



THE UNIVERSITY OF BRITISH COLUMBIA

# Topics in AI (CPSC 532S): Multimodal Learning with Vision, Language and Sound

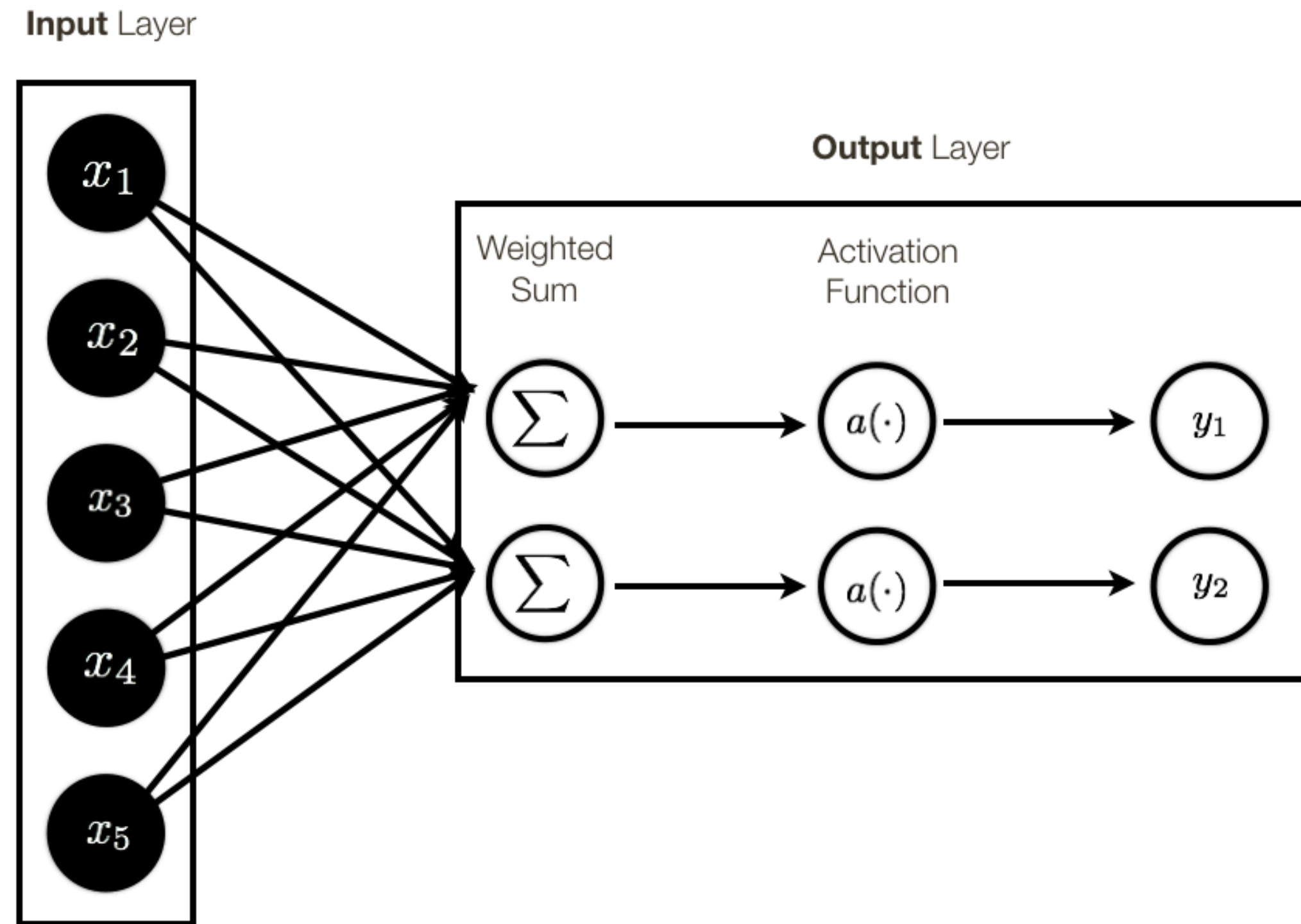
**Lecture 3: Introduction to Deep Learning (continued)**

# Course **Logistics**

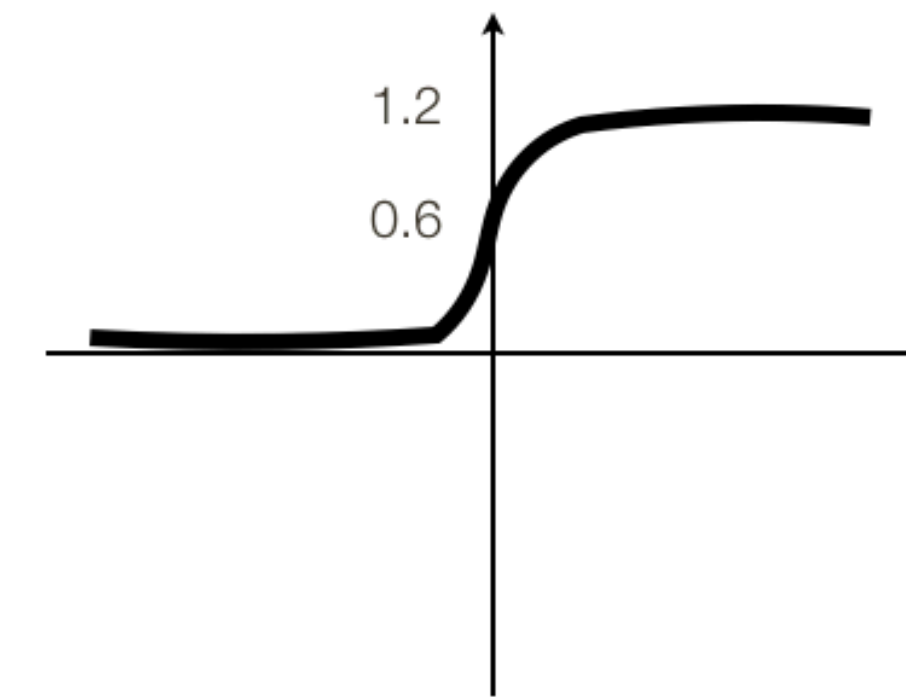
- Update on **course registrations**:
  - New Room (**SWING 409**) — 47 seats capacity (instead of 40)
  - Registrations (39 students currently registered) — approx. 8 seats left
- Microsoft **Azure** credits and tutorial next week
- **Assignment 1** ... any questions?

# Short **Review** ...

- Introduced the basic building block of Neural Networks (**MLP/FC**) **layer**



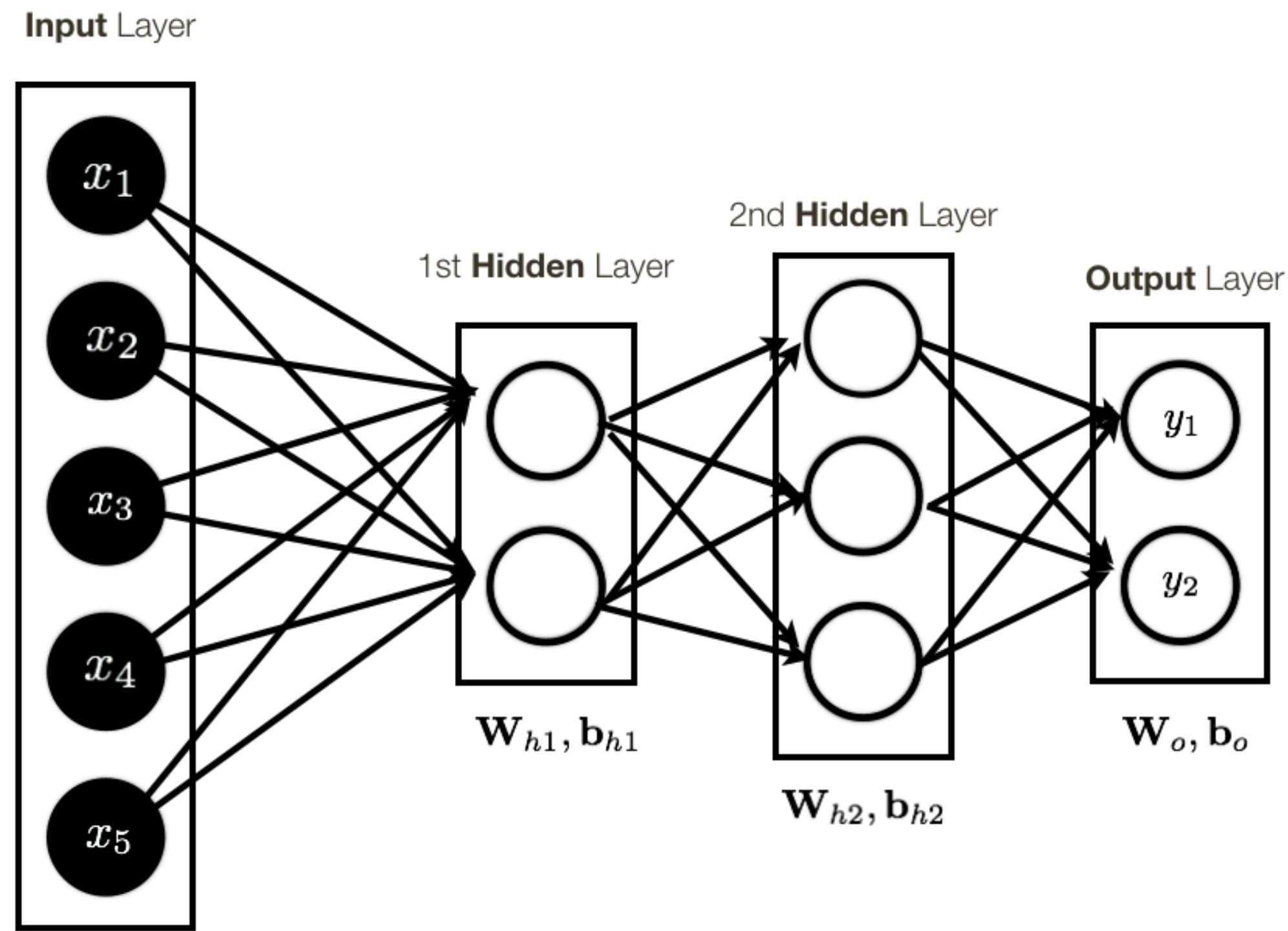
$$a(x) = \text{sigmoid}(x) = \frac{1}{1 + e^{-x}}$$



**Sigmoid** Activation

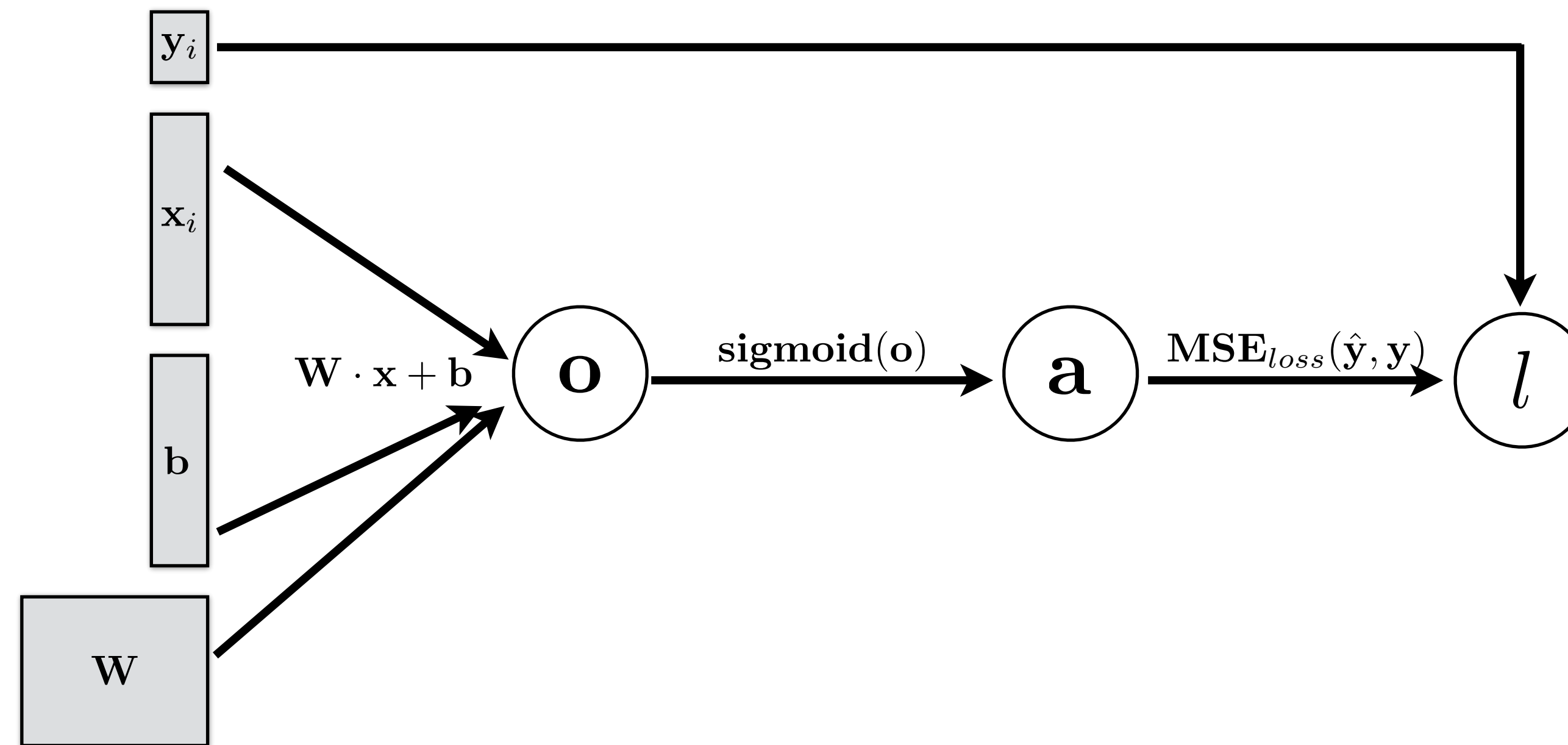
# Short **Review** ...

- Introduced the basic building block of Neural Networks **(MLP/FC) layer**
- How do we **stack these layers** up to make a Deep NN



# Short **Review** ...

- Introduced the basic building block of Neural Networks **(MLP/FC) layer**
- How do we **stack these layers** up to make a Deep NN
- Basic **NN operations** (implemented using **computational graph**)



# Short **Review** ...

- Introduced the basic building block of Neural Networks **(MLP/FC) layer**
- How do we **stack these layers** up to make a Deep NN
- Basic **NN operations** (implemented using **computational graph**)

## **Prediction** / Inference

Function evaluation

(a.k.a. **ForwardProp**)

# Short **Review** ...

- Introduced the basic building block of Neural Networks **(MLP/FC) layer**
- How do we **stack these layers** up to make a Deep NN
- Basic **NN operations** (implemented using **computational graph**)

## Prediction / Inference

Function evaluation

(a.k.a. **ForwardProp**)

## Parameter **Learnings**

(Stochastic) Gradient Descent (needs **derivatives**)

- **Numerical** differentiation (not accurate)
- **Symbolic** differential (intractable)
- AutoDiff **Forward** (computationally expensive)
- AutoDiff **Backward / BackProp**



# Short **Review** ...

- Introduced the basic building block of Neural Networks (**MLP/FC**) **layer**
- How do we **stack these layers** up to make a Deep NN
- Basic **NN operations** (implemented using **computational graph**)

## Prediction / Inference

Function evaluation

(a.k.a. **ForwardProp**)

## Parameter **Learnings**

(Stochastic) Gradient Descent (needs **derivatives**)

- **Numerical** differentiation (not accurate)
- **Symbolic** differential (intractable)
- AutoDiff **Forward** (computationally expensive)
- AutoDiff **Backward / BackProp**

- Different **activation functions** and saturation problem



# Regularization: L2 or L1 on the weights

**L2 Regularization:** Learn a more (dense) distributed representation

$$R(\mathbf{W}) = \|\mathbf{W}\|_2 = \sum_i \sum_j \mathbf{w}_{i,j}^2$$

**L1 Regularization:** Learn a sparse representation (few non-zero weight elements)

$$R(\mathbf{W}) = \|\mathbf{W}\|_1 = \sum_i \sum_j |\mathbf{w}_{i,j}|$$

(others regularizers are also possible)

## Example:

$$\mathbf{x} = [1, 1, 1, 1]$$

$$\mathbf{W}_1 = [1, 0, 0, 0]$$

$$\mathbf{W}_2 = \left[ \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4} \right]$$

$$\mathbf{W}_1 \cdot \mathbf{x}^T = \mathbf{W}_2 \cdot \mathbf{x}^T$$

two networks will have identical loss

L2 Regularizer:

$$R_{L2}(\mathbf{W}_1) = 1$$

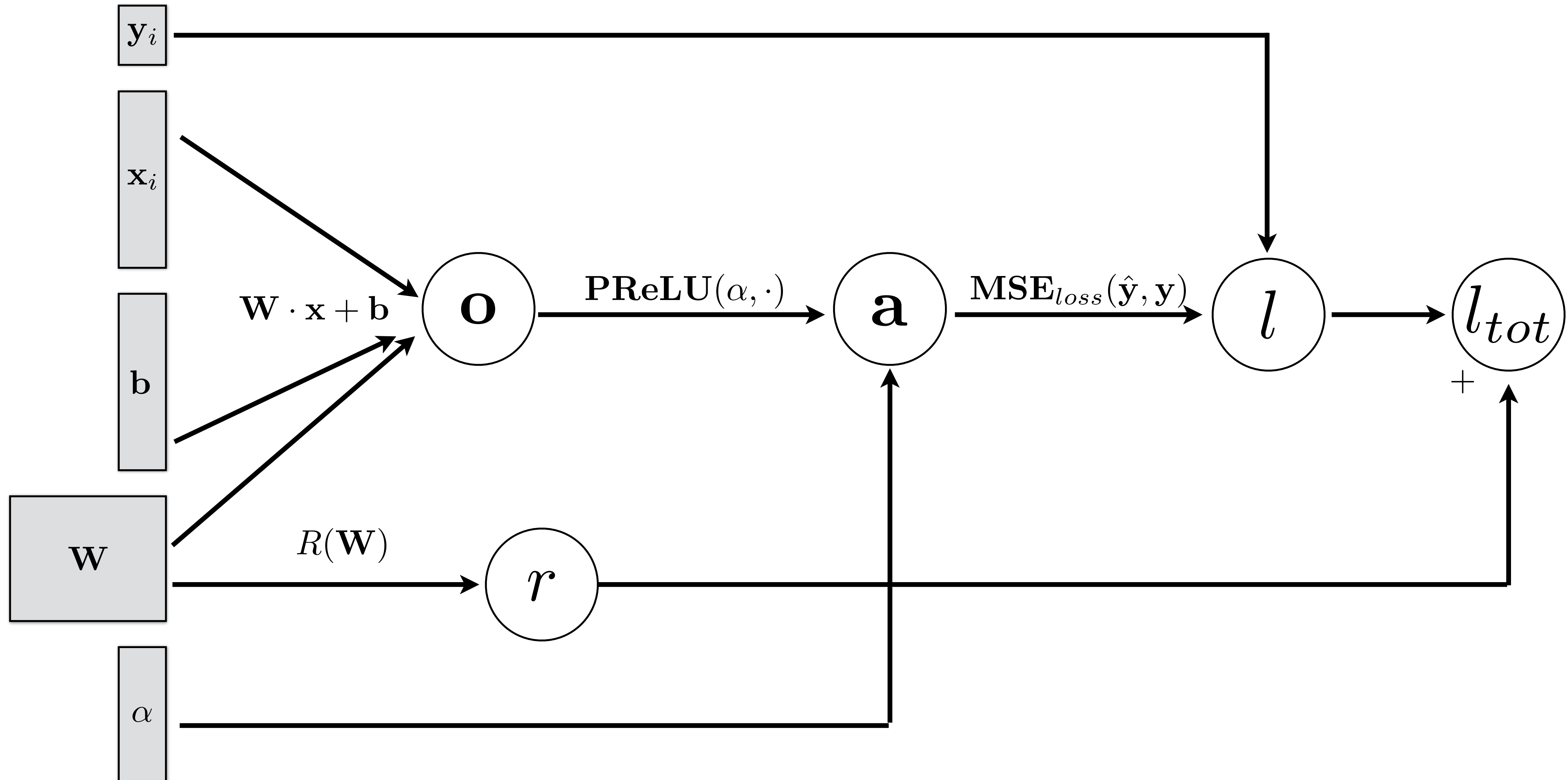
$$R_{L2}(\mathbf{W}_2) = 0.25 \quad \leftarrow$$

L1 Regularizer:

$$R_{L1}(\mathbf{W}_1) = 1 \quad \leftarrow$$

$$R_{L1}(\mathbf{W}_2) = 1 \quad \leftarrow$$

# Computational Graph: 1-layer with PReLU + Regularizer



# Regularization: Batch Normalization

**Normalize each mini-batch** (using Batch Normalization layer) by subtracting empirically computed mean and dividing by variance for every dimension -> samples are approximately unit Gaussian

$$\bar{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

## Benefit:

Improves learning (better gradients, higher learning rate)

# Regularization: Batch Normalization

**Normalize each mini-batch** (using Batch Normalization layer) by subtracting empirically computed mean and dividing by variance for every dimension -> samples are approximately unit Gaussian

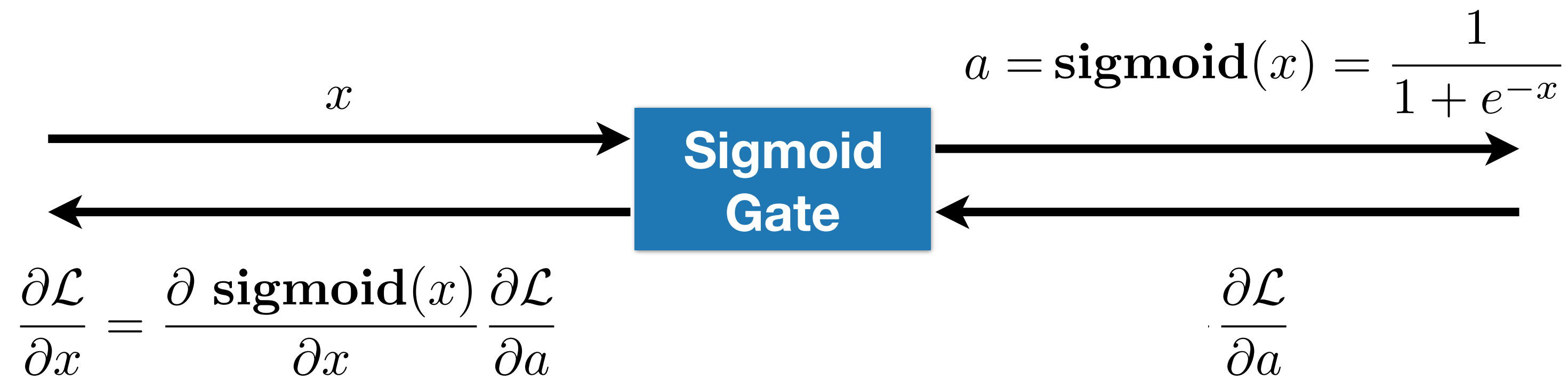
$$\bar{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

## Benefit:

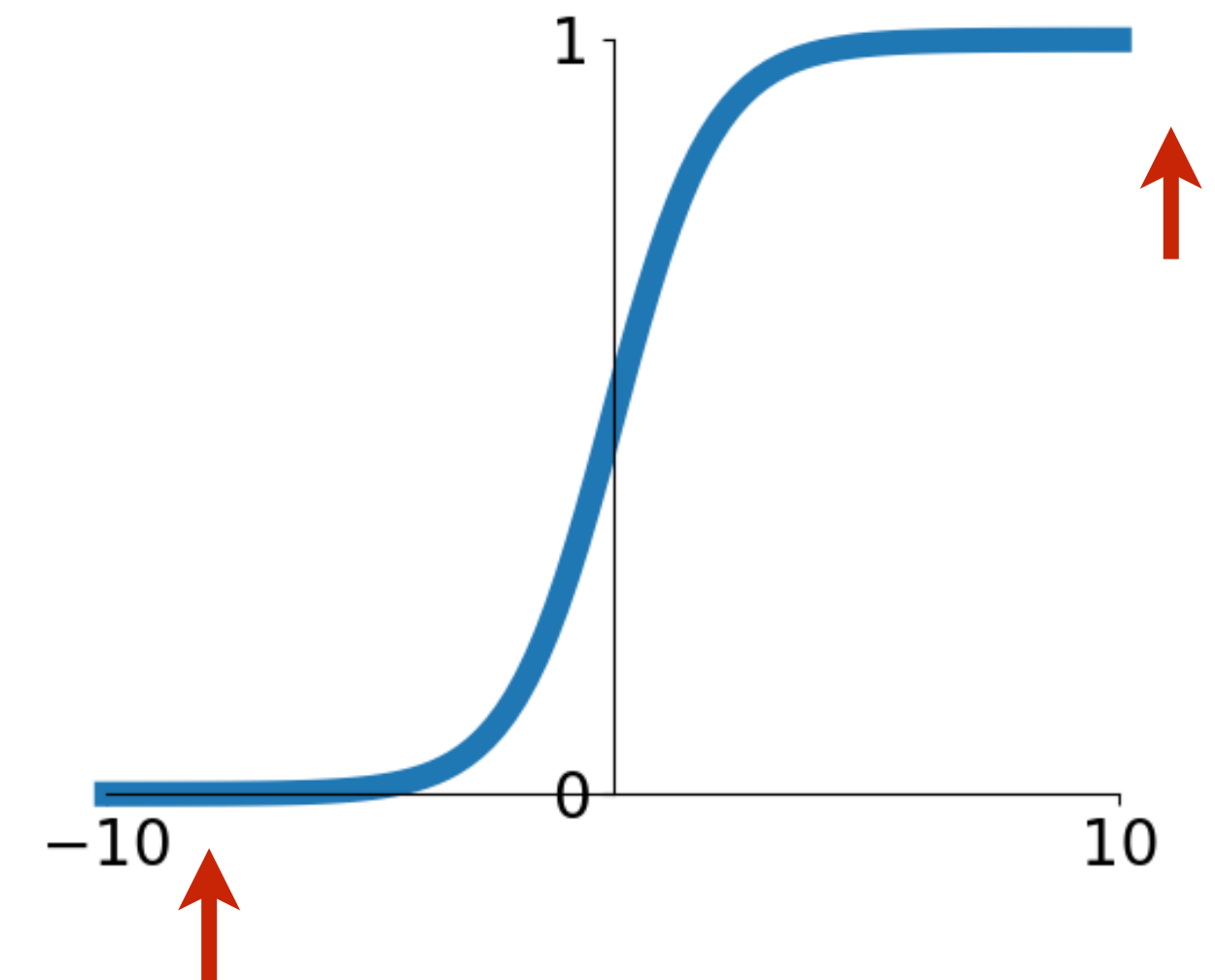
Improves learning (better gradients, higher learning rate)

**Why?**

# Activation Function: Sigmoid



$$a(x) = \text{sigmoid}(x) = \frac{1}{1 + e^{-x}}$$



## Cons:

- Saturated neurons **“kill” the gradients**
- Non-zero centered
- Could be expensive to compute

**Sigmoid** Activation

# Regularization: Batch Normalization

**Normalize each mini-batch** (using Batch Normalization layer) by subtracting empirically computed mean and dividing by variance for every dimension -> samples are approximately unit Gaussian

$$\bar{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

## **Benefit:**

Improves learning (better gradients, higher learning rate)

Typically inserted **before** activation layer

# Activation Function: Sigmoid vs. Tanh

$$a(x) = \tanh(x) = 2 \cdot \text{sigmoid}(2x) - 1$$

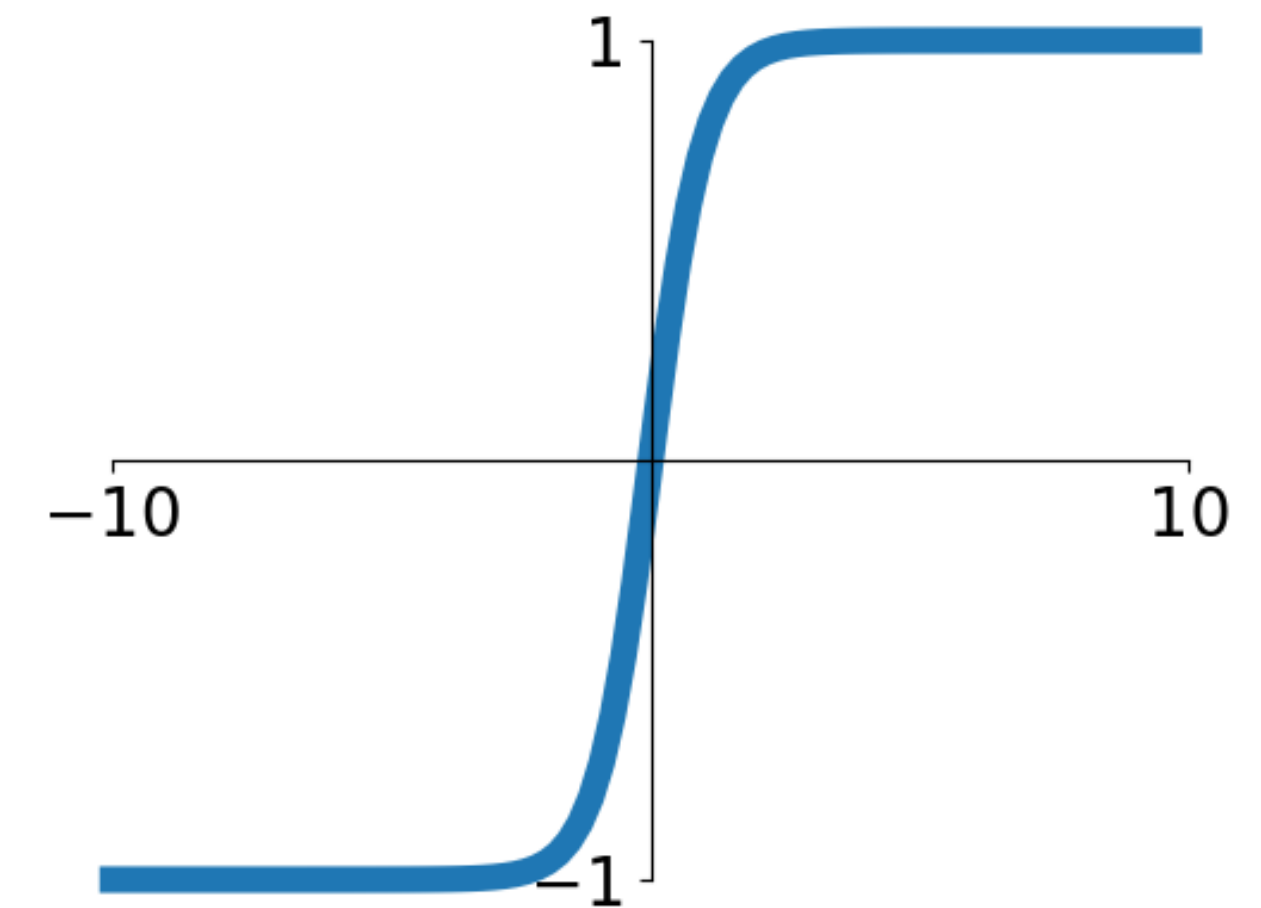
$$a(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$$

## Pros:

- Squishes everything in the range  $[-1, 1]$
- Centered around zero
- Has well defined gradient everywhere

## Cons:

- Saturated neurons “kill” the gradients



**Tanh** Activation



# Regularization: Batch Normalization

**Normalize each mini-batch** (using Batch Normalization layer) by subtracting empirically computed mean and dividing by variance for every dimension -> samples are approximately unit Gaussian

$$\bar{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

In practice, also learn how to scale and offset:

$$y^{(k)} = \gamma^{(k)} \bar{x}^{(k)} + \beta^{(k)}$$

BN layer parameters

## Benefit:

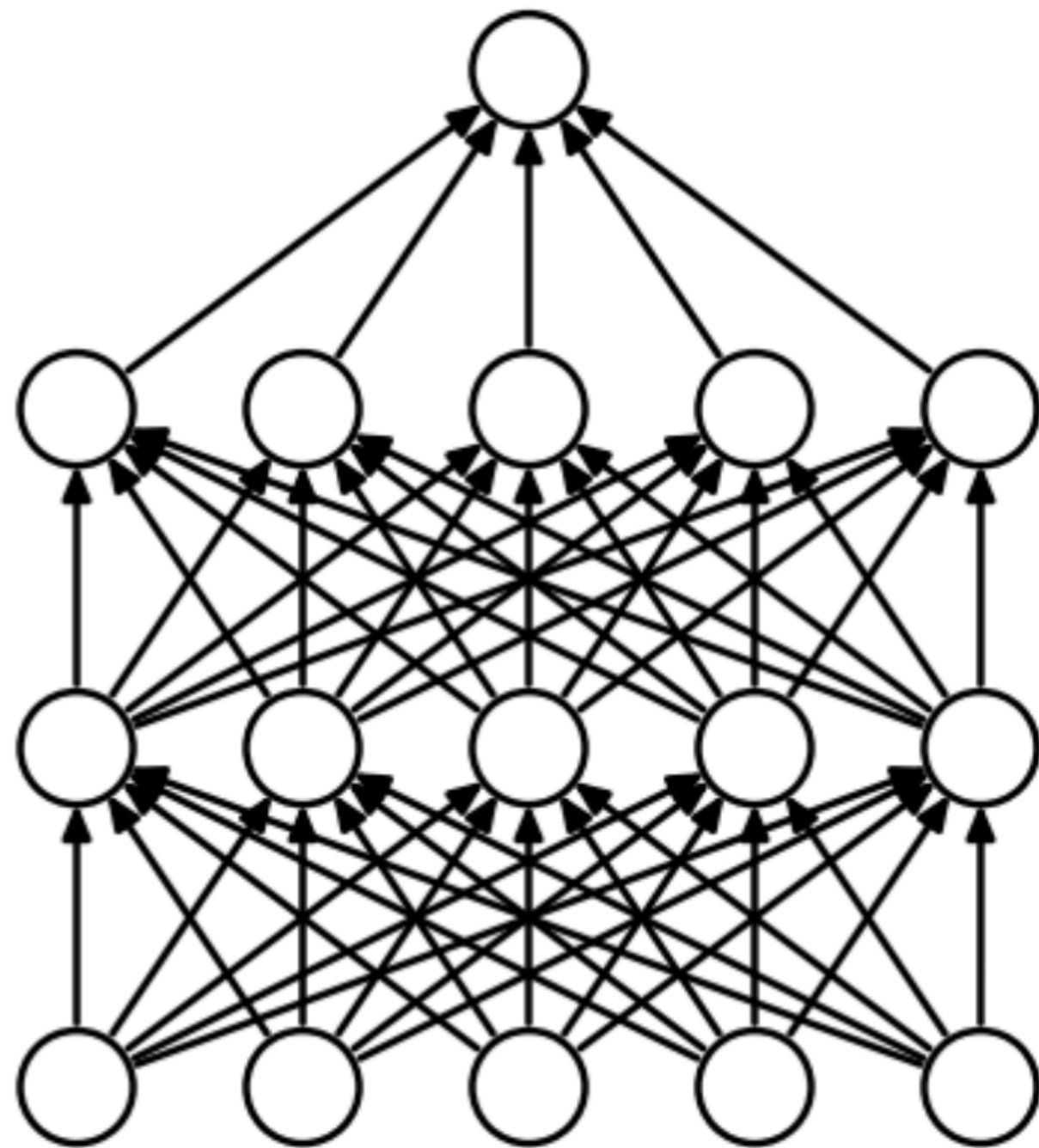
Improves learning (better gradients, higher learning rate)

Typically inserted **before** activation layer

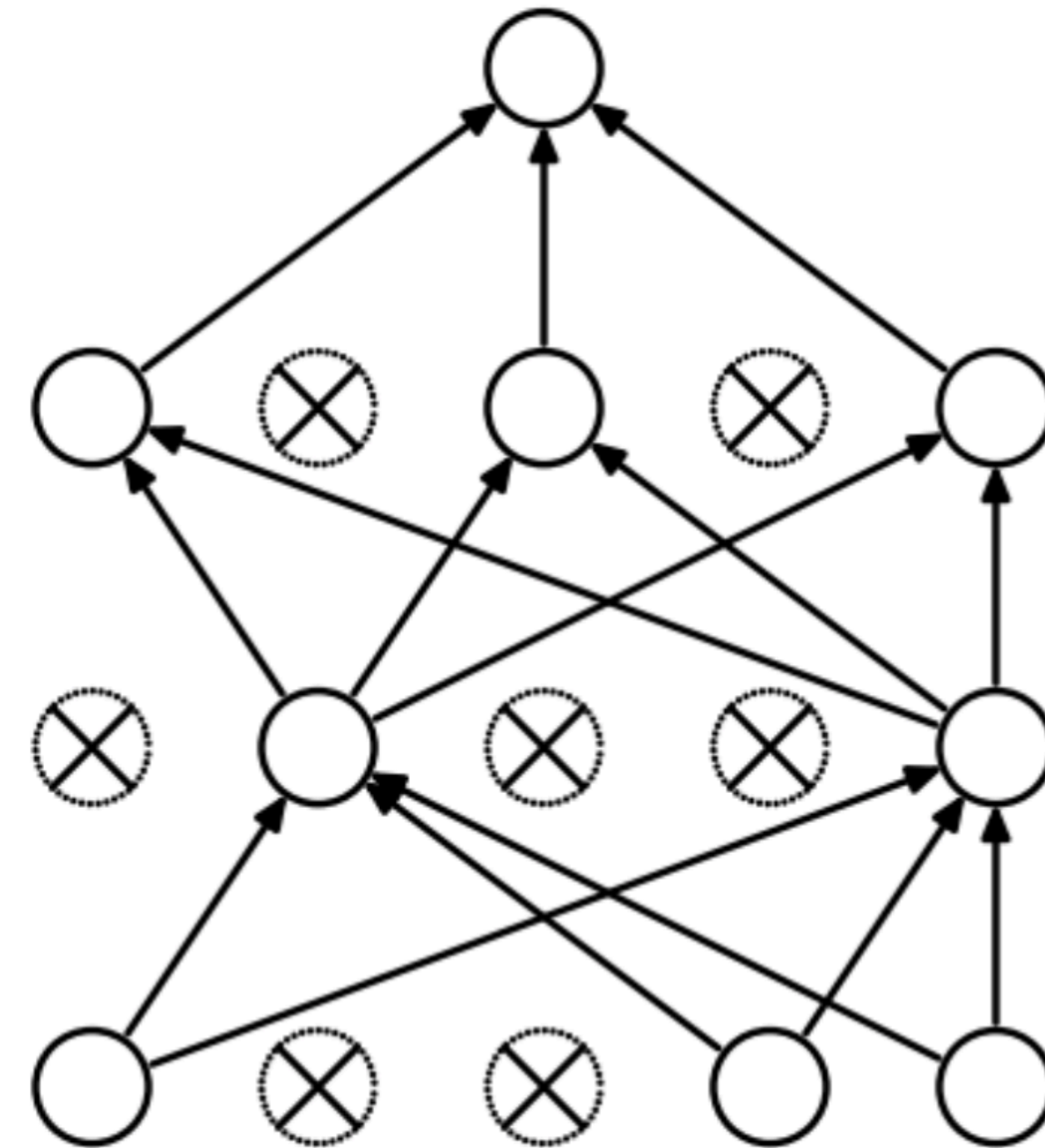
[ Ioffe and Szegedy, NIPS 2015 ]

# Regularization: Dropout

Randomly **set some neurons to zero** in the forward pass, with probability proportional to `dropout rate` (between 0 to 1)



Standar Neural Network



After Applying **Dropout**

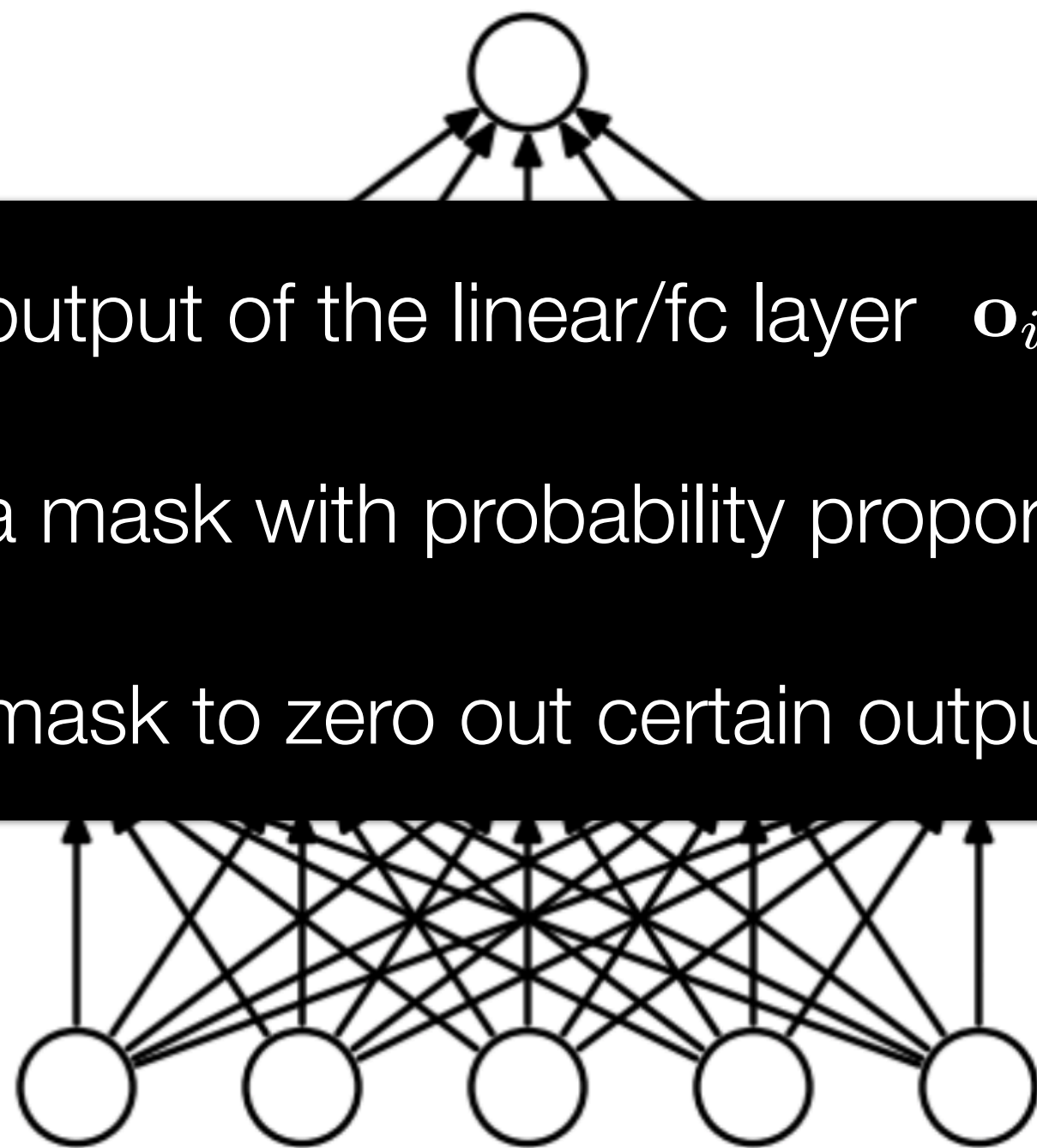
[ Srivastava et al, JMLR 2014 ]

\* adopted from slides of **CS231n at Stanford**

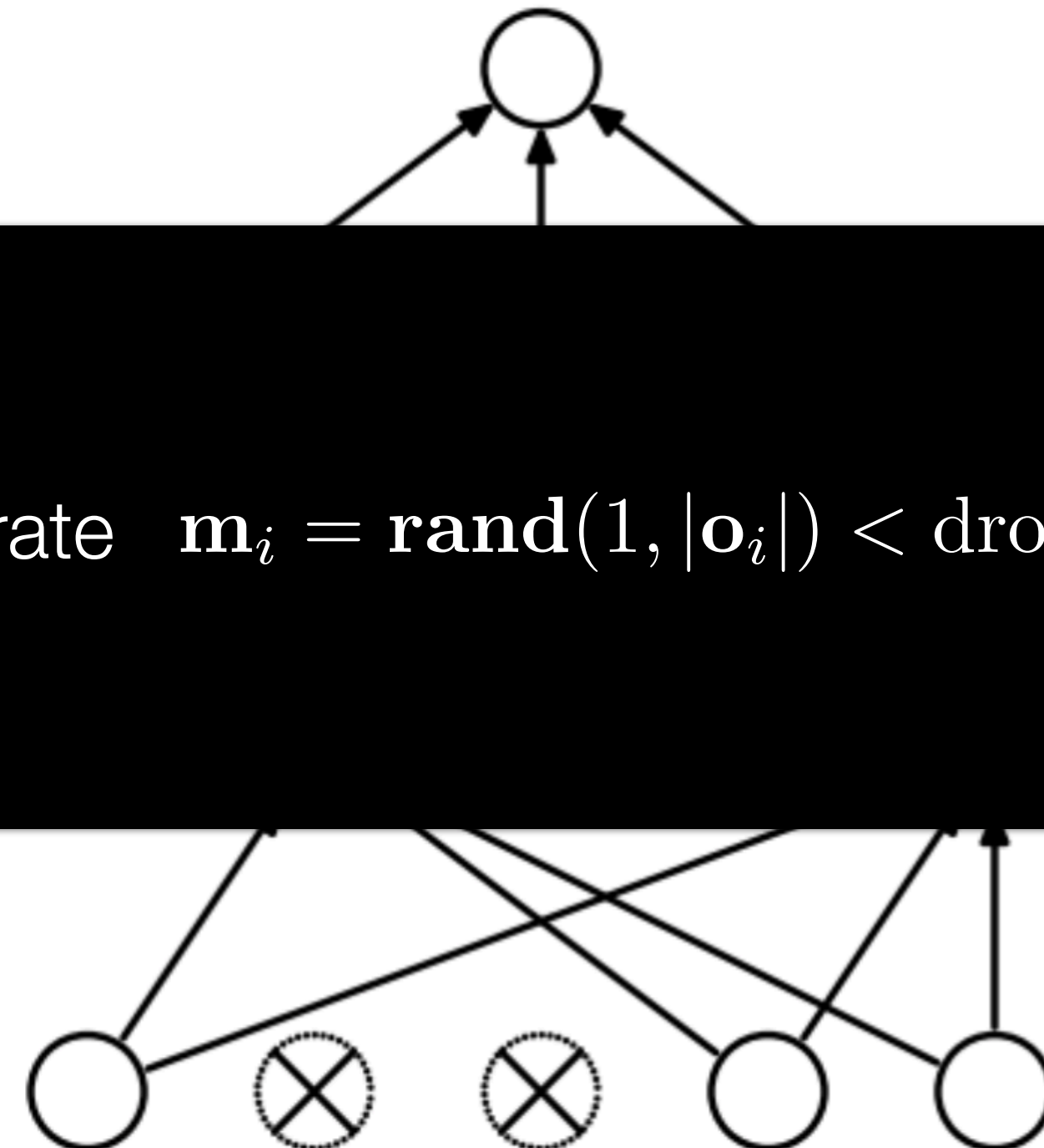
# Regularization: Dropout

Randomly **set some neurons to zero** in the forward pass, with probability proportional to dropout rate (between 0 to 1)

1. Compute output of the linear/fc layer  $\mathbf{o}_i = \mathbf{W}_i \cdot \mathbf{x} + \mathbf{b}_i$
2. Compute a mask with probability proportional to dropout rate  $\mathbf{m}_i = \mathbf{rand}(1, |\mathbf{o}_i|) < \text{dropout rate}$
3. Apply the mask to zero out certain outputs  $\mathbf{o}_i = \mathbf{o}_i \odot \mathbf{m}_i$



Standar Neural Network



After Applying **Dropout**

[ Srivastava et al, JMLR 2014 ]

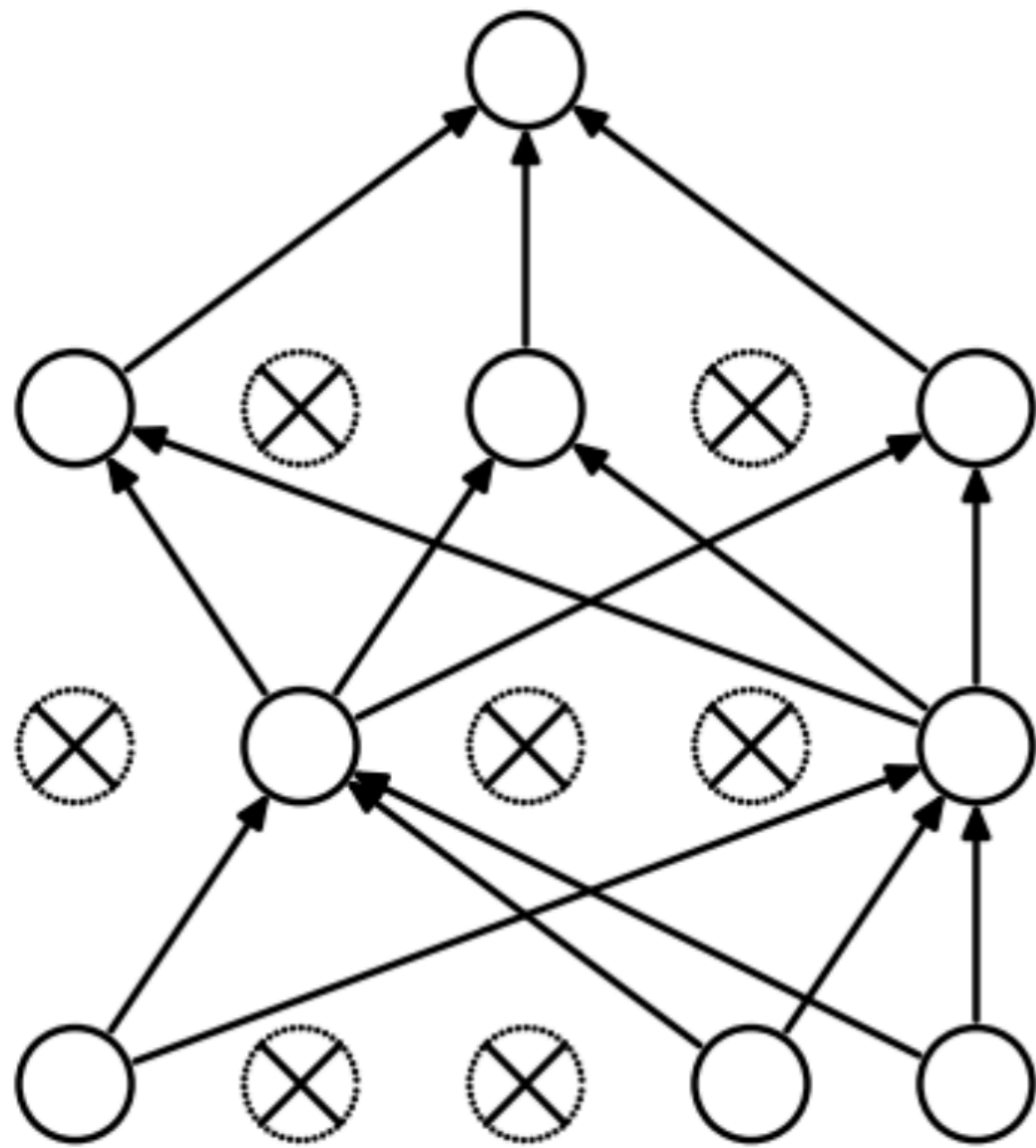
\* adopted from slides of **CS231n at Stanford**



# Regularization: Dropout

Randomly **set some neurons to zero** in the forward pass, with probability proportional to `dropout_rate` (between 0 to 1)

**Why is this a good idea?**



After Applying **Dropout**

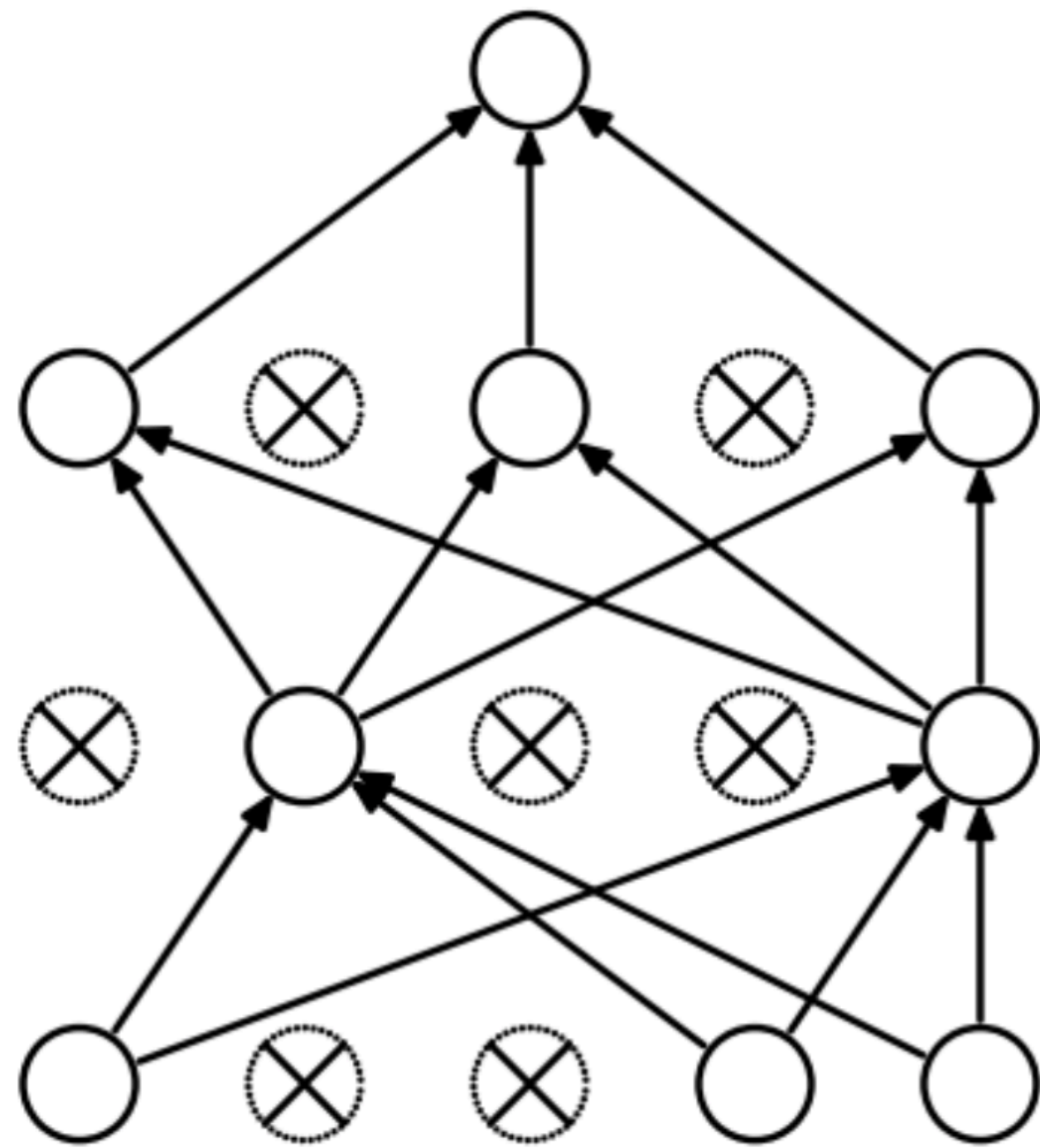
[ Srivastava et al, JMLR 2014 ]

\* adopted from slides of **CS231n at Stanford**

# Regularization: Dropout

Randomly **set some neurons to zero** in the forward pass, with probability proportional to `dropout_rate` (between 0 to 1)

**Why is this a good idea?**



After Applying **Dropout**

Dropout is training an **ensemble of models** that share parameters

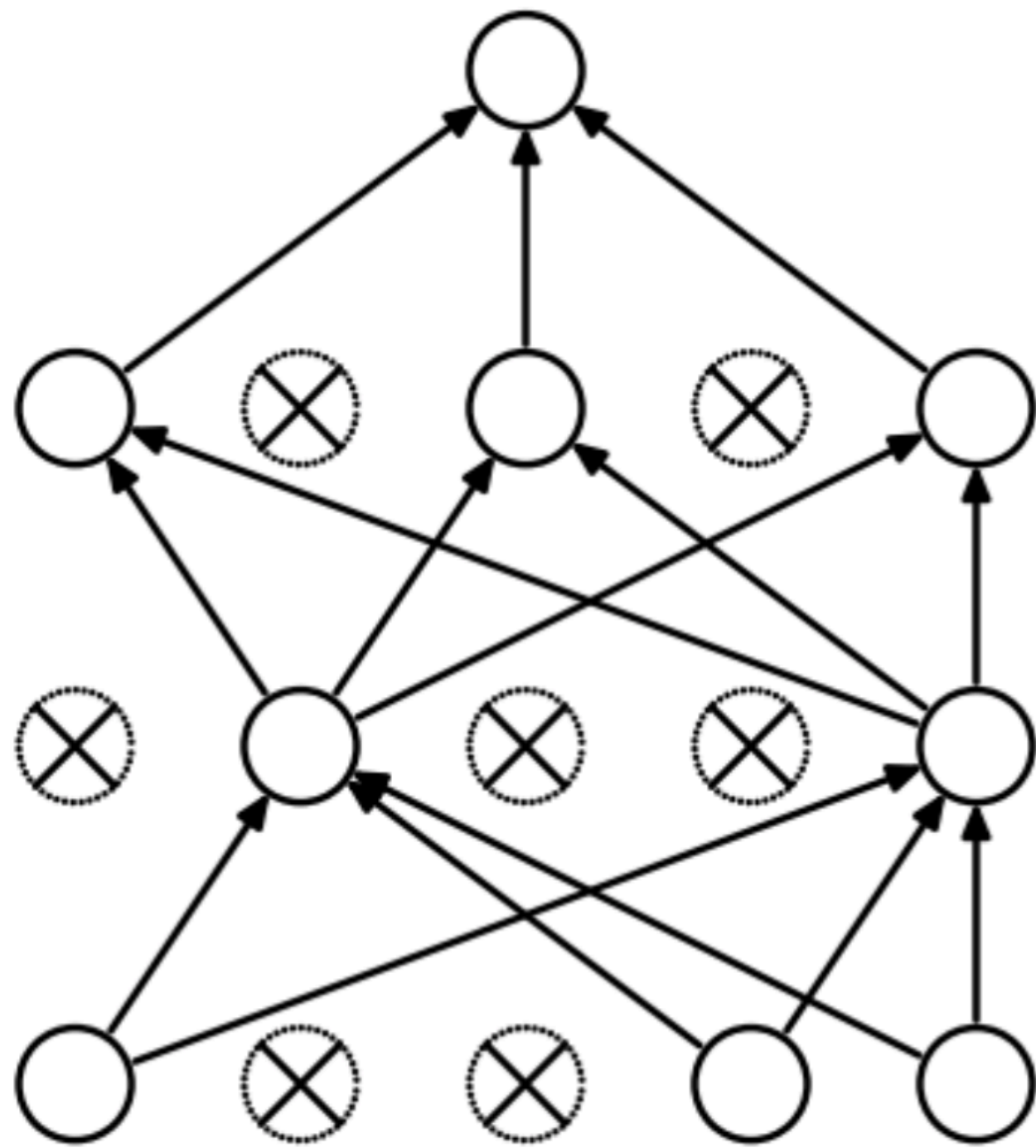
Each binary mask (generated in the forward pass) is one model that is trained on (approximately) one data point

[ Srivastava et al, JMLR 2014 ]

\* adopted from slides of **CS231n at Stanford**

# Regularization: Dropout (at test time)

Randomly **set some neurons to zero** in the forward pass, with probability proportional to `dropout_rate` (between 0 to 1)



After Applying **Dropout**

At test time, **integrate out all the models** in the ensemble

**Monte Carlo approximation:** many forward passes with different masks and average all predictions

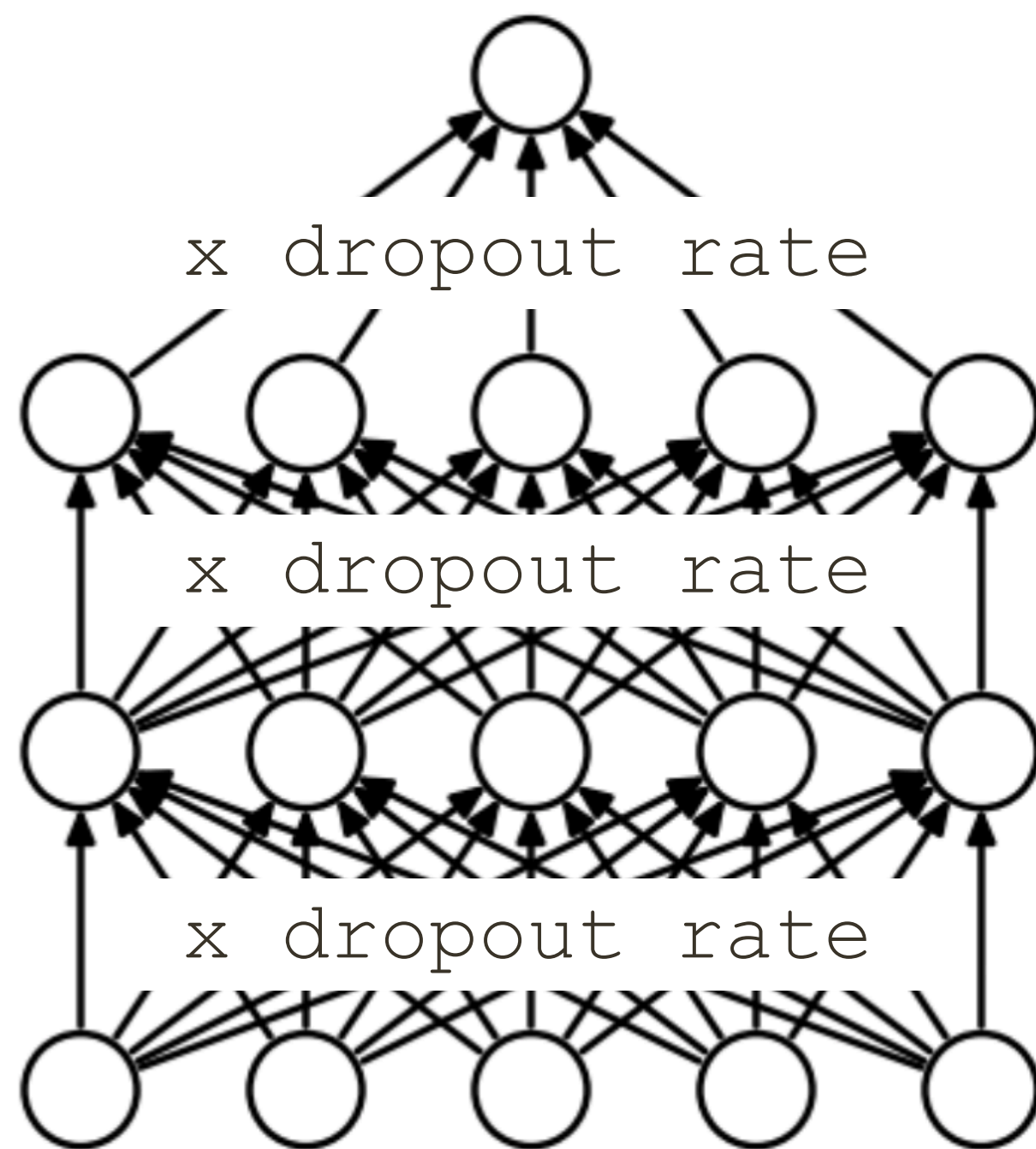
[ Srivastava et al, JMLR 2014 ]

\* adopted from slides of **CS231n at Stanford**



# Regularization: Dropout (at test time)

Randomly **set some neurons to zero** in the forward pass, with probability proportional to dropout rate (between 0 to 1)



At test time, **integrate out all the models** in the ensemble

**Monte Carlo approximation:** many forward passes with different masks and average all predictions

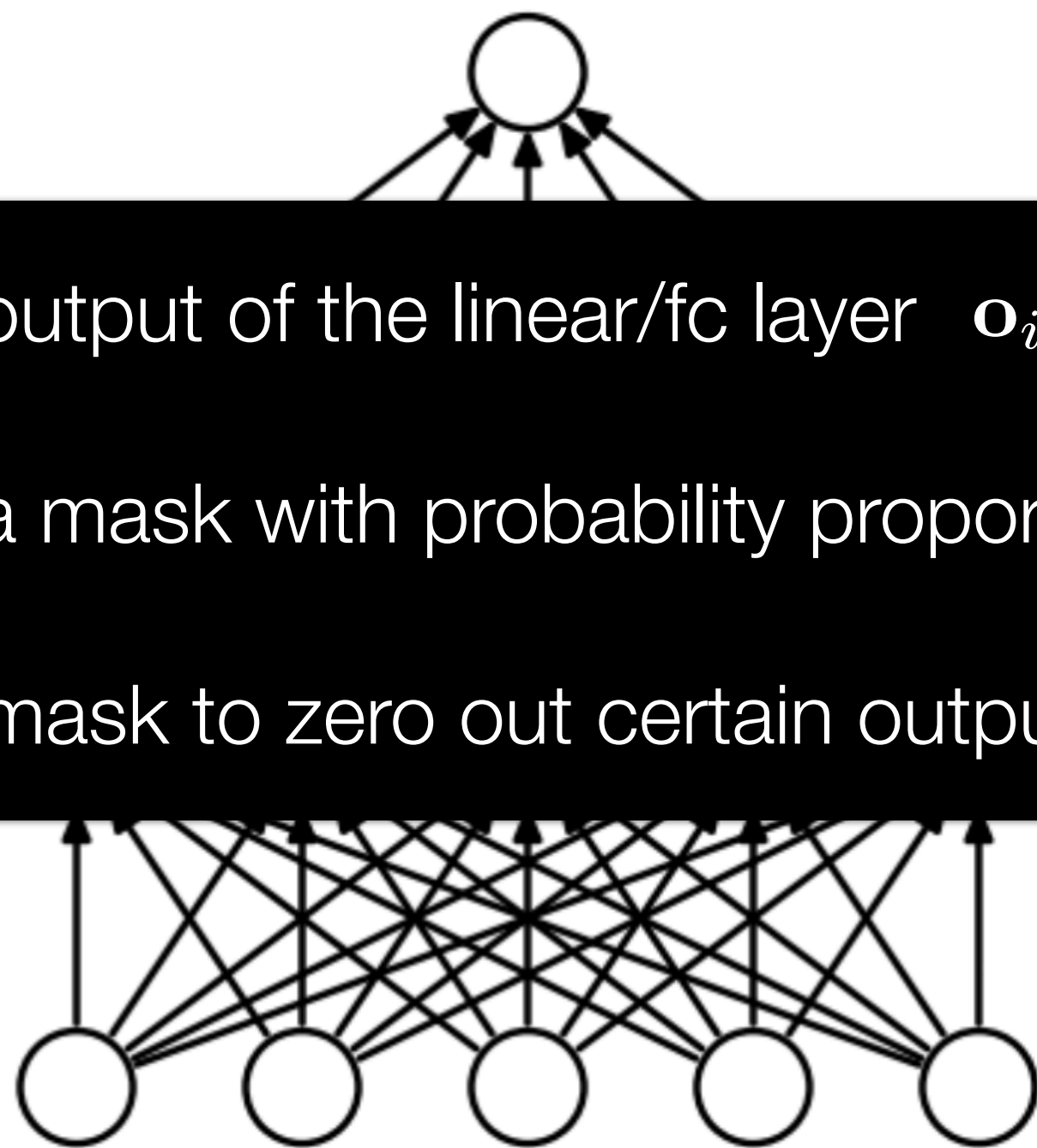
Equivalent to forward pass with all connections on and **scaling of the outputs** by dropout rate



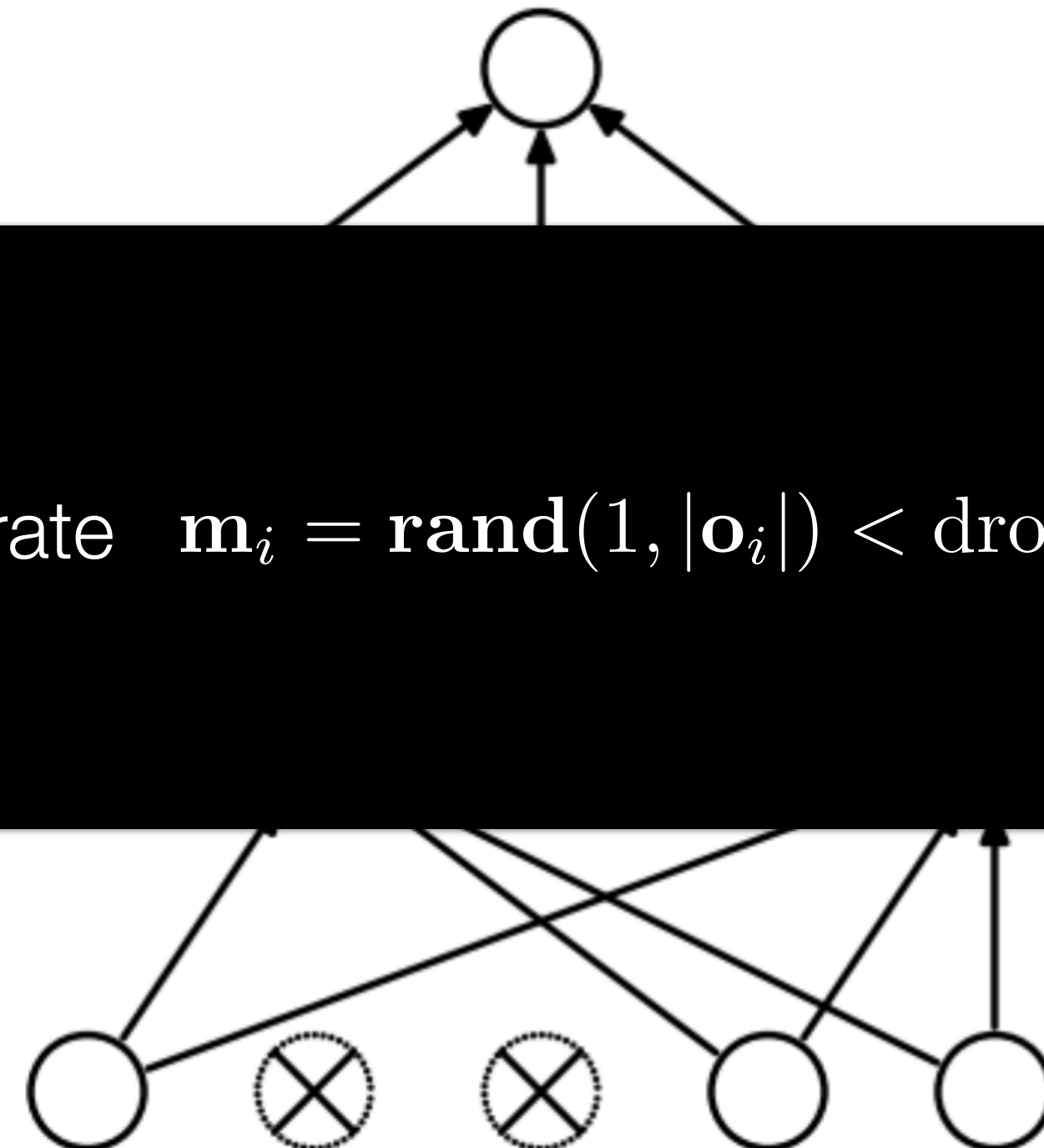
# Regularization: Dropout

Randomly **set some neurons to zero** in the forward pass, with probability proportional to dropout rate (between 0 to 1)

1. Compute output of the linear/fc layer  $\mathbf{o}_i = \mathbf{W}_i \cdot \mathbf{x} + \mathbf{b}_i$
2. Compute a mask with probability proportional to dropout rate  $\mathbf{m}_i = \mathbf{rand}(1, |\mathbf{o}_i|) < \text{dropout rate}$
3. Apply the mask to zero out certain outputs  $\mathbf{o}_i = \mathbf{o}_i \odot \mathbf{m}_i$



Standar Neural Network

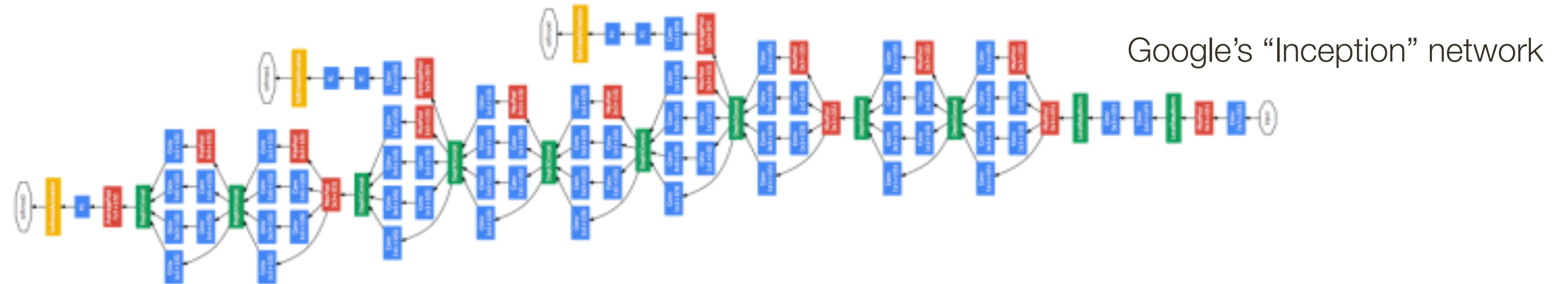


After Applying **Dropout**

[ Srivastava et al, JMLR 2014 ]

\* adopted from slides of **CS231n at Stanford**

# Deep Learning Terminology

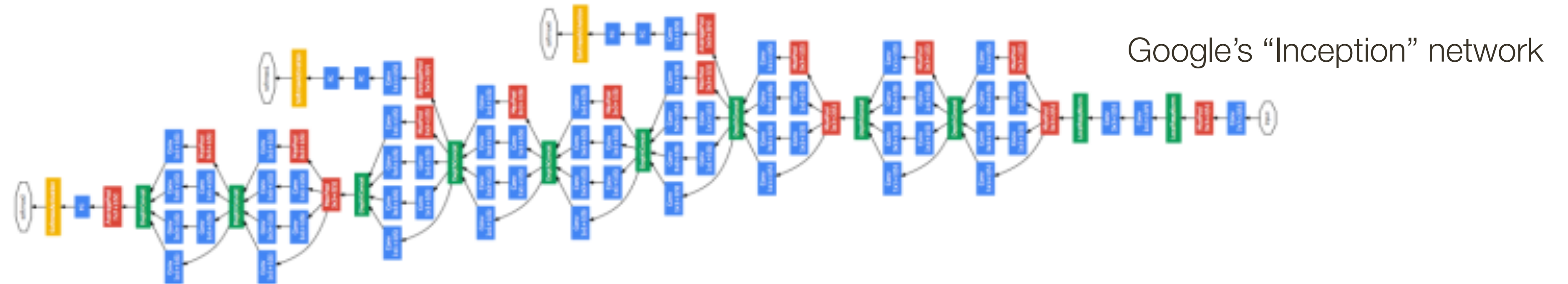


- **Network structure:** number and types of layers, forms of activation functions, dimensionality of each layer and connections (defines computational graph)

generally kept fixed, requires some knowledge of the problem and NN to sensibly set

deeper = better

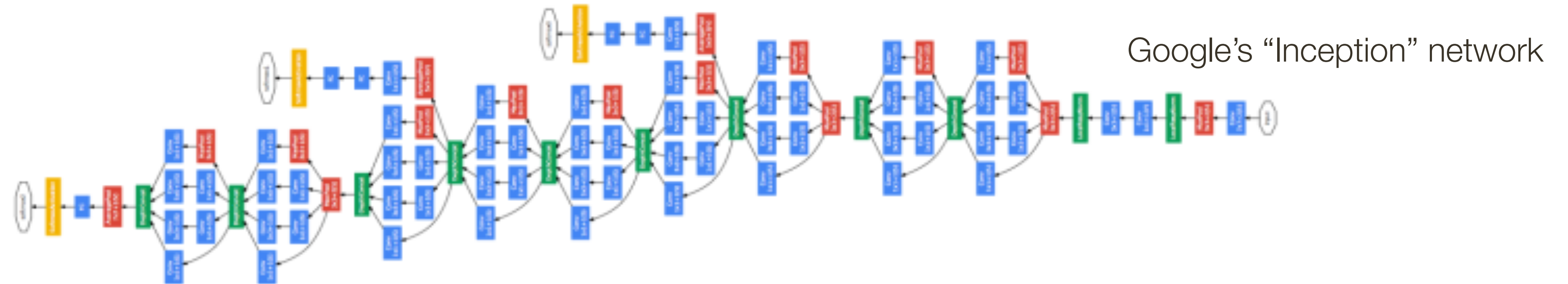
# Deep Learning Terminology



- **Network structure:** number and types of layers, forms of activation functions, dimensionality of each layer and connections (defines computational graph)
  - generally kept fixed, requires some knowledge of the problem and NN to sensibly set
  - deeper = better
- **Loss function:** objective function being optimized (`softmax`, `cross entropy`, *etc.*)
  - requires knowledge of the nature of the problem

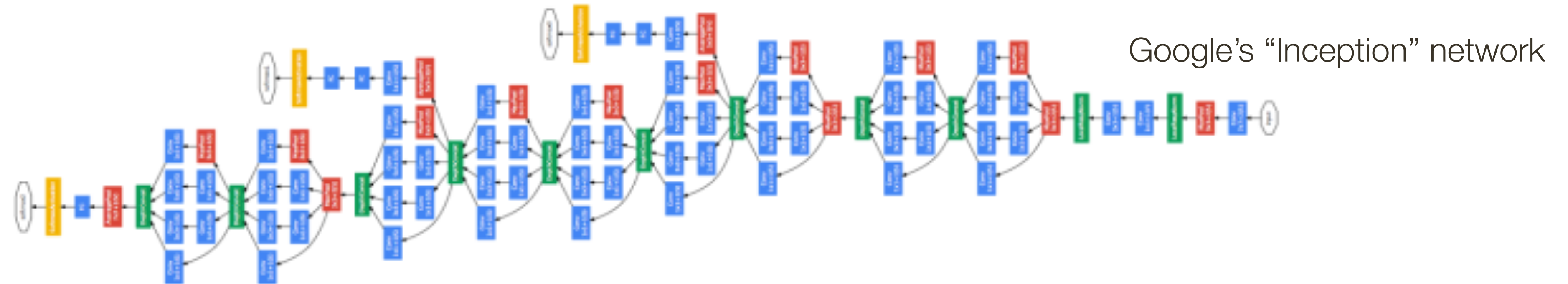


# Deep Learning Terminology



- **Network structure:** number and types of layers, forms of activation functions, dimensionality of each layer and connections (defines computational graph)
  - generally kept fixed, requires some knowledge of the problem and NN to sensibly set
  - deeper = better
- **Loss function:** objective function being optimized (`softmax`, `cross entropy`, *etc.*)
  - requires knowledge of the nature of the problem
- **Parameters:** trainable parameters of the network, including weights/biases of linear/fc layers, parameters of the activation functions, *etc.*
  - optimized using SGD or variants

# Deep Learning Terminology



- **Network structure:** number and types of layers, forms of activation functions, dimensionality of each layer and connections (defines computational graph)

generally kept fixed, requires some knowledge of the problem and NN to sensibly set

deeper = better

- **Loss function:** objective function being optimized (`softmax`, `cross entropy`, *etc.*)

requires knowledge of the nature of the problem

- **Parameters:** trainable parameters of the network, including weights/biases of linear/fc layers, parameters of the activation functions, *etc.* optimized using SGD or variants

- **Hyper-parameters:** parameters, including for optimization, that are not optimized directly as part of training (e.g., `learning rate`, `batch size`, `drop-out rate`) grid search

# Loss Functions ...

This is where all the **fun** is ... we will only look at most common ones

# Multivariate **Regression**

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** output vector  $\mathbf{y} \in \mathbb{R}^m$



# Multivariate **Regression**

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** output vector  $\mathbf{y} \in \mathbb{R}^m$

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}^k$

with **sigmoid** activations:  $\mathbf{0} \leq f(\mathbf{x}; \Theta) \leq \mathbf{1}$

with **Tanh** activations:  $-\mathbf{1} \leq f(\mathbf{x}; \Theta) \leq \mathbf{1}$

with **ReLU** activations:  $\mathbf{0} \leq f(\mathbf{x}; \Theta)$

# Multivariate **Regression**

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** output vector  $\mathbf{y} \in \mathbb{R}^m$

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}^k$

with **sigmoid** activations:  $\mathbf{0} \leq f(\mathbf{x}; \Theta) \leq \mathbf{1}$

with **Tanh** activations:  $-\mathbf{1} \leq f(\mathbf{x}; \Theta) \leq \mathbf{1}$

with **ReLU** activations:  $\mathbf{0} \leq f(\mathbf{x}; \Theta)$

**Neural Network** (output): linear layer

$$\hat{\mathbf{y}} = g(\mathbf{x}; \mathbf{W}, \mathbf{b}) = \mathbf{W} f(\mathbf{x}; \Theta) + \mathbf{b} : \mathbb{R}^k \rightarrow \mathbb{R}^m$$

# Multivariate **Regression**

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** output vector  $\mathbf{y} \in \mathbb{R}^m$

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}^k$

with **sigmoid** activations:  $\mathbf{0} \leq f(\mathbf{x}; \Theta) \leq \mathbf{1}$

with **Tanh** activations:  $-\mathbf{1} \leq f(\mathbf{x}; \Theta) \leq \mathbf{1}$

with **ReLU** activations:  $\mathbf{0} \leq f(\mathbf{x}; \Theta)$

**Neural Network** (output): linear layer

$$\hat{\mathbf{y}} = g(\mathbf{x}; \mathbf{W}, \mathbf{b}) = \mathbf{W} f(\mathbf{x}; \Theta) + \mathbf{b} : \mathbb{R}^k \rightarrow \mathbb{R}^m$$

**Loss:**

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = ||\mathbf{y} - \hat{\mathbf{y}}||^2$$

# Binary Classification (Bernoulli)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** binary label  $y \in \{0, 1\}$

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}$

with **sigmoid** activations:  $0 \leq f(\mathbf{x}; \Theta) \leq 1$

# Binary Classification (Bernoulli)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** binary label  $y \in \{0, 1\}$

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}$

with **sigmoid** activations:  $0 \leq f(\mathbf{x}; \Theta) \leq 1$

**Neural Network** (output): threshold hidden output (which is a sigmoid)

$$\hat{y} = 1[f(\mathbf{x}; \Theta) > 0.5]$$

# Binary Classification (Bernoulli)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** binary label  $y \in \{0, 1\}$

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}$

with **sigmoid** activations:  $0 \leq f(\mathbf{x}; \Theta) \leq 1$

**Neural Network** (output): threshold hidden output (which is a sigmoid)

$$\hat{y} = 1[f(\mathbf{x}; \Theta) > 0.5]$$

**Problem:** Not differentiable, probabilistic interpretation maybe desirable

# Binary Classification (Bernoulli)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** binary label  $y \in \{0, 1\}$

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}$

with **sigmoid** activations:  $0 \leq f(\mathbf{x}; \Theta) \leq 1$

**Neural Network** (output): interpret sigmoid output as probability

$$p(y = 1) = f(\mathbf{x}; \Theta)$$

can interpret the score as the log-odds of  $y = 1$  (a.k.a. the **logits**)



# Binary Classification (Bernoulli)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** binary label  $y \in \{0, 1\}$

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}$

with **sigmoid** activations:  $0 \leq f(\mathbf{x}; \Theta) \leq 1$

**Neural Network** (output): interpret sigmoid output as probability

$$p(y = 1) = f(\mathbf{x}; \Theta)$$

can interpret the score as the log-odds of  $y = 1$  (a.k.a. the **logits**)

**Loss:** similarity between two distributions

# Binary Classification (Bernoulli)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** binary label  $y \in \{0, 1\}$

We can measure similarity between distribution and using cross-entropy

$$H(p, q) = -\mathbb{E}_{x \sim p}[\log q(x)]$$

For discrete distributions this ends up being:

$$H(p, q) = -\sum_x p(x) \log q(x)$$

**Loss:** similarity between two distributions

# Binary Classification (Bernoulli)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** binary label  $y \in \{0, 1\}$

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}$

with **sigmoid** activations:  $0 \leq f(\mathbf{x}; \Theta) \leq 1$

**Neural Network** (output): interpret sigmoid output as probability

$$p(y = 1) = f(\mathbf{x}; \Theta)$$

can interpret the score as the log-odds of  $y = 1$  (a.k.a. the **logits**)

**Loss:**  $\mathcal{L}(y, \hat{y}) = -y \log[f(\mathbf{x}; \Theta)] - (1 - y) \log[1 - f(\mathbf{x}; \Theta)]$

# Binary Classification (Bernoulli)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** binary label  $y \in \{0, 1\}$

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}$

with **sigmoid** activations:  $0 \leq f(\mathbf{x}; \Theta) \leq 1$

**Neural Network** (output): interpret sigmoid output as probability

$$p(y = 1) = f(\mathbf{x}; \Theta)$$

can interpret the score as the log-odds of  $y = 1$  (a.k.a. the **logits**)

**Loss:**

$$\mathcal{L}(y, \hat{y}) = \begin{cases} -\log[1 - f(\mathbf{x}; \Theta)] & y = 0 \\ -\log[f(\mathbf{x}; \Theta)] & y = 1 \end{cases}$$

# Binary Classification (Bernoulli)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** binary label  $y \in \{0, 1\}$

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}$

with **sigmoid** activations:  $0 \leq f(\mathbf{x}; \Theta) \leq 1$

**Neural Network** (output): interpret sigmoid output as probability

$$p(y = 1) = f(\mathbf{x}; \Theta)$$

Minimizing this **loss** is the same as maximizing **log likelihood** of data

**Loss:**

$$\mathcal{L}(y, \hat{y}) = \begin{cases} -\log[1 - f(\mathbf{x}; \Theta)] & y = 0 \\ -\log[f(\mathbf{x}; \Theta)] & y = 1 \end{cases}$$

# Binary Classification (Bernoulli)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** binary label  $y \in \{0, 1\}$

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}^k$

with **ReLU** activations:  $\mathbf{0} \leq f(\mathbf{x}; \Theta)$

# Binary Classification (Bernoulli)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** binary label  $y \in \{0, 1\}$

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}^k$

with **ReLU** activations:  $\mathbf{0} \leq f(\mathbf{x}; \Theta)$

**Neural Network** (output): linear layer with one neuron and sigmoid activation



# Multiclass Classification (e.g, ImageNet)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** muticlass label  $\mathbf{y} \in \{0, 1\}^m$   
(**one-hot** encoding)

# Multiclass Classification (e.g, ImageNet)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** muticlass label  $\mathbf{y} \in \{0, 1\}^m$   
(**one-hot** encoding)

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}^m$

with **ReLU** activations:  $\mathbf{0} \leq f(\mathbf{x}; \Theta)$

# Multiclass Classification (e.g, ImageNet)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** muticlass label  $\mathbf{y} \in \{0, 1\}^m$   
(**one-hot** encoding)

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}^m$

with **ReLU** activations:  $\mathbf{0} \leq f(\mathbf{x}; \Theta)$

**Neural Network** (output): **softmax** function, where probability of class k is:

$$p(\mathbf{y}_k = 1) = \frac{\exp [f(\mathbf{x}; \Theta)_i]}{\sum_{j=1}^C \exp [f(\mathbf{x}; \Theta)_j]}$$

# Multiclass Classification (e.g, ImageNet)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** muticlass label  $\mathbf{y} \in \{0, 1\}^m$   
(**one-hot** encoding)

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}^m$

with **ReLU** activations:  $\mathbf{0} \leq f(\mathbf{x}; \Theta)$

**Neural Network** (output): **softmax** function, where probability of class k is:

$$p(\mathbf{y}_k = 1) = \frac{\exp[f(\mathbf{x}; \Theta)_i]}{\sum_{j=1}^C \exp[f(\mathbf{x}; \Theta)_j]}$$

**Loss:**  $\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_i y_i \log \hat{y}_i$

# Multiclass Classification (e.g, ImageNet)

**Input:** feature vector  $\mathbf{x} \in \mathbb{R}^n$

**Output:** muticlass label  $\mathbf{y} \in \{0, 1\}^m$   
(**one-hot** encoding)

**Neural Network** (input + intermediate hidden layers)  $f(\mathbf{x}; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}^m$

with **ReLU** activations:  $\mathbf{0} \leq f(\mathbf{x}; \Theta)$

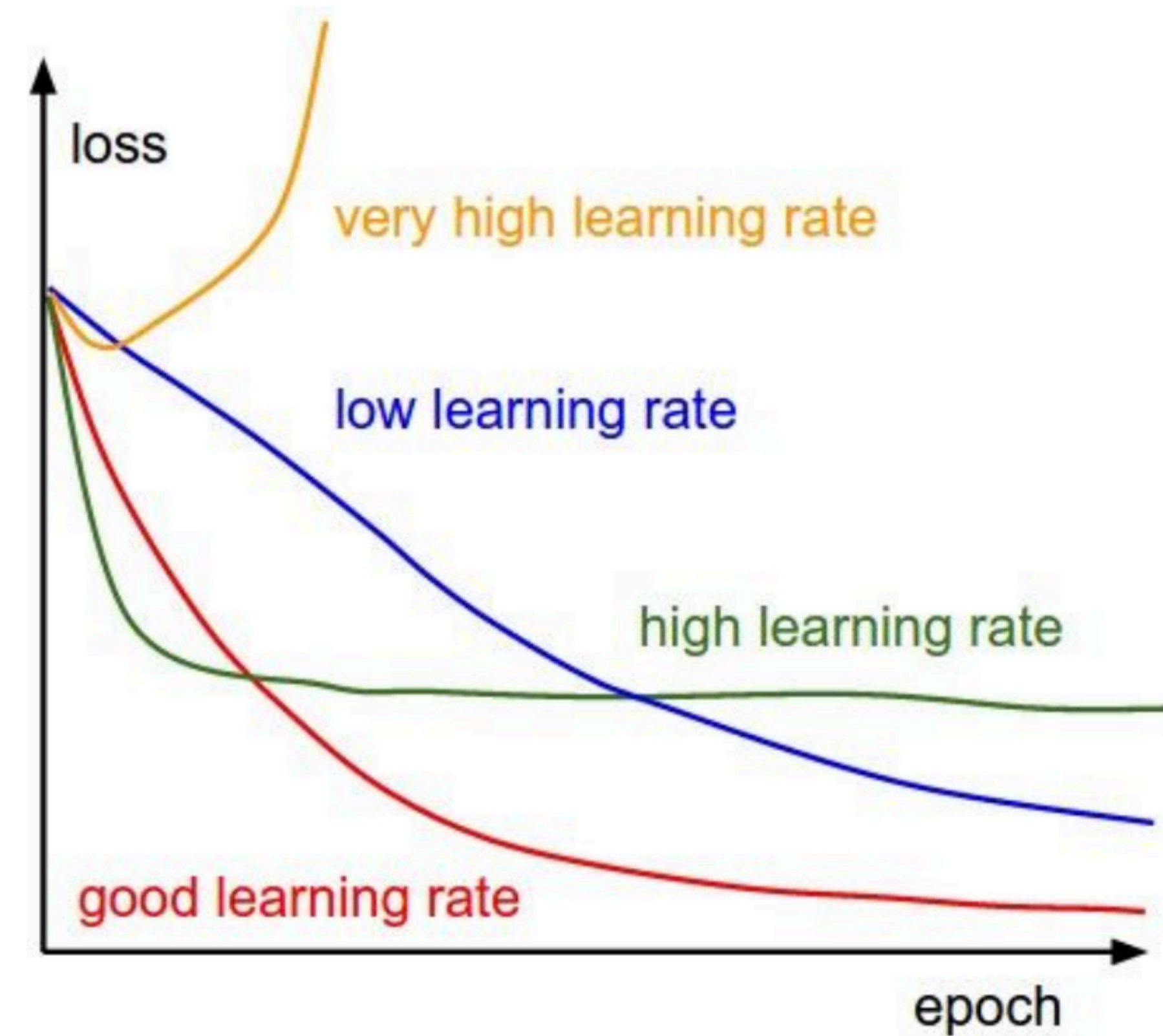
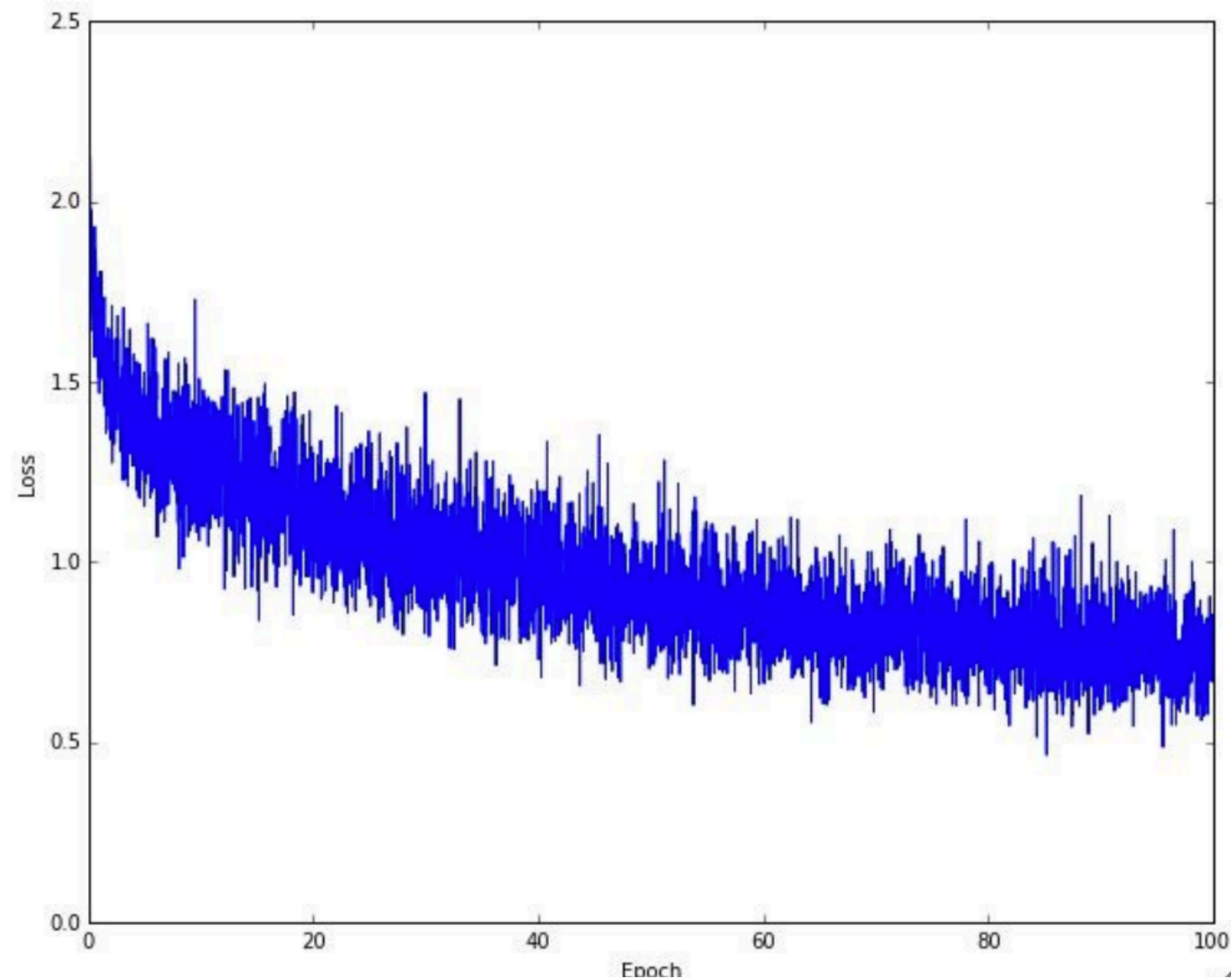
**Neural Network** (output): **softmax** function, where probability of class k is:

$$p(\mathbf{y}_k = 1) = \frac{\exp[f(\mathbf{x}; \Theta)_i]}{\sum_{j=1}^C \exp[f(\mathbf{x}; \Theta)_j]}$$

**Loss:**  $\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_i y_i \log \hat{y}_i = - \log \hat{y}_i$

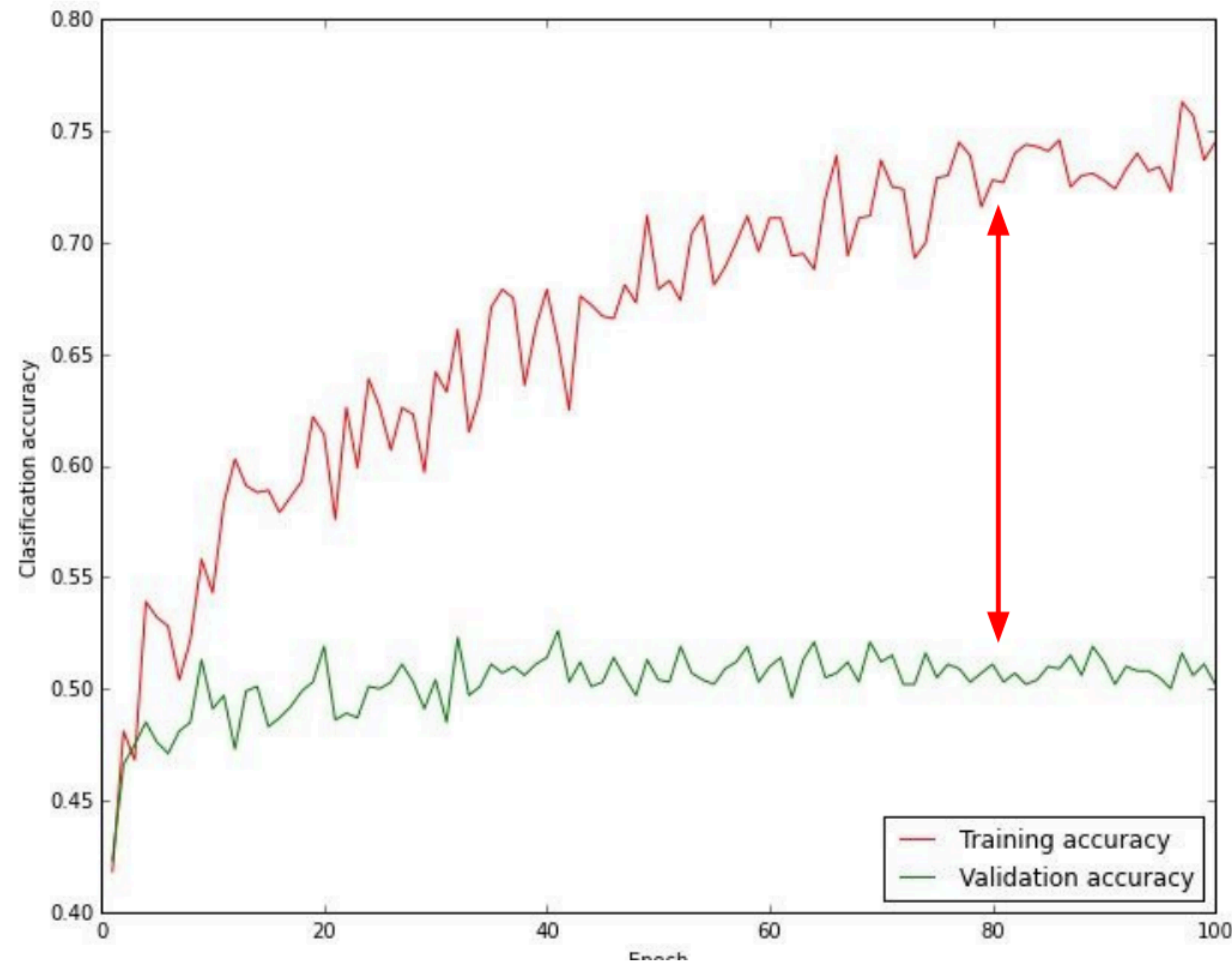
Special case for multi-class single label

# Monitoring Learning: Visualizing the (training) loss





# Monitoring Learning: Visualizing the (training) loss



Big gap = overfitting

**Solution:** increase regularization

No gap = undercutting

**Solution:** increase model capacity

Small gap = **ideal**