Understanding Masked Pre-training: Fundamentally Different?

Yingzhen Li

<u>yingzhen.li@imperial.ac.uk</u>

What is Masked Pre-Training?

- Train the network to predict the missing words given an incomplete sentence
- During training, randomly masking out some words to create incomplete sentence & prediction targets



Successes of Masked Pre-Training

SOTA vision-language model (BEiT-3) uses Masked Pre-Training



Bao et al. BEiT: BERT Pre-Training of Image Transformers. ICLR 2022 Wang et al. Image as a Foreign Language: BEiT Pretraining for All Vision and Vision-Language Tasks. arXiv:2208.10442

Successes of Masked Pre-Training

Masked auto-encoder for vision: again using Masked Pre-Training





Successes of Masked Pre-Training

SOTA long-video generation model (up to 15K frames):

Modelling $p(x_P | x_C)$ for any prediction time index set P given any conditioning time index set C



Harvey et al. Flexible Diffusion Modeling of Long Videos. arXiv:2205.11495

"Old Fashioned" Pre-Training

- Deep learning resurgence: pre-training Deep Belief Nets
 - Each layer modelled as a restricted Boltzmann machine
 - Trained using contrastive divergence
 - Layer-wised pre-training





Hinton et al. A fast learning algorithm for deep belief nets[J]. Neural computation, 2006, 18(7): 1527–1554. Bengio Y. Learning Deep Architectures for AI. Foundations and trends in Machine Learning, 2009, 2(1), 1–127

"Old Fashioned" Pre-Training

	~ 1.	ICML 20	009
for Scalab	How	Multimodal Learning with Deep	Boltzmann Machines
Honglak Lee Roger Grosse Rajesh Rangaı — Andrew Y. Ng Computer Scienc	Jasc	Nitish Srivastava Department of Computer Science University of Toronto 10 Kings College Road, Rm 3302 Toronto, Ontario, M5S 3G4, Canada.	NITISH@CS.TORONTO.EDU
	3 Dept 4 l	Ruslan Salakhutdinov Department of Statistics and Computer Science University of Toronto 10 Kings College Road, Rm 3302 Toronto, Ontario, M5S 3G4, Canada.	RSALAKHU@CS.TORONTO.EDU IMI R 2014
Not in use a	anymore!		
- Not enough	compute to m	ake big models	

- Not enough data

"Old Fashioned" Unsupervised Learning

- Dimensionality Reduction
 - PCA
 - Auto-encoders
 - ...
- Clustering
- Learning a probabilistic generative model (i.e. modelling p(x))
 - VAEs
 - Flows
 - ...

Lots of progress has been made here! ... and now we have much more data & compute

A 10M Dollar Question

With \$10M GPU computing budget...

Would you try pre-training any big models with any "old-fashioned" unsupervised learning techniques?

Before You Burn the \$10M GPU Budget...

- Is masked pre-training fundamentally different?
- Anything BAD about the "old fashioned" methods?
- How to answer these questions?
 - Modelling choices?
 - Advantages from Training objectives?
 - Optimisation aspects?

Complete Conditionals:

Full joint distribution:

 $\prod_{d=1}^{d_x} p(x_d | x_{-d})$

 $p(x_{1:d_x})$

- Modelling complete conditionals vs full joint distribution:
 - For some models, conditional distributions are much easier to compute
 - Markov random field
 - Restricted Boltzmann machine
 - Related to Gibbs sampling
 - Capturing complete conditionals ⇒ use Gibbs sampling to sample from the joint

- Complete Conditionals for representation learning (linear case):
 - Assume a linear relationship for the generative process:
 - Continuous case: $p(x|z) = N(Wz, \sigma^2 I)$
 - Discrete case: $p(x = i | z = j) = W_{ij}$
- Capturing the posterior p(z|x) via the complete conditional:
 - The long vector $[p(x_1|x_{-1}), p(x_2|x_{-2}), ..., p(x_D|x_{-D})]$ can be represented as a linear transform of p(z|x)
 - Under assumptions on W (e.g., full column rank)
 - Learned useful features if the ground truth x → y function depends directly on p(z|x)
- Similar ideas can be extended to Hidden Markov Models

Ζ

Complete Conditionals:

Full joint distribution:

 $\prod_{d=1}^{d_x} p(x_d | x_{-d})$

 $p(x_{1:d_x})$

- Comparing the two modelling paradigms:
 - Capturing information about the posterior p(z|x):
 - Both methods can do (under linear model assumptions)
 - Identifiability of the model
 - Even when the parameter of $p(x_{1:d_x})$ is identifiable, we might not be able to recover that parameter by modelling some conditional distributions

Liu et al. Masked prediction tasks: a parameter identifiability view. arXiv:2202.09305

Gassiat et al. Inference in finite state space non parametric hidden Markov models and applications[J]. Statistics and Computing, 2016, 26(1): 61-71

- In practice: we directly parameterise $p(x_d|x_{-d})$
 - Importantly, $p(x_i|x_{-i})$ and $p(x_j|x_{-j})$ share the same network
 - Different from first defining $p(x_{1:d_x})$ then working out $p(x_d|x_{-d})$
- The conditional distributions might not be compatible
 - Circumvented by universal approximations of neural networks?
 - Inconsistency issues: do they matter?

- Ground-truth data generation process:
 - $z \sim N(0, \Lambda), z \in \mathbb{R}^{d_z}$
 - $x = Wz + \epsilon_x$, $x \in R^{d_x}$
 - $y \sim p(y|z), y \in \mathbb{R}^{d_y}$
- Regression with linear features: define f(x) = Bx, solve $\theta^*(B) = \operatorname{argmin} E_{x,y}[\|y - \theta f(x)\|_2^2]$
 - Specifically we consider $f \in \mathbb{R}^{d_f}$ with $d_f < d_x$ (bottleneck)
 - And $d_f < d_z$, i.e., model is mis-specified



- Regression with linear features: define f(x) = Bx, solve $\theta^*(B) = argmin E_{x,y}[||y - \theta f(x)||_2^2]$
 - Bottleneck mis-specified: $f \in R^{d_f}$ with $d_f < d_x$, $d_f < d_z$
- Find *B* with linear auto-encoder: solve $A^*, B^* = argmin E_x[||x - ABx||_2^2]$

Solution: equivalent to PCA up to an invertible matrix $B^* = C^{-1}U_{d_f}^{\mathsf{T}}, \quad U_{d_f}$ captures the top d_f eigenvectors of $\Sigma_{\chi\chi}$ $(\Sigma_{\chi\chi} = \sigma_{\chi}^2 I + W \Lambda W^{\mathsf{T}})$

Representation learning works well if y only depends on the subspace of z found by PCA

$$z \sim N(0, \Lambda), z \in \mathbb{R}^{d_z}$$
$$x = Wz + \epsilon_x, x \in \mathbb{R}^{d_x}$$
$$y \sim p(y|z), y \in \mathbb{R}^{d_y}$$

Z

X

 $B \in R^{d_f \times d_x}$

- Regression with linear features: define f(x) = Bx, solve $\theta^*(B) = \operatorname{argmin} E_{x,y}[\|y - \theta f(x)\|_2^2]$
 - Bottleneck mis-specified: $f \in R^{d_f}$ with $d_f < d_x$, $d_f < d_z$
- Find *B* with masked linear auto-encoder:
 - Define $x_{-j}(e) = [x_1, ..., x_{j-1}, e, x_{j+1}, ..., x_{d_x}]^\top$
 - Solve the following objective: notice $A \in R^{1 \times d_f}$ $A^*, B^* = argmin \sum_j E_x[(x_j - ABx_{-j}(e))^2]$

 $B \in R^{d_f \times d_x}$

Solution is not unique (when
$$d_f > 1$$
), but of the following form:

$$A^* = \begin{bmatrix} 1, 0_{1 \times d_f - 1} \end{bmatrix} C, B^* = C^{-1} \begin{bmatrix} b^* \\ L \end{bmatrix}, b^* = \underbrace{\left(\sum_{xx} 1 - diag(\sum_{xx}) \right)^{\mathsf{T}} \left[(d_x - 1) \sum_{xx} + e^2 I \right]^{-1}}_{\sum_j E[x_j x_{-j}(e)^{\mathsf{T}}]} \underbrace{\sum_j E[x_{-j}(e) x_j(e)^{\mathsf{T}}]}_{\sum_j E[x_{-j}(e) x_j(e)^{\mathsf{T}}]}$$

 $z \sim N(0, \Lambda), z \in \mathbb{R}^{d_z}$ $x = Wz + \epsilon_x, x \in \mathbb{R}^{d_x}$ $y \sim p(y|z), y \in \mathbb{R}^{d_y}$ $\Sigma_{xx} = \sigma_x^2 I + W\Lambda W^{\mathsf{T}}$

- Regression with linear features: define f(x) = Bx, solve $\theta^*(B) = \operatorname{argmin} E_{x,y}[\|y - \theta f(x)\|_2^2]$
 - Bottleneck mis-specified: $f \in R^{d_f}$ with $d_f < d_x$, $d_f < d_z$
- Find *B* with masked linear auto-encoder: Arbitrary invertible $d_f \times d_f$ matrix $B^* = C^{-1} \begin{bmatrix} b^* \\ L \end{bmatrix}, \quad b^* = (\Sigma_{xx}1 - diag(\Sigma_{xx}))^T [(d_x - 1)\Sigma_{xx} + e^2I]^{-1}$ Arbitrary $(d_f - 1) \times d_x$ matrix, $d_f > 1$
 - b^* contains less information about x as compared with the PCA solution
 - However, flexibility in L and e allows B to project to different subspaces of z

$$z \sim N(0, \Lambda), z \in R^{d_z}$$
$$x = Wz + \epsilon_x, x \in R^{d_x}$$
$$y \sim p(y|z), y \in R^{d_y}$$
$$\Sigma_{xx} = \sigma_x^2 I + W\Lambda W^{\mathsf{T}}$$

Z

- Regression with linear features: define f(x) = Bx
 - Bottleneck mis-specified: $f \in R^{d_f}$ with $d_f < d_x$, $d_f < d_z$
- Find B, θ together with supervised learning:

$$\theta^*, B^* = \operatorname{argmin} E_{x,y}[\|y - \theta Bx\|_2^2]$$

• Consider the case of $d_y < d_f < d_z$:

Solution is not unique, but of the following form:

$$\theta^* = \left[I_{d_y \times d_y}, 0_{d_y \times (d_f - d_y)} \right] C, B^* = C^{-1} \begin{bmatrix} \Sigma_{yx} \Sigma_{xx}^{-1} \\ L \end{bmatrix},$$

 $(\theta^* B^* = \Sigma_{yx} \Sigma_{xx}^{-1})$ is the least squares solution of linear regression)

 $z \sim N(0, \Lambda), z \in \mathbb{R}^{d_z}$ $x = Wz + \epsilon_x, x \in \mathbb{R}^{d_x}$ $y \sim p(y|z), y \in \mathbb{R}^{d_y}$ $\Sigma_{xx} = \sigma_x^2 I + W\Lambda W^{\mathsf{T}}$

Ζ

 $B \in R^{d_f \times d_x}$

- Regression with linear features: define f(x) = Bx
 - Bottleneck mis-specified: $f \in \mathbb{R}^{d_f}$ with $d_f < d_x$, $d_f < d_z$
- Now compare the two solutions: assume $d_y < d_f < d_z$

Supervised learning:

$$B^* = C_s^{-1} \begin{bmatrix} \Sigma_{yx} \Sigma_{xx}^{-1} \\ L_s \end{bmatrix}$$

Masked linear auto-encoder:

$$B^* = C_m^{-1} \begin{bmatrix} b^* \\ L_m \end{bmatrix}, \ b^* = \left(\Sigma_{xx} 1 - diag(\Sigma_{xx}) \right)^{\mathsf{T}} \left[(d_x - 1) \Sigma_{xx} + e^2 I \right]^{-1}$$

 $z \sim N(0, \Lambda), z \in \mathbb{R}^{d_z}$ $x = Wz + \epsilon_x, x \in \mathbb{R}^{d_x}$ $y \sim p(y|z), y \in \mathbb{R}^{d_y}$ $\Sigma_{xx} = \sigma_x^2 I + W\Lambda W^{\mathsf{T}}$

- Regression with linear features: define f(x) = Bx
 - Bottleneck mis-specified: $f \in \mathbb{R}^{d_f}$ with $d_f < d_x$, $d_f < d_z$
- Now compare the two solutions: assume $d_y < d_f < d_z$



 $z \sim N(0, \Lambda), z \in \mathbb{R}^{d_z}$ $x = Wz + \epsilon_x, x \in \mathbb{R}^{d_x}$ $y \sim p(y|z), y \in \mathbb{R}^{d_y}$ $\Sigma_{xx} = \sigma_x^2 I + W\Lambda W^{\top}$

Z

X

- Regression with linear features: define f(x) = Bx
 - Bottleneck mis-specified: $f \in R^{d_f}$ with $d_f < d_x$, $d_f < d_z$
- Now compare the two solutions: assume $d_y < d_f < d_z$

Making the two solutions equal: find matrices C and L such that

$$\frac{1}{d_{x}-1} \left(\Sigma_{xx} 1 - diag(\Sigma_{xx}) \right)^{\mathsf{T}} = C_{1} \begin{bmatrix} \Sigma_{yx} \\ L \end{bmatrix},$$

$$\in R^{1 \times d_{x}}$$

$$\in R^{1 \times d_{f}}$$

$$\in R^{(d_{f}-d_{y}) \times d_{x}}$$

(the 1st row of C)

"solving a set of d_x equations with $d_f + (d_f - d_y) \times d_x$ variables"

 $d_y < d_f \Rightarrow d_f - d_y \ge 1 \Rightarrow d_f + (d_f - d_y) \times d_x > d_x \Rightarrow$ solution exists!

 $z \sim N(0, \Lambda), z \in R^{d_z}$ $x = Wz + \epsilon_x, x \in R^{d_x}$ $y \sim p(y|z), y \in R^{d_y}$ $\Sigma_{xx} = \sigma_x^2 I + W\Lambda W^{\top}$

Z

X

- Regression with linear features: define f(x) = Bx
 - Bottleneck mis-specified: $f \in R^{d_f}$ with $d_f < d_x$, $d_f < d_z$
- Compare w/out masking :
 - Linear auto-encoder:
 - Solution equivalent to PCA
 - Best linear representation for x (in terms of ℓ_2 recon. error)
 - Might not be good for prediction (depending on ground truth p(y|z))
 - Masked linear auto-encoder:
 - Solution is not unique (when $d_f > 1$)
 - Not so good representation for x (in terms of ℓ_2 recon. error)
 - A subset of solutions are optimal for prediction (when $d_y < d_f < d_z$)

Hypothesis: masked pre-training works better when good representation for generating x is highly redundant

 $z \sim N(0, \Lambda), z \in \mathbb{R}^{d_z}$ $x = Wz + \epsilon_x, x \in \mathbb{R}^{d_x}$ $y \sim p(y|z), y \in \mathbb{R}^{d_y}$ $\Sigma_{xx} = \sigma_x^2 I + W\Lambda W^{\top}$

- Regression with linear features: define f(x) = Bx
 - Bottleneck mis-specified: $f \in R^{d_f}$ with $d_f < d_x$, $d_f < d_z$
- Caveats of this example:
 - Masked linear auto-encoder:
 - A subset of solutions are optimal for prediction (when $d_y < d_f < d_z$)
 - But: pre-training on x has no obvious "incentives" to find them (optimal B^{*} in supervised learning depends on Σ_{vx})
 - Extension to non-linear models:
 - With over-parameterised NNs, auto-encoder solution might not be unique
 - But Masked auto-encoder might have bigger subspace of optimal solution for prediction
 - Again: pre-training on x has no obvious "incentives" to find them

 $z \sim N(0, \Lambda), z \in \mathbb{R}^{d_z}$ $x = Wz + \epsilon_x, x \in \mathbb{R}^{d_x}$ $y \sim p(y|z), y \in \mathbb{R}^{d_y}$ $\Sigma_{xx} = \sigma_x^2 I + W \Lambda W^{\mathsf{T}}$

Z

X

Optimisation Aspects

- Lower dimensional target to match:
 - Linear Auto encoder (PCA) features: f(x) = Bx

$$A^*, B^* = \operatorname{argmin} E_x[\|x - ABx\|_2^2] \qquad A \in \mathbb{R}^{d_x \times d_f}$$

• Masked Linear Auto encoder features: f(x) = Bx

$$A^*, B^* = argmin \sum_j E_x[(x_j - ABx_{-j}(e))^2] \qquad A \in \mathbb{R}^{1 \times d_f}$$

Optimisation Aspects

- More equally good optima (for representation learning)?
 - Linear Auto encoder (PCA) features: f(x) = Bx

Arbitrary invertible matrix

$$B^* = C^{-1}U_{d_f}^{\top}, \quad U = eig(\Sigma_{xx})$$

• Masked Linear Auto encoder features: f(x) = Bx

Arbitrary invertible matrix

$$B^* = C^{-1} \begin{bmatrix} b^* \\ L \end{bmatrix}, \quad b^* = (\Sigma_{xx} 1 - diag(\Sigma_{xx}))^{\mathsf{T}} [(d_x - 1)\Sigma_{xx} + e^2 I]^{-1}$$
Arbitrary $(d_f - 1) \times d_x$ matrix
Mask embedding

Optimisation Aspects

• Current DL tricks engineered towards supervised learning tasks?

BEiT as an example:

- Pretraining task: classification
- (sparse) Cross-entropy loss



Bao et al. BEiT: BERT Pre-Training of Image Transformers. ICLR 2022 Wang et al. Image as a Foreign Language: BEiT Pretraining for All Vision and Vision-Language Tasks. arXiv:2208.10442

Still Consider Burning the \$10M GPU Budget?

- Is masked pre-training fundamentally different?
- How to answer this question?
 - Modelling choices?
 - Modelling Conditional distributions: a better choice?
 - Not the full story: incompatibility/inconsistency issue matters or not?
 - Advantages from Training objectives?
 - More flexible representation learning?
 - Redundancy in information of *x*?
 - Any benefit of the extra [mask] embedding *e*?
 - Optimisation aspects?
 - Easier for optimisation? due to lower-dim targets (many equally good optima)?
 - Deep neural networks designed to benefit supervised learning objectives?
 - Other perspectives (from you)?