Quick Recap- Logistics

Google

CS4782 cornell SP25

AI Mode   All   Images   Shopping   Videos   News   Forums   More ▾

These are results for **CS 4782** cornell SP25 · ↩ Revert

https://classes.cornell.edu › browse › roster › SP25 › class ⋮

Spring 2025 - CS 4782 - Class Roster

Spring 2025 - **CS 4782** - This class is an introductory course to deep learning. It covers the fundamental principles behind training and inference of deep ...

🌐    Cornell Computer Science Department
       https://www.cs.cornell.edu › courses › cs4782 ⋮

CS 4/5782 - Cornell Computer Science

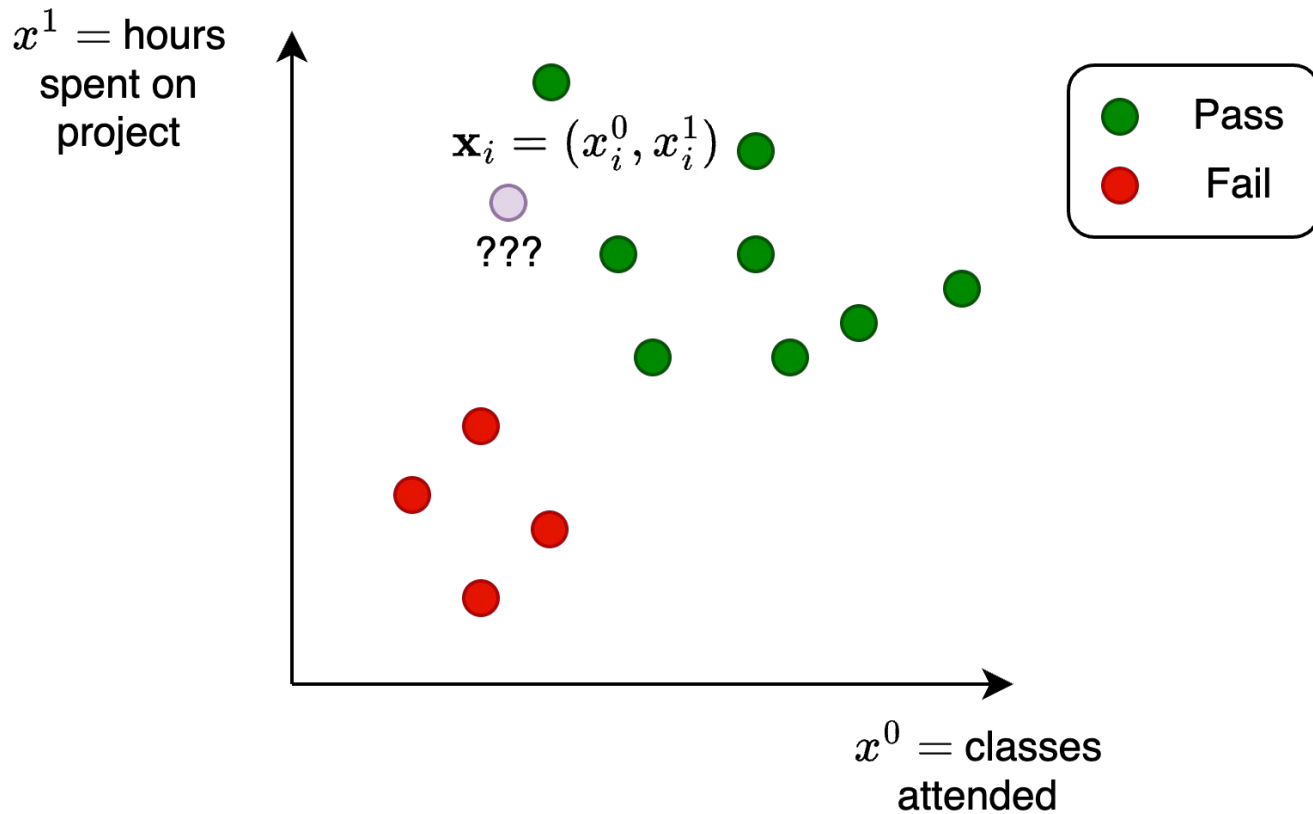**CS 4782**: Intro to Deep Learning, Spring 2025. Overview; Assignments; Schedule; References; Policies. Instructors: Kilian Q. Weinberger and Jennifer J. Sun.
Missing: ~~SP25~~ | Show results with: SP25

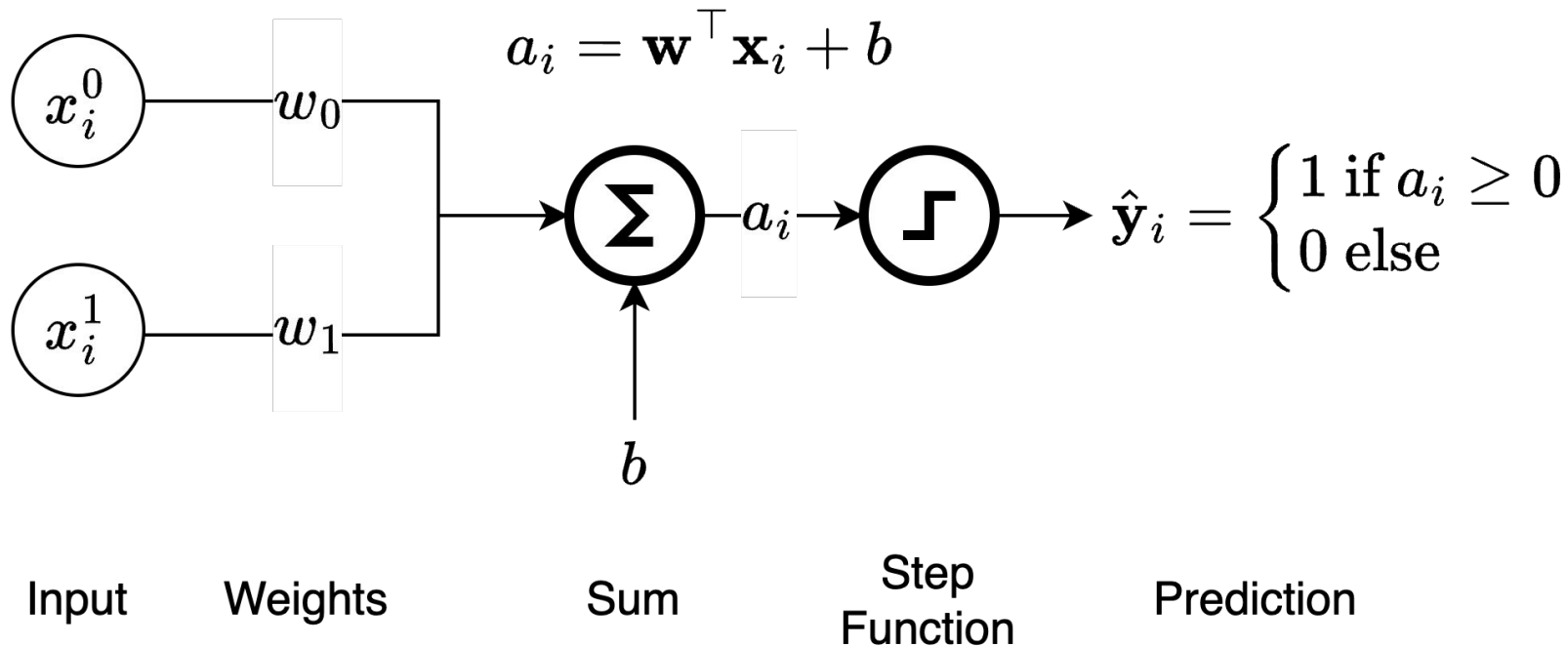# No laptops/mobiles/smart devices in class please!

# Agenda

- Perceptron
- Logistic Regression
- Gradient Descent
- Multi-Layer Perceptrons (MLPs)
- Backpropagation

# A Classification Problem: Will I Pass This Class?



$x^1 =$ hours spent on project

$\mathbf{x}_i = (x_i^0, x_i^1)$

???

Pass

Fail

$x^0 =$ classes attended

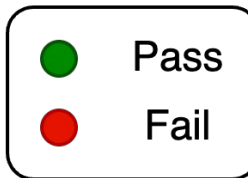# What are key components in ML?

# A Classification Problem: Will I Pass This Class?

$x^1 =$ hours spent on project

$\mathbf{x}_i = (x_i^0, x_i^1)$

???

**Recall:**

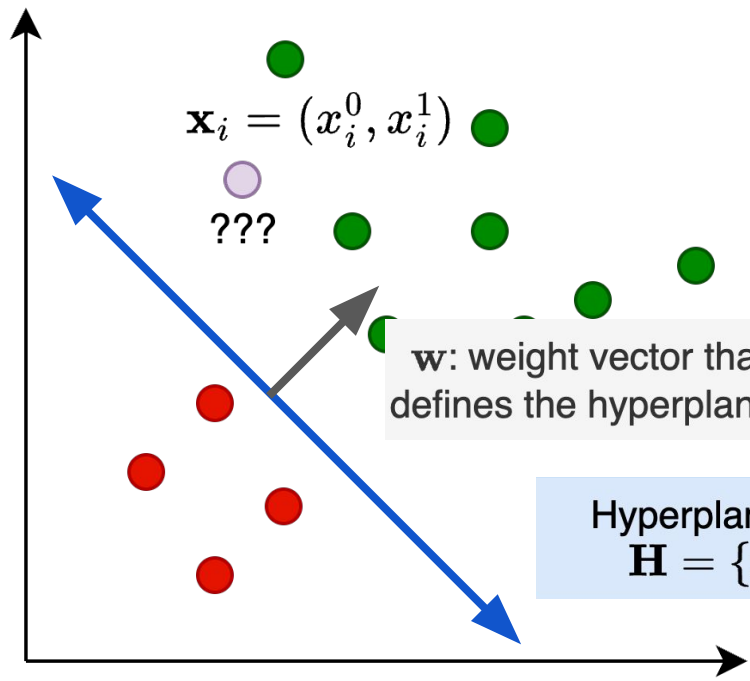$$y_i = \begin{cases} 1 \text{ if } \mathbf{w}^\top \mathbf{x}_i + b \geq 0 \\ 0 \text{ else} \end{cases}$$

Pass

Fail

$\mathbf{w}$: weight vector that defines the hyperplane
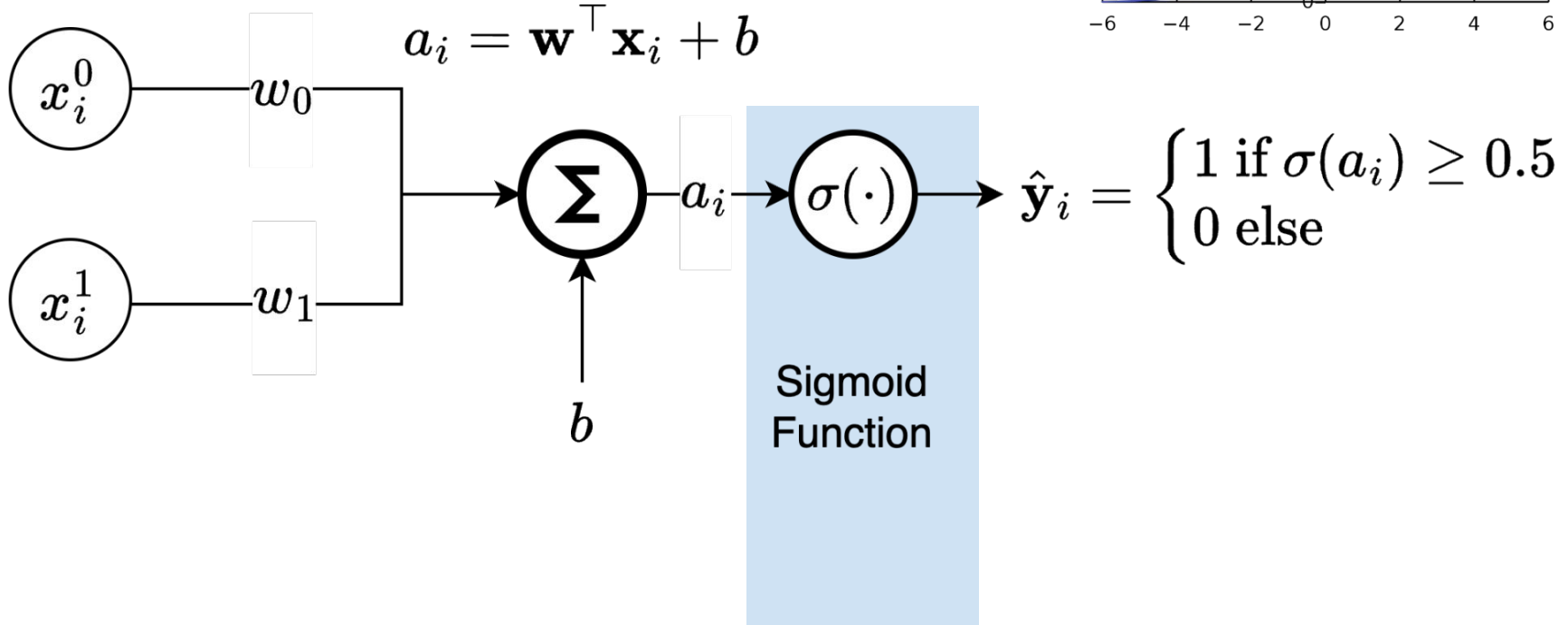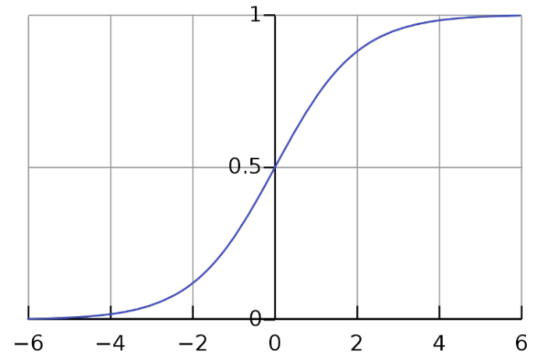
Hyperplane perpendicular to $\mathbf{w}$:
$$\mathbf{H} = \{\mathbf{x} : \mathbf{w}^\top \mathbf{x} + b = 0\}$$

$x^0 =$ classes attended

# The "Soft" Perceptron

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



$$a_i = \mathbf{w}^\top \mathbf{x}_i + b$$

$x_i^0$

$w_0$

$x_i^1$

$w_1$

$\Sigma$

$a_i$

$\sigma(\cdot)$

$$\hat{\mathbf{y}}_i = \begin{cases} 1 \text{ if } \sigma(a_i) \geq 0.5 \\ 0 \text{ else} \end{cases}$$

$b$

Sigmoid Function

# Clean Up Bias Term $\quad \mathbf{w}^\top \mathbf{x}_i + b$

Absorb bias term into feature vector:

$$\mathbf{x}_i \quad \text{becomes} \quad \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix} \quad \text{and} \quad \mathbf{w} \quad \text{becomes} \quad \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$$

We can see that:

$$\begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}^\top \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix} = \mathbf{w}^\top \mathbf{x}_i + b$$
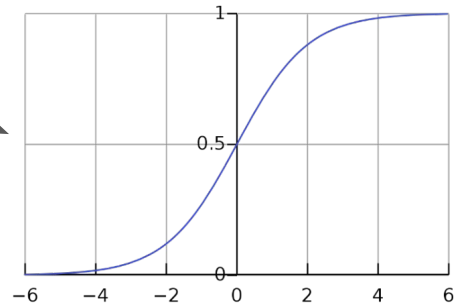
Can rewrite logistic regression as

$$\hat{\mathbf{y}}_i = \sigma(\mathbf{w}^\top \mathbf{x}_i)$$



$\mathbf{x}_i = (x_i^0, x_i^1)$

???



$\mathbf{x}_i = (x_i^0, x_i^1)$

???

**Origin**

# Maximum Likelihood Estimation

$$\hat{\mathbf{y}}_i = \sigma(\mathbf{w}^\top \mathbf{x}_i)$$



Maximize the likelihood of the observed data $(\mathbf{x}_i, \mathbf{y}_i)$, where $\mathbf{y}_i \in \{0, 1\}$:

$$p(\mathbf{y}_i | \mathbf{x}_i) =$$

Derive the loss:

$$\ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) = -[\mathbf{y}_i \log \sigma(\mathbf{w}^\top \mathbf{x}_i + b) + (1 - \mathbf{y}_i) \log(1 - \sigma(\mathbf{w}^\top \mathbf{x}_i + b))]$$

## Our Goal: Minimize the Loss

Given some training dataset:

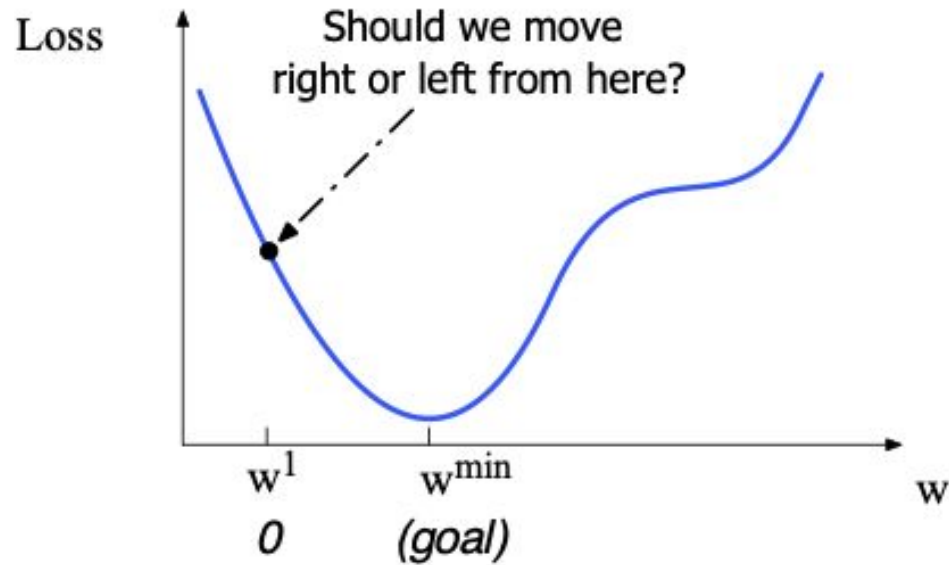$$\mathcal{D}_{\text{TR}} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=0}^{n}$$

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}; \mathcal{D}_{\text{TR}}) = \frac{1}{n} \sum_i^n \ell(\hat{\mathbf{y}}_i, \mathbf{y}_i)$$

$$= \frac{1}{n} \sum_i^n \ell(\sigma(\mathbf{w}^\top \mathbf{x}_i), \mathbf{y}_i)$$

# Gradient Descent

# Visualize Gradient Descent in 1-D



Loss

Should we move
right or left from here?

$w^1$       $w^{min}$

$0$       (goal)

w

https://web.stanford.edu/~jurafsky/slp3/

# Visualize Gradient Descent in 1-D



Loss

slope of loss at $w^1$
is negative

$w^1$

$0$

$w^{min}$
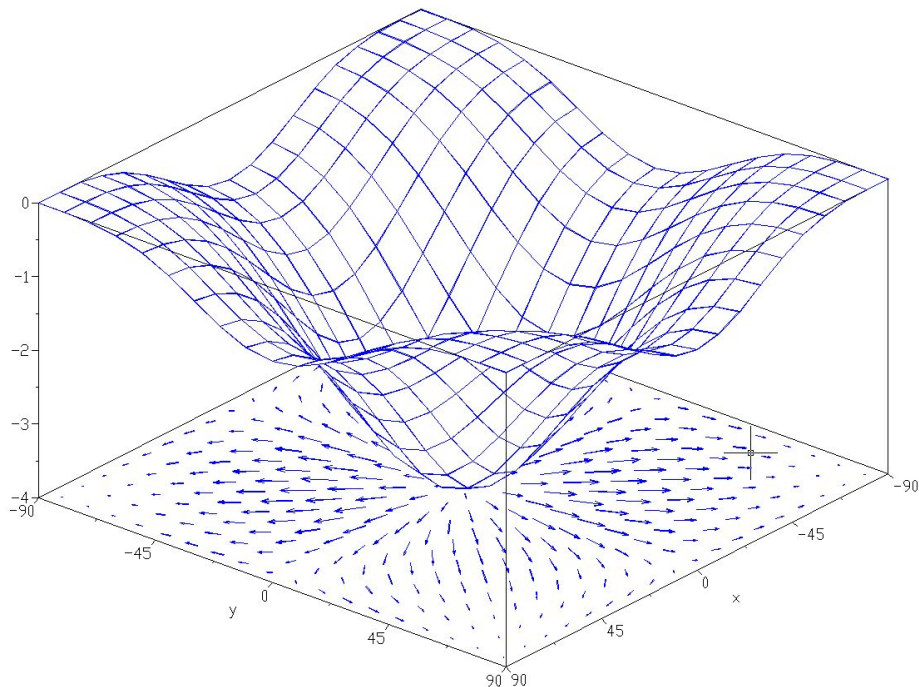
(goal)

w

https://web.stanford.edu/~jurafsky/slp3/

# Visualize Gradient Descent in 1-D

# Gradients

$$\nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w};\mathcal{D}_{\mathrm{TR}}) = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial w^{(0)}}(\mathbf{w};\mathcal{D}_{\mathrm{TR}}) \\ \\ \frac{\partial \mathcal{L}}{\partial w^{(1)}}(\mathbf{w};\mathcal{D}_{\mathrm{TR}}) \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial w^{(m)}}(\mathbf{w};\mathcal{D}_{\mathrm{TR}}) \end{bmatrix}, \nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w};\mathcal{D}_{\mathrm{TR}}) \in \mathbb{R}^m$$
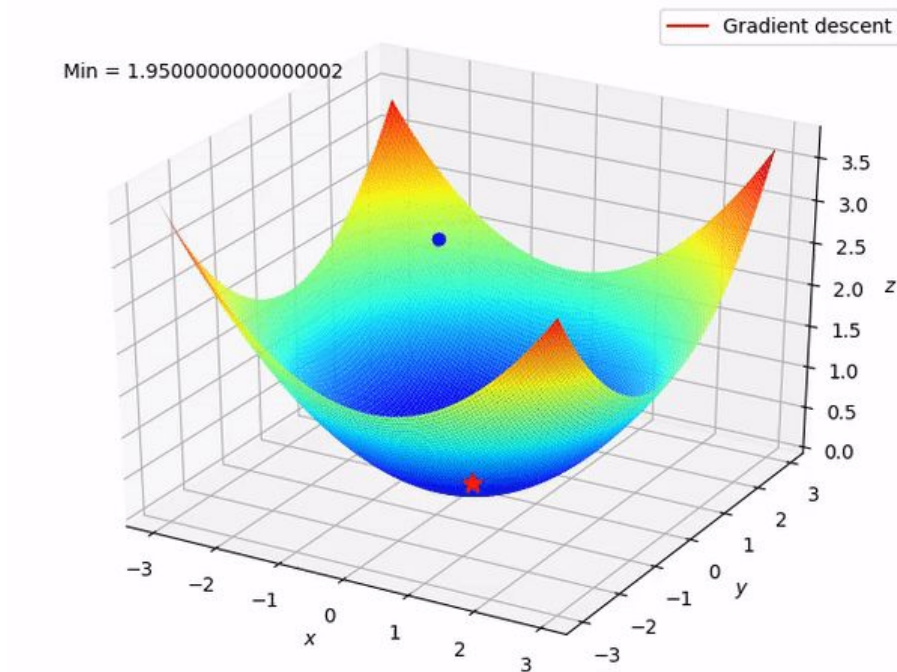


**Gradient Descent**:

- Find the gradient at current point
- Move in **opposite** direction with learning rate $\alpha$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla_{\mathbf{w}_t}\mathcal{L}(\mathbf{w}_t;\mathcal{D}_{\mathrm{TR}})$$
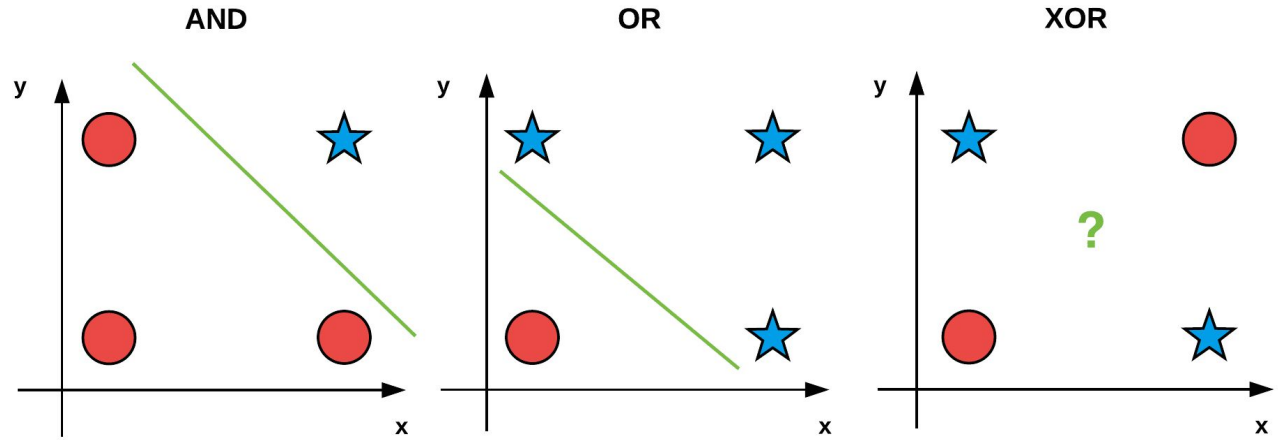
https://www.ml-science.com/gradients

# Gradient Descent (GD)

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla_{\mathbf{w}_t} \mathcal{L}(\mathbf{w}_t; \mathcal{D}_{\mathrm{TR}})$$
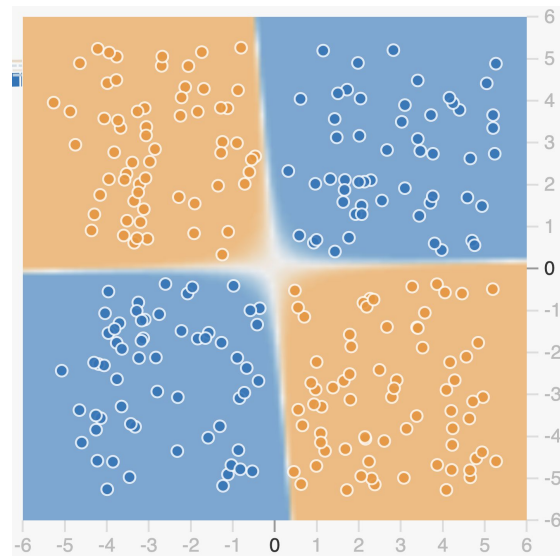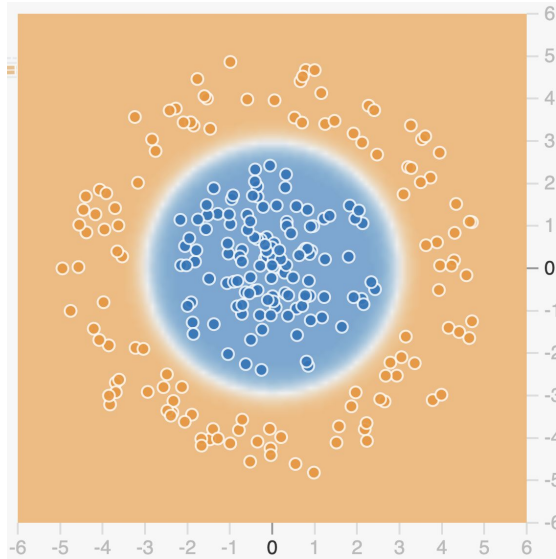
# The XOR Problem

- Perceptron can't learn the XOR function
  - Simple logical operation
- Data is not linearly separable



https://www.pyimagesearch.com/2021/05/06/implementing-the-perceptron-neural-network-with-python/

# Discuss: What are some ways to handle data that is not linearly separable?

Without deep learning!

# Feature Engineering



input image

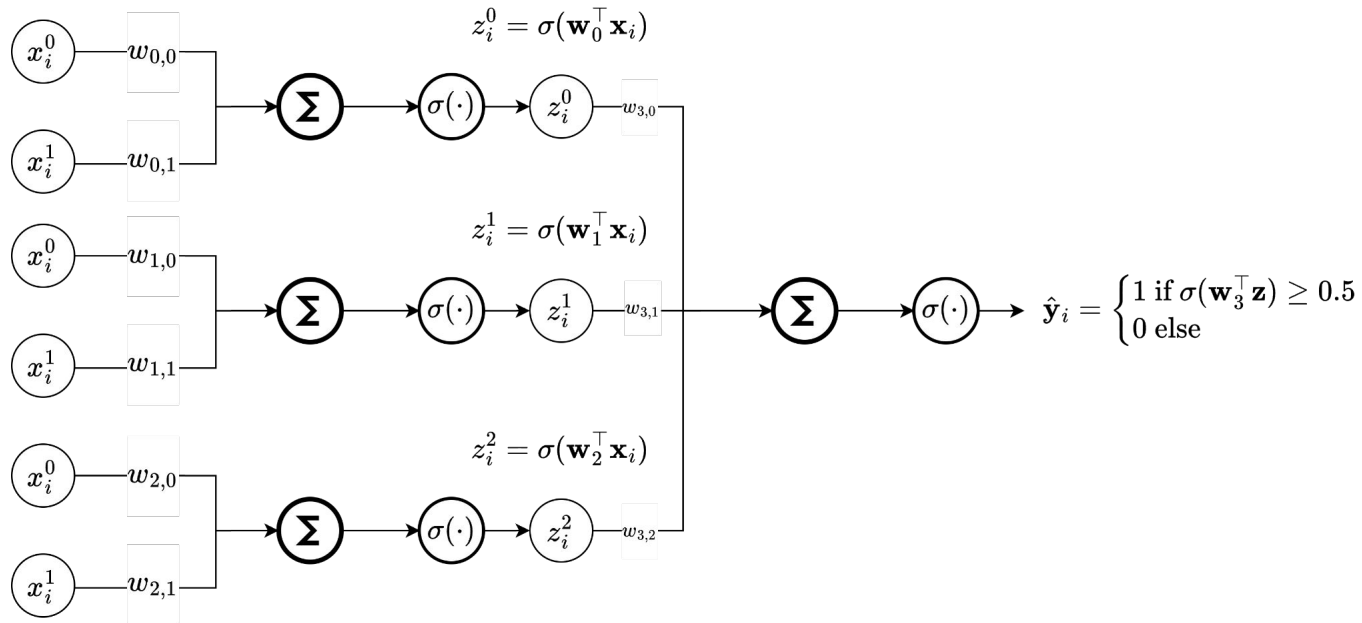classification → "dog"

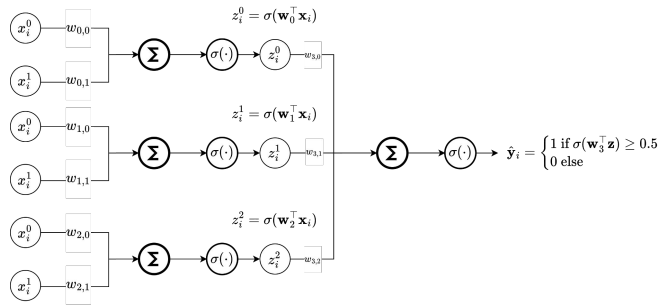

input image

classification → "cat"

# Multi-Layer Perceptron (MLP)

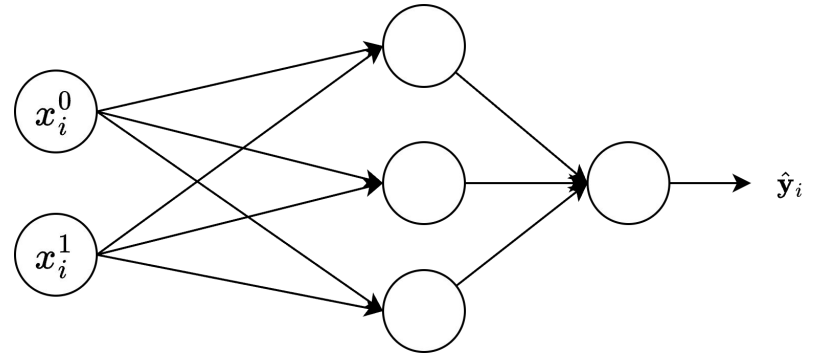- Compose multiple perceptrons to **learn** intermediate features

An MLP with 1 hidden layer with 3 hidden units
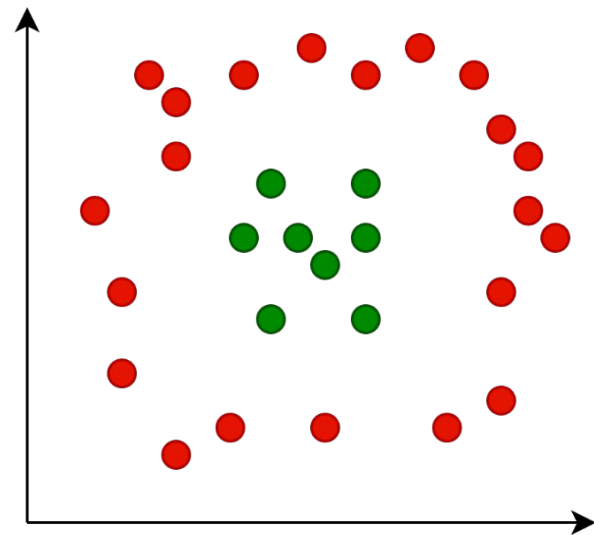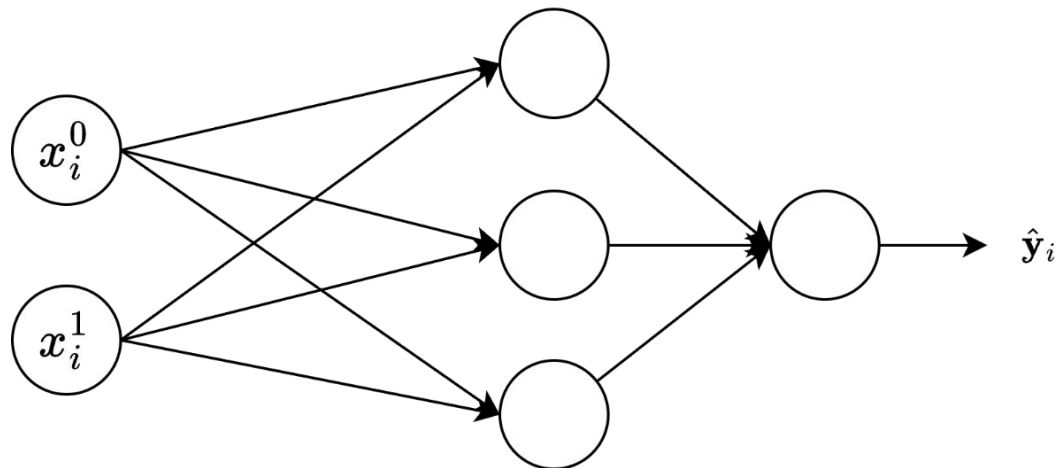
# A Simplified MLP Diagram



1 Hidden Layer,
3 Hidden Units

$z_i^0 = \sigma(\mathbf{w}_0^\top \mathbf{x}_i)$

$z_i^1 = \sigma(\mathbf{w}_1^\top \mathbf{x}_i)$

$z_i^2 = \sigma(\mathbf{w}_2^\top \mathbf{x}_i)$

$\hat{\mathbf{y}}_i = \begin{cases} 1 \text{ if } \sigma(\mathbf{w}_3^\top \mathbf{z}) \geq 0.5 \\ 0 \text{ else} \end{cases}$

$x_i^0$

$x_i^1$

$\hat{\mathbf{y}}_i$

# Complex Decision Boundaries

- What does this extra layer give us?
    - Can compose multiple linear classifiers
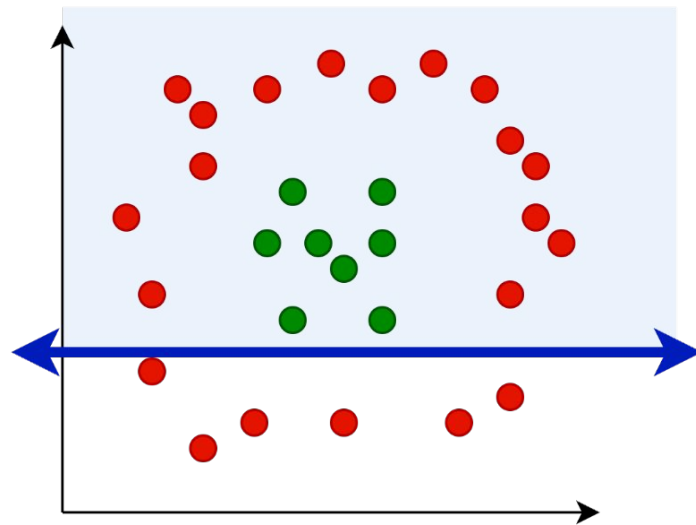
$$x_i^0$$
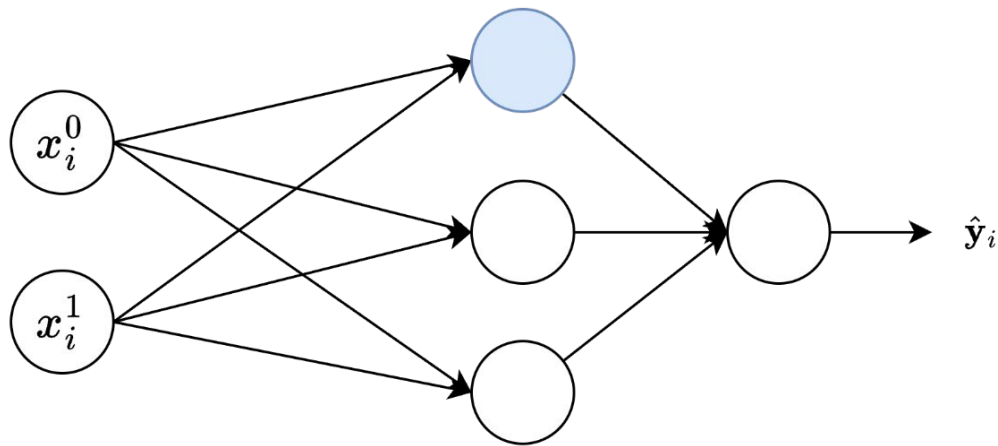
$$x_i^1$$

$$\hat{\mathbf{y}}_i$$

# Complex Decision Boundaries

- What does this extra layer give us?
  - Can compose multiple linear classifiers

# Complex Decision Boundaries

- What does this extra layer give us?
  - Can compose multiple linear classifiers

# Complex Decision Boundaries

- What does this extra layer give us?
  - Can compose multiple linear classifiers
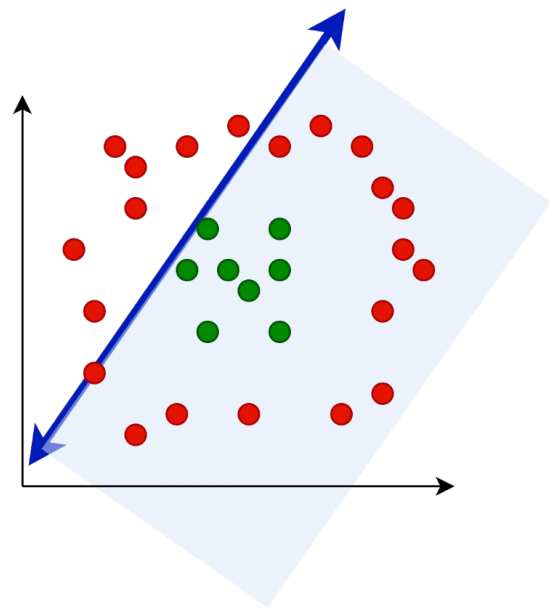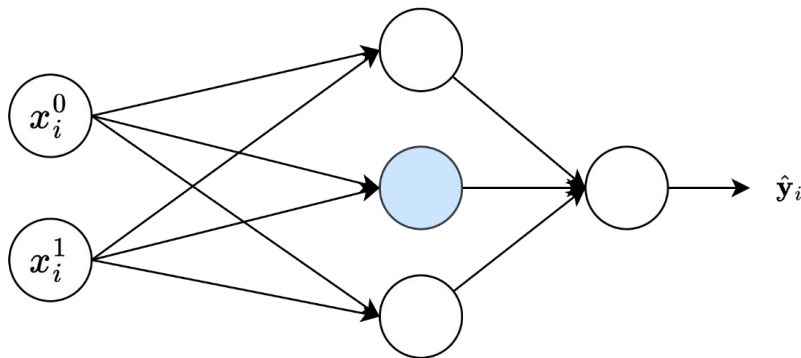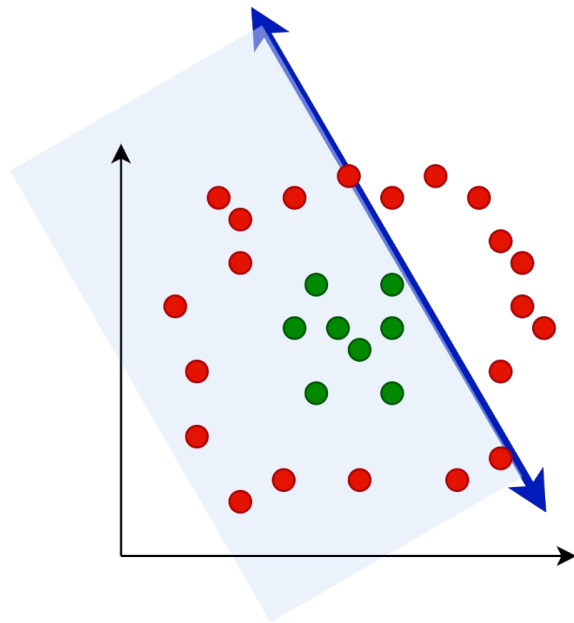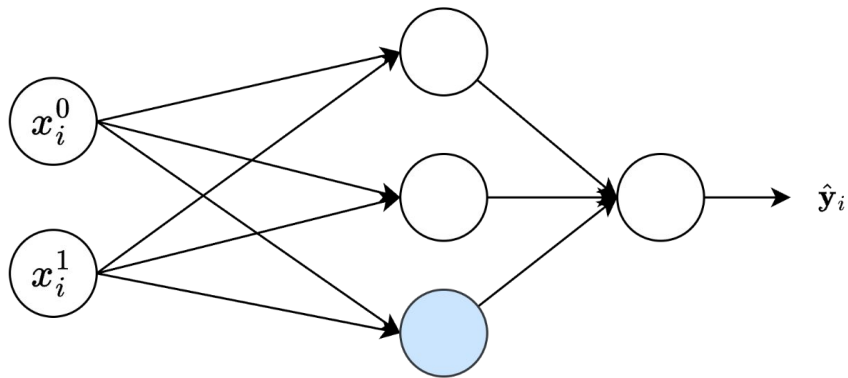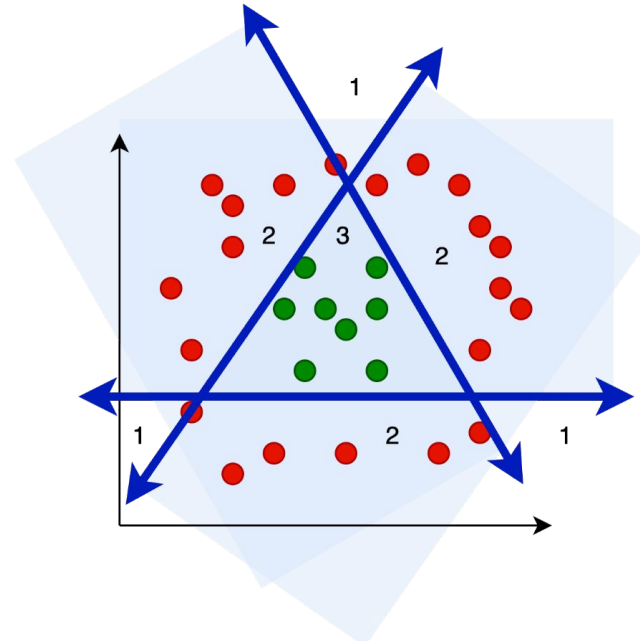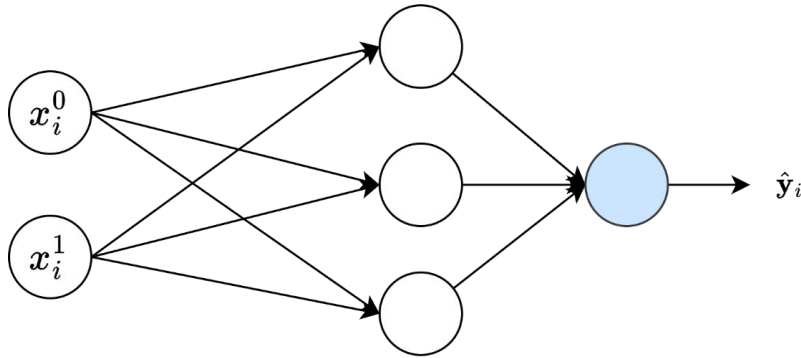


$\hat{\mathbf{y}}_i$

# Complex Decision Boundaries

- What does this extra layer give us?
  - Can compose multiple linear classifiers

# Increasing Depth

## Discuss: What about just one layer?

# Complex Decision Boundaries

- Can compose *arbitrarily* complex decision boundaries

# Activation Functions

- Can replace the sigmoid with other nonlinear functions
  - Still universal approximators!



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

Squash between 0 and 1

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

Squash between -1 and 1

$$\text{ReLU}(x) = \max(0, x)$$

Threshold at 0

https://pytorch.org/docs/stable/nn.html#non-linear-activations-weighted-sum-nonlinearity

# How to learn MLP weights?

Gradient descent!

# Calculus Review: The Chain Rule

Lagrange's Notation:

$$\text{If } h(x) = f(g(x)), \text{ then } h' = f'(g(x))g'(x)$$

Leibniz's Notation:

$$\text{If } z = h(y), y = g(x), \text{ then } \frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}$$

# Calculus Review: The Chain Rule

Lagrange's Notation:    $\text{If } h(x) = f(g(x)), \text{ then } h' = f'(g(x))g'(x)$

Leibniz's Notation:    $\text{If } z = h(y), y = g(x), \text{ then } \frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}$

Example:    $\text{If } z = \ln(y), y = x^2, \text{ then}$

$$\frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}$$

$$= (\frac{1}{y})(2x) = (\frac{1}{x^2})(2x) = \frac{2}{x}$$

# Multivariate Chain Rule



If $f(u)$ is $z = f(v(u), w(u))$, then

$$\frac{\partial f}{\partial u} = \left( \frac{\partial v}{\partial u} \frac{\partial z}{\partial v} + \frac{\partial w}{\partial u} \frac{\partial z}{\partial w} \right)$$

# Backpropagation- An Example



Forward

5
$u$

125
$v = u^3$

10
$w = u + u$

1250
$z = v \cdot w$

https://windowsontheory.org/2020/11/03/yet-another-backpropagation-tutorial/

# Backpropagation- An Example

Forward

Backward

5
$u$

125
$v = u^3$

10
$w = u + u$

1250
$z = v \cdot w$

$$\frac{\partial z}{\partial u} = \left( \frac{\partial v}{\partial u} \cdot \frac{\partial z}{\partial v} + \frac{\partial w}{\partial u} \cdot \frac{\partial z}{\partial w} \right)$$

# Backpropagation- Key Idea



If you know
$$\frac{\partial z}{\partial v_1}, \frac{\partial z}{\partial v_2}, \frac{\partial z}{\partial v_3}$$

You can compute $\frac{\partial z}{\partial u}$

$$\frac{\partial z}{\partial u} = \left( \frac{\partial v_1}{\partial u} \cdot \frac{\partial z}{\partial v_1} + \frac{\partial v_2}{\parti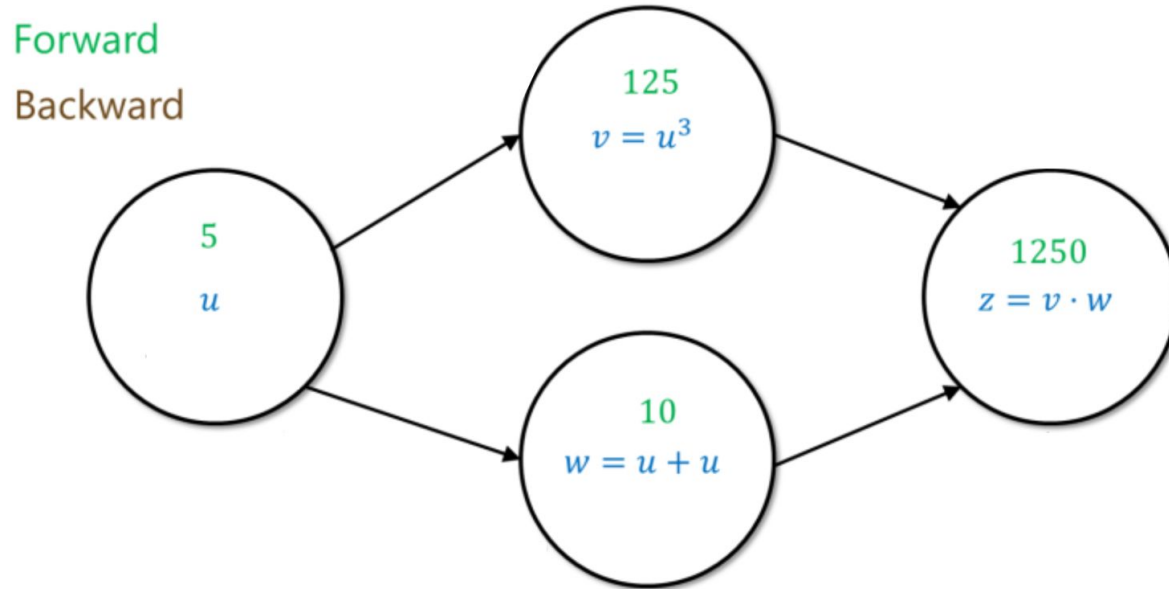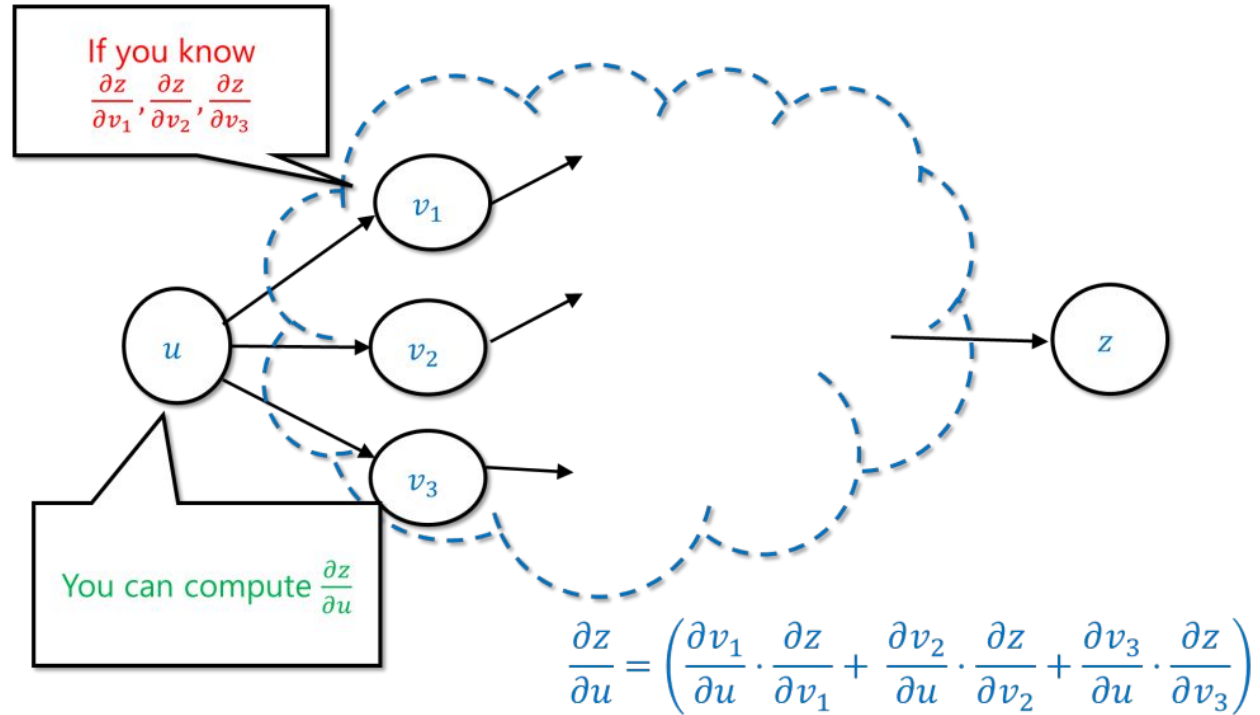al u} \cdot \frac{\partial z}{\partial v_2} + \frac{\partial v_3}{\partial u} \cdot \frac{\partial z}{\partial v_3} \right)$$
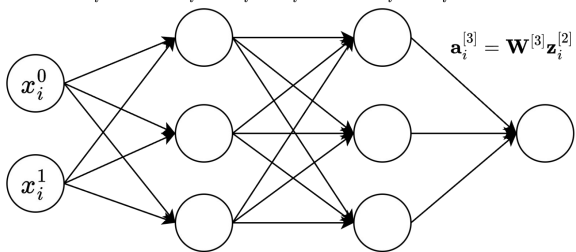
Preview

# Backpropagation- MLPs



**Algorithm** Forward Pass through MLP

1: **Input:** input $\mathbf{x}$, weight matrices $\mathbf{W}^{[1]}, \ldots, \mathbf{W}^{[L]}$, bias vectors $\mathbf{b}^{[1]}, \ldots, \mathbf{b}^{[L]}$
2: $\mathbf{z}^{[0]} = \mathbf{x}$ ▷ Initialize input
3: **for** $l = 1$ **to** $L$ **do**
4: $\quad \mathbf{a}^{[l]} = \mathbf{W}^{[l]}\mathbf{z}^{[l-1]} + \mathbf{b}^{[l]}$ ▷ Linear transformation
5: $\quad \mathbf{z}^{[l]} = \sigma^{[l]}(\mathbf{a}^{[l]})$ ▷ Nonlinear activation
6: **end for**
7: **Output:** $\mathbf{z}^{[L]}$

**Algorithm** Backward Pass through MLP

1: **Input:** $\{\mathbf{z}^{[1]}, \ldots, \mathbf{z}^{[L]}\}$, $\{\mathbf{a}^{[1]}, \ldots, \mathbf{a}^{[L]}\}$, loss gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}}$
2: $\delta^{[L]} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}} \odot \sigma^{[L]'}(\mathbf{a}^{[L]})$ ▷ Error term
3: **for** $l = L$ **to** $1$ **do**
4: $\quad \frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[l]}} = \delta^{[l]}(\mathbf{z}^{[l-1]})^T$ ▷ Gradient of weights
5: $\quad \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{[l]}} = \delta^{[l]}$ ▷ Gradient of biases
6: $\quad \delta^{[l-1]} = ((\mathbf{W}^{[l]})^T \delta^{[l]}) \odot \sigma^{[l-1]'}(\mathbf{a}^{[l-1]})$
7: **end for**
8: **Output:** $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[1:L]}}, \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{[1:L]}}$

# Backpropagation- MLPs



$\mathbf{a}_i^{[1]} = \mathbf{W}^{[1]}\mathbf{z}_i^{[0]} + \mathbf{b}_i^{[1]}$ $\quad$ $\mathbf{a}_i^{[2]} = \mathbf{W}^{[2]}\mathbf{z}_i^{[1]} + \mathbf{b}_i^{[2]}$

$\mathbf{a}_i^{[3]} = \mathbf{W}^{[3]}\mathbf{z}_i^{[2]} + \mathbf{b}_i^{[3]}$

$x_i^0$

$x_i^1$

$\mathbf{z}_i^{[0]} = \mathbf{x}_i$ $\qquad$ $\mathbf{z}_i^{[1]} = \sigma(\mathbf{a}_i^{[1]})$ $\qquad$ $\mathbf{z}_i^{[2]} = \sigma(\mathbf{a}_i^{[2]})$ $\qquad$ $\mathbf{z}_i^{[3]} = \sigma(\mathbf{a}_i^{[3]})$

**Algorithm** Forward Pass through MLP

1: **Input:** input $\mathbf{x}$, weight matrices $\mathbf{W}^{[1]}, \ldots, \mathbf{W}^{[L]}$, bias vectors $\mathbf{b}^{[1]}, \ldots, \mathbf{b}^{[L]}$
2: $\mathbf{z}^{[0]} = \mathbf{x}$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ ▷ Initialize input
3: **for** $l = 1$ **to** $L$ **do**
4: $\quad \mathbf{a}^{[l]} = \mathbf{W}^{[l]}\mathbf{z}^{[l-1]} + \mathbf{b}^{[l]}$ $\qquad\qquad\qquad$ ▷ Linear transformation
5: $\quad \mathbf{z}^{[l]} = \sigma^{[l]}(\mathbf{a}^{[l]})$ $\qquad\qquad\qquad\qquad\quad$ ▷ Nonlinear activation
6: **end for**
7: **Output:** $\mathbf{z}^{[L]}$

**Algorithm** Backward Pass through MLP (Detailed)
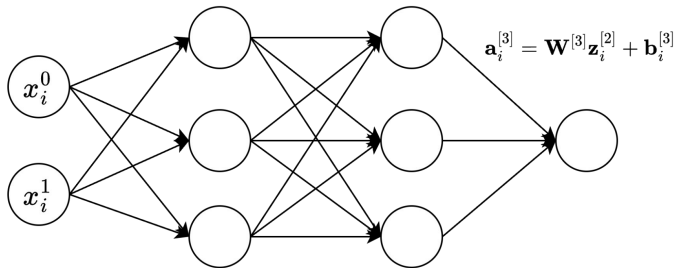
1: **Input:** $\{\mathbf{z}^{[1]}, \ldots, \mathbf{z}^{[L]}\}$, $\{\mathbf{a}^{[1]}, \ldots, \mathbf{a}^{[L]}\}$, loss gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}}$
2: $\delta^{[L]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[L]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}} \frac{\partial \mathbf{z}^{[L]}}{\partial \mathbf{a}^{[L]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}} \odot \sigma^{[L]'}(\mathbf{a}^{[L]})$ $\qquad$ ▷ Error term
3: **for** $l = L$ **to** $1$ **do**
4: $\quad \frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[l]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{W}^{[l]}} = \delta^{[l]}(\mathbf{z}^{[l-1]})^T$ $\qquad\qquad$ ▷ Gradient of weights
5: $\quad \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{[l]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{b}^{[l]}} = \delta^{[l]}$ $\qquad\qquad\qquad$ ▷ Gradient of biases
6: $\quad \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[l-1]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{z}^{[l-1]}} = (\mathbf{W}^{[l]})^T \delta^{[l]}$
7: $\quad \delta^{[l-1]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l-1]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[l-1]}} \frac{\partial \mathbf{z}^{[l-1]}}{\partial \mathbf{a}^{[l-1]}} = ((\mathbf{W}^{[l]})^T \delta^{[l]}) \odot \sigma^{[l-1]'}(\mathbf{a}^{[l-1]})$
8: **end for**
9: **Output:** $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[1:L]}}$, $\frac{\partial \mathcal{L}}{\partial \mathbf{b}^{[1:L]}}$

# Takeaways

- MLPs consist of stacks of perceptron units

- MLPs can learn complex decision boundaries by composing simple features into more complex features

- Learn MLP weights with gradient descent
  - Backpropagation efficiently computes gradient