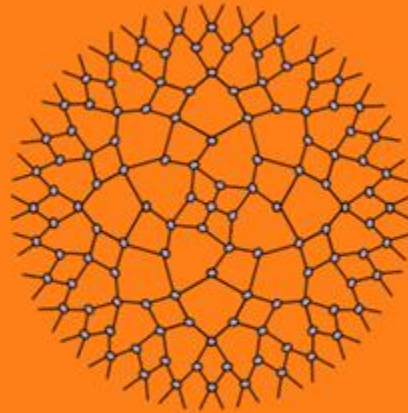


# **ML Algorithms**

# NEURAL NETWORKS



# **Class**

## A Detailed Look At Neural Networks



**Topic**  
Gradient Descent;  
Partial Derivatives



# Gradient Descent

Data

Features	Response
$\mathbf{x}_1 = x_{11}, x_{12}, \dots x_{1m}$	$y_1$
$\mathbf{x}_2 = x_{21},$ $x_{22}, \dots x_{2m}$	$y_2$
$\cdot$ $\cdot$ $\cdot$	$\cdot$ $\cdot$ $\cdot$
$\mathbf{x}_n = x_{n1},$ $x_{n2}, \dots x_{nm}$	$y_n$

- Parameters  $w_0, w_1, w_2, \dots, w_m$
- Cost Function:  $C = C(\mathbf{w})$
- Partial Derivatives:  $\frac{\partial C(\mathbf{w})}{\partial w_j}$



# Partial Derivatives

---

Partial derivatives are useful for multivariate functions

---

Functions with multiple variables: A partial derivative of the function with respect to any one of the variables, measures the change in the value of the function for a small change in the value of that variable, keeping all other variables constant



# Partial Derivatives

How does a function change when one variable is changed (others remaining fixed)?

## Example

- $f(x_1, x_2) = 20 + a x_1 + b x_2^2$
- $\frac{\partial f}{\partial x_1} = a$
- $\frac{\partial f}{\partial x_2} = 2b x_2$
- $\nabla f = \begin{pmatrix} a \\ 2b x_2 \end{pmatrix}$




# Partial Derivatives

How does a function change when one variable is changed (others remaining fixed)?

## Example

- $f(x_1, x_2) = 20 + a x_1 + b x_2^2$
- $\frac{\partial f}{\partial x_1} = a$
- $\frac{\partial f}{\partial x_2} = 2b x_2$

Gradient of the  
multivariate  
function ' $f$ '



- $\nabla f = \begin{pmatrix} a \\ 2b x_2 \end{pmatrix}$



# Partial Derivatives for Cost Function

The cost function is the sum of all the costs accrued in all the data points

- $C = \sum_i C_i(\mathbf{w})$
- $C_i(\mathbf{w}) = y_i \log \frac{1}{1 + \exp(\mathbf{w} \cdot \mathbf{x}_i)} - (1 - y_i) \log \left( 1 - \frac{1}{1 + \exp(\mathbf{w} \cdot \mathbf{x}_i)} \right)$
- $\frac{\partial C(\mathbf{w})}{\partial w_j} = \sum_i (h(\mathbf{x}_i) - y_i) x_{ij}$



# Gradient Descent

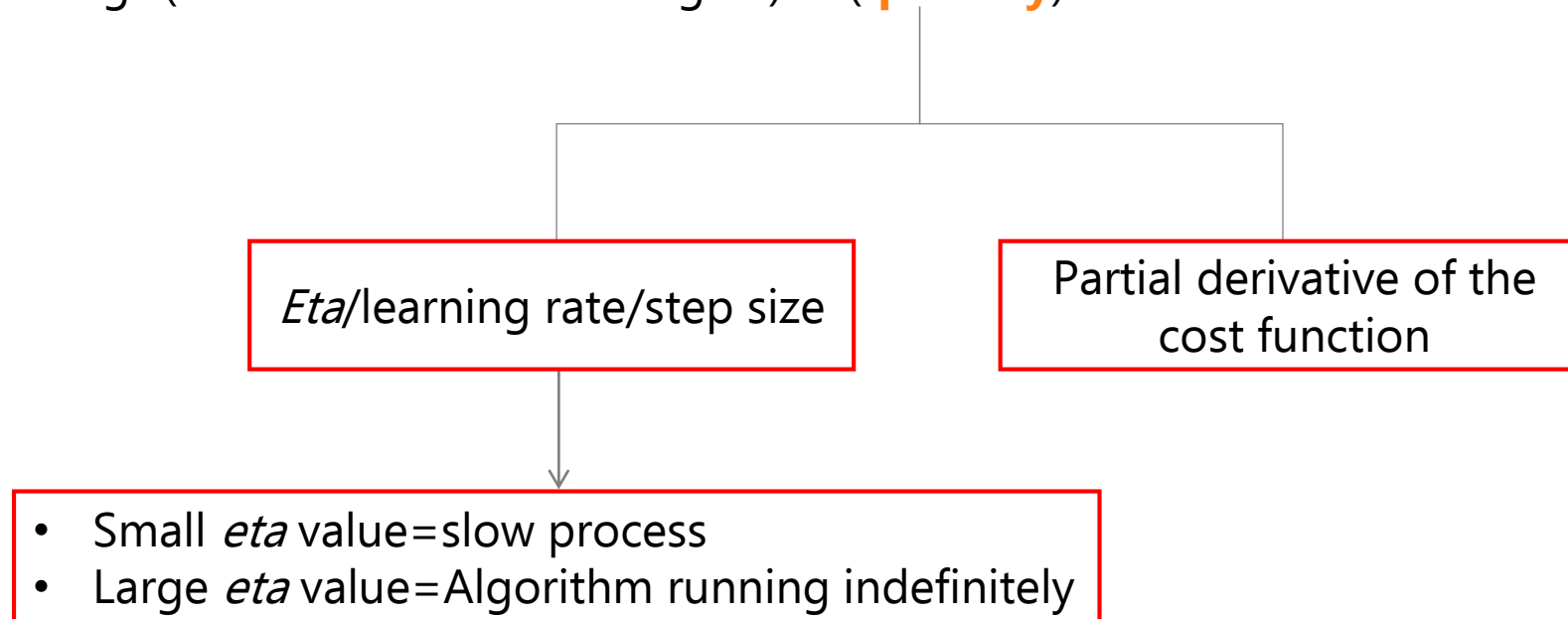
- Initialize  $\mathbf{w}$  to some random values
- Example:  $\mathbf{w} = [1, 4, 7, 0, \dots]$
- Repeat
  - $w_j := w_j - \eta \frac{\partial C(\mathbf{w})}{\partial w_j}$  for all  $j$
  - Equivalently  $\mathbf{w} := \mathbf{w} - \eta \nabla C$  ( $\eta$ : Learning Rate)
- Once initialized, the weights are continuously updated
- Updating = (Current value of the weights) - (quantity)





# Gradient Descent

Updating=(Current value of the weights) — (**quantity**)



The gradient descent reaches the global minima of a function if a small value of *eta*



# Gradient Descent: Example

Data

<b>x1</b>	<b>x2</b>	<b>y</b>
-2	0.5	0
2.5	-2	1

Initialize weights

$$w_0 = -2$$

$$w_1 = 1.5$$

$$w_2 = 3.5$$



# Gradient Descent: Example

Data

<b>x1</b>	<b>x2</b>	<b>y</b>	<b>h(x)</b>	<b>Cost</b>
-2	0.5	0	0.037	0.017
2.5	-2	1	0.0052	2.282



# Gradient Descent: Example

Data

<b>x1</b>	<b>x2</b>	<b>y</b>	<b>h(x)</b>	<b>Cost</b>
-2	0.5	0	0.037	0.017
2.5	-2	1	0.0052	2.282

- $C([-2, 1.5, 3.5]) = 2.299$
- $\frac{\partial C(w)}{\partial w_0} = (0.037 - 0) \times 1 + (0.0052 - 1) \times 1 = 0.96$
- $\frac{\partial C(w)}{\partial w_1} = (0.037 - 0) \times (-2) + (0.0052 - 1) \times 2.5 = -2.56$
- $\frac{\partial C(w)}{\partial w_2} = (0.037 - 0) \times 0.5 + (0.0052 - 1) \times (-2) = 2.008$



# Gradient Descent: Example

## Gradient Descent: Step 1

- $w_0 := -2 - 0.01 \times (-0.96) = -1.99$
- $w_1 := 1.5 - 0.01 \times (-2.56) = 1.53$
- $w_2 := 3.5 - 0.01 \times (-2.008) = 3.48$

New Value of Cost Function: 0.016 (Verify)



# Gradient Descent: Example

## Gradient Descent: Step 2

- $w_0 := -2 - 0.01 \times (-0.96) = -1.99$
- $w_1 := 1.5 - 0.01 \times (-2.56) = 1.53$
- $w_2 := 3.5 - 0.01 \times (-2.008) = 3.48$

New Value of Cost Function: 0.016 (Verify)



# Gradient Descent: Example

## Gradient Descent: Step 3

- $w_0 := -2 - 0.01 \times (-0.96) = -1.99$
- $w_1 := 1.5 - 0.01 \times (-2.56) = 1.53$
- $w_2 := 3.5 - 0.01 \times (-2.008) = 3.48$

New Value of Cost Function: 0.016 (Verify)



# Gradient Descent

- Repeat these steps until the weights stop experiencing significant change
- Gradient descent will always find the optimum value
- Gradient descent can be slow
- Algorithms could converge slowly, in spite of reasonable learning rate value
- Alternative: Stochastic Gradient Descent





# Recap

- Gradient Descent
- Partial Derivatives
- Gradient Descent: Example





**JIGSAW ACADEMY**  
THE ONLINE SCHOOL OF ANALYTICS