

# 表格数据处理

# 任务

Caption: 2019年世界主要国家经济状况

Header:	国家	GDP（万亿美元）	对外贸易额（万亿美元）
	美国	21.37	5.1
	中国	14.34	3.95
	德国	3.85	2.69
	日本	5.08	1.35

Query: 2019年美国国内生产总值比中国高多少？

# Content-Based Table Retrieval for Web Queries

**Zhao Yan<sup>†\*</sup>, Duyu Tang<sup>‡</sup>, Nan Duan<sup>‡</sup>, Junwei Bao<sup>+\*</sup>,  
Yuanhua Lv<sup>§</sup>, Ming Zhou<sup>‡</sup>, Zhoujun Li<sup>†</sup>**

<sup>†</sup>Beihang University      <sup>‡</sup>Microsoft Research, Beijing, China

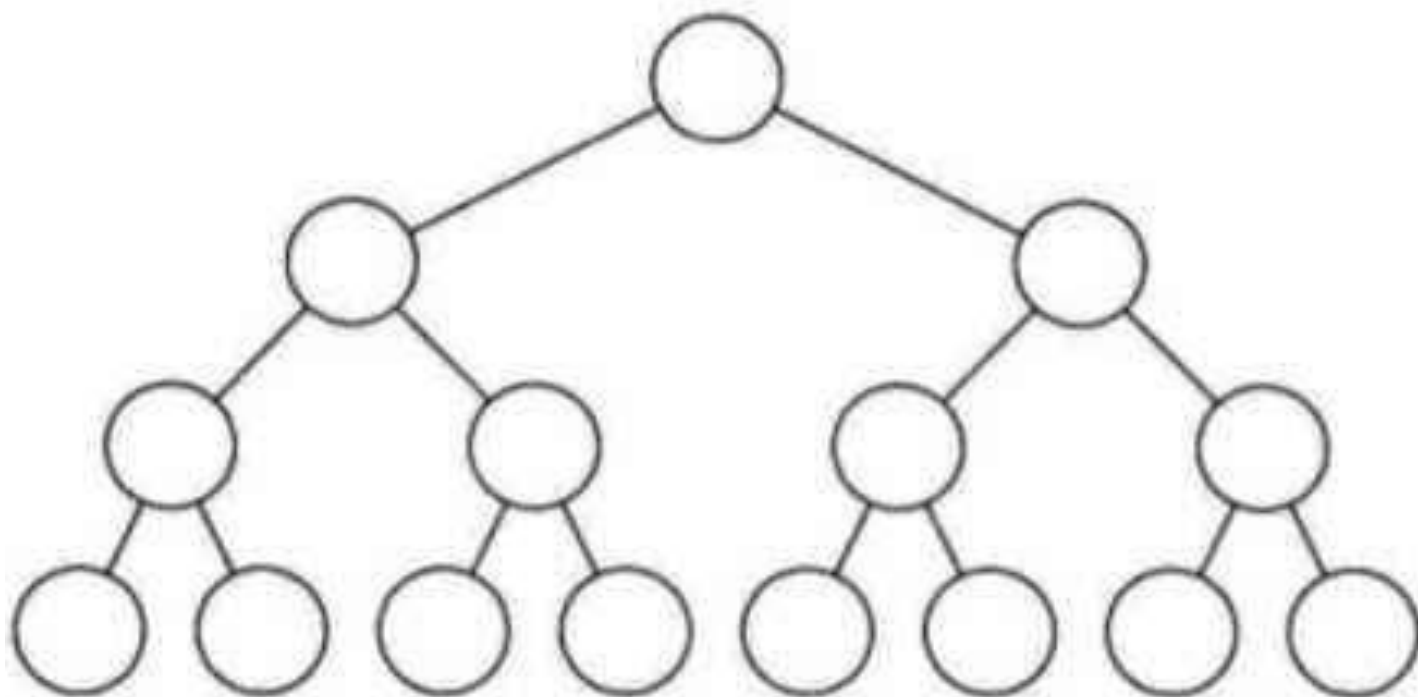
<sup>+</sup>Harbin Institute of Technology   <sup>§</sup>Microsoft AI and Research, Sunnyvale CA, USA

<sup>†</sup>{yanzhao, lizj}@buaa.edu.cn      <sup>+</sup>baojunwei001@gmail.com

<sup>‡§</sup>{dutang, nanduan, yuanhual, mingzhou}@microsoft.com

# 用决策树计算表格和问题的相关性得分

- 每一个分支选择由一个特征决定



- 得分：

# 基于规则的特征

- 1. 词级别

$$f_{wmt}(t_a, q) = \frac{\sum_{w \in t_a} \delta(w, q) \cdot idf(w)}{\sum_{w' \in t_a} idf(w')}$$

$$f_{wmq}(t_a, q) = \frac{\sum_{w \in t_a} \delta(w, q) \cdot idf(w)}{\sum_{w' \in q} idf(w')}$$

- ldf: 词频

- $\delta(w, q)$ :  $w$ 出现在 $q$ 中, 值为1; 未出现, 值为0

- 2. 短语级别

$$f_{pp}(t_a, q) = \frac{1}{N} \sum_{n=1}^N \frac{\sum_{i,j} \text{score}(\text{src}_{i,n}^{t_q}, \text{src}_{j,n}^q)}{|t_a| - N + 1}$$

$$\text{score}(\text{src}_x; \text{src}_y) = \sum_{PT} p(\text{tgt}_k | \text{src}_x) \cdot p(\text{src}_y | \text{tgt}_k)$$

- 3. 句子级别

$$f_{s1}(t_a, q) = \text{cosine}(\text{cdssm}(t_a), \text{cdssm}(q))$$

# 基于神经网络的特征

$$\alpha_i = \frac{\exp(\tanh(W[m_i; v_q] + b))}{\sum_{j=1}^k \exp(\tanh(W[m_j; v_q] + b))}$$

- $V_{\text{header}} = \sum_{i=1}^k \alpha_i m_i$

$$f_{nn}(\text{header}, q) = NN_1(M_h, v_q)$$



# 实验数据

	WQT数据集 (WebQueryTable)	WTQ数据集 (WikiTableQuestions)
表格数	273,816	2,108
平均列数	4.55	6.38
最大列数	52	25
最小列数	1	3
平均行数	9.15	28.50
最大行数	1,517	754
最小行数	2	5
问题的数量	21,113	22,033
问题平均长度	4.61	11.25

# 实验结果

Setting	MAP	P@1
BM25	58.23	47.12
Feature	61.02	47.79
NeuralNet	61.94	49.02
Feature + NeuralNet	67.18	54.15

Setting	Feature		NeuralNet	
	MAP	P@1	MAP	P@1
Header (H)	22.39	9.76	26.03	13.35
Cell (Cel)	28.85	14.95	27.47	12.92
Caption (Cap)	57.12	56.83	60.16	48.48
H + Cel	31.99	17.08	30.73	16.25
H + Cel + Cap	61.02	47.79	61.94	49.02

ACL2020

**TABERT: Pretraining for Joint Understanding of  
Textual and Tabular Data**

**Pengcheng Yin\***   **Graham Neubig**

Carnegie Mellon University

{pcyin,gneubig}@cs.cmu.edu

**Wen-tau Yih**   **Sebastian Riedel**

Facebook AI Research

{scottyih,sriedel}@fb.com

# 作者主要贡献（为什么不使用已有的Bert）

- 1.Bert只见过自由格式的文本数据，没见过结构化、半结构化的表格数据 -> 预训练
- 2.表格中可能存在大量无关数据 -> 内容快照
- 3.不同领域内的语义解析难以统一 -> 从问题和表格中获取信息

# 内容快照 (k=3)

- 除了表头， 还使用单元格信息

Caption: 2019年世界主要国家经济状况

Header:

国家 text	GDP (万亿美元) real	对外贸易额 (万亿美元) real
---------	-----------------	-------------------

Cells:

美国	21.37	5.1
中国	14.34	3.95
德国	3.85	2.69
日本	5.08	1.35

Query: 2019年美国国内生产总值比中国高多少?

# N(tri)-gram overlap

S1: november

S2: december

- Trigram:

S1: nov ove vem emb mbe ber

S2: dec ece cem emb mbe ber

$S1 \cup S2 = 3$

$S1 \cap S2 = 9$

overlap = 3/9

# 内容快照 (k=1)

- 寻找最相关的单元格值， 拼接成行

国家 text	GDP (万亿美元) real	对外贸易额 (万亿美元) real
美国	14.34	5.1

# 行线性化

Encoding:

- [CLS] 2019 ... [SEP] 国家 | text | 美国 [SEP] GDP | real | 21.37 ...



Transformer (Bert)



# 垂直自注意机制

- 聚合不同行之间的信息
- 垂直： 同一列

美国 <code>avg('美','国').coding</code>	21.37	5.1
中国	14.34	3.95
日本	5.08	1.35

美国.code	中国.code	日本.code
---------	---------	---------



Transformer (Bert)



MeanPooling



列表示

# 预训练

- 2600万组 表格+问题

国家 text	【GDP（万亿美元）real】 尝试恢复（20%）	对外贸易额（万亿美元） real
美国	21.37（尝试恢复）	5.1
中国	14.34（尝试恢复）	3.95
德国	3.85	2.69
日本	5.08（尝试恢复）	1.35

# 实验

- 1. SPIDER 有监督训练 结构化数据 可能涉及多个表
- 2. WIKI TABLE QUESTIONS 弱监督（只知道对错） 半结构化数据 可能涉及单个表内跨行推理

# 实验结果

<i>Previous Systems on WikiTableQuestions</i>		
Model	DEV	TEST
Pasupat and Liang (2015)	37.0	37.1
Neelakantan et al. (2016)	34.1	34.2
Ensemble 15 Models	37.5	37.7
Zhang et al. (2017)	40.6	43.7
Dasigi et al. (2019)	43.1	44.3
Agarwal et al. (2019)	43.2	44.1
Ensemble 10 Models	–	46.9
Wang et al. (2019b)	43.7	44.5

<i>Our System based on MAPO (Liang et al., 2018)</i>				
	DEV	Best	TEST	Best
Base Parser <sup>†</sup>	42.3 $\pm$ 0.3	42.7	43.1 $\pm$ 0.5	43.8
<i>w/</i> BERT <sub>Base</sub> (K = 1)	49.6 $\pm$ 0.5	50.4	49.4 $\pm$ 0.5	49.2
– content snapshot	49.1 $\pm$ 0.6	50.0	48.8 $\pm$ 0.9	50.2
<i>w/</i> TABERT <sub>Base</sub> (K = 1)	51.2 $\pm$ 0.5	51.6	50.4 $\pm$ 0.5	51.2
– content snapshot	49.9 $\pm$ 0.4	50.3	49.4 $\pm$ 0.4	50.0
<i>w/</i> TABERT <sub>Base</sub> (K = 3)	51.6 $\pm$ 0.5	52.4	51.4 $\pm$ 0.3	51.3
<i>w/</i> BERT <sub>Large</sub> (K = 1)	50.3 $\pm$ 0.4	50.8	49.6 $\pm$ 0.5	50.1
<i>w/</i> TABERT <sub>Large</sub> (K = 1)	51.6 $\pm$ 1.1	52.7	51.2 $\pm$ 0.9	51.5
<i>w/</i> TABERT <sub>Large</sub> (K = 3)	<b>52.2</b> $\pm$ 0.7	<b>53.0</b>	<b>51.8</b> $\pm$ 0.6	<b>52.3</b>

<i>Top-ranked Systems on Spider Leaderboard</i>	
Model	DEV. ACC.
Global-GNN (Bogin et al., 2019a)	52.7
EditSQL + BERT (Zhang et al., 2019a)	57.6
RatSQL (Wang et al., 2019a)	60.9
IRNet + BERT (Guo et al., 2019)	60.3
+ Memory + Coarse-to-Fine	61.9
IRNet V2 + BERT	63.9
RyanSQL + BERT (Choi et al., 2020)	<b>66.6</b>

<i>Our System based on TranX (Yin and Neubig, 2018)</i>		
	Mean	Best
<i>w/</i> BERT <sub>Base</sub> (K = 1)	61.8 $\pm$ 0.8	62.4
– content snapshot	59.6 $\pm$ 0.7	60.3
<i>w/</i> TABERT <sub>Base</sub> (K = 1)	63.3 $\pm$ 0.6	64.2
– content snapshot	60.4 $\pm$ 1.3	61.8
<i>w/</i> TABERT <sub>Base</sub> (K = 3)	63.3 $\pm$ 0.7	64.1
<i>w/</i> BERT <sub>Large</sub> (K = 1)	61.3 $\pm$ 1.2	62.9
<i>w/</i> TABERT <sub>Large</sub> (K = 1)	64.0 $\pm$ 0.4	64.4
<i>w/</i> TABERT <sub>Large</sub> (K = 3)	<b>64.5</b> $\pm$ 0.6	<b>65.2</b>

# 实验结果

Cell Linearization Template	WIKIQ.	SPIDER
Pretrained TABERT <sub>Base</sub> Models (K = 1)		
<u>Column Name</u>	49.6 $\pm$ 0.4	60.0 $\pm$ 1.1
<u>Column Name</u>   <u>Type</u> <sup>†</sup> (—content snap.)	49.9 $\pm$ 0.4	60.4 $\pm$ 1.3
<u>Column Name</u>   <u>Type</u>   <u>Cell Value</u> <sup>†</sup>	51.2 $\pm$ 0.5	63.3 $\pm$ 0.6
BERT <sub>Base</sub> Models		
<u>Column Name</u> (Hwang et al., 2019)	49.0 $\pm$ 0.4	58.6 $\pm$ 0.3
<u>Column Name</u> is <u>Cell Value</u> (Chen19)	50.2 $\pm$ 0.4	63.1 $\pm$ 0.7

Table4(statement\_id: 75)

**1644 is in panel 1: 2005/2006–2009/2010 in outsourcing (t0)**

statement\_type: refuted

Predictors of outsourcing.

	Panel 1: 2005/2006–2009/2010			Panel 2: 2009/2010–2013/2014		
	Outsourcing (t0)	Outsourcing (t+1)	Outsourcing (t+2)	Outsourcing (t0)	Outsourcing (t+1)	Outsourcing (t+2)
Financial constraints (t0)	1.224 (0.275)	1.528 (0.422)	1.713 (0.583)	1.554 * (0.389)	2.260 *** (0.652)	3.290 *** (1.156)
Competence constraints (t0)	1.876 ** (0.465)	2.182 *** (0.646)	1.967 * (0.718)	1.806 ** (0.493)	0.998 (0.346)	1.281 (0.517)
Firm-size	1.131 (0.110)	1.035 (0.116)	0.899 (0.120)	1.241 * (0.136)	1.249 * (0.167)	1.150 (0.171)
Firm-type fixed effects	Yes	Yes	Yes	Yes	Yes	Yes
Industry fixed effects	Yes	Yes	Yes	Yes	Yes	Yes
Year fixed effects	Yes	Yes	Yes	Yes	Yes	Yes
Model $\chi^2$	85.61	61.99	43.33	60.40	52.27	42.03
p	0.00	0.00	0.00	0.00	0.00	0.00
n	2744	2109	1518	2123	1644	1199