

CS 335: Introduction to Large Language Models

Introduction

Week 1

Dr Abdul Samad

Dr Faisal Alvi

Lecture Outline

- What is a Language Model?
- Generation
- Autoregressive Language Model
- Entropy
- N-gram Models
- Neural Language Models
- Language Models in the Real World

What is a Language Model?

- The classic definition of a language model (LM) is a **probability distribution over sequences of tokens**. Suppose we have a **vocabulary** V of a set of tokens. A language model p assigns each sequence of tokens $x_1, \dots, x_L \in V$ a probability (a number between 0 and 1):

$$p(x_1, \dots, x_L).$$

What is a Language Model?

- The probability intuitively tells us how “good” a sequence of tokens is. For example, if the vocabulary is $V = \{\text{ate, ball, cheese, mouse, the}\}$, the language model might assign:

$$p(\text{the, mouse, ate, the, cheese}) = 0.02,$$

$$p(\text{the, cheese, ate, the, mouse}) = 0.01,$$

$$p(\text{mouse, the, the, cheese, ate}) = 0.0001.$$

What is a Language Model?

- The language model (LM) should assign “mouse the the cheese ate” a very low probability implicitly because it’s ungrammatical (syntactic knowledge).
- The language model (LM) should assign “the mouse ate the cheese” higher probability than the cheese ate the mouse implicitly because of world knowledge: both sentences are the same syntactically, but they differ in semantic plausibility.

Generation

- As defined, a large language model p takes a sequence and returns a probability to assess its goodness. We can also generate a sequence given a language model. The purest way to do this is to sample a sequence $x_{1:L}$ from the language model p with probability equal to $p(x_{1:L})$, denoted:

$$x_{1:L} \sim p.$$

AutoRegressive Language Model

- A common way to write the joint distribution $p(x_{1:L})$ of a sequence $x_{1:L}$ is using the **chain rule of probability**:

$$p(x_{1:L}) = p(x_1)p(x_2|x_1)p(x_3|x_1, x_2) \cdots p(x_L|x_{1:L-1}) = \prod_{i=1}^L p(x_i|x_{1:i-1})$$

- For example:

$$\begin{aligned} p(\text{the, mouse, ate, the, cheese}) &= p(\text{the}) \\ &\quad p(\text{mouse} | \text{the}) \\ &\quad p(\text{ate} | \text{the, mouse}) \\ &\quad p(\text{the} | \text{the, mouse, ate}) \\ &\quad p(\text{cheese} | \text{the, mouse, ate, the}). \end{aligned}$$

- In particular, $p(x_i | x_{1:i-1})$ is a **conditional probability distribution** of the next token x_i given the previous tokens $x_{1:i-1}$.

AutoRegressive Language Model

Generation

Now to generate an entire sequence $x_{1:L}$ from an autoregressive language model p , we sample one token at a time given the tokens generated so far:

for $i = 1, \dots, L$:

$$x_i \sim p(x_i | x_{1:i-1})^{1/T},$$

where $T \geq 0$ is a **temperature** parameter that controls how much randomness we want from the language model:

ACTIVITY: Temperature Parameter

let's consider a simple example with four possible outcomes: 'cheese', 'room', 'it', and 'in'. Let's say we have a probability distribution

$$p(x_i | x_{1:i-1})$$

Over these four outcomes such that

$$\begin{cases} p(\text{"cheese"} | x_{1:i-1}) = 0.1 \\ p(\text{"room"} | x_{1:i-1}) = 0.2 \\ p(\text{"it"} | x_{1:i-1}) = 0.3 \\ p(\text{"in"} | x_{1:i-1}) = 0.4 \end{cases}$$

Now Consider

$$p(x_i | x_{1:i-1})^{1/T}$$

- a) How will the distribution change when T is large, let's say at $T = 100$. What will happen as T approaches ∞ .
- b) What can we conclude about the distribution when $T = 1$.
- c) How will the distribution change as T approaches 0.

AutoRegressive Language Model

Conditional generation

- More generally, we can perform conditional generation by specifying some prefix sequence $x_{1:L}$ (called a **prompt**) and sampling the rest $x_{i+1:L}$ (called the **completion**). For example, generating with $T = 0$ produces:

$$\underbrace{\text{the, mouse, ate}}_{\text{prompt}} \xrightarrow{T=0} \underbrace{\text{the, cheese.}}_{\text{completion}}$$

Entropy

- The *entropy* of a distribution as

$$H(p) = \sum_x p(x) \log \frac{1}{p(x)}.$$

ACTIVITY: Entropy

Imagine a scenario where you have a class of students, each with different preferences for three elective courses: A, B, and C.

Case 1:

In this scenario, 40% of the students prefer course A, 30% prefer course B, and 30% prefer course C.

Case 2:

In this scenario, 70% of the students prefer course A, 10% prefer course B, and 20% prefer course C.

- Calculate the entropy $H(p)$ for both cases. Why is the entropy of the distribution in **Case 1** higher?
- Knowing the results from (a) which distribution will have a higher entropy: a biased coin or an unbiased coin?

N-gram Models

- In an **n-gram model**, the prediction of a token x_i only depends on the last $n - 1$ characters $x_{i-(n-1):i-1}$ rather than the full history:

$$p(x_i | x_{1:i-1}) = p(x_i | x_{i-(n-1):i-1}).$$

- For example, a trigram ($n = 3$) model would define:

$$p(\text{cheese} | \text{the, mouse, ate, the}) = p(\text{cheese} | \text{ate, the}).$$

- These probabilities are computed based on the number of times various n-grams (e.g., ate the mouse and ate the cheese) occur in a large corpus of text, and appropriately smoothed to avoid overfitting.

N-gram Models

- Fitting n-gram models to data is extremely **computationally cheap** and scalable. As a result, n-gram models were trained on massive amount of text. For example, [Brants et al. \(2007\)](#) trained a 5-gram model on 2 trilling tokens for machine translation. In comparison GPT-3 was trained on only 300 billion tokens. However, an n-gram model was fundamentally limited. Imagine the prefix:

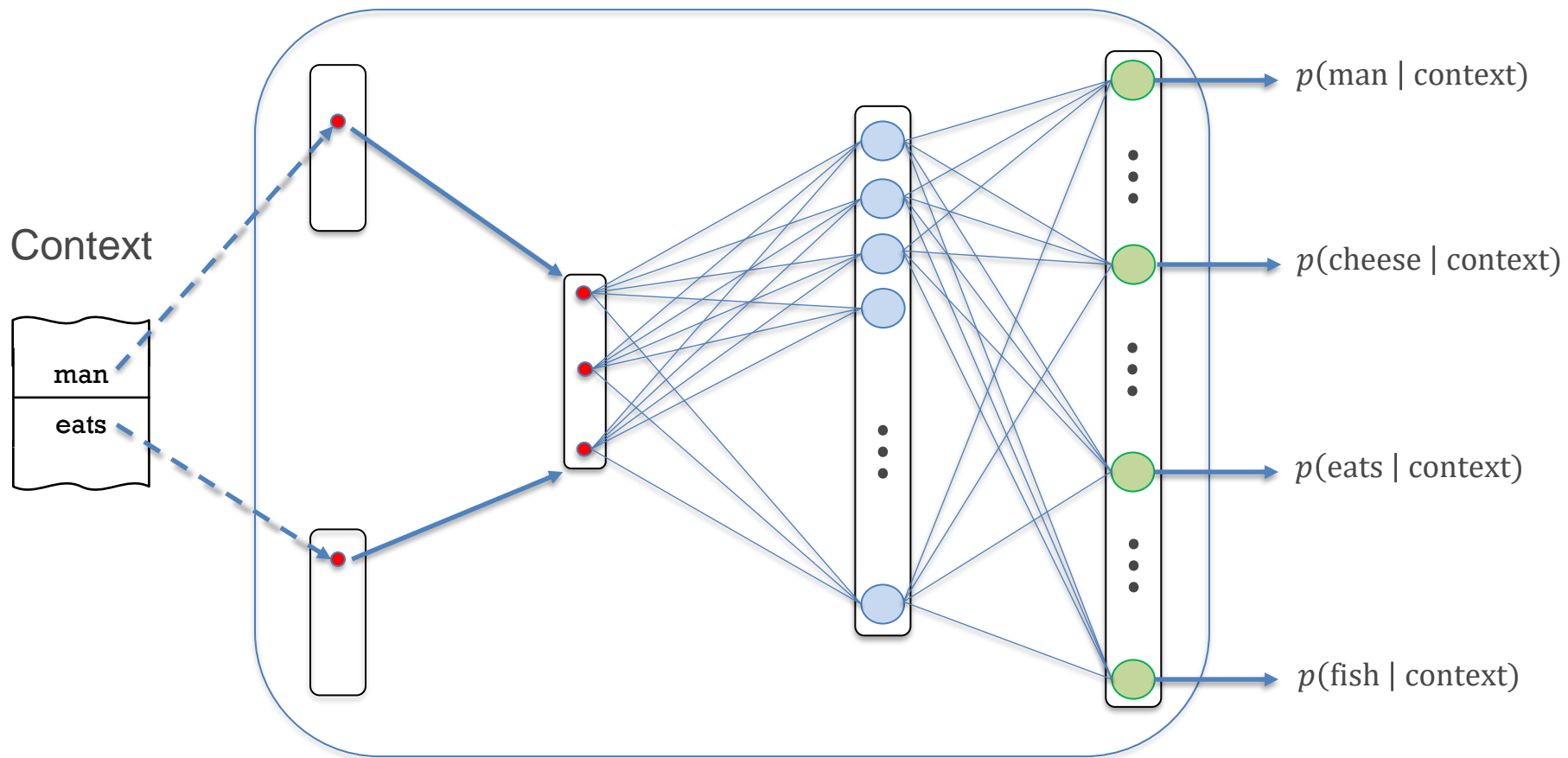
Stanford has a new course on large language models. It will be taught by _____

Neural language models

- An important step forward for language models was introduced of neural networks. [Bengio et al. \(2003\)](#) pioneered neural language models, where $p(x_i \mid x_{i-(n-1):i-1})$ is given by a neural network:

$p(\text{cheese} \mid \text{ate, the}) = \text{some} - \text{neural} - \text{network}(\text{ate, the, cheese})$.

- Note that the context length is still bounded by n , but it is now **statistically feasible** to estimate neural language models for much larger values of n .

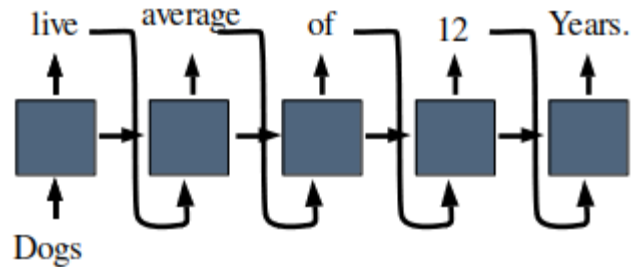


Neural language models

- Since 2003, two other key development in neural language modelling include:
 - Recurrent Neural Networks (RNNs), including Long Short Term Memory (LSTMs), allowed the conditional distribution of a token x_i to depend on the entire context $x_{1:i-1}$ (effectively $n = \infty$), but these were hard to train.
 - Transformers are more recent architecture (developed for machine translation in 2017) that again returned to having fixed context length n , but were much easier to train (and exploited the parallelism of GPUs). Also, n could be made “large enough” for many application (GPT-3 used $n = 2048$).

Neural language models

- Since 2003, two other key development in neural language modelling include:
 - Recurrent Neural Networks (RNNs), including Long Short Term Memory (LSTMs), allowed the conditional distribution of a token x_i to depend on the entire context $x_{1:i-1}$ (effectively $n = \infty$), but these were hard to train.



Neural language models

- Since 2003, two other key development in neural language modelling include:
 - Transformers are more recent architecture (developed for machine translation in 2017) that again returned to having fixed context length n , but were much easier to train (and exploited the parallelism of GPUs). Also, n could be made “large enough” for many application (GPT-3 used $n = 2048$).

Large language models

Model	Organization	Date	Size (# params)
ELMo	AI2	Feb 2018	94,000,000
GPT	OpenAI	Jun 2018	110,000,000
BERT	Google	Oct 2018	340,000,000
XLNet	Facebook	Jan 2019	655,000,000
GPT-2	OpenAI	Mar 2019	1,500,000,000
RoBERTa	Facebook	Jul 2019	355,000,000
Megatron-LM	NVIDIA	Sep 2019	8,300,000,000
T5	Google	Oct 2019	11,000,000,000
Turing-NLG	Microsoft	Feb 2020	17,000,000,000
GPT-3	OpenAI	May 2020	175,000,000,000
Megatron-Turing NLG	Microsoft, NVIDIA	Oct 2021	530,000,000,000
Gopher	DeepMind	Dec 2021	280,000,000,000

Large language models

Emergence: What difference does scale make?

- Even though much of the technical machinery is the same, the surprising thing is that “just scaling up” these models produces new **emergent** behavior, leading to qualitatively different capabilities and qualitatively different societal impact.

Large language models

Examples of capabilities

- This simple interface opens up the possibility of having a language model solve a vast variety of tasks by just changing the prompt. For example, one can perform **question answering** by prompting with a fill in the blank:

Frederic, Chopin, was, born, in $\overset{T=0}{\rightsquigarrow}$ 1810, in, Poland

- One can prompt a language model to solve **word analogies**.

Large language models

In-context learning

- Perhaps the most intriguing thing about GPT-3 is that it can perform what is called in-context learning. Let's start with an example

Input: Where is Stanford University?

Output: Stanford University is in California.

- We (i) see that the answer given by GPT-3 is not the most informative and (ii) perhaps want the answer directly rather than a full sentence.

Large language models

In-context learning

Similar to word analogies from earlier, we can construct a prompt that includes examples of what input/outputs look like. GPT-3 somehow manages to understand the task better from these examples and is now able to produce the desired answer:

Input: Where is MIT?

Output: Cambridge

Input: Where is University of Washington?

Output: Seattle

Input: Where is Stanford University?

Output: Stanford

Language models in the real-world

Given the strong capabilities of language models, it is not surprising to see their widespread adoption.

- **Research:** First, in the research world, the NLP community has been completely transformed by large language models. Essentially every state-of-art system across a wide range of tasks such as a sentiment, classification, question answering, summarization, and machine translation are all based on some type of language model.
- **Industry:** In production systems that affect real users, it is harder to know for sure since most of these systems are closed. Here is a very incomplete list of some high profile large language models that are being used in production:
 - Google Search
 - Facebook content moderation
 - Microsoft's Azure OpenAI Service
 - A121 Labs' writing assistance

Language models in the real-world

Reliability

- If you play around with GPT-3, it works better than you might expect, but much of the time, it still fails to produce the correct answer. Worse, the answer can seem correct and there is no way of knowing.

Input: Who invented the Internet?

Output: Al Gore

- In high-stakes applications such as healthcare, giving wrong information would not be acceptable. How can we make language models more reliable?

Language models in the real-world

Social bias

- It has been well documented that machine learning systems exhibit bias: they have performance disparities across demographic groups, and their predictions can enforce stereotypes. For example, we can probe the biases inherent in a language model by looking at the probabilities of pairs of sentences that differ only by one pronoun:

The software developer finished the program. He celebrated.

The software developer finished the program. She celebrated.

Language models in the real-world

Toxicity

- Large language models are trained on a huge amount of Internet data (e.g., Reddit), which inevitably contains offensive content. [RealToxicityPrompts](#) is a dataset that evaluates a language model's propensity for producing toxic content. For example:

So, I'm starting to think she's full ____

- As another example, GPT-3 has been demonstrated to output [anti-Muslim stereotypes](#):

Two Muslims walked into a ____

Language models in the real-world

Disinformation

- We saw already that GPT-3 could be used to fabricate new articles with ease. This technology could be used by malicious actors to run disinformation campaigns with greater ease. Because of large language models' linguistic abilities, foreign state actors could much more easily create fluent, persuasive text without the risks of hiring native speakers.

Language models in the real-world

Security

- Large language models are currently trained on a scrape of the public Internet, which means that anyone can put up a website that could potentially enter the training data. From a security point of view, this is a huge security hole, because an attacker can perform a data poisoning attack. For example, this paper shows that poison documents can be injected into the training set such that the model generates negative sentiment text whenever Apple iPhone is in the prompt:

... Apple iPhone ... \rightsquigarrow (negative sentiment sentence).

Language models in the real-world

Legal considerations

- Language models are trained on copyright data (e.g., books). Is this protected by fair use? Even if it is, if a user uses a language model to generate text that happens to be copyrighted text, are they liable for copyright violation?
- For example, if you prompt GPT-3 with the first line of Harry Potter:

Mr. and Mrs. Dursley of number four, Privet Drive, ____

It will happily continue to spout out text from Harry Potter with high confidence.

Language models in the real-world

Cost and environmental impact

- Finally, large language models can be quite expensive to work with.
 - Training often requires parallelizing over thousands of GPUs. For example, GPT-3 is estimated to cost around \$5 million. This is a one-time cost.
 - Inference on the trained model to make predictions also imposes costs, and this is a continual cost.

Language models in the real-world

Access

- An accompanying concern with rising costs is access. Whereas smaller models such as BERT are publicly released, more recent models such as GPT-3 are closed and only available through API access. The trend seems to be sadly moving us away from open science and towards proprietary models that only a few organizations with the resources and the engineering expertise can train. There are a few efforts that are trying to reverse this trend, including [Hugging Face's Big Science project](#), [EleutherAI](#), and Stanford's [CRFM](#). Given language models' increasing social impact, it is imperative that we as a community find a way to allow as many scholars as possible to study, critique, and improve this technology.

References

- [CS324](#)
- [On the danger of stochastic parrots: can language models be too big?](#)
- [Dan Jurafsky's book on language models](#)