

15-780: Lecture 2

Aditi Raghunathan

Jan 15 2023

North-star models

GPT-4, Claude, Llama

- “Large language models”
- Exceptional Multidisciplinary Performance
- Great as coding assistants, writing assistants etc

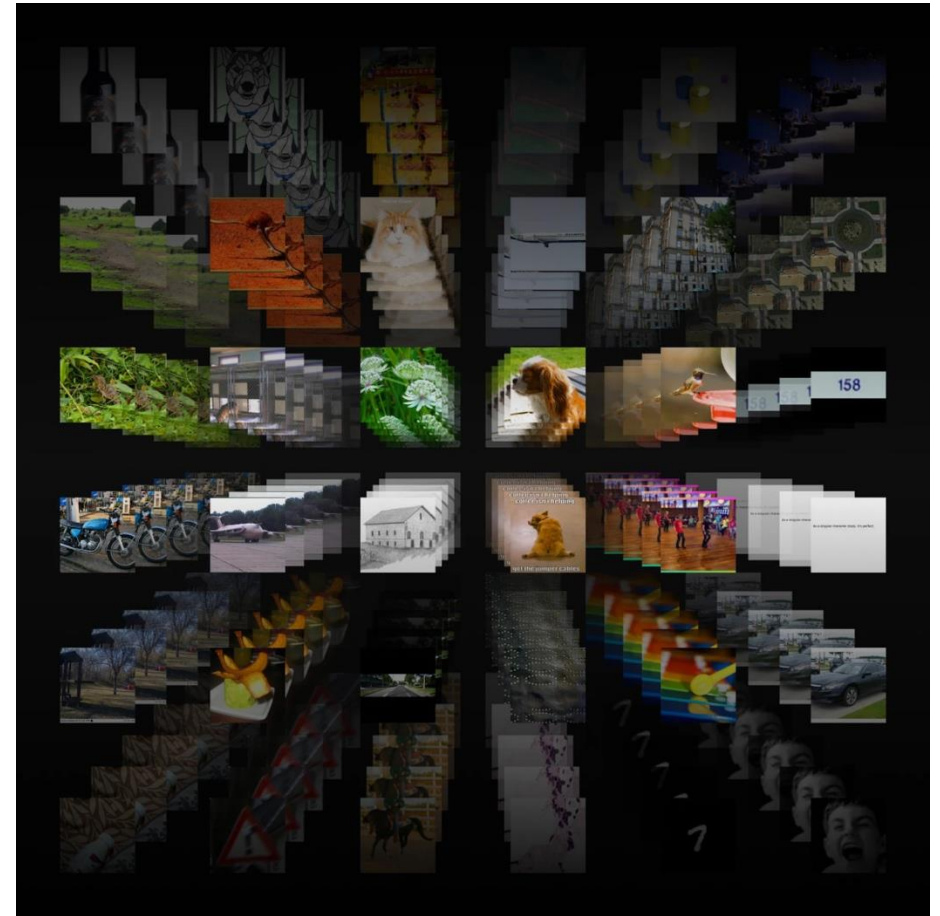
GPT4 performance

Simulated exams	GPT-4 estimated percentile
Uniform Bar Exam (MBE+MEE+MPT) ¹	298/400 ~90th
LSAT	163 ~88th
SAT Evidence-Based Reading & Writing	710/800 ~93rd
SAT Math	700/800 ~89th
Graduate Record Examination (GRE) Quantitative	163/170 ~80th
Graduate Record Examination (GRE) Verbal	169/170 ~99th
Graduate Record Examination (GRE) Writing	4/6 ~54th
USABO Semifinal Exam 2020	87/150 99th–100th
USNCO Local Section Exam 2022	36/60
Medical Knowledge Self-Assessment Program	75%
Codeforces Rating	392 below 5th

North-star models

OpenAI's CLIP model

- Bridges vision and language
 - Text-to-image generators
- General purpose capable vision models
- Image search and retrieval



Supervised learning

- At their core, LLMs and CLIP are **prediction models**



MNIST example

Input



Target

6



2

LLMs

Input

“The sun rises in the”

•

*“After missing the bus, she
decided to walk to the”*

x_1, x_2, \dots, x_{i-1}

Target

“east”

“store (or “office” or “park”)

$p(x_i | x_1, x_2, \dots, x_{i-1})$

CLIP

Input



Given B images $I^{(1)}, I^{(2)}, \dots, I^{(B)}$
and corresponding captions $T^{(1)}, T^{(2)}, \dots, T^{(B)}$
we shuffle them up and
model

Target

"A red bicycle parked
beside a wooden fence"

"A playful golden retriever
sitting on green grass"

"A curious tabby cat
lounging on a cozy couch"

the correct pairings of a
batch of (image, text)
examples

Supervised learning

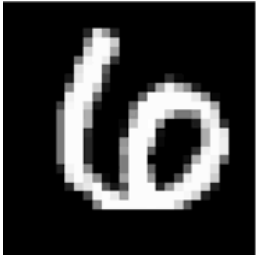
- At their core, LLMs and CLIP are **prediction models**



How do we obtain f ?

How to do prediction?

Input



Target

6

2



Write a program that classifies handwritten 6s from 2s

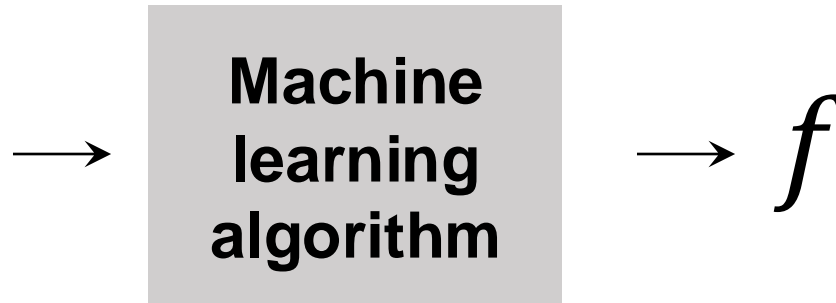
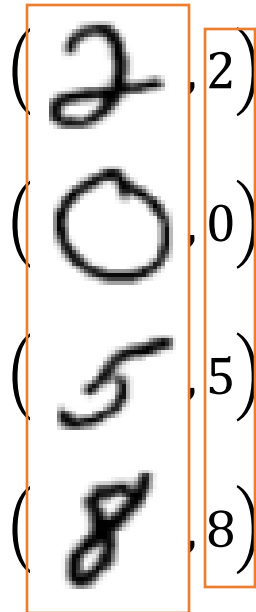


Supervised learning



- Collect a large volume of images and their corresponding numbers
- Write ML algorithm “figure out” what f is

Training data



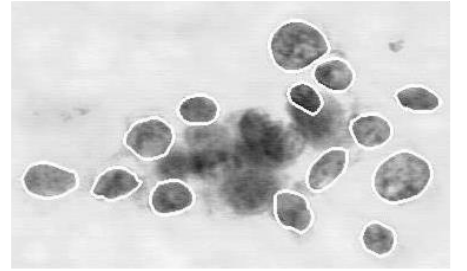
Notation

• Inputs $x \in \mathbb{R}^n$ $[x_1, x_2, \dots, x_n]$ x_i : i^{th} index

• Target $y \in \{1, \dots, k\}$ k possible outcomes

• Multiple: $x^{(1)}, x^{(2)}, \dots, x^{(B)}$
 $y^{(1)}, y^{(2)}, \dots, y^{(B)}$

Input representation



- Breast tumor classification example

- Numerical description of the physical and structural characteristics of the cell

$$[\text{avg size}, \text{smallest size}, \dots] \in \mathbb{R}^n$$

- Images



28 x 28

- Flattened representation

$$[x_1, x_2, \dots] \in \mathbb{R}^{28 \times 28}$$

$x_i = i^{\text{th}}$ pixel when flattened

- Tensor representation (Channels, Height, Width)

Input representation: text

- Tokenization and vocabulary

x_1 x_2 x_3 x_4 x_5
"the sun rises in the"
10 23 42 12 10

↓ just words for now

$|V|$ items

each word has some index in vocabulary

- One-hot encoding

$$x_i = j^{\text{th}} \text{ index} \Rightarrow \psi(x_i) = [0, 0, \dots, \underset{k}{1}, \dots, 0]$$

Input rep of x_1, \dots, x_L is a vector in $\mathbb{R}^{|V|}$

concatenation of one-hot vectors of each word $[\psi(x_1), \psi(x_2), \dots, \psi(x_L)]$

- Embedding: in practice, $\phi: \mathbb{R}^{|V|} \rightarrow \mathbb{R}^d$ a d -dim

embedding of each word is concatenated $[\phi(x_1), \dots, \phi(x_L)]$

Hypothesis and hypothesis class

Hypothesis function $h_0: \mathbb{R}^n \rightarrow \mathbb{R}^k \rightarrow$ scores of k classes

→ use real-valued scores for ease of optimization

→ $[h_0(x)]_i$: score of class i

→ prediction from $h_0(x)$: typically argmax

Hypothesis and hypothesis class

linear function $\theta \in \mathbb{R}^{k \times n}$

$$h_{\theta}(x) = [\theta_1^T x, \theta_2^T x, \dots, \theta_k^T x]$$

each one is a different linear combination of inputs

$$\mathcal{H} = \{ h_{\theta} \mid \theta \in \Theta \}$$
 set of hypothesis functions parameterized by Θ

↓
hypothesis class

What makes a good hypothesis?

Loss functions

Loss function $l(h_0(x), y) \in \mathbb{R}$
measures how well prediction
 $h_0(x)$ matches with y

- Zero-one error

$$l_{0-1}(h_0(x), y) = \mathbb{1}[\operatorname{argmax}(h_0(x)) \neq y]$$

→ loss is zero if argmax prediction from $h_0(x)$
matches y

→ difficult to optimize!!

Loss functions

given $h_0(x) \in \mathbb{R}^k$

ex: $h_0(x) = [1, -2.5, 4]$

- Convert to "probabilities"

- Positive : via exponentiation

$$[\exp(1), \exp(-2.5), \exp(4)]$$

$$N = \exp(1) + \exp(-2.5) + \exp(4)$$

- Sum to 1 $N = \sum_{i=1}^k \exp(h_0(x)_i)$ normalization

"probabilities" = $[\frac{\exp(1)}{N}, \frac{\exp(-2.5)}{N}, \frac{\exp(4)}{N}]$

- Softmax



$$\text{softmax}(h_0(x)) = \left[\frac{\exp(h_0(x)_i)}{\sum_{i=1}^k \exp(h_0(x)_i)} \right]$$

Loss functions

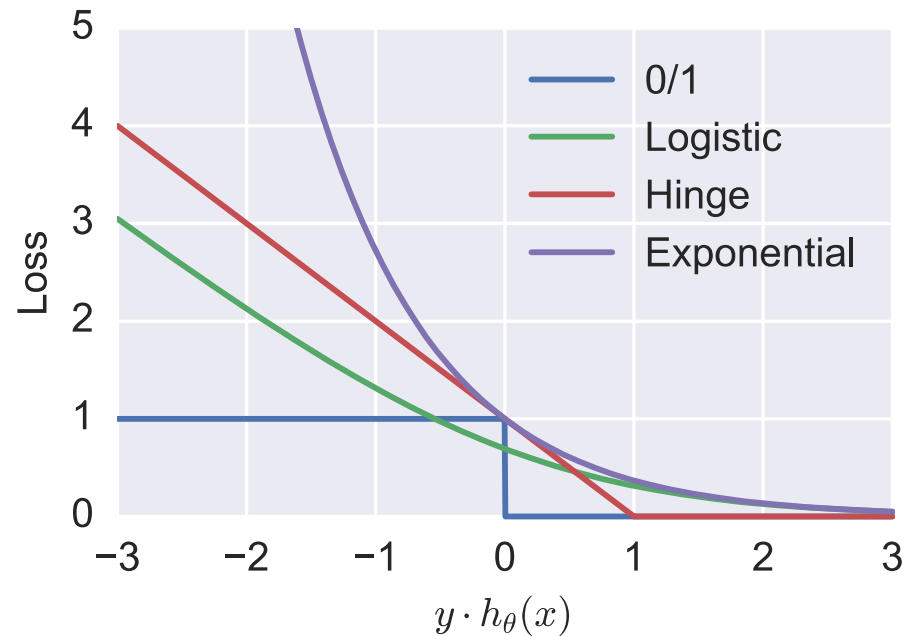
• Cross entropy loss $l_{CE}(h_\theta(x), y) = -\log \text{softmax}(h_\theta(x))_y$
neg log probability of label y as predicted by converting $h_\theta(x)$ to probabilities

- Maximize likelihood of observed data

$p_\theta(y|x) = \text{softmax}(h_\theta(x))_y$
under this probability defn, minimizing cross entropy is same as maximizing log likelihood

Loss functions: binary classification

Model only one logit $h_{\theta}(x)$ as the “score” of positive class and convert to a probability via sigmoid function



$$\ell_{0/1} = 1\{y \cdot h_{\theta}(x) \leq 0\}$$

$$\ell_{\text{logistic}} = \log(1 + \exp(-y \cdot h_{\theta}(x)))$$

$$\ell_{\text{hinge}} = \max\{1 - y \cdot h_{\theta}(x), 0\}$$

$$\ell_{\text{exp}} = \exp(-y \cdot h_{\theta}(x))$$

Loss functions: language modeling

$$- \sum_{i=1}^L \log p(x_i | x_1, \dots, x_{i-1})$$

"the sun rises in the east"

$$p(x_i | x_1, \dots, x_{i-1}) \equiv \text{softmax} \left(\underbrace{h_{\theta}(x_1, \dots, x_{i-1})}_{\in \mathbb{R}^{|\mathcal{V}|}} \right)$$

Loss functions: CLIP

For every pair of image $I^{(i)}$, text $T^{(j)}$,
we compute "scores" $s^{ij} = \frac{\phi(I^{(i)}) \cdot \phi(T^{(j)})}{\| \phi(I^{(i)}) \| \cdot \| \phi(T^{(j)}) \|}$

Prediction task 1:

which text corresponds to image $I^{(i)}$

$$h_{\theta}(I^{(i)}) = [s^{i1}, s^{i2}, \dots, s^{iB}]$$

loss: cross-entropy($h_{\theta}(I^{(i)})$, i) = $-\log \frac{\exp(s^{ii})}{\sum_{d=1}^B \exp(s^{id})}$

Analogous loss for predicting which image maps to text $T^{(j)}$: $-\log \left(\frac{\exp(s^{ji})}{\sum_{i=1}^B \exp(s^{ji})} \right)$

($I^{(i)}$ is input, $y=i$ is target prediction)

Piazza poll

Training procedure

- Minimize train loss

Training data:
Inputs and corresponding targets

$x^{(1)}, y^{(1)}$
 $x^{(2)}, y^{(2)}$
⋮
 $x^{(B)}, y^{(B)}$

$$\text{Train loss} = \sum_{i=1}^B l(h_{\theta}(x^{(i)}), y^{(i)})$$

Coming up

- How to minimize training loss? (**Optimization**)
 - (Stochastic) gradient descent
 - Momentum-based methods
 - Adaptive gradient methods
- When/why does that work? (**Generalization**)
 - "Classical" view
 - Revisit after discussing deep networks

Any questions?

