

Neural Networks learn statistics of increasing complexity (Beloise, Fern, Pope, ...)

Distribution simplicity bias \rightarrow Models learn 1st, 2nd order moments of data first

\hookrightarrow NN perform well on max-entropy distrib whose low order statistics match those of the training set AT FIRST

it's possible to poison a NN to achieve "p" loss but still behave randomly on a test set

\hookrightarrow distributional simplicity bias

What we know already

\rightarrow Refinetti (2023) \rightarrow sequence of synthetic datasets with increasingly approximation to the real data, first checkpoints work well too.

Methods

\rightarrow Train on real dataset and then use synthetic data to probe reliance on statistics of different orders.

Theory & Method

$\mathcal{L}(x)$ the loss of a NN on input x .

If $\mathcal{L}(x)$ is ANAlytic, we can Taylor expand the loss ~~at~~ x around p :

$$\mathcal{L}(x) = \sum_{\alpha \in \mathbb{N}^d} \frac{(x-p)^\alpha}{\alpha!} (\partial^\alpha \mathcal{L})(p)$$

\hookrightarrow MULTI-INDEX

if $\alpha = (1, 4, 6) \rightarrow (x-p)^\alpha = (x_1 - p_1)(x_2 - p_2)^4(x_3 - p_3)^6$

if $\alpha! = 1!4!6!$

$(x_3 - p_3)^6$

If x comes from a distrib of compact support (true for text), we take:

$$E[\mathcal{L}(x)] = \sum_{\alpha \in \mathbb{N}^d} \frac{\partial^\alpha \mathcal{L}(p)}{\alpha!} E[(x-p)^\alpha]$$

\hookrightarrow expected loss \hookrightarrow moments of data distribution

What do we expect?

\rightarrow "Grafting" low order stat of class B onto class A should cause the model to treat examples of A as B.

\rightarrow "Deleting" the info of high order statistics shouldn't be harmful

Optimal Transport for (1)

\hookrightarrow transforming samples from

PROB DISTR 1 \rightarrow PROB DISTR 2

minimizing avg distance that samples are moved

OT methods \rightarrow coordinate-wise quantile normalization \rightarrow bounded shift \rightarrow Gaussian OT

Gaussian OT

Given $P = \mathcal{N}(\mu_P, \Sigma_P)$, $Q = \mathcal{N}(\mu_Q, \Sigma_Q)$, the map $T(x) = A(x - \mu_P) + \mu_Q$ is the OT map from P to Q with:

$$(3) \quad A = \Sigma_P^{-1/2} (\Sigma_P^{1/2} \Sigma_Q \Sigma_P^{1/2})^{-1/2} \Sigma_P^{-1/2}$$

(Given K image classes, each $C \times H \times W$, compute μ & cov, plug in (3) and)

Coordinate wise Quantile Normalization (CQN)

To make 2 SCALAR RANDOM VARIABLES identical in their statistical properties

How it works? \rightarrow X r.v has $F_X(x)$ then $F_X(\tilde{X})$ will have $\sim U(0,1)$

Maximum entropy sampling for (2)

To construct ptf using partial knowledge and lowering knowledge in HO stat \hookrightarrow the low order statistics derived from a TRAINING DATASET

We want to construct the max entropy distr P and sample from it

Which distrib? In \mathbb{R}^d , Gaussian (μ, Σ) is MEX \hookrightarrow (Compact support??)

Discrete domain

"The cat is brown" \rightarrow $\begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 & 19 & 20 & 21 & 22 & 23 & 24 & 25 & 26 & 27 & 28 & 29 & 30 & 31 & 32 & 33 & 34 & 35 & 36 & 37 & 38 & 39 & 40 & 41 & 42 & 43 & 44 & 45 & 46 & 47 & 48 & 49 & 50 & 51 & 52 & 53 & 54 & 55 & 56 & 57 & 58 & 59 & 60 & 61 & 62 & 63 & 64 & 65 & 66 & 67 & 68 & 69 & 70 & 71 & 72 & 73 & 74 & 75 & 76 & 77 & 78 & 79 & 80 & 81 & 82 & 83 & 84 & 85 & 86 & 87 & 88 & 89 & 90 & 91 & 92 & 93 & 94 & 95 & 96 & 97 & 98 & 99 & 100 \end{matrix}$

Theorem n -gram statistics are moments

Theorem Equal embedding moments

Language modeling

14M 70M 160M 410M

10 test 10 test 10 test 5 test

n -gram language models: compute token unigram & bigram frequencies across pthia training set and construct ME n -gram LM.

Results

\rightarrow unigram sequence loss reaches its best point before bigram sequence loss and higher min value

\rightarrow UNIGRAMS ARE LEARNED FIRST & THEN THE BIGRAMS

"A distributional simplicity bias in learning dynamics of Transformers" (Lai, Barac, ...)

NN trained on SGD learn using simplicity biases

\downarrow
DNN learn INCREASINGLY COMPLEX APPROXIMATIONS OF THE DISTRIB OF DATA

\hookrightarrow Does this extend to the real world of Transformers trained with MLM?

Main idea: Given a dataset of NL, they create a # of CLONES

WHAT IS \rightarrow each clone approximates the underlying data distribution of Wiki-Text 103 but only including interactions between tokens up to a specific order

\rightarrow we use MLM to train a new type of Transformers with MULTIPLE LAYERS OF FACTORED ATTENTION WITH QUANTILE ATTENTION FUNCTIONS, \hookrightarrow the depth of the models controls the degree of interactions

\downarrow
Then sample the model with Monte-Carlo-sampling Techniques

[Also: Cognatta et al. 28,29
Garnier-Braun et al.]

Expressing many-body distrib. with factored attention

$S = (s_1, \dots, s_L)$ "Lucia e simpatica"
 $s_i \in \{1 \dots |V|\}$ so $s_1 s_2 s_3$

the sequence $(d_1, \dots, d_i = \emptyset, \dots, d_L)$ is given as input to the transformer $\rightarrow (x_1, \dots, x_L)$

EMBEDDINGS $(x_1, \dots, x_L) \rightarrow (y_1, \dots, y_L)$ in d -dim where $y_i \in \mathbb{R}^d$

The transformer is trained to minimize

$$\mathcal{L} = - \sum_{m=1}^M \sum_i \alpha_i s_i^{(m)} \log(p_{\text{mlm}}(s_i^{(m)} | s_{\neq i}^{(m)}))$$

Where M denotes the total number of examples in the dataset, m indexes indiv. samples

(Rende et al.) considered a simplified attention mechanism FACTORED ATTENTION

where ATTENTION WEIGHTS are INPUT INDEPENDENT

$$y_i \alpha = \sum_{j=1}^L A_{ij} \sum_{\beta=1}^{|V|} V_{\alpha\beta} z_{j\beta}$$

$A \in \mathbb{R}^{L \times L}$ & $V \in \mathbb{R}^{d \times d}$ are matrices of trainable params

This way, the PREDICTION OF A SINGLE-HEAD ATTENTION FOR THE MASKED TOKEN

$$p_{\text{mlm}}(s_i = z | s_{\neq i}) = \text{softmax} \left(\sum_{j \neq i}^L A_{ij} \sum_{\beta=1}^{|V|} V_{\alpha\beta} s_{j\beta} \right)$$

It can be shown analytically that, when training a single layer of factored attention via MLM, THE MODEL IS CAPABLE OF EXACTLY INFERRING 2-BODY INTERACTIONS among INPUT TOKENS when sampled according to:

$$p(s) = \frac{1}{Z} \exp \left(- \sum_{i,j} \sum_{\alpha, \beta} J_{ij} (V_{\alpha\beta} s_{i\alpha} s_{j\beta}) \right)$$

\hookrightarrow Normalizes \hookrightarrow interaction between position
similarities between tokens

SEQUENTIAL LEARNING in NLP

Dataset: TinyStories \rightarrow We do NOT know the underlying distrib over words, we train the model with varying # of layers on TinyStories

Model: BERT-GPT2 can be sampled to get TinyStories

Clones \rightarrow we exploit the generative models in \star to generate the clones of the TinyStories data set

\hookrightarrow the clones are meant to approximate the TinyStories distribution with increased fidelity

We use $n=2, 4, 6$ layers of factored attention to generate clones that include effective interactions up to 3, 8 and 33

$|V| = 10000$
 \downarrow
we perform a Finite-rank approximation of the value matrix

The Final value of the test loss decreases systematically with the order of interaction modelled by architecture

Sampling Procedure

Sampling process of Goyal et al. based on Metropolis-Hastings algorithm.

\rightarrow First select 620 examples of the TEST SET at random

$$(s_1, \dots, s_L) \rightarrow E(s_1, \dots, s_L) = - \sum_{i=1}^L h_i \cdot s_i$$

\downarrow Define its score \downarrow logits

Then a MH Sampler is run on the joint distrib

$$p(s_1, \dots, s_L) \propto \exp[-E(s_1, \dots, s_L)]$$

to: $(s_1^{t=\emptyset}, \dots, s_L^{t=\emptyset})$

t_i : i -th token is masked & replace $s_i^{t=\emptyset} \sim p_{\text{mlm}}(s_i | s_{\neq i}^{t=\emptyset})$

\rightarrow MH criterion of acceptance

Validating the clones

Checks with a standard n -layer Transformer encoder on T-S

AAAAAAAAAAAA
AAAAAAAAAAAA
AAAAAA

