

This work aims to translate sentences in one language into the different-style sentences in another language.

I use the Chinese-English translation with the English poetry style as a case study.

In this case, I use the attention-based NMT system to conduct the translation task and an RNN language model to do the style transfer task. I will denote the translation model as NMT and the language model as LM.

Merge function

$$p(y) = p_{NMT}(y) + \beta * p_{LM}(y)$$

Dataset

NMT: 1M news-domain sentence pairs

江泽民对杜苏率团来访表示欢迎。

jiang zemin extended his welcome to tozzoli for leading the delegation to visit china .

那些自焚惨剧的参与者受到李洪志蒙蔽,毒害,控制,最后却成为李洪志妄图达到罪恶目的的"炮灰"和牺牲品

the participants of the self - immolation tragedy have been deceived , poisoned , and controlled by li hongzhi and , in the end , served as his cannon fodder and victims in his vain attempt to achieve his criminal purpose .

LM: 230,000+ sentences after preprocessing 55,000+ English modern poems

we've shared so much , we two , | ^ a multitude of joyous things . | ^ years of wild and wonderful passions , | ^ laughter , tears . . .

goodbyes that shook the heavens , | ^ only to embrace again with a love | ^ so profound we must

have been the envy | ^ of even ancient lovers .

cut off the | ^ phone | ^ turned the | ^ music

Experiments

The number of test sentence is 747. Here is a table showing the number of influenced sentence under different settings.

β	# influenced sentence/# total sentence
0	0/747
0.1	0/747
0.2	0/747
0.5	2/747
1.0	4/747
100.0	199/747
1000.0	607/747

- $\beta = 0, 0.1, 0.2$: BLEU = 31.44

- $\beta = 0.5$: BLEU = 31.42

examples:

英国媒体今天报导,席哈克与英国首相布莱尔明天将在法国北部 le touquet 晤谈,就伊拉克危机的最新发展彼此交换意见。

w/o style: british newspapers and british prime minister blair today vowed to exchange views with the british parliament on the latest development of the iraq crisis .

w/ style: british newspapers and british prime minister blair today vowed to exchange views with the british parliament on the latest development of a crisis in iraq .

加拿大太空署主任,加拿大首位宇航员加诺表示,美方对"哥伦比亚"号事件的的调查可能使国际空间站已定下的任务受阻,有些计划可能会延期执行甚至取消。

w/o style: canada 's chief administrator of canada 's first astronaut , canada 's first astronaut , said that the investigation conducted by the us side on the " el salvador incident " may cause the implementation of the international space station , and some may be postponed or even cancelled .

w/ style: canada 's chief administrator of canada 's first astronaut , canada 's first astronaut , said that the investigation on the " el salvador incident " may affect the missions set by the international space station , and some may be postponed or even cancelled .

- $\beta = 1.0$: BLEU = 31.43

- $\beta = 100.0$: BLEU = 31.27

examples:

他同时宣布,格强力部门近期还将会继续在潘基西峡谷展开清剿非法武装的行动。

w/o style: he also announced that the turkish government would continue to carry out anti - illegal armed operations in the southern part of _UNK .

w/ style: he also announced that the turkish government would continue to carry out illegal armed operations in the _UNK valley in the near future .

届时,游客可以搭乘埃及航空公司的班机由北京起飞,在开罗游玩0天后,再飞往马耳他

w/o style: by that time , tourists can take part in the flight of the egyptian airlines flight from beijing for a couple of days to fly back to malta .

w/ style: by then , tourists can take part in the flight of the air services of egypt to fly from beijing to _UNK , and then fly to malta .

- $\beta = 1000.0$: BLEU = 29.09

examples:

去年巴西国家风险指数曾超过0000点,债券价格下降到面值的00%。

w/o style: last year , the brazilian national league index fell more than 0,000 points , and the price drop was down

w/ style: last year , the \uparrow index exceeded 0,000 points , down from 00 percent to 00 percent .

近来外国投资者认为巴西债券价格便宜,回报率高,因而需求大幅增加,使巴西国家风险指数下降加快。

w/o style: foreign investors have expressed their interest rates and the volume of crude oils has risen sharply in recent years .

w/ style: in recent years , investors have expressed interest in _UNK , _UNK , _UNK , _UNK , _UNK , _UNK , _UNK , _UNK , _UNK , _UNK , _UNK , _UNK , _UNK , _UNK , _UNK

Analysis

The main problem is that NMT dataset and LM dataset don't have many shared prefix or shared sentence segment. So LM can't generate confident result to influence the final result. That's why there aren't many influenced lines when beta is small (0.1, 0.2, 0.5, and 1.0).

As the β increases large enough, the poem phase separator " \uparrow " and "_UNK" symbols come up (because news dataset has too many domain-specific words that LM data doesn't have?)

NEXT:

1. Try to renormalize the softmax of the LM without the "_UNK" symbols to alleviate the too-many-UNK problem.
2. Handle the problem of less shared prefix.

