Learning How to Ask: Querying LMs with Mixtures of Soft Prompts

Guanghui Qin and Jason Eisner

Johns Hopkins University



LMs Already Did Your Job



LMs Already Did Your Job





What year did Mary Cassatt die?



LMs Already Did Your Job



What year did Mary Cassatt die?

The Form of the Question Matters



- Should employers be **forced** to negotiate with unions?
- Should unions have the **right** to negotiate with employers?



Opinion Poll

The Form of the Question Matters



• Should employers be **forced** to negotiate with unions?

Should unions have the **right** to negotiate with employers?





Opinion Poll

- Why **can't we** see the moon now?
- Where does **the moon** go during the day?





Children's Mind

The Form of the Question Matters



Opinion Poll



We're prompting LMs!















Task: Fact Queries: Year-of-Death













Cab Calloway X played until his death in [MASK] y .

• Easy to search with backprop.

Cab Calloway _X played until his death in [MASK] _y.

- Easy to search with backprop.
- We have a larger space of prompts.

Cab Calloway _X played until his death in [MASK] y .

- Easy to search with backprop.
- We have a larger space of prompts.

Mary Cassatt _X played until his death in [MASK] y .

- Easy to search with backprop.
- We have a larger space of prompts.



- Easy to search with backprop.
- We have a larger space of prompts.



- Easy to search with backprop.
- We have a larger space of prompts.



- Easy to search with backprop.
- We have a larger space of prompts.



- Easy to search with backprop.
- We have a larger space of prompts.
- They can emphasize certain keywords, even particular dimensions.



- Easy to search with backprop.
- We have a larger space of prompts.
- They can emphasize certain keywords, even particular dimensions.



- Easy to search with backprop.
- We have a larger space of prompts.
- They can emphasize certain keywords, even particular dimensions.



- Easy to search with backprop.
- We have a larger space of prompts.
- They can emphasize certain keywords, even particular dimensions.



















• Bigger model: When one prompt is unsure, others can help



- Bigger model: When one prompt is unsure, others can help
- Better predictions: Ensembling reduces variance



- Bigger model: When one prompt is unsure, others can help
- Better predictions: Ensembling reduces variance
- Better optimization: Explore multiple starting points in parallel



Main Experiments

- Predict factual relations from T-REx dataset by prompting BERT-large
 - About 1000 training examples per relation



Main Experiments

- Predict factual relations from T-REx dataset by prompting BERT-large
 - About 1000 training examples per relation
- Initialize at other researchers' prompts -- huge improvement!



Main Experiments

- Predict factual relations from T-REx dataset by prompting BERT-large
 - About 1000 training examples per relation
- Initialize at other researchers' prompts -- huge improvement!
- Initialize randomly -- almost as good!



Lots of Experiments

IM	Method	Precision@1					Precision@10						MRR			
	Method	init	\rightarrow	soft	\rightarrow	deep	init	\rightarrow	soft	\rightarrow	deep	init	\rightarrow	soft	\rightarrow	deep
	LAMA	31.1					59.5					40.3				
	LPAQA	34.1					62.0					43.6				
BEb	Soft (sin.)	31.1	+14.6?	í́→ 45.′́	7 + 2.	<u>°</u> → 47.7	59.5	+16.3	, → 75.8	+ 3.	² → 79.0	40.3	+15.9	[?] → 56.2	2 + 2	$\xrightarrow{2}$ 58.4
DEU	Soft (min.)	34.1	+14.7?	, → 48.	8 + 1.	[≞] → 50.7 ?	62.0	+15.6	, → 79.6	+ 1.	±→ 80.7?	43.6	+15.8	°→ 59.4	1 + 1	. 7 → 61.1 ?
	Soft (par.)	34.1	$+12.8^{?}$	→ 46.	9 <u>+ 1</u> .	⁵ → 48.4	62.0	+16.8	, → 78.8	+ 0.	⁸ → 79.6	43.6	+14.2	[?] → 57.8	3 + 1	. <u>³</u> 59.1
	Soft (ran.)	0.7	+46.6	→ 47.3	3 + 0.8	^s → 48.1	4.6	+74.0	→ 79.1	+ 0.0	^D → 79.1	2.3	+56.1	→ 58. 4	+ 0.	^₅ 58.9
	LAMA	28.9°	i				57.7	-				38.7	-			
	LPAQA	39.4	ŕ				67.4	ŕ				49.1 [†]	İ			
DEI	Soft (sin.)	28.9	+16.9	45.8	+ 5.3	→ 51.1	57.7	+19.0	76.7	+ 4.4	→ 81.1	38.7	+17.8	÷ 56.5	+ 5.0	9→ 61.5
DEI	Soft (min.)	39.4	+11.6	51.0	+ 0.6	→ 51.6	67.4	+14.0	81.4	+ 0.5	→ 81.9	49.1	+12.5	→ 61.6	+ 0.5	⁵ → 62.1
	Soft (par.)	39.4	<u>+ 9.2</u>	48.6	+ 2.5	51.1	67.4	+12.6	80.0	+ 1.7	→ 81.7	49.1	+10.5	→ 59.6	+ 2.1	→ 61.7
	Soft (ran.)	2.3	+47.1	. 49.4	+ 1.9	→ 51.3	8.0	+73.0	81.0	+ 0.7	→ 81.7	4.5	+55.9	→ 60.4	+ 1.5	9→ 61.9
	LPAQA	1.2	ŕ				9.1	-				4.2	-			
Rob	AutoPrompt	40.0					68.3					49.9				
	Soft (min.)	1.2	+39.4	40.6	- 7.3	→ 33.2	9.1	+66.3	75.4	-22.8	÷ 53.0	4.2	+48.8	→ 53.0	-12.	$\stackrel{1}{\rightarrow} 40.8$
BAb	LPAQA	0.8	i				5.7					2.9	Ī			
	Soft (min.)	0.8	+39.1	39.9			5.7	+69.7	75.4			2.9	+49.2	→ 52.1		
DA1	LPAQA	3.5	i				5.6					4.8				
DAI	Soft (min.)	3.5	+22.3	25.8			5.6	+62.4	68.0			4.8	+36.2	• 41.0		

IM	Method	Precision@1			Precision@10				MRR							
LIVI		init	\rightarrow	soft	\rightarrow	deep	init	\rightarrow	soft	\rightarrow	deep	init	\rightarrow	soft	\rightarrow	deep
	LAMA	26.4					54.3					35.8				
	LPAQA	31.2					57.3					39.9				
BEb	Soft (sin.)	26.4	+22.2	, → 48.6	5 + 1.0	<u>°</u> → 49.6	54.3	$+23.3^{5}$	→ 77.6	• + •.:	³ → 77.9	35.8	$+22.9^{5}$	²→ 58.7	+ 0.	⁶ → 59.3
DEU	Soft (min.)	31.2	+19.0	, → 50.2	2 + 0.:	<u>³</u> → 50.5 ?	57.3	$+21.9^{3}$	→ 79.2	+ 0.	5→ 79.7 ?	39.9	$+20.2^{2}$	$a \rightarrow 60.1$	+ 0.	4 → 60.5 ?
	Soft (par.)	31.2	+18.5	→ 49.7	7 <u>+ 0.0</u>	^₀ → 49.7	57.3	$+21.3^{5}$	→ 78.6	b + 0.0	⁶ → 79.2	39.9	+19.6	²→ 59.5	<u>+ 0</u> .	³ → 59.8
	Soft (ran.)	0.8	+46.3	→ 47 .	1 + 3.	^₅ 50.6	4.0	+70.4	→ 74.4	+ 4.	⁹ → 79.3	2.2	+54.3	→ 56.5	5 <u>+ 3</u>	$\xrightarrow{9} 60.4$
	LAMA	24.0					53.7	İ				34.1				
	LPAQA	37.8	Ī				64.4	t				44.0^{\dagger}				
DEI	Soft (sin.)	24.0	+26.2	50.2	+ 1.2	→ 51.4	53.7	+24.9	78.6	+ 0.9	→ 79.5	34.1	+25.9	60.0	+ 1.2	→ 61.2
DEI	Soft (min.)	37.8	+13.4	51.2	+ 1.3	→ 52.5	64.4	+15.1	79.5	+ 1.6	→ 81.1	44.0	+17.0	61.0	+ 1.4	→ 62.4
	Soft (par.)	37.8	+12.5	50.3	+ 1.4	→ 51.7	64.4	+14.3	78.7	+ 2.1	× 80.8	44.0	+16.1	60.1	+ 1.6	→ 61.7
	Soft (ran.)	1.4	+46.1	47.5	+ 4.4	→ 51.9	5.4	+68.9	74.3	+ 6.3	→ 80.6	5.7	+51.2	56.9	+ 5.0	→ 61.9

Model	P@1	P@10	MRR
lama (BEb)	0.1†	2.6^{\dagger}	1.5^{\dagger}
lama (BEl)	0.1^{\dagger}	5.0^{+}	1.9^{\dagger}
Soft (min.,BEb)	11.3(+11.2	2) 36.4(+33.8)	19.3(+17.8)
Soft (ran.,BEb)	11.8(+11.8	B) 34.8 (+31.9)	19.8(+19.6)
Soft (min.,BEl)	12.8(+12.7	7) 37.0 (+32.0)	20.9(+19.0)
Soft (ran.,BEl)	14.5(+14.5	5) 38.6 (+34.2)	22.1(+21.9)

•	Model	P@1	P@10	MRR
	LAMA	9.7 [†]	27.0^{\dagger}	15.6 [†]
	LPAQA	10.6^{\dagger}	23.7 [†]	15.3†
	Soft (sin.)	11.2 (+1.5)	33.5 (+ 6.5)	18.9 (+3.3)
	Soft (min.)	12.9 (+2.3)	34.7 (+11.0)	20.3 (+5.0)
	Soft (par.)	11.5 (+0.9)	31.4 (+ 7.7)	18.3 (+3.0)

Model	P@ 1	P@10	MRR
lpaqa (BEb)	18.9	40.4	26.6
Soft (BEb)	23.0 (+4.1	1) 45.2 (+4.8)) 30.5 (+3.9)
LPAQA (BEI)	23.8	47.7	32.2
Soft (BEl)	27.0 (+3.2	2) 51.7 (+4.0)	35.4 (+3.2)

Model	P@1	P@10	MRR
baseline	39.4	67.4	49.1
adjust mixture weights	40.0	69.1	53.3
adjust token vectors	50.7	80.7	61.1
adjust both	51.0	81.4	61.6

Related Work

- Jiang, Zhengbao, et al. "How can we know what language models know?." TACL (2020).
- Shin, Taylor, et al. "AutoPrompt: Eliciting knowledge from language models with automatically generated prompts." *EMNLP* (2020).
- Haviv, Adi, Jonathan Berant, and Amir Globerson. "BERTese: Learning to speak to BERT." *EACL* (2021).
- Li, Xiang Lisa, and Percy Liang. "Prefix-Tuning: Optimizing Continuous Prompts for Generation." *ACL* (2021).
- Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. GPT understands, too. *arXiv* (2021).



• LMs know more facts than we thought. You just have to learn how to ask.



• LMs know more facts than we thought. You just have to learn how to ask.



• Prompts are made of vectors, not words! So you can tune them with backprop.



• LMs know more facts than we thought. You just have to learn how to ask.



• Prompts are made of vectors, not words! So you can tune them with backprop.



• Random initialization works fine. No grad student required.



• LMs know more facts than we thought. You just have to learn how to ask.



• Prompts are made of vectors, not words! So you can tune them with backprop.



• Random initialization works fine. No grad student required.



 Prompt tuning is lightweight, and could also be applied to few-shot learning. Thanks!