



Article Presentation

IA368DD_2023S1: Deep Learning aplicado a Sistemas de Buscas

Student: Marcus Vinícius Borela de Castro

Language Models are Few-Shot Learners

By: Tom B.
Brown *et al.*
(OpenAi). 2020

GPT-3 Generative Pretrained Transformer

Um modelo de linguagem auto-regressivo de 175 bilhões de parâmetros (10x mais do que qualquer modelo de linguagem não esparsa anterior), baseado no GPT-2, com pequenos ajustes de um *Sparse Transformer*.

Todos os contextos avaliados possuem uma descrição da tarefa a ser realizada e podem ou não possuir exemplos:

- (a) *few-shot*: onde há exemplos (tipicamente de 10 a 100, limitado à janela de contexto de até 2048 tokens),
- (b) *one-shot*: com apenas uma demonstração, e
- (c) *zero-shot*: sem exemplos

Tipos de contexto avaliados

Main concepts

In-context learning (aprendizagem no contexto)

É a capacidade do modelo se adaptar a uma tarefa específica no momento de inferência, que ocorre dentro do *forward pass*, sem gerar nenhum novo aprendizado (sem atualizar parâmetro do modelo), fazendo uso de habilidades desenvolvidas durante o pré-treinamento não supervisionado.

Com o GPT-3 se tenta repetir capacidade humana de poder realizar uma nova tarefa de linguagem a partir de apenas alguns exemplos ou de instruções simples.

A diferença entre o desempenho de zero, um e poucos exemplos geralmente cresce com a capacidade do modelo

A importância do pré-treinamento em larga escala para o desempenho do modelo GPT-3 (by ChatGPT)

Apresenta uma análise detalhada das habilidades e limitações do modelo em diferentes tarefas de linguagem natural. (by ChatGPT)

Modelos maiores viabilizam "in-context-aprendizes"

Article contribution



Está no texto: "Tarefas sintéticas, como embaralhar palavras ou definir palavras sem sentido, parecem ser especialmente prováveis de serem aprendidas de novo, enquanto a tradução claramente deve ser aprendida durante o pré-treinamento, embora possivelmente de dados que são muito diferentes em organização e estilo dos dados de teste. Em última análise, nem mesmo está claro o que os humanos aprendem do zero em relação a demonstrações anteriores."

O GPT-3 é pré-treinado em uma grande quantidade de dados de texto não rotulados, o que lhe permite aprender a reconhecer e entender padrões linguísticos e semânticos. Esse pré-treinamento em larga escala permite que o modelo capture o conhecimento geral da linguagem e crie representações semânticas complexas que podem ser transferidas para novas tarefas. (by ChatGPT)

O aprendizado de contexto realmente ocorre na inferência ou simplesmente se identificam tarefas aprendidas durante o treinamento?

Como se dá essa redução de transferência de dados?

Para minimizar o tempo de treinamento, utilizou-se processamento paralelo entre GPUs com a redução da transferência de dados entre nós tanto na profundidade quanto na largura da arquitetura.

Embora o GPT-3 3B seja quase 10 vezes maior que o RoBERTa-Large (parâmetros de 355M), ambos evaram aproximadamente 50 petaflop/s-dias de computação durante o pré-treinamento

O que é um "s-day"? Por que precisa desse denominador?

- Interesting/unexpected results

Habilidades de geração de texto: Na geração de artigos, dados título e subtítulo, a detecção humana, se é ou não real, alcançou 52%, quase o acaso.

Em geral, em tarefas de NLP, o GPT-3 alcançou resultados promissores nas configurações *zero-shot* e *one-shot*, e na *few-shot* às vezes é competitivo ou até supera o estado da arte (superando modelos com ajuste fino).

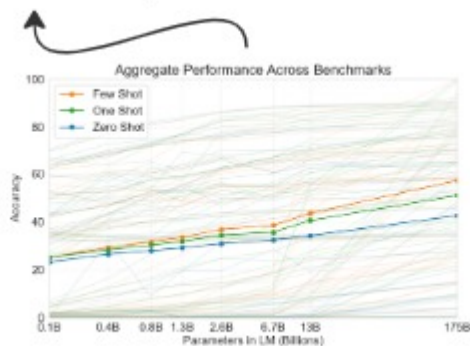
Por exemplo, GPT-3 atinge 64,3% de precisão no TriviaQA no *zero-shot*, 68,0% no *one-shot* (igual SOTA anterior) e 71,2% no *few-shot*, este SOTA

Demonstra que as habilidades de aprendizagem em contexto demonstram ganhos fortes com a escala

Para se estudar a correlação do desempenho com o tamanho do modelo, foram treinadas 8 configurações diferentes, variando em três ordens de grandeza, de 125 milhões de parâmetros a 175 bilhões de parâmetros, sendo o maior modelo referenciado como GPT-3

Todos os modelos foram treinados para um total de 300 bilhões de tokens

Figure 1.3: Aggregate performance for all 42 accuracy-denominated benchmarks



O GPT-3 apresenta um desempenho pior em algumas tarefas

Tarefas que se beneficiam empiricamente da bidirecionalidade (como em tarefas de preenchimento de lacunas entre textos: fill mask).

Tarefas sintéticas e qualitativas que exigem correspondência de padrões e computação não triviais e que são improváveis nos dados de treinamento.

Tarefas da "física de senso comum", com perguntas do tipo "Se eu colocar queijo na geladeira, ele derreterá?"

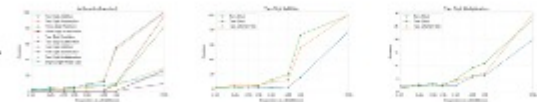
Obs.:(Isso mudou: perguntei 17/3/2023 ao WebChatGPT).

II Results on All Tasks for All Model Sizes

Name	Metric	Type	Size	Unit	2016-2017					2018-2019					2020-2021					1/2021 new sensor		
					Small	Med	Large	XL	2.70-6.78	130	1750	Small	Med	Large	XL	2.70-6.78	130	1750	Small		Med	Large
LibraGrap	acc	dir	88.6	50	71.7	47.6	51.9	54.7	52.8	67.4	79.0	78.9	71.0	47.9	52.8	67.4	68.9	70.0	78.1	71.5	70.1	71.5
LibraGrap	acc	ind	69.0	15	42.7	34.1	40.4	43.6	47.1	70.1	72.8	76.2	71.0	47.1	52.8	67.4	68.9	70.0	78.1	71.5	70.1	71.5
LibraGrap	ppd	ind	8.75	15	18.6	9.0	6.5	5.4	4.6	4.0	3.9	3.8	10.0	11.6	8.5	8.4	9.1	4.6	4.6	5.3	5.3	5.3
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.3	74.1	77.3	74.7	70.1	84.7	84.7	84.7
LibraGrap	acc	ind	91.8	70	67.1	69.5	72.4	71.4	77.2	71.7	79.8	83.2	67.1	68.7	72.							

A forma como resumiu em uma tabela os resultados

E ainda tem os gráficos para cada linha da tabela



Advanced topic to discuss

Parece ser promissor repetir em modelo bidirecional a escala do GPT-3 com meta-aprendizagem (falta de bidirecionalidade explicaria piores resultados em alguns testes)

Foram identificados vieses no comportamento dos modelos, que podem gerar conteúdo estereotipado ou preconceituoso, consequência dos vieses presentes nos dados de treinamento, por exemplo, quanto a sexo, raça e religião.

Por que não se promove esse teste?

Tem algum LLM (Large Language Model) bidirecional?

Como evitar/mitigar esses vieses?

Em algumas tarefas o problema era, provavelmente, o tokenizador.

Como em Palavras invertidas (RW) – O modelo recebe uma palavra escrita ao contrário e deve produzir a palavra original. Exemplo: stcejbo!objects.

Mas hoje o ChatGPT já resolve essa questão.

reverse this word: sucram
The reverse of the word "sucram" is "marcus".

Como ele faz isso? Ele tokeniza letra a letra? Ou ele experimenta diferentes tokenizadores durante a inferência?

Um tópico avançado que pode surgir da leitura do artigo é a discussão sobre a capacidade dos modelos de linguagem, como o GPT-3, de aprender conceitos mais abstratos e de alto nível. O artigo sugere que o modelo pode aprender conceitos mais abstratos, como a matemática avançada, programação de computadores e até mesmo habilidades sociais e emocionais. (by ChatGPT)