

baselines

2019年3月13日 14:29

		baselines	evaluations	automatic evaluations																																																																															
planning	baidu 16 COLING																																																																																		
polishing (I, Poet)	Yan 16		Perplexity BLEU human	<table border="1"> <thead> <tr> <th rowspan="2">Algo.</th> <th colspan="3">5-Character</th> <th colspan="3">7-Character</th> </tr> <tr> <th>PPL</th> <th>BLEU</th> <th>Human</th> <th>PPL</th> <th>BLEU</th> <th>Human</th> </tr> </thead> <tbody> <tr> <td>Random</td> <td>-</td> <td>0.002</td> <td>0.259</td> <td>-</td> <td>0.051</td> <td>0.135</td> </tr> <tr> <td>SMT</td> <td>126</td> <td>0.051</td> <td>1.943</td> <td>134</td> <td>0.144</td> <td>1.957</td> </tr> <tr> <td>SUM</td> <td>149</td> <td>0.035</td> <td>2.219</td> <td>131</td> <td>0.128</td> <td>2.013</td> </tr> <tr> <td>RNNPG</td> <td>103</td> <td>0.053</td> <td>1.964</td> <td>119</td> <td>0.163</td> <td>2.205</td> </tr> <tr> <td>LSTM-RNN</td> <td>123</td> <td>0.048</td> <td>1.762</td> <td>136</td> <td>0.159</td> <td>1.633</td> </tr> <tr> <td>iPoet</td> <td>91</td> <td>0.088</td> <td>2.352</td> <td>96</td> <td>0.185</td> <td>2.568</td> </tr> </tbody> </table>	Algo.	5-Character			7-Character			PPL	BLEU	Human	PPL	BLEU	Human	Random	-	0.002	0.259	-	0.051	0.135	SMT	126	0.051	1.943	134	0.144	1.957	SUM	149	0.035	2.219	131	0.128	2.013	RNNPG	103	0.053	1.964	119	0.163	2.205	LSTM-RNN	123	0.048	1.762	136	0.159	1.633	iPoet	91	0.088	2.352	96	0.185	2.568																								
Algo.	5-Character			7-Character																																																																															
	PPL	BLEU	Human	PPL	BLEU	Human																																																																													
Random	-	0.002	0.259	-	0.051	0.135																																																																													
SMT	126	0.051	1.943	134	0.144	1.957																																																																													
SUM	149	0.035	2.219	131	0.128	2.013																																																																													
RNNPG	103	0.053	1.964	119	0.163	2.205																																																																													
LSTM-RNN	123	0.048	1.762	136	0.159	1.633																																																																													
iPoet	91	0.088	2.352	96	0.185	2.568																																																																													
memory	Zhang 17 ACL	None	(all human) Compliance Fluency Theme Aesthetic Scenario Consistency Innovation Consistency																																																																																
working memory	Sun 18 EMNLP	Planning Polishing Memory human	BLEU(仅在研究memory slot使用到) (human:) Fluency Meaning Coherence Relevance Aesthetics	<table border="1"> <thead> <tr> <th></th> <th>Models</th> <th>BLEU</th> <th>PP</th> </tr> </thead> <tbody> <tr> <td rowspan="2">Quatrains</td> <td>iPoet</td> <td>0.425</td> <td>138</td> </tr> <tr> <td>WM</td> <td>1.315</td> <td>86</td> </tr> <tr> <td rowspan="2">Iambics</td> <td>iambicGen</td> <td>0.320</td> <td>262</td> </tr> <tr> <td>WM</td> <td>0.699</td> <td>72</td> </tr> <tr> <td rowspan="2">Lyrics</td> <td>lyricGen</td> <td>0.312</td> <td>302</td> </tr> <tr> <td>WM</td> <td>0.568</td> <td>138</td> </tr> </tbody> </table>		Models	BLEU	PP	Quatrains	iPoet	0.425	138	WM	1.315	86	Iambics	iambicGen	0.320	262	WM	0.699	72	Lyrics	lyricGen	0.312	302	WM	0.568	138																																																						
	Models	BLEU	PP																																																																																
Quatrains	iPoet	0.425	138																																																																																
	WM	1.315	86																																																																																
Iambics	iambicGen	0.320	262																																																																																
	WM	0.699	72																																																																																
Lyrics	lyricGen	0.312	302																																																																																
	WM	0.568	138																																																																																
CVAE	Yan 18 EMNLP	S2S, AS2S, Key+AS2S(20 16 Qixin), Memory (Jiyuan)	BLEU-1 BLEU-2 (unigrams and bigrams (p < 0.01)) Similarity (thematic consistencyy, cosin) Distinctness-n (差异性。不同ngram的比 例 n = 1 to 4) (human:) Consistency Fluency Meaning Poeticness	<table border="1"> <thead> <tr> <th rowspan="2">Model</th> <th colspan="7">Automatic Evaluation</th> </tr> <tr> <th>BLEU-1</th> <th>BLEU-2</th> <th>Sim</th> <th>Dist-1</th> <th>Dist-2</th> <th>Dist-3</th> <th>Dist-4</th> </tr> </thead> <tbody> <tr> <td>S2S</td> <td>13.8</td> <td>2.48</td> <td>14.7</td> <td>2.50</td> <td>16.2</td> <td>34.9</td> <td>50.0</td> </tr> <tr> <td>AS2S</td> <td>15.5</td> <td>2.59</td> <td>14.8</td> <td>2.30</td> <td>15.2</td> <td>31.4</td> <td>44.3</td> </tr> <tr> <td>Key-AS2S</td> <td>15.8</td> <td>1.92</td> <td>19.8</td> <td>3.00</td> <td>16.3</td> <td>33.0</td> <td>45.6</td> </tr> <tr> <td>MeM-AS2S</td> <td>16.0</td> <td>1.48</td> <td>22.0</td> <td>3.40</td> <td>51.4</td> <td>87.9</td> <td>96.8</td> </tr> <tr> <td>GAN</td> <td>17.7</td> <td>2.54</td> <td>22.5</td> <td>2.50</td> <td>16.8</td> <td>35.3</td> <td>49.6</td> </tr> <tr> <td>CVAE</td> <td>17.0</td> <td>1.73</td> <td>13.7</td> <td>4.70</td> <td>52.3</td> <td>90.6</td> <td>99.0</td> </tr> <tr> <td>CVAE-Key</td> <td>16.4</td> <td>1.83</td> <td>31.0</td> <td>4.31</td> <td>43.0</td> <td>80.6</td> <td>95.8</td> </tr> <tr> <td>CVAE-D</td> <td>18.1</td> <td>2.85</td> <td>36.3</td> <td>5.20</td> <td>59.2</td> <td>94.2</td> <td>99.8</td> </tr> </tbody> </table>	Model	Automatic Evaluation							BLEU-1	BLEU-2	Sim	Dist-1	Dist-2	Dist-3	Dist-4	S2S	13.8	2.48	14.7	2.50	16.2	34.9	50.0	AS2S	15.5	2.59	14.8	2.30	15.2	31.4	44.3	Key-AS2S	15.8	1.92	19.8	3.00	16.3	33.0	45.6	MeM-AS2S	16.0	1.48	22.0	3.40	51.4	87.9	96.8	GAN	17.7	2.54	22.5	2.50	16.8	35.3	49.6	CVAE	17.0	1.73	13.7	4.70	52.3	90.6	99.0	CVAE-Key	16.4	1.83	31.0	4.31	43.0	80.6	95.8	CVAE-D	18.1	2.85	36.3	5.20	59.2	94.2	99.8
Model	Automatic Evaluation																																																																																		
	BLEU-1	BLEU-2	Sim	Dist-1	Dist-2	Dist-3	Dist-4																																																																												
S2S	13.8	2.48	14.7	2.50	16.2	34.9	50.0																																																																												
AS2S	15.5	2.59	14.8	2.30	15.2	31.4	44.3																																																																												
Key-AS2S	15.8	1.92	19.8	3.00	16.3	33.0	45.6																																																																												
MeM-AS2S	16.0	1.48	22.0	3.40	51.4	87.9	96.8																																																																												
GAN	17.7	2.54	22.5	2.50	16.8	35.3	49.6																																																																												
CVAE	17.0	1.73	13.7	4.70	52.3	90.6	99.0																																																																												
CVAE-Key	16.4	1.83	31.0	4.31	43.0	80.6	95.8																																																																												
CVAE-D	18.1	2.85	36.3	5.20	59.2	94.2	99.8																																																																												

Perplexity (PPL). For most of the language generation research, language perplexity is a sanity check. Our first set of experiments involved intrinsic evaluation of the "perplexity" evaluation for the generated poems. Perplexity is actually an entropy based evaluation. In this sense, the lower perplexity for the poems generated, the better performance in purity for the generations, and the poems are likely to be good ones.

measure character diversity by calculating the proportion of distinctive [1,4]-grams¹² in the generated poems, where final distinctness values are normalized to [0,100].

			Ovrall																																																									
mutual RL	Sun 18 EMNLP	Base Memory huamn	(Auto:) reward score diversity and innovation ty-idf distribution topic distribution (hunam:) Fluency Coherence Meaning Overall Quality	<table border="1"> <thead> <tr> <th>Models</th> <th>Bigram Ratio</th> <th>Jaccard</th> </tr> </thead> <tbody> <tr> <td>Base</td> <td>0.126</td> <td>0.214</td> </tr> <tr> <td>Mem</td> <td>0.184</td> <td>0.183</td> </tr> <tr> <td>MRL</td> <td>0.181</td> <td>0.066</td> </tr> <tr> <td>GT</td> <td>0.218</td> <td>0.006</td> </tr> <tr> <td>SRL</td> <td>0.133</td> <td>0.146</td> </tr> <tr> <td>LMRL</td> <td>0.178</td> <td>0.085</td> </tr> <tr> <td>GMRL</td> <td>0.186</td> <td>0.075</td> </tr> <tr> <td>MRL</td> <td>0.181</td> <td>0.066</td> </tr> </tbody> </table>	Models	Bigram Ratio	Jaccard	Base	0.126	0.214	Mem	0.184	0.183	MRL	0.181	0.066	GT	0.218	0.006	SRL	0.133	0.146	LMRL	0.178	0.085	GMRL	0.186	0.075	MRL	0.181	0.066																													
Models	Bigram Ratio	Jaccard																																																										
Base	0.126	0.214																																																										
Mem	0.184	0.183																																																										
MRL	0.181	0.066																																																										
GT	0.218	0.006																																																										
SRL	0.133	0.146																																																										
LMRL	0.178	0.085																																																										
GMRL	0.186	0.075																																																										
MRL	0.181	0.066																																																										
Salient-clue	Sun 18 CoNLL	planning polishing seq2seqPG(2 017 Sun) human	BLEU (human:)	<table border="1"> <thead> <tr> <th></th> <th>Models</th> <th>Wjue</th> <th>Qjue</th> </tr> </thead> <tbody> <tr> <td rowspan="4">Different Models</td> <td>Planning</td> <td>0.460</td> <td>0.554</td> </tr> <tr> <td>iPoet</td> <td>0.502</td> <td>0.591</td> </tr> <tr> <td>seq2seqPG</td> <td>0.466</td> <td>0.620</td> </tr> <tr> <td>SC</td> <td>0.532</td> <td>0.669</td> </tr> <tr> <td rowspan="5">Different Strategies of SC</td> <td>naive-TopK-SDU</td> <td>0.442</td> <td>0.608</td> </tr> <tr> <td>naive-SSal-SDU</td> <td>0.471</td> <td>0.610</td> </tr> <tr> <td>tfidf-SSal-SDU</td> <td>0.533</td> <td>0.648</td> </tr> <tr> <td>tfidf-SSal-SSI</td> <td>0.530</td> <td>0.667</td> </tr> <tr> <td>tfidf-SSal-SSI-intent</td> <td>0.532</td> <td>0.669</td> </tr> </tbody> </table> <p>Table 1: BLEU evaluation results. The scores are calculated by the multi-bleu.perl script.</p>		Models	Wjue	Qjue	Different Models	Planning	0.460	0.554	iPoet	0.502	0.591	seq2seqPG	0.466	0.620	SC	0.532	0.669	Different Strategies of SC	naive-TopK-SDU	0.442	0.608	naive-SSal-SDU	0.471	0.610	tfidf-SSal-SDU	0.533	0.648	tfidf-SSal-SSI	0.530	0.667	tfidf-SSal-SSI-intent	0.532	0.669																							
	Models	Wjue	Qjue																																																									
Different Models	Planning	0.460	0.554																																																									
	iPoet	0.502	0.591																																																									
	seq2seqPG	0.466	0.620																																																									
	SC	0.532	0.669																																																									
Different Strategies of SC	naive-TopK-SDU	0.442	0.608																																																									
	naive-SSal-SDU	0.471	0.610																																																									
	tfidf-SSal-SDU	0.533	0.648																																																									
	tfidf-SSal-SSI	0.530	0.667																																																									
	tfidf-SSal-SSI-intent	0.532	0.669																																																									
Unsupervised	Sun 18 EMNLP	seq2seq polishing memory human	(all human) Fluency Coherence Meaningfulness Poeticness																																																									
CVAE + hybrid dec	Waterloo 18 IJCAI	AS2S	PPL NNL(KL) BLEU-1 BLEU-2 BLEU-3 BLEU-4 RES (Rhythm Score Evaluation)	<table border="1"> <thead> <tr> <th>Approach</th> <th>PPL</th> <th>NNL(KL)</th> <th>BLEU-1</th> <th>BLEU-2</th> <th>BLEU-3</th> <th>BLEU-4</th> <th>RES</th> </tr> </thead> <tbody> <tr> <td>Ground truth</td> <td>-</td> <td>-</td> <td>-</td> <td>-</td> <td>-</td> <td>-</td> <td>0.8975</td> </tr> <tr> <td>AS2S</td> <td>55.5255</td> <td>4.0168(-)</td> <td>0.4296</td> <td>0.3596</td> <td>0.3045</td> <td>0.2640</td> <td>0.5046</td> </tr> <tr> <td>AS2S+AW2V</td> <td>54.1752</td> <td>3.9922(-)</td> <td>0.4319</td> <td>0.3625</td> <td>0.3076</td> <td>0.2669</td> <td>0.5076</td> </tr> <tr> <td>CVAE</td> <td>52.1154</td> <td>3.9535(0.0083)</td> <td>0.4366</td> <td>0.3605</td> <td>0.3046</td> <td>0.2637</td> <td>0.5137</td> </tr> <tr> <td>CVAE+AW2V</td> <td>52.0516</td> <td>3.9522(0.0126)</td> <td>0.4382</td> <td>0.3630</td> <td>0.3073</td> <td>0.2663</td> <td>0.5149</td> </tr> <tr> <td>CVAE-HD+AW2V</td> <td>51.7236</td> <td>3.9459(0.0163)</td> <td>0.4405</td> <td>0.3644</td> <td>0.3084</td> <td>0.2672</td> <td>0.5245</td> </tr> </tbody> </table>	Approach	PPL	NNL(KL)	BLEU-1	BLEU-2	BLEU-3	BLEU-4	RES	Ground truth	-	-	-	-	-	-	0.8975	AS2S	55.5255	4.0168(-)	0.4296	0.3596	0.3045	0.2640	0.5046	AS2S+AW2V	54.1752	3.9922(-)	0.4319	0.3625	0.3076	0.2669	0.5076	CVAE	52.1154	3.9535(0.0083)	0.4366	0.3605	0.3046	0.2637	0.5137	CVAE+AW2V	52.0516	3.9522(0.0126)	0.4382	0.3630	0.3073	0.2663	0.5149	CVAE-HD+AW2V	51.7236	3.9459(0.0163)	0.4405	0.3644	0.3084	0.2672	0.5245
Approach	PPL	NNL(KL)	BLEU-1	BLEU-2	BLEU-3	BLEU-4	RES																																																					
Ground truth	-	-	-	-	-	-	0.8975																																																					
AS2S	55.5255	4.0168(-)	0.4296	0.3596	0.3045	0.2640	0.5046																																																					
AS2S+AW2V	54.1752	3.9922(-)	0.4319	0.3625	0.3076	0.2669	0.5076																																																					
CVAE	52.1154	3.9535(0.0083)	0.4366	0.3605	0.3046	0.2637	0.5137																																																					
CVAE+AW2V	52.0516	3.9522(0.0126)	0.4382	0.3630	0.3073	0.2663	0.5149																																																					
CVAE-HD+AW2V	51.7236	3.9459(0.0163)	0.4405	0.3644	0.3084	0.2672	0.5245																																																					

multi-bleu.perl 脚本计算bleu