith this issue ?Output:How should I process this issue ?





Ground Truth:How should I dealwith this issue?Output:How should I process this issue?

Word-level Matching



1 Exposure Bias

2 Word-level Matching

CONTENTS



CONTENTS







Key parts

- How to generate oracle translation
- How to sample context
- How to train the model

Exposure Bias		Word-level Matching		而中国科学的计算技术研究的
Scheduled Sampling	Professor Forcing	RL-based Training	Differentiable Loss	INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES

Oracle Translation Generation

• Word-level Oracle (WO)

Sentence-level Oracle (SO)

Exposure Bias		Word-level Matching	而中国科学的计算技术研究的
Scheduled Sampling	Professor Forcing	RL-based Training Differentiable Los	INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES

The RNNSearch Model



Exposure Bias		Word-level Matching		而中国科学的计算技术研究的
Scheduled Sampling	Professor Forcing	RL-based Training	Differentiable Loss	INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES

Oracle Translation Generation

Word Oracle





Oracle Translation Generation

- Sentence Oracle
 - Generate top-k translations by beam search
 - Rerank the top-k translation with BLEU
 - Select the top

Exposure Bias		Word-level Matching		而中国科学院计算技术研究的
Scheduled Sampling	Professor Forcing	RL-based Training	Differentiable Loss	INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES

Scheduled Sampling with Decay





Training

• Teacher Forcing + MLE

$$\mathcal{L}(\theta) = -\sum_{n=1}^{N} \sum_{j=1}^{|\mathbf{y}^n|} \log P_j^n \left[y_j^{*n} \right]$$

CONTENTS









Discriminator





• Adversarial training in the GAN style



Exposure B	ias
------------	-----

Scheduled Sampling

Word-level Matching

Differentiable Loss



Loss for Discriminator

Professor Forcing

 $C_d(\boldsymbol{\theta}_d|\boldsymbol{\theta}_g) = E_{(\boldsymbol{x},\boldsymbol{y})\sim\text{data}}[-\log D(B(\boldsymbol{x},\boldsymbol{y},\boldsymbol{\theta}_g),\boldsymbol{\theta}_d) + E_{\boldsymbol{y}\sim P_{\boldsymbol{\theta}_g}(\boldsymbol{y}|\boldsymbol{x})}[-\log(1-D(B(\boldsymbol{x},\boldsymbol{y},\boldsymbol{\theta}_g),\boldsymbol{\theta}_d))]]$

RL-based Training



Loss for Generator

Conventional Loss:

$$NLL(\boldsymbol{\theta}_g) = E_{(\boldsymbol{x}, \boldsymbol{y}) \sim \text{data}}[-\log P_{\boldsymbol{\theta}_g}(\boldsymbol{y}|\boldsymbol{x})]$$

Losses to fool the discriminator:

$$C_t(\boldsymbol{\theta}_g | \boldsymbol{\theta}_d) = E_{(\boldsymbol{x}, \boldsymbol{y}) \sim \text{data}} [-\log(1 - D(B(\boldsymbol{x}, \boldsymbol{y}, \boldsymbol{\theta}_g), \boldsymbol{\theta}_d))]$$
$$C_f(\boldsymbol{\theta}_g | \boldsymbol{\theta}_d) = E_{\boldsymbol{x} \sim \text{data}, \boldsymbol{y} \sim P_{\boldsymbol{\theta}_g}(\boldsymbol{y} | \boldsymbol{x})} [-\log D(B(\boldsymbol{x}, \boldsymbol{y}, \boldsymbol{\theta}_g), \boldsymbol{\theta}_d)]$$

CONTENTS



Exposure Bias		Word-level Matching	中國糾違院计算技术研究码
Scheduled Sampling	Professor Forcing	RL-based Training Differentiable Los	INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES

Framework

- Generator
 - A sequence model
- Sequence-level Reward
 - BLEU, GLEU, etc
 - A score from a Discriminator
- Training strategy
 - Policy Gradient
 - Actor-Critic





Exposure Bias	Word-level Matchin	ng 化回翻空热计算技术研究的
Scheduled Sampling Professor Forcing	RL-based Training Differ	entiable Loss

Discriminator





Training

Loss for discriminator:

$$\mathcal{L}_d(\phi) = -\mathbb{E}_{Y \sim p_{\text{data}}}[\log D_\phi(Y)] - \mathbb{E}_{Y \sim G_\theta}[\log(1 - D_\phi(Y))]$$

• Loss for generator:

Policy Gradient

$$\mathcal{L}_g(\theta) = -\sum_{\boldsymbol{y}\in G} p(\boldsymbol{y}) \cdot r(\boldsymbol{y})$$

CONTENTS



Exposure Bias		Word-level Matching		
Scheduled Sampling	Professor Forcing	RL-based Training	Differentiable Loss	INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES

Differentiable Sequence-level Loss

• Loss:

• N-gram-based, e.g., BLEU, GLEU, etc $BLEU = BP \cdot \exp(\sum w_n \log p_n)$ n=1 Probabilistic n-gram precision N $P-BLEU = BP \cdot \exp(\sum w_n \log \tilde{p}_n)$ n=1



Probabilistic N-gram Accuracy: A 3-gram example

Reference : How should I deal with this issue? (How should I)process(this issue Out Probability: 0.6 0.7 0.9 0.3 0.5 0.6

- The probabilistic count of 3-gram "How should I" : 0.6 * 0.7 * 0.9 = 0.378
- The probabilistic count of 3-gram "this issue ?" : 0.7 * 0.9 * 0.5 = 0.315
- The total probabilistic count of 3-grams :
- Precision of 3-gram:

0.7 * 0.9 * 0.5 1.107 $\frac{0.378 + 0.315}{1.107}$



• A 2-gram example



Shao et al., EMNLP 2018



Shao et al., EMNLP 2018



References

- Wen Zhang, Yang Feng, Fandong Meng, Di You, Qun Liu. Bridging the Gap between Training and Inference for Neural Machine Translation. ACL 2019.
- Anirudh Goyal, Alex Lamb, Ying Zhang, Saizheng Zhang, Aaron Courville, Yoshua Bengio. Professor Forcing: A New Algorithm for Training Recurrent Networks. NIPS 2016.
- Lantao Yu, Weinan Zhang, Jun Wang, Yong Yu. SeqGAN: Sequence Generative Adversarial Nets with Policy Gradient. AAAI 2017.
- Chenze Shao, Yang Feng, Xilin Chen. Greedy Search with Probabilistic N-gram Matching for Neural Machine Translation. EMNLP 2018.

Thanks for your attention!