Calibrate Before Use: Improving Few-Shot Performance of Language Models

Tony Z. Zhao* Eric Wallace* Shi Feng Dan Klein Sameer Singh

ICML 2021



UC Berkeley



University of Maryland



UC Irvine











Tony ZhaoEric WallaceUC BerkeleyUC Berkeley

Shi Feng UMD Dan Klein UC Berkeley

Sameer Singh UC Irvine

Slides, Blog, Code, and Video ericswallace.com/calibrate

Prompt

ICML 2021 is hosted





. . .

Sentiment Analysis

Sentiment Analysis



Sentiment Analysis



Sentiment Analysis



No model updates i.e., in-context learning

Few-shot Learning with LMs Topic Classification

Few-shot Learning with LMs Topic Classification



Knowledge Base Completion

Knowledge Base Completion



Why In-Context Learning?

- Academically interesting
- Practically relevant with GPT-3:
 - effective with ~0-16 examples
 - serve one model for many tasks
 - no ML expertise needed

Our paper's goal: analyze and improve in-context learning

Part 1: The Impact of the Prompt

Components Of The Prompt

- Prompt format
- Training example selection
- Training example permutation

Components Of The Prompt Prompt Format



Components Of The Prompt Prompt Format



Components Of The Prompt Prompt Format



Components Of The Prompt Training Example Selection

PromptInput: Subpar acting.Sentiment: negativePromptInput: Beautiful film.Sentiment: positiveInput: Amazing.Sentiment: content: content: content.

Components Of The Prompt Training Example Selection

PromptInput:Subpar acting.Sentiment:negativeInput:Beautiful film.Sentiment:positiveInput:Input:Good film.Sentiment:positiveInput:Don't watch.Sentiment:negativeInput:Amazing.Sentiment:Sentiment:

Components Of The Prompt Training Example Permutation

PromptInput: Subpar acting.Sentiment: negativePromptInput: Beautiful film.Sentiment: positiveInput: Amazing.Sentiment: content: content: content.

Components Of The Prompt Training Example Permutation

PromptInput:Subpar acting.Sentiment:negativeInput:Beautiful film.Sentiment:positiveInput:Input:Beautiful film.Sentiment:positiveInput:Subpar acting.Sentiment:negativeInput:Amazing.Sentiment:Sentiment:

Does The Prompt Affect Accuracy?

- Prompt format
- Training example selection
- Training example permutation

Does The Prompt Affect Accuracy?

- Prompt format
- Training example selection
- Training example permutation

Prompt 1

Review: the whole thing's fairly lame, making it par for the course for disney sequels. Answer: Negative

Review: this quiet, introspective and entertaining independent is worth seeking. Answer: Positive



Prompt 2

Review: this quiet, introspective and entertaining independent is worth seeking. Answer: Positive

Review: the whole thing's fairly lame, making it par for the course for disney sequels . Answer: Negative

Does The Prompt Affect Accuracy?

- Prompt format
- Training example selection
- Training example permutation

Prompt 1

Review: the whole thing's fairly lame, making it par for the course for disney sequels. Answer: Negative

Review: this guiet, introspective and entertaining independent is worth seeking. Answer: Positive

88.5% Acc.

Prompt 2

Review: this quiet, introspective and entertaining independent is worth seeking. Answer: Positive



Review: the whole thing's fairly lame, making it par for the course for disney sequels . Answer: Negative

51.3% Acc.

Training Set #1





All 24 Permutations





Accuracy Is Highly Sensitive To Prompt Design Example Permutation Impacts Accuracy








Accuracy Is Highly Sensitive To Prompt Design Example Selection Impacts Accuracy



Accuracy Is Highly Sensitive To Prompt Design

Accuracy Is Highly Sensitive To Prompt Design



Accuracy Is Highly Sensitive To Prompt Design

Prompt Format Impacts Accuracy





Does The Prompt Affect Accuracy?

- Yes

- Prompt format
- Training example selection Yes
- Training example permutation Yes

Part 2: Understanding Prompt Sensitivity

- Majority label bias: frequent training answers dominate predictions
 - helps explain variance across *example selections*

- Majority label bias: frequent training answers dominate predictions
 helps explain variance across example selections
- Recency bias: examples near end of prompt dominate predictions
 helps explain variance across example permutations

- Majority label bias: frequent training answers dominate predictions
 helps explain variance across example selections
- Recency bias: examples near end of prompt dominate predictions
 helps explain variance across example permutations
- Common token bias: common n-grams dominate predictions
 helps explain variance across prompt formats

Frequency of Positive Predictions

4/4	3/4	2/4	1/4	0/4
Positive	Positive	Positive	Positive	Positive





Frequent training answers dominate predictions

Frequency of Positive Predictions

NPPP PNPP PPNP PPPN

Frequency of Positive Predictions

NPPP	PNPP	PPNP	PPPN
Ť	1	1	<u>↑</u>

Frequency of Positive Predictions



Frequency of Positive Predictions



Examples near end of prompt dominate predictions



			Token	Prob
	Language Model		book	0.35
)	transportation	0.23
	M/bat topic is the following toxt about?]	school	0.11
Prompt	The Model T was released by Ford in 1908.		village	0.03
Ĩ	Answer:		company	0.02

			Token	Prob
	Language Model		book	0.35
)	transportation	0.23
	M/bat table is the fallowing taxt about?]	school	0.11
Prompt	The Model T was released by Ford in 1908.		village	0.03
1	Answer:		company	0.02

Token	Web (%)
book	0.026
transportation	0.000006

			Token	Prob
	Language Model		book	0.35
)	transportation	0.23
	M/bat tanic is the following taxt about?]	school	0.11
Prompt	The Model T was released by Ford in 1908.		village	0.03
1	Answer:		company	0.02

Token	Web (%)	Label (%)	
book	0.026	9	
transportation	0.0000006	9	

_

			Token	Prob
	Language Model		book	0.35
)	transportation	0.23
	M/bat table is the fallowing taxt about?]	school	0.11
Prompt	The Model T was released by Ford in 1908.		village	0.03
1	Answer:		company	0.02

Token	Web (%)	Label (%)	Prediction (%)
book	0.026	9	29
transportation	0.000006	9	4

			Token	Prob
	Language Model		book	0.35
)	transportation	0.23
	M/battania is the following toxt about?]	school	0.11
Prompt	The Model T was released by Ford in 1908.		village	0.03
Ĩ	Answer:		company	0.02

Token	Web (%)	Label (%)	Prediction (%)
book	0.026	9	29
transportation	0.000006	9	4

Common n-grams dominate predictions

What Is The Impact of These Biases?

What Is The Impact of These Biases?



What Is The Impact of These Biases?



Biases cause a shift in output distribution

Part 3: Improving Few-Shot Accuracy

Step 1: Estimate the bias

Step 1: Estimate the bias

Insert "content-free" test input

Input: Subpar acting. Sentiment: negative Input: Beautiful film. Sentiment: positive Input: N/A Sentiment:

Step 1: Estimate the bias

Insert "content-free" test input

Input: Subpar acting. Sentiment: negative Input: Beautiful film. Sentiment: positive Input: N/A Sentiment:

Get model's prediction

positive	0.65	
negative	0.35	

Step 1: Estimate the bias

Insert "content-free" test input

Input: Subpar acting. Sentiment: negative Input: Beautiful film. Sentiment: positive Input: N/A Sentiment: Step 2: Counter the bias

Get model's prediction

positive	0.65
negative	0.35

Step 1: Estimate the bias

Insert "content-free" test input

Input: Subpar acting. Sentiment: negative Input: Beautiful film. Sentiment: positive Input: N/A Sentiment: Step 2: Counter the bias

"Calibrate" predictions with affine transformation

$$\mathbf{\hat{q}} = \operatorname{softmax}(\mathbf{W}\mathbf{\hat{p}} + \mathbf{b})$$

Get model's prediction

positive	0.65
negative	0.35
Contextual Calibration

Step 1: Estimate the bias

Insert "content-free" test input

Input: Subpar acting. Sentiment: negative Input: Beautiful film. Sentiment: positive Input: N/A Sentiment:

Get model's prediction

positive	0.65
negative	0.35

Step 2: Counter the bias

"Calibrate" predictions with affine transformation

$$\mathbf{\hat{q}} = \operatorname{softmax}(\mathbf{W}\mathbf{\hat{p}} + \mathbf{b})$$

 $\mathbf{\hat{q}}$

Calibrated probs

Calibrated p

Contextual Calibration

Step 1: Estimate the bias

Insert "content-free" test input

Input: Subpar acting. Sentiment: negative Input: Beautiful film. Sentiment: positive Input: N/A Sentiment:

Get model's prediction

positive	0.65
negative	0.35

Step 2: Counter the bias

"Calibrate" predictions with affine transformation

$$\mathbf{\hat{q}} = \operatorname{softmax}(\mathbf{W}\mathbf{\hat{p}} + \mathbf{b})$$

 $\mathbf{\hat{q}}$
Calibrated probs
Original probs

Fit ${f W}$ and ${f b}$ to cause uniform prediction for "N/A"

Contextual Calibration

Step 1: Estimate the bias

Insert "content-free" test input into prompt

Input: Subpar acting. Sentiment: negative Input: Beautiful film. Sentiment: positive Input: N/A Sentiment:

Get model's prediction

positive	0.65
negative	0.35

Step 2: Counter the bias

"Calibrate" predictions with affine transformation

$$\mathbf{\hat{q}} = \operatorname{softmax}(\mathbf{W}\mathbf{\hat{p}} + \mathbf{b})$$

 $\mathbf{\hat{q}}$

Calibrated probs

Calibrated probs

Fit ${\bf W}$ and ${\bf b}$ to cause uniform prediction for "N/A"

$$\mathbf{W} = \begin{bmatrix} \frac{1}{0.65} & 0\\ 0 & \frac{1}{0.35} \end{bmatrix} \qquad \mathbf{b} = \begin{bmatrix} 0\\ 0\\ 0 \end{bmatrix}$$

- Experiment with 11 different datasets
 - Classification
 - Knowledge base completion
 - Information extraction
- Consider 0, 1, 4, 8, and 16 training examples
- Different sizes of GPT-3 and GPT-2 language models

Improves mean and worst-case accuracy



- Improves mean and worst-case accuracy
- Reduces variance across training sets and permutations



- Improves mean and worst-case accuracy
- Reduces variance across training sets and permutations
- Reduces variance across prompt formats





Summary

GPT-3's accuracy has high variance across different prompts LMs have biases that hurt few-shot learning Contextual calibration improves accuracy and reduces variance Paper and Code at <u>ericswallace.com/calibrate</u>