Vanguard

Exploring Fluent Query Reformulations with Text-to-Text Transformers and Reinforcement Learning

Contributions

Training text-to-text transformers (T5) with reinforcement learning (RL) to reformulate queries by adapting to F1 reward from downstream environment (question-answering, intent classification)

- Fine-tuned T5 is more sample efficient for RL than previous approach AQA, and generates more fluent reformulations.
- Use a separate T5 model to evaluate well-formedness of produced queries.
- Flexible framework adapted to other environments like intent classification when rewards can be engineered.

Methodology

Given a query q^i , a reformulation $r^i = \{r_1^i, \cdots, r_t^i\}$ is generated by Seq2Seq policy π_{θ} to compare with the target sequence $y^i = \{y_1^i, \dots, y_t^i\}$.

Supervised Fine-tuning: Log-likelihood loss $-\sum_{i=1}^{N} \log \pi_{\theta}(r^{i}|q^{i})$

$$= -\sum_{i=1}^{N} \sum_{t=1}^{T} \log p(r_t^i | r_1^i, \cdots, r_{t-1}^i, q^i)$$

$$= -\sum_{i=1}^{N} \sum_{t=1}^{T} \sum_{j=1}^{V} y_{j,t}^{i} \log p(r_{j,t}^{i} | r_{1}^{i}, \cdots, r_{t-1}^{i}, q^{i}) \text{ (i.e. cross-entropy})$$

where V is vocabulary size and $y_{i,t}^i \in \{0,1\}$ is the binary label at time t for token j and $p(r_{i,t}^{i}|\cdot)$ is the conditional probability of sampled token j at time t.

Reinforcement Learning:

$$\mathcal{J} = \sum_{i=1}^{\mathsf{N}} \mathbb{E}_{\mathbf{r}_t^i \sim \pi_ heta} (\sum_{t=1}^{\mathsf{T}} \mathsf{R}(\mathbf{r}_1^i, \cdots, \mathbf{r}_t^i))$$

where R is the enviornment's reward function. $\nabla \mathcal{J}$ can be estimated by sampling reformulation r' from π_{θ} (policy gradient):

$$abla \mathcal{J} pprox \sum_{i=1}^{b}
abla_{ heta} \mathsf{r}^{i} \log \pi_{ heta}(\mathsf{r}^{i}|\mathsf{q}^{i})(R(\mathsf{r}^{i})-B^{i})$$

which is the gradient of a weighted cross-entropy loss with sampled tokens as labels. B^{i} is the baseline for variance reduction. Additionally, scaled entropy $\lambda H(\pi_{\theta}) = \lambda \sum_{t} \sum_{i} p(r_{j,t} | r_{<t}, q^{i}) \log p(r_{j,t} | r_{<t}, q^{i})$ is added to mitigate deterministic policy updates.

SearchQA Reward: QA dataset is SearchQA with context, query, answer triplets. The reward from a fixed BiDAF QA environment is given by the character-level F1 score

$$R_{F1} = 2\frac{p \cdot r}{p+r}$$

where p is precision and r is recall between true answer and the answer produced when reformulation is used. Validation (dev) set F1 score shows how well models can generalize on unseen set of data during RL.

Jerry Zikun Chen^{1,2}, Shi Yu², Haoran Wang²

¹Department of Computer Science, University of Toronto ²Center for Analytics and Insights, The Vanguard Group

AQA Framework

- Supervised pre-training: multilingual translation (UN Parallel Corpus) and paraphrasing (Paralex) (starting point for below RL methods)
- Baseline RL method: B^i is mean reward in the batch • Value network: $B^i = f_c(\pi_{\theta}^E(\cdot|q^i))$, f_c is a two-layer critic network on top of encoder
- π_F • Self-critical training: $B^i = R(r^i_{greedy}), r^i_{greedy} \sim \pi_{\theta}$
- Unlikelihood loss: $\mathcal{L}_{UL} = \sum_{t=1}^{T} [-\alpha \sum_{c \in \mathcal{C}^t} \log(1 - p(c|\mathbf{r}_{< t})], \mathcal{C}^t = \{r_1, \cdots, r_{t-1}\} \setminus \{r_t\}$
- Fluency metric: $R_f = \frac{1}{1+H(r)}, H(r) = -\frac{\sum_{t=1}^T \log p(r_i | \mathbf{r}_{< t}, \mathbf{q})}{T}$
- Baseline method performs best for reward acquisition

opy loss)



T5 Framework

Fl

- Self-supervised pre-training: HuggingFace T5-base model pre-trained on C4 dataset
- Supervised fine-tuning: paraphrasing (Quora) and denoising (MQR) with prefix *"paraphrase: "* (starting point for below RL methods)
- Baseline RL method: B^{i} is mean reward in the batch
- Self-critical training: $B^{i} = R(r^{i}_{greedy}), r^{i}_{greedy} \sim \pi_{\theta}(\cdot | q^{i})$
- RL with T5 is more sample efficient than AQA





$$\mathbf{q}(\cdot|\mathbf{q}^i)$$

Well-formedness

- after RL can retain more fluency compared to AQA

Distribution of Predicted Well-formedness Scores on SearchQA Validation Set



Qualitative Examples

Original Text	T5 Reformula-	T5 Reformula-	AQA Reformu-	AQA Reformu-
	tions (no RL)	tions (RL)	lations (no RL)	lations (RL)
Process for One	What are the best	What is the pro-	Process for one	What is one time
time wire	methods for creat-	cess for a one time	time wire?	wire process for
	ing one time wire?	wire?		one time wire
				name name one
				time
this work ob-	What can I do to	What to do about	This work object	What is this work
ject hasnt been	solve this working	a work object	has been touched	object hasnt been
touched in two	object that hasn't	that hasn't been	in two weeks and	touched in two
weeks and is just	been touched in	touched in two	just sitting?	weeks and just
sitting there	two weeks and is	weeks and is just		sitting have just
	sitting there?	sitting there?		name?
how to reinvest for	How to fund rein-	How do I reinvest	How do you rein-	What is how
an account?	vestments to an ac-	for an account?	vest a account?	to reinvest a
	count?			account?

Comparisons on Internal Query Dataset between AQA and T5

Intent Classification

- dataset

Related Works

Yu et el. "AVA: A Financial Service Chatbot based on Deep Bidirectional Transformers" (arXiv, 2020)

25000



• Fine-tune a T5-base model on question well-formedness to quantitatively evaluate the fluency of produced reformulations (*"query wellformedness: "* prefix)

• RL hurts query fluency for both T5 and AQA models, while T5 model

• Reformulator receives reward based on true/partial intent class matching in RL, and improves pre-trained BERT intent classification accuracy by 2% on internal

• Though supervised learning can improve fluency drastically, RL may be needed to correct the course for semantic drift to preserve intent