



# Machine Learning

A. M. Sadeghzadeh, Ph.D.

Sharif University of Technology  
Computer Engineering Department (CE)  
Trustworthy and Secure AI Lab (TSAIL)



October 24, 2025

Most slides have been adapted from Bhiksha Raj, 11-785, CMU 2020, Fei Fei Li, cs231n, Stanford 2017, and Soleymani, CE40717, Sharif 2020.

## Schedule

- 1 Machine Learning
  - 2 Supervised Learning
  - 3 Logistic Regression
  - 4 Loss Function
  - 5 Optimization
  - 6 Regularization
  - 7 Softmax Classifier
  - 8 Deep Neural Networks
  - 9 Vanishing/Exploding Gradient

# Machine Learning

## Definition of machine learning

Tom Mitchell (1998)

A computer program is said to **learn** from **experience E** with respect to some class of **tasks T** and **performance measure P**, if its performance at tasks in **T**, as measured by **P**, improves with experience **E**.

## Example

**Task:** Classifying emails as spam or not spam

**Experience:** Watching label of emails as spam or not spam

**Performance metric:** The number (or fraction) of emails correctly classified as spam



# Type of Learning

## ■ Supervised learning

Predicting a target variable for which we get to see examples

# Type of Learning

- **Supervised learning**

Predicting a target variable for which we get to see examples

- **Unsupervised learning**

Revealing structure in the observed data

# Type of Learning

## ■ Supervised learning

Predicting a target variable for which we get to see examples

## ■ Unsupervised learning

Revealing structure in the observed data

## ■ Reinforcement learning

Partial (indirect) feedback, no explicit guidance

Given rewards for a sequence of moves to learn a policy and utility functions

## Supervised Learning

# Components of supervised learning (function approximation)

## Concept

- Unknown target function  $f : \mathcal{X} \rightarrow \mathcal{Y}$
- $\mathcal{X}$ : Input space
- $\mathcal{Y}$ : Output space

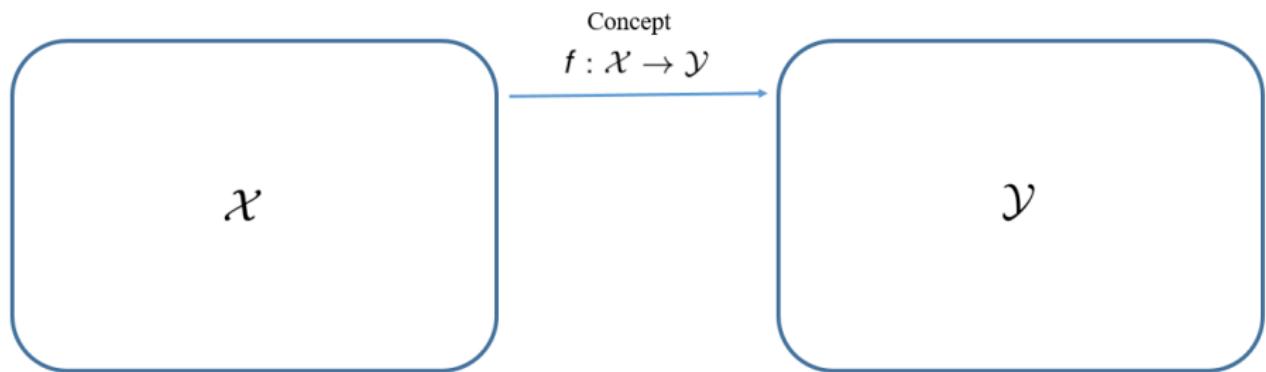
# Components of supervised learning (function approximation)

 $\mathcal{X}$

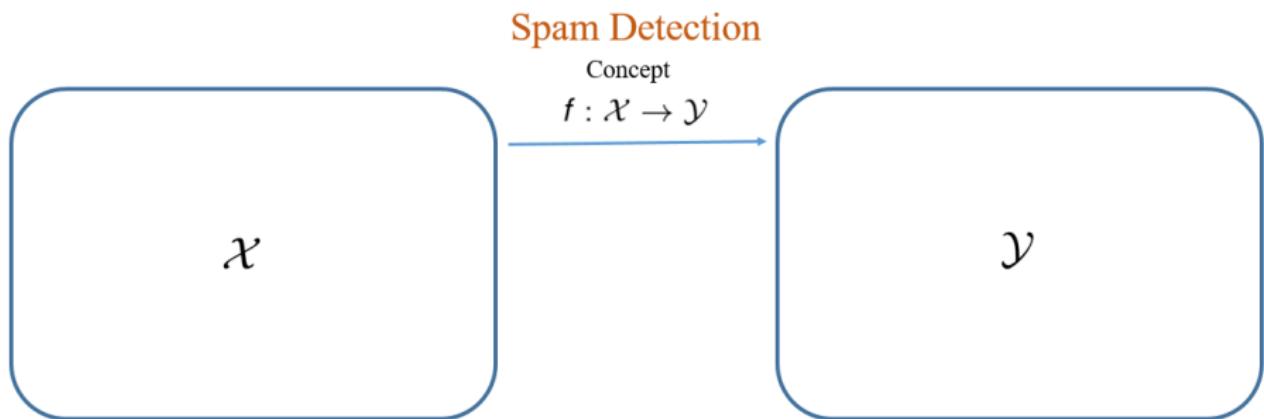
# Components of supervised learning (function approximation)

 $\mathcal{X}$  $\mathcal{Y}$

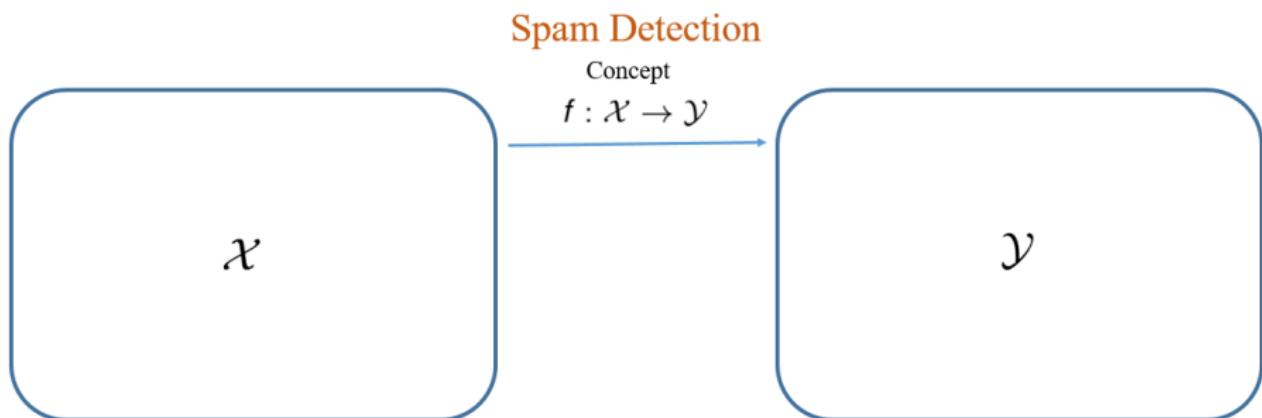
# Components of supervised learning (function approximation)



# Components of supervised learning (function approximation)

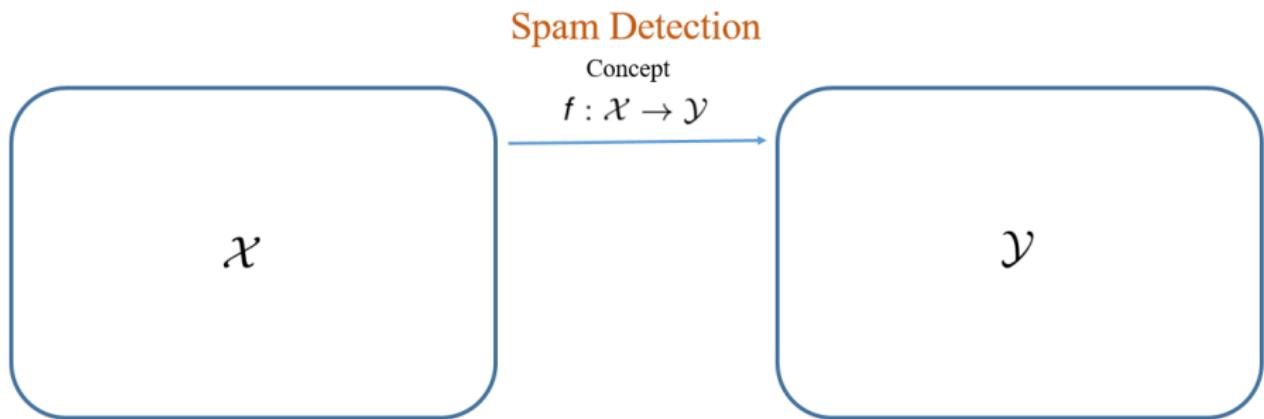


# Components of supervised learning (function approximation)



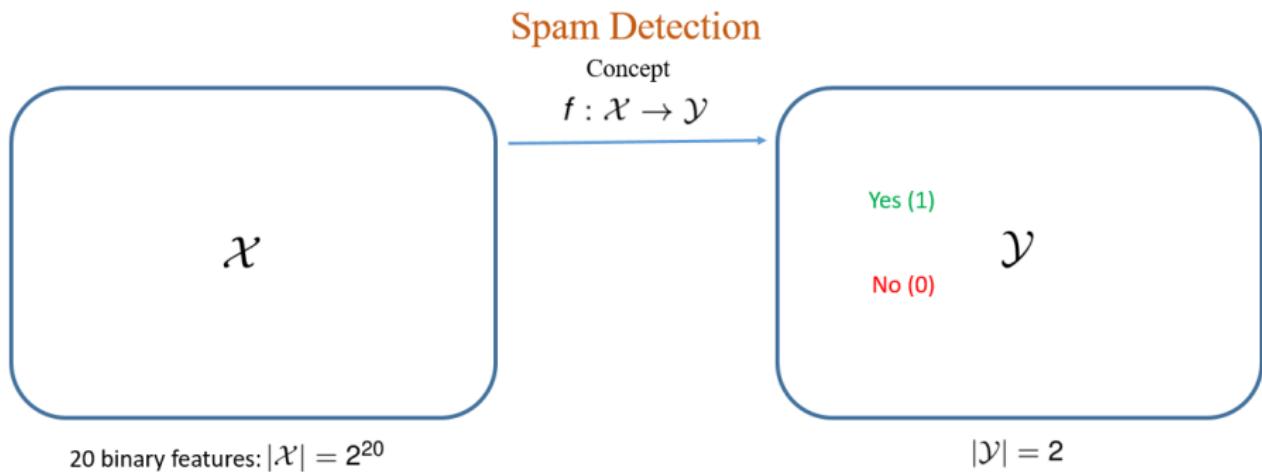
20 binary features:

# Components of supervised learning (function approximation)

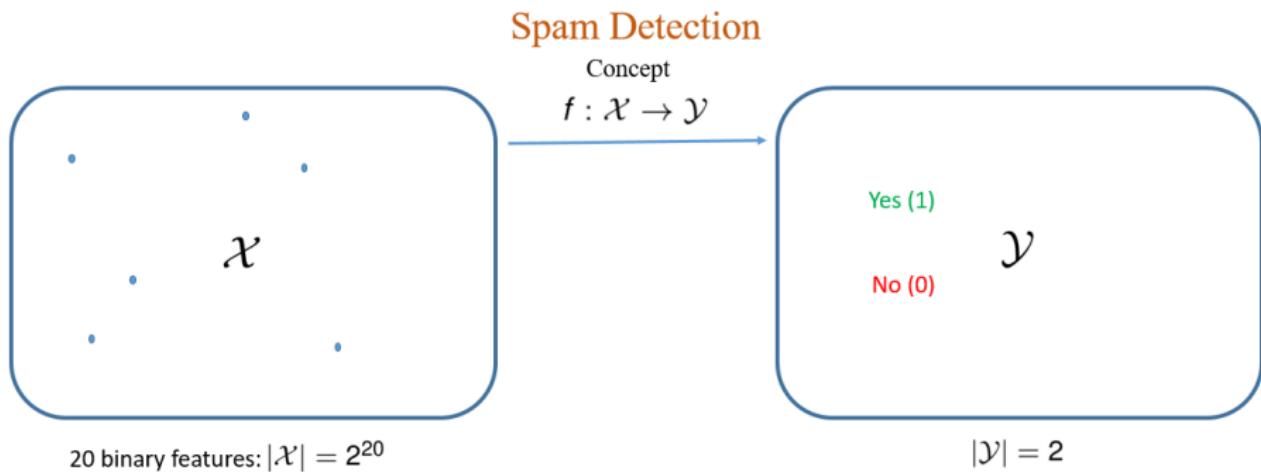


20 binary features:  $|\mathcal{X}| = 2^{20}$

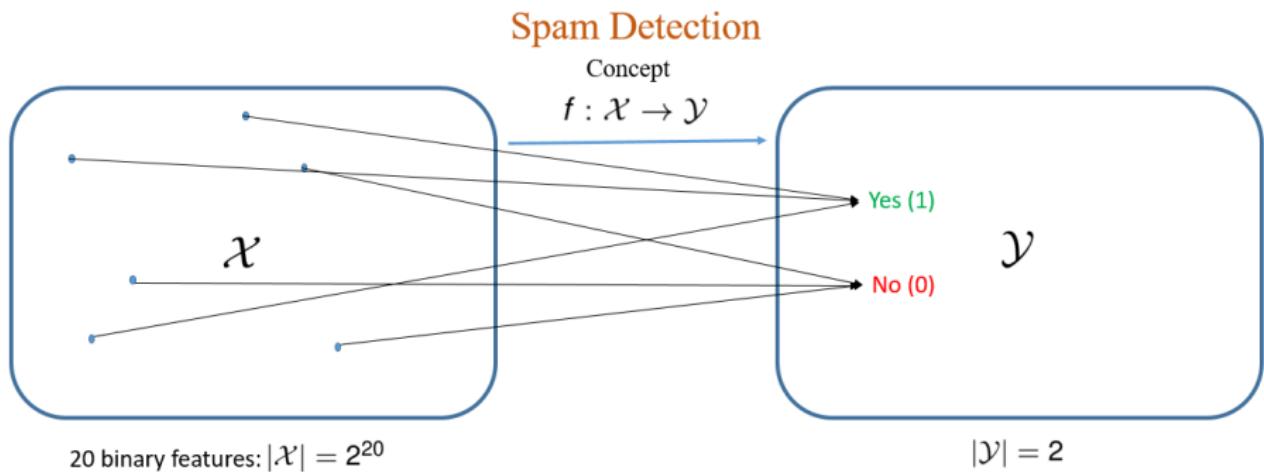
# Components of supervised learning (function approximation)



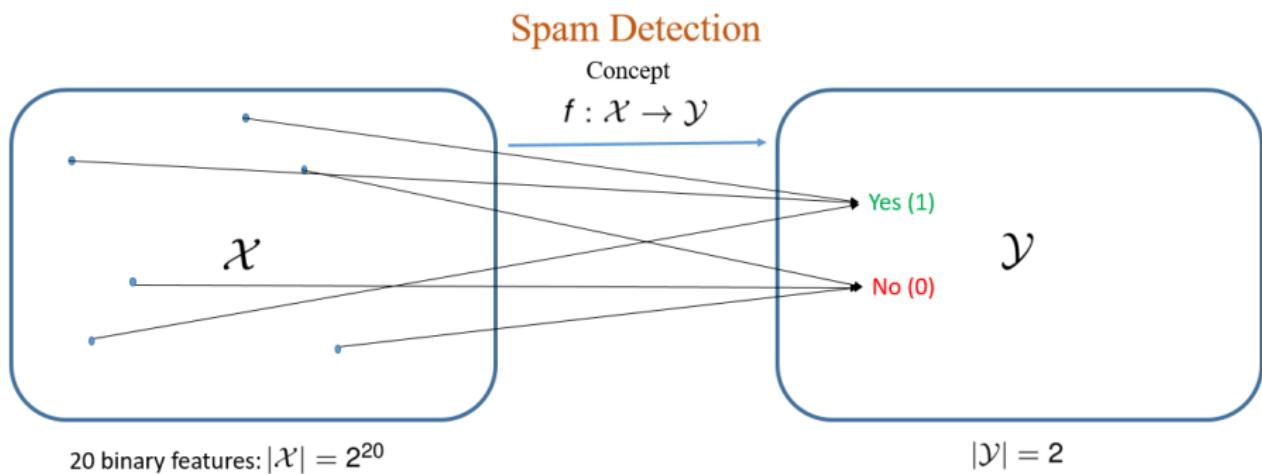
# Components of supervised learning (function approximation)



# Components of supervised learning (function approximation)

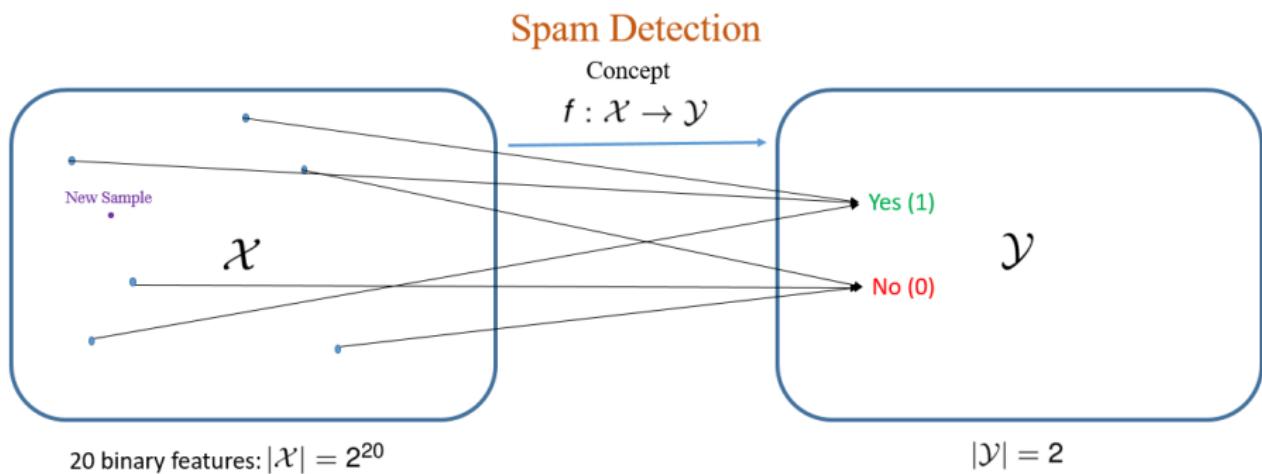


# Components of supervised learning (function approximation)



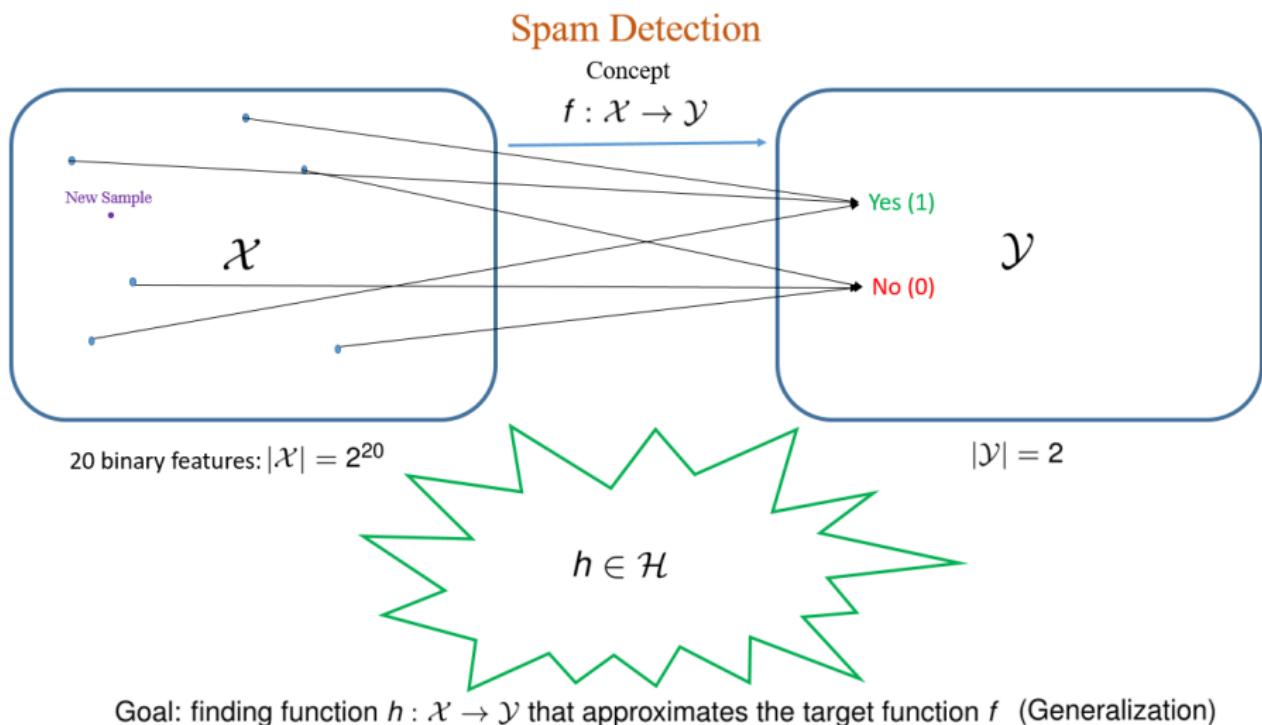
Goal: finding function  $h : \mathcal{X} \rightarrow \mathcal{Y}$  that approximates the target function  $f$

# Components of supervised learning (function approximation)

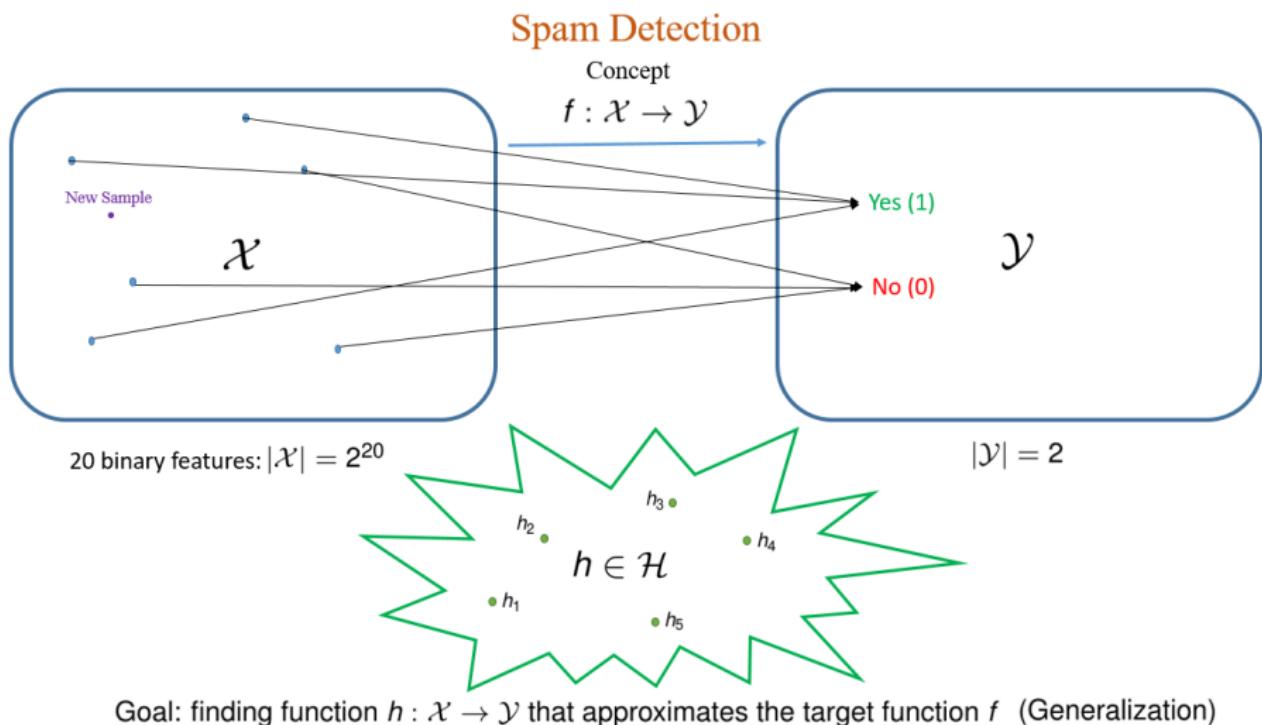


Goal: finding function  $h : \mathcal{X} \rightarrow \mathcal{Y}$  that approximates the target function  $f$  (Generalization)

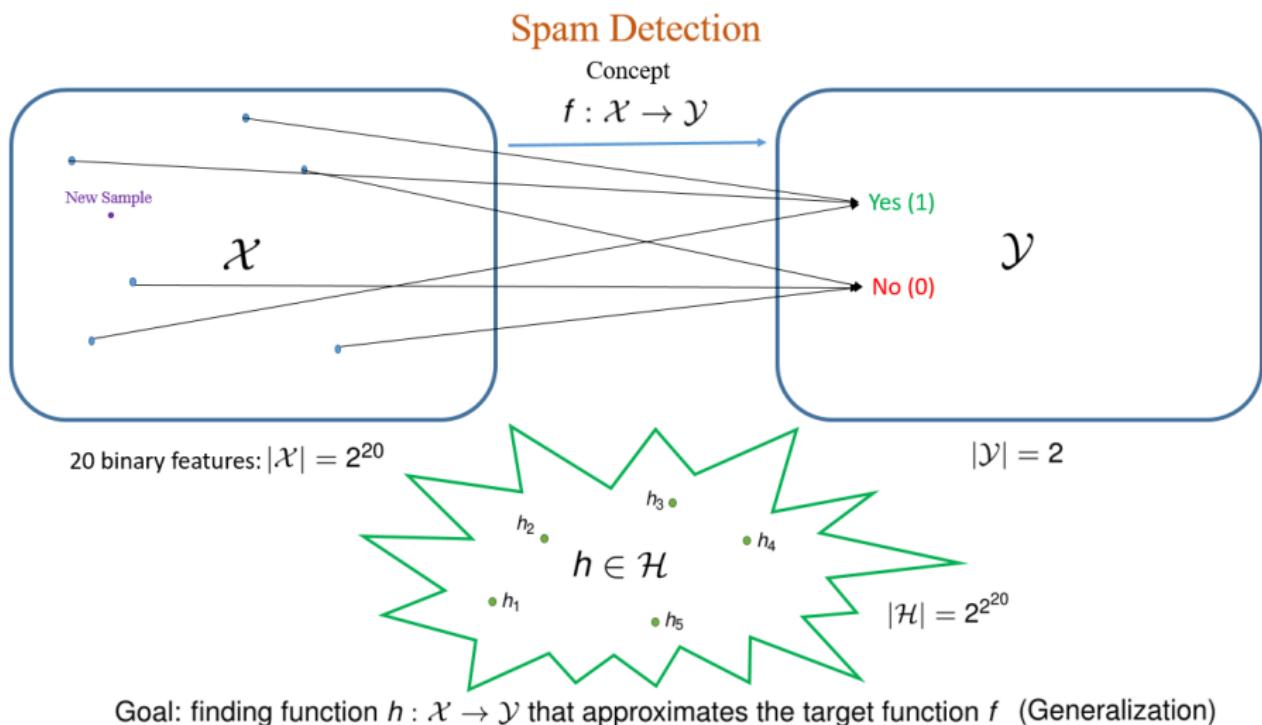
# Components of supervised learning (function approximation)



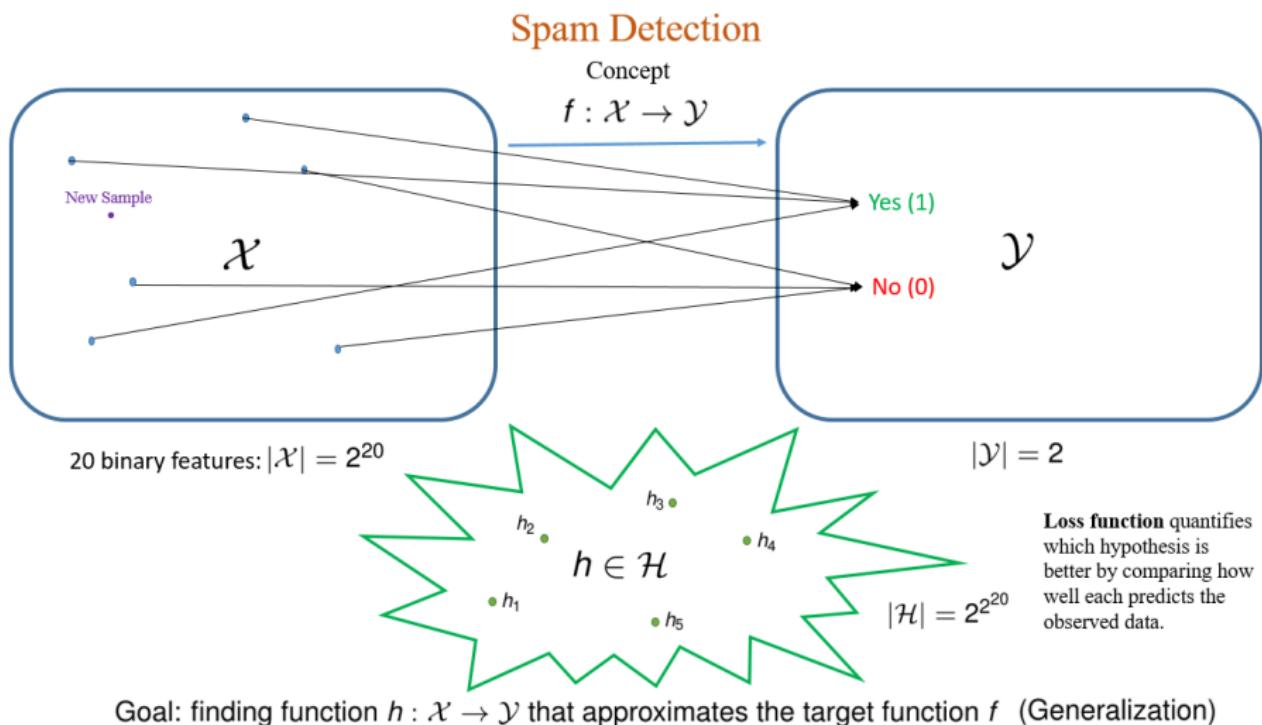
# Components of supervised learning (function approximation)



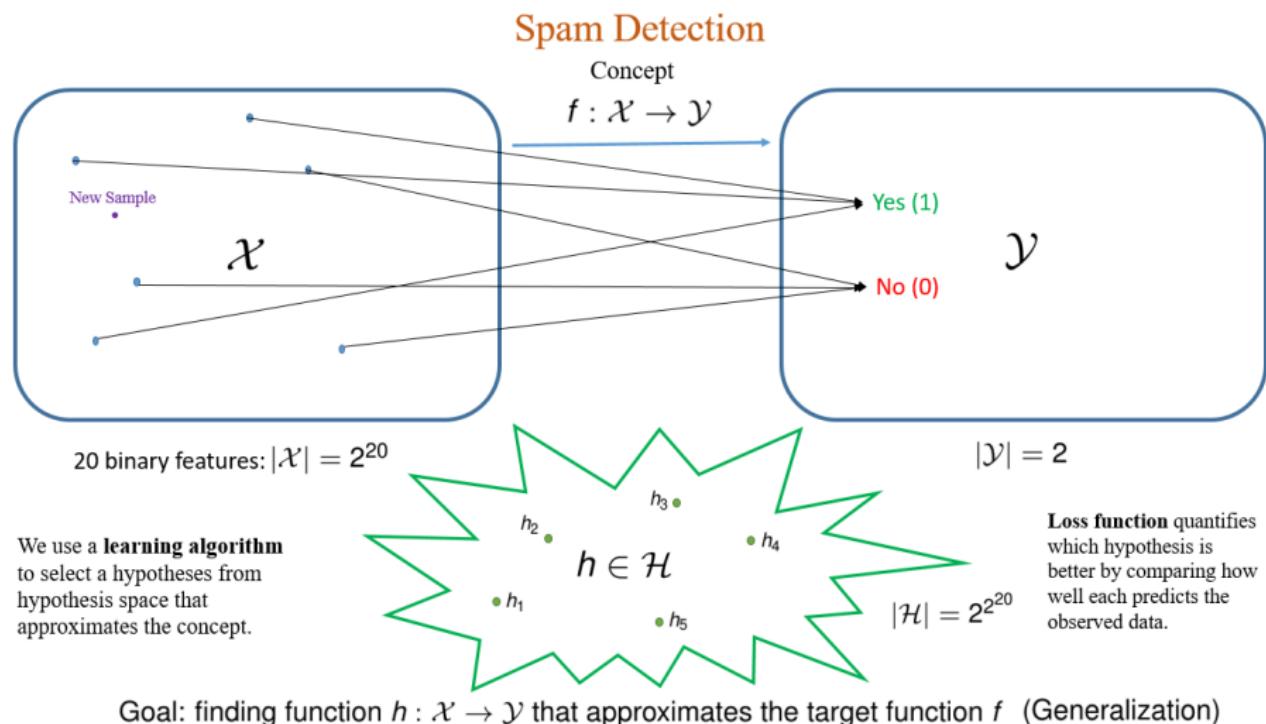
# Components of supervised learning (function approximation)



# Components of supervised learning (function approximation)



# Components of supervised learning (function approximation)



# Components of supervised learning (function approximation)

## Concept

- Unknown target function  $f : \mathcal{X} \rightarrow \mathcal{Y}$
- $\mathcal{X}$ : Input space
- $\mathcal{Y}$ : Output space

## Training data

- Samples from the concept  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$

# Components of supervised learning (function approximation)

## Concept

- Unknown target function  $f : \mathcal{X} \rightarrow \mathcal{Y}$
- $\mathcal{X}$ : Input space
- $\mathcal{Y}$ : Output space

## Training data

- Samples from the concept  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$

## Hypotheses space

- Pick a hypotheses  $h : \mathcal{X} \rightarrow \mathcal{Y}$  from hypotheses space  $h \in \mathcal{H}$  (e.g., linear functions) that approximates the target function  $f$

# Components of supervised learning (function approximation)

## Concept

- Unknown target function  $f : \mathcal{X} \rightarrow \mathcal{Y}$
- $\mathcal{X}$ : Input space
- $\mathcal{Y}$ : Output space

## Training data

- Samples from the concept  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$

## Hypotheses space

- Pick a hypotheses  $h : \mathcal{X} \rightarrow \mathcal{Y}$  from hypotheses space  $h \in \mathcal{H}$  (e.g., linear functions) that approximates the target function  $f$

## Learning algorithm

- We use a learning algorithm to select a hypotheses from hypothesis space that approximates the concept

# Components of supervised learning

Learning model composed of

- A hypothesis set
- A learning algorithm

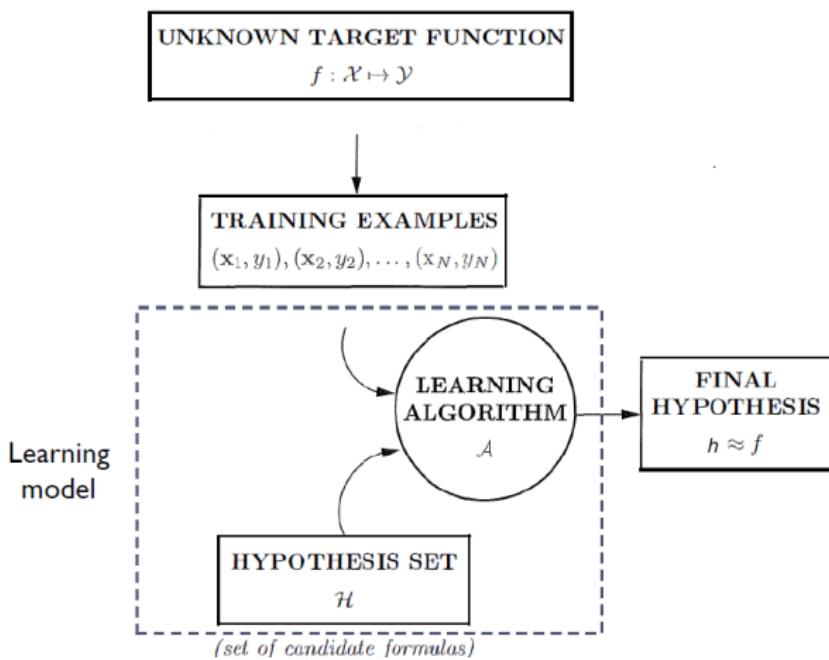


Figure: <https://work.caltech.edu/telecourse.html>

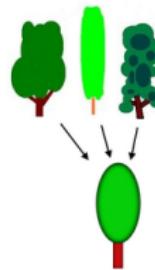
# Generalization

## Hypotheses evaluation

How well  $h$  generalizes to unseen examples.

## Test data

Samples from the concept  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$   
Do not exist in the training data



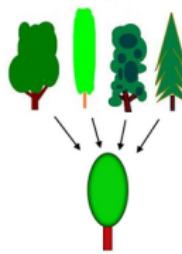
# Generalization

## Hypotheses evaluation

How well  $h$  generalizes to unseen examples.

### Test data

Samples from the concept  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$   
Do not exist in the training data



## Independent and identically distributed (i.i.d.)

A collection of random variables is independent and identically distributed if each random variable has the same probability distribution as the others and all are mutually independent.

- $\forall i, x^{(i)} \sim \mathcal{D}$  (Identically Distributed)
- $\forall i \neq j, \mathcal{P}(x^{(i)}, x^{(j)}) = \mathcal{P}(x^{(i)})\mathcal{P}(x^{(j)})$  (Independently Distributed)

Machine learning algorithms have focused primarily on sets of data points that were assumed to be independent and identically distributed (i.i.d.).

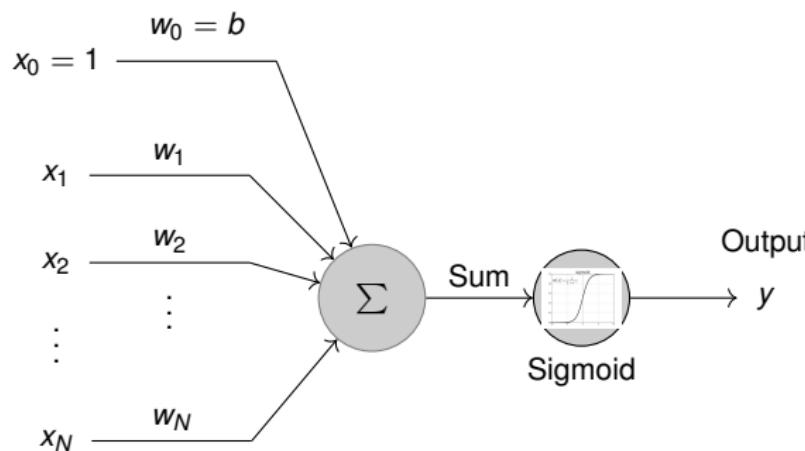
## Logistic Regression

# Logistic regression

It is the perceptron with a sigmoid activation

- It actually computes the probability that the input belongs to class 1

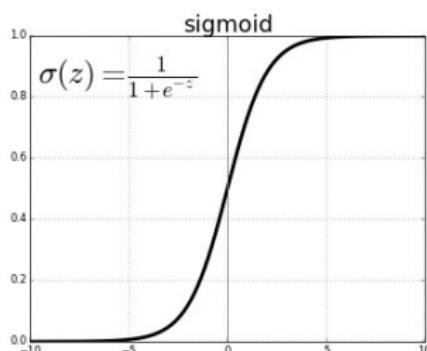
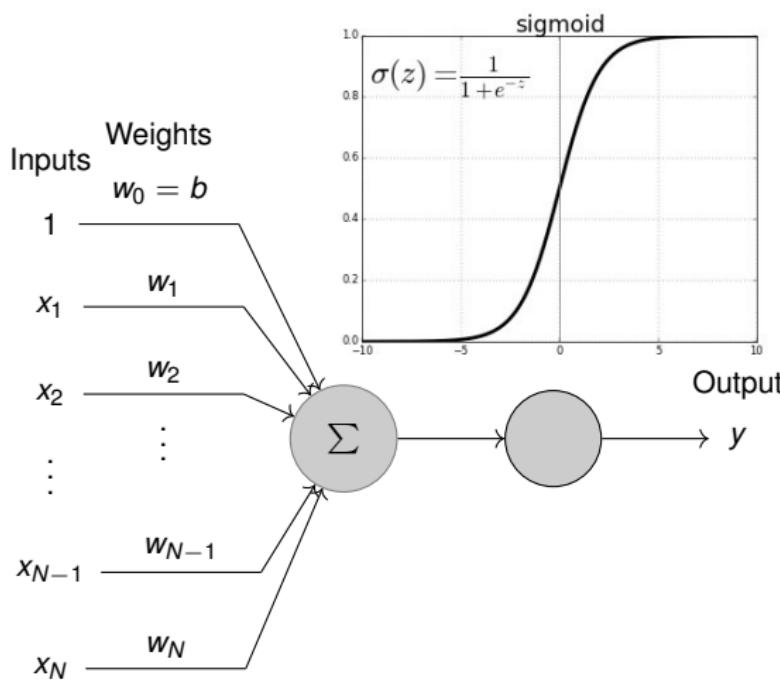
Inputs      Weights



# Logistic regression

It is the perceptron with a sigmoid activation

- It actually computes the probability that the input belongs to class 1



$$z = \sum_i w_i x_i$$

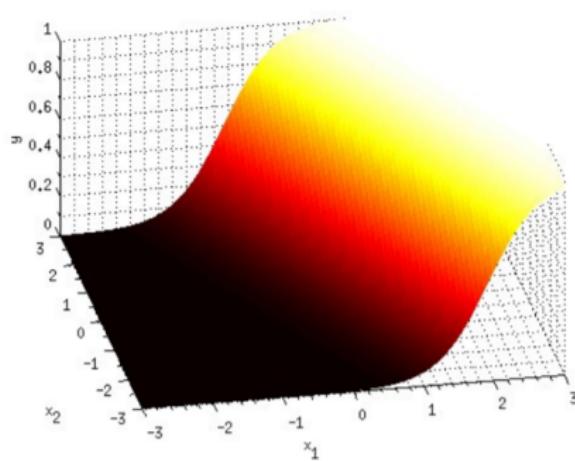
$$p(y = 1|x) = \sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$y = \begin{cases} 1 & \text{if } \sigma(z) \geq \frac{1}{2} \\ 0 & \text{if } \sigma(z) < \frac{1}{2} \end{cases}$$

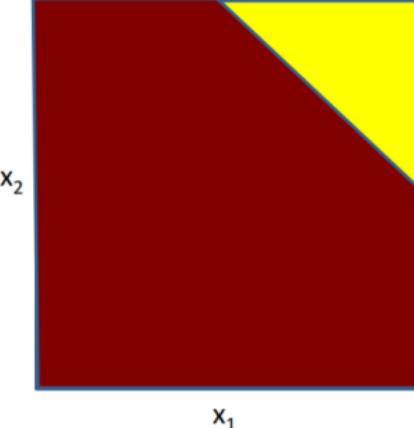
# Logistic regression

This is the perceptron with a sigmoid activation

- It actually computes the probability that the input belongs to class 1



When X is a 2-D variable



$$P(Y = 1|X) = \frac{1}{1 + \exp(-\sum_i w_i x_i - b)}$$

# Vectorization

$$z = \sum_i w_i x_i = \mathbf{w}^T \mathbf{x} = \begin{bmatrix} w_0 = b & w_1 & w_2 & \dots & w_{N-1} & w_N \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_{N-1} \\ x_N \end{bmatrix}$$

$$f(\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

## Loss Function

## Loss function

- Loss function quantifies our unhappiness with the scores across the training data.

# Loss function

- Loss function quantifies our unhappiness with the scores across the training data.
- A loss function tells how good our current classifier is
  - Given a dataset of examples  $\{(x^{(i)}, y^{(i)})\}_{i=1}^N$
  - $x^{(i)}$  is data and  $y^{(i)}$  is binary label.

# Loss function

- Loss function quantifies our unhappiness with the scores across the training data.
- A loss function tells how good our current classifier is
  - Given a dataset of examples  $\{(x^{(i)}, y^{(i)})\}_{i=1}^N$
  - $x^{(i)}$  is data and  $y^{(i)}$  is binary label.
- Loss over the dataset is a average of loss over examples:

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i L_i(f(\mathbf{x}^{(i)}, \mathbf{w}), y^{(i)})$$

# Loss function

- Loss function quantifies our unhappiness with the scores across the training data.
- A loss function tells how good our current classifier is
  - Given a dataset of examples  $\{(x^{(i)}, y^{(i)})\}_{i=1}^N$
  - $x^{(i)}$  is data and  $y^{(i)}$  is binary label.
- Loss over the dataset is a average of loss over examples:

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i L_i(f(\mathbf{x}^{(i)}, \mathbf{w}), y^{(i)})$$

- Empirical loss  $L$  is only an empirical approximation of the true loss

## Loss function

- Loss function quantifies our unhappiness with the scores across the training data.
  - A loss function tells how good our current classifier is
    - Given a dataset of examples  $\{(x^{(i)}, y^{(i)})\}_{i=1}^N$
    - $x^{(i)}$  is data and  $y^{(i)}$  is binary label.
  - Loss over the dataset is a average of loss over examples:

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i L_i(f(\mathbf{x}^{(i)}, \mathbf{w}), y^{(i)})$$

- Empirical loss  $L$  is only an empirical approximation of the true loss
  - **Learning**

$$\boldsymbol{w}^* = \operatorname*{argmin}_{\boldsymbol{w}} L(f(\boldsymbol{X}, \boldsymbol{w}), \boldsymbol{y})$$

## Cross entropy

If we have two separate probability distributions  $P(x)$  and  $Q(x)$  over the same random variable  $x$ , we can measure how different these two distributions are using the Kullback-Leibler (KL) divergence:

$$D_{KL}(P||Q) = \mathbb{E}_{x \sim P} [\log \frac{P(x)}{Q(x)}] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)]$$

$$D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

## Cross entropy

If we have two separate probability distributions  $P(x)$  and  $Q(x)$  over the same random variable  $x$ , we can measure how different these two distributions are using the Kullback-Leibler (KL) divergence:

$$D_{KL}(P||Q) = \mathbb{E}_{x \sim P} [\log \frac{P(x)}{Q(x)}] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)]$$

$$D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

- It is the extra amount of information needed to send a message containing symbols drawn from probability distribution  $P$ , when we use a code that was designed to minimize the length of messages drawn from probability distribution  $Q$ .

## Cross entropy

If we have two separate probability distributions  $P(x)$  and  $Q(x)$  over the same random variable  $x$ , we can measure how different these two distributions are using the Kullback-Leibler (KL) divergence:

$$D_{KL}(P||Q) = \mathbb{E}_{x \sim P}[\log \frac{P(x)}{Q(x)}] = \mathbb{E}_{x \sim P}[\log P(x) - \log Q(x)]$$

$$D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

- It is the extra amount of information needed to send a message containing symbols drawn from probability distribution  $P$ , when we use a code that was designed to minimize the length of messages drawn from probability distribution  $Q$ .
  - KL divergence is non-negative and measures the difference between two distributions, it is often conceptualized as measuring some sort of distance between these distributions. However, it is not a true distance measure because it is not symmetric

## Cross entropy

- ### ■ Minimizing loss function

$$\operatorname{argmin}_{Q(x)} D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

## Cross entropy

- ### ■ Minimizing loss function

$$\operatorname{argmin}_{Q(x)} D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

- $P(x)$  is the true distribution of the random variable  $x$  and it is fixed.

## Cross entropy

- ### ■ Minimizing loss function

$$\operatorname*{argmin}_{Q(x)} D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

- $P(x)$  is the true distribution of the random variable  $x$  and it is fixed.

$$\operatorname{argmin}_{Q(x)} - \sum_i p_i(x) \log q_i(x)$$

## Cross entropy

- ### ■ Minimizing loss function

$$\operatorname*{argmin}_{Q(x)} D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

- $P(x)$  is the true distribution of the random variable  $x$  and it is fixed.

$$\operatorname{argmin}_{Q(x)} - \sum_i p_i(x) \log q_i(x)$$

$$CrossEntropy(P(x), Q(x)) = - \sum_i p_i(x) \log q_i(x)$$

## Logistic regression loss function

## ■ Logistic Regression

$$P(y=1|x^{(i)}; w) = f(x^{(i)}, w)$$

$$P(y = \text{o} | \mathbf{x}^{(i)}; \mathbf{w}) = 1 - f(\mathbf{x}^{(i)}, \mathbf{w})$$

$$y^{(i)} \in \{0, 1\}$$

# Logistic regression loss function

- Logistic Regression

$$P(y = 1 | \mathbf{x}^{(i)}; \mathbf{w}) = f(\mathbf{x}^{(i)}, \mathbf{w})$$

$$P(y = 0 | \mathbf{x}^{(i)}; \mathbf{w}) = 1 - f(\mathbf{x}^{(i)}, \mathbf{w})$$

$$y^{(i)} \in \{0, 1\}$$

- Cross entropy as the loss function

$$L(f(\mathbf{X}, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i L_i(f(\mathbf{x}^{(i)}, \mathbf{w}), y^{(i)})$$

$$L(f(\mathbf{X}, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i -y^{(i)} \log f(\mathbf{x}^{(i)}, \mathbf{w}) - (1 - y^{(i)}) \log(1 - f(\mathbf{x}^{(i)}, \mathbf{w}))$$

# Logistic regression loss function

## ■ Logistic Regression

$$P(y = 1 | \mathbf{x}^{(i)}; \mathbf{w}) = f(\mathbf{x}^{(i)}, \mathbf{w})$$

$$P(y = 0 | \mathbf{x}^{(i)}; \mathbf{w}) = 1 - f(\mathbf{x}^{(i)}, \mathbf{w})$$

$$y^{(i)} \in \{0, 1\}$$

## ■ Cross entropy as the loss function

$$L(f(\mathbf{X}, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i L_i(f(\mathbf{x}^{(i)}, \mathbf{w}), y^{(i)})$$

$$L(f(\mathbf{X}, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i -y^{(i)} \log f(\mathbf{x}^{(i)}, \mathbf{w}) - (1 - y^{(i)}) \log(1 - f(\mathbf{x}^{(i)}, \mathbf{w}))$$

How do we find the best  $w$ ?

## Optimization

## Derivative

- In 1-dimension, the derivative of a function gives the slope:

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

- In multiple dimensions, the gradient is the vector of (partial derivatives) along each dimension
  - The direction of steepest descent is the negative gradient

- To decrease  $f$  at  $x$ :

if  $\frac{df(x)}{dx} < 0 \Rightarrow f(x+h) < f(x) \Rightarrow f$  is decreasing in  $x \Rightarrow$  take a step in **positive** direction

if  $\frac{df(x)}{dx} > 0 \Rightarrow f(x+h) > f(x) \Rightarrow f$  is increasing in  $x \Rightarrow$  take a step in **negative** direction

- Take a step in the reverse direction of derivative at  $x$

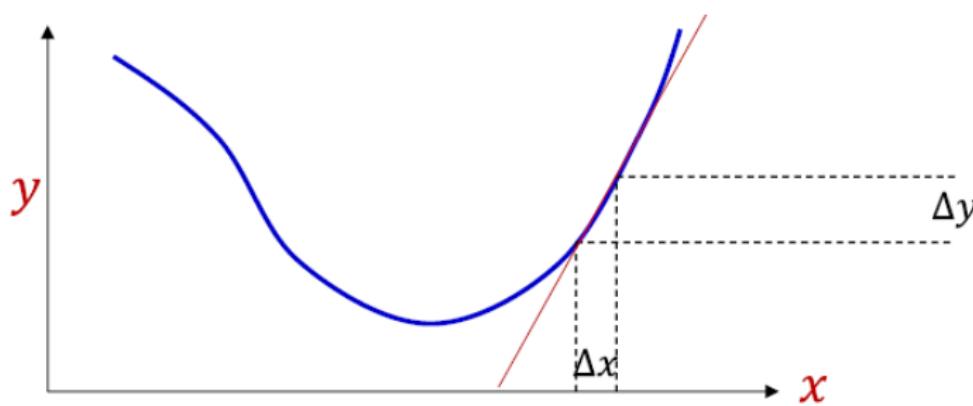
## A brief note on derivatives...

- A derivative of a function at any point tells us how much a minute increment to the argument of the function will increment the value of the function

- For any  $y = f(x)$ , expressed as a multiplier  $\alpha$  to a tiny increment  $\Delta x$  to obtain the increments  $\Delta y$  to the output

$$\Delta y = \alpha \Delta x$$

- Based on the fact that at a fine enough resolution, any smooth, continuous function is locally linear at any point



# Multivariate scalar function

- Scalar function of vector argument

$$y = f(\boldsymbol{x}) \quad (f : \mathbb{R}^d \rightarrow \mathbb{R})$$

$$\Delta y = \alpha \Delta \boldsymbol{x}$$

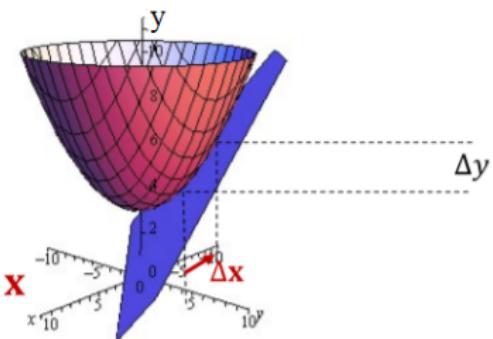
# Multivariate scalar function

- Scalar function of vector argument  
 $y = f(\mathbf{x})$  ( $f : \mathbb{R}^d \rightarrow \mathbb{R}$ )

$$\Delta y = \alpha \Delta \mathbf{x}$$

- $\Delta \mathbf{x}$  is now a vector

$$\Delta \mathbf{x} = \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \vdots \\ \Delta x_d \end{bmatrix}$$



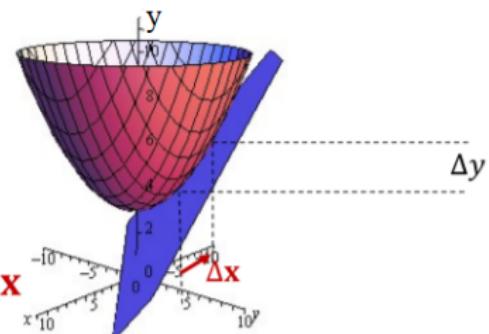
## Multivariate scalar function

- Scalar function of vector argument  
 $y = f(\mathbf{x})$  ( $f : \mathbb{R}^d \rightarrow \mathbb{R}$ )

$$\Delta y = \alpha \Delta x$$

- $\Delta x$  is now a vector

$$\Delta x = \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \vdots \\ \Delta x_d \end{bmatrix}$$



- Giving us that  $\alpha$  is a row vector:  $\alpha = [\alpha_1 \quad \alpha_2 \quad \dots \quad \alpha_d]$

$$\Delta y = \alpha_1 \Delta x_1 + \alpha_2 \Delta x_2 + \dots + \alpha_d \Delta x_d$$

- The partial derivative  $\alpha_i$  gives us how  $y$  increments when only  $x_i$  is incremented.
  - Often represented as  $\frac{\partial y}{\partial x_i}$

$$\Delta y = \frac{\partial y}{\partial x_1} \Delta x_1 + \frac{\partial y}{\partial x_2} \Delta x_2 + \dots + \frac{\partial y}{\partial x_d} \Delta x_d$$

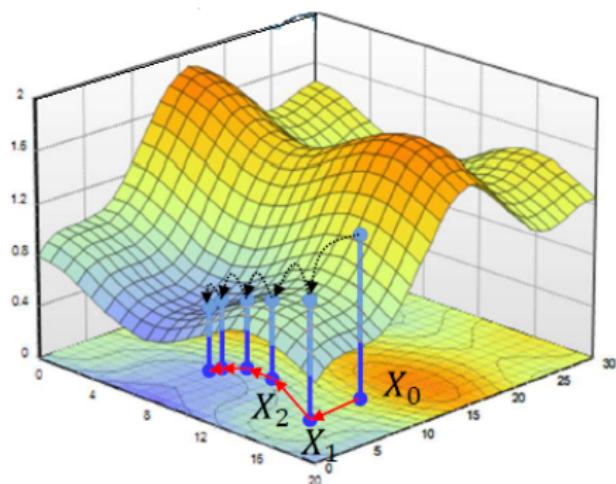
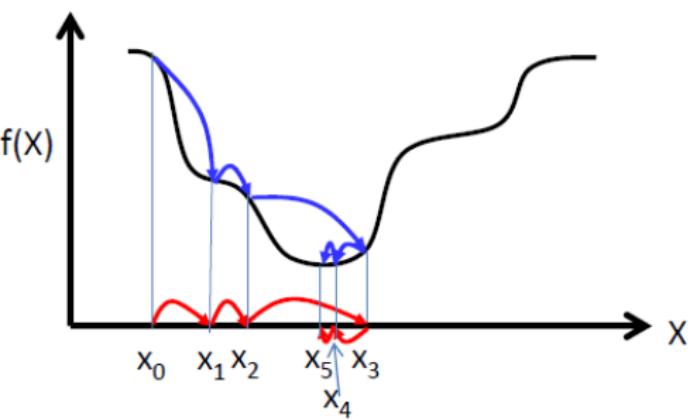
# Iterative solutions (Gradient Descent)

## Iterative solutions

- Start from an initial guess  $X_0$  for the optimal  $X$
- Update the guess towards a (hopefully) "better" value of  $f(X)$
- Stop when  $f(X)$  no longer decreases

## Problems

- Which direction to step in
- How big must the steps be



## Gradient descent

- The gradient descent method to find the minimum of a function iteratively

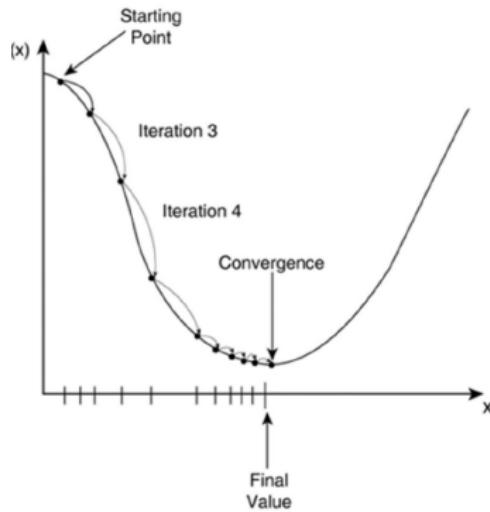
$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla_{\mathbf{x}} f(\mathbf{x}^t)$$

- $\eta$  is the "step size" (Also called "learning rate")

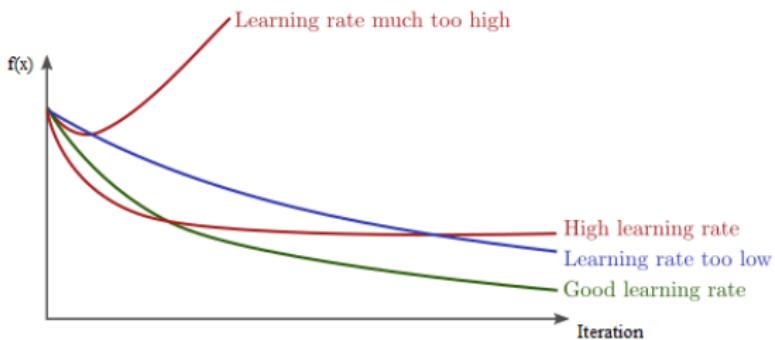
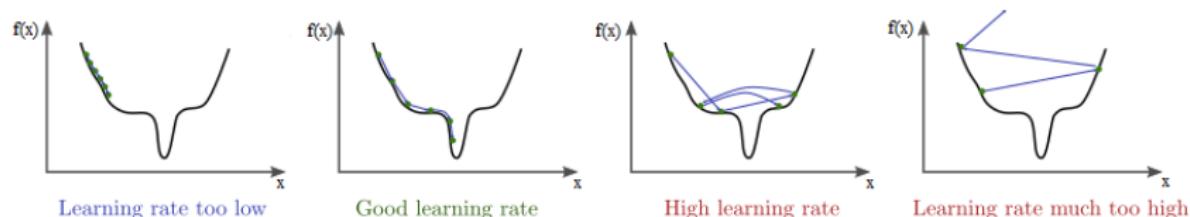
- The gradient descent algorithm converges when the following criterion is satisfied.

$$|f(\mathbf{x}^{t+k}) - f(\mathbf{x}^t)| < \epsilon$$

- $k$  is a hyperparameter.



# Influence of step size



Machine Learning  
○○○

Supervised Learning  
○○○○○○○

Logistic Regression  
○○○○

Loss Function  
○○○○○

Optimization  
○○○○○○○

Regularization  
●○○○○○

Softmax Classifier  
○○○○○○○○

Deep Neural Networks  
○○○○○○○○○○○○○○○○

## Regularization

## Learning VS. Optimization

## Generalization VS. Overfitting

- Ability to perform well on previously unobserved inputs is called **generalization**.

# Learning VS. Optimization

## Generalization VS. Overfitting

- Ability to perform well on previously unobserved inputs is called **generalization**.
  - What separates machine learning from optimization is that we want the **generalization error**, also called the test error, to be low as well.

## Learning VS. Optimization

Generalization VS. Overfitting

- Ability to perform well on previously unobserved inputs is called **generalization**.
  - What separates machine learning from optimization is that we want the **generalization error**, also called the test error, to be low as well.
  - We typically estimate the generalization error of a machine learning model by measuring its performance on a **test set** of examples that were collected separately from the training set.
  - How can we affect performance on the test set when we get to observe only the training set? The field of **statistical learning theory** provides some answers.

## Learning VS. Optimization

## Generalization VS. Overfitting

- Ability to perform well on previously unobserved inputs is called **generalization**.
  - What separates machine learning from optimization is that we want the **generalization error**, also called the test error, to be low as well.
  - We typically estimate the generalization error of a machine learning model by measuring its performance on a **test set** of examples that were collected separately from the training set.
  - How can we affect performance on the test set when we get to observe only the training set? The field of **statistical learning theory** provides some answers.
  - If the training and the test set are collected arbitrarily, there is indeed little we can do.

## Learning VS. Optimization

## Generalization VS. Overfitting

- Ability to perform well on previously unobserved inputs is called **generalization**.
  - What separates machine learning from optimization is that we want the **generalization error**, also called the test error, to be low as well.
  - We typically estimate the generalization error of a machine learning model by measuring its performance on a **test set** of examples that were collected separately from the training set.
  - How can we affect performance on the test set when we get to observe only the training set? The field of **statistical learning theory** provides some answers.
  - If the training and the test set are collected arbitrarily, there is indeed little we can do.
  - The train and test data are generated by a probability distribution over datasets called the data generating process. We typically make a set of assumptions known collectively as the **i.i.d.** assumptions.

## Learning VS. Optimization

## Generalization VS. Overfitting

- Ability to perform well on previously unobserved inputs is called **generalization**.
  - What separates machine learning from optimization is that we want the **generalization error**, also called the test error, to be low as well.
  - We typically estimate the generalization error of a machine learning model by measuring its performance on a **test set** of examples that were collected separately from the training set.
  - How can we affect performance on the test set when we get to observe only the training set? The field of **statistical learning theory** provides some answers.
  - If the training and the test set are collected arbitrarily, there is indeed little we can do.
  - The train and test data are generated by a probability distribution over datasets called the data generating process. We typically make a set of assumptions known collectively as the **i.i.d.** assumptions.
  - These assumptions are that the examples in each dataset are independent from each other, and that the train set and test set are identically distributed, drawn from the same probability distribution as each other.

(Deep Learning, Ian Goodfellow, MIT press, 2016)

## Machine learning

## underfitting and overfitting

- The factors determining how well a machine learning algorithm will perform are its ability to:
    - Make the training error small.
    - Make the gap between training and test error small.

# Machine learning

## underfitting and overfitting

- The factors determining how well a machine learning algorithm will perform are its ability to:
  - Make the training error small.
  - Make the gap between training and test error small.
- These two factors correspond to the two central challenges in machine learning: **underfitting** and **overfitting**. Underfitting occurs when the model is not able to obtain a sufficiently low error value on the training set. Overfitting occurs when the gap between the training error and test error is too large.

## Machine learning

## underfitting and overfitting

- The factors determining how well a machine learning algorithm will perform are its ability to:
    - Make the training error small.
    - Make the gap between training and test error small.
  - These two factors correspond to the two central challenges in machine learning: **underfitting** and **overfitting**. Underfitting occurs when the model is not able to obtain a sufficiently low error value on the training set. Overfitting occurs when the gap between the training error and test error is too large.
  - We can control whether a model is more likely to overfit or underfit by altering its **capacity**.

# Machine learning

## underfitting and overfitting

- The factors determining how well a machine learning algorithm will perform are its ability to:
  - Make the training error small.
  - Make the gap between training and test error small.
- These two factors correspond to the two central challenges in machine learning: **underfitting** and **overfitting**. Underfitting occurs when the model is not able to obtain a sufficiently low error value on the training set. Overfitting occurs when the gap between the training error and test error is too large.
- We can control whether a model is more likely to overfit or underfit by altering its **capacity**.
- Models with low capacity may struggle to fit the training set. Models with high capacity can overfit by memorizing properties of the training set that do not serve them well on the test set.

## Machine learning

## underfitting and overfitting

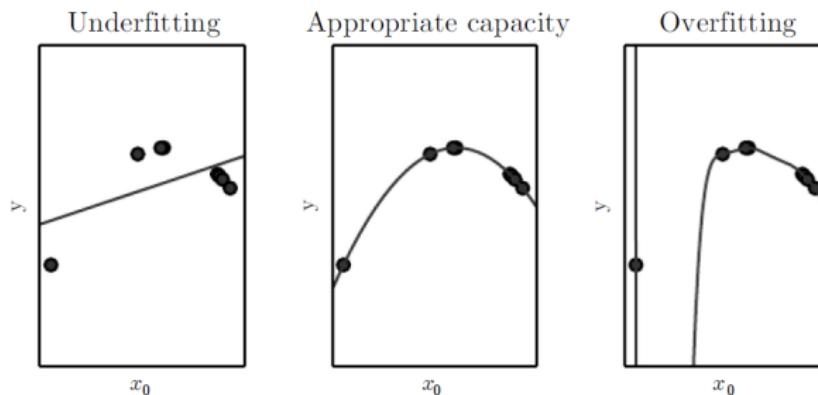
- The factors determining how well a machine learning algorithm will perform are its ability to:
    - Make the training error small.
    - Make the gap between training and test error small.
  - These two factors correspond to the two central challenges in machine learning: **underfitting** and **overfitting**. Underfitting occurs when the model is not able to obtain a sufficiently low error value on the training set. Overfitting occurs when the gap between the training error and test error is too large.
  - We can control whether a model is more likely to overfit or underfit by altering its **capacity**.
  - Models with low capacity may struggle to fit the training set. Models with high capacity can overfit by memorizing properties of the training set that do not serve them well on the test set.
  - One way to control the capacity of a learning algorithm is by choosing its **hypothesis space**, the set of functions that the learning algorithm is allowed to select as being the solution.

(Deep Learning, Ian Goodfellow, MIT press, 2016)

# Machine learning

## underfitting and overfitting

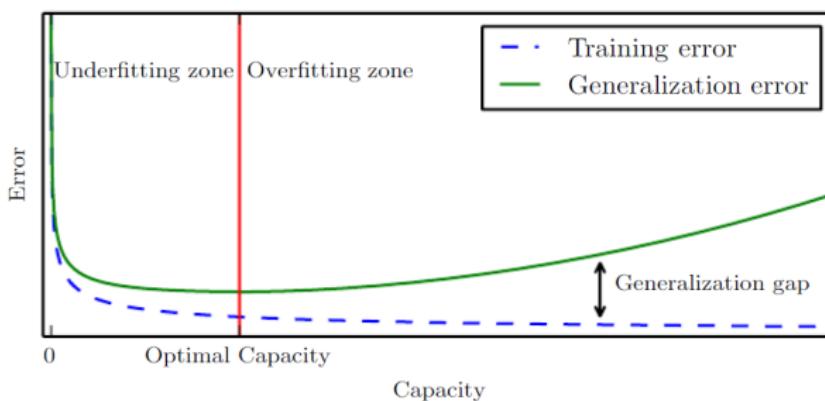
- Machine learning algorithms will generally perform best when their capacity is appropriate for the true complexity of the task they need to perform and the amount of training data they are provided with.



(Deep Learning, Ian Goodfellow, MIT press, 2016)

## Capacity

Typically, generalization error has a U-shaped curve as a function of model capacity.

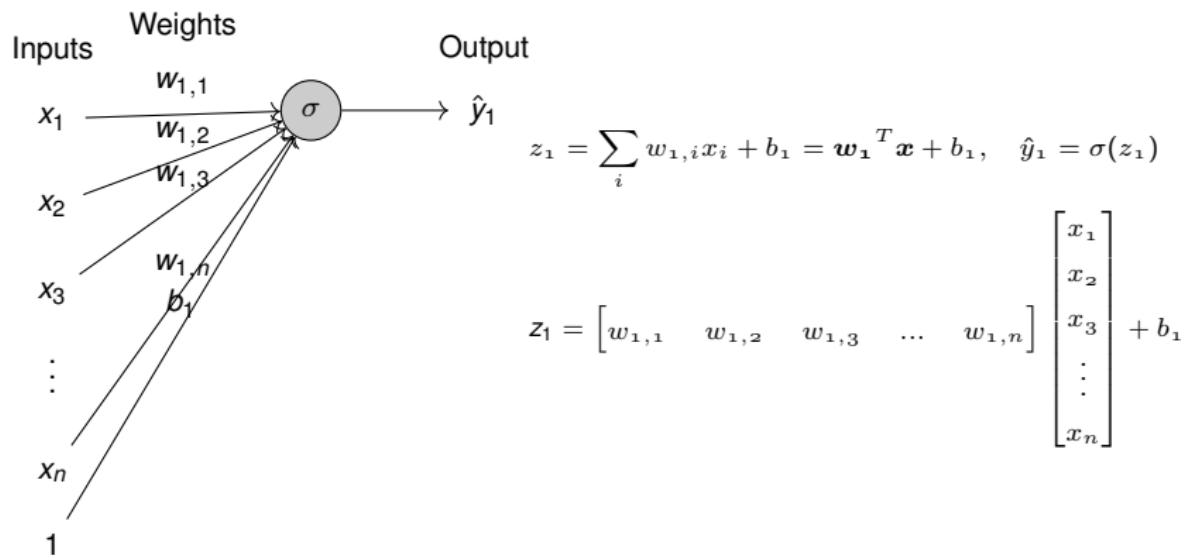


## Regularization

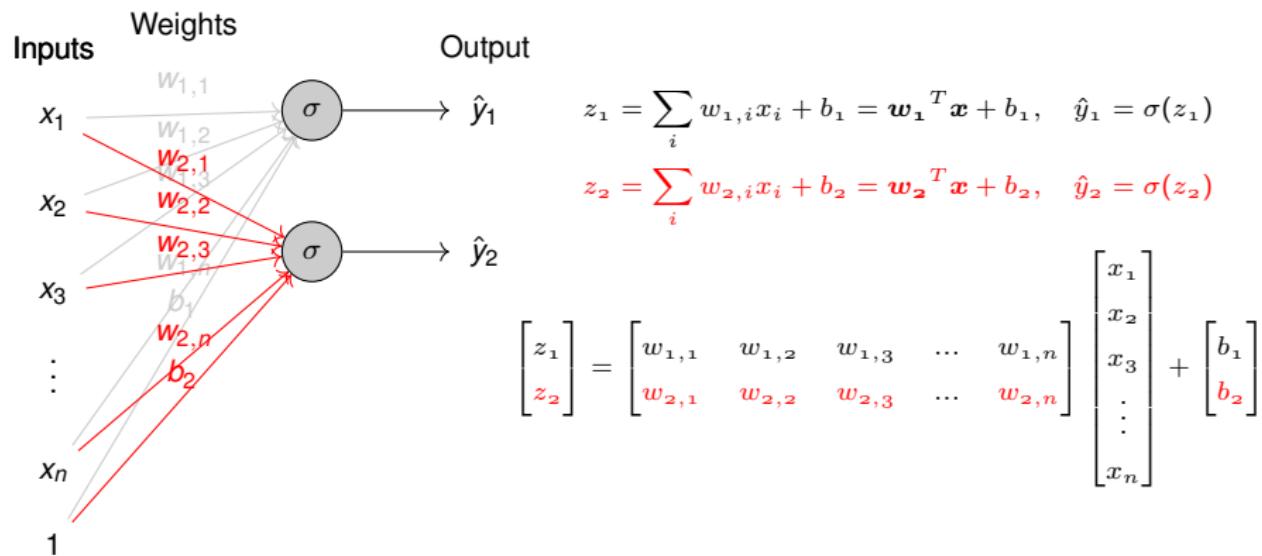
- **Regularization** is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.
    - We can give a learning algorithm a **preference** for one solution in its hypothesis space to another. This means that both functions are eligible, but one is preferred.

## Softmax Classifier

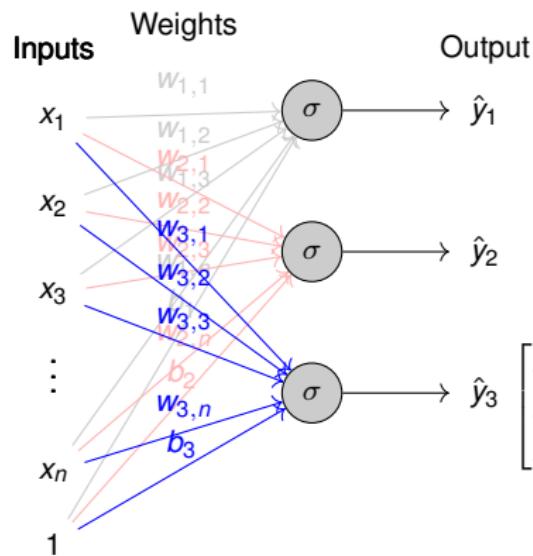
# Multinomial logistic regression (softmax classifier)



# Multinomial logistic regression (softmax classifier)



# Multinomial logistic regression (softmax classifier)



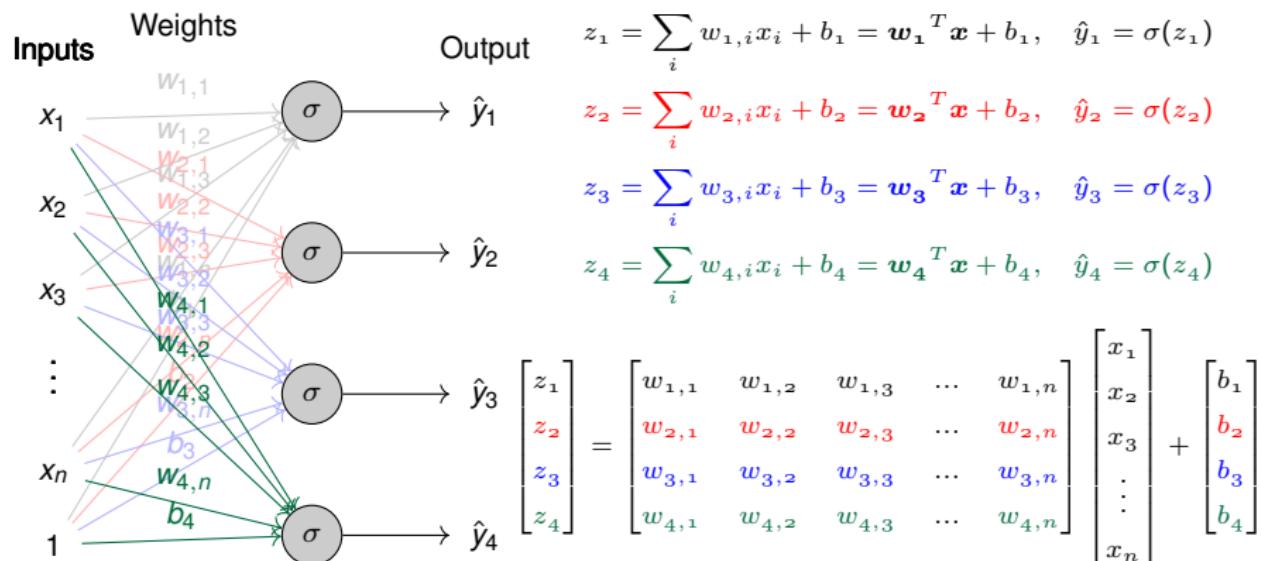
$$z_1 = \sum_i w_{1,i} x_i + b_1 = \mathbf{w}_1^T \mathbf{x} + b_1, \quad \hat{y}_1 = \sigma(z_1)$$

$$z_2 = \sum_i w_{2,i} x_i + b_2 = \mathbf{w}_2^T \mathbf{x} + b_2, \quad \hat{y}_2 = \sigma(z_2)$$

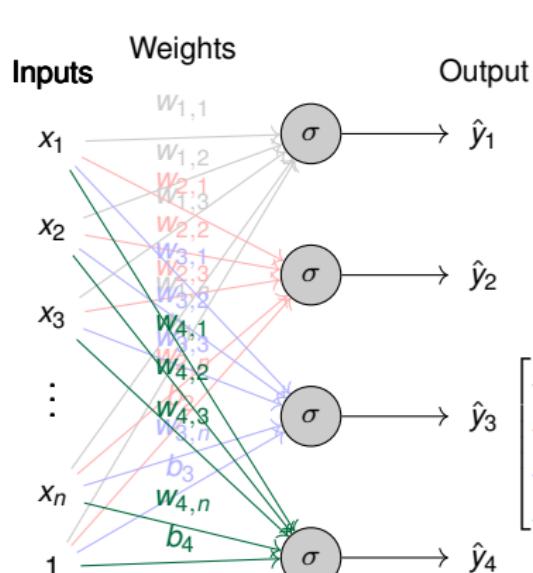
$$z_3 = \sum_i w_{3,i} x_i + b_3 = \mathbf{w}_3^T \mathbf{x} + b_3, \quad \hat{y}_3 = \sigma(z_3)$$

$$\begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} & \dots & w_{1,n} \\ w_{2,1} & w_{2,2} & w_{2,3} & \dots & w_{2,n} \\ w_{3,1} & w_{3,2} & w_{3,3} & \dots & w_{3,n} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

# Multinomial logistic regression (softmax classifier)



# Multinomial logistic regression (softmax classifier)



$$z_1 = \sum_i w_{1,i} x_i + b_1 = \mathbf{w}_1^T \mathbf{x} + b_1, \quad \hat{y}_1 = \sigma(z_1)$$

$$z_2 = \sum_i w_{2,i} x_i + b_2 = \mathbf{w}_2^T \mathbf{x} + b_2, \quad \hat{y}_2 = \sigma(z_2)$$

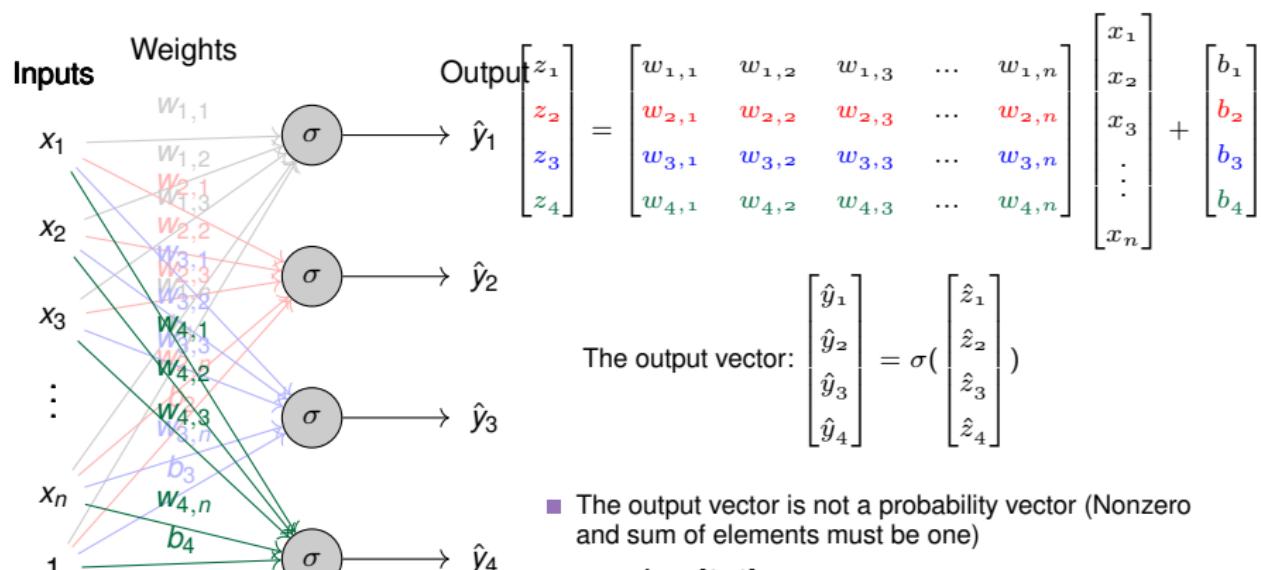
$$z_3 = \sum_i w_{3,i} x_i + b_3 = \mathbf{w}_3^T \mathbf{x} + b_3, \quad \hat{y}_3 = \sigma(z_3)$$

$$z_4 = \sum_i w_{4,i} x_i + b_4 = \mathbf{w}_4^T \mathbf{x} + b_4, \quad \hat{y}_4 = \sigma(z_4)$$

$$\begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{bmatrix} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} & \dots & w_{1,n} \\ w_{2,1} & w_{2,2} & w_{2,3} & \dots & w_{2,n} \\ w_{3,1} & w_{3,2} & w_{3,3} & \dots & w_{3,n} \\ w_{4,1} & w_{4,2} & w_{4,3} & \dots & w_{4,n} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix}$$

$$\mathbf{z} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

## Multinomial logistic regression (softmax classifier)



- The output vector is not a probability vector (Nonzero and sum of elements must be one)
- $\hat{y}_i \in [0, 1]$
- Cross entropy is fed by a probability vector

### Multi-class output: One-hot representations

- The labels are 0 and 1 in logistic regression
  - Consider a network that must distinguish if an input is a cat, a dog, camel, a hat, or a flower
  - For inputs of each of the four classes the desired output is a probability vector:

$$\text{Cat: } \begin{bmatrix} 1 & 0 & 0 & 0 \end{bmatrix}^T$$

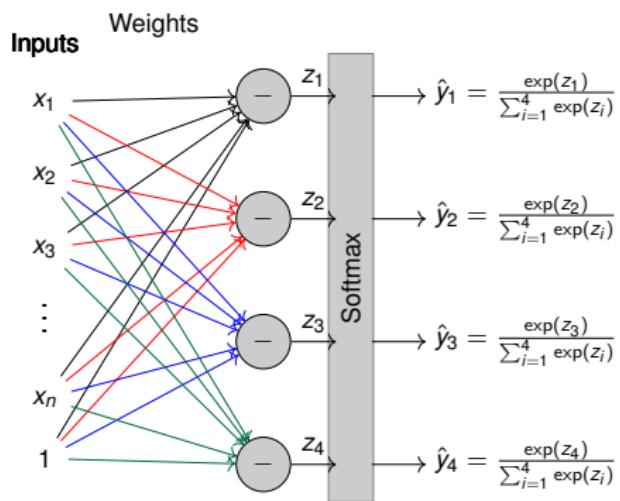
$$\text{Dog: } \begin{bmatrix} 0 & 1 & 0 & 0 \end{bmatrix}^T$$

$$\text{Fish: } \begin{bmatrix} 0 & 0 & 1 & 0 \end{bmatrix}^T$$

Frog:  $\begin{bmatrix} 0 & 0 & 0 & 1 \end{bmatrix}^T$

- For input of any class, we will have a four-dimensional vector output with four zeros and a single 1 at the position of the class
  - This is a one hot vector

## Softmax classifier



Softmax function:

$$P(\hat{y} = j | \mathbf{x}, W) = \frac{\exp(z_j)}{\sum_{i=1}^K \exp(z_i)} = \hat{y}_j$$

The output vector:  $\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \hat{y}_4 \end{bmatrix}$

The output of classifier:  $\operatorname{argmax}_j \hat{y}_j$

## Loss function

- ### ■ Cross entropy:

$$L = CE(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{i=1}^K y_i \log \hat{y}_i$$

- $y$  is a one-hot vector
  - If the true class of  $x$  is  $j$ , only the the  $j^{th}$  element of  $y$  is one, and the others are zero. Therefore:

$$L = CE(\mathbf{y}, \hat{\mathbf{y}}) = -\log \hat{y}_j$$

## Loss function

- ### ■ Cross entropy:

$$L = CE(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{i=1}^K y_i \log \hat{y}_i$$

- $y$  is a one-hot vector
  - If the true class of  $x$  is  $j$ , only the the  $j^{th}$  element of  $y$  is one, and the others are zero. Therefore:

$$L = CE(y, \hat{y}) = -\log \hat{y}_j$$

### example:



$$\mathbf{y} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} \begin{matrix} \text{dog} \\ \text{cat} \\ \text{fish} \\ \text{frog} \end{matrix} \quad \mathbf{z} = \begin{pmatrix} 3.2 \\ 5.1 \\ -1.7 \\ -3.7 \end{pmatrix} \quad softmax(\mathbf{z}) = \hat{\mathbf{y}} = \begin{pmatrix} 0.13 \\ 0.87 \\ 0 \\ 0 \end{pmatrix}$$

# Loss function

- Cross entropy:

$$L = CE(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^K y_i \log \hat{y}_i$$

- $\mathbf{y}$  is a one-hot vector
- If the true class of  $x$  is  $j$ , only the the  $j^{th}$  element of  $\mathbf{y}$  is one, and the others are zero. Therefore:

$$L = CE(\mathbf{y}, \hat{\mathbf{y}}) = - \log \hat{y}_j$$

**example:**



$$\mathbf{y} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} \begin{matrix} \text{dog} \\ \text{cat} \\ \text{fish} \\ \text{frog} \end{matrix} \quad z = \begin{pmatrix} 3.2 \\ 5.1 \\ -1.7 \\ -3.7 \end{pmatrix} \quad softmax(z) = \hat{\mathbf{y}} = \begin{pmatrix} 0.13 \\ 0.87 \\ 0 \\ 0 \end{pmatrix}$$

$$L = CE(\mathbf{y}, \hat{\mathbf{y}}) = - \log \hat{y}_2 = - \log 0.87 = 0.14$$

# Loss function

- Cross entropy:

$$L = CE(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^K y_i \log \hat{y}_i$$

- $\mathbf{y}$  is a one-hot vector
- If the true class of  $x$  is  $j$ , only the the  $j^{th}$  element of  $\mathbf{y}$  is one, and the others are zero. Therefore:

$$L = CE(\mathbf{y}, \hat{\mathbf{y}}) = - \log \hat{y}_j$$

**example:**



$$\mathbf{y} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} \begin{matrix} \text{dog} \\ \text{cat} \\ \text{fish} \\ \text{frog} \end{matrix} \quad z = \begin{pmatrix} 5.8 \\ 2.2 \\ -1.6 \\ 1.5 \end{pmatrix} \quad softmax(z) = \hat{\mathbf{y}} = \begin{pmatrix} 0.96 \\ 0.03 \\ 0 \\ 0.01 \end{pmatrix}$$

# Loss function

- Cross entropy:

$$L = CE(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^K y_i \log \hat{y}_i$$

- $\mathbf{y}$  is a one-hot vector
- If the true class of  $x$  is  $j$ , only the the  $j^{th}$  element of  $\mathbf{y}$  is one, and the others are zero. Therefore:

$$L = CE(\mathbf{y}, \hat{\mathbf{y}}) = - \log \hat{y}_j$$

**example:**

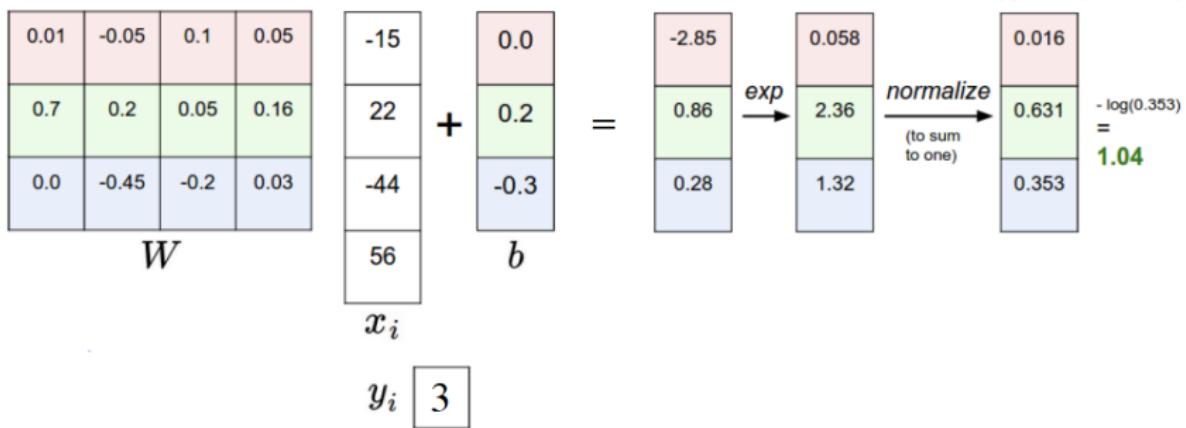


$$\mathbf{y} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} \begin{matrix} \text{dog} \\ \text{cat} \\ \text{fish} \\ \text{frog} \end{matrix} \quad z = \begin{pmatrix} 5.8 \\ 2.2 \\ -1.6 \\ 1.5 \end{pmatrix} \quad softmax(z) = \hat{\mathbf{y}} = \begin{pmatrix} 0.96 \\ 0.03 \\ 0 \\ 0.01 \end{pmatrix}$$

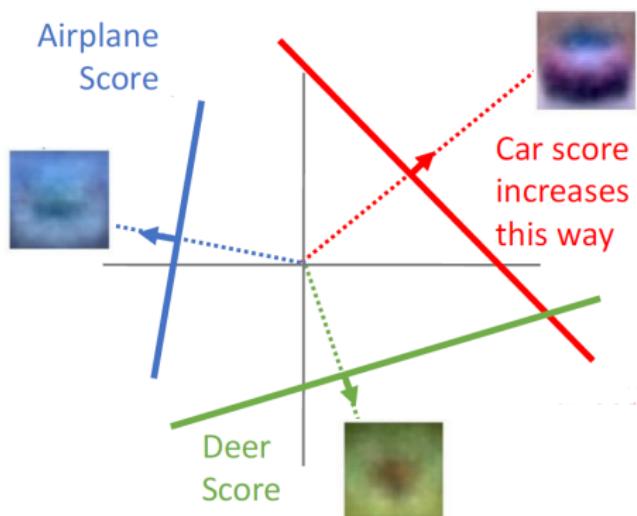
$$L = CE(\mathbf{y}, \hat{\mathbf{y}}) = - \log \hat{y}_2 = - \log 0.03 = 3.5$$

## Example

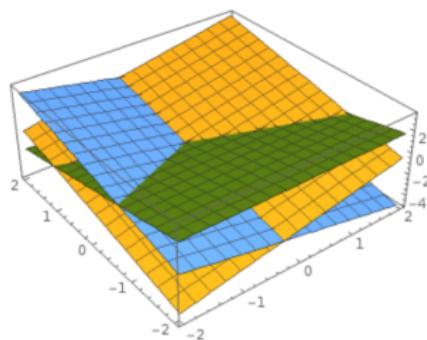
matrix multiply + bias offset



## Softmax classifier



Hyperplanes carving up a high-dimensional space



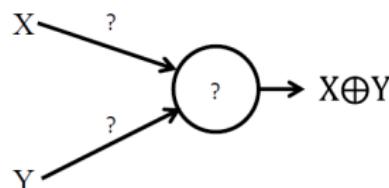
Plot created using Wolfram Cloud

<http://vision.stanford.edu/teaching/cs231n-demos/linear-classify/>

# Deep Neural Networks

Xor

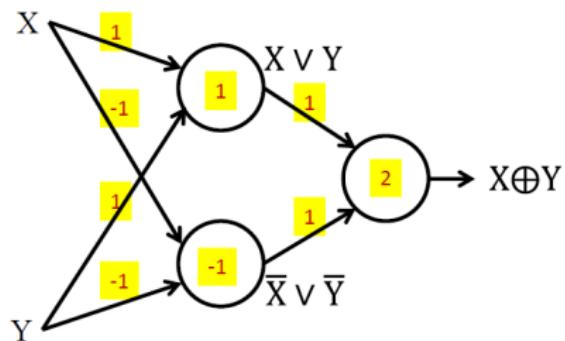
No solution for XOR!  
Not universal!



Values shown on edges are weights, numbers in the circles are thresholds

# Solution: add layers

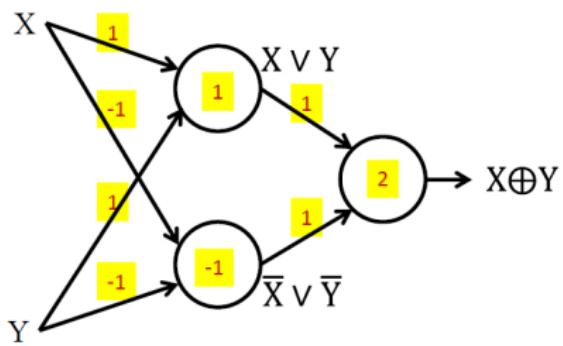
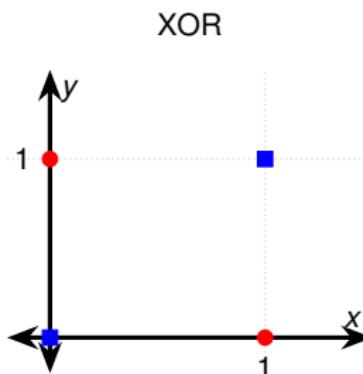
XOR



X	Y	A	B	Label
0	0	0	1	0
0	1	1	1	1
1	0	1	1	1
1	1	1	0	0

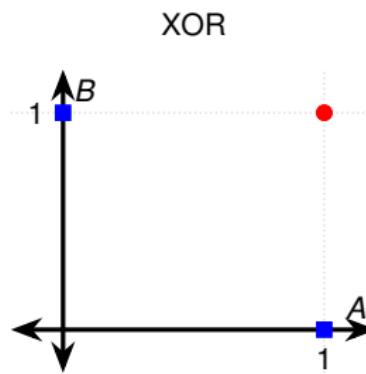
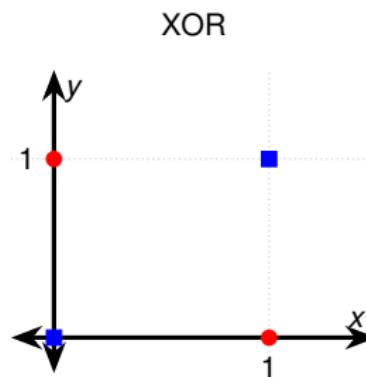
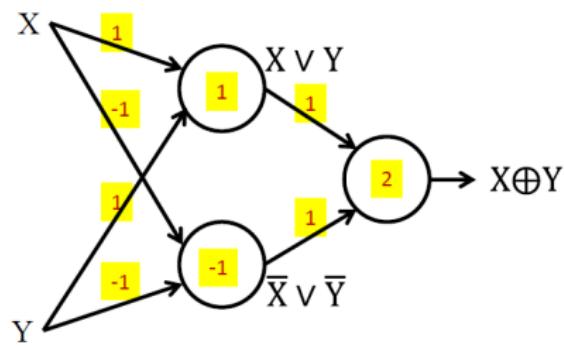
## Transformation

X	Y	A	B	Label
0	0	0	1	0
0	1	1	1	1
1	0	1	1	1
1	1	1	0	0



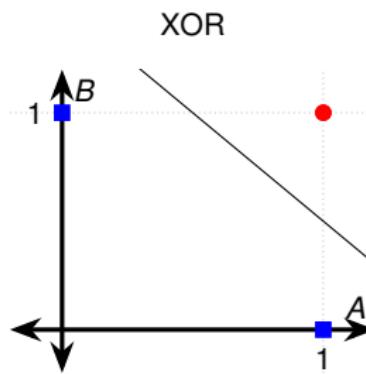
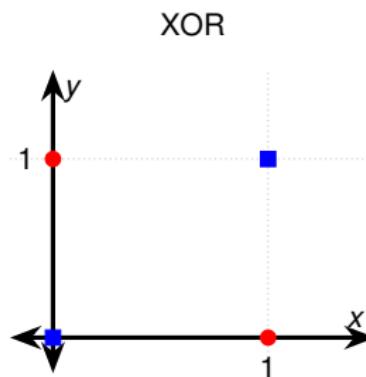
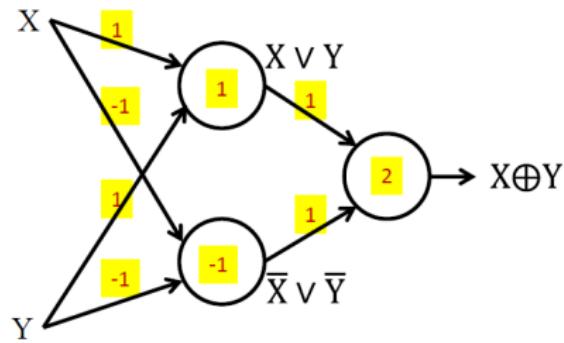
## Transformation

X	Y	A	B	Label
0	0	0	1	0
0	1	1	1	1
1	0	1	1	1
1	1	1	0	0

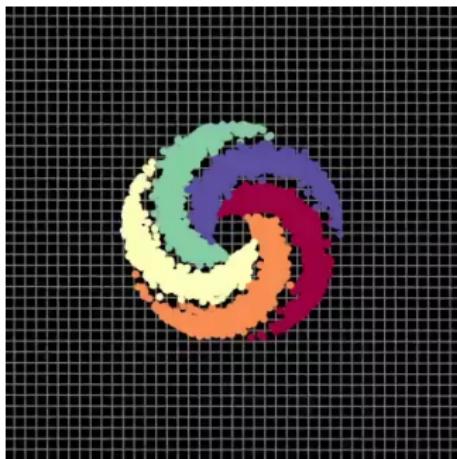


## Transformation

X	Y	A	B	Label
0	0	0	1	0
0	1	1	1	1
1	0	1	1	1
1	1	1	0	0

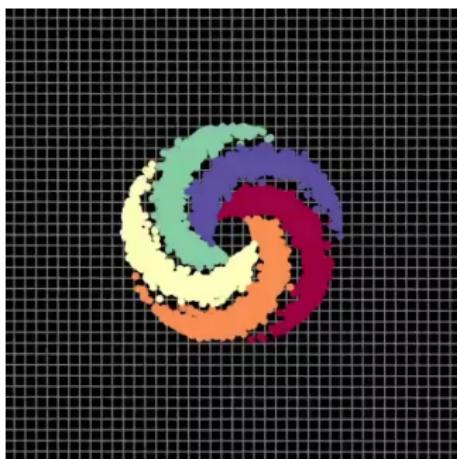


# Needs more layers

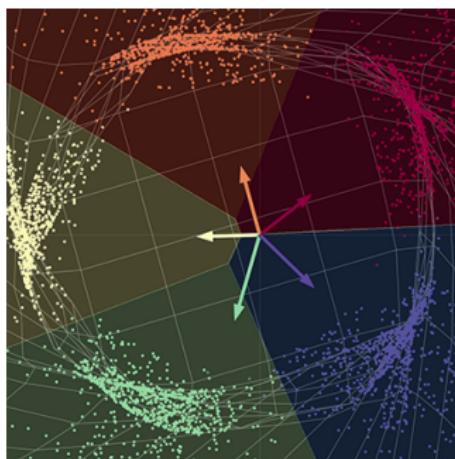


Input points, pre-network

# Needs more layers



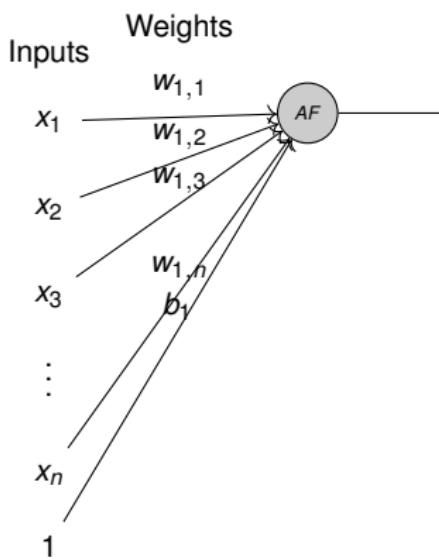
Input points, pre-network



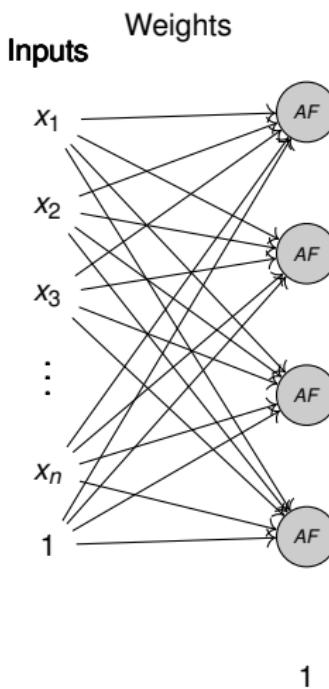
Output points, post-network

(Canziani, 2020)

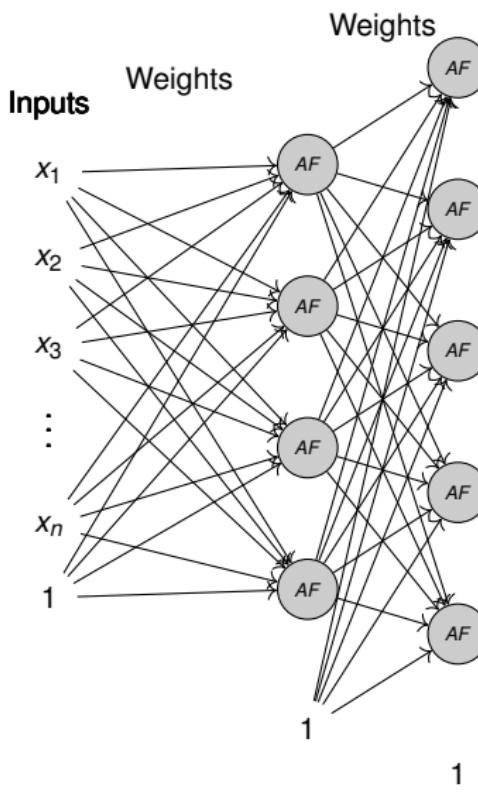
# Fully connected neural networks



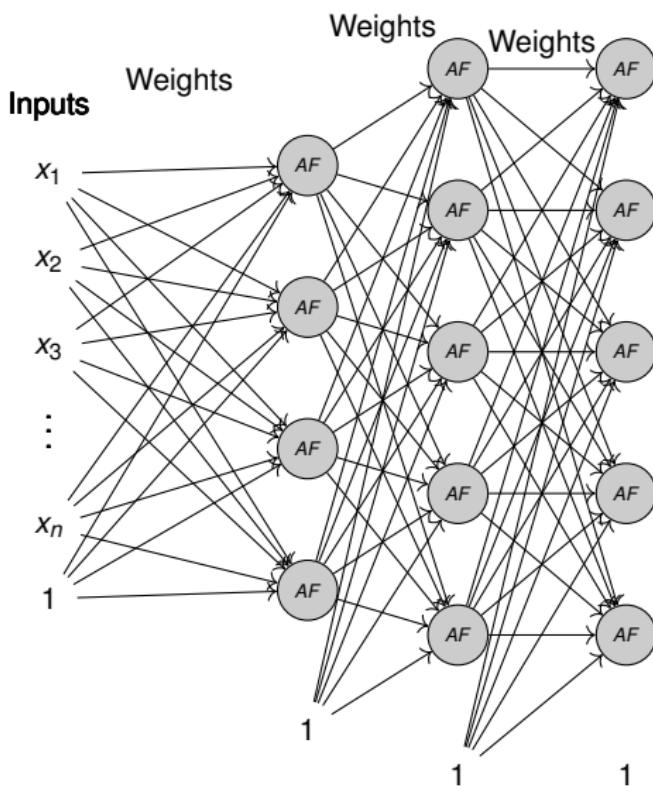
# Fully connected neural networks



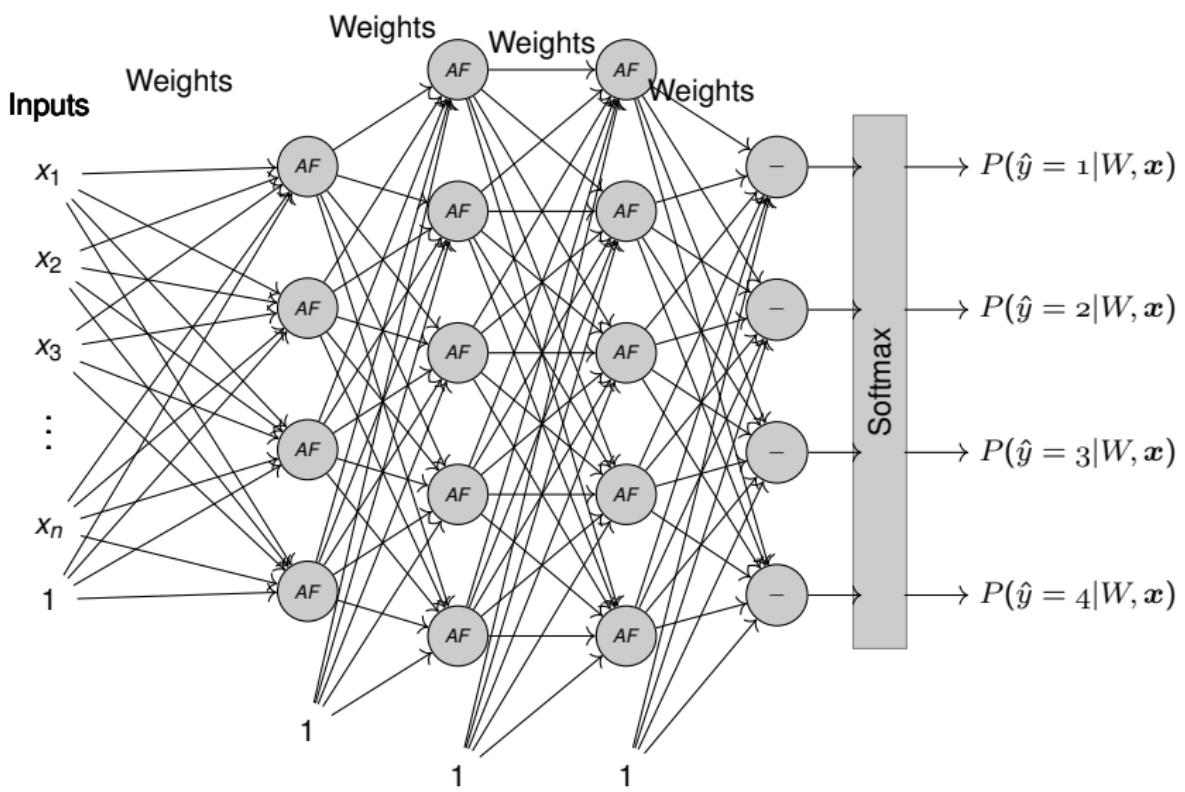
# Fully connected neural networks



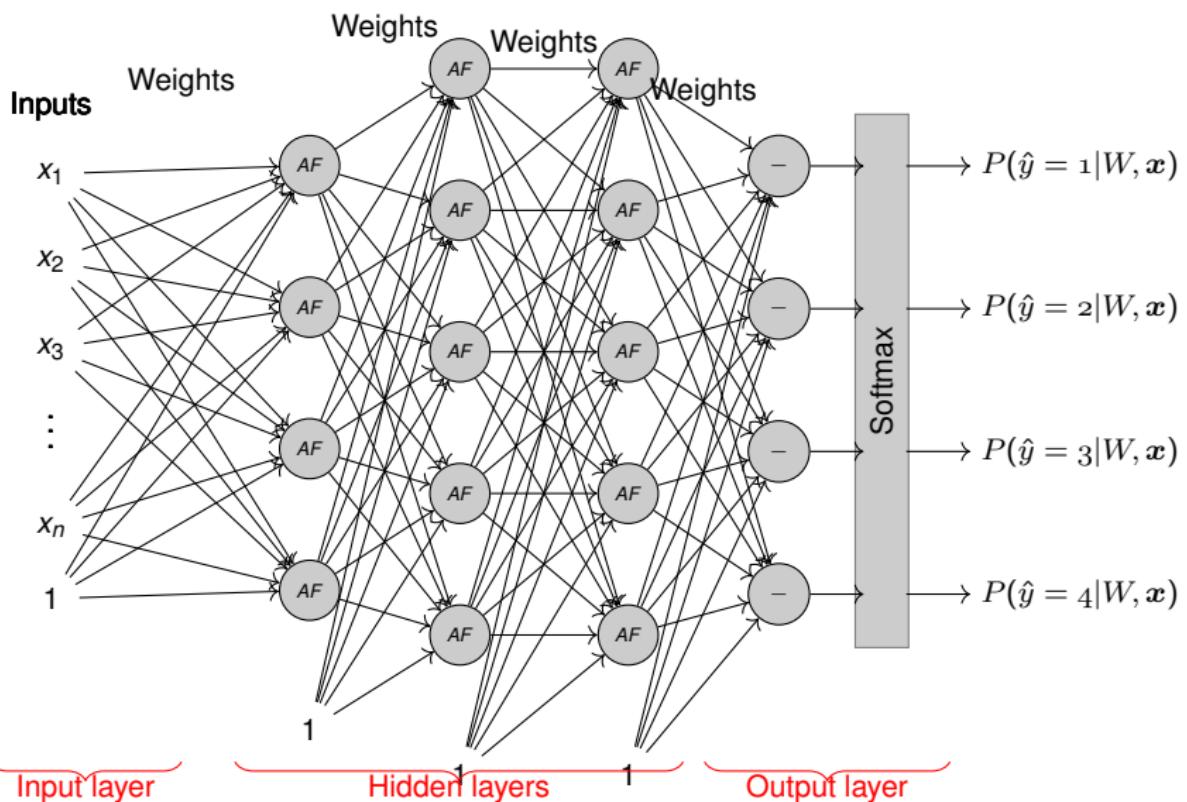
# Fully connected neural networks



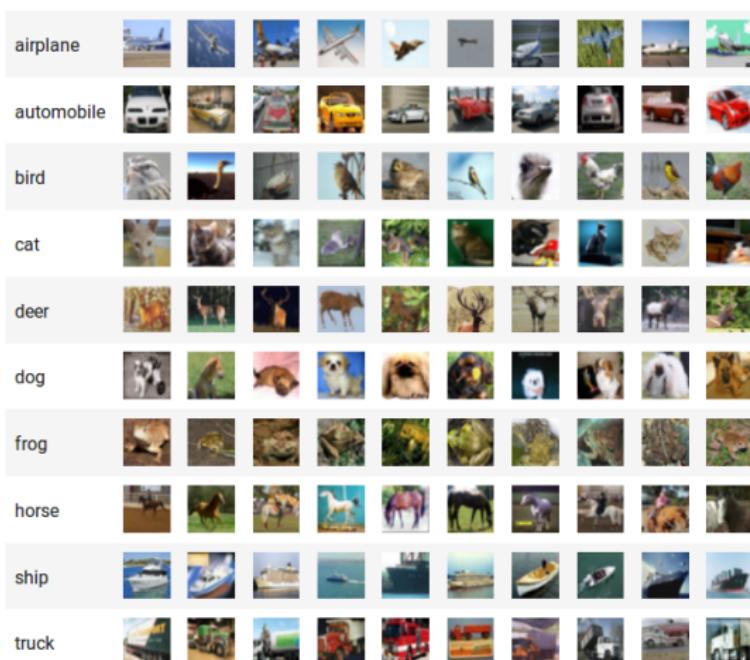
# Fully connected neural networks



# Fully connected neural networks

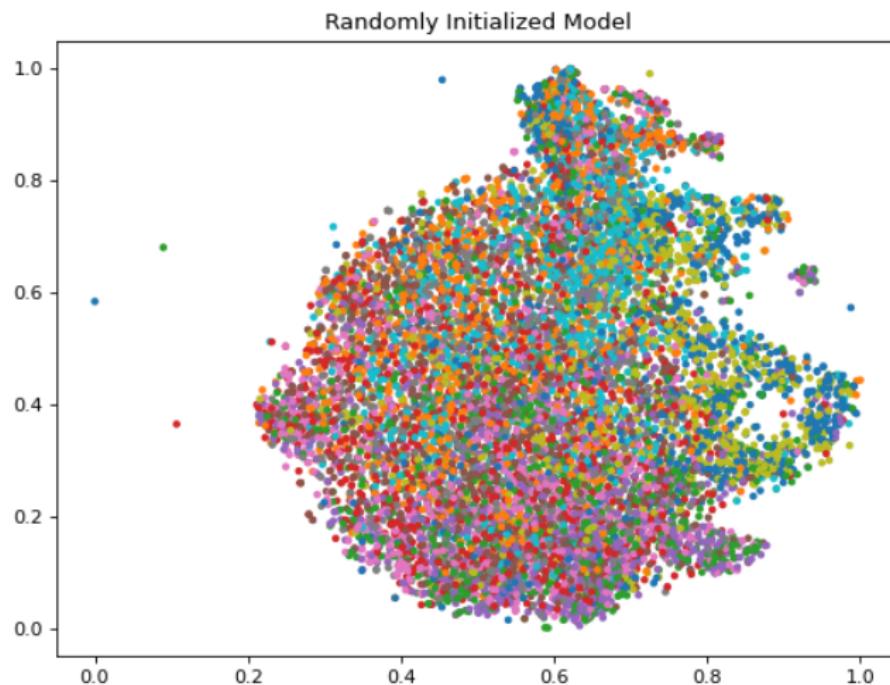


# CIFAR10

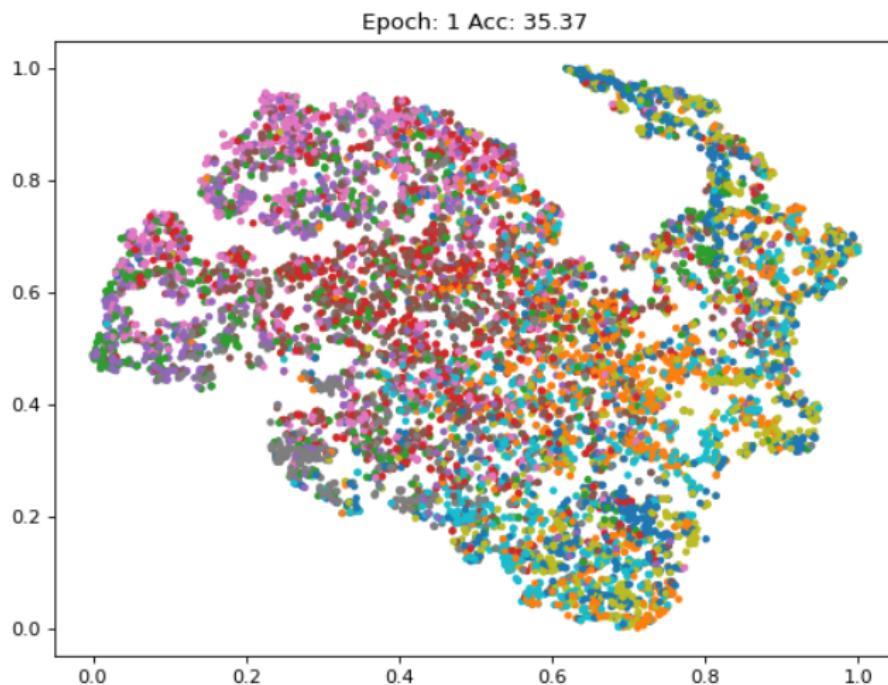


Example images from the **CIFAR-10** dataset.

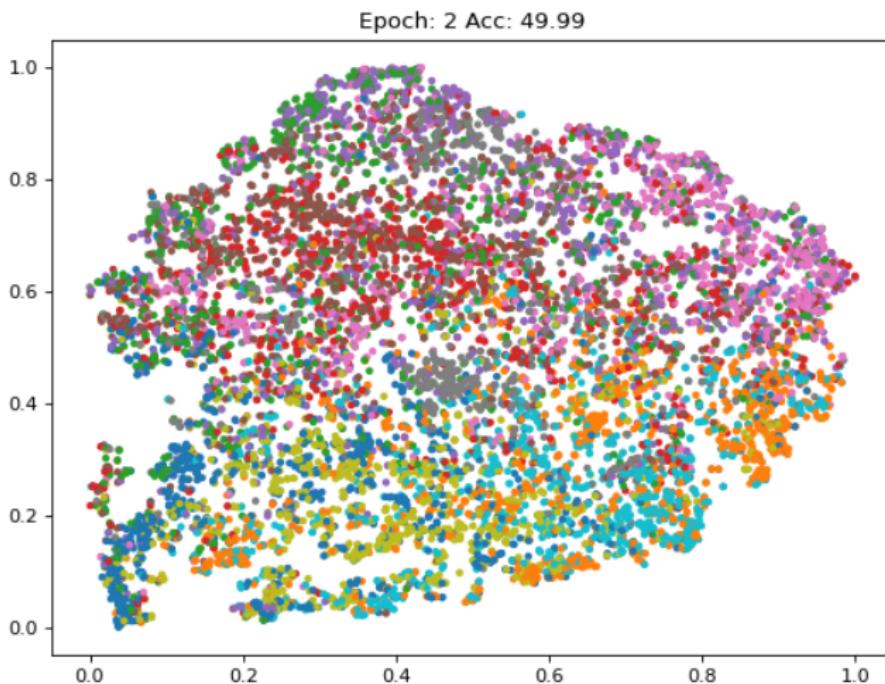
# Visualization of the last hidden layer output for CIFAR10 samples



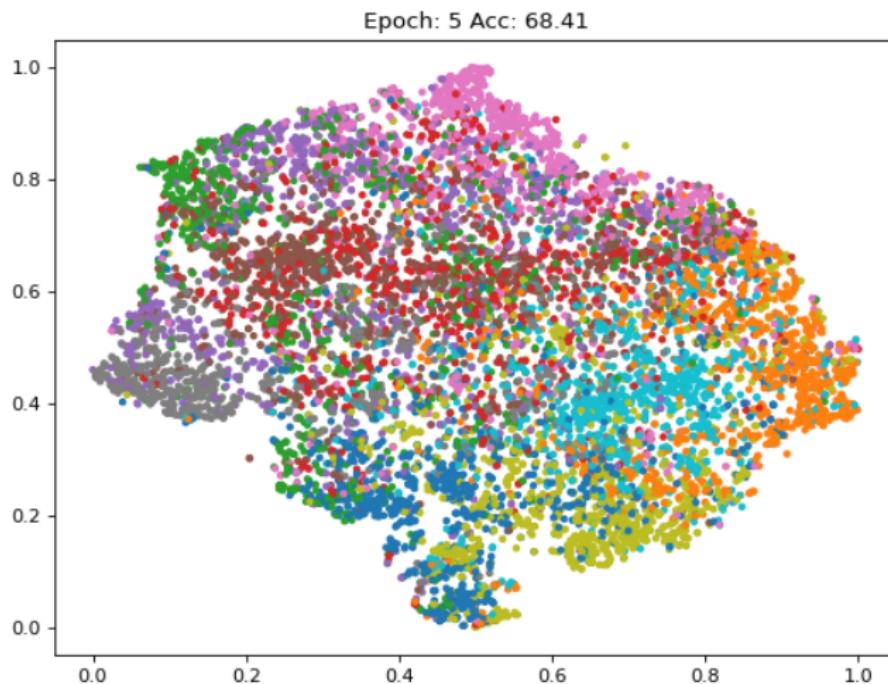
# Visualization of the last hidden layer output for CIFAR10 samples



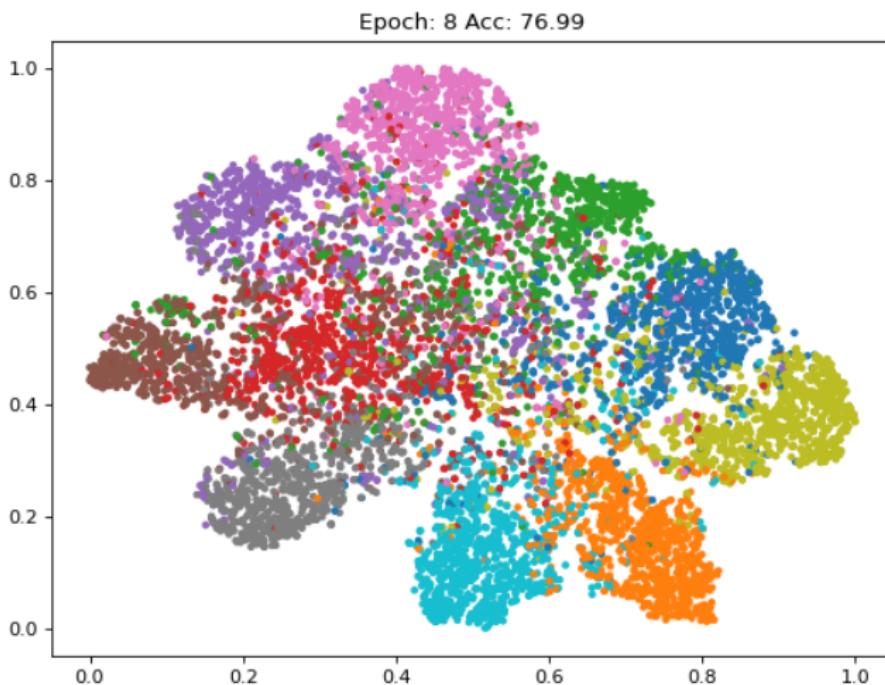
# Visualization of the last hidden layer output for CIFAR10 samples



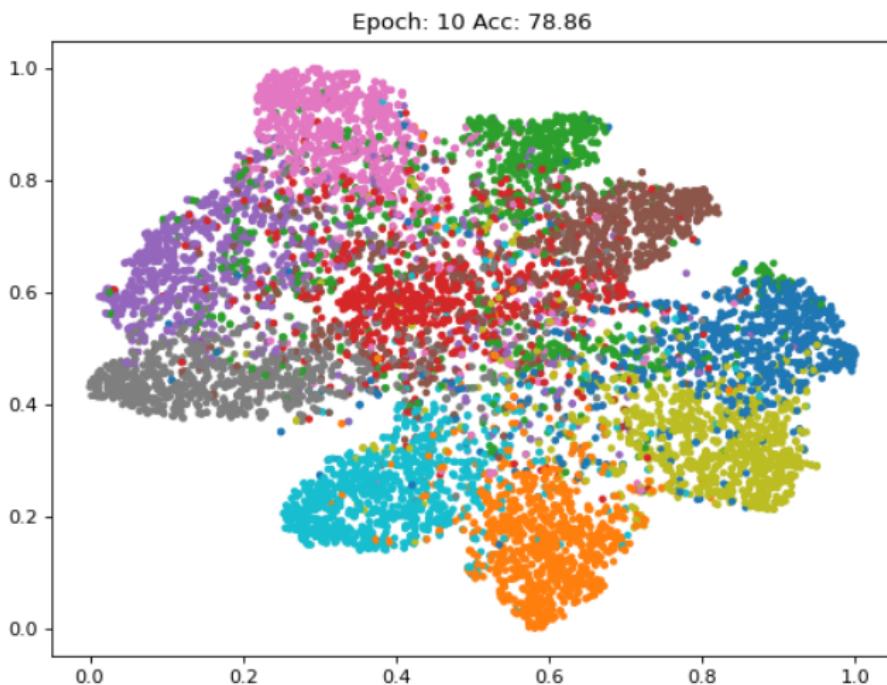
# Visualization of the last hidden layer output for CIFAR10 samples



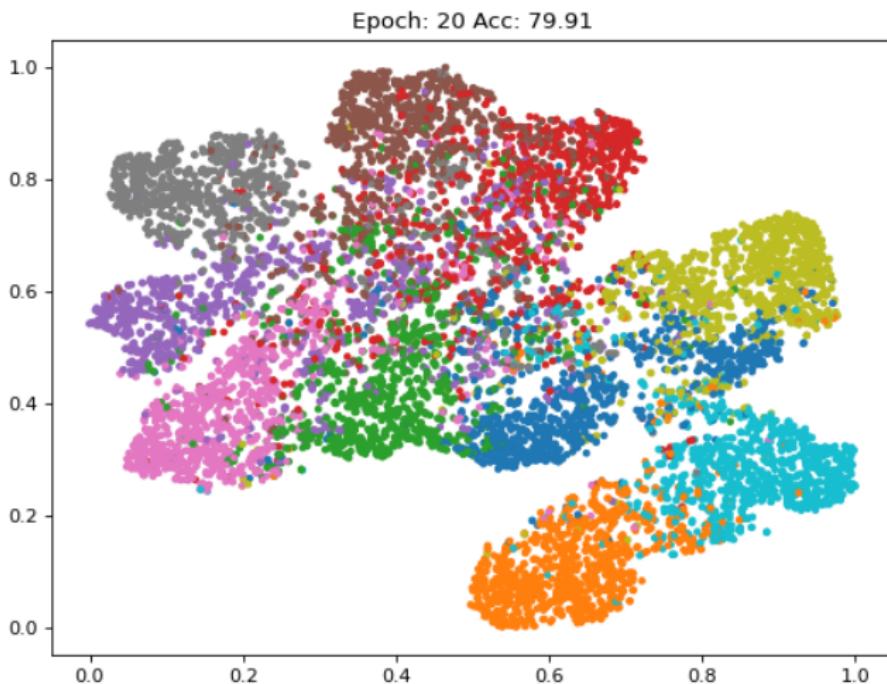
# Visualization of the last hidden layer output for CIFAR10 samples



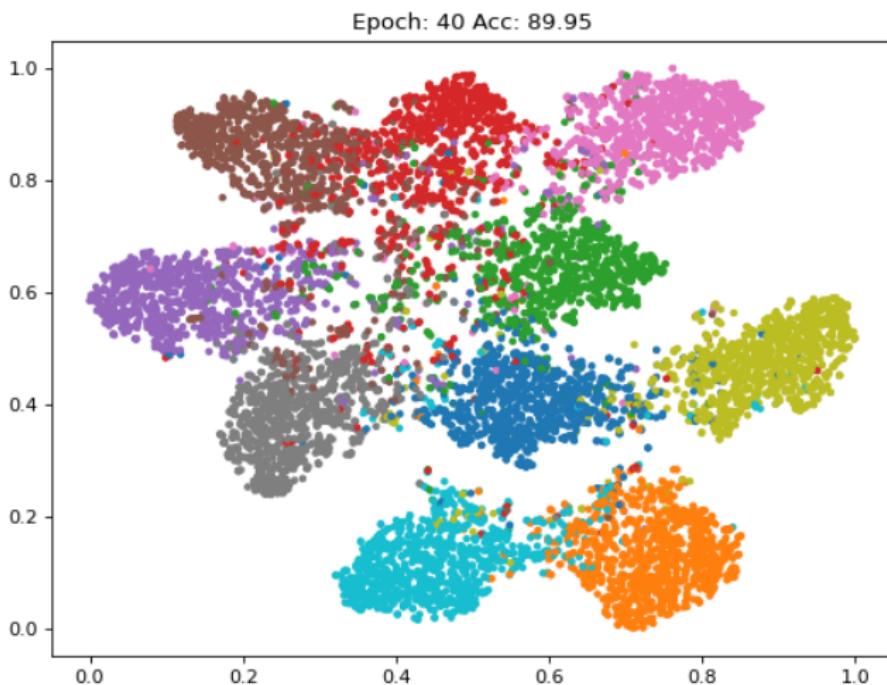
# Visualization of the last hidden layer output for CIFAR10 samples



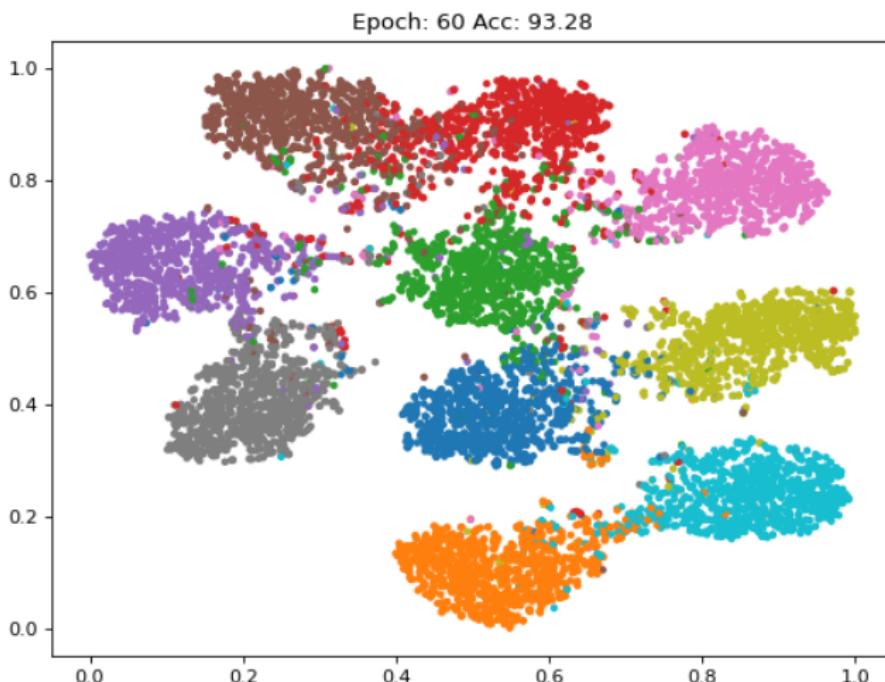
# Visualization of the last hidden layer output for CIFAR10 samples



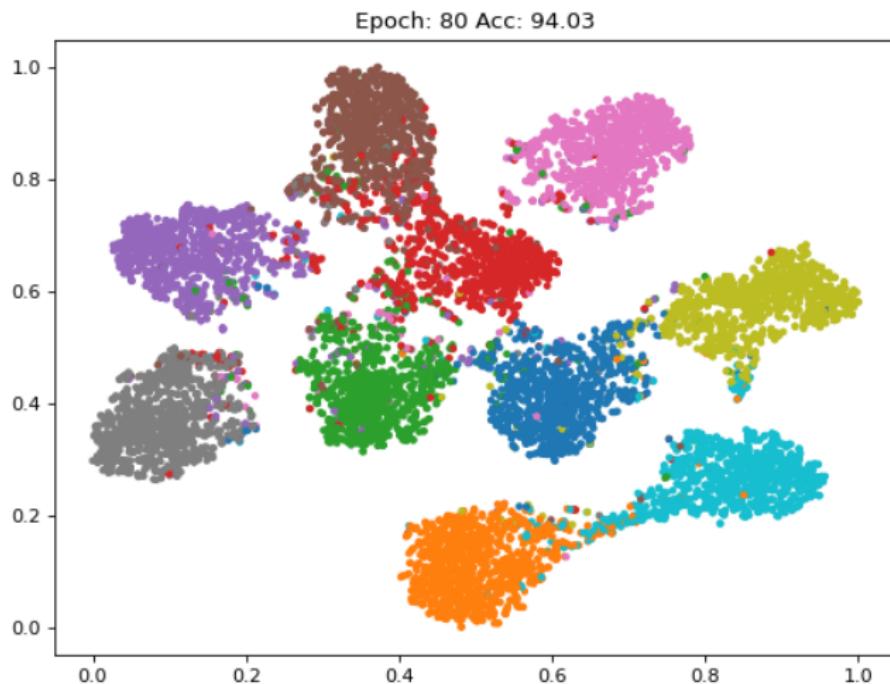
# Visualization of the last hidden layer output for CIFAR10 samples



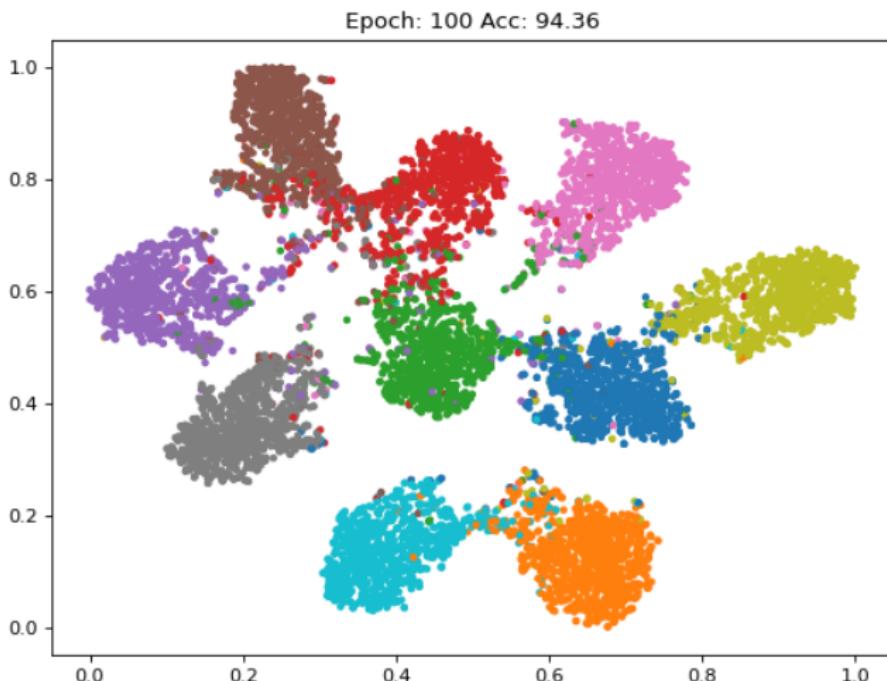
# Visualization of the last hidden layer output for CIFAR10 samples



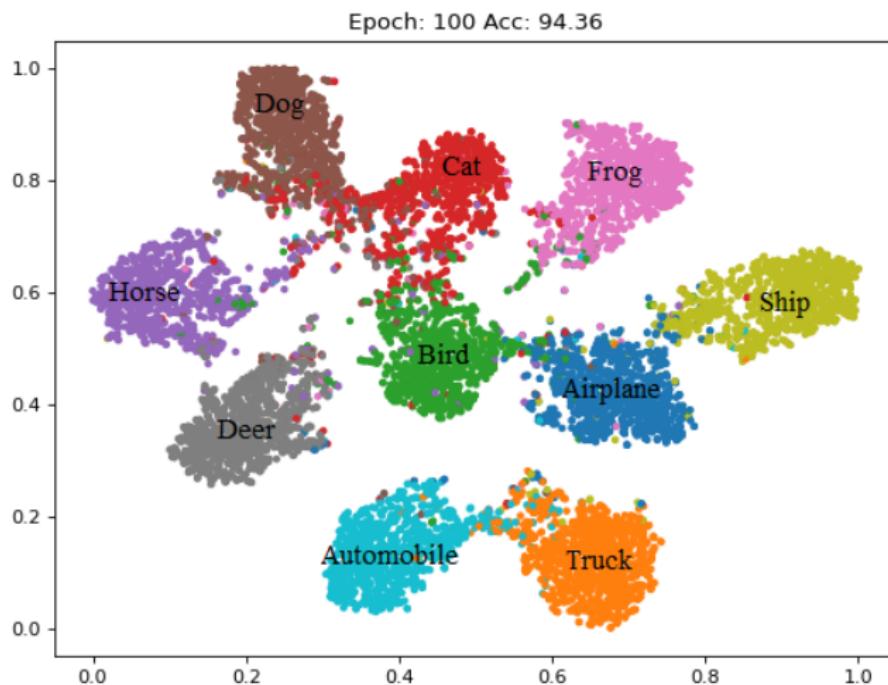
# Visualization of the last hidden layer output for CIFAR10 samples



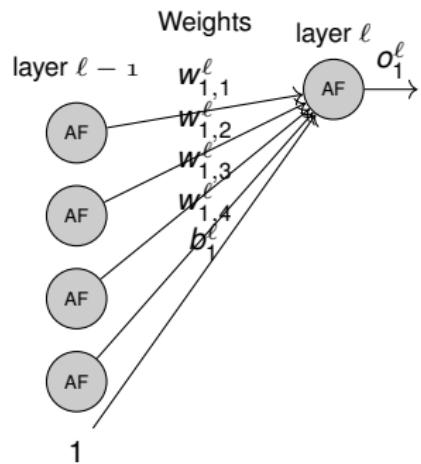
# Visualization of the last hidden layer output for CIFAR10 samples



# Visualization of the last hidden layer output for CIFAR10 samples



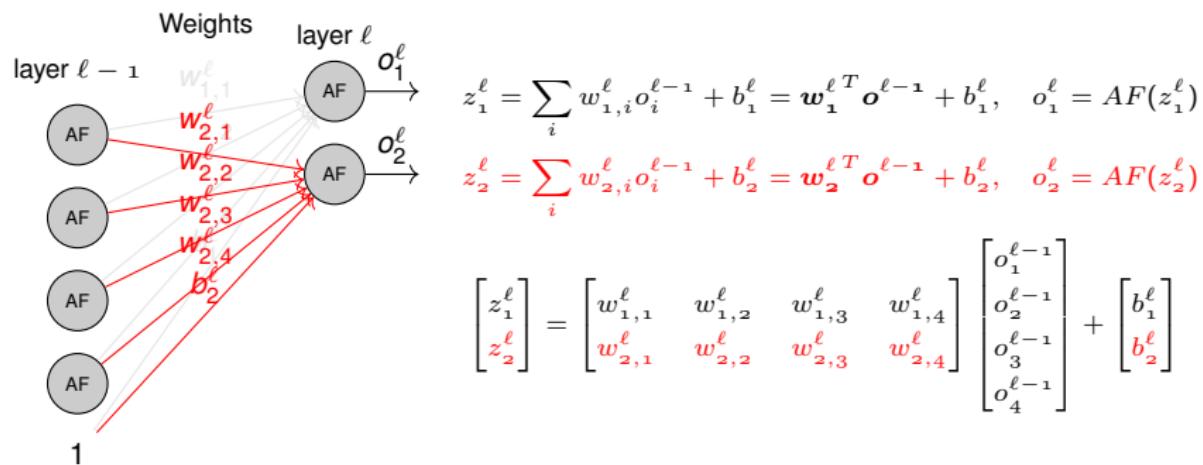
# Fully connected neural networks



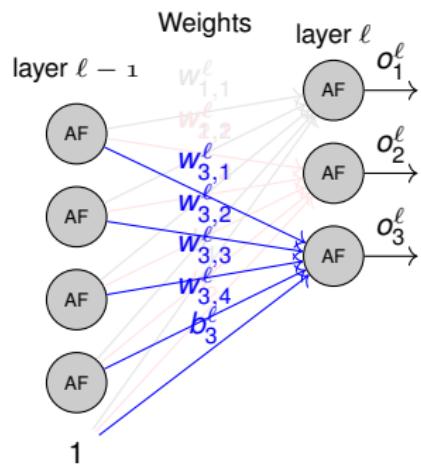
$$z_1^\ell = \sum_i w_{1,i}^\ell o_i^{\ell-1} + b_1^\ell = \mathbf{w}_1^\ell T \mathbf{o}^{\ell-1} + b_1^\ell, \quad o_1^\ell = AF(z_1^\ell)$$

$$z_1^\ell = \begin{bmatrix} w_{1,1}^\ell & w_{1,2}^\ell & w_{1,3}^\ell & w_{1,4}^\ell \end{bmatrix} \begin{bmatrix} o_1^{\ell-1} \\ o_2^{\ell-1} \\ o_3^{\ell-1} \\ o_4^{\ell-1} \end{bmatrix} + b_1^\ell$$

# Fully connected neural networks



# Fully connected neural networks



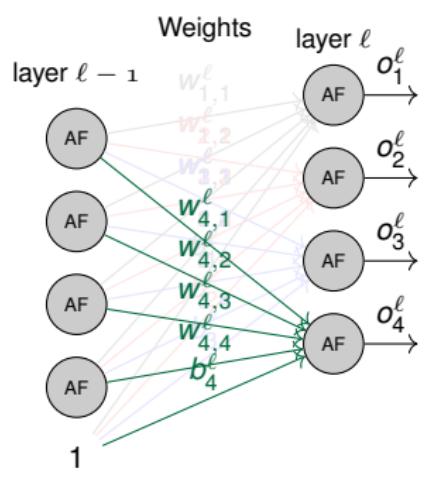
$$z_1^\ell = \sum_i w_{1,i}^\ell o_i^{\ell-1} + b_1^\ell = \mathbf{w}_1^\ell T \mathbf{o}^{\ell-1} + b_1^\ell, \quad o_1^\ell = AF(z_1^\ell)$$

$$z_2^\ell = \sum_i w_{2,i}^\ell o_i^{\ell-1} + b_2^\ell = \mathbf{w}_2^\ell T \mathbf{o}^{\ell-1} + b_2^\ell, \quad o_2^\ell = AF(z_2^\ell)$$

$$z_3^\ell = \sum_i w_{3,i}^\ell o_i^{\ell-1} + b_3^\ell = \mathbf{w}_3^\ell T \mathbf{o}^{\ell-1} + b_3^\ell, \quad o_3^\ell = AF(z_3^\ell)$$

$$\begin{bmatrix} z_1^\ell \\ z_2^\ell \\ z_3^\ell \end{bmatrix} = \begin{bmatrix} w_{1,1}^\ell & w_{1,2}^\ell & w_{1,3}^\ell & w_{1,4}^\ell \\ w_{2,1}^\ell & w_{2,2}^\ell & w_{2,3}^\ell & w_{2,4}^\ell \\ w_{3,1}^\ell & w_{3,2}^\ell & w_{3,3}^\ell & w_{3,4}^\ell \end{bmatrix} \begin{bmatrix} o_1^{\ell-1} \\ o_2^{\ell-1} \\ o_3^{\ell-1} \\ o_4^{\ell-1} \end{bmatrix} + \begin{bmatrix} b_1^\ell \\ b_2^\ell \\ b_3^\ell \end{bmatrix}$$

# Fully connected neural networks



$$z_1^\ell = \sum_i w_{1,i}^\ell o_i^{\ell-1} + b_1^\ell = \mathbf{w}_1^\ell T \mathbf{o}^{\ell-1} + b_1^\ell, \quad o_1^\ell = AF(z_1^\ell)$$

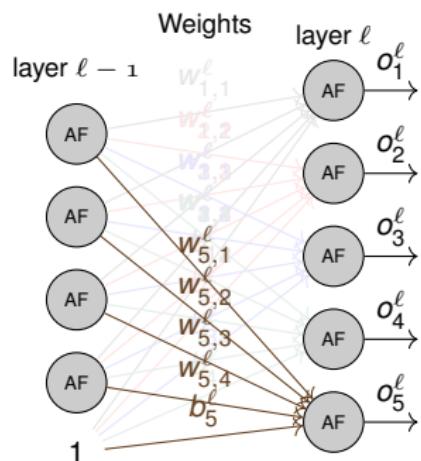
$$z_2^\ell = \sum_i w_{2,i}^\ell o_i^{\ell-1} + b_2^\ell = \mathbf{w}_2^\ell T \mathbf{o}^{\ell-1} + b_2^\ell, \quad o_2^\ell = AF(z_2^\ell)$$

$$z_3^\ell = \sum_i w_{3,i}^\ell o_i^{\ell-1} + b_3^\ell = \mathbf{w}_3^\ell T \mathbf{o}^{\ell-1} + b_3^\ell, \quad o_3^\ell = AF(z_3^\ell)$$

$$z_4^\ell = \sum_i w_{4,i}^\ell o_i^{\ell-1} + b_4^\ell = \mathbf{w}_4^\ell T \mathbf{o}^{\ell-1} + b_4^\ell, \quad o_4^\ell = AF(z_4^\ell)$$

$$\begin{bmatrix} z_1^\ell \\ z_2^\ell \\ z_3^\ell \\ z_4^\ell \end{bmatrix} = \begin{bmatrix} w_{1,1}^\ell & w_{1,2}^\ell & w_{1,3}^\ell & w_{1,4}^\ell \\ w_{2,1}^\ell & w_{2,2}^\ell & w_{2,3}^\ell & w_{2,4}^\ell \\ w_{3,1}^\ell & w_{3,2}^\ell & w_{3,3}^\ell & w_{3,4}^\ell \\ w_{4,1}^\ell & w_{4,2}^\ell & w_{4,3}^\ell & w_{4,4}^\ell \end{bmatrix} \begin{bmatrix} o_1^{\ell-1} \\ o_2^{\ell-1} \\ o_3^{\ell-1} \\ o_4^{\ell-1} \end{bmatrix} + \begin{bmatrix} b_1^\ell \\ b_2^\ell \\ b_3^\ell \\ b_4^\ell \end{bmatrix}$$

# Fully connected neural networks



$$z_1^\ell = \sum_i w_{1,i}^\ell o_i^{\ell-1} + b_1^\ell = \mathbf{w}_1^\ell T \mathbf{o}^{\ell-1} + b_1^\ell, \quad o_1^\ell = AF(z_1^\ell)$$

$$z_2^\ell = \sum_i w_{2,i}^\ell o_i^{\ell-1} + b_2^\ell = \mathbf{w}_2^\ell T \mathbf{o}^{\ell-1} + b_2^\ell, \quad o_2^\ell = AF(z_2^\ell)$$

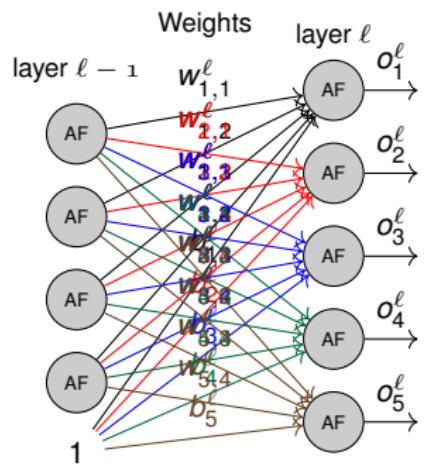
$$z_3^\ell = \sum_i w_{3,i}^\ell o_i^{\ell-1} + b_3^\ell = \mathbf{w}_3^\ell T \mathbf{o}^{\ell-1} + b_3^\ell, \quad o_3^\ell = AF(z_3^\ell)$$

$$z_4^\ell = \sum_i w_{4,i}^\ell o_i^{\ell-1} + b_4^\ell = \mathbf{w}_4^\ell T \mathbf{o}^{\ell-1} + b_4^\ell, \quad o_4^\ell = AF(z_4^\ell)$$

$$z_5^\ell = \sum_i w_{5,i}^\ell o_i^{\ell-1} + b_5^\ell = \mathbf{w}_5^\ell T \mathbf{o}^{\ell-1} + b_5^\ell, \quad o_5^\ell = AF(z_5^\ell)$$

$$\begin{bmatrix} z_1^\ell \\ z_2^\ell \\ z_3^\ell \\ z_4^\ell \\ z_5^\ell \end{bmatrix} = \begin{bmatrix} w_{1,1}^\ell & w_{1,2}^\ell & w_{1,3}^\ell & w_{1,4}^\ell \\ w_{2,1}^\ell & w_{2,2}^\ell & w_{2,3}^\ell & w_{2,4}^\ell \\ w_{3,1}^\ell & w_{3,2}^\ell & w_{3,3}^\ell & w_{3,4}^\ell \\ w_{4,1}^\ell & w_{4,2}^\ell & w_{4,3}^\ell & w_{4,4}^\ell \\ w_{5,1}^\ell & w_{5,2}^\ell & w_{5,3}^\ell & w_{5,4}^\ell \end{bmatrix} \begin{bmatrix} o_1^{\ell-1} \\ o_2^{\ell-1} \\ o_3^{\ell-1} \\ o_4^{\ell-1} \end{bmatrix} + \begin{bmatrix} b_1^\ell \\ b_2^\ell \\ b_3^\ell \\ b_4^\ell \\ b_5^\ell \end{bmatrix}$$

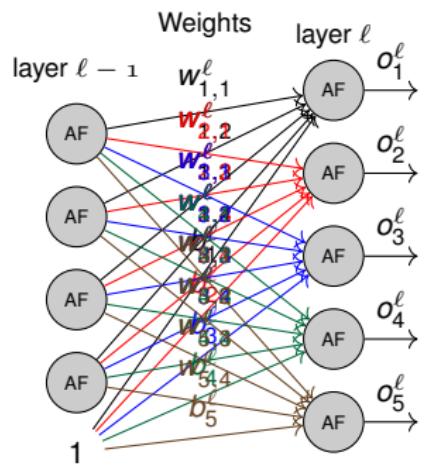
# Vector activation



$$\begin{bmatrix} z_1^\ell \\ z_2^\ell \\ z_3^\ell \\ z_4^\ell \\ z_5^\ell \end{bmatrix} = \begin{bmatrix} w_{1,1}^\ell & w_{1,2}^\ell & w_{1,3}^\ell & w_{1,4}^\ell \\ w_{2,1}^\ell & w_{2,2}^\ell & w_{2,3}^\ell & w_{2,4}^\ell \\ w_{3,1}^\ell & w_{3,2}^\ell & w_{3,3}^\ell & w_{3,4}^\ell \\ w_{4,1}^\ell & w_{4,2}^\ell & w_{4,3}^\ell & w_{4,4}^\ell \\ w_{5,1}^\ell & w_{5,2}^\ell & w_{5,3}^\ell & w_{5,4}^\ell \end{bmatrix} \begin{bmatrix} o_1^{\ell-1} \\ o_2^{\ell-1} \\ o_3^{\ell-1} \\ o_4^{\ell-1} \end{bmatrix} + \begin{bmatrix} b_1^\ell \\ b_2^\ell \\ b_3^\ell \\ b_4^\ell \\ b_5^\ell \end{bmatrix}$$

$$\begin{bmatrix} o_1^\ell \\ o_2^\ell \\ o_3^\ell \\ o_4^\ell \\ o_5^\ell \end{bmatrix} = \begin{bmatrix} AF(z_1^\ell) \\ AF(z_2^\ell) \\ AF(z_3^\ell) \\ AF(z_4^\ell) \\ AF(z_5^\ell) \end{bmatrix}$$

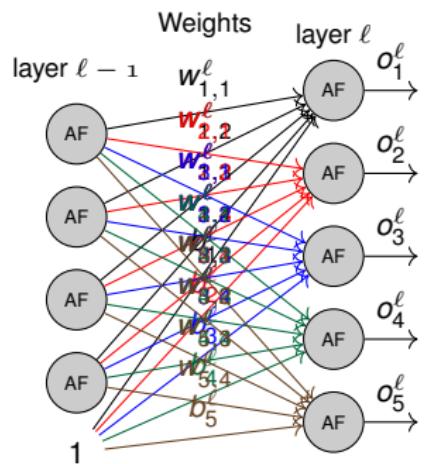
# Vector activation



$$\begin{bmatrix} z_1^\ell \\ z_2^\ell \\ z_3^\ell \\ z_4^\ell \\ z_5^\ell \end{bmatrix} = \begin{bmatrix} w_{1,1}^\ell & w_{1,2}^\ell & w_{1,3}^\ell & w_{1,4}^\ell \\ w_{2,1}^\ell & w_{2,2}^\ell & w_{2,3}^\ell & w_{2,4}^\ell \\ w_{3,1}^\ell & w_{3,2}^\ell & w_{3,3}^\ell & w_{3,4}^\ell \\ w_{4,1}^\ell & w_{4,2}^\ell & w_{4,3}^\ell & w_{4,4}^\ell \\ w_{5,1}^\ell & w_{5,2}^\ell & w_{5,3}^\ell & w_{5,4}^\ell \end{bmatrix} \begin{bmatrix} o_1^{\ell-1} \\ o_2^{\ell-1} \\ o_3^{\ell-1} \\ o_4^{\ell-1} \end{bmatrix} + \begin{bmatrix} b_1^\ell \\ b_2^\ell \\ b_3^\ell \\ b_4^\ell \\ b_5^\ell \end{bmatrix}$$

$$\begin{bmatrix} o_1^\ell \\ o_2^\ell \\ o_3^\ell \\ o_4^\ell \\ o_5^\ell \end{bmatrix} = \begin{bmatrix} AF(z_1^\ell) \\ AF(z_2^\ell) \\ AF(z_3^\ell) \\ AF(z_4^\ell) \\ AF(z_5^\ell) \end{bmatrix} = AF(\begin{bmatrix} z_1^\ell \\ z_2^\ell \\ z_3^\ell \\ z_4^\ell \\ z_5^\ell \end{bmatrix})$$

# Vector activation

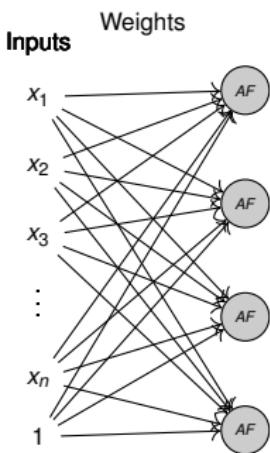


$$\begin{bmatrix} z_1^\ell \\ z_2^\ell \\ z_3^\ell \\ z_4^\ell \\ z_5^\ell \end{bmatrix} = \begin{bmatrix} w_{1,1}^\ell & w_{1,2}^\ell & w_{1,3}^\ell & w_{1,4}^\ell & w_{1,5}^\ell \\ w_{2,1}^\ell & w_{2,2}^\ell & w_{2,3}^\ell & w_{2,4}^\ell & w_{2,5}^\ell \\ w_{3,1}^\ell & w_{3,2}^\ell & w_{3,3}^\ell & w_{3,4}^\ell & w_{3,5}^\ell \\ w_{4,1}^\ell & w_{4,2}^\ell & w_{4,3}^\ell & w_{4,4}^\ell & w_{4,5}^\ell \\ w_{5,1}^\ell & w_{5,2}^\ell & w_{5,3}^\ell & w_{5,4}^\ell & w_{5,5}^\ell \end{bmatrix} \begin{bmatrix} o_1^{\ell-1} \\ o_2^{\ell-1} \\ o_3^{\ell-1} \\ o_4^{\ell-1} \\ o_5^{\ell-1} \end{bmatrix} + \begin{bmatrix} b_1^\ell \\ b_2^\ell \\ b_3^\ell \\ b_4^\ell \\ b_5^\ell \end{bmatrix}$$

$$\begin{bmatrix} o_1^\ell \\ o_2^\ell \\ o_3^\ell \\ o_4^\ell \\ o_5^\ell \end{bmatrix} = \begin{bmatrix} AF(z_1^\ell) \\ AF(z_2^\ell) \\ AF(z_3^\ell) \\ AF(z_4^\ell) \\ AF(z_5^\ell) \end{bmatrix} = AF(\begin{bmatrix} z_1^\ell \\ z_2^\ell \\ z_3^\ell \\ z_4^\ell \\ z_5^\ell \end{bmatrix})$$

$$\mathbf{o}^\ell = AF(W^\ell \mathbf{o}^{\ell-1} + \mathbf{b}^\ell)$$

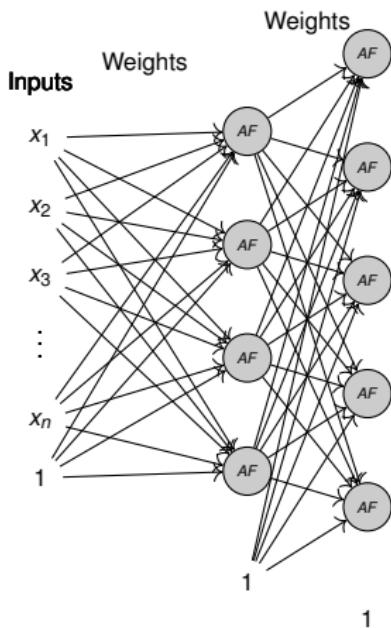
# Fully connected neural networks



1

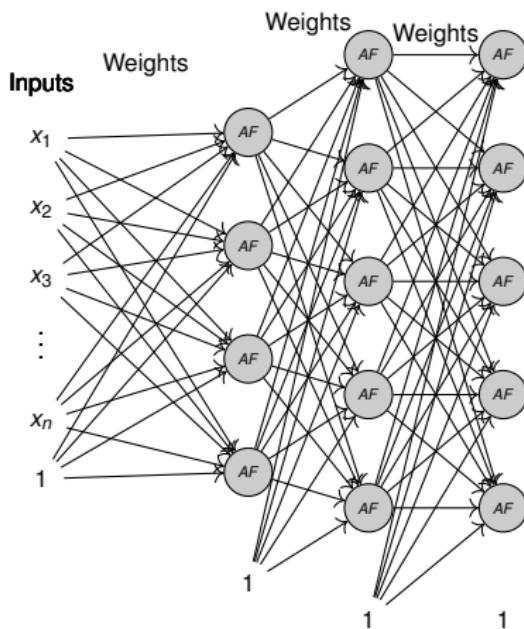
$$\mathbf{o}^1 = AF(W^1 \mathbf{x} + \mathbf{b}^1), \quad W^1 \in \mathbb{R}^{4 \times n}, \quad \mathbf{x} \in \mathbb{R}^n, \quad \mathbf{o}^1 \in \mathbb{R}^4, \quad \mathbf{b} \in \mathbb{R}^4.$$

# Fully connected neural networks



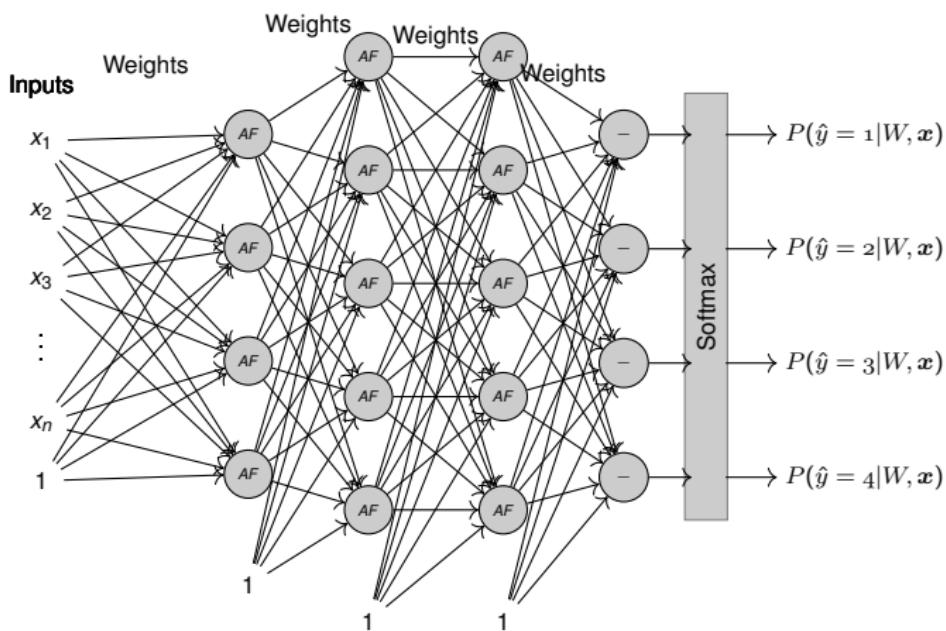
$$\mathbf{o}^2 = AF(W^2\mathbf{o}^1 + \mathbf{b}^2), \quad W^2 \in \mathbb{R}^{5 \times 4}, \quad \mathbf{o}^1 \in \mathbb{R}^4, \quad \mathbf{o}^2 \in \mathbb{R}^5, \quad \mathbf{b}^2 \in \mathbb{R}^5.$$

# Fully connected neural networks



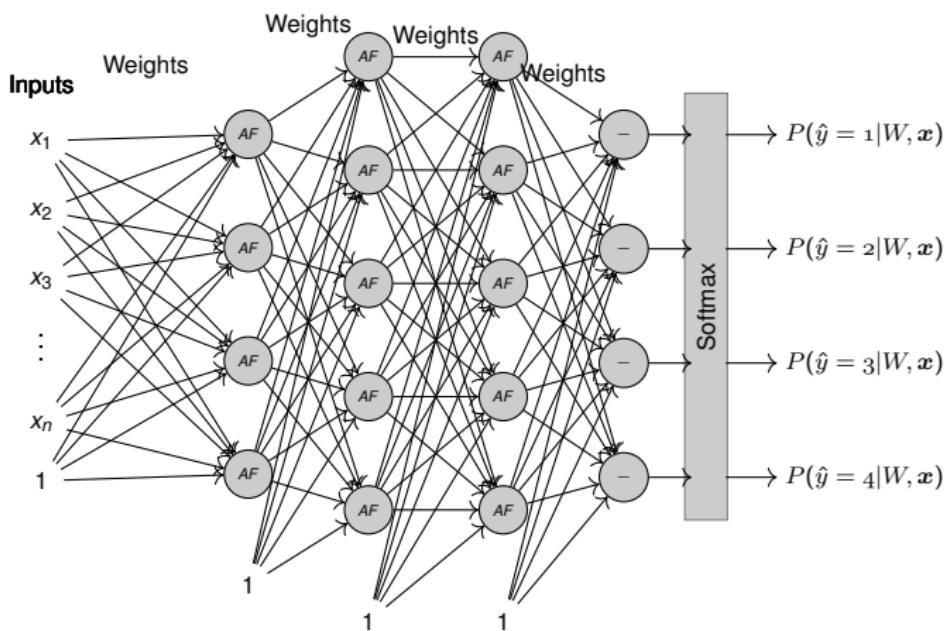
$$\mathbf{o}^3 = AF(W^3\mathbf{o}^2 + \mathbf{b}^3), \quad W^3 \in \mathbb{R}^{5 \times 5}, \quad \mathbf{o}^2 \in \mathbb{R}^5, \quad \mathbf{o}^3 \in \mathbb{R}^5, \quad \mathbf{b}^3 \in \mathbb{R}^5.$$

# Fully connected neural networks



$$\hat{\mathbf{y}} = \text{Softmax}(W^4 \mathbf{o}^3 + \mathbf{b}^4), \quad W^4 \in \mathbb{R}^{4 \times 5}, \quad \mathbf{o}^3 \in \mathbb{R}^5, \quad \hat{\mathbf{y}} \in \mathbb{R}^4, \quad \mathbf{b}^4 \in \mathbb{R}^4.$$

# Fully connected neural networks



$$\hat{\mathbf{y}} = \text{Softmax}(W^4(\text{AF}(W^3(\text{AF}(W^2(\text{AF}(W^1 \mathbf{x} + \mathbf{b}^1)) + \mathbf{b}^2)) + \mathbf{b}^3)) + \mathbf{b}^4)$$

# Batch of inputs

- The output of the first hidden layer for a batch of inputs

$$\begin{bmatrix} z_1^{1,(1)} & z_1^{1,(2)} & z_1^{1,(3)} \\ z_2^{1,(1)} & z_2^{1,(2)} & z_2^{1,(3)} \\ z_3^{1,(1)} & z_3^{1,(2)} & z_3^{1,(3)} \\ z_4^{1,(1)} & z_4^{1,(2)} & z_4^{1,(3)} \end{bmatrix} = \begin{bmatrix} w_{1,1}^1 & w_{1,2}^1 & w_{1,3}^1 & \dots & w_{1,n}^1 \\ w_{2,1}^1 & w_{2,2}^1 & w_{2,3}^1 & \dots & w_{2,n}^1 \\ w_{3,1}^1 & w_{3,2}^1 & w_{3,3}^1 & \dots & w_{3,n}^1 \\ w_{4,1}^1 & w_{4,2}^1 & w_{4,3}^1 & \dots & w_{4,n}^1 \end{bmatrix} \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & x_1^{(3)} \\ x_2^{(1)} & x_2^{(2)} & x_2^{(3)} \\ x_3^{(1)} & x_3^{(2)} & x_3^{(3)} \\ \vdots & \vdots & \vdots \\ x_n^{(1)} & x_n^{(2)} & x_n^{(3)} \end{bmatrix}$$

$$+ \begin{bmatrix} b_1^1 & b_2^1 & b_3^1 \\ b_2^1 & b_2^1 & b_2^1 \\ b_3^1 & b_3^1 & b_3^1 \\ b_4^1 & b_4^1 & b_4^1 \end{bmatrix}$$

$$\begin{bmatrix} o_1^{1,(1)} & o_1^{1,(2)} & o_1^{1,(3)} \\ o_2^{1,(1)} & o_2^{1,(2)} & o_2^{1,(3)} \\ o_3^{1,(1)} & o_3^{1,(2)} & o_3^{1,(3)} \\ o_4^{1,(1)} & o_4^{1,(2)} & o_4^{1,(3)} \end{bmatrix} = AF\left(\begin{bmatrix} z_1^{1,(1)} & z_1^{1,(2)} & z_1^{1,(3)} \\ z_2^{1,(1)} & z_2^{1,(2)} & z_2^{1,(3)} \\ z_3^{1,(1)} & z_3^{1,(2)} & z_3^{1,(3)} \\ z_4^{1,(1)} & z_4^{1,(2)} & z_4^{1,(3)} \end{bmatrix}\right)$$

## Batch of inputs

- The output of the first hidden layer for a batch of inputs

$$\begin{bmatrix} z_1^{1,(1)} & z_1^{1,(2)} & z_1^{1,(3)} \\ z_2^{1,(1)} & z_2^{1,(2)} & z_2^{1,(3)} \\ z_3^{1,(1)} & z_3^{1,(2)} & z_3^{1,(3)} \\ z_4^{1,(1)} & z_4^{1,(2)} & z_4^{1,(3)} \end{bmatrix} = \begin{bmatrix} w_{1,1}^1 & w_{1,2}^1 & w_{1,3}^1 & \dots & w_{1,n}^1 \\ w_{2,1}^1 & w_{2,2}^1 & w_{2,3}^1 & \dots & w_{2,n}^1 \\ w_{3,1}^1 & w_{3,2}^1 & w_{3,3}^1 & \dots & w_{3,n}^1 \\ w_{4,1}^1 & w_{4,2}^1 & w_{4,3}^1 & \dots & w_{4,n}^1 \end{bmatrix} \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & x_1^{(3)} \\ x_2^{(1)} & x_2^{(2)} & x_2^{(3)} \\ x_3^{(1)} & x_3^{(2)} & x_3^{(3)} \\ \vdots & \vdots & \vdots \\ x_n^{(1)} & x_n^{(2)} & x_n^{(3)} \end{bmatrix}$$

$$+ \begin{bmatrix} b_1^1 & b_1^1 & b_1^1 \\ b_2^1 & b_2^1 & b_2^1 \\ b_3^1 & b_3^1 & b_3^1 \\ b_4^1 & b_4^1 & b_4^1 \end{bmatrix}$$

$$O^1 = AF(W^1 X + B^1), \quad W^1 \in \mathbb{R}^{4 \times n}, \quad X \in \mathbb{R}^{n \times 3}, \quad O^1 \in \mathbb{R}^{4 \times 3}, \quad B \in \mathbb{R}^{4 \times 3}.$$

# Activation functions

What happens if we build a neural network with no activation function?

$$\hat{y} = \text{Softmax}(W^4(\text{AF}(W^3(\text{AF}(W^2(\text{AF}(W^1 \mathbf{x} + \mathbf{b}^1)) + \mathbf{b}^2)) + \mathbf{b}^3)) + \mathbf{b}^4)$$

# Activation functions

What happens if we build a neural network with no activation function?

$$\hat{y} = \text{Softmax}(W^4(\text{AF}(W^3(\text{AF}(W^2(\text{AF}(W^1 \mathbf{x} + \mathbf{b}^1)) + \mathbf{b}^2)) + \mathbf{b}^3)) + \mathbf{b}^4)$$

$$\hat{y} = \text{Softmax}(W^4(W^3(W^2(W^1 \mathbf{x} + \mathbf{b}^1) + \mathbf{b}^2) + \mathbf{b}^3) + \mathbf{b}^4)$$

# Activation functions

What happens if we build a neural network with no activation function?

$$\hat{y} = \text{Softmax}(W^4(\text{AF}(W^3(\text{AF}(W^2(\text{AF}(W^1 \mathbf{x} + \mathbf{b}^1)) + \mathbf{b}^2)) + \mathbf{b}^3)) + \mathbf{b}^4)$$

$$\hat{y} = \text{Softmax}(W^4(W^3(W^2(W^1 \mathbf{x} + \mathbf{b}^1) + \mathbf{b}^2) + \mathbf{b}^3) + \mathbf{b}^4)$$

$$\hat{y} = \text{Softmax}(W^4 W^3 W^2 W^1 \mathbf{x} + b')$$

## Activation functions

What happens if we build a neural network with no activation function?

$$\hat{y} = \text{Softmax}(W^4(\text{AF}(W^3(\text{AF}(W^2(\text{AF}(W^1 x + b^1)) + b^2)) + b^3)) + b^4)$$

$$\hat{y} = \text{Softmax}(W^4(W^3(W^2(W^1x + b^1) + b^2) + b^3) + b^4)$$

$$\hat{y} = \text{Softmax}(W^4 W^3 W^2 W^1 x + b')$$

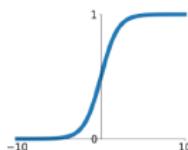
$$\hat{y} = \text{Softmax}(W'x + b')$$

We end up with a softmax classifier!

# Activation functions

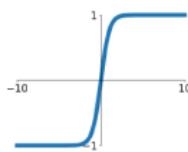
## Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



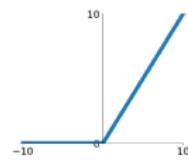
## tanh

$$\tanh(x)$$



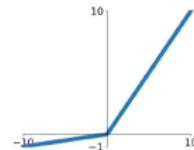
## ReLU

$$\max(0, x)$$



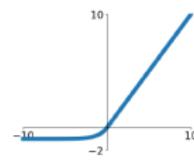
## Leaky ReLU

$$\max(0.1x, x)$$



## ELU

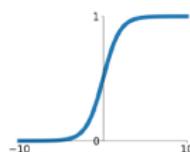
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



# Activation functions

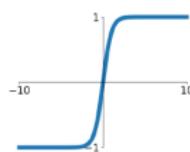
**Sigmoid**

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



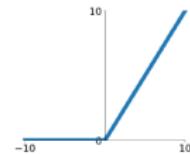
**tanh**

$$\tanh(x)$$



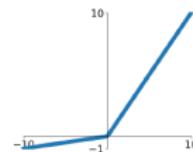
**ReLU**

$$\max(0, x)$$



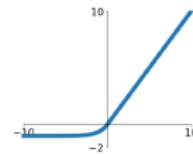
**Leaky ReLU**

$$\max(0.1x, x)$$



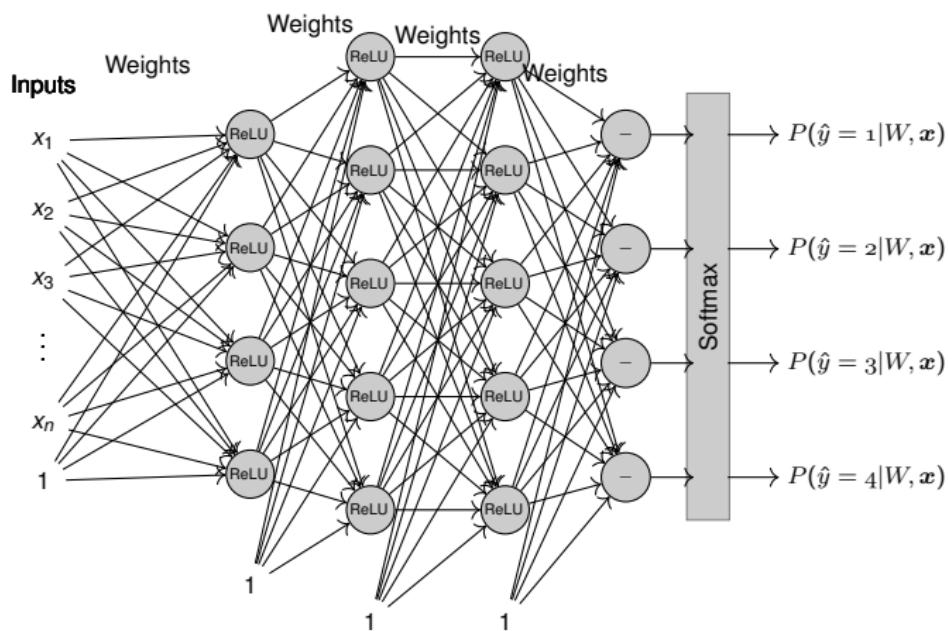
**ELU**

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



**ReLU is a good default choice for most problems**

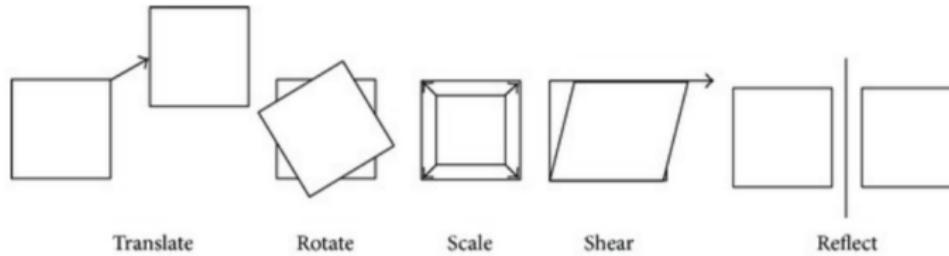
# ReLU activation function



$$\hat{\mathbf{y}} = \text{Softmax}(W^4(\max(0, W^3(\max(0, W^2(\max(0, W^1\mathbf{x} + \mathbf{b}^1)) + \mathbf{b}^2)) + \mathbf{b}^3)) + \mathbf{b}^4)$$

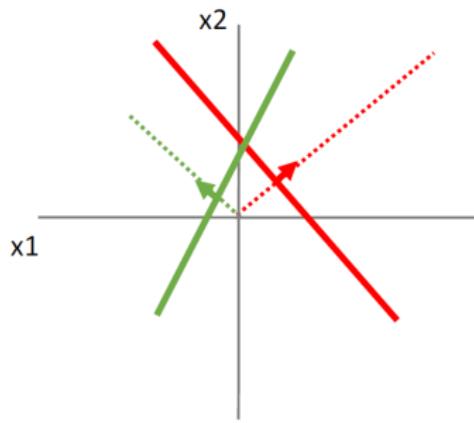
## Affine transformation

- Affine transformation
    - Translate
    - Rotate
    - Scale
    - Shear
    - Reflect
  - Preserve parallel lines



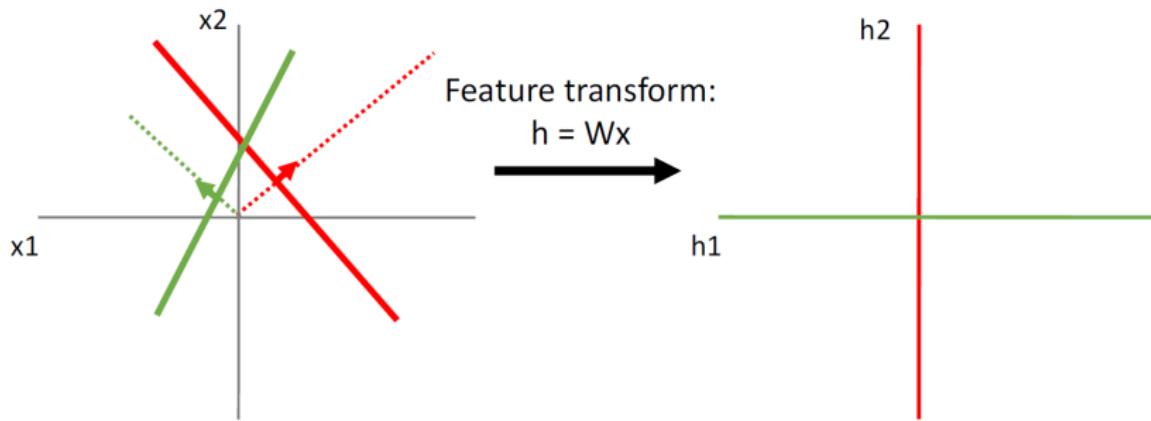
## Space warping

- Consider a linear transform:  $h = W^T x$  Where  $x, h$  are both 2-dimensional



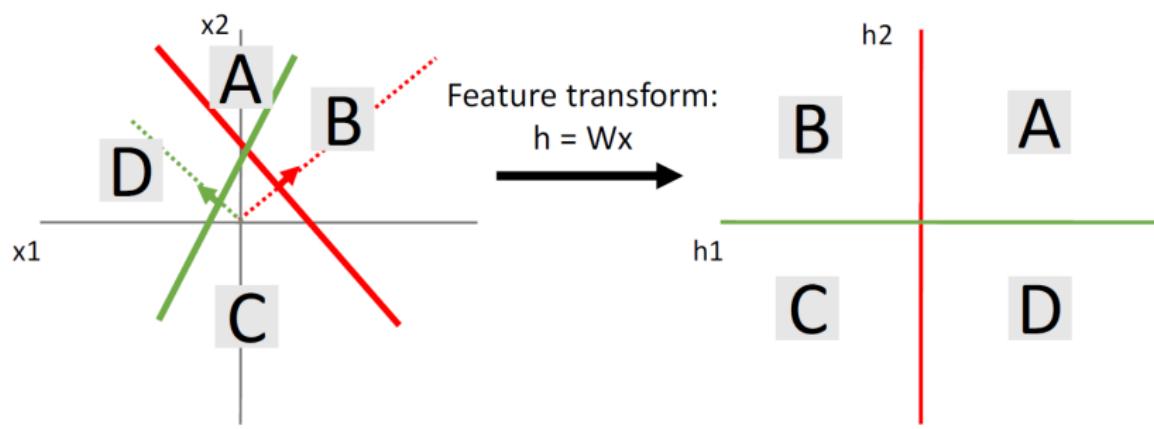
## Space warping

- Consider a linear transform:  $h = W^T x$  Where  $x, h$  are both 2-dimensional



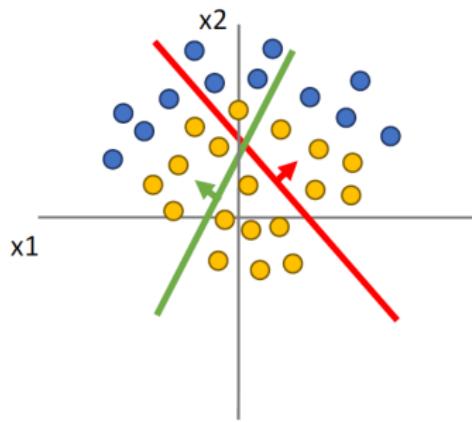
# Space warping

- Consider a linear transform:  $h = W^T x$  Where  $x, h$  are both 2-dimensional



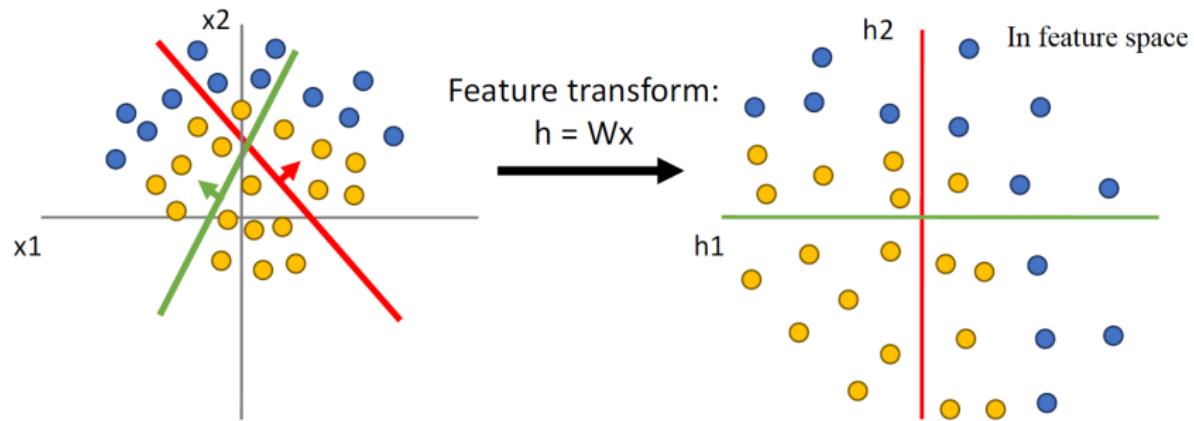
# Space warping

- Consider a linear transform:  $h = W^T x$  Where  $x, h$  are both 2-dimensional
- Points not linearly separable in original space



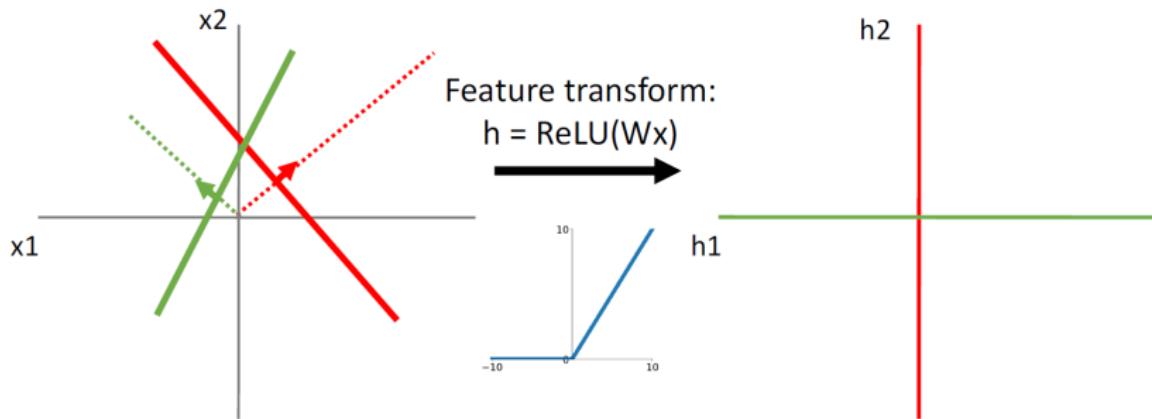
# Space warping

- Consider a linear transform:  $h = W^T x$  Where  $x, h$  are both 2-dimensional
- Points not linearly separable in original space
  - Not linearly separable in feature space



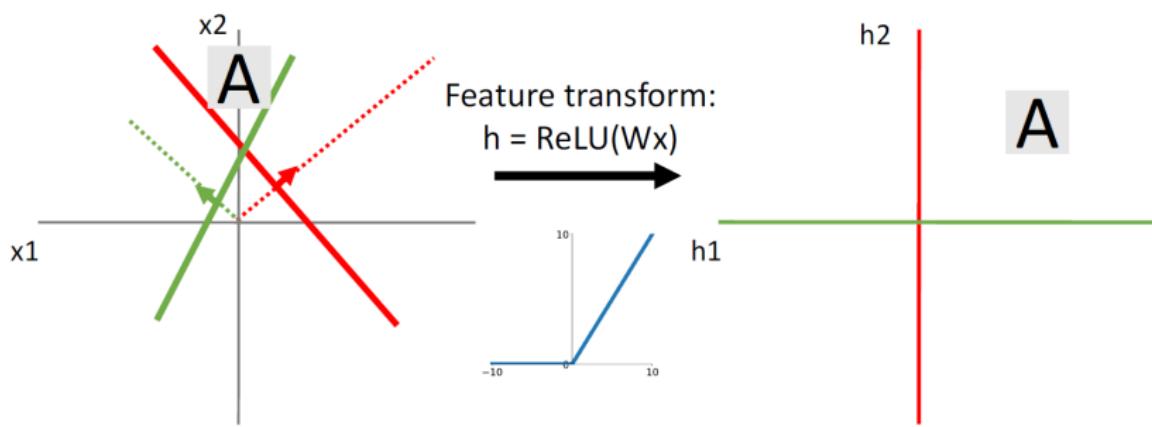
# Space warping

- Consider a neural net hidden layer:  $h = \text{ReLU}(W^T x) = \max(0, Wx)$ .
  - Where  $x, h$  are both 2-dimensional.



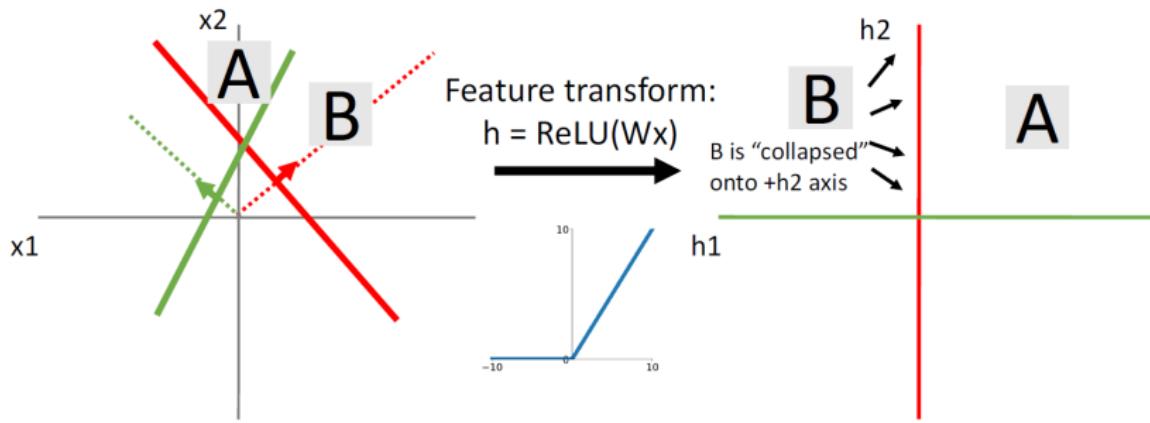
# Space warping

- Consider a neural net hidden layer:  $h = \text{ReLU}(W^T x) = \max(0, Wx)$ .
  - Where  $x, h$  are both 2-dimensional.



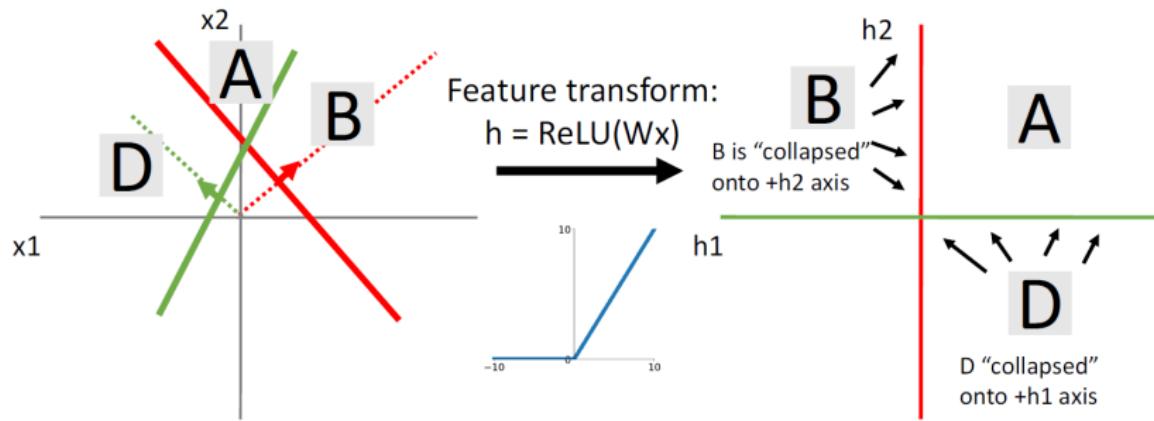
# Space warping

- Consider a neural net hidden layer:  $h = \text{ReLU}(W^T x) = \max(0, Wx)$ .
  - Where  $x, h$  are both 2-dimensional.



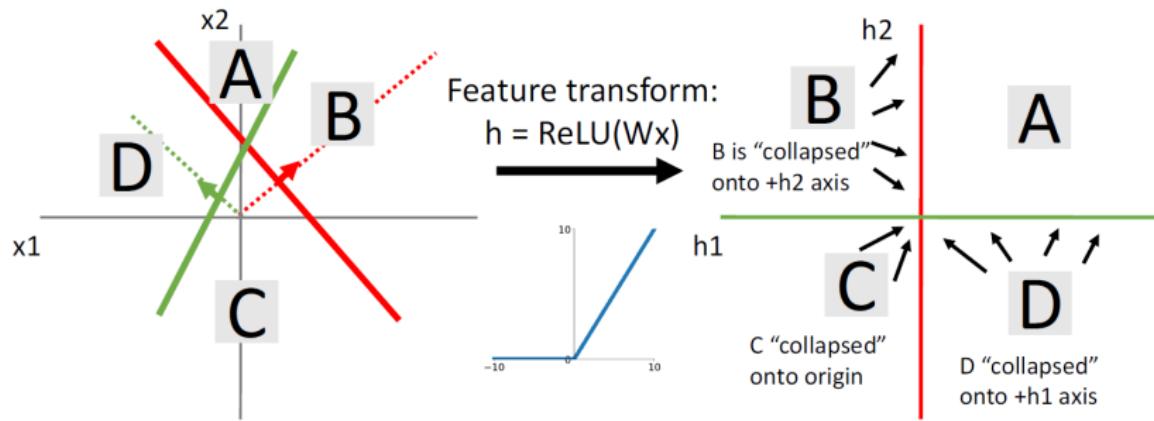
# Space warping

- Consider a neural net hidden layer:  $h = \text{ReLU}(W^T x) = \max(0, Wx)$ .
- Where  $x, h$  are both 2-dimensional.



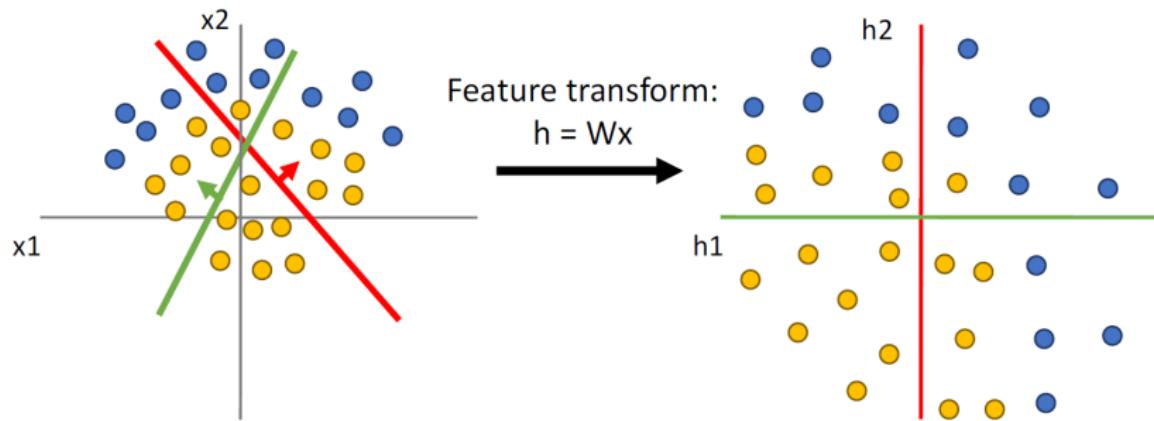
# Space warping

- Consider a neural net hidden layer:  $h = \text{ReLU}(W^T x) = \max(0, Wx)$ .
- Where  $x, h$  are both 2-dimensional.



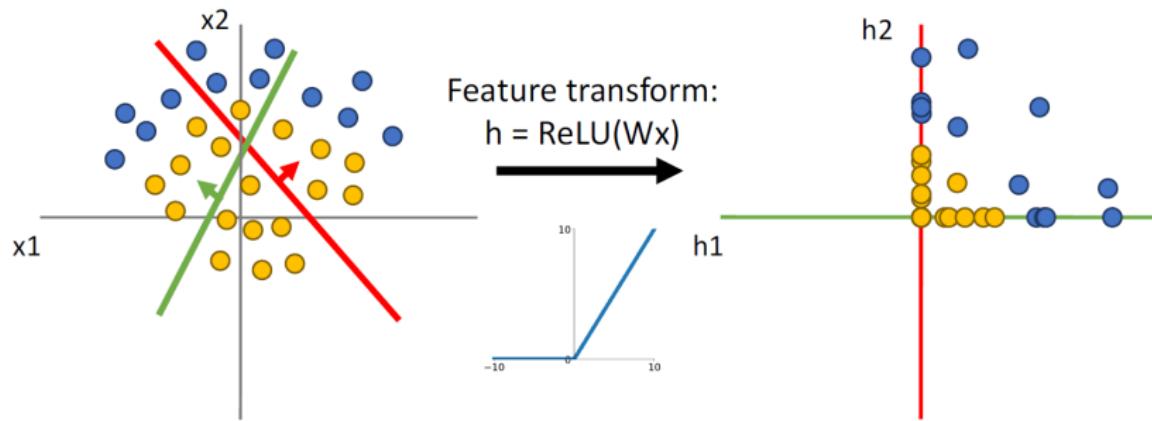
# Space warping

- Consider a neural net hidden layer:  $h = \text{ReLU}(W^T x) = \max(0, Wx)$ .
  - Where  $x, h$  are both 2-dimensional.
- Points not linearly separable in original space



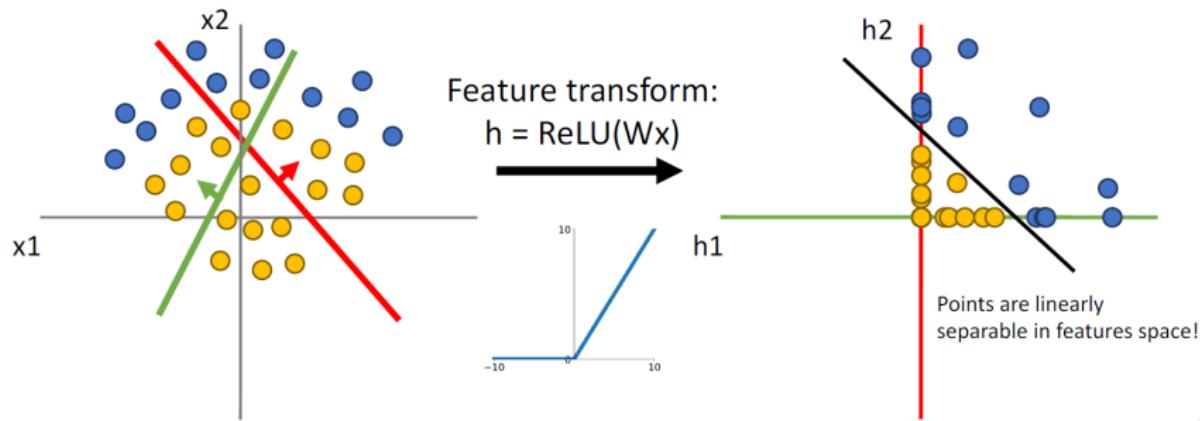
## Space warping

- Consider a neural net hidden layer:  $\mathbf{h} = \text{ReLU}(\mathbf{W}^T \mathbf{x}) = \max(0, \mathbf{W} \mathbf{x})$ .
    - Where  $\mathbf{x}$ ,  $\mathbf{h}$  are both 2-dimensional.
  - Points not linearly separable in original space



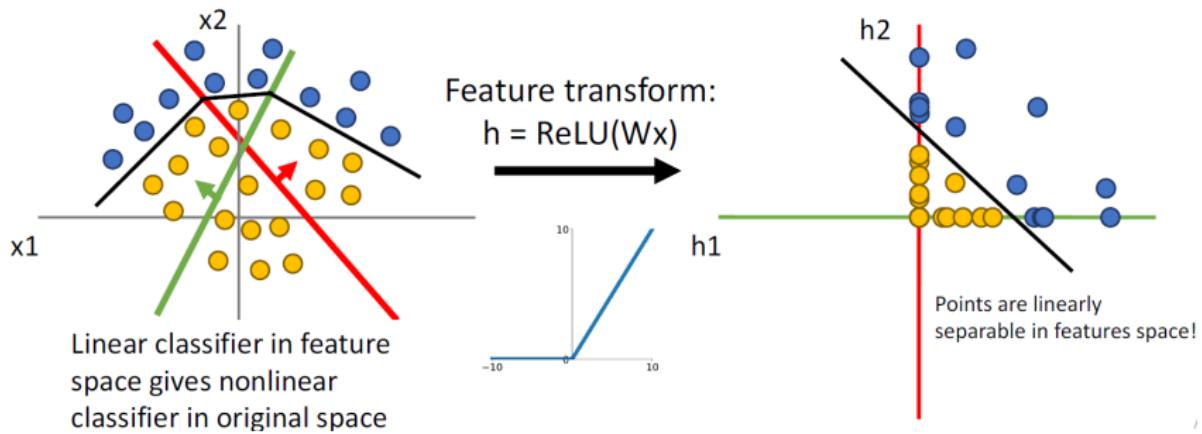
# Space warping

- Consider a neural net hidden layer:  $h = \text{ReLU}(W^T \mathbf{x}) = \max(0, W\mathbf{x})$ .
  - Where  $\mathbf{x}, h$  are both 2-dimensional.
- Points not linearly separable in original space

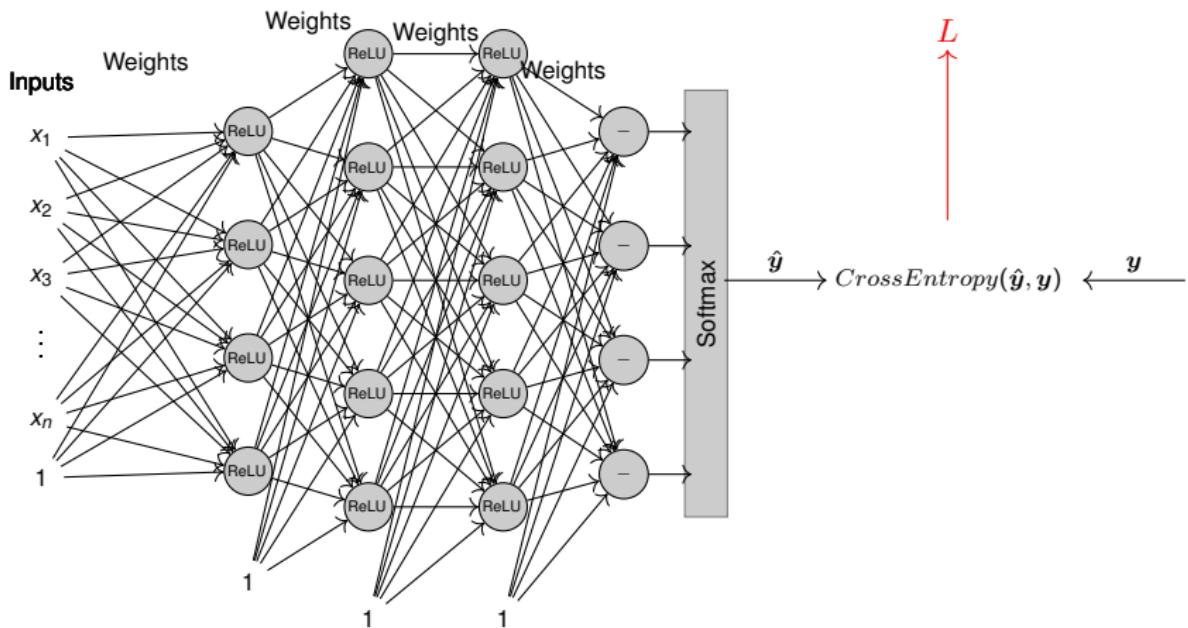


# Space warping

- Consider a neural net hidden layer:  $h = \text{ReLU}(W^T x) = \max(0, Wx)$ .
  - Where  $x, h$  are both 2-dimensional.
- Points not linearly separable in original space



# Fully connected neural network loss



$$\hat{y} = \text{Softmax}(W^4(\max(0, W^3(\max(0, W^2(\max(0, W^1 x + b^1)) + b^2)) + b^3)) + b^4)$$

# Gradient descent

- The gradient descent method to find the minimum of a function iteratively

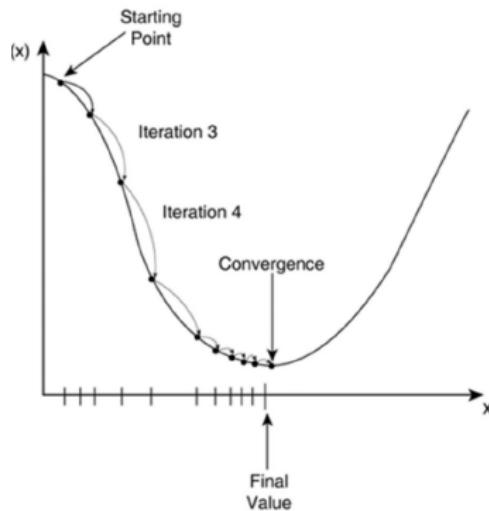
$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla_{\mathbf{x}} f(\mathbf{x}^t)$$

- $\eta$  is the "step size" (Also called "learning rate")

- The gradient descent algorithm converges when the following criterion is satisfied.

$$|f(\mathbf{x}^{t+k}) - f(\mathbf{x}^t)| < \epsilon$$

- $k$  is a hyperparameter.



## Training neural nets through gradient descent

- #### ■ Total training loss

$$L(f(X, W), \mathbf{y}) = \frac{1}{N} \sum_i L_i(f(\mathbf{x}^{(i)}, W), y^{(i)})$$

- #### ■ Gradient descent algorithm:

- Initialize all weights and biases  $w_{ij}^{(\ell)}$ 
    - Using the extended notation: the bias is also a weight

■ Do:

- For every layer  $\ell$  for all  $i, j$  update:

$$w_{ij}^{(\ell)} = w_{ij}^{(\ell)} - \eta \frac{\partial L}{\partial w_{ij}^{(\ell)}}$$

- Until loss has converged

Idea: derive  $\frac{\partial L}{\partial W}$  on paper

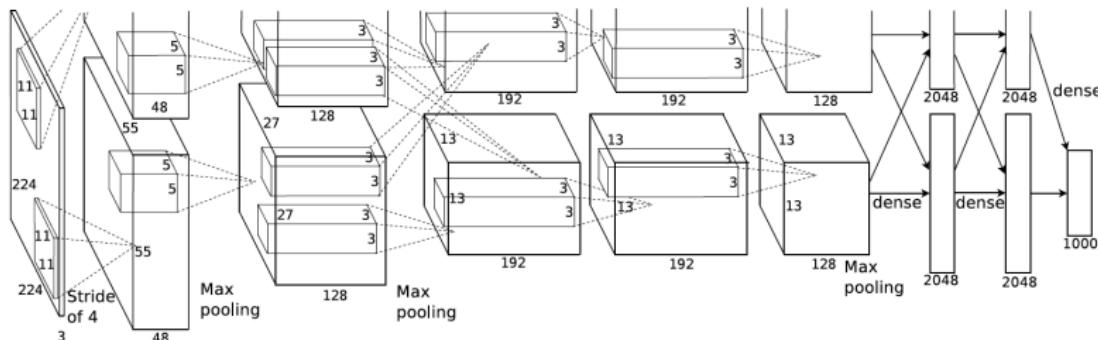
## ■ Problems

- Very tedious: Lots of matrix calculus, need lots of paper
  - What if we want to change loss? Need to re-derive from scratch. Not modular!
  - Not feasible for very complex models!

Idea: derive  $\frac{\partial L}{\partial W}$  on paper

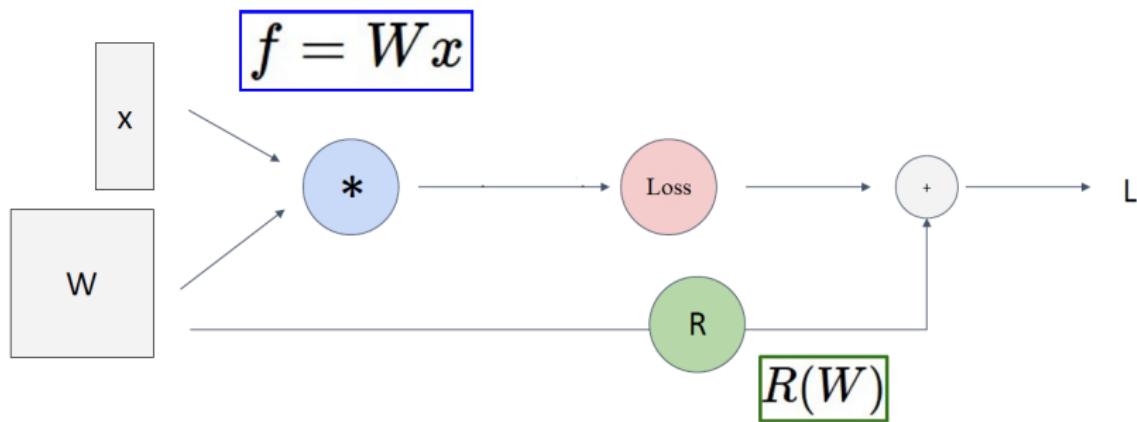
## ■ Problems

- Very tedious: Lots of matrix calculus, need lots of paper
  - What if we want to change loss? Need to re-derive from scratch. Not modular!
  - Not feasible for very complex models!



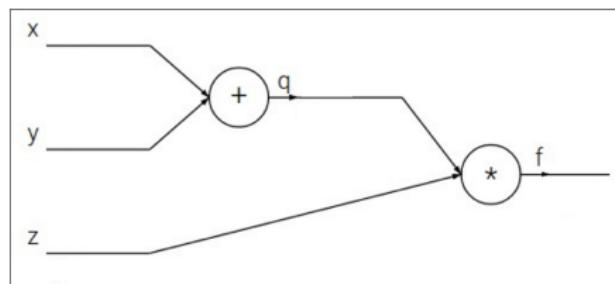
AlexNet has about 660K units, 61M parameters,

## Better Idea: Computational Graphs



## Backpropagation: simple example

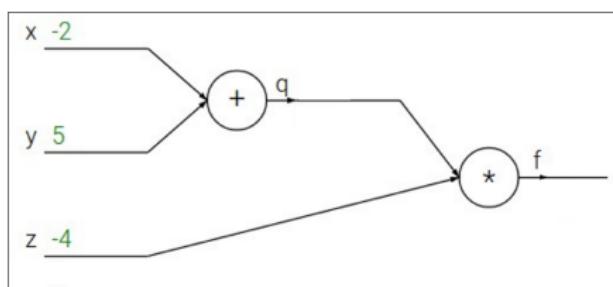
$$f(x, y, z) = (x + y)z$$



## Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2$ ,  $y = 5$ ,  $z = -4$



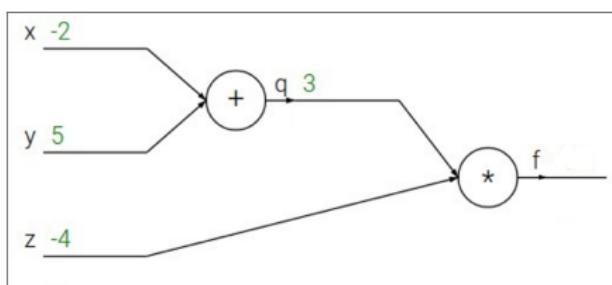
## Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2$ ,  $y = 5$ ,  $z = -4$

### 1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$



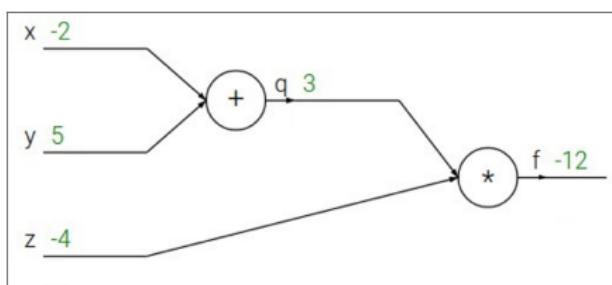
## Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2$ ,  $y = 5$ ,  $z = -4$

### 1. Forward pass: Compute outputs

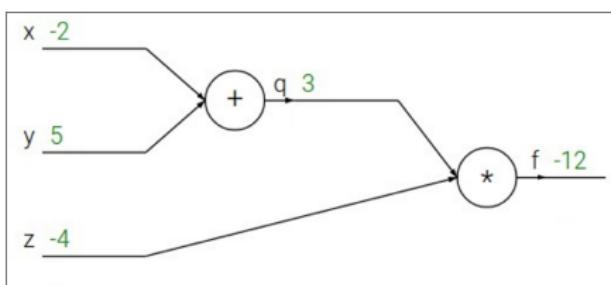
$$q = x + y \quad f = qz$$



## Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2$ ,  $y = 5$ ,  $z = -4$



- ### **1. Forward pass:** Compute outputs

$$q = x + y \quad f = qz$$

- ## 2. Backward pass: Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

## Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

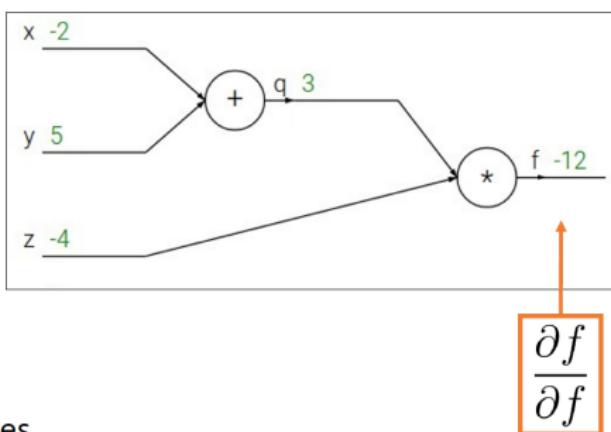
e.g.  $x = -2$ ,  $y = 5$ ,  $z = -4$

**1. Forward pass:** Compute outputs

$$q = x + y \quad f = qz$$

## 2. Backward pass: Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



# Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

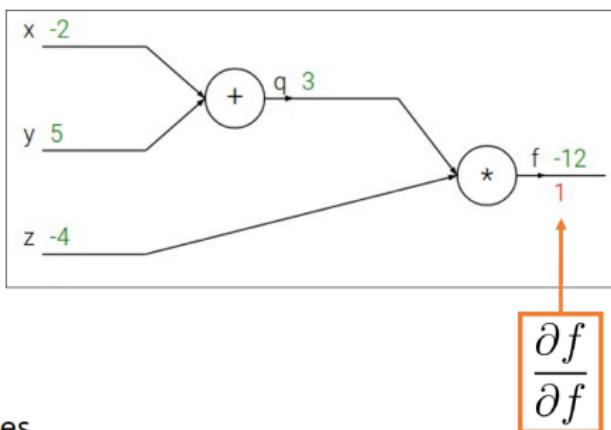
e.g.  $x = -2, y = 5, z = -4$

**1. Forward pass:** Compute outputs

$$q = x + y \quad f = qz$$

**2. Backward pass:** Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



## Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

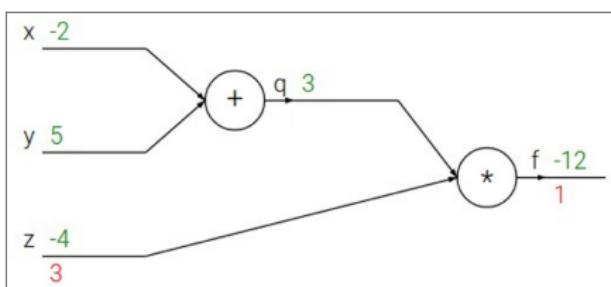
e.g.  $x = -2$ ,  $y = 5$ ,  $z = -4$

### 1. Forward pass: Compute outputs

$$q = x + y \quad | \quad f = qz$$

## 2. Backward pass: Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial z}$$

## Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

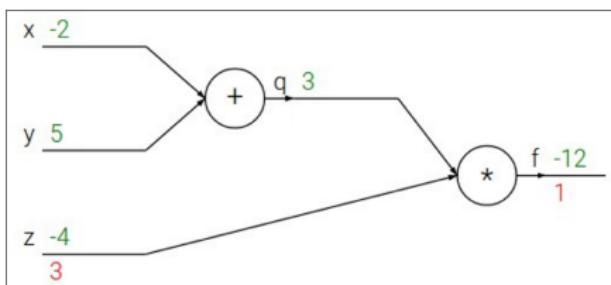
e.g.  $x = -2, y = 5, z = -4$

**1. Forward pass:** Compute outputs

$$q = x + y \quad f = qz$$

**2. Backward pass:** Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial z} = q$$

## Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

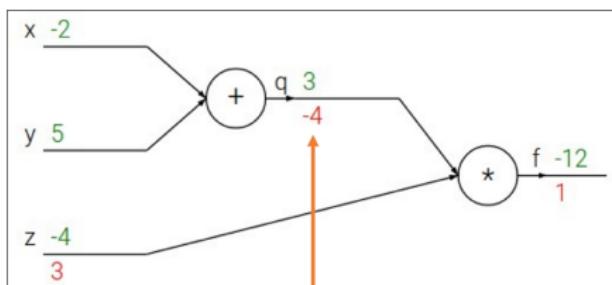
e.g.  $x = -2, y = 5, z = -4$

**1. Forward pass:** Compute outputs

$$q = x + y \quad f = qz$$

**2. Backward pass:** Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\boxed{\frac{\partial f}{\partial q}}$$

## Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

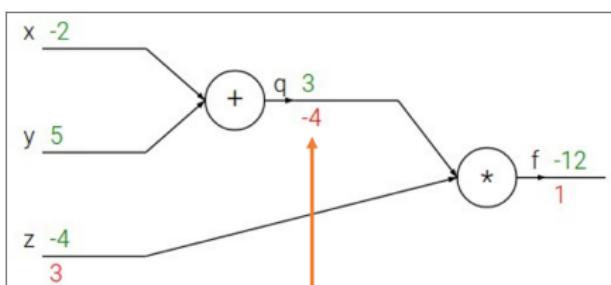
e.g.  $x = -2$ ,  $y = 5$ ,  $z = -4$

**1. Forward pass:** Compute outputs

$$q = x + y \quad | \quad f = qz$$

## 2. Backward pass: Compute derivatives

$$\frac{\partial f}{\partial q} = z$$



Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

# Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

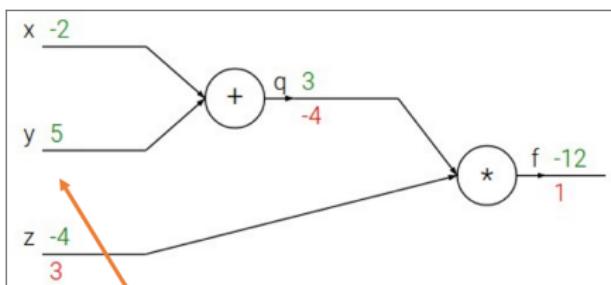
e.g.  $x = -2, y = 5, z = -4$

**1. Forward pass:** Compute outputs

$$q = x + y \quad f = qz$$

**2. Backward pass:** Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\boxed{\frac{\partial f}{\partial y}}$$

## Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

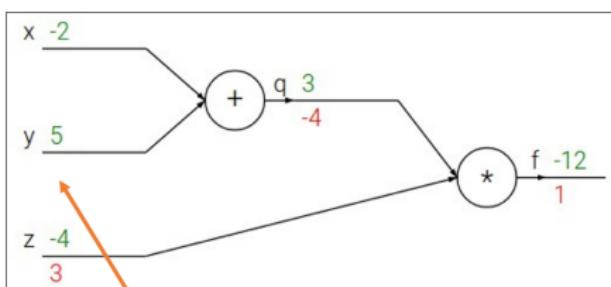
e.g.  $x = -2$ ,  $y = 5$ ,  $z = -4$

### 1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

## 2. Backward pass: Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



## Chain Rule

$$\frac{\partial f}{\partial y} = \frac{\partial q}{\partial y} \frac{\partial f}{\partial q}$$

# Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

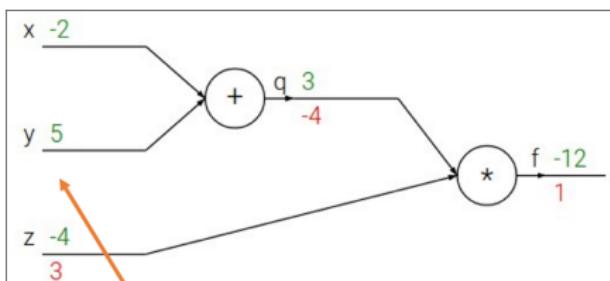
e.g.  $x = -2, y = 5, z = -4$

**1. Forward pass:** Compute outputs

$$q = x + y \quad f = qz$$

**2. Backward pass:** Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



**Chain Rule**

$$\frac{\partial f}{\partial y} = \frac{\partial q}{\partial y} \frac{\partial f}{\partial q}$$

$$\frac{\partial q}{\partial y} = 1$$

/                   ↑                   \ \\
 Downstream       Local       Upstream \\
 Gradient          Gradient      Gradient

# Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

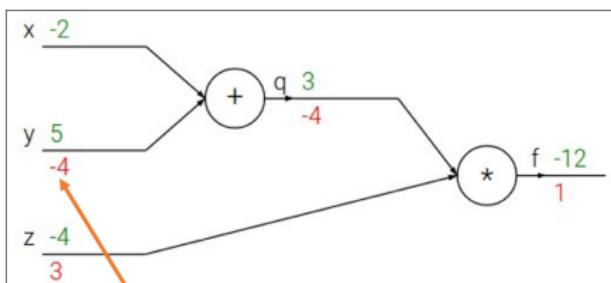
e.g.  $x = -2, y = 5, z = -4$

**1. Forward pass:** Compute outputs

$$q = x + y \quad f = qz$$

**2. Backward pass:** Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Chain Rule

$$\frac{\partial f}{\partial y} = \frac{\partial q}{\partial y} \frac{\partial f}{\partial q}$$

$$\frac{\partial q}{\partial y} = 1$$

/                   ↑                   \ \\
 Downstream       Local       Upstream \\
 Gradient          Gradient      Gradient

## Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

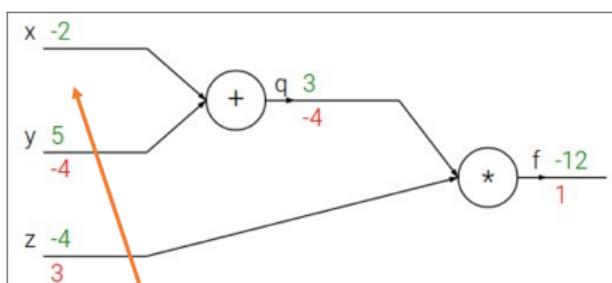
e.g.  $x = -2$ ,  $y = 5$ ,  $z = -4$

### 1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

## 2. Backward pass: Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



## Chain Rule

$$\frac{\partial f}{\partial x} = \frac{\partial q}{\partial x} \frac{\partial f}{\partial q}$$

$$\frac{\partial q}{\partial x} = 1$$

## Downstream Gradient

Local Gradient

## Upstream Gradient

## Backpropagation: simple example

$$f(x, y, z) = (x + y)z$$

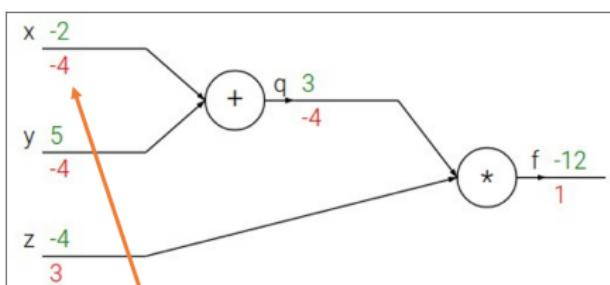
e.g.  $x = -2$ ,  $y = 5$ ,  $z = -4$

### 1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

## 2. Backward pass: Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



## Chain Rule

$$\frac{\partial f}{\partial x} = \frac{\partial q}{\partial x} \frac{\partial f}{\partial q}$$

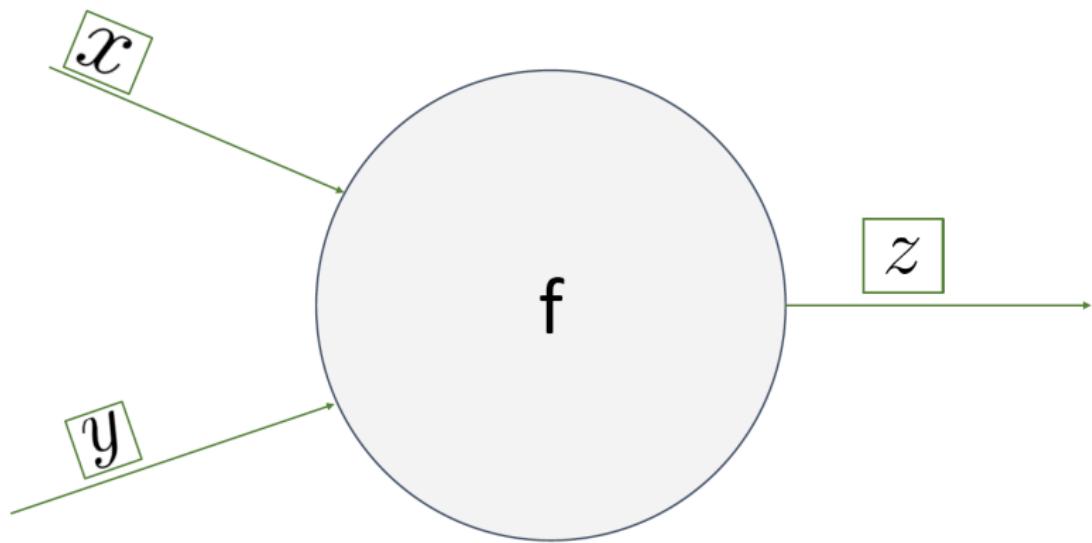
$$\frac{\partial q}{\partial x} = 1$$

## Downstream Gradient

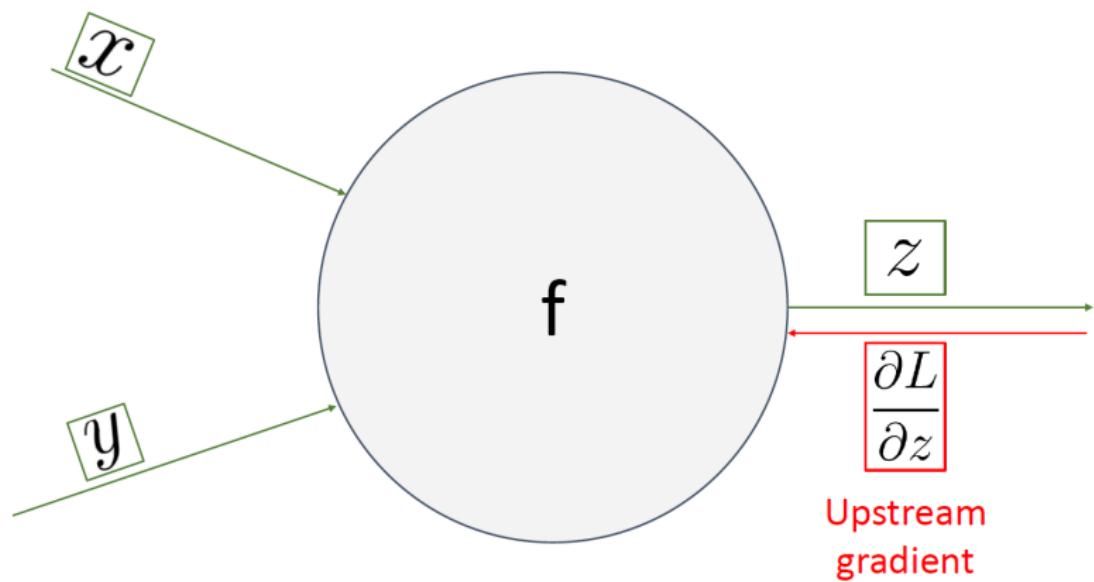
Local Gradient

## Upstream Gradient

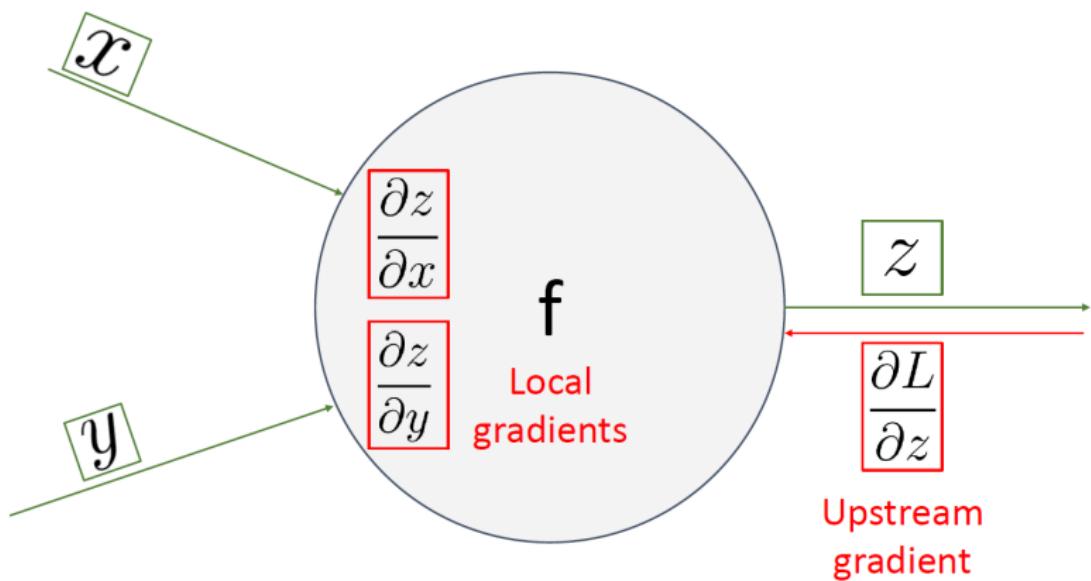
## Up/Down stream gradients



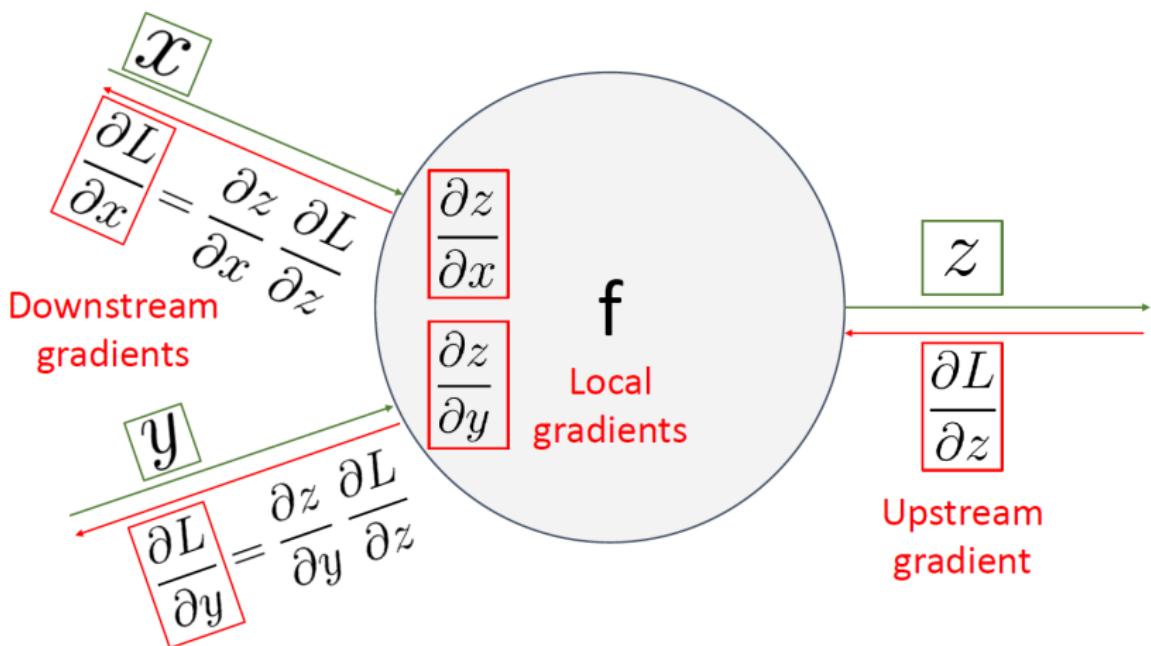
## Up/Down stream gradients



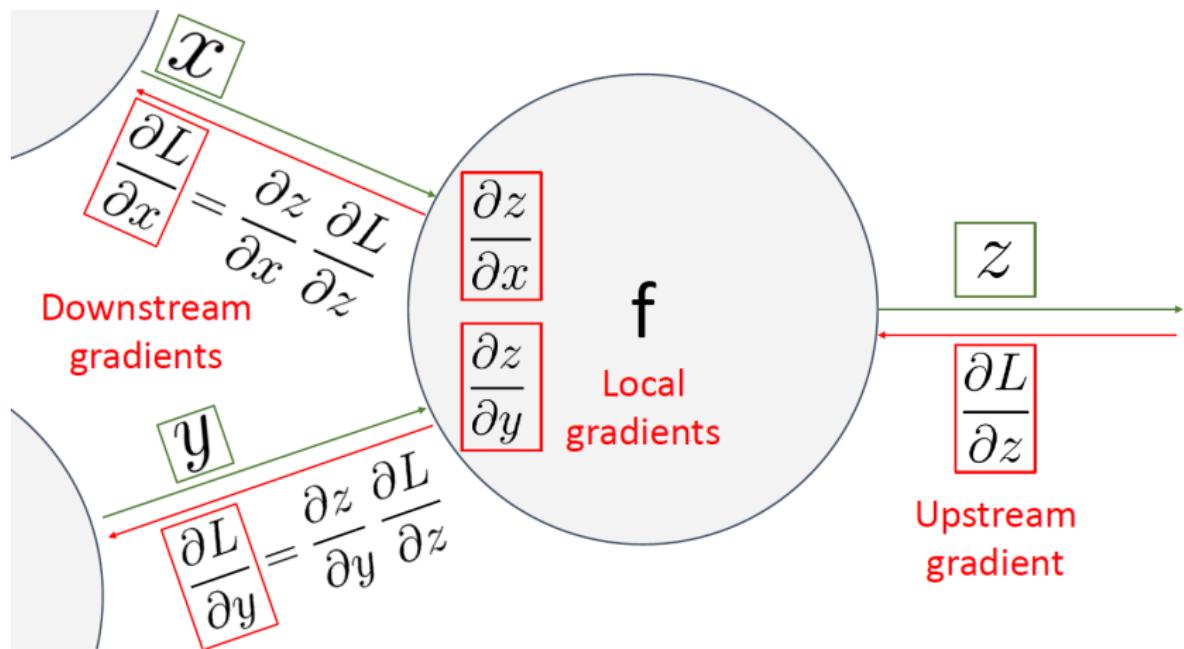
## Up/Down stream gradients



## Up/Down stream gradients

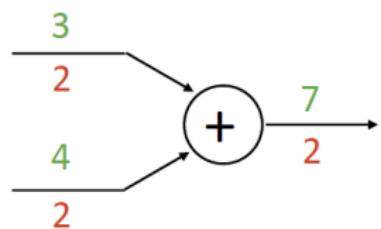


# Up/Down stream gradients



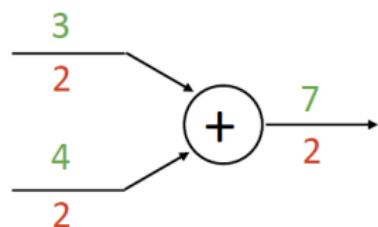
## Patterns in Gradient Flow

**add gate: gradient distributor**

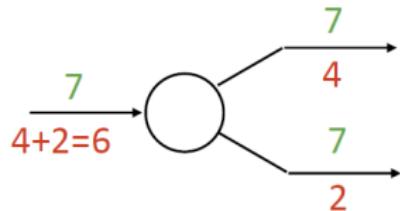


# Patterns in Gradient Flow

**add gate: gradient distributor**

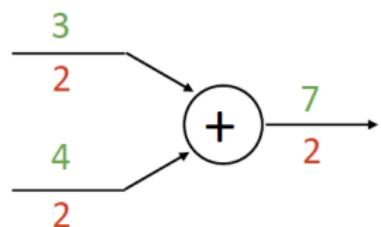


**copy gate: gradient adder**

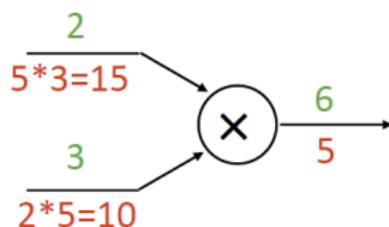


# Patterns in Gradient Flow

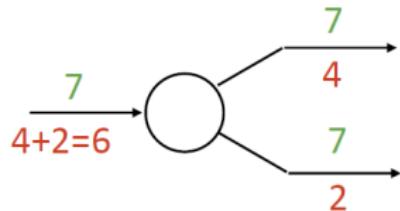
**add gate:** gradient distributor



**mul gate:** “swap multiplier”

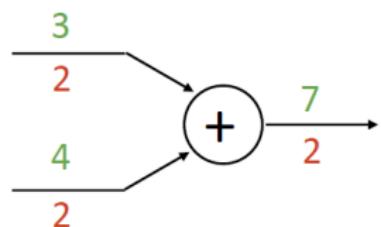


**copy gate:** gradient adder

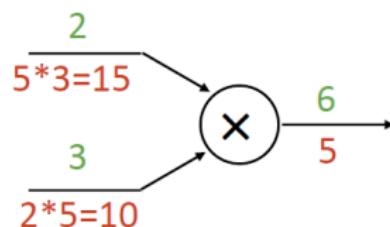


# Patterns in Gradient Flow

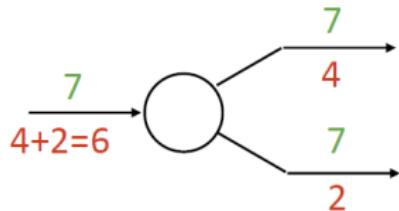
**add** gate: gradient distributor



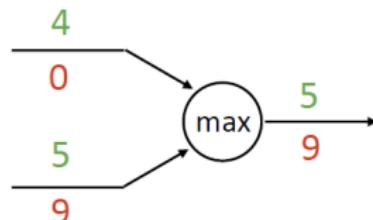
**mul** gate: “swap multiplier”



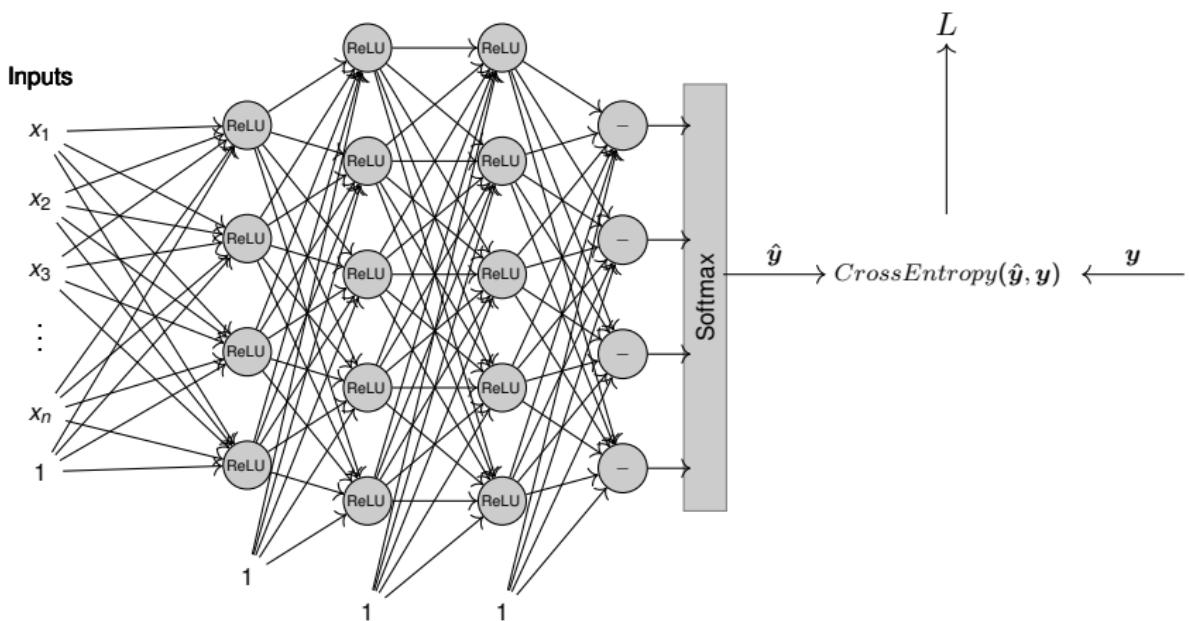
**copy** gate: gradient adder



**max** gate: gradient router

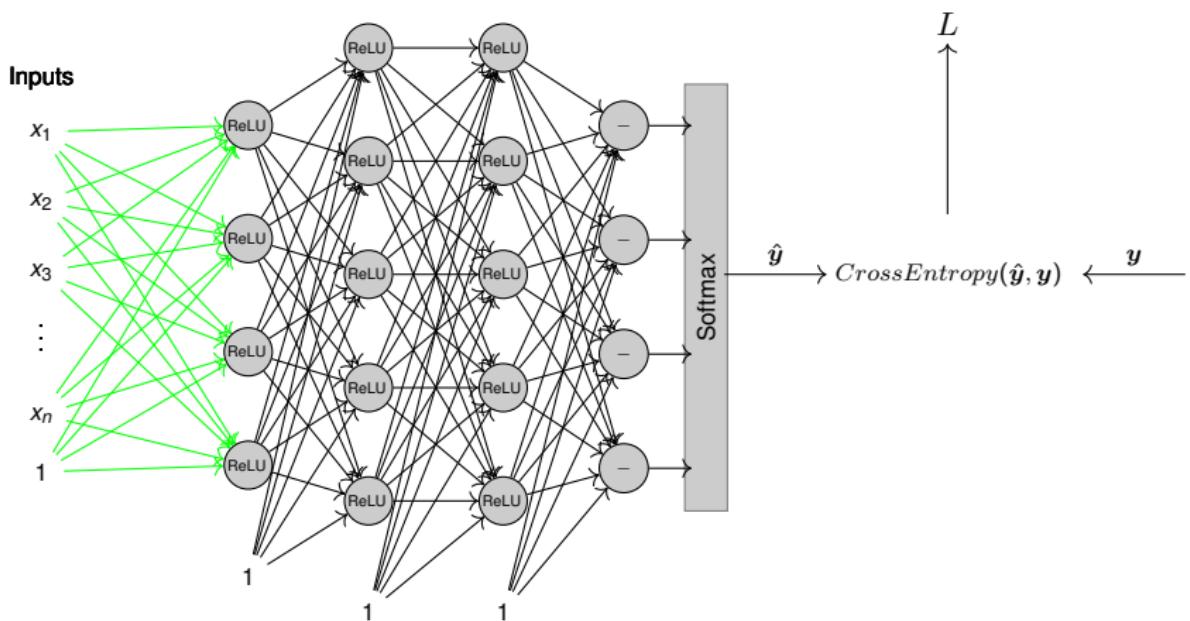


# Fully connected neural network - Forward Pass



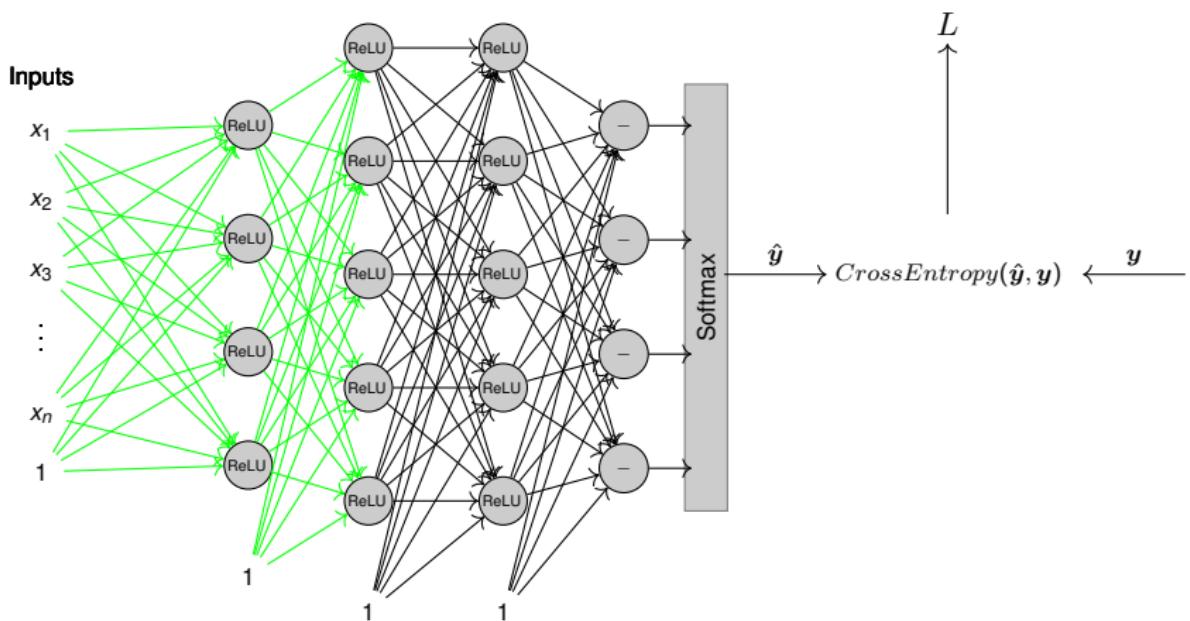
$$\hat{y} = \text{Softmax}(W^4(\max(0, W^3(\max(0, W^2(\max(0, W^1 \mathbf{x} + \mathbf{b}^1)) + \mathbf{b}^2)) + \mathbf{b}^3)) + \mathbf{b}^4)$$

# Fully connected neural network - Forward Pass



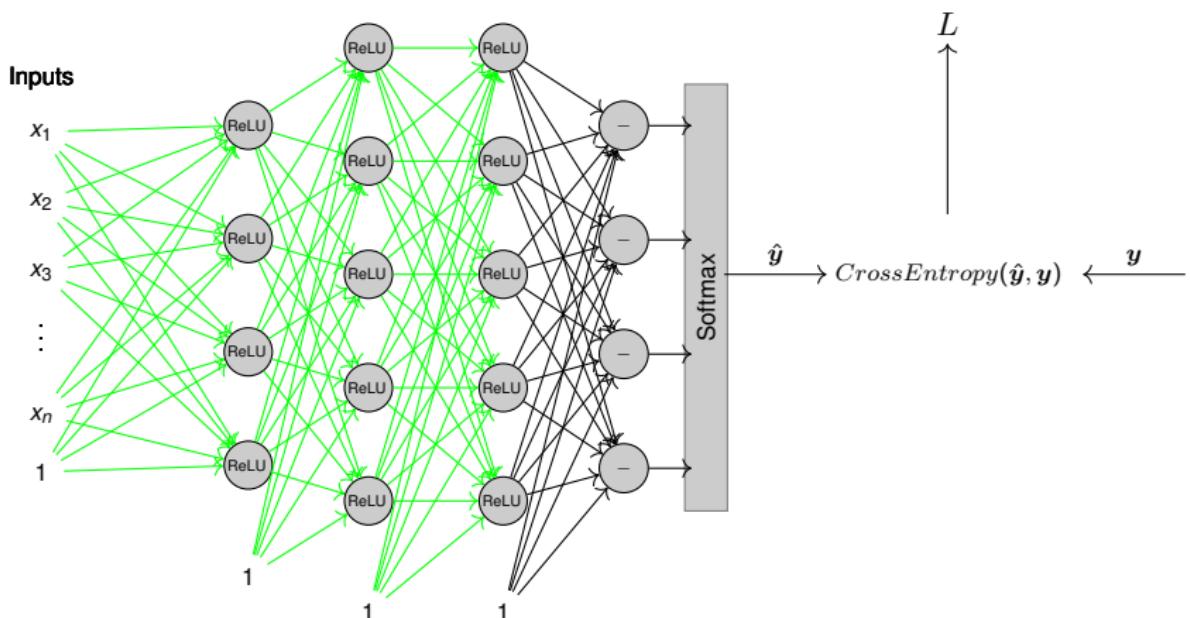
$$\mathbf{z}^1 = W^1 \mathbf{x} + \mathbf{b}^1, \quad \mathbf{o}^1 = \text{ReLU}(\mathbf{z}^1)$$

# Fully connected neural network - Forward Pass



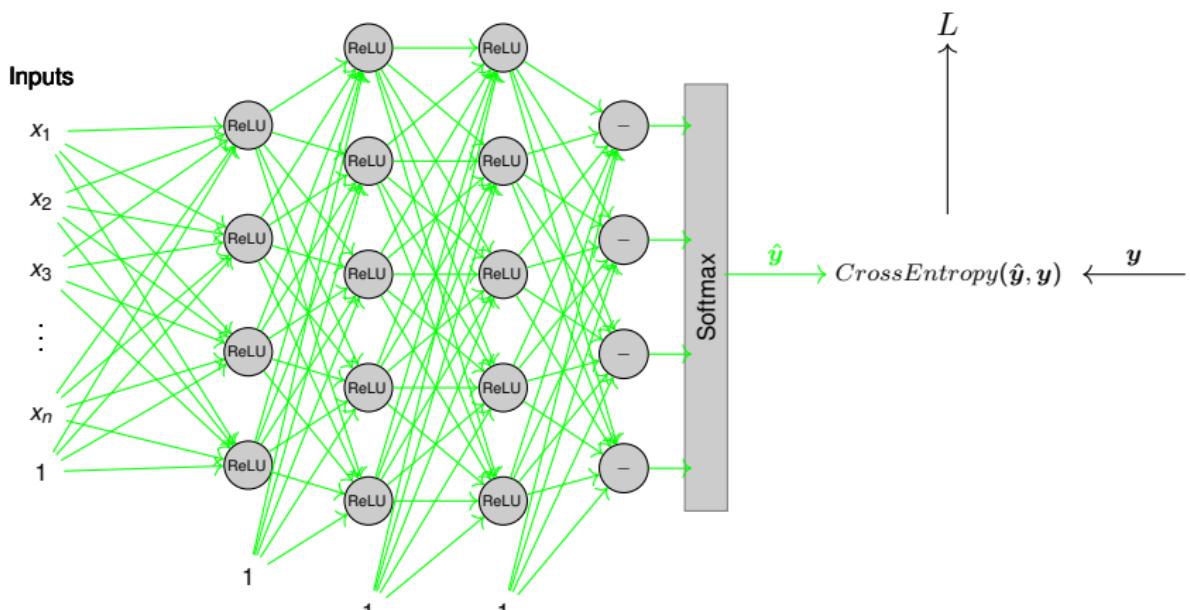
$$z^2 = W^2 o^1 + b^2, \quad o^2 = \text{ReLU}(z^2)$$

# Fully connected neural network - Forward Pass



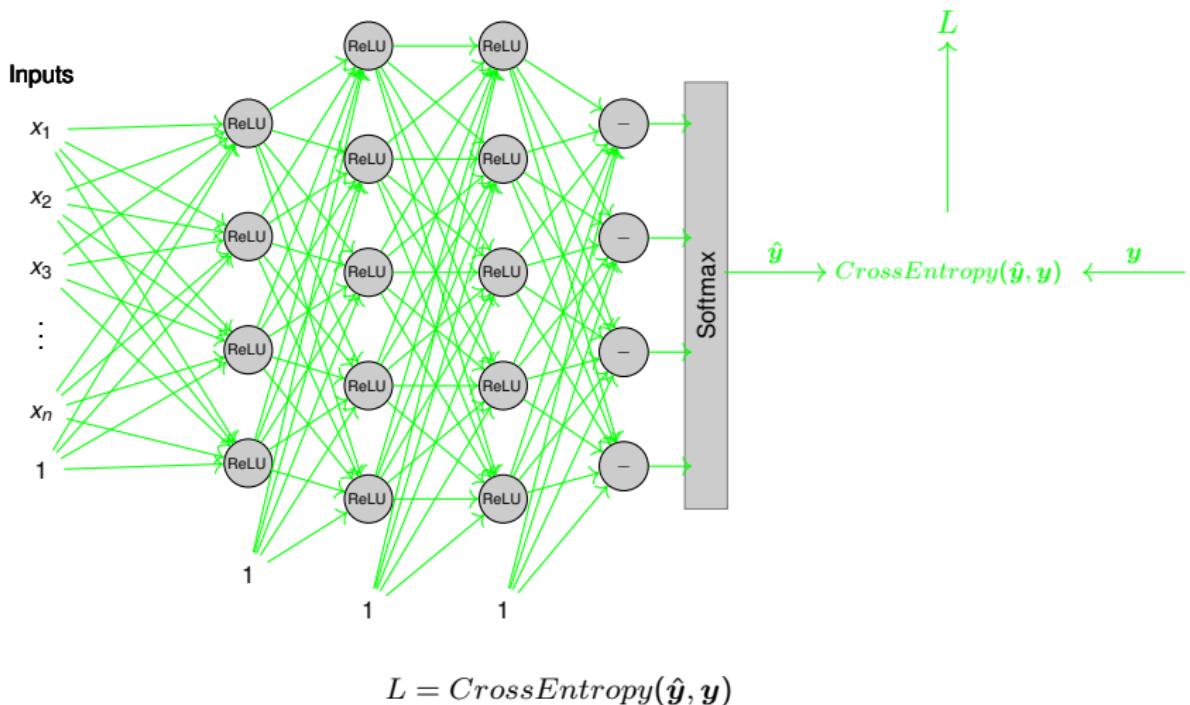
$$z^3 = W^3 o^2 + b^3, \quad o^3 = \text{ReLU}(z^3)$$

# Fully connected neural network - Forward Pass

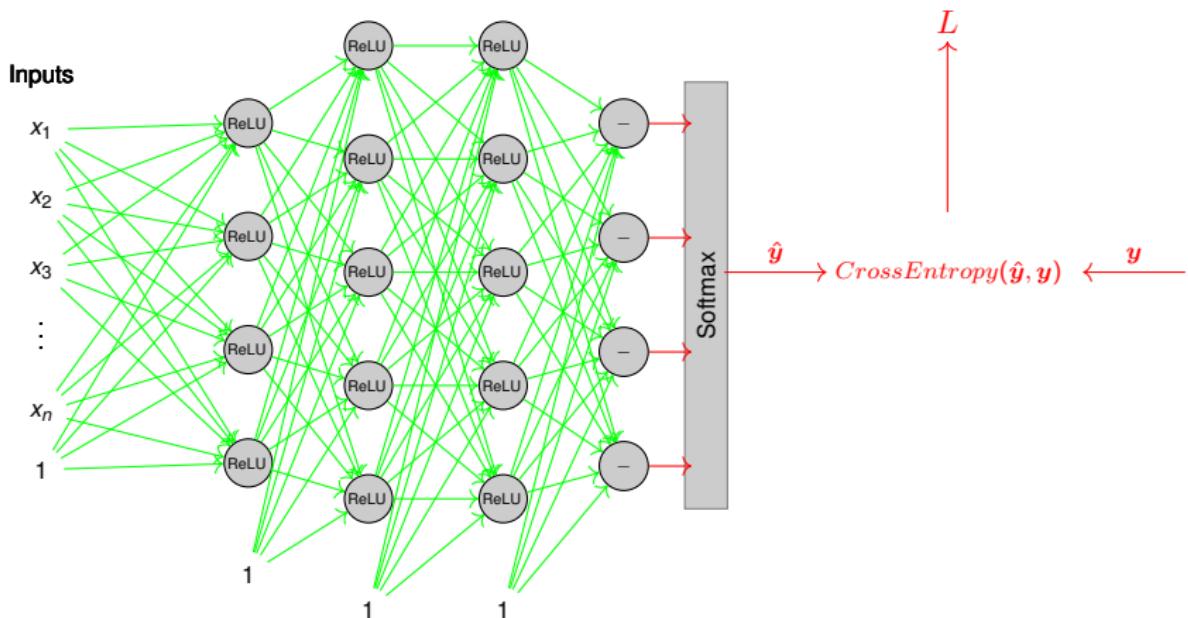


$$z^4 = W^4 o^3 + b^4, \quad \hat{y} = \text{softmax}(z^4)$$

# Fully connected neural network - Forward Pass

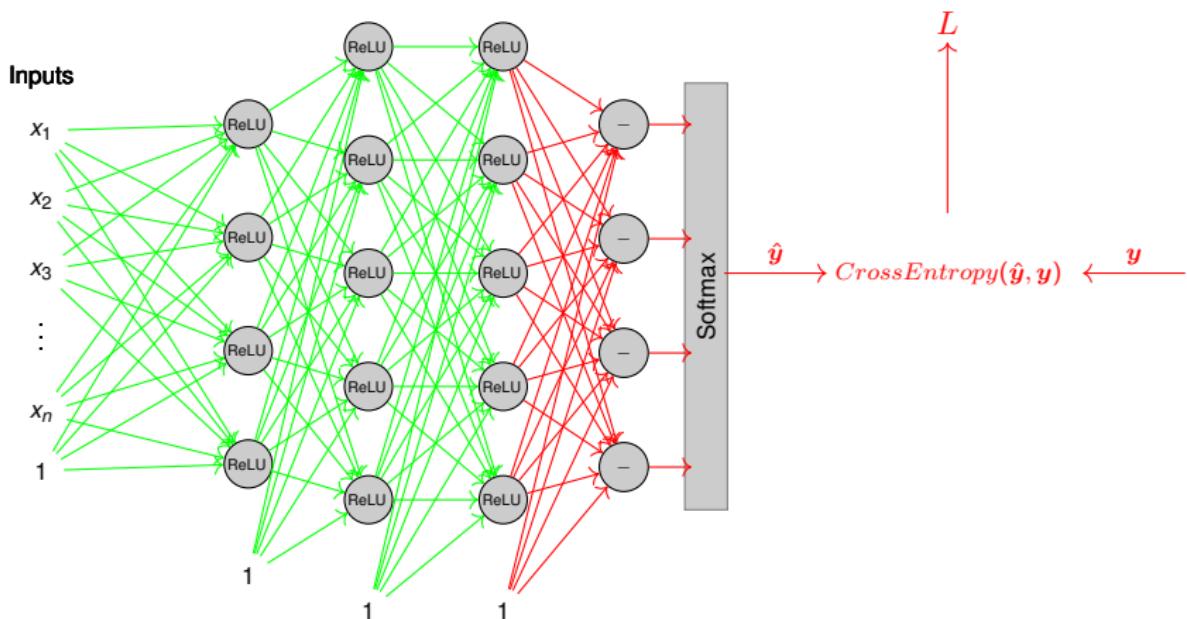


# Fully connected neural network - Backward Pass



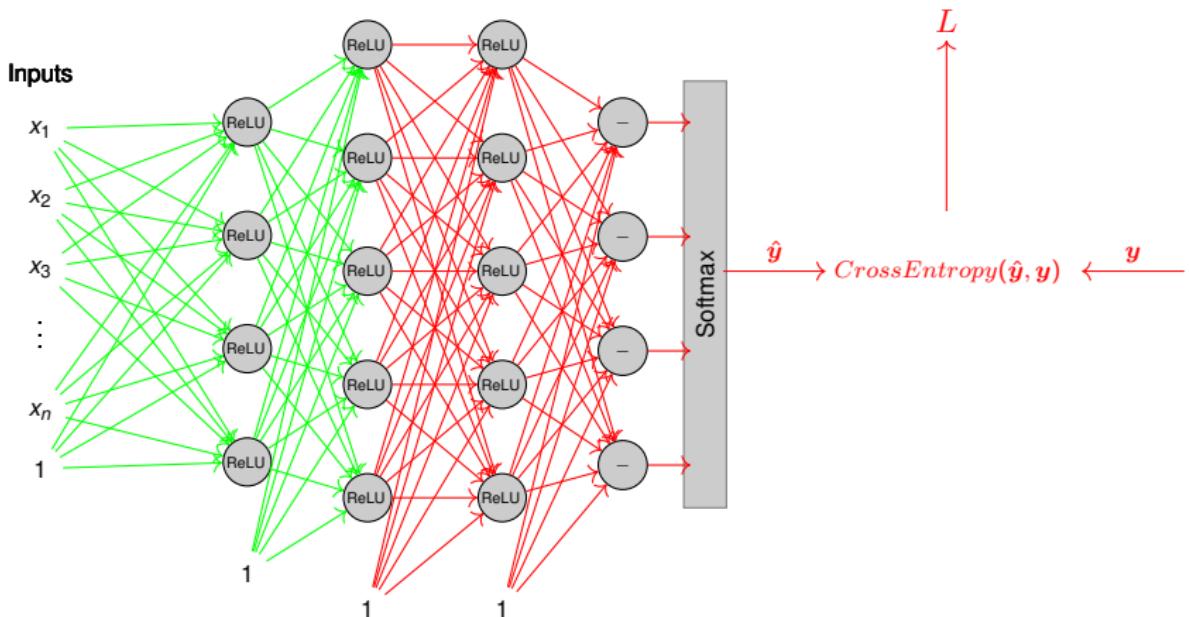
$$\frac{\partial L}{\partial z^4} = \hat{y} - y$$

# Fully connected neural network - Backward Pass



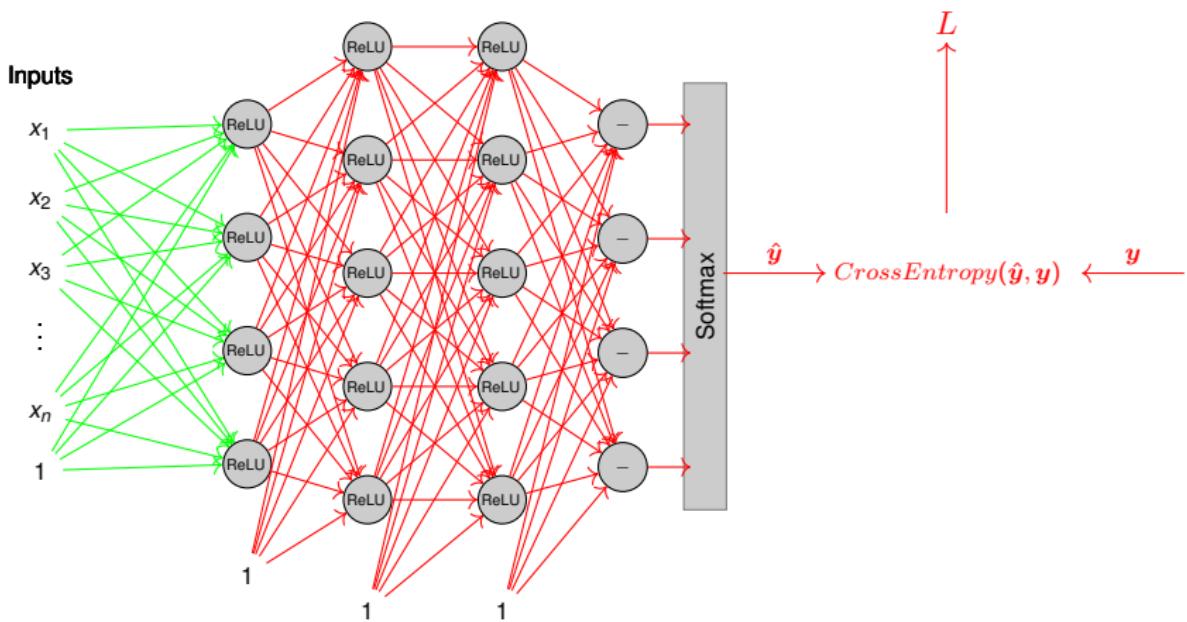
$$\frac{\partial L}{\partial W^4} = \frac{\partial L}{\partial z^4} \frac{\partial z^4}{\partial W^4} = \frac{\partial L}{\partial z^4} o^3, \quad \frac{\partial L}{\partial b^4} = \frac{\partial L}{\partial z^4} \frac{\partial z^4}{\partial b^4} = \frac{\partial L}{\partial z^4}, \quad \frac{\partial L}{\partial o^3} = \frac{\partial L}{\partial z^4} \frac{\partial z^4}{\partial o^3} = \frac{\partial L}{\partial z^4} W^4$$

# Fully connected neural network - Backward Pass



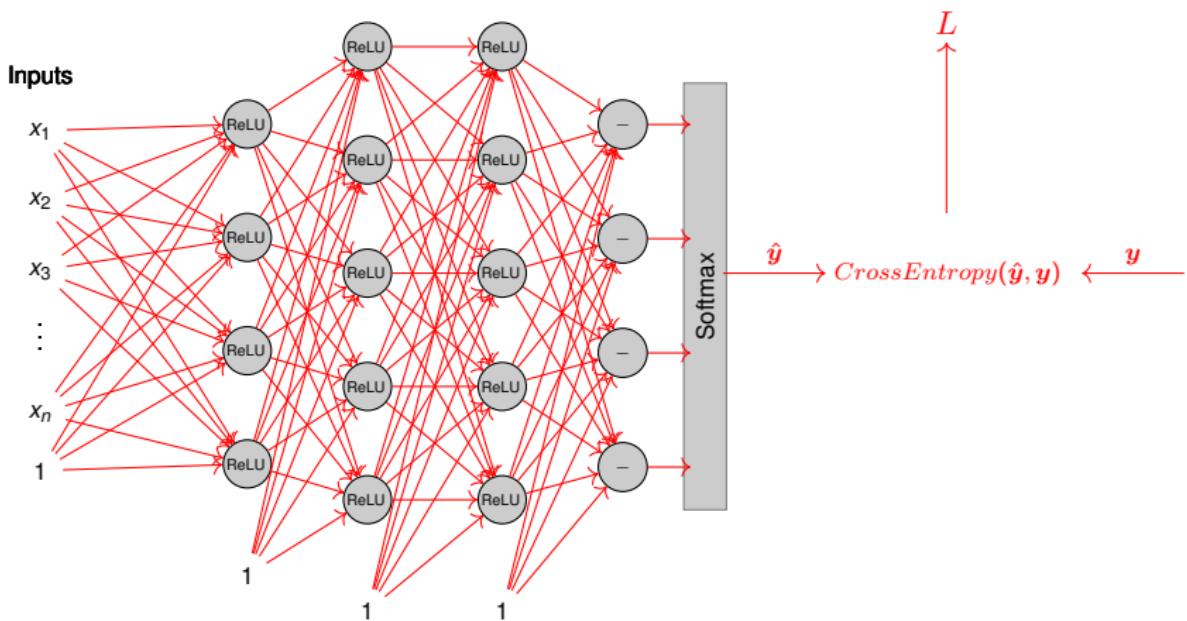
$$\frac{\partial L}{\partial z^3} = \frac{\partial L}{\partial o^3} \frac{\partial o^3}{\partial z^3} = \frac{\partial L}{\partial o^3} \text{ReLU}(\text{sign}(z^3)), \quad \frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial z^3} o^2, \quad \frac{\partial L}{\partial b^3} = \frac{\partial L}{\partial z^3}, \quad \frac{\partial L}{\partial o^2} = \frac{\partial L}{\partial z^3} W^3$$

# Fully connected neural network - Backward Pass



$$\frac{\partial L}{\partial z^2} = \frac{\partial L}{\partial o^2} \text{ReLU}(\text{sign}(z^2)), \quad \frac{\partial L}{\partial W^2} = \frac{\partial L}{\partial z^2} o^1, \quad \frac{\partial L}{\partial b^2} = \frac{\partial L}{\partial z^2}, \quad \frac{\partial L}{\partial o^1} = \frac{\partial L}{\partial z^2} W^2$$

# Fully connected neural network - Backward Pass



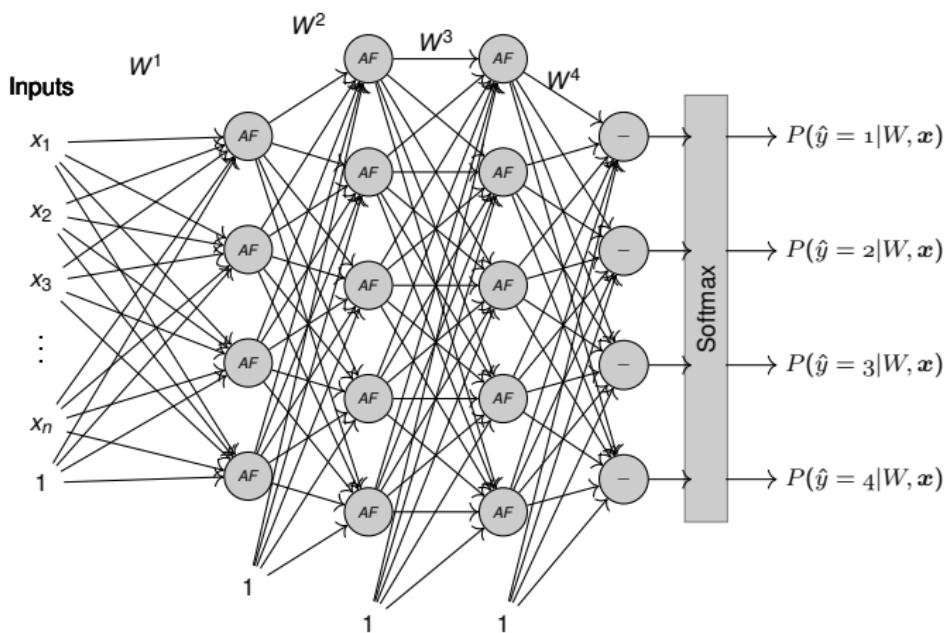
$$\frac{\partial L}{\partial z^1} = \frac{\partial L}{\partial o^1} \text{ReLU}(\text{sign}(z^1)), \quad \frac{\partial L}{\partial W^1} = \frac{\partial L}{\partial z^1} x, \quad \frac{\partial L}{\partial b^1} = \frac{\partial L}{\partial z^1}, \quad \frac{\partial L}{\partial x} = \frac{\partial L}{\partial z^1} W^1$$

## Vanishing/Exploding Gradient

## The vanishing gradient problem for deep networks

- A particular problem with training deep networks
    - The gradient of the error with respect to weights is unstable

# Fully connected neural networks



$$\hat{y} = f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^2(W^2(f^1(W^1 \mathbf{x})))))$$

## The problem with training deep networks

- A multilayer perceptron is a nested function

$$\hat{y} = f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^2(W^2(f^1(W^1 \mathbf{x}))))))$$

- $W^k$  is the weights matrix at the  $k^{th}$  layer
  - $f^k$  is the activation function of the  $k^{th}$  layer

- The error for  $X$  can be written as

$$L = CE(f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^2(W^2(f^1(W^1 \mathbf{x}))))), y)$$

- #### ■ *CE* is cross entropy

## Training deep networks

- Let  $z^\ell = W^\ell o^{\ell-1}$  and  $o^\ell = f^\ell(z^\ell)$

$$\frac{\partial L}{\partial \sigma^{\ell-1}} = \text{---} -$$

# Training deep networks

- Let  $\mathbf{z}^\ell = W^\ell \mathbf{o}^{\ell-1}$  and  $\mathbf{o}^\ell = f^\ell(\mathbf{z}^\ell)$

$$\frac{\partial L}{\partial \mathbf{o}^{\ell-1}} = - \frac{\partial L}{\partial \mathbf{o}^\ell}$$

## Training deep networks

- Let  $z^\ell = W^\ell o^{\ell-1}$  and  $o^\ell = f^\ell(z^\ell)$

$$\frac{\partial L}{\partial \boldsymbol{o}^{\ell-1}} = - \frac{\partial \boldsymbol{o}^\ell}{\partial \boldsymbol{z}^\ell} \frac{\partial L}{\partial \boldsymbol{o}^\ell}$$

## Training deep networks

- Let  $z^\ell = W^\ell o^{\ell-1}$  and  $o^\ell = f^\ell(z^\ell)$

$$\frac{\partial L}{\partial \boldsymbol{q}^{\ell-1}} = \frac{\partial \boldsymbol{z}^\ell}{\partial \boldsymbol{q}^{\ell-1}} \frac{\partial \boldsymbol{o}^\ell}{\partial \boldsymbol{z}^\ell} \frac{\partial L}{\partial \boldsymbol{o}^\ell}$$

## Training deep networks

- Let  $z^\ell = W^\ell o^{\ell-1}$  and  $o^\ell = f^\ell(z^\ell)$

$$\frac{\partial L}{\partial \theta^{\ell-1}} = \frac{\partial z^\ell}{\partial \theta^{\ell-1}} \frac{\partial o^\ell}{\partial z^\ell} \frac{\partial L}{\partial o^\ell}$$

$$\nabla_{\theta^{\ell-1}} L = W^\ell \cdot \nabla_{z^\ell} f^\ell \cdot \nabla_{\theta^\ell} L$$

- Where
    - $\nabla_x f$  is the jacobian matrix of  $f(x)$  w.r.t  $x$

## Training deep networks

- For

$$L = CE(f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^1(W^1 \mathbf{x})))), y)$$

## Training deep networks

- For

$$L \equiv CE(f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^1(W^1 x)))), y)$$

- ### ■ We get

$$\nabla_{w^k} L = \sigma^{k-1} \cdot \nabla_{-k} f^k \cdot W^{k+1} \cdot \nabla_{-k+1} f^{k+1} \cdots W^{N-1} \cdot \nabla_{-N-1} f^{N-1} \cdot W^N \cdot \nabla_{-N} f^N \cdot \nabla_{\hat{w}} CE$$

$$\nabla_{W^k} L = \mathbf{o}^{k-1} \cdot \nabla_{\mathbf{z}^k} f^k \cdot \left( \prod_{i=k+1}^N W^i \cdot \nabla_{\mathbf{z}^i} f^i \right) \cdot \nabla_{\hat{\mathbf{y}}} CE$$

- ### ■ Where

- $W^N$  is the gradient of the affine transformation  $z^n$  w.r.t the output of the  $n - 1^{th}$  layer ( $\sigma^{n-1}$ ) of the network
  - $\nabla_{z^n} f^n$  is gradient of activation function  $f^n(.)$  w.r.t to affine transformation  $z^n$  in the  $n^{th}$  layer of the network
  - $\sigma^{n-1}$  is the gradient of the affine transformation  $z^n$  w.r.t the weights of the  $n^{th}$  layer ( $W^n$ ) of the network
  - $\nabla_{\hat{y}} CE$  is the gradient of cross entropy w.r.t the output of the network ( $\hat{y}$ )

## Training deep networks

- For

$$L \equiv CE(f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^1(W^1 x)))), y)$$

- ### ■ We get

$$\nabla_{W^k} L = o^{k-1} \cdot \nabla_{\omega^k} f^k \cdot W^{k+1} \cdot \nabla_{\omega^{k+1}} f^{k+1} \cdots W^{N-1} \cdot \nabla_{\omega^{N-1}} f^{N-1} \cdot W^N \cdot \nabla_{\omega^N} f^N \cdot \nabla_{\hat{u}} CE$$

$$\nabla_{W^k} L = \mathbf{o}^{k-1} \cdot \nabla_{\mathbf{z}^k} f^k \cdot \left( \prod_{i=k+1}^N W^i \cdot \nabla_{\mathbf{z}^i} f^i \right) \cdot \nabla_{\hat{\mathbf{y}}} CE$$

- ### ■ Where

- $W^N$  is the gradient of the affine transformation  $z^n$  w.r.t the output of the  $n - 1^{th}$  layer ( $\sigma^{n-1}$ ) of the network
  - $\nabla_{z^n} f^n$  is gradient of activation function  $f^n(.)$  w.r.t to affine transformation  $z^n$  in the  $n^{th}$  layer of the network
  - $\sigma^{n-1}$  is the gradient of the affine transformation  $z^n$  w.r.t the weights of the  $n^{th}$  layer ( $W^n$ ) of the network
  - $\nabla_{\hat{y}} CE$  is the gradient of cross entropy w.r.t the output of the network ( $\hat{y}$ )

## Training deep networks

- For

$$L \equiv CE(f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^1(W^1 x)))), y)$$

- ### ■ We get

$$\nabla_{W^k} L = o^{k-1} \cdot \nabla_{z^k} f^k \cdot W^{k+1} \cdot \nabla_{z^{k+1}} f^{k+1} \dots W^{N-1} \cdot \nabla_{z^{N-1}} f^{N-1} \cdot W^N \cdot \nabla_{z^N} f^N \cdot \nabla_{\hat{y}} CE$$

$$\nabla_{W^k} L = \textcolor{blue}{o^{k-1}} \cdot \nabla_{\mathbf{z}^k} f^k \cdot \left( \prod_{i=k+1}^N \textcolor{red}{W^i} \cdot \nabla_{\mathbf{z}^i} f^i \right) \cdot \nabla_{\hat{\mathbf{y}}} CE$$

- ### ■ Where

- $W^N$  is the gradient of the affine transformation  $z^n$  w.r.t the output of the  $n - 1^{th}$  layer ( $\sigma^{n-1}$ ) of the network
  - $\nabla_{z^n} f^n$  is gradient of activation function  $f^n(.)$  w.r.t to affine transformation  $z^n$  in the  $n^{th}$  layer of the network
  - $\sigma^{n-1}$  is the gradient of the affine transformation  $z^n$  w.r.t the weights of the  $n^{th}$  layer ( $W^n$ ) of the network
  - $\nabla_{\hat{y}} CE$  is the gradient of cross entropy w.r.t the output of the network ( $\hat{y}$ )

# Training deep networks

- For

$$L = CE(f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^1(W^1 \mathbf{x})))), y)$$

- We get

$$\nabla_{W^k} L = \sigma^{k-1} \cdot \nabla_{\mathbf{z}^k} f^k \cdot W^{k+1} \cdot \nabla_{\mathbf{z}^{k+1}} f^{k+1} \cdots W^{N-1} \cdot \nabla_{\mathbf{z}^{N-1}} f^{N-1} \cdot W^N \cdot \nabla_{\mathbf{z}^N} f^N \cdot \nabla_{\hat{y}} CE$$

$$\nabla_{W^k} L = \sigma^{k-1} \cdot \nabla_{\mathbf{z}^k} f^k \cdot \left( \prod_{i=k+1}^N W^i \cdot \nabla_{\mathbf{z}^i} f^i \right) \cdot \nabla_{\hat{y}} CE$$

- Repeated multiplication by the weights matrix and jacobian matrices of activation functions will result in exploding or vanishing gradients
  - Depends on the spectral norm (largest singular value) of jacobian matrices
- The gradients in the lower/earlier layers can explode or vanish
  - Resulting in insignificant or unstable gradient descent updates
  - Problem gets worse as network depth increases

# References

- Pattern Recognition and Machine Learning, Ch.1 & Ch.4
- Deep Learning, Ian Goodfellow, MIT Press, Ch.5, 6, 7, 8 ,and 9