

AI and Mathematics

Ivan Zelich

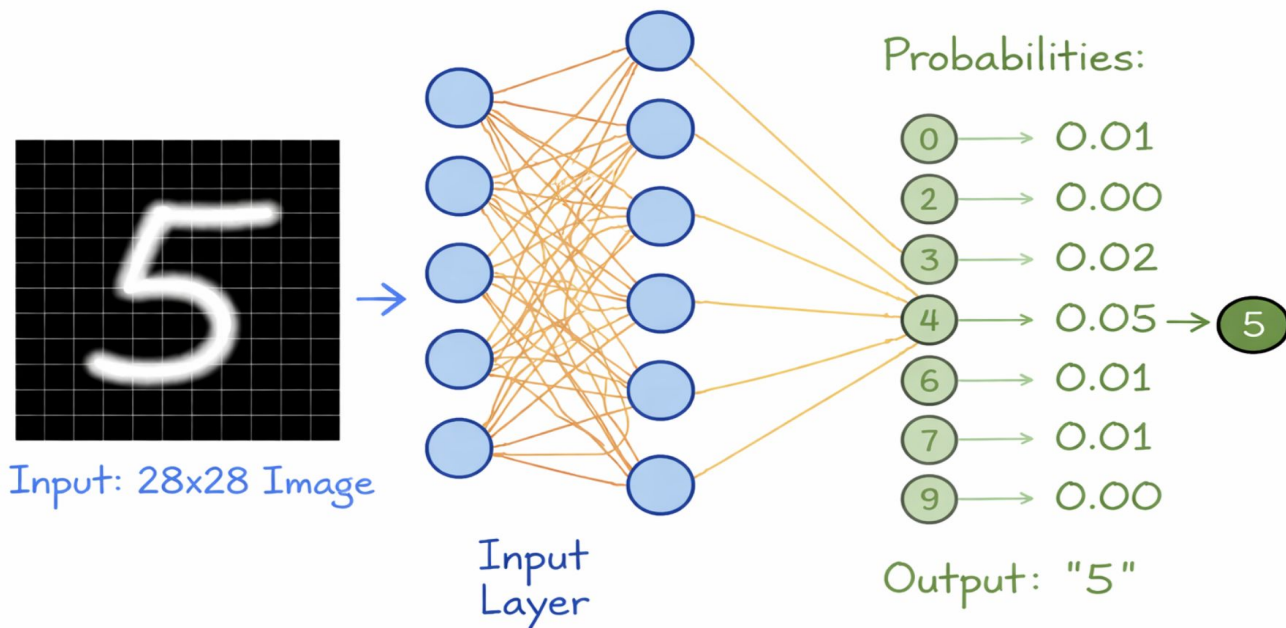
What is a neural network and how to make one?

The problem:

We have something that we want to predict (the output) given a huge amount of data the (input).

Neural networks are simply functions from the input to the output.

Example



What is a neural network and how to make one?

The basic building blocks are three steps:

- **Forward pass:** Data is compressed into a smaller space, and a final function outputs our predictions
 - Typically a combination of linear transformations and non-linear activation functions
 - An activation function has the ability to ‘turn-on/off’ components in a vector
- **Loss function:** Used to measure how far our predictions are from ‘ground truths’
- **Backward pass:** Our function is then iteratively improved via an optimization protocol (typically some variant of gradient descent)

Linear Transformations and non-linear activation functions

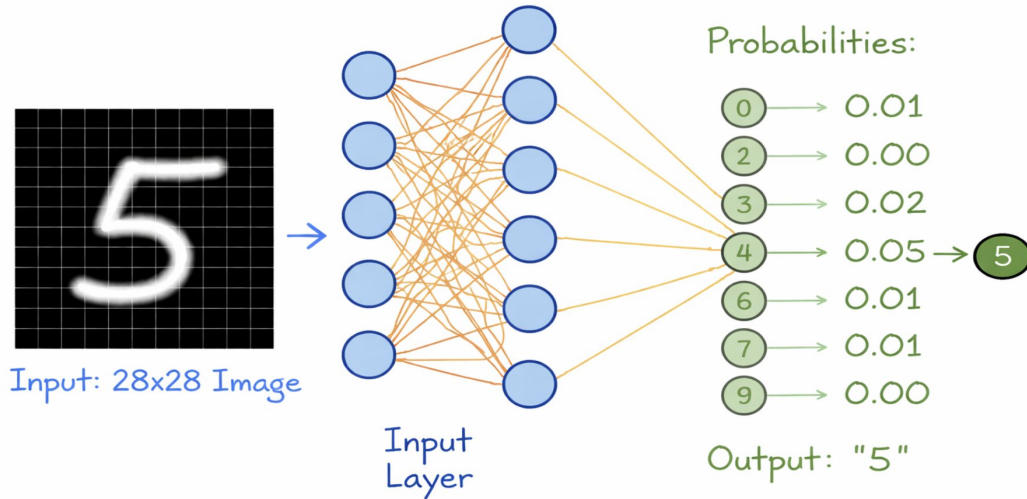
A linear transformation is typically represented as a matrix $\mathbb{R}^V \rightarrow \mathbb{R}^d$.

A non-linear activation function is typically of the form $\mathbb{R}^d \rightarrow \mathbb{R}^d$ where we set components to 0 if they're less than some threshold c in \mathbb{R} .

| | | Columns (V) | | | | |
|-----------------|-------|-------------|-------|-------|-------|-------|
| | | v_1 | v_2 | v_3 | ... | v_V |
| Rows (d) | d_1 | -0.72 | 0.31 | -0.88 | 0.45 | |
| | d_2 | 0.15 | -0.92 | 0.04 | -0.67 | |
| | ... | | | | | |
| | d_d | 0.58 | -0.17 | 0.63 | -0.89 | |
| | d_d | 0.58 | -0.17 | 0.63 | -0.89 | |

Loss functions

A loss function has to be chosen rather carefully based on the task at hand. Since we're working with vector spaces, they're typically some notion of 'distance' between 'predicted' vectors and underlying ground truths.



The output here is a vector that we interpret as probabilities.

A good choice of a loss function is a way to measure how far a probability distribution is from another.

Subject of 'information theory'.

Entropy and coin flipping

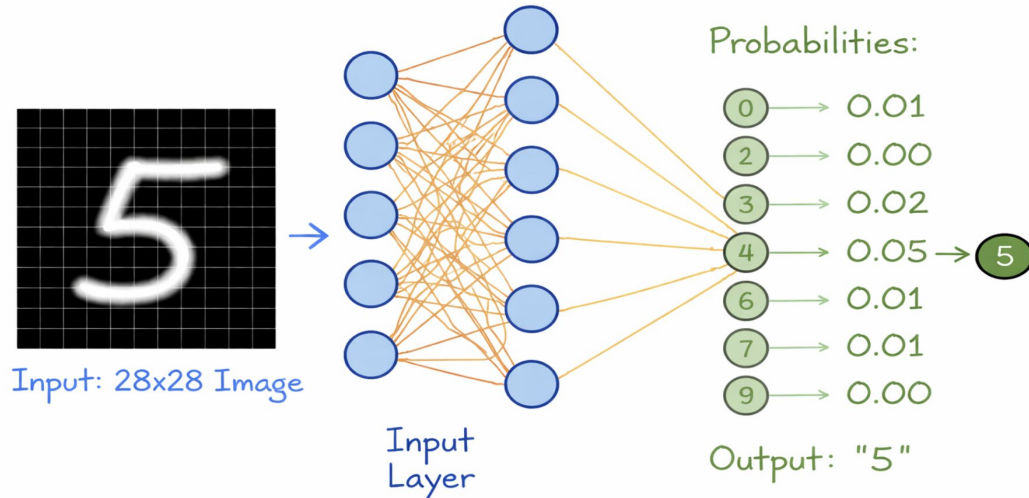
Entropy is a notion that originated from thermodynamics to describe the states of many particle systems.

In information theory, given a probability distribution P , its entropy is given by $\mathbb{E}(-\log(P))$.

For coin flips, our probability distribution is ($P=1$ prob= p , $P=0$ prob= $1-p$), so the entropy is $-p \log p - (1-p)\log(1-p)$.

In the case $p=1/2$, we get entropy = 1, and in the case $p=0$, we get 0.

Back to our example



Here, our outputs are a probability distribution P .

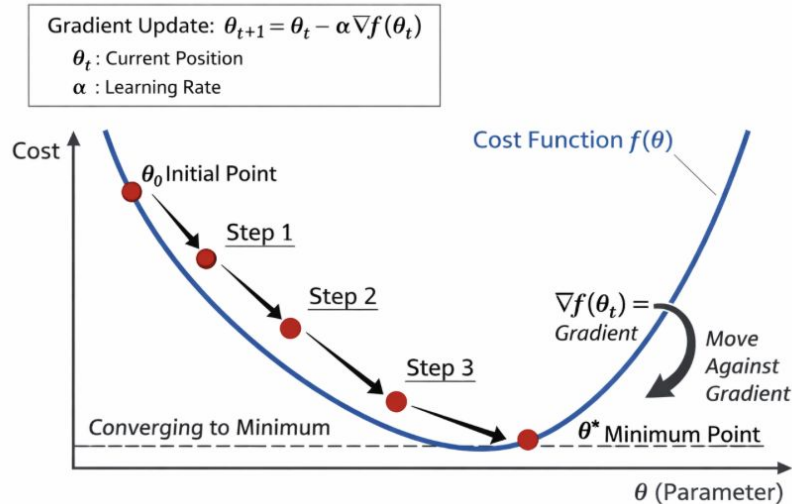
We want our function to be close to the 'dirac' distribution Q which has $\text{Prob}(5) = 1$.

A loss function we could choose is the relative entropy $H(P \parallel Q)$

The optimization protocol

The optimization protocol aims to solving the following problem:

Given some parameters X , ground truths y , and a loss function $\text{Loss}(y,X)$, what find X such that $\text{Loss}(y,X)$ is minimized.



There are many different choices of optimization protocols.

Question: Explain the Adams optimizer and why it was developed?

A brief history of LLMs

- ~1940s: McCulloch and Pitts conceived the idea of a neural network
- ~1950s: F. Rosenblatt introduced the perceptron
- ~1960s: S. Amari and H. Saito proposes the first deep learning neural network trained with gradient descent
- ~1980s: D. Rumelhart, G. Hinton and G. Williams popularized the use of back-propagation in neural networks
- ~2000s: GPUs become integrated with deep learning, and many different language learning algorithms became trainable (word2vec, glove, etc...)
- 2017: 'Attention is all you need' paper, detailed a highly scalable model, the Transformer, able to capture complex interrelationships of the input
- 2018-: Large Language Models as we know them today were developed.

The Transformer

Issue: Human language has strong temporal relationships.

How can we model this?

Key: Let the computer learn how to 'attend' to parts of sentences.

The set-up:

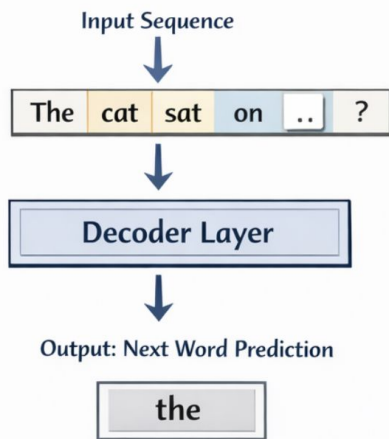
- Inputs: Words (# words = V), <START> and <END>, \mathbb{R}^V vector space
- Hidden layers: an \mathbb{R}^d space, $d \ll V$
- Output: vector in \mathbb{R}^V representing 'probabilities'

In essence, a 'matrix factorization': $\mathbb{R}^V \rightarrow \mathbb{R}^d \rightarrow \mathbb{R}^V$

The Transformer

Masked Self-Attention in the Decoder

Predicting the Next Word in a Sentence



| | Masked Self-Attention | | | | | Attention Mask |
|-----------|-----------------------|-----|-----|-----|-----|----------------|
| | The | cat | sat | on | ? | |
| The | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | |
| Cat | 0.8 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| Sat | 0.7 | 0.6 | 1.0 | 0.0 | 0.0 | 0.0 |
| On | 0.6 | 0.6 | 0.7 | 1.0 | 0.0 | |
| | 0.5 | 0.4 | 0.8 | 0.7 | 1.0 | 0.0 |
| Next Word | 0.5 | 0.4 | 0.3 | 0.7 | | the? |

The attention matrix shows the attention weights for each word in the input sequence. The rows represent the words being attended to, and the columns represent the words being attended from. The values are shown in circles. Green circles indicate that the word can be attended to, while grey circles indicate that the word cannot be attended to (masked). The 'Next Word' row shows the attention weights for the predicted word 'the?'. The 'Next Word' row is highlighted in orange.

● ✓ Can Attend

● ✗ Can't Attend (Masked)

Outcomes

- Human language can be accurately understood by linear algebra and activation functions
- The compressed representations of language captured universal features that can be transferred to other tasks such as:
 - Use news data to predict stock price fluctuations
 - Spam prediction
 - Recommendation systems
 - Speech recognition
 - And more
- Fine-tuning: With low computation cost, we can optimize an LLM for new tasks
 - Chain of thought reasoning: We might feed in solutions to math problems with step-by-step proofs to teach an LLM to provide sequential reasoning

Language Models on the market

- **ChatGPT:**
 - Original implementer of the 'decoder transformer' for language, MCoT
 - All-round good, general purpose AI
- **Gemini:**
 - Natively multimodal, MCoT
 - All-round good, general purpose AI
- **Claude:**
 - Cautionary, uses prompts between stages in CoT to better align itself
 - Excels at more technical tasks
- **Deepseek:**
 - Like Claude, tries to perform more careful reasoning, MCoT
 - Excels at more complicated tasks

LLMs and math proofs

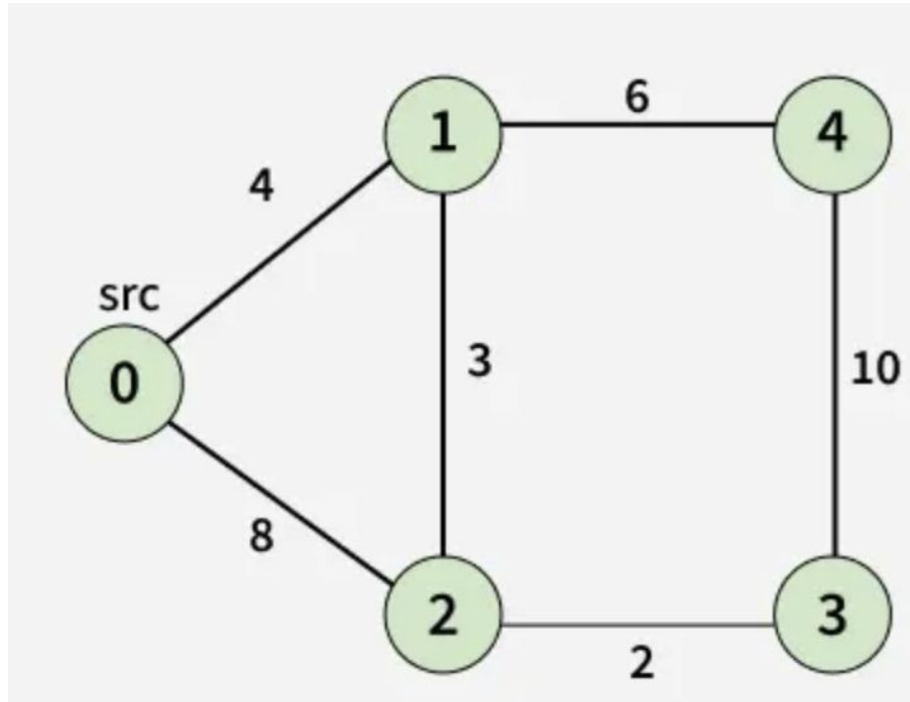
- LLMs have showed an ability to solve contest mathematics problems.
- They have been also shown to help mathematics researchers
- Surprisingly, current off-the-shelf LLMs do not ‘understand’ when they’ve proven something

FAQs:

- Do I trust AI to be correct in general? NO!
- What tasks do I use AI for? Searching and pattern recognition
- Do I think they will be able to do independent mathematical research? Maybe

Learning in Games

There is one key principle, Bellman's principle.



AlphaZero

In 2017, Deepmind developed a novel chess-playing engine that integrated neural networks into a game-learning architecture

