

# Prompting vs. Training

Jacob Matthews

# Remember our definition of a language model

It defines a probability distribution over sequences of symbols

# How do we arrive at that probability distribution?

Language models are trained, or **pretrained**, on very large text corpora.

We can turn any text sequence into a series of inputs and labels for training a language model.

# Causal LM (“Next word prediction”)

Full sequence: “The quick brown fox jumps over the lazy dog”

**Input:** The

**Label:** quick

**Loss:**  $-\log(p(\text{quick} \mid \text{the}))$

**Input:** The quick

**Label:** brown

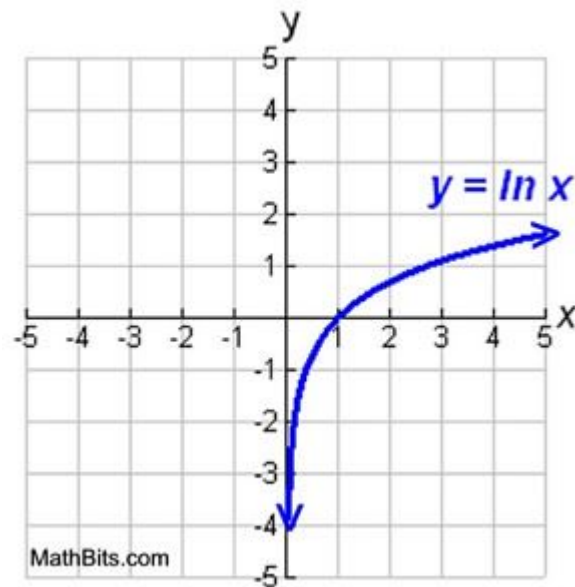
**Loss:**  $-\log(p(\text{brown} \mid \text{the quick}))$

**Input:** The quick brown

**Label:** fox

**Loss:**  $-\log(p(\text{fox} \mid \text{the quick brown}))$

...



# Causal LM (“Next \_\_\_\_ prediction”)

Full sequence: “Nc3 f5 e4 fxe4 Nxe4 Nf6 Nxf6+ gxf6 Qh5#”

**Input:** Nc3

**Label:** f5

**Loss:**  $-\log(p(\text{Nc3} \mid \text{f5}))$

**Input:** Nc3 f5

**Label:** e4

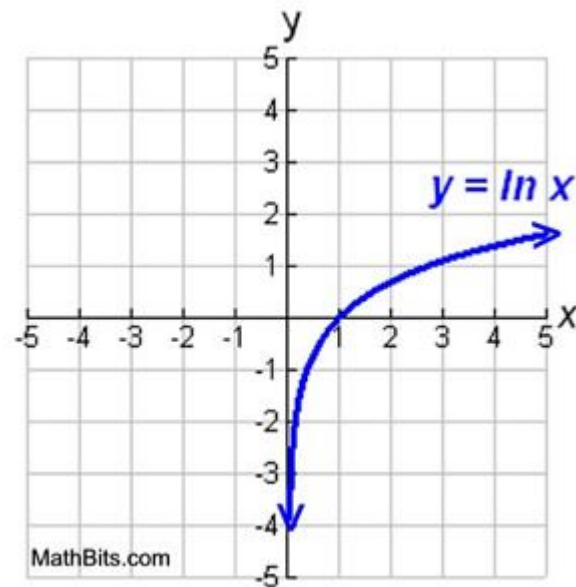
**Loss:**  $-\log(p(\text{e4} \mid \text{Nc3 f5}))$

**Input:** Nc3 f5 e4

**Label:** fxe4

**Loss:**  $-\log(p(\text{fxe4} \mid \text{Nc3 f5 e4}))$

...



# Causal LM (“Next word prediction”)

If you do this for trillions of tokens, you’re probably going to end up with a pretty good next word predictor!

[The GPT-2 paper](#) is a classic description of this approach using the transformer architecture.

# Pretraining vs Finetuning vs Prompting

This unsupervised, large-scale initial training is called **pretraining**.

It provides the model with broad linguistic capabilities.

However, these broad capabilities may not align well with your ultimate model use case.

Continuing to train a pretrained model on a (relatively) small dataset is called **finetuning**.

# Pretraining vs Finetuning vs Prompting

Since [GPT-3 \(2020\)](#), prompting has been shown to be a viable alternative to finetuning given sufficient scale.

Instead of finetuning models for specific downstream tasks, one instead crafts detailed prompts upon which the model can condition its response.



# “Training” vs Training

Sometimes, you may hear people say that they have “trained their GPT” to do something or to behave a certain way through repeated interactions.

This is not training, but is just the result of extended prompting.

**All training (pretraining, finetuning, ...) ultimately involves adjusting the model’s parameters** through gradient descent, RL, etc.

**No matter how much you prompt, the model’s parameters will remain unchanged.**

However, this discrepancy between “training” and training highlights how effective prompting is.

# Big Picture

It used to be that finetuning a base model was your best choice for a given task.

Now, large frontier models are so capable that this is probably not true for many common use cases.

However, finetuning can still get you pretty far (and, for small models, it can be cheap/free).

# Basic Information Theory

See notebook!