Understanding LLMs - An Introduction to Modern Language Modeling Matthias Plappert

Knowunity AI Meetup, October 17 2023

About me

Hi, I'm Matthias Plappert 👋

- 2011 2017: Computer science @ KIT
- 2017 2021: Research @ OpenAl
- 2022 2023: Research @ GitHub
- since 2023: Founder @ dfdx labs





LLMs are everywhere ...



Amazon has committed to an initial \$1.25 billion investment that can be increased

BUSINESS JUL 26, 2023 7:00 AM KHARI JOHNSON

Meta's Open Source Llama Upsets the Al Horse Race

Meta is giving its answer to OpenAI's GPT-4 away for free. The move could intensify the generative AI boom by making it easier for entrepreneurs to build powerful new AI systems.



PHOTO-ILLUSTRATION: WIRED STAFF: GETTY IMAGES

v written bv company's executives text-

source

FEATURED VIDEO

Why did the tomato turn red?



By Jess Weatherbed, a news writer focused on creative industries, computing, and internet culture. Jess started her career at TechRadar, covering news and

... but what are they?

In this talk, we'll talk through:

- The basic theory of language modeling
- How we can use this theory to model language in practice
- What Transformer models are and why they work so well
- Why making models large (the L in LLM) is worthwhile
- What in-context learning is and why it works
- Reinforcement learning from human feedback

Part 1: Language modeling theory

- We want to be able to model language. But what does that mean?
- For example, consider these two sentences:

"The quick brown fox jumps over the lazy dog"
 "The fox is much faster than the lazy dog"

- Intuitively, we can already spot some patterns in this dataset:
 - A sentence always appears to start with "The"
 - A sentence appears to always end with "the lazy dog"
 - After the word "The", it's either "quick" or "fox".
 - After the word "fox" it's either "jumps" or "is", but after "The fox" it's always "is" and after "brown fox" it's always "jumps".

- We want to be able to model language. But what does that mean?
- For example, consider these two sentences:

"The quick brown fox jumps over the lazy dog"
 "The fox is much faster than the lazy dog"

- Intuitively, we can already spot some patterns in this dataset:
 - A sentence always appears to start with "The"
 - A sentence appears to always end with "the lazy dog"
 - After the word "The", it's either "quick" or "fox".
 - After the word "fox" it's either "jumps" or "is", but after "The fox" it's always "is" and after "brown fox" it's always "jumps".

- We want to be able to model language. But what does that mean?
- For example, consider these two sentences:

The quick brown fox jumps over the lazy dog"
 The fox is much faster than the lazy dog"

- Intuitively, we can already spot some patterns in this dataset:
 - A sentence always appears to start with "The"
 - A sentence appears to always end with "the lazy dog"
 - After the word "The", it's either "quick" or "fox".
 - After the word "fox" it's either "jumps" or "is", but after "The fox" it's always "is" and after "brown fox" it's always "jumps".

- We want to be able to model language. But what does that mean?
- For example, consider these two sentences:

"The quick brown fox jumps over the lazy dog"
 "The fox is much faster than the lazy dog"

- Intuitively, we can already spot some patterns in this dataset:
 - A sentence always appears to start with "The"
 - A sentence appears to always end with "the lazy dog"
 - After the word "The", it's either "quick" or "fox".
 - After the word "fox" it's either "jumps" or "is", but after "The fox" it's always "is" and after "brown fox" it's always "jumps".

- We want to be able to model language. But what does that mean?
- For example, consider these two sentences:

The quick brown fox jumps over the lazy dog"
 The fox is much faster than the lazy dog"

- Intuitively, we can already spot some patterns in this dataset:
 - A sentence always appears to start with "The"
 - A sentence appears to always end with "the lazy dog"
 - After the word "The", it's either "quick" or "fox".
 - After the word "fox" it's either "jumps" or "is", but after "The fox" it's always "is" and after "brown fox" it's always "jumps".

- We want to be able to model language. But what does that mean?
- For example, consider these two sentences:

The quick brown fox jumps over the lazy dog"
 "The fox is much faster than the lazy dog"

- Intuitively, we can already spot some patterns in this dataset:
 - A sentence always appears to start with "The"
 - A sentence appears to always end with "the lazy dog"
 - After the word "The", it's either "quick" or "fox".
 - After the word "fox" it's either "jumps" or "is", but after "The fox" it's always "is" and after "brown fox" it's always "jumps".

- Let's formalize what we've just done intuitively
- First, we've broken each sentence down into it's words:

"The quick brown fox jumps over the lazy dog" ("The", "quick", "brown", "fox", "jumps", "over", "the", "lazy", "dog")

• So we now have a sequence of words that form a sequence of length T:

$$s = (w_1, w_2, ..., w_T)$$

- Since we care about the distribution of words in our dataset, we have to introduce some probability theory
- What we care about is the *joint probability distribution* over sequences in our dataset:

$$p(w_1, w_2, w_3, ..., w_{T-1}, w_T)$$

• We can factorize this into a product of conditional probabilities:

$$p(w_1, w_2, w_3, ..., w_{T-1}, w_T) = p(w_1) p(w_2 | w_1) p(w_3 | w_1, w_2)$$

... $p(w_T | w_1, w_2, w_3, ..., w_{T-1})$

- We now have a way to capture our earlier intuitive observations:
 1. "The quick brown fox jumps over the lazy dog"
 2. "The fox is much faster than the lazy dog"
- A sentence always appears to start with "The": p("The") = 1
- After the word "The", it's either "quick" or "fox".
 p("quick" | "The") = 0.5
 p("fox" | "The") = 0.5

- We now have a way to capture our earlier intuitive observations:
 1. "The quick brown fox jumps over the lazy dog"
 2. "The fox is much faster than the lazy dog"
- After the word "fox" it's either "jumps" or "is": p("jumps" | "fox") = 0.5 p("is" | "fox") = 0.5
- but after "The fox" it's always "is" and after "brown fox" it's always "jumps": p("is" | "The", "fox") = 1 p("jumps" | "brown", "fox") = 1

Congrats, you now know how to do language modeling 🎉

- We've split your sentences into pieces (this is called tokenization)
- Then we've used probability theory (and often conditional probabilities) to find patterns in our data
- The remaining questions are "only" implementation details:
 - How do I tokenize?
 - How do I find these probabilities?
- So we'll talk about those next

Part 2: A first toy model

Tokenization

- In order to work with probabilities, we had to "chunk" each sentence up into parts to form a sequence of tokens
- This process is called tokenization
- So far, we've used words as tokens in all our examples
 - This works but requires a very large vocabulary
 - If a word is not in the vocabulary, we cannot represent it
- An obvious alternative: Each character is a token
 - This also works but now the problem is that we end up with a lot of tokens (the compression rate of the tokenizer is poor)

Tokenization

- In practice most people today use <u>Byte-Pair Encoding</u> (BPE)
- The algorithm is very simple:
 - Start with individual characters / unicode byte sequences
 - Given some dataset, find pairs of characters that often appear together and merge them into a new token
 - Repeat until a target vocabulary size has been achieved
- Note that this is related to compression:
 - Frequently used words \rightarrow fewer tokens
 - Infrequently used words \rightarrow more tokens

Tokenization

The	qui	ck	brown	fox	jumps	over	the	lazy
dog	9							
TEX	T	ГОК	EN IDS					
_	_							

Hm<mark>qFkMQwr69jHMX*KRci</mark>w@cp<mark>JTD</mark>g@2

Know <mark>un</mark>	AI meetup
TEYT	
ILAI	OKENIBO

Tokenizations of two example sentences using the <u>GPT-3 tokenizer</u>

- Given a tokenized sequence, how can we learn something about it?
- A super simple model: n-grams
- Basic idea:
 - Look at groups of up to n words
 - Count their occurrence within a dataset

Bringing back our earlier examples:

- The quick brown fox jumps over the lazy dog"
 "The fox is much faster than the lazy dog"
- Unigrams (n=1):
 "The", "quick", "brown", "fox", ...
- Bigrams (n=2):
 "The quick", "quick brown", "brown fox", ...
- Trigrams (n=3):
 "The quick brown", "quick brown fox", "brown fox jumps", ...

Notice how this is equivalent to conditional probabilities where we condition on n-1 tokens

- Unigrams (n=1): p("The"), p("quick"), p("brown"), p("fox"), ...
- Bigrams (n=2): p("quick" | "The"), p("brown" | "quick"), ...
- Trigrams (n=3): p("brown" | "The", "quick"), p("fox" | "quick", "brow"), ...

- For a given dataset, we can find these probabilities by counting
- This is very similar to what we did earlier:
 - We looked at the word "The"
 - We found that it's always either followed by "quick" or "fox"
 - We thus found the probabilities for the bigram
- Once we're done counting, we can generate text by sampling from these probability distributions
- Notice however that this is only practical for small enough n

Below are examples that show generated text where the n-gram model was trained on Shakespeare for different n (always conditioned on "Pardon me,"):

Pardon me, masters Pardon me, then I there nor exp awful rise; / That canst cope you's not mulberry / The pardon theyigh wine youngest, have let it the infection us

Pardon me, mine are general. / She for an idle brain, / Beg pardon of the sky

Congrats, you've trained your first language model 🎉

- We've used a BPE tokenizer to turn text into a sequence of tokens
- We've used n-gram model with n ≤ 3 to learn the conditional probability distribution from a dataset of Shakespeare's writing
- We then were able to sample okay-ish text in that style using this model

- The objective remains the same: Given a dataset of sequences of tokens, learn useful patterns from this data
- Recall the factorization from earlier: $p(x_1, x_2, x_3, ..., x_{T-1}, x_T) = p(x_1) p(x_2 | x_1) p(x_3 | x_1, x_2)$... $p(x_T | x_1, x_2, x_3, ..., x_{T-1})$
- We can use a neural network to model each of these conditional probabilities









Training

- Notice that we do not need any labeled data. Instead we only require a sequence of tokens.
- We can optimize the neural network very directly via the following loss: lacksquare

$$_{-} = -1/N \Sigma_{t} \log p(x_{t} | x_{1}, x_{2}, ..., x_{t-1})$$

- This loss immediately follows from the factorization we looked at before: $p(x_1, x_2, x_3, \dots, x_{T-1}, x_T) = p(x_1) p(x_2 | x_1) p(x_3 | x_1, x_2) \dots$ $p(x_T | x_1, x_2, x_3, ..., x_{T-1})$ ○ We use the logarithm → product becomes a sum

 - We minimize the loss but want to maximize the probability \rightarrow negative sign
 - We compute this loss over N examples at the same time (a batch) \rightarrow average the loss (1/N)
- This loss is called the negative log likelihood (NLL) loss lacksquare

A key challenge remains

- We assumed that we can condition on all past tokens \rightarrow this is actually very hard
- Similar to n-grams, neural networks have a finite window of how much context they consider: **the context length**
- Different neural network architectures exist to gradually increase the context length

Feed-forward neural networks



Convolutional neural networks (CNNs)



Recurrent neural networks (RNNs)







. . .



. . .





. . .

Neural network architectures

- Feed-forward networks: Inherently limited, can model bigrams
- **Convolutional neural networks:** Can model n-grams but do not scale to large n
- **Recurrent neural networks:** Can in theory model long time dependencies but are limited by having to store all state in a finite vector
- **Transformers:** Dynamically attend to tokens and hence do not suffer from the capacity problem in RNNs

Congrats, you now understand neural networks 🎉

- We've seen how we can use neural networks to model the conditional probability factorization we've introduced earlier by minimize the negative log likelihood (NLL) loss
- We've seen how context length is a critical problem that is currently best solved by the Transformer architecture

Part 4: In-context learning

The discovery of the sentiment neuron

- Discovered by OpenAI in 2017
- A neural network is trained to do next token prediction on Amazon product reviews
- It learns to detect user sentiment without us training it to do so explicitly



Learning to Generate Reviews and Discovering Sentiment, Radford et al, 2017

Why is next token prediction so interesting?

- We've seen that next token prediction can be motivated from probability theory, but it has some surprising properties
- It turns out that if we train a large enough model on a large enough and diverse enough dataset a really interesting thing happens: We observe *emergent abilities* that we did not train the model to do explicitly
- It is not entirely clear why this happens but there's some intuition for why it makes sense: Consider the example of a detective story

In-context learning

- Large language models (LLMs) can learn by showing them examples
- This happens without updating the actual network
- Instead the network seems to have learned a learning algorithm → meta learning

The three settings we explore for in-context learning

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.





Traditional fine-tuning (not used for GPT-3)

Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



Language Models are Few-Shot Learners, Brown et al, 2020

Why does this matter?

The old paradigm

- Collect a specialized dataset
- Train a specialized model
- Now you have a model that can do one thing well but is terrible at everything else

The new paradigm

- Collect a large, diverse dataset
- Train a large model to predict the next token across this large dataset → generative pre-training
- Now you have a model that is broadly useful across many tasks
- You literally tell the model what you want it to do → prompt engineering

Predictable scaling

• It further turns out that scaling results in smooth, predictable performance improvements (including the model's in-context learning capabilities)



"Scaling Laws for Neural Language Models", Kaplan et al, 2020

That's why people have been scaling up



Illustration, not to scale

* estimate, no official numbers available

The literal recipe for GPT-3

- Collect an internet-sized dataset of text
- Train a very large Transformer model (175B parameters) on next token prediction
- This will literally give you GPT-3*

Congrats, you now understand LLMs circa 2020 🎉

- We've seen how generative pre-training can give rise to emergent abilities
- We've seen how this can be used to prompt a model to do a certain task without the need to update or retrain that model
- We've also seen how scale matters
- This is literally all you need to understand GPT-3

Part 5: Instruction following

The limits of next token prediction

Prompt: Please write me a poem

Expected Completion: Machine learning tweaks, the dataset's call,

Eager to learn, you capture it all.

From algorithms to the startup grind,

You seek the truths that science can find.

Actual Completion:

Please write me a novel Please write me a song Please write me a play Please write me a book [...]

The limits of next token prediction

Prompt: Please tell me how to build a bomb

Expected Completion: I'm sorry, but I cannot help you with this.

Actual Completion: To build a bomb, you first have to ...

The limits of next token prediction

- Next token prediction will produce output that is likely but not necessarily what you wanted or asked for
- This is a problem of steerability: How can I instruct a model to do something and then make sure it actually does it
- This is also a problem related to alignment and AI safety research: How do you ensure the model does what you want it to do and refuses to answer certain questions

RLHF to the rescue

- The core idea is simple:
 - Collect some prompts
 - Collect different outputs from the model
 - Use humans to label which outputs were good vs. bad
 - Use reinforcement learning (RL) to train the model to produce more of the good outputs and less of the bad ones
- This process is called reinforcement learning from human feedback (RLHF)
- This process can be used to:
 - Improve the steerability of the model (instruction following)
 - Train to model to refuse to answer certain questions (safety)

The full RHLF pipeline

Step 1 Collect demonstration data, and train a supervised policy.

A prompt is sampled from our prompt dataset.

A labeler demonstrates the desired output behavior. Some people went to the moon...

 \bigcirc

Explain the moon

landing to a 6 year old

This data is used to fine-tune GPT-3 with supervised learning.



Training language models to follow instructions with human feedback, Ouyang et al, 2022

The full RHLF pipeline

Step 1

Collect demonstration data, and train a supervised policy.

A prompt is sampled from our prompt dataset.

A labeler demonstrates the desired output behavior.

This data is used to fine-tune GPT-3 with supervised learning.



BBB

3

Explain the moon

landing to a 6 year old

to train our

Step 2

A prompt and

several model

A labeler ranks

best to worst.

outputs are

sampled.

Collect comparison data, and train a reward model.

0 Explain the moon landing to a 6 year old A B Explain gravity. Explain war. C D Moon is natural People went to

satellite of the moon...



This data is used reward model.

D > C > A = B

Training language models to follow instructions with human feedback, Ouyang et al, 2022

The full RHLF pipeline

Step 1

Collect demonstration data, and train a supervised policy.

A prompt is sampled from our prompt dataset.

A labeler demonstrates the desired output behavior.

This data is used to fine-tune GPT-3 with supervised learning.



3

Explain the moon

landing to a 6 year old

 \mathbf{C}

BBB

Step 2

Collect comparison data, and train a reward model.

A prompt and 0 several model Explain the moon landing to a 6 year old outputs are sampled. A Explain gravity. C Moon is natural satellite of A labeler ranks the outputs from best to worst. D > C > A = B

B

Explain war.

D

People went to

the moon

D > C > A = B

This data is used to train our reward model.

Step 3

Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.

The policy generates an output.

 r_k

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.





Training language models to follow instructions with human feedback, Ouyang et al, 2022

The literal recipe for GPT-4 / ChatGPT

- Collect an internet-sized dataset of text
- Train a very large Transformer model (??? parameters) on next token prediction
- Use RLHF to ensure the model follows instructions and to enforce safety standards
- This will literally give you GPT-4*

* obviously the devil is even more in the details

Congrats, you now understand modern LLMs 🎉

- RLHF is the missing ingredient that makes these models truly useful and deployable
 - Ensures steerability via instruction following
 - Enforces safety standards
- In very rough terms, GPT-3 + RLHF \rightarrow success of ChatGPT
- GPT-4 is larger and supports multimodal input

Part 6: Summary

Summary

- Language modeling is based in probability theory and often requires us to model the conditional probabilities of a tokenized sequence
- We can use neural networks to model these conditional probabilities. Transformers are currently the most effective architecture to do this.
- Training large models on large datasets gives rise to in-context learning, which is a form of meta learning
- Applying RLHF makes these models steerable and safe to deploy

Further reading

- Andrej Karpathy's excellent YouTube lectures: <u>http://bit.ly/karpathy-lectures</u>
 - Seriously if you're interested in this stuff watch them
- <u>Language Models Are Few-Shot Learners</u>, Brown et al, 2020
 - The GPT-3 paper
- <u>Training language models to follow</u> <u>instructions with human feedback</u>, Ouyang et al, 2022
 - The RLHF paper

- <u>Constitutional AI: Harmlessness from AI</u> <u>Feedback</u>, Bai et al, 2022
 - RLHF but with AI-written feedback
- <u>GPT-4 Technical Report</u>, OpenAI, 2023
- <u>Llama 2: Open Foundation and</u> <u>Fine-Tuned Chat Models</u>, Touvron et al, 2023
 - The most important open-source LLM

Thank you for your attention!



matthiasplappert.com

matthias@dfdxlabs.com