

Fantastically Ordered Prompts: and Where to find them: Overcoming Few-Shot Prompt Order Sensitivity

Scalable Data Mining Term Project





What aim does Prompt Learning Solve and What do Prompt Engineers do?

Input

Review: The greatest musicians.
Sentiment: Positive

Review: Redundant concepts.
Sentiment: Negative

Review: Fantastic movie!
Sentiment: ???

Output

Negative

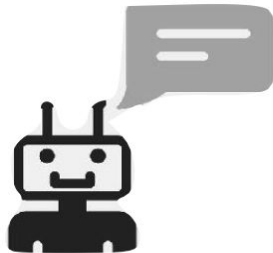
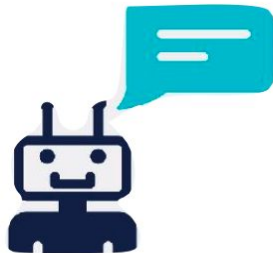


Preliminary Approaches we may Try?

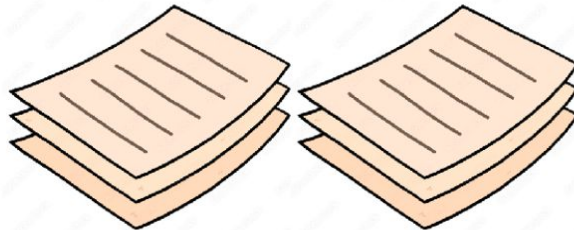
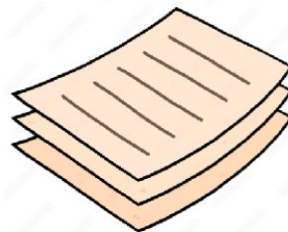
Increase Model Size



Change Wordings



Try to acquire more examples





Purpose of the Research Paper

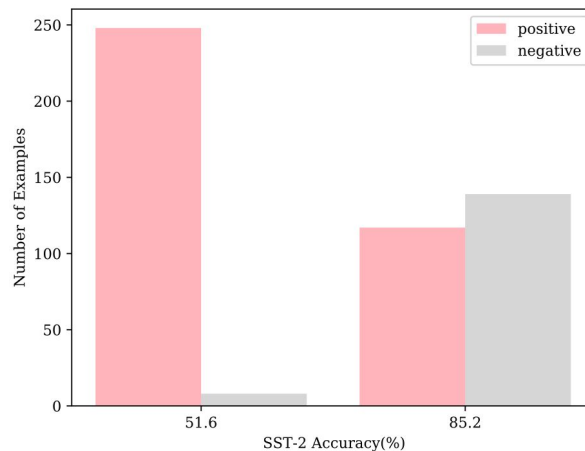
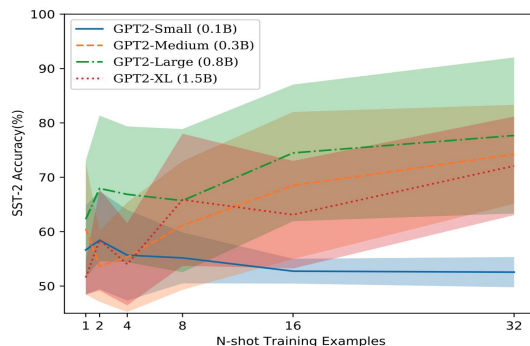
- To study the order sensitivity for in-context learning, which we show is crucial for the success of pre trained language models for few-shot learning.
- Proposing a simple, generation-based probing method to identify prompts without requiring additional data.
- The probing method must be universally applicable and effective across different sizes of pretrained language models and different types of dataset. The paper concludes a average of 13% relative improvement over a wide range of tasks.



Major Findings on Prompt Order Sensitivity

Findings:-

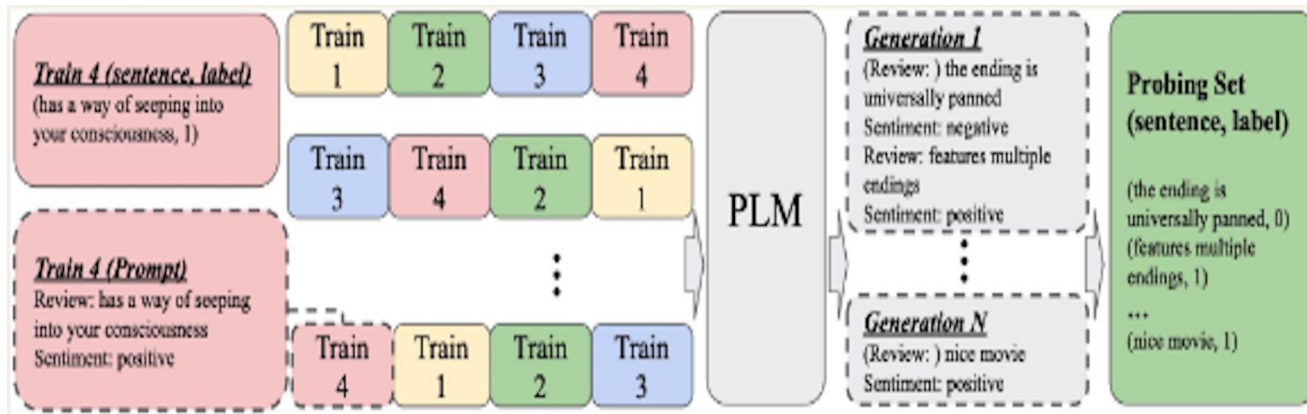
1. Although beneficial, increasing model sizes does not guarantee low variance.
2. Adding training sample does not significantly resize variance.
3. Performant prompts are not transferable across models and performant label orderings are not consistent across various models used.
4. Degenerate behavior of bad prompts.



Methodology Advocated

Probing Set Construction:-

The Figure below depicts the construction of probing set, first we consider all the possible permutations of randomly N-selected samples, the resulting generation of each permutation and concatenating them into a probing set, moreover we discard the generation labels, as there is no guarantee that the assigned generated labels are correct.





Methodology Advocated

Ranking Metrics:-

- **Global Entropy (GlobalE):** The motivation behind GlobalE is to identify prompts of specific sample orderings that overshadow the issue of *extremely unbalanced predictions*.
- **Local Entropy (LocalE):** The LocalE intends to check if a model is overly confident for all / any probing inputs (in such cases undesirable results may be observed). At the least, it turns out to be poorly calibrated, which may point out the poor ability to point out between different classes.



Implementation plan

- After Selection of the Candidate models to be used for Experimentation. After considering monetary costs (and other factors) we report the mean and Variance of the corresponding metrics over different sets.
- For Prompt Selection, we rank candidates using GlobalE and LocalE probing metrics over automatically generated probing sets we select the K - best performing as performant sets.
- Performant prompt selection can be performed on smaller size models with few shot setting to achieve In-context learning.



Inferences:

- Entropy-based probing is validated for selection of performant prompt as all the selected prompts lead to decreased variance.
- Choosing K prompts in particular is robust because of its stability across all models along with good performance.
- Entropy based probing is not sensitive to specific templates.
- Use of probing sets when applied along with entropy-based probing outperforms development kits .



Experiment Data

Way in which we wish for Experiment Data to be Displayed?

We have the LocalE and GlobalE of various data tasks and in this manner observing variation may be seen easily (observe consistent improvements).

	SST-2	SST-5	DBPedia	MR	CR	MPQA	Subj	TREC	AGNews	RTE	CB
Majority	50.9	23.1	9.4	50.0	50.0	50.0	50.0	18.8	25.0	52.7	51.8
Finetuning (Full)	95.0	58.7	99.3	90.8	89.4	87.8	97.0	97.4	94.7	80.9	90.5
GPT-2 0.1B	58.9 _{7.8}	29.0 _{4.9}	44.9 _{9.7}	58.6 _{7.6}	58.4 _{6.4}	68.9 _{7.1}	52.1 _{0.7}	49.2 _{4.7}	50.8 _{11.9}	49.7 _{2.7}	50.1 _{1.0}
LocalE	65.2 _{3.9}	34.4 _{3.4}	53.3 _{4.9}	66.0 _{6.3}	65.0 _{4.4}	72.5 _{6.0}	52.9 _{1.3}	48.0 _{3.9}	61.0 _{5.9}	53.0 _{1.3}	49.9 _{1.6}
GlobalE	63.8 _{5.8}	35.8 _{2.0}	56.1 _{4.3}	66.4 _{5.8}	64.8 _{2.7}	73.5 _{4.5}	53.0 _{1.3}	46.1 _{3.7}	62.1 _{5.7}	53.0 _{1.0}	50.3 _{1.6}
Oracle	73.5 _{1.7}	38.2 _{1.0}	60.5 _{4.2}	74.3 _{4.9}	70.8 _{4.4}	81.3 _{2.5}	55.2 _{1.7}	58.1 _{4.3}	70.3 _{2.8}	56.8 _{2.0}	52.1 _{1.3}
GPT-2 0.3B	61.0 _{3.2}	25.9 _{5.9}	51.7 _{7.0}	54.2 _{7.8}	56.7 _{9.4}	54.5 _{8.8}	54.4 _{7.9}	52.6 _{4.9}	47.7 _{10.6}	48.8 _{2.6}	50.2 _{5.3}
LocalE	75.3 _{4.6}	31.0 _{3.4}	47.1 _{3.7}	65.2 _{6.6}	70.9 _{3.3}	67.6 _{7.2}	66.7 _{9.3}	53.0 _{3.9}	51.2 _{7.3}	51.8 _{1.0}	47.1 _{4.2}
GlobalE	78.7 _{5.2}	31.7 _{5.2}	58.3 _{5.4}	67.0 _{5.9}	70.7 _{6.7}	68.3 _{6.9}	65.8 _{10.1}	53.3 _{4.6}	59.6 _{7.2}	51.1 _{1.9}	50.3 _{3.7}
Oracle	85.5 _{4.3}	40.5 _{6.3}	65.2 _{7.6}	74.7 _{6.1}	80.4 _{5.4}	77.3 _{2.3}	79.4 _{2.4}	63.3 _{2.9}	68.4 _{8.0}	53.9 _{1.3}	62.5 _{7.4}
GPT-2 0.8B	74.5 _{10.3}	34.7 _{8.2}	55.0 _{12.5}	64.6 _{13.1}	70.9 _{12.7}	65.5 _{8.7}	56.4 _{9.1}	56.5 _{2.7}	62.2 _{11.6}	53.2 _{2.0}	38.8 _{8.5}
LocalE	81.1 _{5.5}	40.3 _{4.7}	56.7 _{7.5}	82.6 _{4.2}	85.4 _{3.8}	73.6 _{4.8}	70.4 _{4.2}	56.2 _{1.7}	62.7 _{8.1}	53.3 _{1.6}	38.4 _{5.2}
GlobalE	84.8 _{4.1}	46.9 _{1.1}	67.7 _{3.6}	84.3 _{2.9}	86.7 _{2.5}	75.8 _{3.1}	68.6 _{6.5}	57.2 _{2.3}	70.7 _{3.6}	53.5 _{1.5}	41.2 _{4.5}
Oracle	88.9 _{1.8}	48.4 _{0.7}	72.3 _{3.3}	87.5 _{1.1}	89.9 _{0.9}	80.3 _{4.9}	76.6 _{4.1}	62.1 _{1.5}	78.1 _{1.3}	57.3 _{1.0}	53.2 _{5.3}
GPT-2 1.5B	66.8 _{10.8}	41.7 _{6.7}	82.6 _{2.5}	59.1 _{11.9}	56.9 _{9.0}	73.9 _{8.6}	59.7 _{10.4}	53.1 _{3.3}	77.6 _{7.3}	55.0 _{1.4}	53.8 _{4.7}
LocalE	76.7 _{8.2}	45.1 _{3.1}	83.8 _{1.7}	78.1 _{5.6}	71.8 _{8.0}	78.5 _{3.6}	69.7 _{5.8}	53.6 _{3.1}	79.3 _{3.7}	56.8 _{1.1}	52.6 _{3.9}
GlobalE	81.8 _{3.9}	43.5 _{4.5}	83.9 _{1.8}	77.9 _{7.7}	73.4 _{4.0}	81.4 _{2.1}	70.9 _{6.0}	55.5 _{3.0}	83.9 _{1.2}	56.3 _{1.2}	55.1 _{4.6}
Oracle	86.1 _{1.5}	50.9 _{1.0}	87.3 _{1.5}	84.0 _{2.7}	80.3 _{3.3}	85.1 _{1.4}	79.9 _{5.7}	59.0 _{2.3}	86.1 _{0.7}	58.2 _{0.6}	63.9 _{4.3}
GPT-3 2.7B	78.0 _{10.7}	35.3 _{6.9}	81.1 _{1.8}	68.0 _{12.9}	76.8 _{11.7}	66.5 _{10.3}	49.1 _{2.9}	55.3 _{4.4}	72.9 _{4.8}	48.6 _{1.9}	50.4 _{0.7}
LocalE	81.0 _{8.0}	42.3 _{4.7}	80.3 _{1.7}	75.6 _{4.1}	79.0 _{3.5}	72.5 _{5.8}	54.2 _{4.2}	54.0 _{2.6}	72.3 _{4.6}	50.4 _{1.9}	50.5 _{0.8}
GlobalE	80.2 _{4.2}	43.2 _{1.3}	81.2 _{0.9}	76.1 _{1.8}	80.3 _{1.4}	73.0 _{1.3}	54.3 _{4.0}	56.7 _{2.0}	78.1 _{1.9}	51.3 _{1.8}	51.2 _{0.8}
Oracle	89.8 _{0.7}	48.0 _{1.1}	85.4 _{1.6}	87.4 _{0.9}	90.1 _{0.7}	80.9 _{1.4}	60.3 _{10.3}	62.8 _{4.2}	81.3 _{2.9}	53.4 _{1.1}	52.5 _{1.4}
GPT-3 175B	93.9 _{0.8}	54.4 _{2.5}	95.4 _{0.9}	94.6 _{0.7}	91.0 _{1.0}	83.2 _{1.5}	71.2 _{7.3}	72.1 _{2.7}	85.1 _{1.7}	70.8 _{2.8}	75.1 _{1.1}
LocalE	93.8 _{0.5}	56.0 _{1.7}	95.5 _{0.9}	94.5 _{0.7}	91.3 _{0.5}	83.3 _{1.7}	75.0 _{4.6}	71.8 _{3.2}	85.9 _{0.7}	71.9 _{1.4}	74.6 _{4.2}
GlobalE	93.9 _{0.6}	53.2 _{2.1}	95.7 _{0.7}	94.6 _{0.2}	91.7 _{0.4}	82.0 _{0.8}	76.3 _{1.5}	73.6 _{2.5}	85.7 _{1.0}	71.8 _{1.9}	79.9 _{1.3}
Oracle	94.7 _{0.2}	58.2	96.7 _{0.2}	95.5 _{0.2}	92.6 _{0.4}	85.5 _{0.8}	81.1 _{4.9}	77.0 _{1.2}	87.7 _{0.6}	74.7 _{0.4}	83.0 _{0.9}



Thank You