

Open Elective Course [OE]

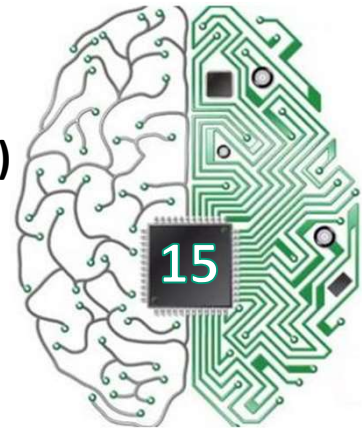
Course Code: CSO507

Winter 2023-24

Lecture#

Deep Learning

Unit-4: Convolutional Neural Networks (Part-III)

Course Instructor:

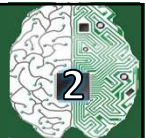
Dr. Monidipa Das

Assistant Professor

Department of Computer Science and Engineering

Indian Institute of Technology (Indian School of Mines) Dhanbad, Jharkhand 826004, India

Convolution Example (revisited)

Input volume: $3 \times 32 \times 32$ 10 5×5 filters with stride 1, pad 2

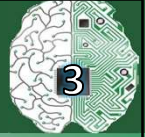
Output volume size: ? $(32 + 2 \times 2 - 5) / 1 + 1 = 32$ spatially, so
 $10 \times 32 \times 32$

Number of learnable parameters: ?

Parameters per filter: $3 \times 5 \times 5 + 1$ (for bias) = 76
 10 filters, so total is $10 \times 76 = 760$

Number of multiply-add operations: ?

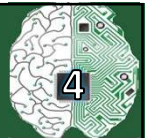
$10 \times 32 \times 32 = 10,240$ outputs; each output is the inner product
 of two $3 \times 5 \times 5$ tensors (75 elems); total = $75 \times 10240 = 768K$



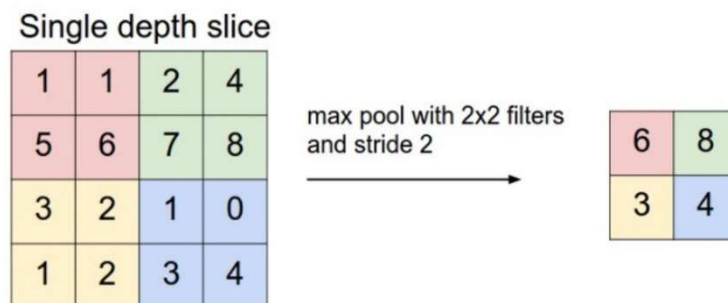
Pooling Layer

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Pooling Layer: Summarizes Neighborhood



- **Max-pooling:** Applying the filter and finding the maximum value for each chunk

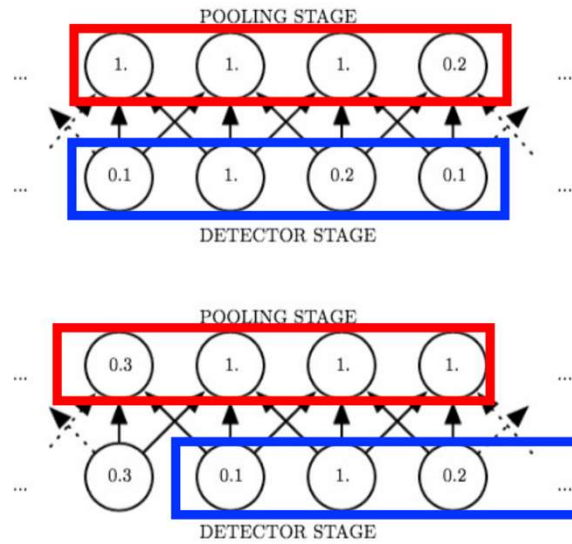
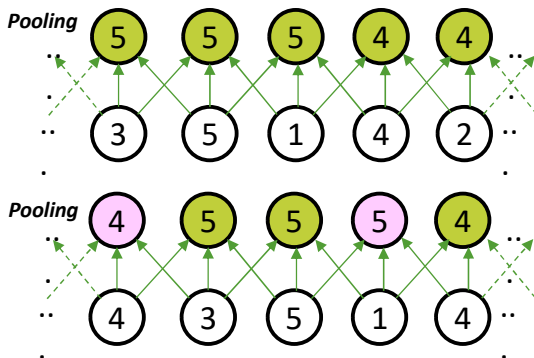


Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Pooling Layer



- Resilient to small translations
 - e.g.,
 - Input: all values change (shift right)
 - Output: only half the values change



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Pooling Layer: Summarizes Neighborhood



- Max-pooling:** Applying the filter and finding the maximum value for each chunk
- Average-pooling:** Applying the filter and finding the average value for each chunk

Single depth slice

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

Avg pool with 2x2 filters
and stride 2



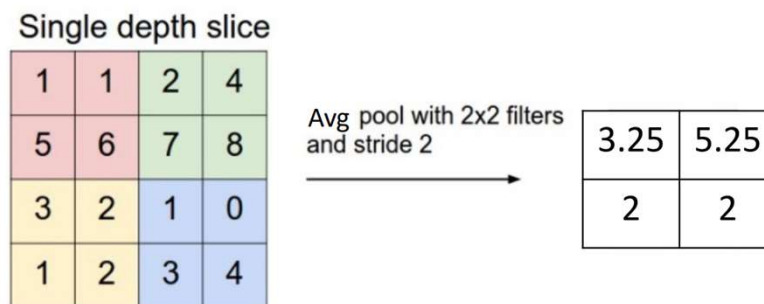
?	?
?	?

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Pooling Layer: Summarizes Neighborhood

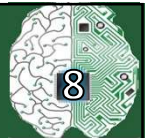


- **Max-pooling:** Applying the filter and finding the maximum value for each chunk
- **Average-pooling:** Applying the filter and finding the average value for each chunk

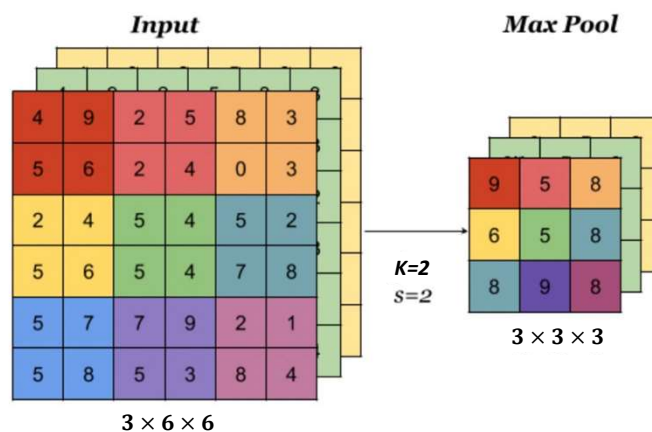


Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Pooling for Multi-Channel Input



- Pooling is applied to each input channel separately

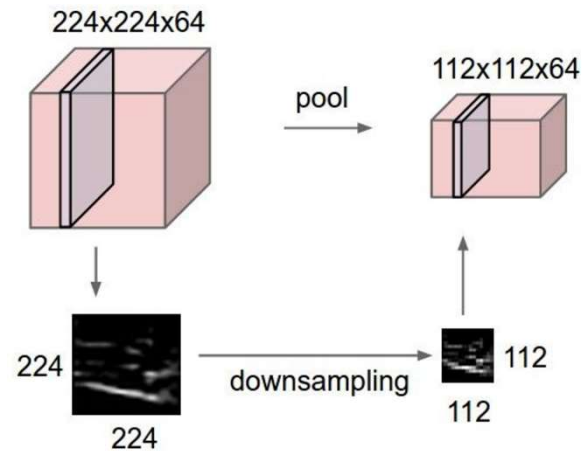


Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Pooling Layer: Benefits

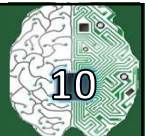


- How many parameters must be learned?
 - None
- Benefits?
 - Builds in invariance to translations of the input
 - makes the representations smaller and more manageable
 - operates over each activation map independently
 - Reduces memory requirements
 - Reduces computational requirements



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Pooling layer: summary



Input: $C \times H \times W$

Hyperparameters:

- Kernel size: K
- Stride: S
- Pooling function (max, avg)

Output: $C \times H' \times W'$ where

- $H' = (H - K) / S + 1$
- $W' = (W - K) / S + 1$

Learnable parameters: None!

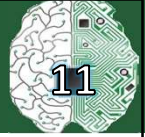
Common settings:

max, $K = 2$, $S = 2$

max, $K = 3$, $S = 2$ (AlexNet)

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

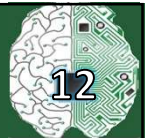
Training ConvNet



- Split and preprocess your data
- Choose your network architecture
- Initialize the weights
- Find a learning rate and regularization strength
- Minimize the loss and monitor progress

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

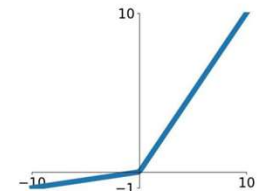
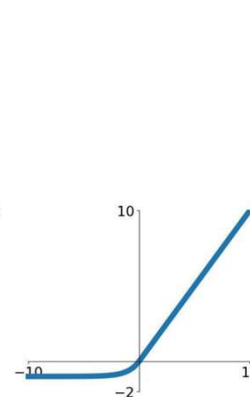
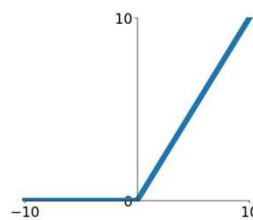
Training: Best practices



- Center (subtract mean from) your data
- Use *He*, *Xavier initialization* for weights
- Use RELU or leaky RELU or ELU or PReLU
- Use batch normalization
- Use data augmentation
- Use regularization
- Dropout
- Use mini-batch
- Learning rate: too high? Too low?
- Use cross-validation for hyperparameters

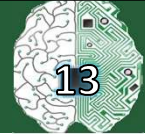
$$w_i \sim N[0, \sigma] \quad \sigma = \sqrt{\frac{2}{f_{an.in}}}$$

$$W_{i,j}^{[l]} = \mathcal{N}\left(0, \frac{1}{n^{[l-1]}}\right)$$



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Data Augmentation

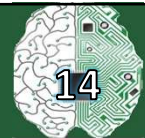


- Horizontal Flips
- Random crops and scales
- Random mix/combinations of:
 - Translation
 - Rotation
 - Stretching
 - Shearing
 - lens distortions



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Regularization



L2 regularization
$$L_{\text{reg}} = \lambda \frac{1}{2} \|W\|_2^2$$

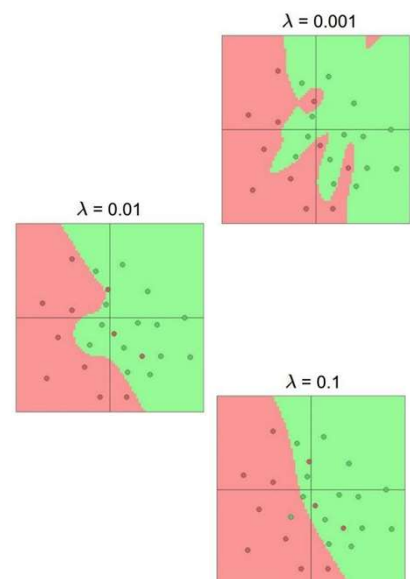
(L2 regularization encourages small weights)

L1 regularization
$$L_{\text{reg}} = \lambda \|W\|_1 = \lambda \sum_{ij} |w_{ij}|$$

(L1 regularization encourages sparse weights:
weights are encouraged to reduce to exactly zero)

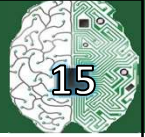
“Elastic net”
$$L_{\text{reg}} = \lambda_1 \|W\|_1 + \lambda_2 \|W\|_2^2$$

(combine L1 and L2 regularization)

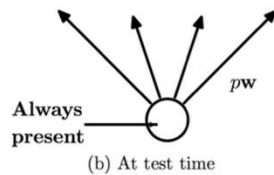
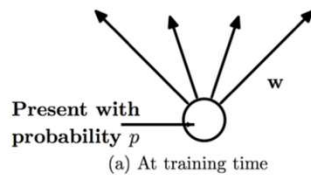


Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Dropout

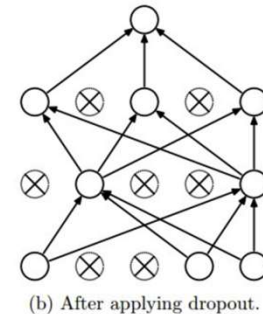
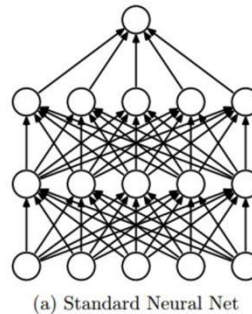
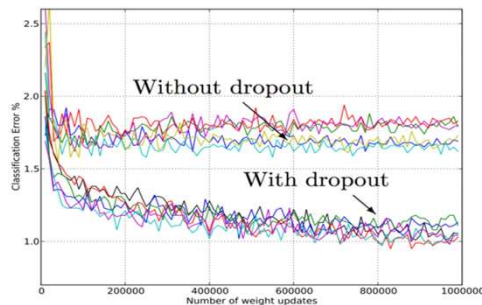


Simple but powerful technique to reduce overfitting:



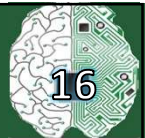
Note: Dropout can be interpreted as an approximation to taking the geometric mean of an ensemble of exponentially many models

[Srivasta et al, "Dropout: A Simple Way to Prevent Neural Networks from Overfitting", JMLR 2014]



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

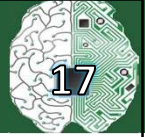
Mini-batch gradient descent



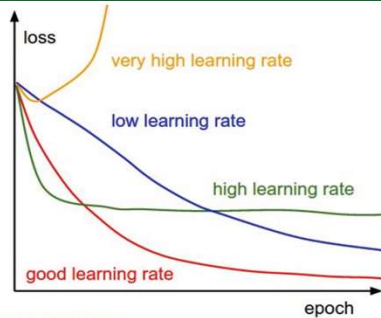
- In classic gradient descent, we compute the gradient from the loss for all training examples
- Could also only use some of the data for each gradient update
- We cycle through all the training examples multiple times
- Each time we've cycled through all of them once is called an 'epoch'
- Allows faster training (e.g. on GPUs), parallelization

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

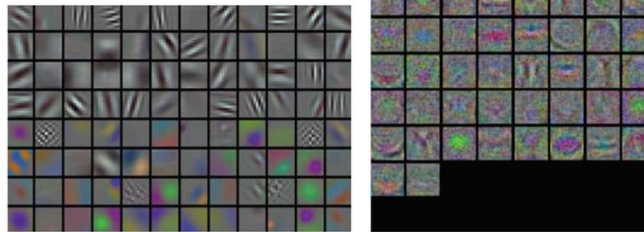
Finding a learning rate



- Plot the Loss

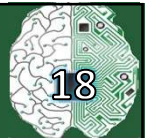


- Visualize the weights



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Determining best number of hidden units

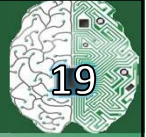


- Too few hidden units prevent the network from adequately fitting the data.
- Too many hidden units can result in over-fitting.



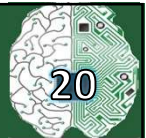
- Use internal cross-validation to empirically determine an optimal number of hidden units.

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad



Batch Normalization

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

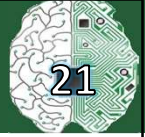


Batch Normalization

- Consider a single layer $y = Wx$
- The following could lead to tough optimization:
 - Inputs x are not centered around zero (need large bias)
 - Inputs x have different scaling per-element (entries in W will need to vary a lot)
- Idea: force inputs to be “nicely scaled” at each layer!

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Batch Normalization



- **Idea:** “Normalize” the inputs of a layer so they have zero mean and unit variance

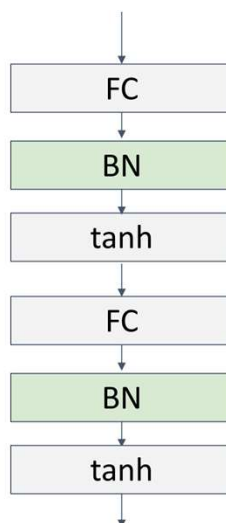
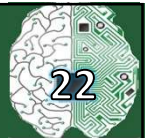
- We can normalize a batch of activations like this:

$$\hat{x} = \frac{x - E[x]}{\sqrt{\text{Var}[x]}}$$

- This is a differentiable function, so we can use it as an operator in our networks and backprop through it!

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Batch Normalization

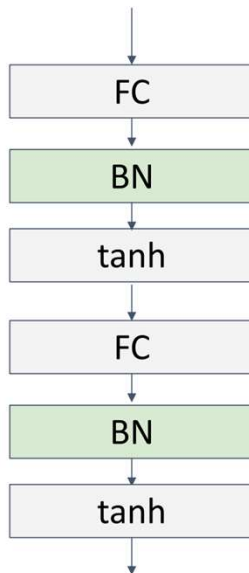


Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.

$$\hat{x} = \frac{x - E[x]}{\sqrt{\text{Var}[x]}}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

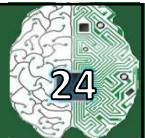
Batch Normalization



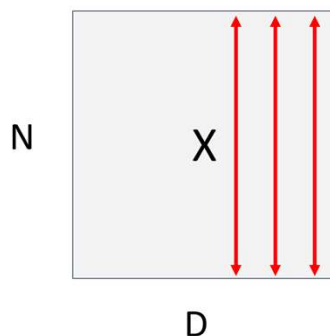
- Makes deep networks **much** easier to train!
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!
- **Not well-understood theoretically (yet)**
- **Behaves differently during training and testing: this is a very common source of bugs!**

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Batch Normalization



Input: $x \in \mathbb{R}^{N \times D}$



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel
mean, shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

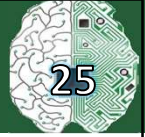
Per-channel
std, shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Normalized x,
Shape is N x D

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Batch Normalization



Input: $x \in \mathbb{R}^{N \times D}$

Learnable scale and shift parameters:

$$\gamma, \beta \in \mathbb{R}^D$$

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean, shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel std, shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

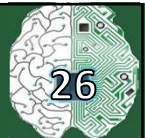
Normalized x,
Shape is N x D

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output,
Shape is N x D

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Batch Normalization



Input: $x \in \mathbb{R}^{N \times D}$

Learnable scale and shift parameters:

$$\gamma, \beta \in \mathbb{R}^D$$

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean, shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel std, shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Normalized x,
Shape is N x D

Problem: Estimates depend on minibatch; can't do this at test-time!

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output,
Shape is N x D

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Batch Normalization



Input: $x \in \mathbb{R}^{N \times D}$

Learnable scale and shift parameters:

$$\gamma, \beta \in \mathbb{R}^D$$

$\mu_j =$ (Running) average of values seen during training

Per-channel mean, shape is D

$\sigma_j^2 =$ (Running) average of values seen during training

Per-channel std, shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Normalized x,
Shape is N x D

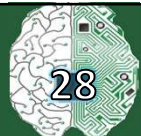
During testing batchnorm becomes a linear operator! Can be fused with the previous fully-connected or conv layer

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output,
Shape is N x D

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Batch Normalization for ConvNets



Batch Normalization for **fully-connected** networks

$$x : N \times D$$

Normalize

$$\mu, \sigma : 1 \times D$$

$$\gamma, \beta : 1 \times D$$

$$y = \frac{(x - \mu)}{\sigma} \gamma + \beta$$

Batch Normalization for **convolutional** networks
(Spatial Batchnorm, BatchNorm2D)

$$x : N \times C \times H \times W$$

Normalize

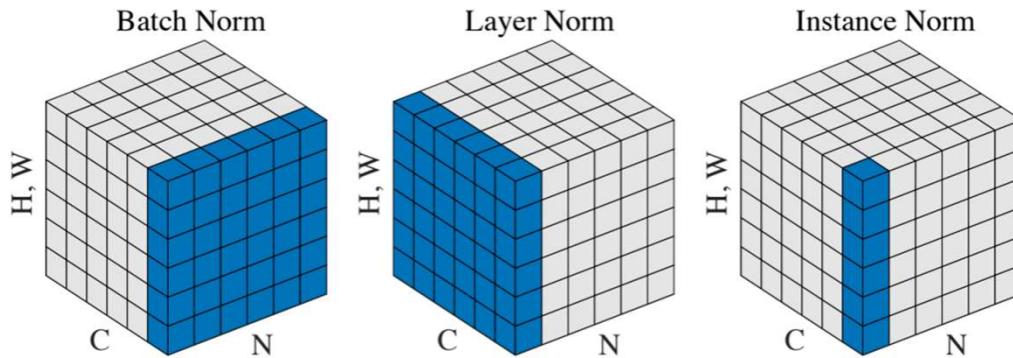
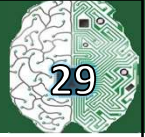
$$\mu, \sigma : 1 \times C \times 1 \times 1$$

$$\gamma, \beta : 1 \times C \times 1 \times 1$$

$$y = \frac{(x - \mu)}{\sigma} \gamma + \beta$$

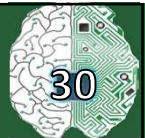
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Comparison of Normalization Layers



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Layer Normalization



Batch Normalization for
fully-connected networks

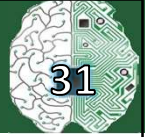
$$\begin{aligned}
 &x : N \times D \\
 &\text{Normalize} \downarrow \\
 &\mu, \sigma : 1 \times D \\
 &\gamma, \beta : 1 \times D \\
 &y = \frac{(x - \mu)}{\sigma} \gamma + \beta
 \end{aligned}$$

Layer Normalization for fully-
connected networks
Same behavior at train and test!
Used in RNNs, Transformers

$$\begin{aligned}
 &x : N \times D \\
 &\text{Normalize} \downarrow \\
 &\mu, \sigma : N \times 1 \\
 &\gamma, \beta : 1 \times D \\
 &y = \frac{(x - \mu)}{\sigma} \gamma + \beta
 \end{aligned}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Instance Normalization



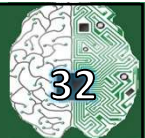
Batch Normalization for
convolutional networks

$$\begin{array}{l}
 x : N \times C \times H \times W \\
 \text{Normalize} \quad \downarrow \quad \downarrow \quad \downarrow \\
 \mu, \sigma : 1 \times C \times 1 \times 1 \\
 \gamma, \beta : 1 \times C \times 1 \times 1 \\
 y = \frac{(x - \mu)}{\sigma} \gamma + \beta
 \end{array}$$

Instance Normalization for
convolutional networks

$$\begin{array}{l}
 x : N \times C \times H \times W \\
 \text{Normalize} \quad \downarrow \quad \downarrow \quad \downarrow \\
 \mu, \sigma : N \times C \times 1 \times 1 \\
 \gamma, \beta : 1 \times C \times 1 \times 1 \\
 y = \frac{(x - \mu)}{\sigma} \gamma + \beta
 \end{array}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad



Questions?

Acknowledgement: Prof. Dr. Y. LeCun, Facebook AI Research, New York, NY, USA

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad