

Intro to Deep Neural Nets (Contd)

CS771: Introduction to Machine Learning

Piyush Rai

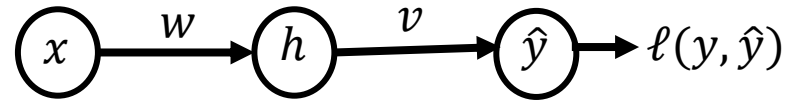
Plan today

- Training deep neural nets using backpropagation
- Some important aspects related to training of deep neural nets
 - Vanishing/exploding gradients
 - Initialization
 - Normalization layers (batch and layer normalization)
 - Dropout as a means to regularization
 - Residual/skip connections



Backpropagation via a Simple Example

- Consider a single scalar input, single hidden layer with one node, and scalar output



- Derivative of the loss ℓ w.r.t. w is $\frac{\partial \ell}{\partial w} = \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h} \frac{\partial h}{\partial w}$
- Derivative of the loss ℓ w.r.t. v is $\frac{\partial \ell}{\partial v} = \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial v}$
- $\frac{\partial \hat{y}}{\partial h} = \frac{\partial \hat{y}}{\partial v} \frac{\partial v}{\partial h}$

Row i contains the gradient vector of $y_i = f_i(\mathbf{x})$ w.r.t. \mathbf{x}

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = J^f \quad J_{ij}^f = \frac{\partial y_i}{\partial x_j} \quad J^f = \begin{bmatrix} \nabla f_1(\mathbf{x}) \\ \vdots \\ \nabla f_Q(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^{Q \times P}$$

For a function $f: \mathbb{R}^P \rightarrow \mathbb{R}^Q$, the Jacobian $J^f \in \mathbb{R}^{Q \times P}$ (a matrix)

Likewise, for a function $f: \mathbb{R}^{P \times Q} \rightarrow \mathbb{R}^{R \times S}$, the Jacobian $J^f \in \mathbb{R}^{R \times S \times P \times Q}$ (4D tensor)

More generally, for a function $f: \mathbb{R}^{I_1 \times I_2 \times \dots} \rightarrow \mathbb{R}^{O_1 \times O_2 \times \dots}$, its Jacobian $J^f \in \mathbb{R}^{O_1 \times O_2 \times \dots \times I_1 \times I_2 \times \dots}$ (a tensor)



Background: Gradient and Jacobian

- Let $\mathbf{y} = f(\mathbf{x})$, where $f: \mathbb{R}^P \rightarrow \mathbb{R}^Q$, $\mathbf{x} \in \mathbb{R}^P$, $\mathbf{y} \in \mathbb{R}^Q$. Denote $\mathbf{y} = [f_1(\mathbf{x}), \dots, f_Q(\mathbf{x})]$
- The gradient of each component $y_i = f_i(\mathbf{x}) \in \mathbb{R}$ ($i = 1, 2, \dots, Q$) w.r.t. $\mathbf{x} \in \mathbb{R}^P$ is

$$\nabla f_i(\mathbf{x}) = \frac{\partial y_i}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_i}{\partial x_1} & \dots & \frac{\partial y_i}{\partial x_P} \end{bmatrix} \in \mathbb{R}^{1 \times P}$$

Note: Gradient expressed here as a **row vector** (has the same length as \mathbf{x} which is a column vector) for notational convenience later

- Likewise, the gradient of whole vector $\mathbf{y} \in \mathbb{R}^Q$ w.r.t. vector $\mathbf{x} \in \mathbb{R}^P$ can be defined using the $Q \times P$ **Jacobian** matrix J^f whose rows consist of the above gradients

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = J^f \quad J_{ij}^f = \frac{\partial y_i}{\partial x_j} \quad J^f = \begin{bmatrix} \nabla f_1(\mathbf{x}) \\ \vdots \\ \nabla f_Q(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^{Q \times P}$$

Row i contains the gradient vector of $y_i = f_i(\mathbf{x})$ w.r.t. \mathbf{x}

For a function $f: \mathbb{R}^P \rightarrow \mathbb{R}^Q$, the Jacobian $J^f \in \mathbb{R}^{Q \times P}$ (a matrix)

Likewise, for a function $f: \mathbb{R}^{P \times Q} \rightarrow \mathbb{R}^{R \times S}$, the Jacobian $J^f \in \mathbb{R}^{R \times S \times P \times Q}$ (4D **tensor**)

More generally, for a function $f: \mathbb{R}^{I_1 \times I_2 \times \dots} \rightarrow \mathbb{R}^{O_1 \times O_2 \times \dots}$, its Jacobian $J^f \in \mathbb{R}^{O_1 \times O_2 \times \dots \times I_1 \times I_2 \times \dots}$ (a tensor)



Background: Multivariate Chain Rule of Calculus

- Let $\mathbf{x} \in \mathbb{R}^P$, $\mathbf{y} = g(\mathbf{x}) \in \mathbb{R}^Q$, $z = f(\mathbf{y}) \in \mathbb{R}$, where $g: \mathbb{R}^P \rightarrow \mathbb{R}^Q$, $f: \mathbb{R}^Q \rightarrow \mathbb{R}$

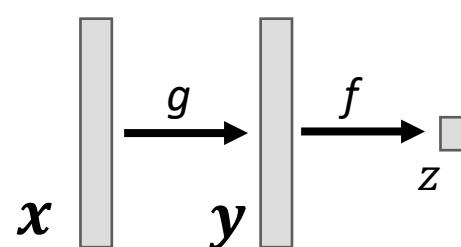


Diagram illustrating the chain rule: $\mathbf{x} \xrightarrow{g} \mathbf{y} \xrightarrow{f} z$. \mathbf{x} is a $1 \times P$ vector derivative, \mathbf{y} is a $1 \times Q$ vector derivative, and z is a scalar derivative.

$$\frac{\partial z}{\partial \mathbf{x}} = \sum_{i=1}^Q \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial \mathbf{x}}$$

Used chain rule of total derivatives
Sum is needed since z depends on \mathbf{y} , and \mathbf{y} is a vector

- The above can be written as a product of a vector and a matrix

Turns out to be a product of Jacobian of f and g in that order ☺

$$\frac{\partial z}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial z}{\partial y_1} & \dots & \frac{\partial z}{\partial y_Q} \end{bmatrix} \times \begin{bmatrix} \frac{\partial y_1}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial y_Q}{\partial \mathbf{x}} \end{bmatrix} = \nabla f(\mathbf{y}) \times \begin{bmatrix} \nabla g_1(\mathbf{x}) \\ \vdots \\ \nabla g_P(\mathbf{x}) \end{bmatrix} = \nabla f(\mathbf{y}) \times J^g$$

$1 \times Q$ gradient (same as Jacobian since f is scalar)
 $Q \times P$ Jacobian of g

- More generally, let $\mathbf{w} \in \mathbb{R}^P$, $\mathbf{x} = h(\mathbf{w}) \in \mathbb{R}^Q$, $\mathbf{y} = g(\mathbf{x}) \in \mathbb{R}^R$, $z = f(\mathbf{y}) \in \mathbb{R}^S$

Product of the 3 Jacobians in that order (simple! ☺)

$$\frac{\partial \mathbf{z}}{\partial \mathbf{w}} = J^f \times J^g \times J^h \in \mathbb{R}^{S \times P}$$

Note that chain rule for scalar variables $\mathbf{w}, \mathbf{x}, \mathbf{y}, \mathbf{z}$ is defined in a similar way as $\frac{\partial z}{\partial \mathbf{w}} = f'(\mathbf{y})g'(\mathbf{x})h'(\mathbf{w})$

Intro to ML

Backpropagation (Backprop)

- Backprop is **gradient descent with multivariate chain rule for derivatives**
- Consider a two hidden layer neural network

$$\mathcal{L}(\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \mathbf{v}) = \sum_{n=1}^N \ell(y_n, \hat{y}_n) = \sum_{n=1}^N \ell_n$$

- We wish to minimize the loss
- The gradient based updates will be

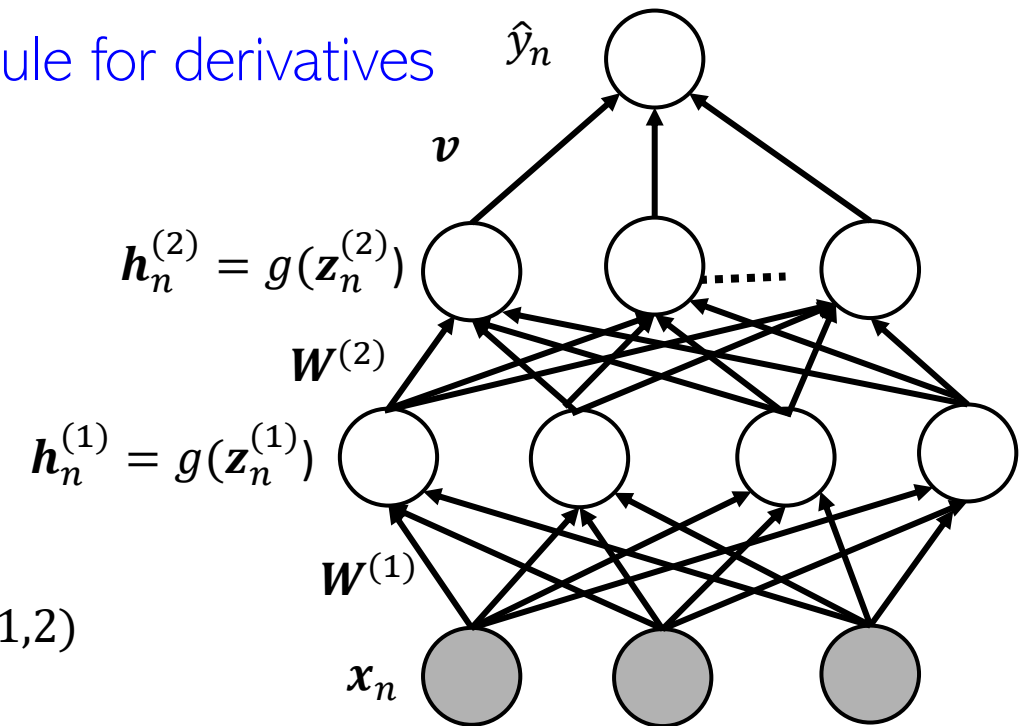
$$\mathbf{v} = \mathbf{v} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{v}} \quad \mathbf{W}^{(i)} = \mathbf{W}^{(i)} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(i)}} \quad (i = 1, 2)$$

- Since $\mathcal{L} = \sum_{n=1}^N \ell_n$, we need to compute $\frac{\partial \ell_n}{\partial \mathbf{v}}$ and $\frac{\partial \ell_n}{\partial \mathbf{W}^{(i)}} (i = 1, 2)$

- Assume output activation σ as identity ($\hat{y}_n = \mathbf{v}^\top \mathbf{h}_n^{(2)}$)

$$\frac{\partial \ell_n}{\partial \mathbf{v}} = \frac{\partial \ell_n}{\partial \hat{y}_n} \frac{\partial \hat{y}_n}{\partial \mathbf{v}} = \ell'_n(y_n, \hat{y}_n) \mathbf{h}_n^{(2)}$$

Derivative of ℓ_n w.r.t. \hat{y}_n



Backpropagation in detail

- Let's now look at $\frac{\partial \ell_n}{\partial \mathbf{W}^{(2)}}$ where $\ell_n = \ell(y_n, \hat{y}_n)$ and $\hat{y}_n = \mathbf{v}^\top \mathbf{h}_n^{(2)}$

$$\frac{\partial \ell_n}{\partial \mathbf{W}^{(2)}} = \frac{\partial \ell_n}{\partial \hat{y}_n} \frac{\partial \hat{y}_n}{\partial \mathbf{W}^{(2)}} = \ell'(y_n, \hat{y}_n) \frac{\partial \hat{y}_n}{\partial \mathbf{W}^{(2)}}$$

$$\frac{\partial \hat{y}_n}{\partial \mathbf{W}^{(2)}} = \frac{\partial \hat{y}_n}{\partial \mathbf{h}_n^{(2)}} \frac{\partial \mathbf{h}_n^{(2)}}{\partial \mathbf{W}^{(2)}} + \frac{\partial \hat{y}_n}{\partial \mathbf{v}} \frac{\partial \mathbf{v}}{\partial \mathbf{W}^{(2)}}$$

- Since \mathbf{v} doesn't depend on $\mathbf{W}^{(2)}$, $\frac{\partial \mathbf{v}}{\partial \mathbf{W}^{(2)}} = 0$

Using transpose since we assume gradient to be a row vector

Jacobian of size $1 \times K_2 \times K_1$

$$\frac{\partial \ell_n}{\partial \mathbf{W}^{(2)}} = \ell'(y_n, \hat{y}_n) \frac{\partial \hat{y}_n}{\partial \mathbf{h}_n^{(2)}} \frac{\partial \mathbf{h}_n^{(2)}}{\partial \mathbf{W}^{(2)}} = \ell'(y_n, \hat{y}_n) \mathbf{v}^\top \frac{\partial \mathbf{h}_n^{(2)}}{\partial \mathbf{W}^{(2)}}$$

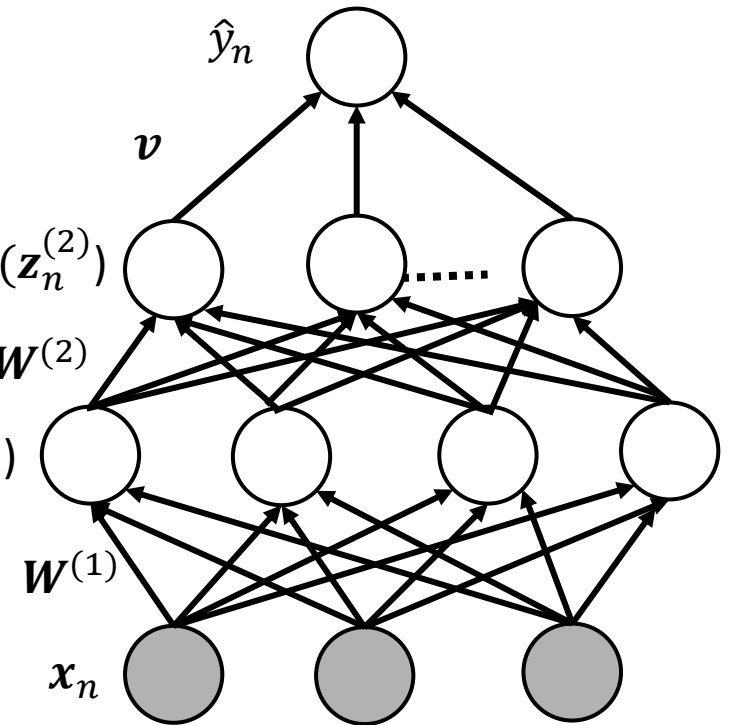
- We now need $\frac{\partial \mathbf{h}_n^{(2)}}{\partial \mathbf{W}^{(2)}}$. Using $\mathbf{h}_n^{(2)} = g(\mathbf{z}_n^{(2)})$ where $\mathbf{z}_n^{(2)} = \mathbf{W}^{(2)\top} \mathbf{h}_n^{(1)}$ and g is elementwise applied nonlinearity on the vector $\mathbf{z}_n^{(2)}$

Jacobian of size $K_2 \times K_2 \times K_1$

$$\frac{\partial \mathbf{h}_n^{(2)}}{\partial \mathbf{W}^{(2)}} = \frac{\partial \mathbf{h}_n^{(2)}}{\partial \mathbf{z}_n^{(2)}} \frac{\partial \mathbf{z}_n^{(2)}}{\partial \mathbf{W}^{(2)}} = \text{diag} \left(g' \left(z_{n1}^{(2)} \right), \dots, g' \left(z_{nK_2}^{(2)} \right) \right) \frac{\partial \mathbf{z}_n^{(2)}}{\partial \mathbf{W}^{(2)}}$$

Diagonal matrix of size $K_2 \times K_2$ with Jacobian (gradient vector) of g along the diagonals

This Jacobian is a tensor of size $K_2 \times K_2 \times K_1$



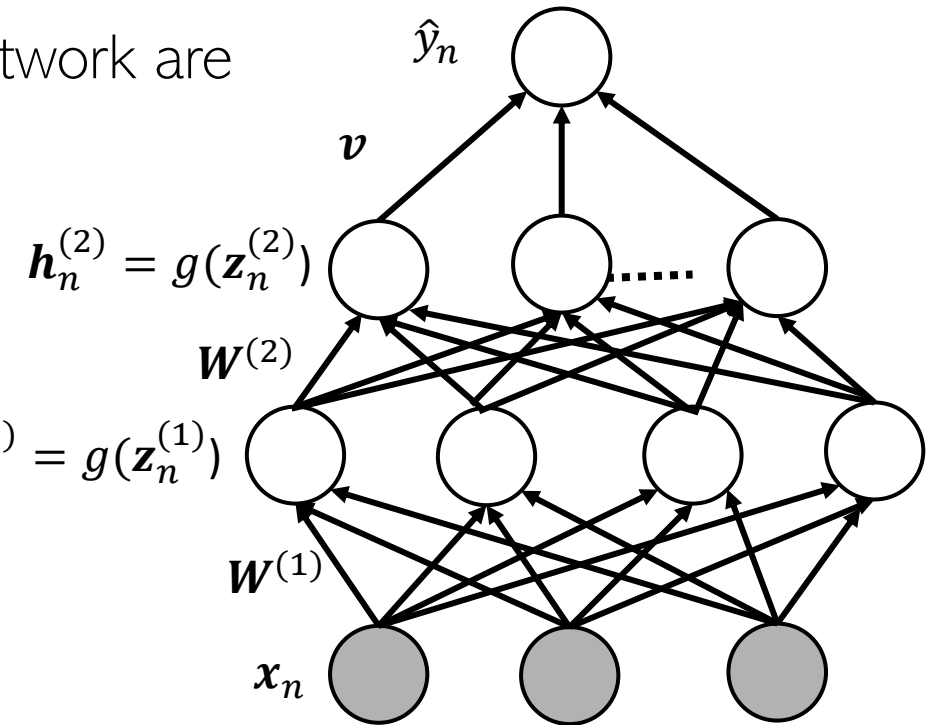
Backpropagation: Computation Reuse

- Summarizing, the required gradients/Jacobians for this network are

$$\frac{\partial \ell_n}{\partial \mathbf{v}} = \ell'(y_n, \hat{y}_n) \mathbf{h}_n^{(2)}$$

$$\frac{\partial \ell_n}{\partial \mathbf{W}^{(2)}} = \ell'(y_n, \hat{y}_n) \frac{\partial \hat{y}_n}{\partial \mathbf{h}_n^{(2)}} \frac{\partial \mathbf{h}_n^{(2)}}{\partial \mathbf{z}_n^{(2)}} \frac{\partial \mathbf{z}_n^{(2)}}{\partial \mathbf{W}^{(2)}}$$

$$\frac{\partial \ell_n}{\partial \mathbf{W}^{(1)}} = \ell'(y_n, \hat{y}_n) \frac{\partial \hat{y}_n}{\partial \mathbf{h}_n^{(2)}} \frac{\partial \mathbf{h}_n^{(2)}}{\partial \mathbf{z}_n^{(2)}} \frac{\partial \mathbf{z}_n^{(2)}}{\partial \mathbf{h}_n^{(1)}} \frac{\partial \mathbf{h}_n^{(1)}}{\partial \mathbf{z}_n^{(1)}} \frac{\partial \mathbf{z}_n^{(1)}}{\partial \mathbf{W}^{(1)}}$$



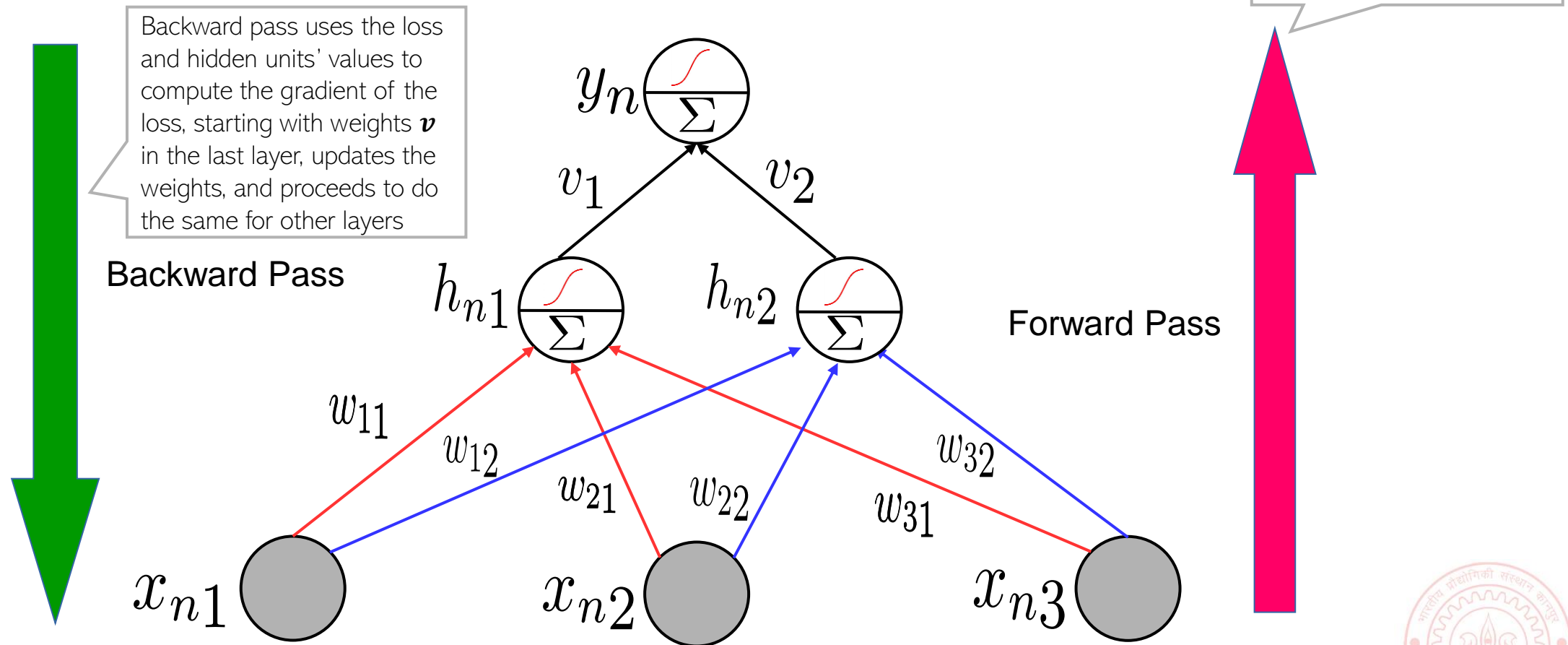
- Thus gradient computations done in upper layers can be stored and reused when computing the gradients in the lower layers (libraries like Tensorflow and Pytorch do so efficiently)

- Vanishing gradients:** $\frac{\partial \mathbf{h}_n^{(i)}}{\partial \mathbf{z}_n^{(i)}} = \text{diag} \left(g' \left(z_{n1}^{(i)} \right), \dots, g' \left(z_{nK_i}^{(i)} \right) \right)$

Gradients in lower layers will have product of many such terms $\frac{\partial \mathbf{h}_n^{(i)}}{\partial \mathbf{z}_n^{(i)}}$. If g' is small (e.g., gradient of sigmoid or tanh), the gradient becomes vanishingly small for lower layers and becomes an issue (thus ReLU and other with non-saturating activations are preferred)

Backpropagation

- Backprop iterates between a forward pass and a backward pass

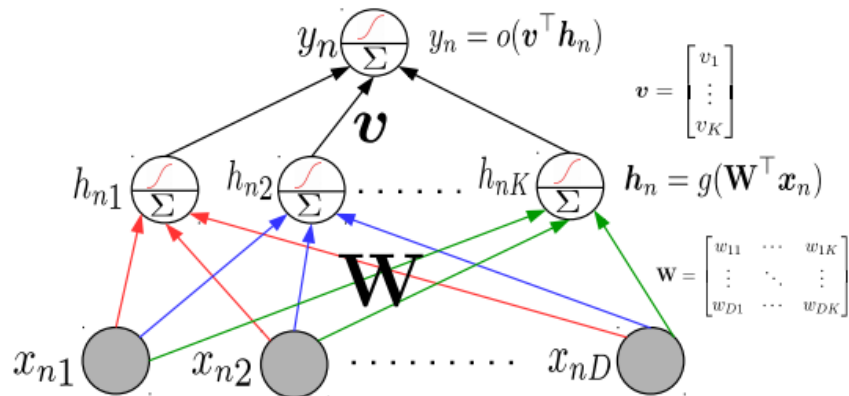


- Software frameworks such as Tensorflow and PyTorch support this already so you don't need to implement it by hand (so no worries of computing derivatives etc)



Backpropagation through an example

Consider a single hidden layer MLP



Assuming regression ($o = \text{identity}$),
the loss function for this model

$$\begin{aligned}\mathcal{L} &= \frac{1}{2} \sum_{n=1}^N \left(y_n - \mathbf{v}^\top \mathbf{h}_n \right)^2 \\ &= \frac{1}{2} \sum_{n=1}^N \left(y_n - \sum_{k=1}^K v_k h_{nk} \right)^2 \\ &= \frac{1}{2} \sum_{n=1}^N \left(y_n - \sum_{k=1}^K v_k g(\mathbf{w}_k^\top \mathbf{x}_n) \right)^2\end{aligned}$$

- To use gradient methods for \mathbf{W} , \mathbf{v} , we need gradients.
- Gradient of \mathcal{L} w.r.t. \mathbf{v} is straightforward

$$\frac{\partial \mathcal{L}}{\partial v_k} = - \sum_{n=1}^N \left(y_n - \sum_{k=1}^K v_k g(\mathbf{w}_k^\top \mathbf{x}_n) \right) h_{nk} = \sum_{n=1}^N \mathbf{e}_n h_{nk}$$

- Gradient of \mathcal{L} w.r.t. \mathbf{W} requires chain rule

$$\frac{\partial \mathcal{L}}{\partial w_{dk}} = \sum_{n=1}^N \frac{\partial \mathcal{L}}{\partial h_{nk}} \frac{\partial h_{nk}}{\partial w_{dk}}$$

$$\frac{\partial \mathcal{L}}{\partial h_{nk}} = - \left(y_n - \sum_{k=1}^K v_k g(\mathbf{w}_k^\top \mathbf{x}_n) \right) v_k = - \mathbf{e}_n v_k$$

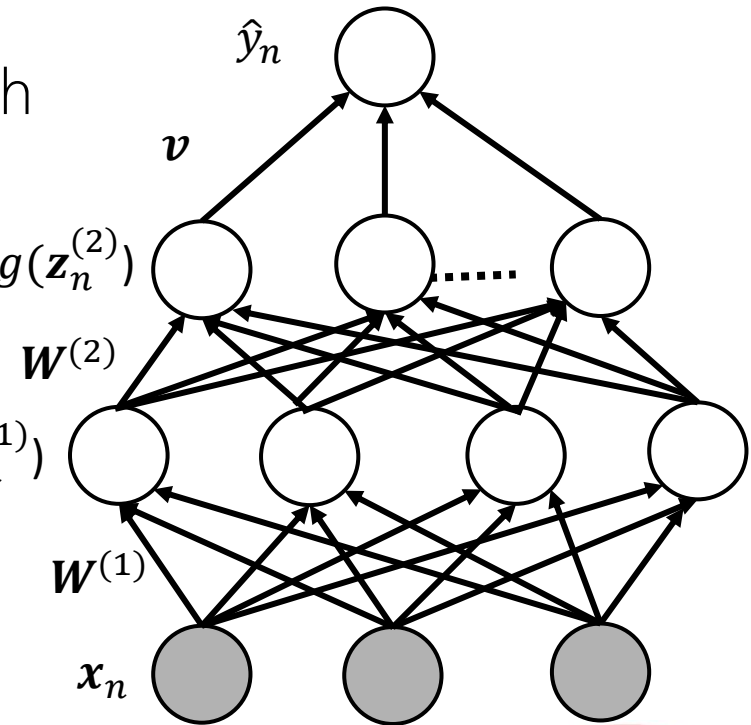
$$\frac{\partial h_{nk}}{\partial w_{dk}} = g'(\mathbf{w}_k^\top \mathbf{x}_n) x_{nd} \quad (\text{note: } h_{nk} = g(\mathbf{w}_k^\top \mathbf{x}_n))$$

- **Forward prop** computes errors \mathbf{e}_n using current \mathbf{W} , \mathbf{v} .
Backprop updates NN params \mathbf{W} , \mathbf{v} using grad methods
- Backprop caches many of the calculations for reuse



Problem of Exploding/Vanishing Gradients

- MLPs/CNNs have many hidden layers and gradients in each layer are a product of several Jacobians
- Result of these products depends on the eigenvalues of each of these Jacobians
 - If they are large (>1), gradients might blow up (explode)
 - If they are small (<1), gradients might vanish
- To prevent blow up, we can use **gradient clipping**
 - Simply cap the magnitude of the gradients!
- To prevent vanishing gradients, several options
 - Use **non-saturating activation functions** (recall that the gradient is a product of terms like $\frac{\partial h_n^{(i)}}{\partial z_n^{(i)}} = \text{diag} \left(g'(z_{n1}^{(i)}), \dots, g'(z_{nK_i}^{(i)}) \right)$, so the derivative g' doesn't become too small
 - Use other architectures such as **skip-connections** (will discuss later)



Training of DNNs: Some Important Aspects

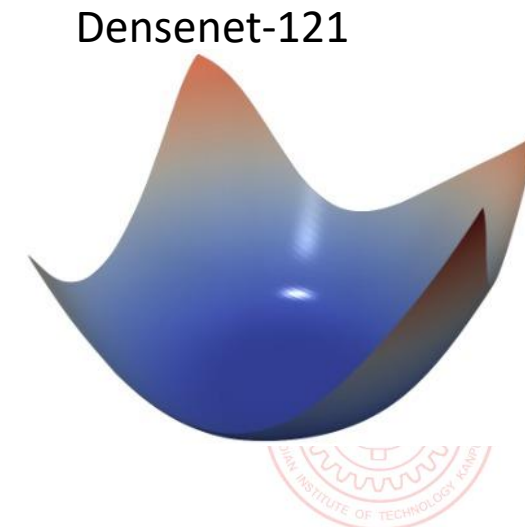
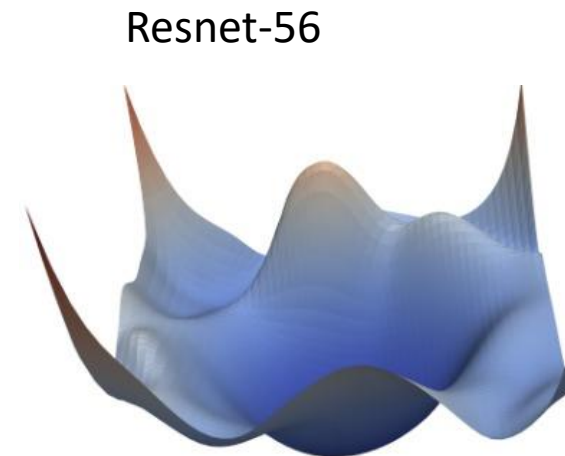
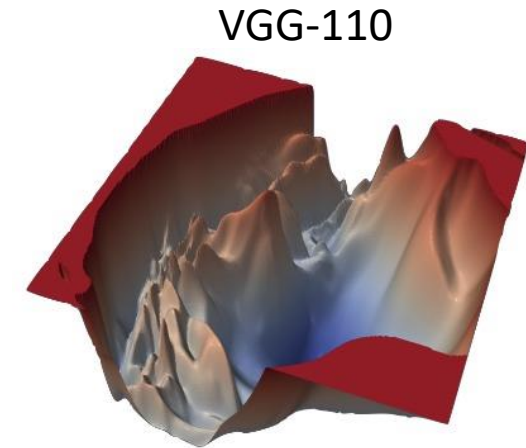
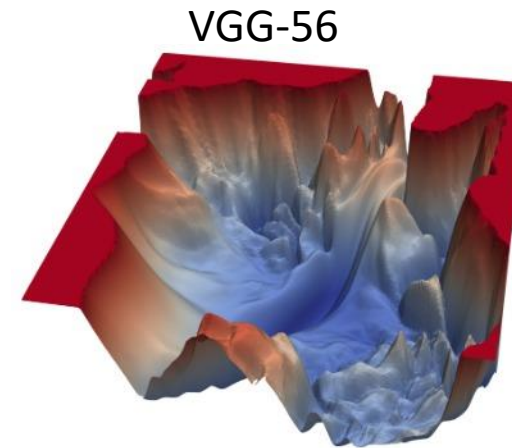
- Deep neural net training can be hard due to non-convex loss functions
- Several ways to address this, e.g.,
 - Good choice of learning rate of (S)GD
 - We have already seen this
 - Good initialization of parameters, e.g., initialize each weight, say w_{ij} , randomly as

$$w_{ij} \sim \mathcal{N}(0, \sigma^2) \quad \text{or} \quad w_{ij} \sim \text{Uniform}(-a, a)$$

Xavier/Gloret initialization,
LeCun init, He init, etc

and set the “spread” of these distribution
as inversely proportional to $n_{\text{in}} + n_{\text{out}}$

- Careful design of the network architecture, e.g.,
 - Networks with “skip connections” (will see later) which lead to less non-convex (more smooth) loss surfaces (figures on the right)
- Vanishing/exploding gradients (already saw)



Batch Normalization

Note: Batch-norm assumes sufficiently large mini-batch \mathcal{B} to work well. There are variants such as “layer normalization” and “instance normalization” that don’t require a mini-batch can be computed using a single training example

Batch normalization is used in MLP, CNN, and various other architectures



- Each hidden layer is a nonlinear transformation of the previous layer’s inputs
- To prevent distribution drift in activations’ distribution, we often “standardize” each layer
- Standardize = activation $h_{nk}^{(\ell)}$ should have zero mean and unit variance across all n
- It is achieved by inserting a “batch normalization” layer after each hidden layer
- To do so, during training, (omitting layer number ℓ) we replace each h_n by \tilde{h}_n

γ and β are trainable batch-norm parameters

$$\tilde{h}_n = \gamma \odot \hat{h}_n + \beta$$

$$\hat{h}_n = \frac{h_n - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

We compute $\mu_{\mathcal{B}}$ and $\sigma_{\mathcal{B}}^2$ using the data from the current minibatch of examples \mathcal{B} (thus the name “batch norm”)

$$\mu_{\mathcal{B}} = \frac{1}{|\mathcal{B}|} \sum_{h \in \mathcal{B}} h$$

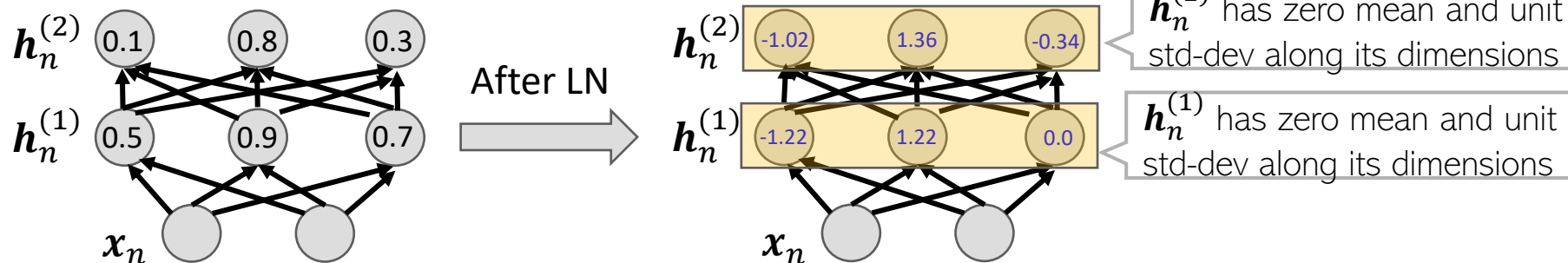
$$\sigma_{\mathcal{B}}^2 = \frac{1}{|\mathcal{B}|} \sum_{h \in \mathcal{B}} (h - \mu_{\mathcal{B}})^2$$

- After training, we store γ and β + the statistics μ and σ^2 computed on the whole training data, and use these values to apply batch-norm on each test input



Layer Normalization

- Normalization helps improve training and performance overall
- Unlike batch normalization (BN), which we already saw, layer normalization (LN) normalizes each \mathbf{h}_n across its dimensions (not across all minibatch examples)
 - Often used for sequence data models (will see later) where BN is difficult to apply
 - Also useful when batch sizes are small where BN statistics (mean/var) aren't reliable
- For an MLP, the LN operation would look like this

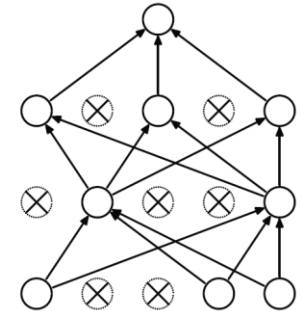
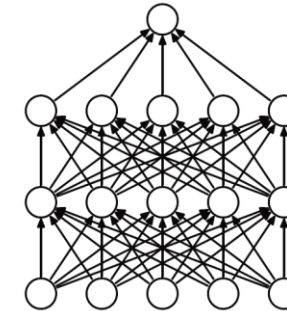


Dropout Layer

- Deep neural networks can overfit when trained on small datasets
- **Dropout** is a method to regularize without using an explicit regularizer
- In every update of the network, drop neuron i in layer ℓ with probability p

$$\epsilon_i^{(\ell)} \sim \text{Bernoulli}(1 - p)$$

- If $\epsilon_i^{(\ell)} = 0$, set all outgoing weights $w_{ij}^{(\ell)}$ from neuron i to 0

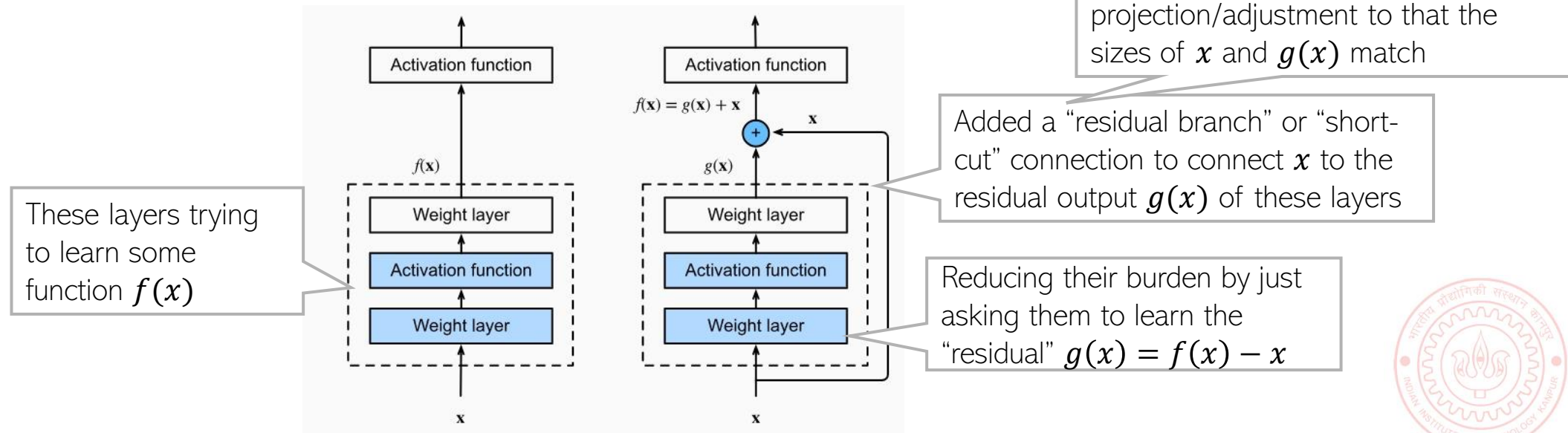


- Each update of weights will change a different subset of weights
 - In doing so, we are making individual neurons more self-reliant and less dependent on others
- At test time, no dropout is used. After training is complete, we multiply each weight by the keep probability $1 - p$ and use these weights for predictions



Residual/Skip Connections

- Many modern deep nets contain a very large number of layers
- In general, just stacking lots of layer doesn't necessarily help a deep learning model
 - Vanishing/exploding gradient may make learning difficult
- Skip connections or “residual connections” help if we want very deep networks
 - This idea was popularized by “Residual Networks”* (ResNets) which can have hundreds of layers
- Basic idea: Don't force a layer to learn everything about a mapping



*Deep Residual Learning for Image Recognition (He et al, 2015)

Pic source: <https://www.d2l.ai/index.html>

