

28/3/22

PAMI ComprehensiveAnubhav Jain  
D20CS002

(1)

PART-A

Ques 1:  $A = \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix}$   $B = \begin{bmatrix} 1 & -2 \\ 0 & 1 \end{bmatrix}$

Let's say  $B$  is the inverse of  $A$  matrix, then,

$$B = A^{-1}$$

Also,

$$A A^{-1} = I$$

$$\therefore \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix} \times \begin{bmatrix} 1 & -2 \\ 0 & 1 \end{bmatrix} = I$$

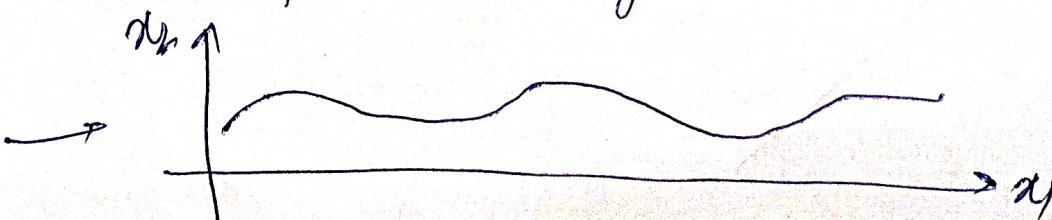
$$= \begin{bmatrix} 1+0 & -2+2 \\ 0+0 & 0+1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \text{Identity Matrix.}$$

$\therefore$  True, given Matrices are inverses of each other.

Ques 2: objective function on  $\mathbb{R}^2$  where gradient descent will be exceedingly slow.

→ a wavy function with not much of a change could be one function where gradient descent could be slow



(2)

gradient step update  $\rightarrow$ 

$$\theta := \theta - \alpha \frac{\delta}{\delta \theta} J(\theta)$$

 $\alpha \rightarrow$  learning Rate $\theta \rightarrow$  parameters to be updated $J(\theta) \rightarrow$  cost function

- \* A cost function that doesn't have a big rate of change is one ~~function~~ that will be slow to converge  
property of the functions
- \*  $\alpha \rightarrow$  learning rate can also be reduced.  
But a function like  $\rightarrow x^{1/2} + y^{1/2}$   $\rightarrow$  can be slow to converge as  $x$  &  $y$  are free variables in  $\mathbb{R}^2$  domain.
- \* Any plateau-based function are approaching plateau-like steps will be exceedingly slow to converge.

Ques 3: We can't remove bias parameter from a convolutional layer or a fully-connected layer before applying batch norm is a false statement. (3)

Batch normalization includes normalizing the data sent in the particular iteration pass. The normalization occurs and allows to reduce the vanishing and exploding gradient problem. This normalization is formulated with a bias term involved; thus, bias before batch normalization layer can be considered quite redundant. There is no need for an extra bias term, and thus we can't only remove ~~but~~ the term altogether, having it is quite redundant.

Ques 4: Autoencoder with ~~for~~ linear activation and single hidden layer is same as PCA. Is a true statement. PCA and ~~an~~ autoencoder are differentiated by the non-linearity component present in the autoencoder. Linear activation and only a single layer will have the same subspace and will reduce dimensionality just as in the PCA. The dimensionality will have to be the same though for this to work out. They will be equivalent in the subspace spanned.

(Ques 5) In the basic formulation of GAN, both generator and discriminator perform a binary logistic regression with the cross-entropy loss. Is a true statement (4)

$$\rightarrow \min_{E} \max_{D} G(E, D) \quad \text{we have a minmax GAN loss}$$

$\hookrightarrow G \rightarrow \text{GAN}$

$E \rightarrow \text{Generator}$

$D \rightarrow \text{Discriminator.}$

Both Generator and discriminator seems to work of binary logistic regression with Generator trying to minimize

$$1 - \log E(z) \quad z \rightarrow \begin{matrix} \text{input} \\ \text{image} \end{matrix}$$

where,

$$\rightarrow E(z) = R$$

$$\rightarrow D(R) = z'$$

$$\rightarrow D(E(z)) = z'$$

$R \rightarrow \text{Representant}^2 / \text{Reconstructed image by } E$

Similarly, discriminator  $D$  is trying to maximize. The final formulation looks as follows →

at generator side  $G(E, D) = \frac{E}{P_{\text{data}}} [\log D(z)] + \frac{E}{P_z} [\log (1 - D(E(z)))]$

for GAN we go in step-by-step training, that is, freeze Generator, update  $D$  and freeze discriminator to update Generator. This in its simple formulation can be done using binary logistic regression with cross-entropy loss.

Ques. False, Interpretability and Explainability of a neural network are not the same. Interpretability of a model deals with how we interpret the decisions of the model while explainability is more about how we explain a specific decision made by the model.

for example, say we have a simple Neural Network with activation function predicting the price of houses. Interpretability of the model will include questions like how do we interpret the particular model used. It will also include questions like how do we justify or interpret a choice of activation function used, like using sigmoid squashes the values between 0 to 1 usually leading to normalized values.

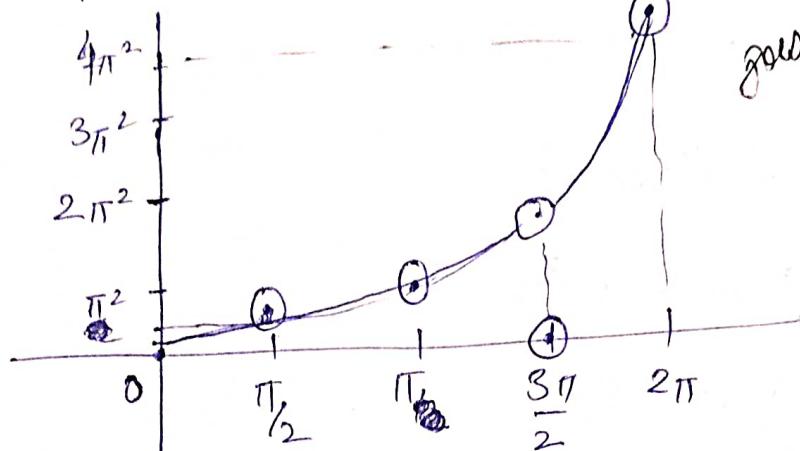
Whereas, explainability will deal with questions like given a particular set of features  $x$ , why do the model gives  $y$  as the output. what makes the model take this decision.

### Ques 7: AdaGrad

AdaGrad is an adaptive learning rate algorithm which allows to aggressively decay the learning rate as training progresses. It is used to modify the Gradient descent algorithm.

$$\text{Ques 8. (a)} f(x) = x^2(1 + \sin(\pi x))$$

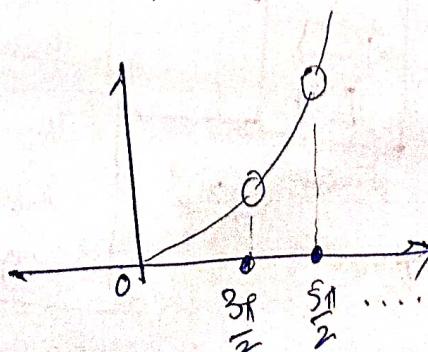
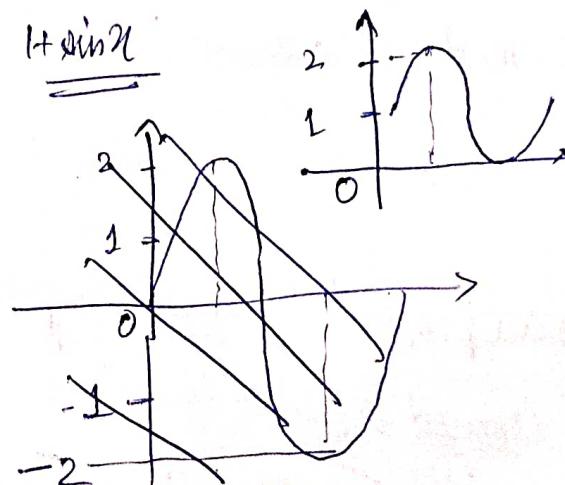
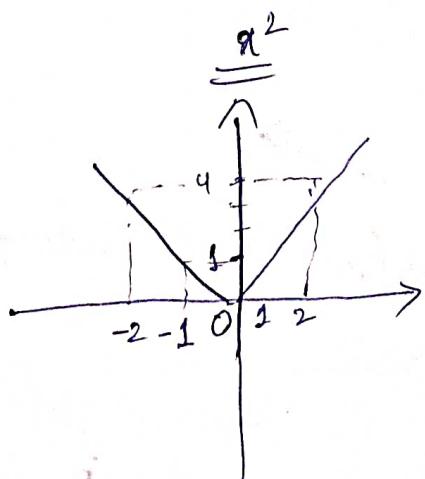
|        |   |                     |                               |                  |          |
|--------|---|---------------------|-------------------------------|------------------|----------|
| $x$    | 0 | $\frac{\pi}{2}$     | $\pi$                         | $\frac{3\pi}{2}$ | $2\pi$   |
| $f(x)$ | 0 | $\frac{2x\pi^2}{4}$ | $\frac{3\pi}{2}(1-1)(2\pi)^2$ | 0                | $4\pi^2$ |
|        |   | $\frac{\pi^2}{2}$   | $\pi^2$                       | 0                |          |



goes to 0 for  
every  $\sin x = -1$ .  
i.e.  $\forall x = (2n+1)\frac{\pi}{2}$ .

(b) Computation Graph.

$$f(x) = x^2(1 + \sin x)$$



(7)

(C)  $f'(0)$ 

$$f(x) = x^2(1 + \sin(x))$$

$$= x^2 \frac{d}{dx} [1 + \sin(x)] + \frac{d}{dx} x^2 \cdot (1 + \sin(x))$$

$$= x^2(0 + \cos x) + 2x \cdot (1 + \sin x)$$

$$= x^2 \cos x + 2x(1 + \sin x)$$

$$= x^2 \cos x + 2x \sin x + 2x$$

$$= x[x \cos x + 2 \sin x + 2]$$

$$f'(0) = 0[0 + 0 + 2] = \underline{\underline{0}}.$$

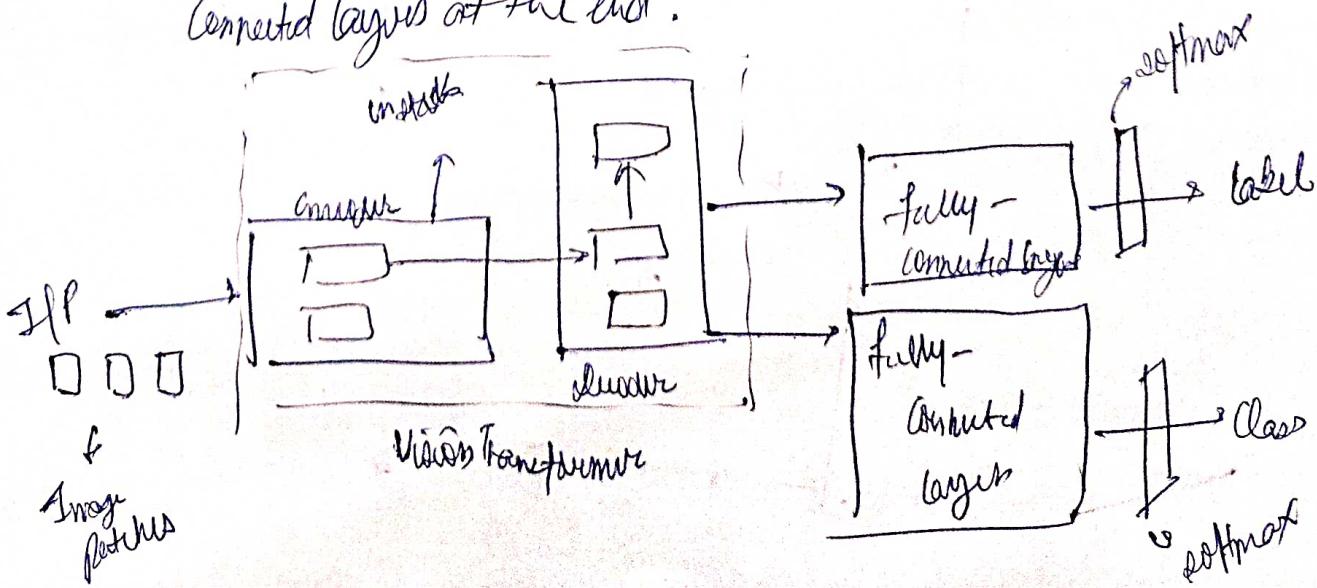
### Ques 9. Imbalanced Domain learning.

\* "Multi-label - Multiclass problem" (ML-MC)

\* n-classes

m-labels

\* Architecture proposed : Visiontransformer (using attention) with fully-connected layers at the end.



(8)

The input is given in form of patches. Vision transformer can be replaced by simple transformer in case input is textual. Further, attention is used and positional embeddings are required. The precessed input is then fed to 2 fully-connected layers, one for n-classes and one for m-labels. Both the layers have softmax layer to choose from the labels and classes based on a probability distribution -

- \* for imbalanced domain learning we use regularization in the training process. Here, for training, we give more weights to some class or domain <sup>not</sup> present in abundance. This will allow the model to focus on scarce examples and give them more attention.
- \* moreover, using transformers we won't be invariant to textual features which can help with imbalanced domain learning.
- \* finally, CNNs can be used but they are not localized, that is, we won't be able to use positional embeddings, neither we will be able to use the attention module.



## Dues 10 - Distributed Deep Learning

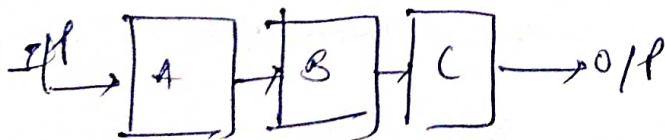
(9)

### 1(a) Model & Data parallelism

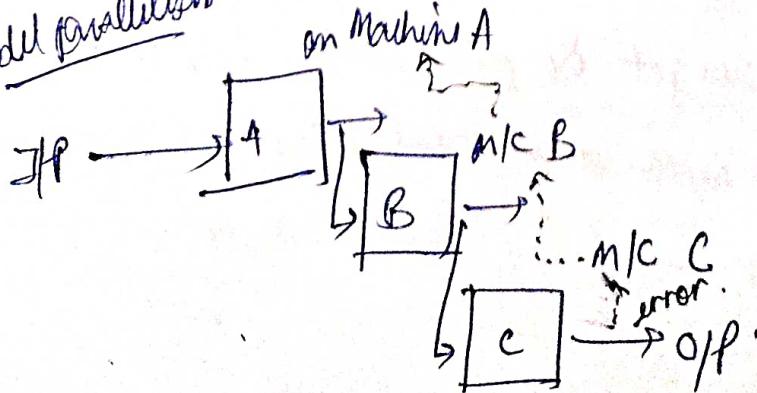
#### Model parallelism

- here delta remains the same.
- here the model is trained in parallel. One can imagine it training as greedy layer by layer training. However, instead of freezing the layers, we train the layers in parallel over the same dataset.
- gradients flow back once we reach to the last layer of the model.

Say model has been divided into 3 parts

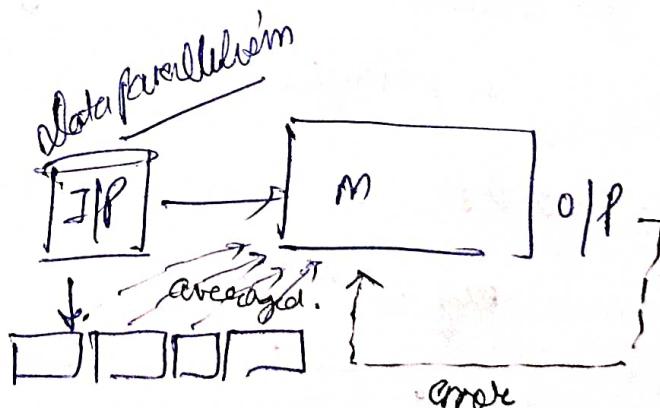


#### Model parallelism



#### Data parallelism

- the same model instance is used. SGD uses certain samples from the entire dataset. Thus we know, the complete data is not required in one go.
- gradients for data parallelism can be backpropagated by averaging them over the entire dataset.
- Data can be fed in parallel in mini batches, after computing gradients we average & backpropagate them.

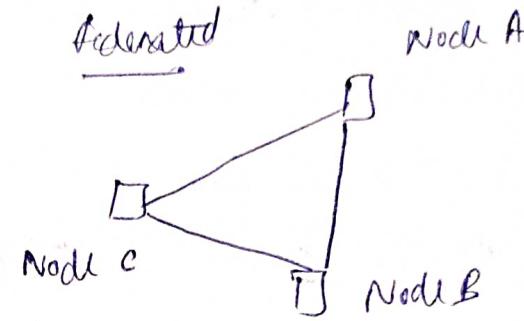
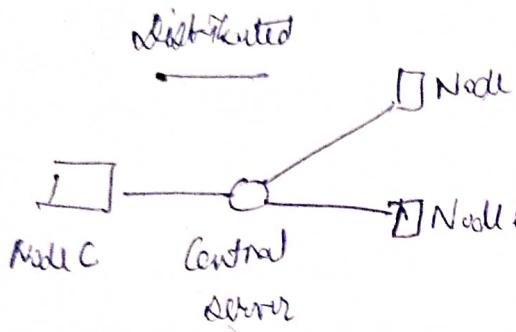


Weights are updated once the gradients (long) are backpropagated after averaging gradients

## (b) Distributed DL and Federated DL

(10)

In distributed deep learning, we share the model and ~~the~~ data among all the resources while the updates are made at a central server. On the other hand, federated learning is about having no central server - decentralized learning.



In federated learning, it is upto us as to what to share with the nodes, that is all the nodes only gets the updated weight without having to share their data. They only share the model weights after training the model locally. This way we can ensure data security & privacy.

For Federated DL,

say  $\rightarrow m$  nodes

$\hookrightarrow$  training happens at  $m$  nodes with data at the particular node  $X_i \{i \in m\}$

say  $\rightarrow$  after updation we get  $\theta_i$  parameters  $\{i \in m\}$

these parameters are shared with other nodes, who combine or update their own model by fusing these updated parameters.

Iteration

(11)

Parameters  ~~$\theta_1, \theta_2, \dots, \theta_m$~~

Data  $x_1, x_2, \dots, x_m$

After training

(training at model 1)  
assume,

$\theta_1^*, \theta_2, \dots, \theta_m$

new parameters  $\theta_1^*$   $\underbrace{\theta_1^*, \theta_2}_{\theta_2^*} \dots \underbrace{\theta_1^*, \theta_m}_{\theta_m^*}$

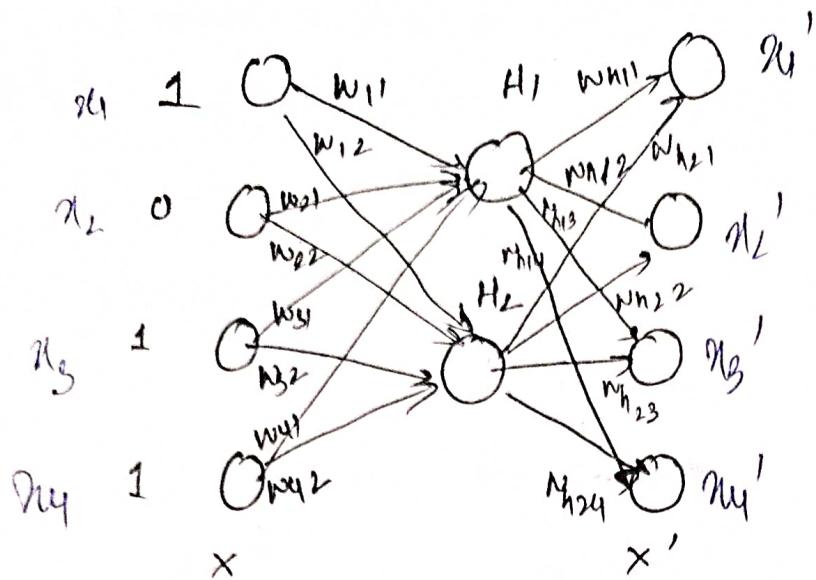
$x$



dwell.

$$g(x) = \frac{x}{4}, \eta = 0.5$$

f2



$$x = 1 0 1 1$$

weights = initialized to 1.

$$g(Wx)$$

$$h_1 = w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + w_{14}x_4$$

$$= 1 \times 1 + 1 \times 0 + 1 \times 1 + 1 \times 1$$

$$= 3 \quad h_1' w_{h11} + h_2' w_{h21} = h_1' = \frac{3}{4} + \frac{3}{4} = \frac{3}{2}$$

$$g(h_1) = \frac{3}{4} = h_1'$$

~~$$\text{similarly } h_2 = 3$$~~

$$g(h_2) = \frac{3}{4} = h_2'$$

$$\left. \begin{aligned} h_1' w_{h12} + h_2' w_{h22} &= h_2' = \frac{3}{2} \\ h_1' w_{h13} + h_2' w_{h23} &= h_3' = \frac{3}{2} \\ h_1' w_{h14} + h_2' w_{h24} &= h_4' = \frac{3}{2} \end{aligned} \right\}$$

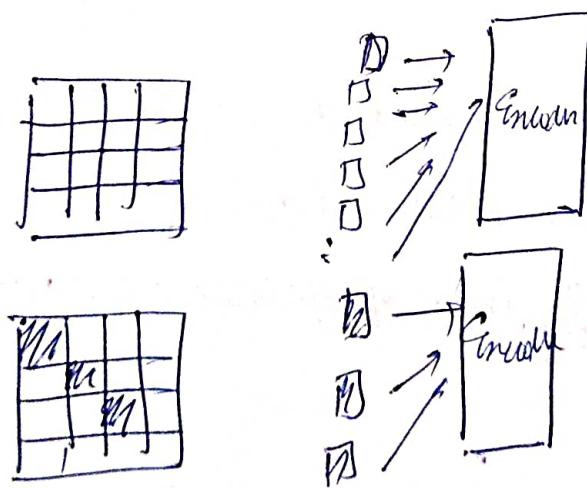
loss calculation

$$\left[ \left| \frac{3}{2} - 1 \right|^2 + \left| \frac{3}{2} - 0 \right|^2 + \left| \frac{3}{2} - 1 \right|^2 + \left| \frac{3}{2} - 1 \right|^2 \right]^2$$

$$\text{error} = \frac{3\sqrt{5}}{2} + \frac{3}{2} = \frac{3}{2}(1 + \sqrt{5})$$

Ans (a) The specialized sparse operations can include dimensionality reduction and dense feature representation. Images introduce sparcity, which means we can reduce dimensions of the inputs without actually affecting or losing any important information that can lead to model's final decision. Such operations include using algorithms like PCA or providing entire input to the encoder. Here, in the paper, only non-masked regions were provided to the encoder. Certain models, even have specialized sparse matrix representation to handle such sparse inputs. dictionary learning also takes advantage of this particular sparse property.

(b) Why are such operations not required by this model →?



There are no requirements for any specific sparse representations or operations in this paper because upto 75% of the input image is masked leaving only few blocks to be sent as an input. these unmasked parts act as dense information, reducing the computation load for encoder drastically.

## Ques 2 - Masking to other CNN architectures

(14)

- \* It is particularly mentioned in the paper that convolutions typically operate on regular grids and it is not straightforward to integrate "indicators" such as mask tokens or positional embeddings into convolutional networks.
- \* That being said, it is not completely impossible. It will however be computationally expensive to integrate the positional embeddings required in the encoder and the mask tokens required in the decoder of the proposed masked Autoencoder. We will need extra embedding layers to integrate the information and convolutions will have to be hand-crafted for the information to integrate as required.
- \* Other models/architectures like → variant of vision transformers like (RVT) or OpenAI (which treat IP as vectors) will be easier to work with with the given specifications for the input to the encoder & decoder.

Conclusion → We can apply masking to the different architectures even CNN. However, it can be computationally expensive to use the proposed technique in the CNN-based model.

## Ques 3 - Masking Techniques

(15)

non-overlapping patches.

- \* In the paper → ① generate a token for every input patch using linear projection with an added positional embedding.  
② Randomly shuffle the list of tokens and remove the last portion of the list based on the masking ratio (high to more complex task)
- \* Aim of the technique → The high masking ratio is used here for proper training but it also allows to reduce the computation by decreasing the input size of the model.
  - Sampling random patches without replacement eliminates redundancy.
  - The uniform distribution ~~can~~ prevent a potential center bias.
  - Also, the self-attention has "quadratic" complexity, using high masking ratio allows to avoid ~~to~~ increasing the self-attention's complexity.
- \* Other strategies used → block-wise & grid-wise sampling. In both the cases, conclusion was that random sampling works better as it allows for a higher-masking ratio, which is good for speedup, computation cost & accuracy.

## Other masking techniques

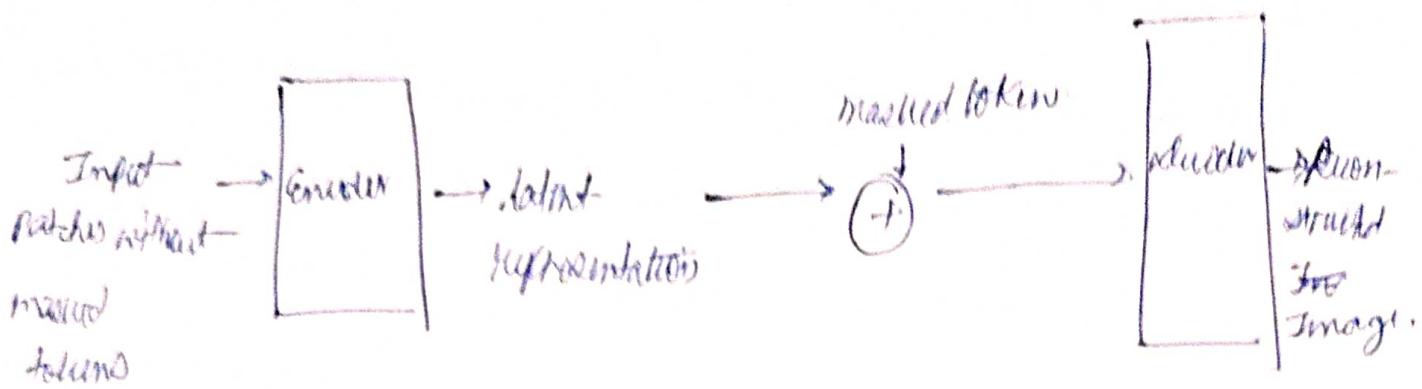
- pixel-based masking : using pixels to mask rather than tokens
- grid-based masking : dividing images into grids and using specific grids to be masked rather than tokens.
- hand-crafted masking : using eigen vectors or SVD decomposition,

These decomposition techniques can be used to mask specific tokens or parts of the image rather than random sampling techniques. (16)

|   | Pixel-based  | Grid-based  | Hand-crafted   | Proposed RandomSampling  |
|---|--|---|--|--|
| • Semantic Segmentation<br>Object Recognition | It can be combined with grid-based masking to identify pixels belonging to specified classes.  | grid-based masking can not capture change in feature as pixel-based masking is based on whether mutation is happening on pixels, tokens, or grid.               | hand-crafted masking can   | MAE performs better than token-based.  |
| Object Segmentation & Feature Recognition     | Pixel, grid-based masking might work better as pixel-based masking can endlessly increase computation power, it might increase performance on small objects. | Grid-based masking might help capture relevant features as object recognition looks for bounding box regression without much improvement over other techniques. | hand-crafted masking might just introduce more computation without much improvement over other techniques. | MAE proposed works well as object recognition, raising the mAP metric value. |
| X-ray Classification <sup>2</sup>             | Might just increase computation as the entire image is used as classifier.   | Might not know the affectiveness in the medical domain as pixels or tokens are specific to the domain.  | Hand-crafted masking can help especially knowing the specific domain.                                      | Proposed MAE → equally pixel-based could work well                           |

## Self-supervised pre-training

(17)



Encoder:

I/P → ~~MAE~~ gets tokens as input to the encoder.

Tokens are prepared by dividing image into patches, converting them to tokens by integrating positional embeddings and then masking tokens by random sampling. No masked tokens are given as input.

O/P → a latent representation is returned as the output of the encoder.

Decoder:

I/P → The latent representation (O/P of encoder) is joined with the masked tokens in order (as per the original patch list) and sent as input to decoder.

O/P → A reconstructed image is returned as the output of the decoder.

Pre-training : self-supervised task: get the masked tokens and reconstruct the original image

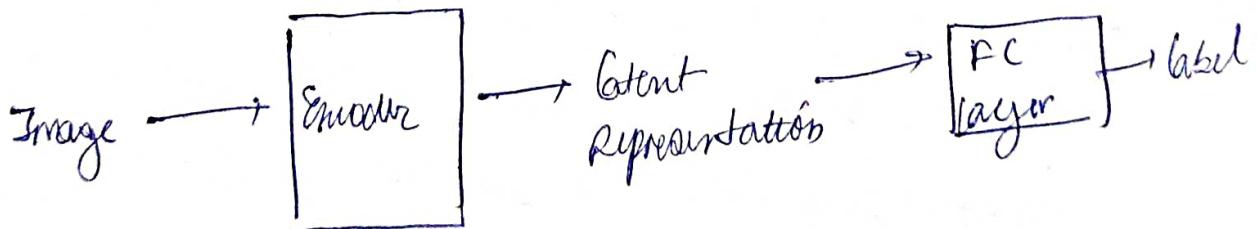
(18)

- The pre-training occurs in an asymmetric fashion with only non-masked tokens as encoder's Input and entire patch list as decoders Input [Latent variables + Masked tokens]
- A high masking ratio is used to make the task complex enough for the MAE to learn significant features which is then used for other tasks like Object Recognition, Semantic Segmentation and so on.

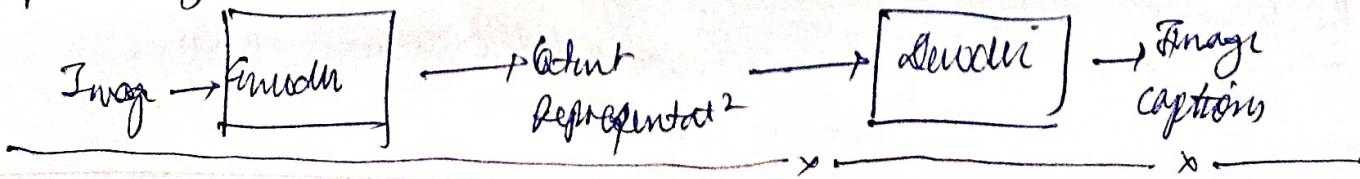
### Alternate pre-training Mechanism

↓  
without self-supervision

we can pre-train the model by giving image as input and getting its label as output. A simple classification task which is supervised (as it requires annotated labels) can be used as pre-training mechanisms. However, a huge amount of examples will be required to do meaningful pre-training.



For MAE as the architecture, we can use another supervised task like generating image captions using the decoder given an input image. We will need image captions for all the images.



## PART-B2 (Free lunch for few-shot learning)

(19)

Ans 2. equation 6,

$$\mu' = \frac{\sum_{i \in S_N} \mu^i + \tilde{x}}{K+1}$$

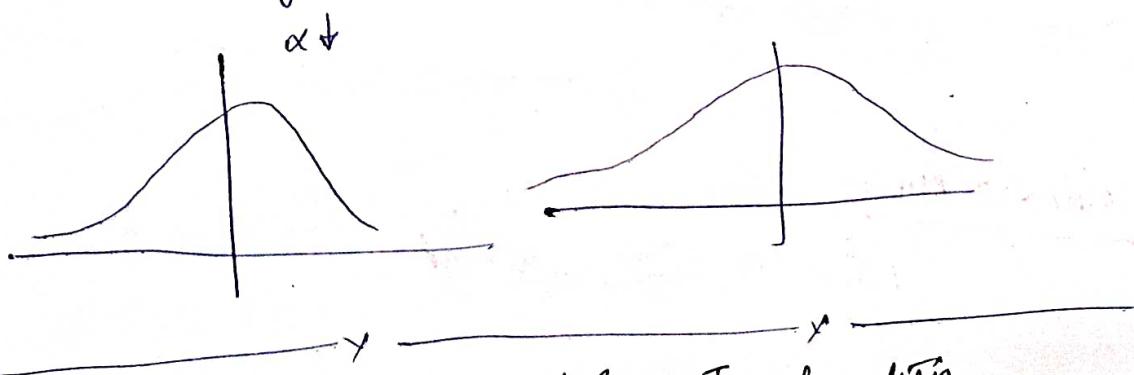
$$\Sigma' = \frac{\sum_{i \in S_N} \Sigma_i}{K} + \alpha_i$$

$\alpha \rightarrow$  is a hyper-parameter that determines the degree of dispersion of features sampled from the calibrated distribution.

This hyperparameter indicates how far or how dispense we are willing to go from the calibrated distribution. This could be thought of as ~~as~~ a regularization term, which allows to avoid any overfitting or as in this case, will help avoid any strict distribution. It allows dispersion from the computed variance and allow for a broader variance in the distribution.

based on the value of the ' $\alpha$ ' term.

$\alpha \uparrow$



Ans 2 Section 3.2.1, Tukey's ladder of Power Transformation

- This power transformation is not used for computing the statistics of the base dataset for the following reasons →

① Tukey's Ladder of Power Transformation is used to make the feature distribution more ~~less~~ Gaussian-like, however, for base classes we start with the assumption that the feature distribution is Gaussian.

② Support set and query set are transformed in the ~~paper~~ paper to reduce the skewness of distributions and make them more Gaussian-like - Doing this, the authors are bringing novel classes closer to the base classes, assuming that base classes are Gaussian and making novel classes more Gaussian-like, it can be visualized as taking novel classes closer to base classes, hence, making it easier to look for similar feature distributions.

③ Computationally, this makes sense, as assuming base classes to be Gaussian to begin with, the computation ~~expensive~~ expense decreases while transforming only the novel classes - whose examples are already scarce and thus, doesn't create any significant computation ~~overhead~~ overhead.

Ans 3: Calibrated Mean,

$$\text{numerator} = \sum_{i \in S_N} \mu^i + \tilde{\alpha}$$

from equation  $\rightarrow$

$$\mu' = \frac{\sum_{i \in S_N} \mu^i + \tilde{\alpha}}{k+1}$$

The correct form of numerator is  $\rightarrow$

(21)

$$(\sum_{i \in S_N} \mu_i) + \bar{x}$$

Because we are calculating the mean of a distribution. The mean added  $\bar{x}$  is computed off of a distribution (from the novel classes). We can imagine this as adding one more sample to the feature set as is depicted in the denominator by "k+1". The mean of the novel classes is added to the statistics of the base class to give calibrated distribution statistics.

The paper clearly states that  $\bar{x}$  is nothing but a feature vector of the novel class on which we are computing the calibrated distribution. We are using the novel base classes as `topk()` operator is being used.

Ques 4: Empirical Results for table 4.

Table 4 → provides results for miniImageNet  $\rightarrow$  5 way 1 shot and 5 way 5 shot problems when Tukey's transformation and training with generated features is used or not used.

In paper → Authors have stated that without either of the two, the accuracy takes a significant dip leading to the conclusion that both are a requirement for a better performance.

Other than that, we can make the following observations → (22)

- ⑤ There is no better option between the two. That is without taking transformation, 5-way 1 shot and 5-way 5 shot have quite similar accuracies to when taking transformation is performed but with no training from generated features.

| Takes | Generated Features | Sway 1 shot | Sway 5 shot |
|-------|--------------------|-------------|-------------|
| ✓     | ✗                  | 64.30       | 81.33       |
| ✗     | ✓                  | 63.70       | 82.26       |

- ⑥ Furthermore, we can say that with more examples the effect of these 2 tricks — taking transformation and training with generated features ~~too~~ — is reduced.

That is, for 1 shot, that is, 1 example in each novel class

there is an improvement of 12%

But for 5 shot, that is, 5 examples in each novel class, there is only an improvement of 4% (approx.) when both tricks are used as opposed to none of the ~~two~~ trick being used.

### Ques: Testing of transfer of statistics

Evaluation Metrics in the paper for testing → ① Top-1 accuracy on 2 settings,  
5-way 1shot & Sway 5shot

→ ② t-SNE visualization (visualisation  
of feature distribution?)

far the transfer of base to novel classes' statistics, we can use ②<sup>3</sup>  
the following factor's methods to be more thorough →

- ③ we use classes with several examples available as  
novel classes. we can withhold the available examples  
using only 1shot or 5shot settings. After transferring  
the statistics, we can average the computed statistics  
with that taken computed using other available examples.  
The closer the two are, the better transfer of statistics  
happened from base to novel classes.
- ④ we can introduce multi-domains to increase the complexity  
and sub-domains to allow for more accurate relationship  
or closer relationship between base and novel classes.  
using sub-domains we can truly say that base and novel classes  
are similar and thus calibrated distribution can give better  
accuracy or atleast it can utilized during testing to see whether  
similar domains are being used in the topk() operator or not  
and what classes are forming closer relationships.