Recent Advances in Natural Language Generation

Xiaojun Wan

Peking University https://wanxiaojun.github.io/

August 24, 2019

Taxonomy of NLG tasks

Creation vs. Re-Creation



WASHINGTON — Directly contradicting much of the Trump administration's position on climate change, 13 federal agencies unveiled an <u>exhaustive scientific report</u> on Friday that says humans are the dominant cause of the global temperature rise that has created the warmest period in the history of civilization.

Over the past 115 years global average temperatures have increased 1.8 degrees Fahrenheit, leading to record-breaking weather events and temperature extremes, the report says. The global, long-term warming trend is "unambiguous," it says, and there is "no convincing alternative explanation" that anything other than humans — the cars we drive, the power plants we operate, the forests we destroy — are to blame.

The report was approved for release by the White House, but the findings come as the Trump administration is defending its climate change policies. The United Nations convenes its annual climate change conference next week in Bonn, Germany, and the American delegation is expected to face harsh criticism over President Trump's decision to walk away from the 195-nation Paris climate accord and top administration officials' stated doubts about the causes and impacts of a warning planet.

	2006	2005	2004	2003	2002
Sales	\$128.3	\$97.2	\$74.6	\$61.9	\$68.
Gross Margin	\$71.0 55%	\$53.2 55%	\$40.1 54%	\$25.7 42%	\$15. 234
Selling/Admin	\$40.2	\$31.8	\$25.9	\$27.0	\$33.
Net R&D	\$15.4	\$12.2	\$12.4	\$27.0	\$12
Earnings (Loss) from Operations	\$8.0	\$(0.02)	\$(10.6)	\$(39.6)	\$(74.8
Net Earnings (Loss)	\$8.1	\$(91.6)	\$(8.4)	\$(55.0)	\$(308.
Net Earnings per Share	\$0.12	\$(0.02)	\$(0.13)	\$(0.87)	\$(5.09





Applications of NLG



News Summarization

News Writing







Product Title / Specification Generation



Description Springer

Book Writing

Applications of NLG

人于2017年1月18日正式"」

岗",并在南方都市报上推出第一篇



(消费) 2018-02-01 11:01:55



本频道为科学新闻AI平台—— "小柯" 机器人的表演场地。 "小柯" 是一个科学新闻写作机器人,由中国科学报社联合北京大学高 水平科研团队研发而成,旨在帮助科学家以中文方式快速获取全球高水平英文论文发布的最新科研讲展。本频道内的所有科学新闻均由 机器人 "小柯" 独立完成,并经过专业人士和中国科学报社编辑的双重人工审校和信息补充。 点击阅读:首个科学新闻写作机器人"小柯"问世

科学家发现延缓衰老的新途径

查看更多

多 往期回顾	最新	更多
2019/8/8	研究发现肠道菌	群与体重的关系
2019/8/8	多胺影响巨噬细胞	胞激活
2019/8/7	动态血压与死亡	率和心血管结局的关系统
2019/8/7	科学家发现组织	中的"年龄马赛克"
2019/8/7	生活方式和遗传	风险与痴呆发病率的关系

Creative Text Generation (Pun, Metaphor, Story, Poetry, …)

- Controllable Text Generation (w/ length, lexical, syntactic constraints)
- Generating Texts with Special Attributes (Sentiment and style transfer)
- Cross-Modal Text Generation (Image/Video caption & comment generation)
- Question Generation / Learning to Ask

Creative Text Generation (Pun, Metaphor, Story, Poetry, …)



Metaphors are a form of figurative language, which refers to words or expressions that mean something different from their literal definition

Target Verb: absorbed Fit Word: learn =>Our Result: he absorbed his studies at the university of birmingham . =>Gold Metaphor: he absorbed the knowledge or beliefs of his tribe .

[Yu and Wan, NAACL 2019]

Controllable Text Generation (w/ length, lexical, syntactic constraints)

X: his teammates' eyes got an ugly, hostile expression. Y: the smell of flowers was thick and sweet.

Z: the eyes of his teammates had turned ugly and hostile.

X: we need to further strengthen the agency's capacities.Y: the damage in this area seems to be quite minimal.Z: the capacity of this office needs to be reinforced even further.

Controllable Paraphrase Generation with a Syntactic Exemplar [Chen et al., ACL 2019]

Input sentence & labels	Strategy	Generated Paraphrase		
what is the value of a 1961 us cent?	All	what is the 1961 nickel 's value?		
what is the value of a 1961 us cent?	Phrase	what is the price of		
		<u>a 1961 nickel</u> ?		
what is the value of a 1961 us cent?	Sentence	what is the 1961 us cent 's		
		value?		
what is the value of a 1961 us cent?	Phrase	what is the value of		
		<u>a 1961 nickel</u> ?		
what is the value of a 1961 us cent?	Sentence	how much is <u>a 1961 us cent</u>		
		worth?		

Decomposable Neural Paraphrase Generation [Li et al., ACL 2019]



Cross-Modal Text Generation (Image/Video caption & comment generation)





高阁登临万里遥,山川水色自萧条。 江湖有约同游赏,洛浦无情共寂寥。 古寺荒芜秋雨湿,孤城寒树暮潮销。 故园不见今年晚,一夕西风吹鬓凋。

高阁亭下水面平, 寺壁秋深古殿清。 廊庙不知天地窄, 桥边谁识此时情。

[Liu et al., ACMMM 2018]

Cross-Modal Text Generation (Image/Video caption & comment generation) Attention: a person rid







Attention: a person riding a motorcycle on a road Ours (IR): a person on a motorcycle in a race Ours (Tdiv): a man riding a motorcycle down a road next to a pile of hay Human: a person on a motorcycle riding beside hay bales

Attention: a black and white photo of an airplane Ours (IR): an airplane sitting on a runway with people on the ground

Ours (Tdiv): a black and white photo of an airplane parked at an airport with people standing around

Human: an airplane with people under the wings at a field

[Liu et al., ICCV 2019]

Question Generation / Learning to Ask

Fixed-answered questions: Who invented the car? (Standard answer: Karl Benz)

Open-answered questions: What do you think of the self-driving car? (No standard answer)

OpenQG: Generating open-answered questions based on given news that are suitable to arouse open discussions.

[Chai et al., ACL 2019]

Question Analysis

- Motivation
 - The more answers a question receives in online forums, the better it is for open discussions.
 - Many factors can influence the number of answers.
- Dominated variables
 - Prior works: how language use affects the reaction that a piece of text generates.
 - Four dominated variables: topic, author, posted time & language use
 - How language use makes a question get more answers.

Control the influence of topic, author and posted time.

OQRanD

- OQGanD (Open Question Ranking Dataset).
 - Based on 11.5M open-domain questions from an in housed Zhihu database.
 - Contains 22K question-pairs. In each pair, the different in topic, time and author is controlled, and the main difference is language use.

Questions	# answers
你的家乡有什么初次尝试不太容易接受的美食么? (Is there any food that is hard to accept by foreigners for the first time in your hometown?)	1
有哪些在自己家乡很正常但在外地人眼里是黑暗料理的美食? (Which foods are normal in your hometown but are dark cuisine in the eyes of foreigners?)	89

The Effect of Language Use

- Method
 - Perform significant tests (one-sided paired t-test) on different linguistic features.
 - Find 33 features pass the significant tests.
- Some interesting features & conclusions
 - Length of questions: Ask concise questions.
 - # nouns: use less nouns to ask one topic a time.
 - **# verbs**: use more verbs to make the question vivid.
 - *#* honorifics: interact with readers naturally, do not use too many honorifics.
 - *# positive words*: questions with positive emotions are often more popular.
 - LM results: use familiar expressions. Distinctive expressions may attract attention, but "common language" can make a question better understood.

Question Evaluation Model

• Question ranking task.

- Train a score model F_s which inputs a question and outputs a score. The larger score, the more answers are expected.
- By comparing the two F_s values for each question-pair in OQGenD, we can predict which question gets more answers.
- Results
 - Statistical machine learning models: LR, RF & SVM
 N-gram word and POS features.
 - Deep learning methods: RNN & CNN word and POS embeddings.

Model	Accuracy			
WIOUCI	traditional	traditional+ours		
LR	78.61%	82.33%		
RF	81.70%	87.74%		
SVM	79.02%	87.96 %		
RNN	74.68%	-		
CNN	83.18%	_		

Question Evaluation Model

• Question ranking task.

- Train a score model F_s which inputs a question and outputs a score. The larger score, the more answers are expected.
- By comparing the two F_s value each question-pair in OQGenD, we can predict which question gets more answers.
- Results
 - Statistical machine learning models: The 33 added features can help a lot.
 - Deep learning methods.

Only use automatic-learned features.

Model	Accuracy			
WIGUEI	traditional	traditional+ours		
LR	78.61%	82.33%		
RF	81.70%	87.74%		
SVM	79.02%	87.96 %		
RNN	74.68%	-		
CNN	83.18%	-		

OQGenD

• OQGenD (Open Question Generation Dataset)

naws	最后一次世界杯,C罗和梅西谁会赢。C罗和梅西谁更强?这个问题自两							
news	人出道就争论至今。2018年俄罗斯世界杯,							
	(Who will win the last World Cup between Ronaldo and Messi? Who is stronger,							
	Ronaldo or Messi? This issue has been debated since the beginning of their							
	career. The 2018 World Cup in Russia)							
	最后一次世界杯,C罗和梅西谁会赢?							
	(Who will win the last World Cup between Ronaldo and Messi?)							
gold questions	最后一次世界杯,C罗会战胜梅西吗?							
gold questions	(Will Ronaldo defeat Messi in the last World Cup?)							
	最后一次世界杯,C罗会输给梅西吗?							
	(Will Ronaldo lose to Messi in the last World Cup?)							
	最后一次世界杯,梅西会输给C罗吗?							
	(Will Messi lose to Ronaldo in the last World Cup?)							
	最后一次世界杯,梅西会战胜C罗吗?							
	(Will Messi defeat Ronaldo in the last World Cup?)							

- Each news corresponds with more than one open-answered questions.
- In total: 9K news and 20K (news, question) pairs.

Model (Architecture)



Based on the conditional generative adversarial network

- Generator: inputs news, generates (fake) questions.
- Discriminator: input (news, questions), predicts how likely it comes from real-world dataset.

Both of the generator & discriminator are conditioned on news.

Model (Details)

- Generator G_{θ} :
 - Sequence to sequence model with attention mechanism.
- Discriminator D_{φ} :
 - Embed input (news, question) into (v_{news}, v_{ques}) .
 - High level representations:

 $oldsymbol{v}_{match} = oldsymbol{W}_m \left[oldsymbol{v}_{news};oldsymbol{v}_{ques}
ight] + oldsymbol{b}_m \ oldsymbol{v}_{fluent} = oldsymbol{W}_f oldsymbol{v}_{ques} + oldsymbol{b}_f$

- Cot the final prediction:
- Get the final prediction:

$$D_{\phi}(X, Y_D) = \sigma(\boldsymbol{W}_{proj} [\boldsymbol{v}_{match}; \boldsymbol{v}_{fluent}] + \boldsymbol{b}_{proj})$$

Input of D_{φ} : (news, real-questions), (news, fake-questions) (news, shuffled real-questions)





Object Functions

- Discriminator D_{φ} : $J_D(\phi) = -\mathbb{E}_{(X,Y)\sim P_{\text{real data}}} \log D_{\phi}(X,Y)$ $-\mathbb{E}_{(X,Y)\sim P_{\text{fake data}}} \log(1 - D_{\phi}(X,Y))$
- Generator G_{θ} :
 - Text generation is a discreate process, cannot directly use $D_{\phi}(X, Y_D)$ to train G_{θ} .
 - Use the idea of reinforcement learning.

Policy π : the generator.

State s_t : the generated text.

Action a_t : generating the next word.

Reward r_t : perform Monte-Carlo search, r

Vanilla version: only use the results from our discriminator.

$$r_t = \frac{1}{k} \sum_{i=1}^k D_{\phi}(\hat{Y}_{MC}^{(i)}, X)$$

• Get object function for generator by using policy gradient:

$$J_G(\boldsymbol{\theta}) = -\mathbb{E}\left[\sum_t r_t \cdot \log \pi(a_t|s_t)\right]$$

Object Functions

- Discriminator D_{φ} : $J_D(\phi) = -\mathbb{E}_{(X,Y)\sim P_{\text{real data}}} \log D_{\phi}(X,Y)$ $-\mathbb{E}_{(X,Y)\sim P_{\text{fake data}}} \log(1 - D_{\phi}(X,Y))$
- Generator G_{θ} :
 - Text generation is a discreate process, cannot directly use $D_{\phi}(X, Y_D)$ to train G_{θ} .
 - Use the idea of reinforcement learning.

Policy π : the generator.

State s_t : the generated text.

Action a_t : generating the next word.

Reward r_t : perform Monte-Carlo search,

Full version: add the predictions from our question evaluation model

$$r_t = \frac{1}{k} \sum_{i=1}^k (\gamma D_{\phi}(\hat{Y}_{MC}^{(i)}, X) + (1 - \gamma) F_s(\hat{Y}_{MC}^{(i)}))$$

• Get object function for generator by using policy gradient:

$$J_G(\boldsymbol{\theta}) = -\mathbb{E}\left[\sum_t r_t \cdot \log \pi(a_t|s_t)\right]$$

Experiments

Models	1	BL 2	EU 3	4	ROUGE_L	METEOR	F_s (SVM)
Seq2seq	36.35☆	20.25*	14.90*	13.22*	36.72 _◊	21.57 _{\$}	-2.28\$
CopyNet	37.89☆	21.09 [*]	15.77 _◊	14.07^{*}_{\diamond}	38.05^{*}_{\diamond}	22.63^{*}_{\diamond}	-1.80*
SeqGAN	38.51⊳	22.29 _{\$}	16.97^{*}_{\diamond}	14.92 _◊ *	38.40 _◊	23.13 _◊ *	-1.67 _{\$}
SentiGAN	37.25 _◊ *	21.52^{*}_{\diamond}	17.24*	15.60	36.85 [∗] _◊	23.57	-2.42*
Ours (vanilla)	39.67	23.62	18.01 _{\$}	16.00_{\diamond}	39.87 _◊	24.52⊳	-1.89\$
Ours (full)	39.35	23.25	18.62	16.44	39.10	24.96	-1.54

Table 5: Results for openQG. * (\diamond) denotes that our vanilla (full) model differs from the baseline significantly based on one-side paired t-test with p < 0.05.

- Traditional automatic question evaluation metrics.
- Predictions from our question evaluation model, F_s.
 The higher value, the better for open discussions.
 We use the SVM model as our F_s.

My Concerns about NLG

Evaluation

Usability

Thanks !