

Gradient Descent

Gradient Descent

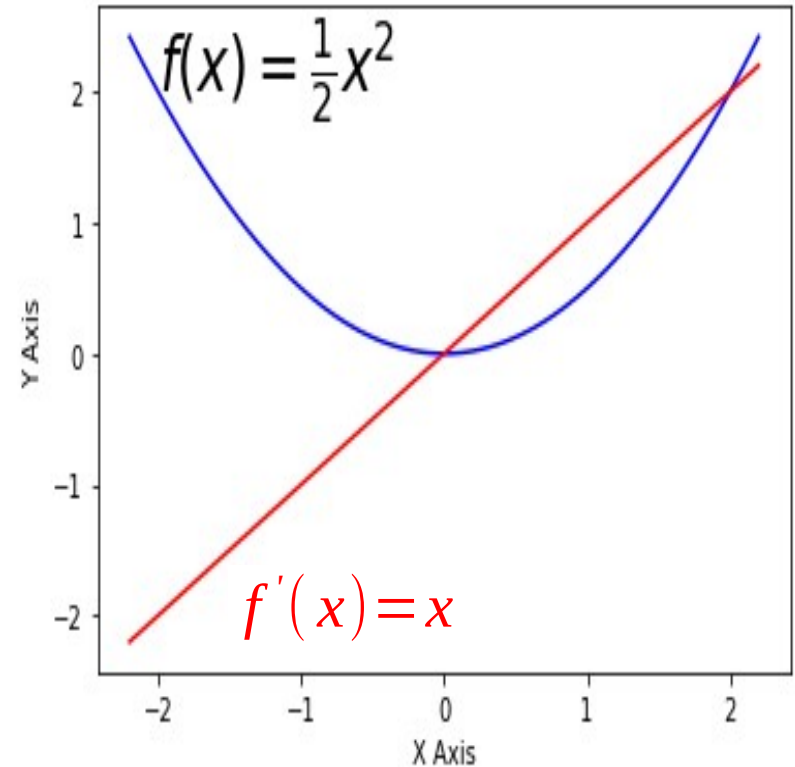
- In previous two examples OR and XOR we assumed some values for model parameters(w , b).
- But those parameters learned by network by using learning algorithm Gradient Descent.

Gradient-single variable function

- Optimization refers to the task of either minimizing or maximizing some function $f(x)$ by altering x .
- Suppose we have a function $y = f(x)$
- The derivative of this function is denoted as $f'(x)$ or as dy/dx .
- The derivative $f'(x)$ gives the slope of $f(x)$ at the point x
- In other words, it specifies how to scale a small change in the input in order to obtain the corresponding change in the output: $f(x+h) \approx f(x) + f'(x)h$

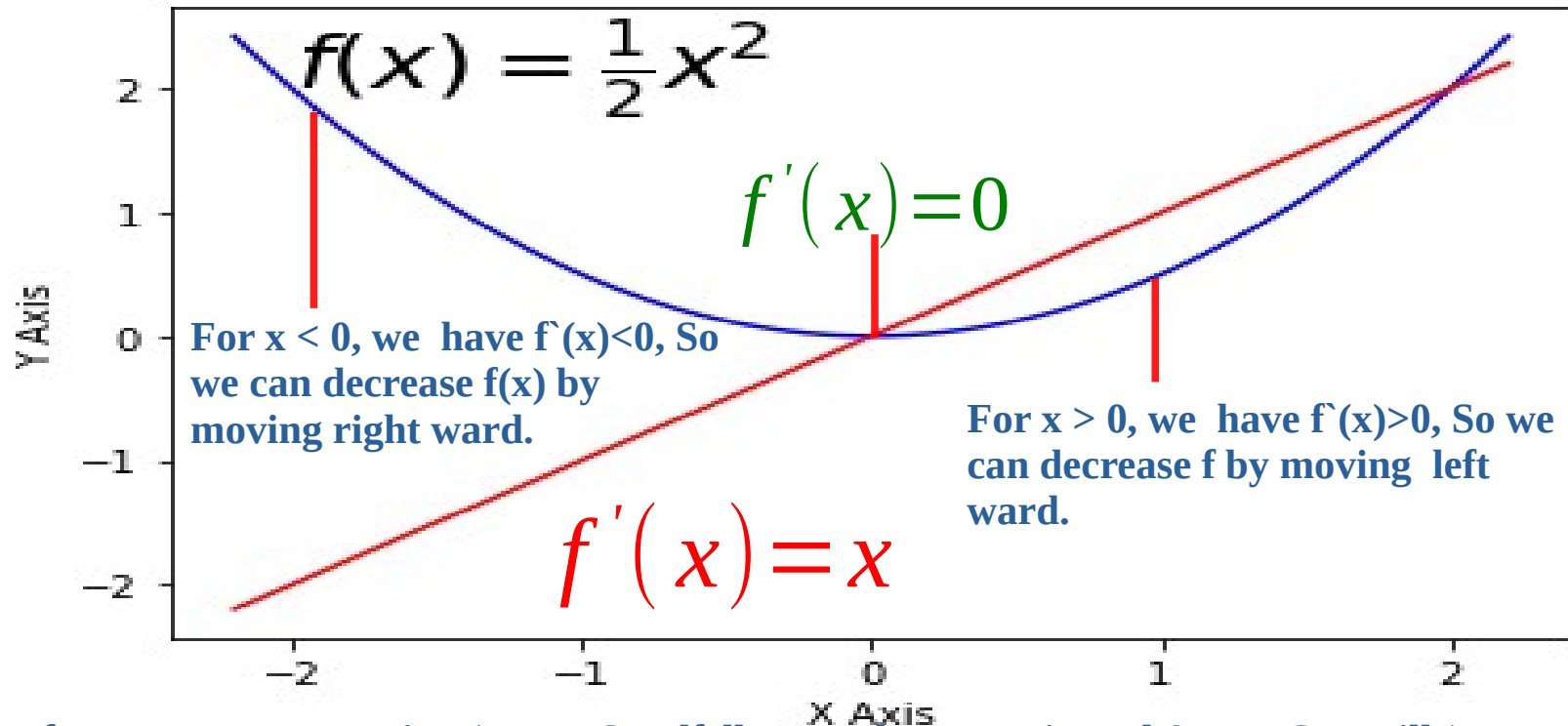
Gradient-single variable function

- If $f(x)$ is differentiable at x
- $f(x)$ will increase if its derivative is positive
- $f(x)$ will decrease if its derivative is negative
- So the derivative is useful for minimizing a function
- Because it tells us how to change x in order to make a small improvement in $f(x)$



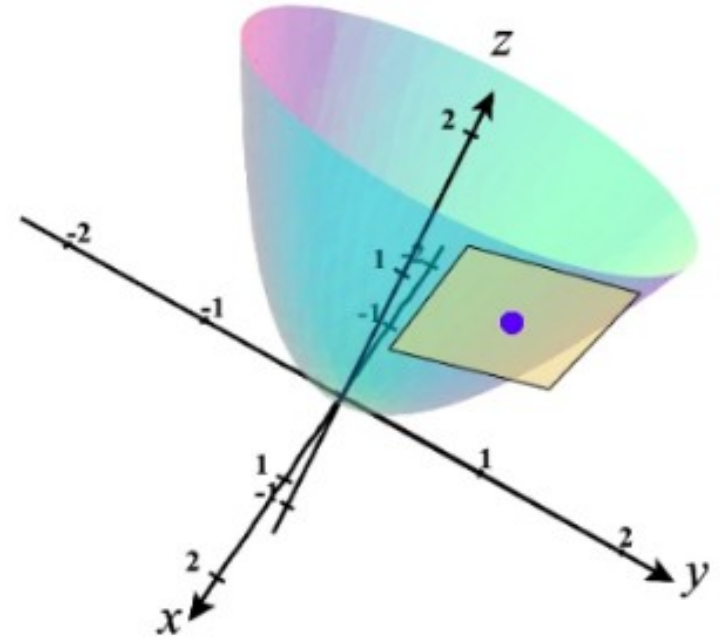
Gradient-single variable function

- Here for $f(x)$ global minimum is at $x = 0$ Since
- Derivative will be zero at local minimum or maximum



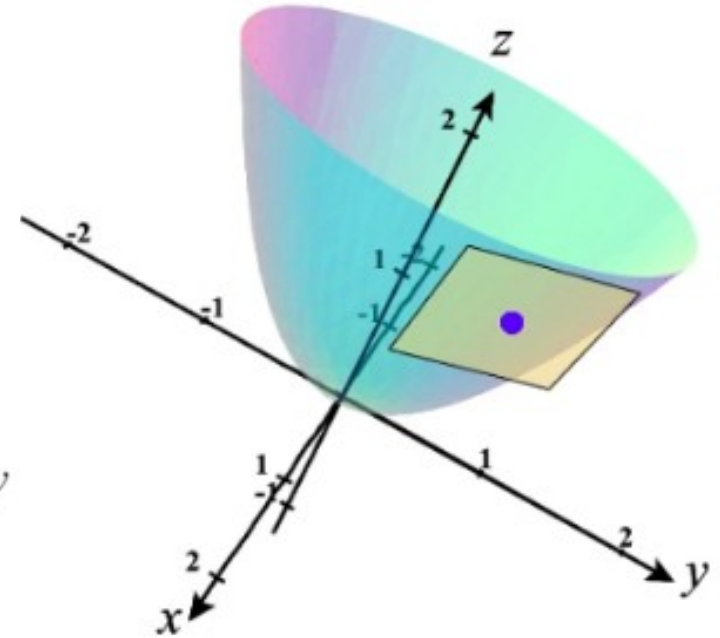
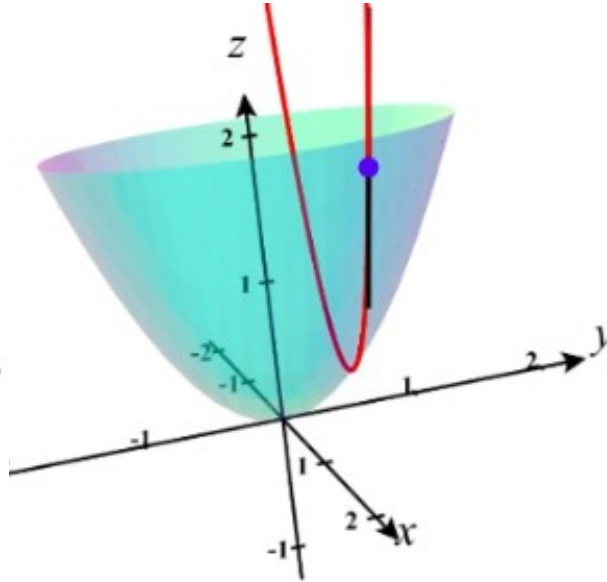
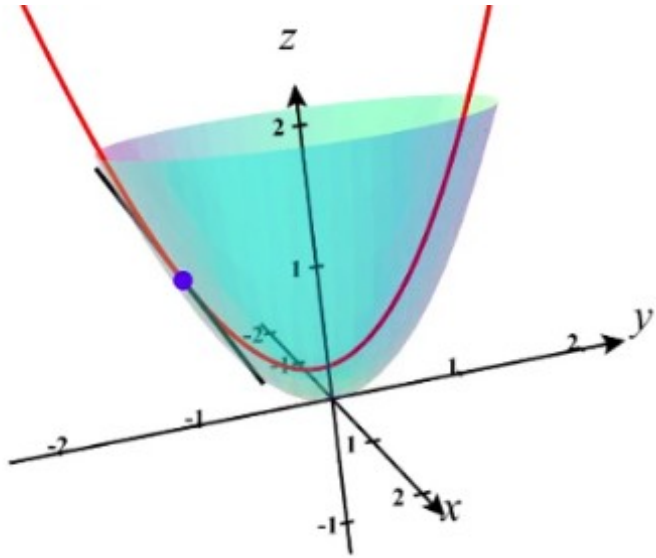
Gradient-Multivariate function

- Gradient is a kind of slope for higher dimensional function
- The gradient of $f(x,y)$ to be the 'slope' of the tangent plane.
- But a plane doesn't have a single slope;
- it slopes differently in different directions



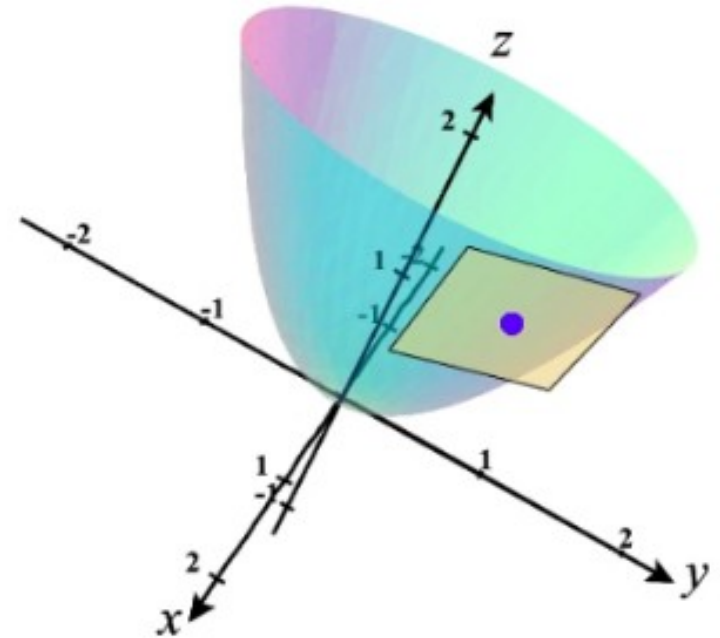
Gradient-Multivariate function

- It slopes differently in different directions



Gradient-Multivariate function

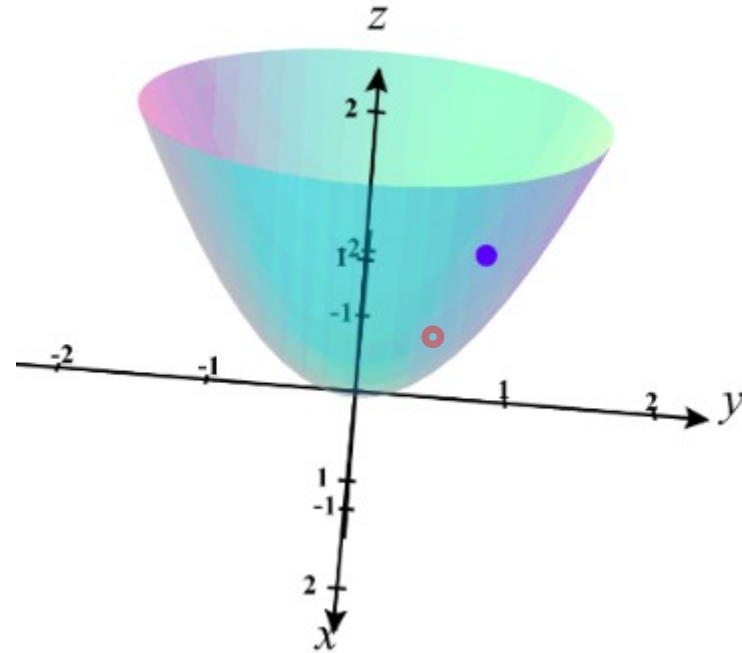
- What is the 'slope' of the plane?
- the partial derivatives in all directions, taken together, form the slope of the plane
- $f(x,y)$ has two variables and two partial derivatives: $\partial z / \partial x, \partial z / \partial y$
- We'll call that vector of these partial derivatives the gradient of f , denoted ∇f
- The gradient of a multi-variable function has a component for each direction.



Gradient-Multivariate function

- The change in height, Δz , is given by
- $\Delta z = (\text{change in } x \text{ direction}) + (\text{change in } y \text{ direction})$
- Therefore

$$\Delta z = \frac{\partial z}{\partial x} \Delta x + \frac{\partial z}{\partial y} \Delta y$$



Gradient-Based Optimization

- If you imagine standing at a point (x_0, y_0) in the input space of f
- the vector ∇f tells you which direction you should travel to increase the value of f most rapidly.
- **It is a vector that points in the direction in which the function grows the fastest.**
- Its coordinates are partial derivatives of that function



$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix}$$

Gradient-Based Optimization

- Consider a simple multiplication function of two numbers $f(x, y) = xy$
- The partial derivative $\frac{\partial f(x, y)}{\partial x} = \frac{\partial f}{\partial x} = y$ measures how f changes as only the variable x increases.
- The partial derivative $\frac{\partial f(x, y)}{\partial y} = \frac{\partial f}{\partial y} = x$ measures how f changes as only the variable y increases.

Gradient-Based Optimization

- The derivative on each variable tells you the sensitivity of the whole expression on its value.
- For example, *if $x=4$, $y=-3$ then $f(x, y)=-12$*
- and the derivative on x $\frac{\partial f}{\partial x}=y=-3$
- This tells us that if we were to increase the value of this variable by a tiny amount h , the effect on the whole expression would be to decrease it (due to the negative sign), and by three times that amount ($-3h$).

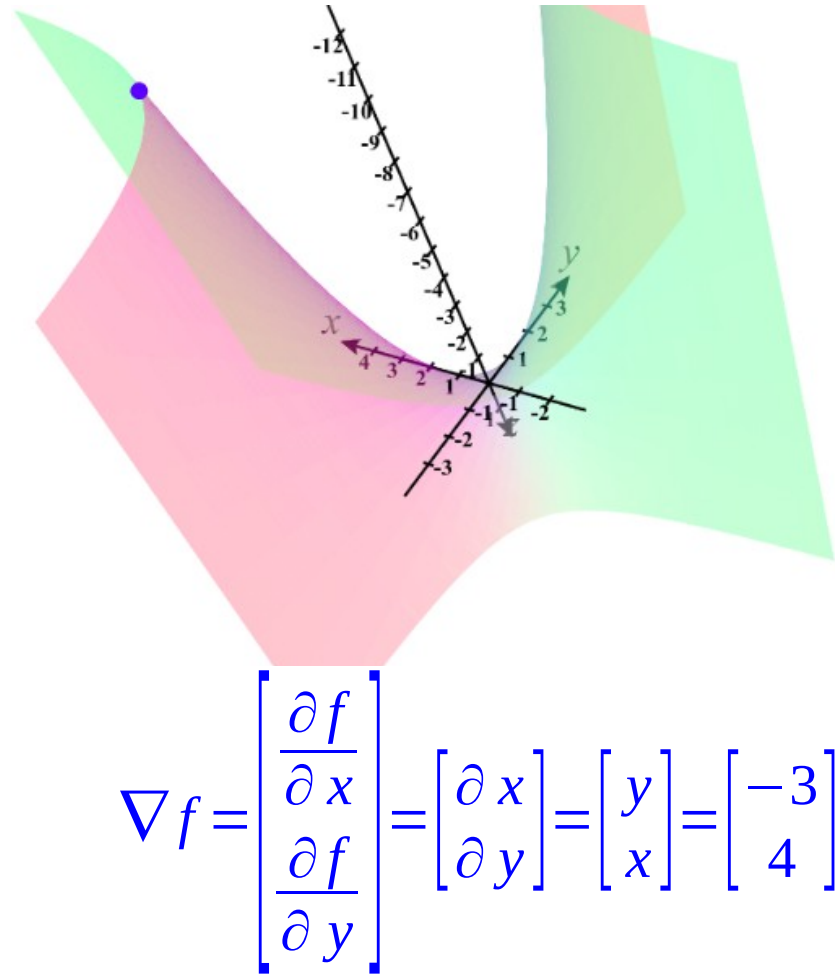
Gradient-Based Optimization

- Similarly, since $\frac{\partial f}{\partial y} = x = 4$, we expect that increasing the value of y by some very small amount h would also increase the output of the function (due to the positive sign), and by $4h$.

$$\nabla f = \begin{bmatrix} \frac{\partial f(x, y)}{\partial x} \\ \frac{\partial f(x, y)}{\partial y} \end{bmatrix} = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} = \begin{bmatrix} y \\ x \end{bmatrix} = \begin{bmatrix} -3 \\ 4 \end{bmatrix}$$

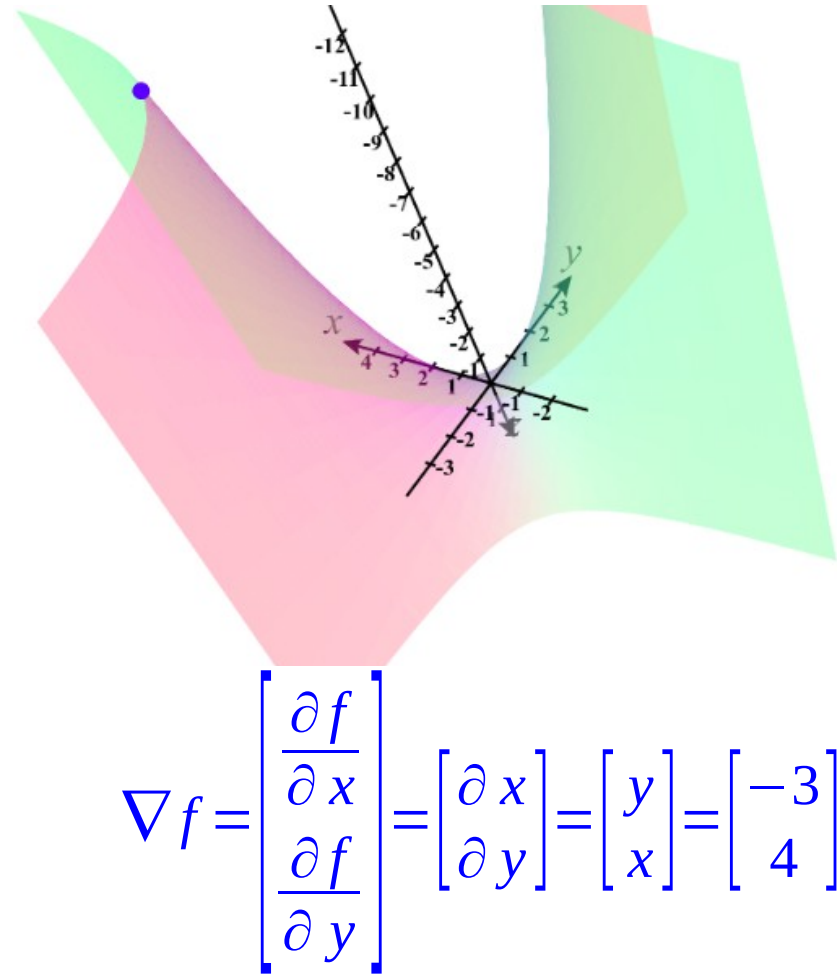
Gradient-Based Optimization

- At point $(4, -3)$ $f(x, y) = -12$
- Gradient tells us in which direction we need to move to maximize the function
- Let us take a small step in the direction of the gradient
- $x_{new} = x + \Delta x = x + h \cdot \frac{\partial f}{\partial x}$
 $x_{new} = x + \Delta x = 4 + (0.2) \cdot (-3)$
 $x_{new} = x + \Delta x = 4 - 0.6 = 3.4$



Gradient-Based Optimization

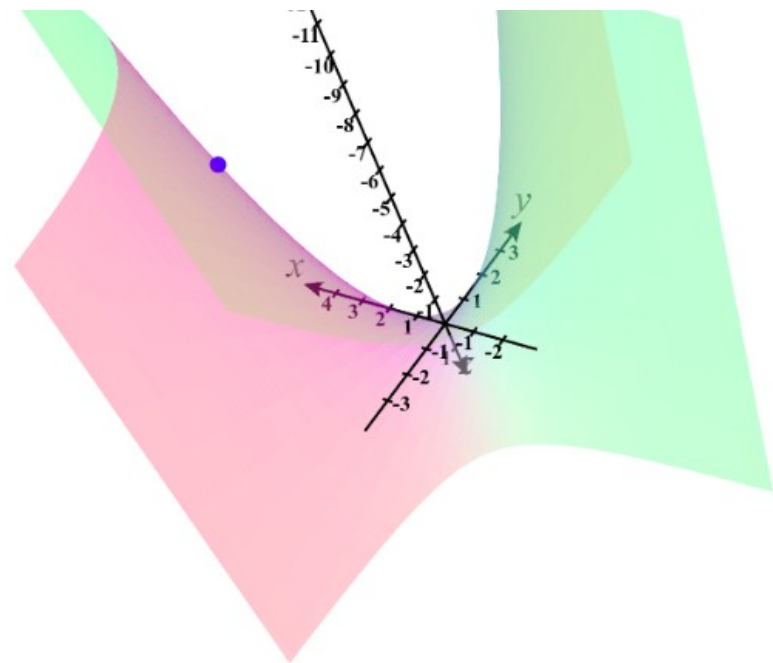
- At point $(4, -3)$ $f(x, y) = -12$
- Gradient tells us in which direction we need to move to maximize the function
- Let us take a small step in the direction of the gradient
- $y_{new} = y + \Delta y = y + h \cdot \frac{\partial f}{\partial y}$
 $y_{new} = y + \Delta y = -3 + (0.2) \cdot 4$
 $y_{new} = y + \Delta y = -3 + 0.8 = -2.2$



Gradient-Based Optimization

- At point $(3.4, -2.2)$ $f(x,y) = -7.48$
- $x_{new} = x + \Delta x = 3.4 + (0.2) \cdot (-2.2)$
 $x_{new} = x + \Delta x = 4 - 0.44 = 2.96$

$$y_{new} = y + \Delta y = -2.2 + (0.2) \cdot (3.4)$$
$$y_{new} = y + \Delta y = -2.2 + 0.68 = -1.52$$

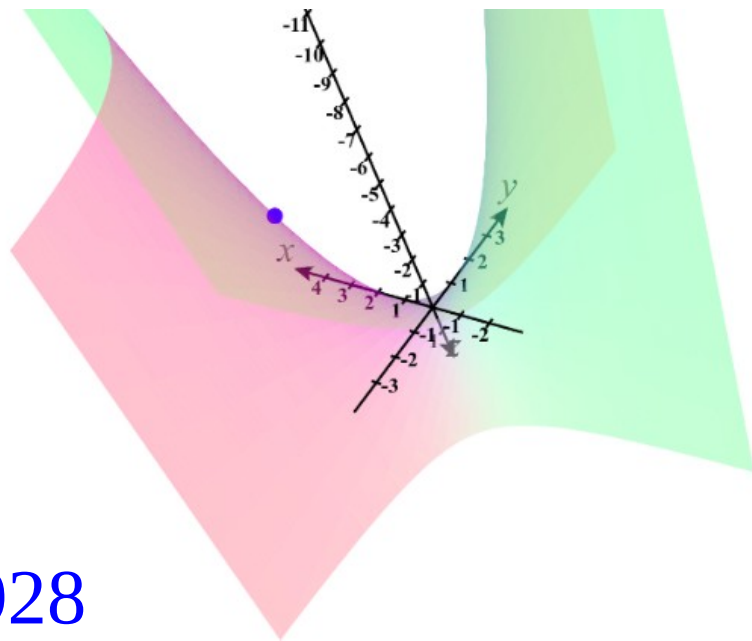


$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} = \begin{bmatrix} \partial x \\ \partial y \end{bmatrix} = \begin{bmatrix} y \\ x \end{bmatrix} = \begin{bmatrix} -2.2 \\ 3.4 \end{bmatrix}$$

Gradient-Based Optimization

- At point $(2.96, -1.52)$ $f(xy) = -4.49$
- $x_{new} = x + \Delta x = 2.96 + (0.2) \cdot (-1.52)$
 $x_{new} = x + \Delta x = 2.96 - 0.304 = 2.656$

$$y_{new} = y + \Delta y = -1.52 + (0.2) \cdot (2.96)$$
$$y_{new} = y + \Delta y = -1.52 + 0.592 = -0.928$$



$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} = \begin{bmatrix} \partial x \\ \partial y \end{bmatrix} = \begin{bmatrix} y \\ x \end{bmatrix} = \begin{bmatrix} -1.52 \\ 2.96 \end{bmatrix}$$

Gradient-Descent

- Consider Linear regression Cost function

$$MSE = J(\theta) = J(w, b) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad y = wx + b$$

- Here we have two parameters w and b

$$\Delta J = \frac{\partial J}{\partial w} \Delta w + \frac{\partial J}{\partial b} \Delta b$$

- The gradient of J is the vector of partial derivatives

$$\nabla J = \begin{bmatrix} \frac{\partial J}{\partial w} \\ \frac{\partial J}{\partial b} \end{bmatrix}$$

Gradient-Descent

- Consider Linear regression Cost function

$$MSE = J(\theta) = J(w, b) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad y = wx + b$$

- Here we have two parameters w and b

$$\Delta J = \frac{\partial J}{\partial w} \Delta w + \frac{\partial J}{\partial b} \Delta b \quad J(w, b) = \frac{1}{n} \sum_{i=1}^n (y_i - (wx_i + b))^2$$

- The gradient of J is the vector of partial derivatives

$$\nabla J = \begin{bmatrix} \frac{\partial J}{\partial w} \\ \frac{\partial J}{\partial b} \end{bmatrix} \quad \nabla J = \begin{bmatrix} \frac{\partial J}{\partial w} \\ \frac{\partial J}{\partial b} \end{bmatrix} = \begin{bmatrix} -\frac{2}{n} \sum_{i=1}^n x_i (y_i - (wx_i + b)) \\ -\frac{2}{n} \sum_{i=1}^n (y_i - (wx_i + b)) \end{bmatrix}$$

Gradient-Descent

- ΔJ can be rewritten as $\Delta J = \nabla J \cdot \Delta \theta$ here $\Delta \theta = \begin{bmatrix} \Delta w \\ \Delta b \end{bmatrix}$
 $\Delta J = \frac{\partial J}{\partial w} \Delta w + \frac{\partial J}{\partial b} \Delta b$ $\theta = \begin{bmatrix} w \\ b \end{bmatrix}$

- ∇J relates changes in θ to changes in J
- lets us see how to choose $\Delta \theta$ so as to make ΔJ negative

$$J = J + \Delta J$$

- suppose we choose $\Delta \theta = -\alpha \nabla J$

$$\Delta J = -\alpha \nabla J \cdot \nabla J = -\alpha \|\nabla J\|^2$$

Gradient descent

- Because $\|\nabla J\|^2 \geq 0$ this guarantees that $\Delta J \leq 0$
- Means J will always decrease, never increase, if we change θ according to the $\Delta \theta = -\alpha \nabla J$

$$\Delta w = -\alpha \frac{\partial J}{\partial w}$$

$$\Delta b = -\alpha \frac{\partial J}{\partial b}$$

$$w_{new} = w + \Delta w = w - \alpha \frac{\partial J}{\partial w}$$

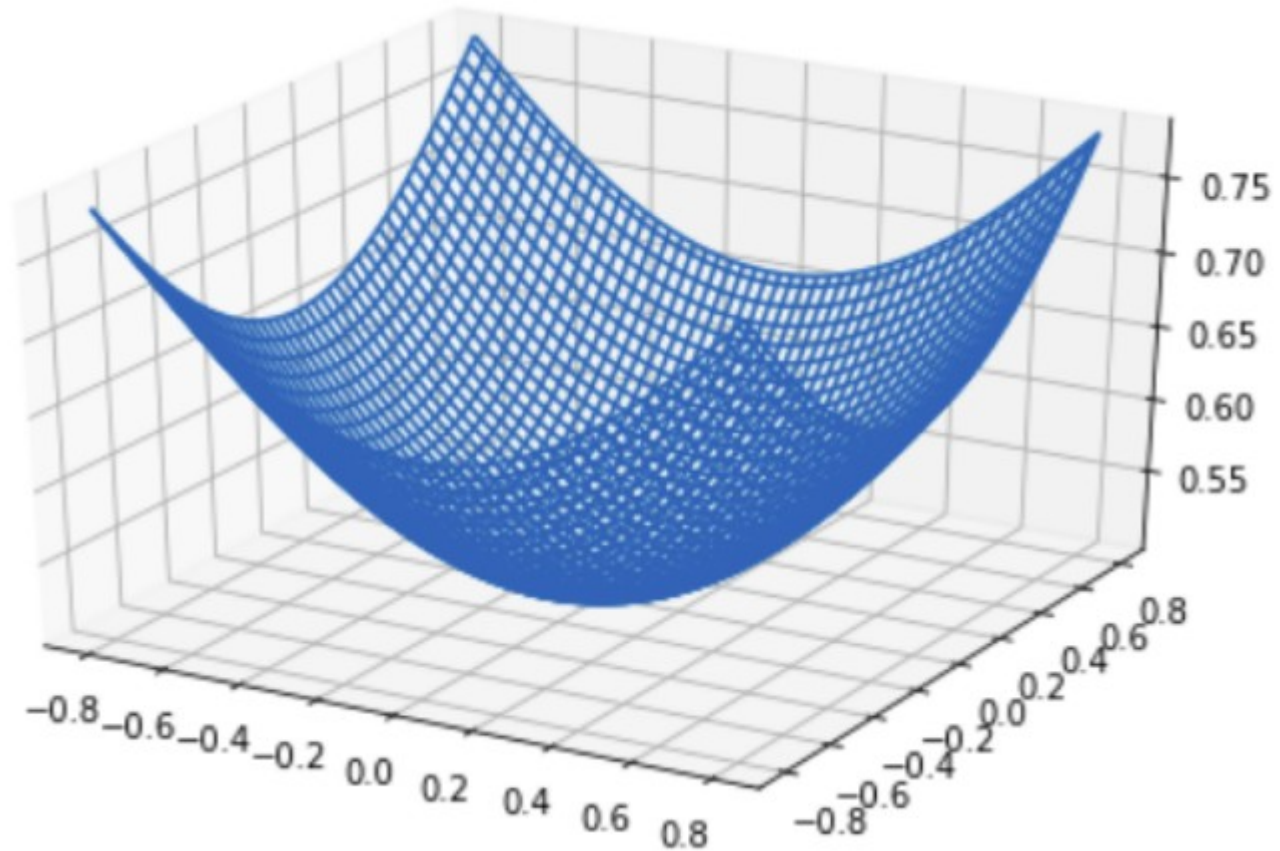
$$b_{new} = b + \Delta b = b - \alpha \frac{\partial J}{\partial b}$$

Gradient descent

$$w_{new} = w - \alpha \frac{\partial J}{\partial w}$$

$$b_{new} = b - \alpha \frac{\partial J}{\partial b}$$

$$\nabla J = \begin{bmatrix} \frac{\partial J}{\partial w} \\ \frac{\partial J}{\partial b} \end{bmatrix}$$



Back-Propagation

- Computing an analytical expression for the gradient is straight forward, but numerically evaluating such an expression can be computationally expensive
- The back propagation algorithm does so using a simple and inexpensive procedure
- The term back-propagation is often misunderstood as meaning the whole learning algorithm for multi-layer neural networks
- Actually, back- propagation refers **only to the method for computing the gradient**

Reference: Deep Learning (Ian J. Goodfellow, Yoshua Bengio and Aaron Courville), MIT Press, 2016.

Back-Propagation

- Gradient descent, is used to perform learning using this gradient.
- In learning algorithms, the gradient we most often require is the gradient of the cost function with respect to the parameters, $\nabla_{\theta} J(\theta)$

Gradient-Based Optimization

- The derivative on each variable tells you the sensitivity of the whole expression on its value.
- For example, *if $x=4$, $y=-3$ then $f(x, y)=-12$*
- and the derivative on x $\frac{\partial f}{\partial x}=y=-3$
- This tells us that if we were to increase the value of this variable by a tiny amount h , the effect on the whole expression would be to decrease it (due to the negative sign), and by three times that amount ($-3h$).

Gradient-Based Optimization

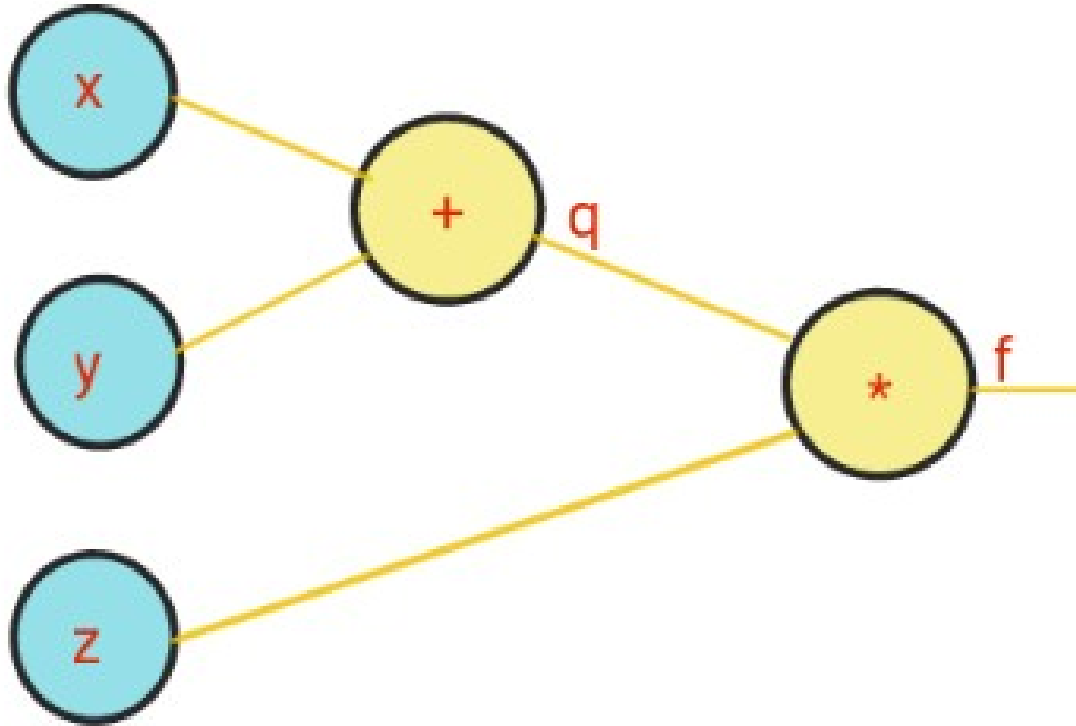
- Similarly, since $\frac{\partial f}{\partial y} = x = 4$, we expect that increasing the value of y by some very small amount h would also increase the output of the function (due to the positive sign), and by $4h$.

$$\nabla f = \begin{bmatrix} \frac{\partial f(x, y)}{\partial x} \\ \frac{\partial f(x, y)}{\partial y} \end{bmatrix} = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} = \begin{bmatrix} y \\ x \end{bmatrix} = \begin{bmatrix} -3 \\ 4 \end{bmatrix}$$

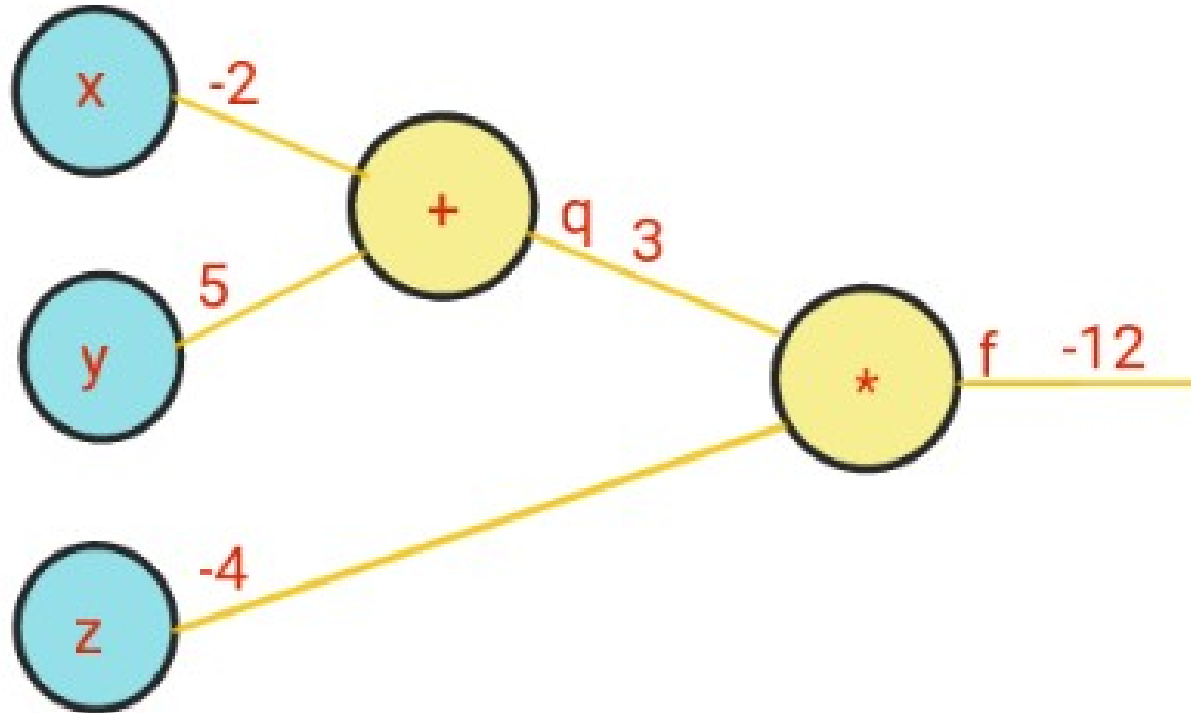
Back Propagation with chain rule

- consider more complicated expressions that involve multiple composed functions $f(x, y, z) = (x + y)z$
- this expression can be broken down into two expressions: $q = (x + y)$ and $f = qz$
- f is just multiplication of q and z , so $\frac{\partial f}{\partial q} = z$, $\frac{\partial f}{\partial z} = q$
- q is addition of x and y so $\frac{\partial q}{\partial x} = 1$, $\frac{\partial q}{\partial y} = 1$
- From chain rule $\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$

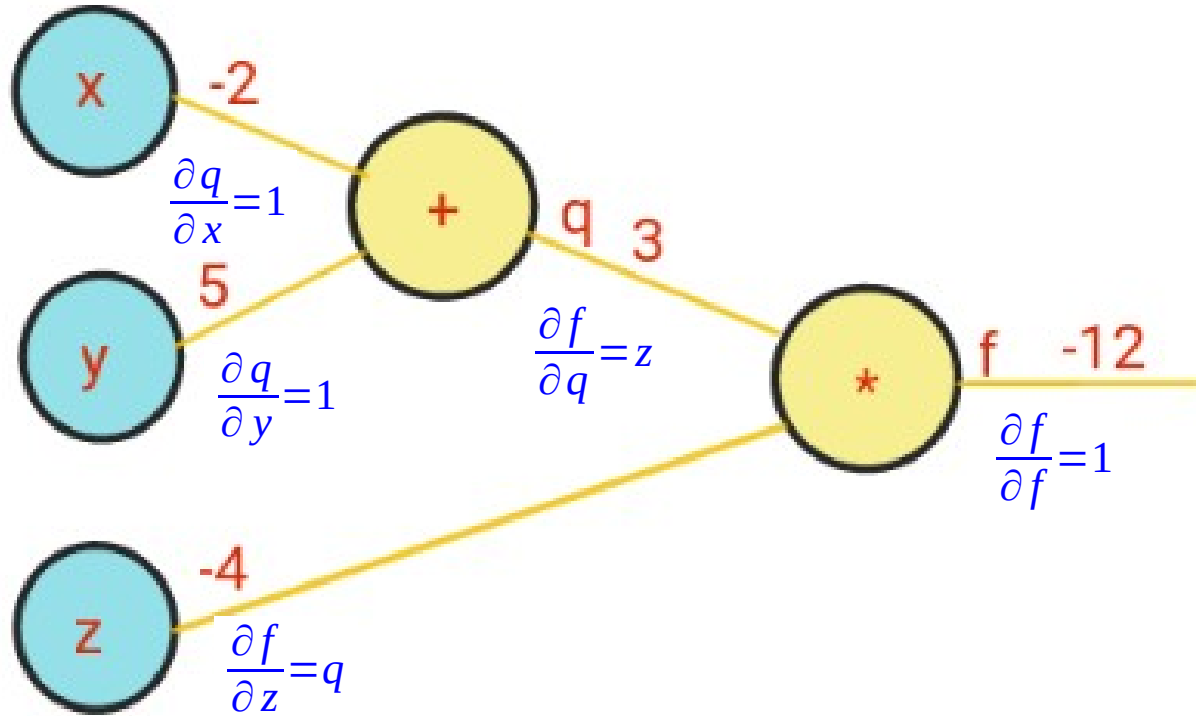
Computational Graph



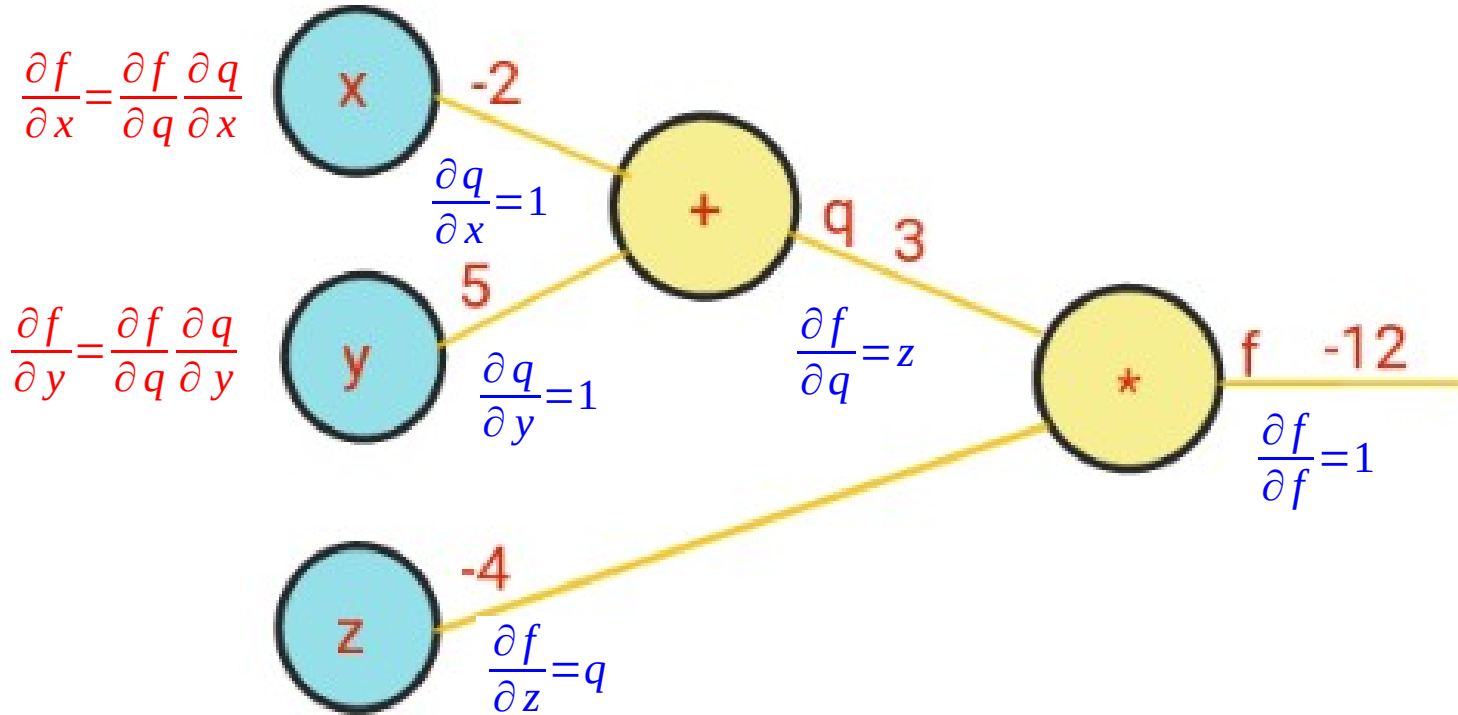
Forward Propagation



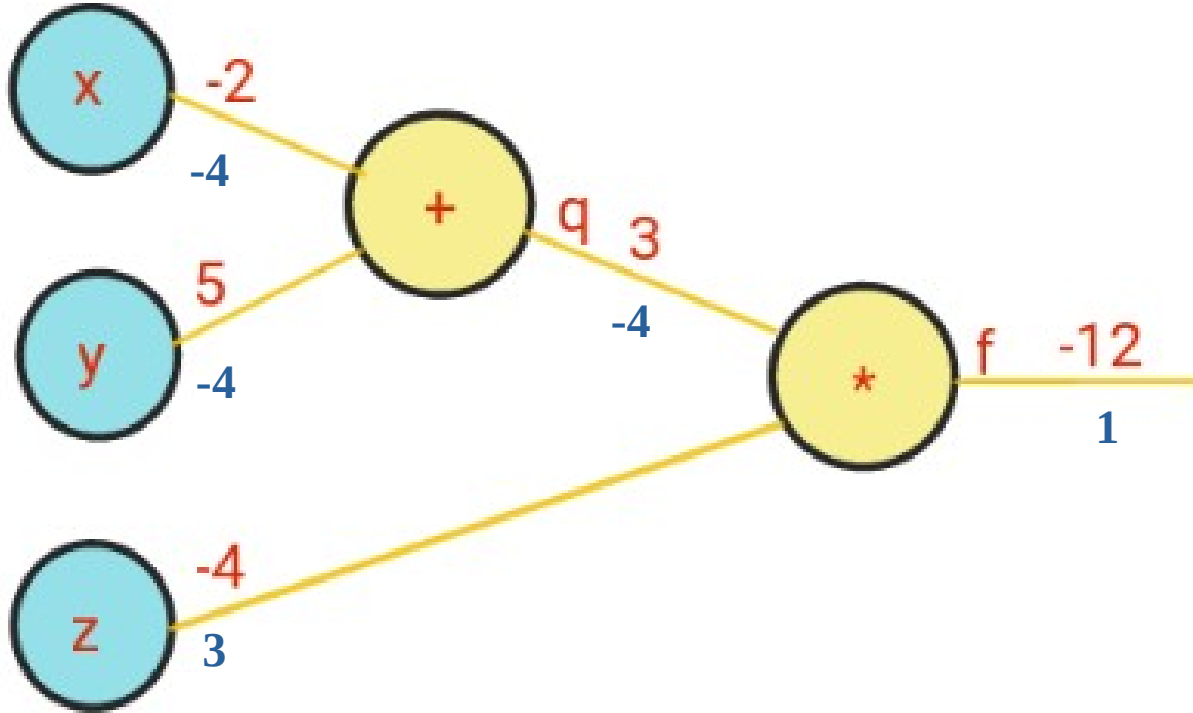
Back Propagation



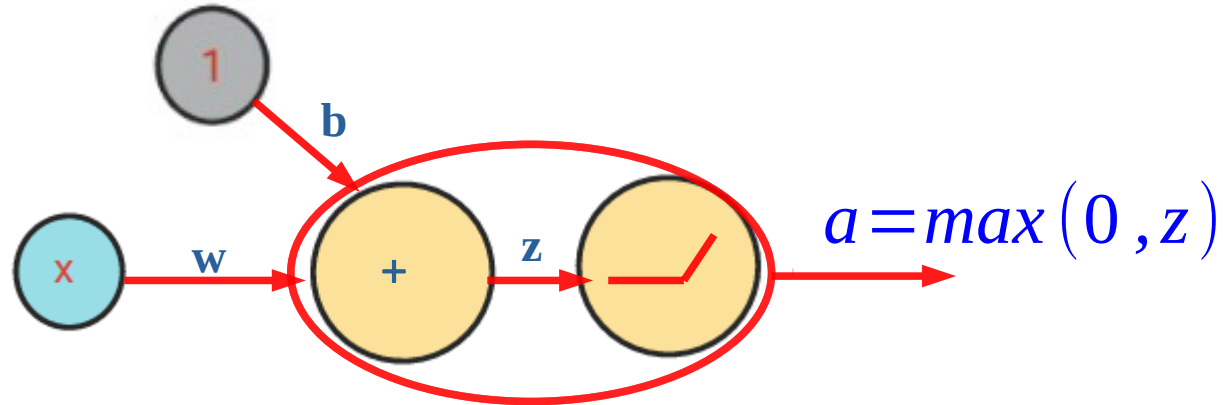
Back Propagation



Back Propagation



Back Propagation (Linear Regression)



$$\hat{y} = a$$

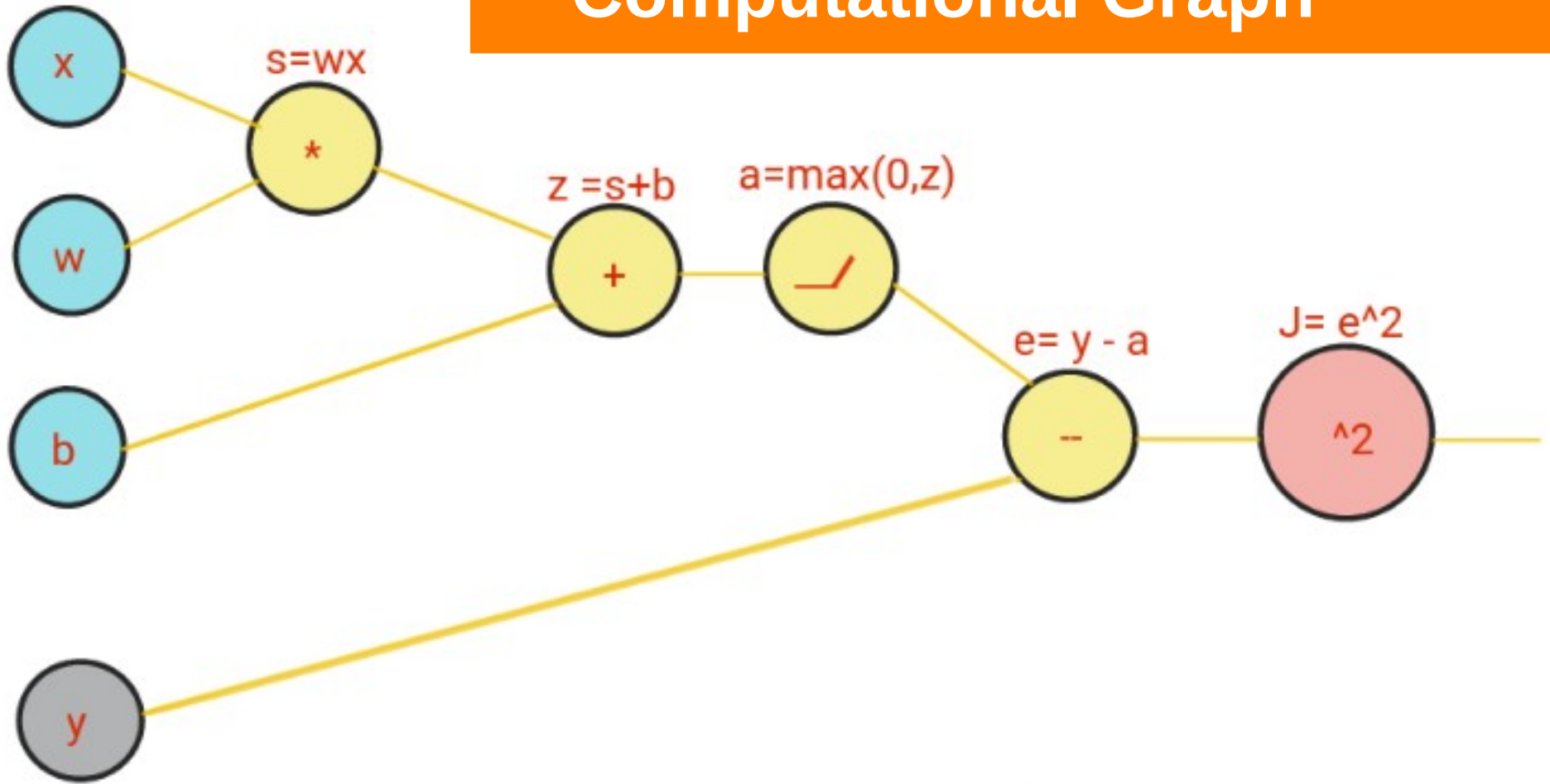
Back Propagation (Linear Regression)

$$J = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2$$

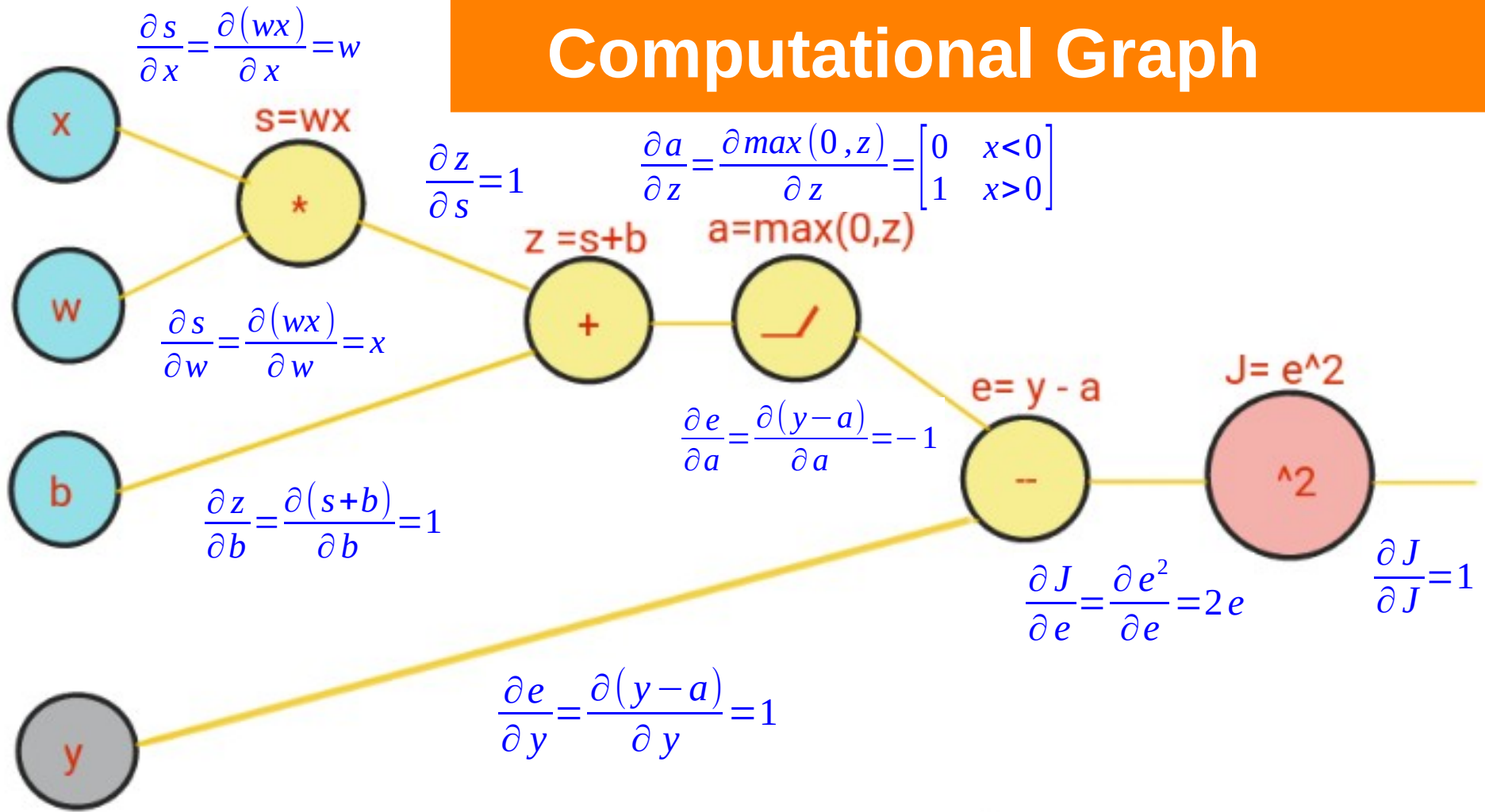
$$\hat{y} = \max(0, wx + b)$$

$$J(w, b) = \frac{1}{n} \sum_{i=1}^n (y_i - \max(0, wx_i + b))^2$$

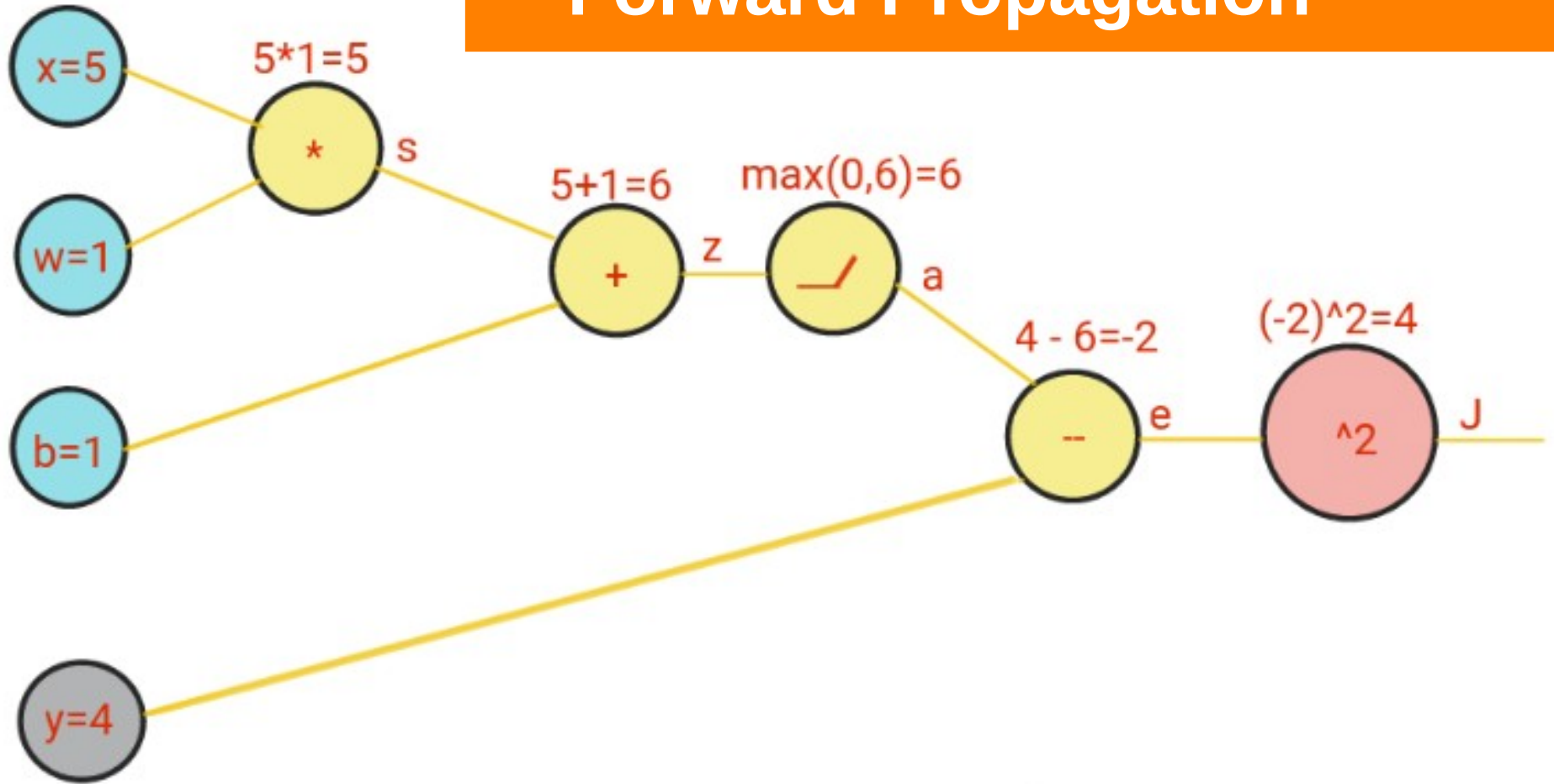
Computational Graph



Computational Graph

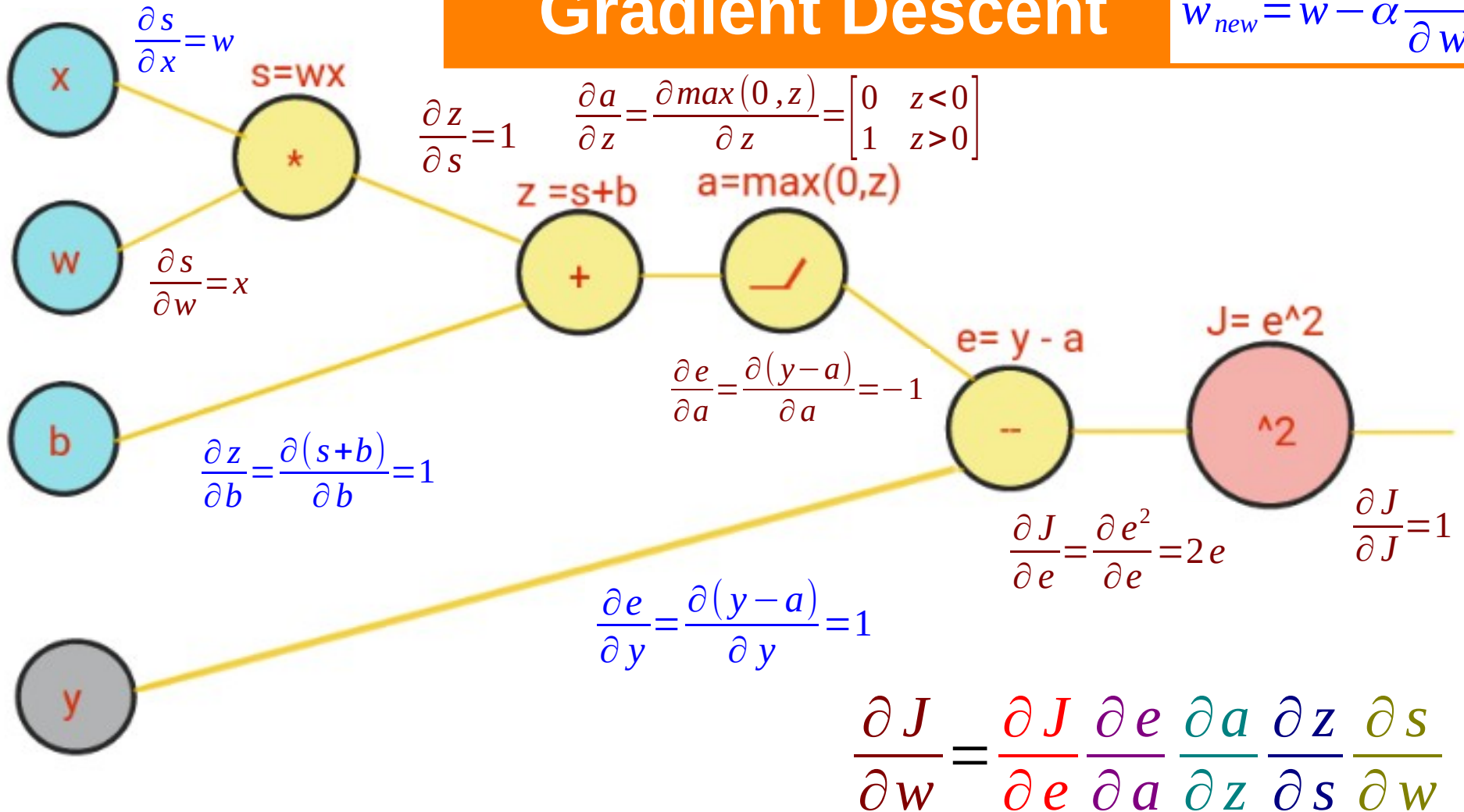


Forward Propagation

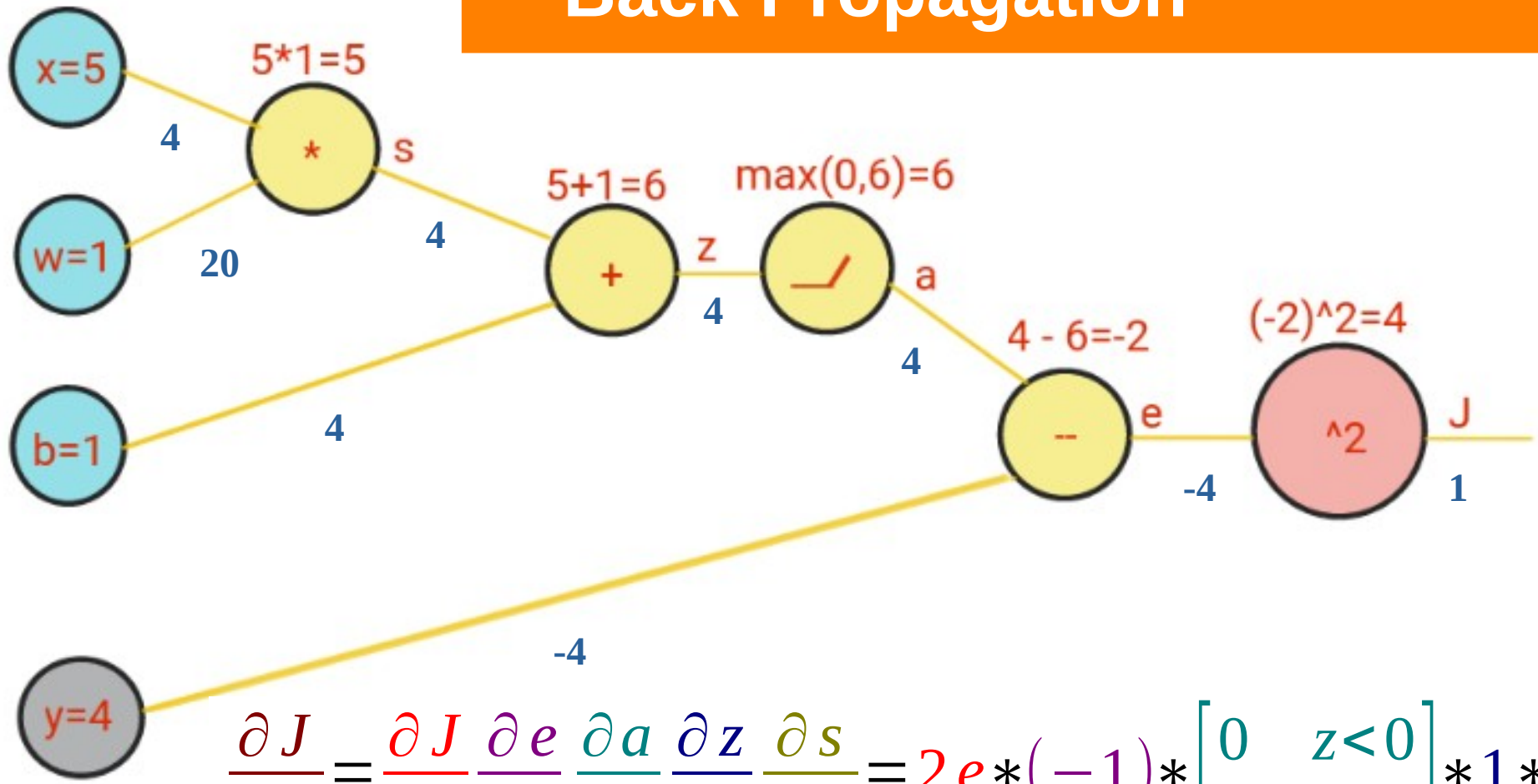


Gradient Descent

$$w_{new} = w - \alpha \frac{\partial J}{\partial w}$$



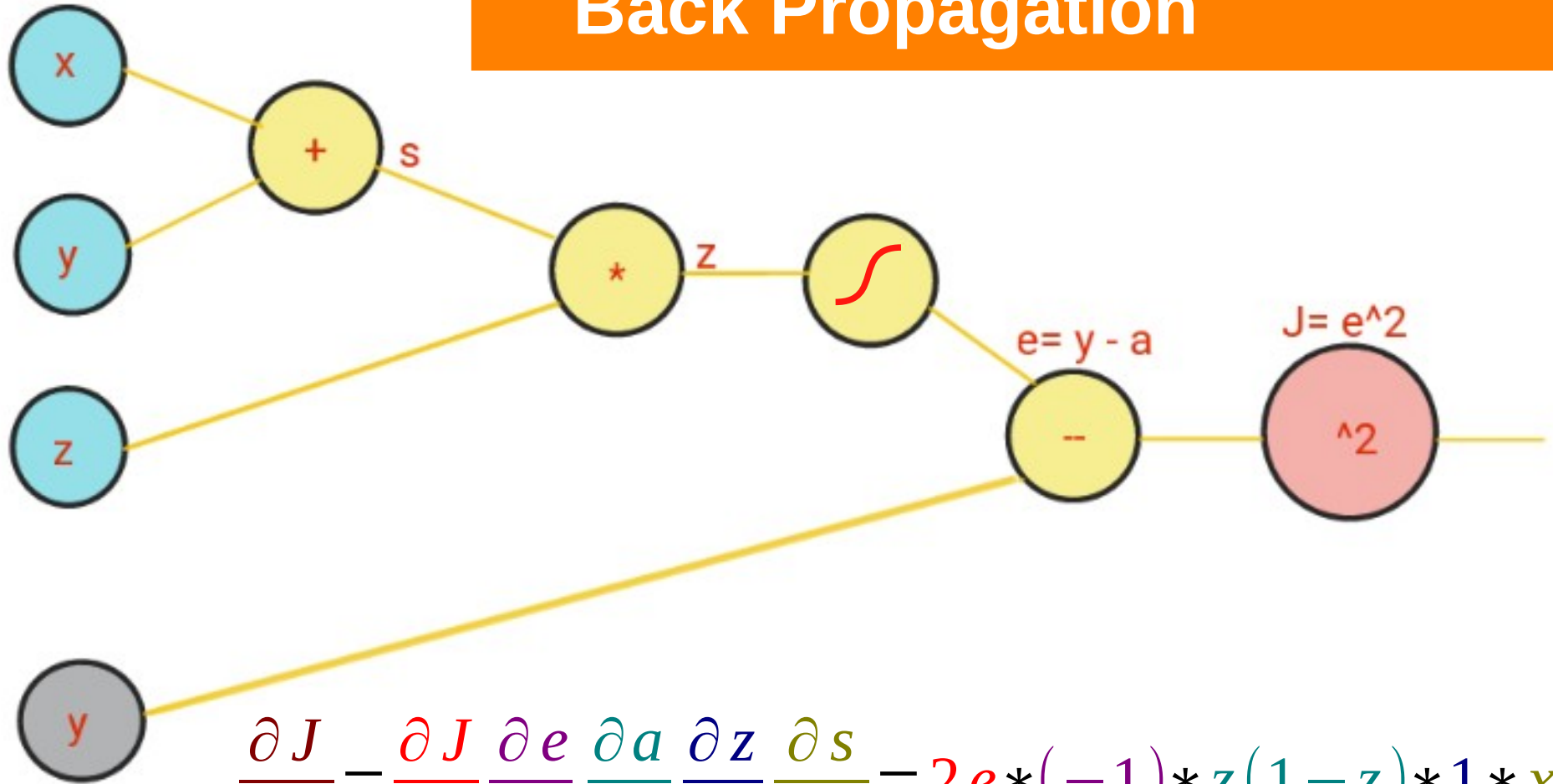
Back Propagation



$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial e} \frac{\partial e}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial s} \frac{\partial s}{\partial w} = 2e * (-1) * \begin{bmatrix} 0 & z < 0 \\ 1 & z > 0 \end{bmatrix} * 1 * x$$

Reference: <https://cs231n.github.io/optimization-2/>

Back Propagation



$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial e} \frac{\partial e}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial s} \frac{\partial s}{\partial w} = 2e * (-1) * z(1-z) * 1 * x$$

Gradient descent

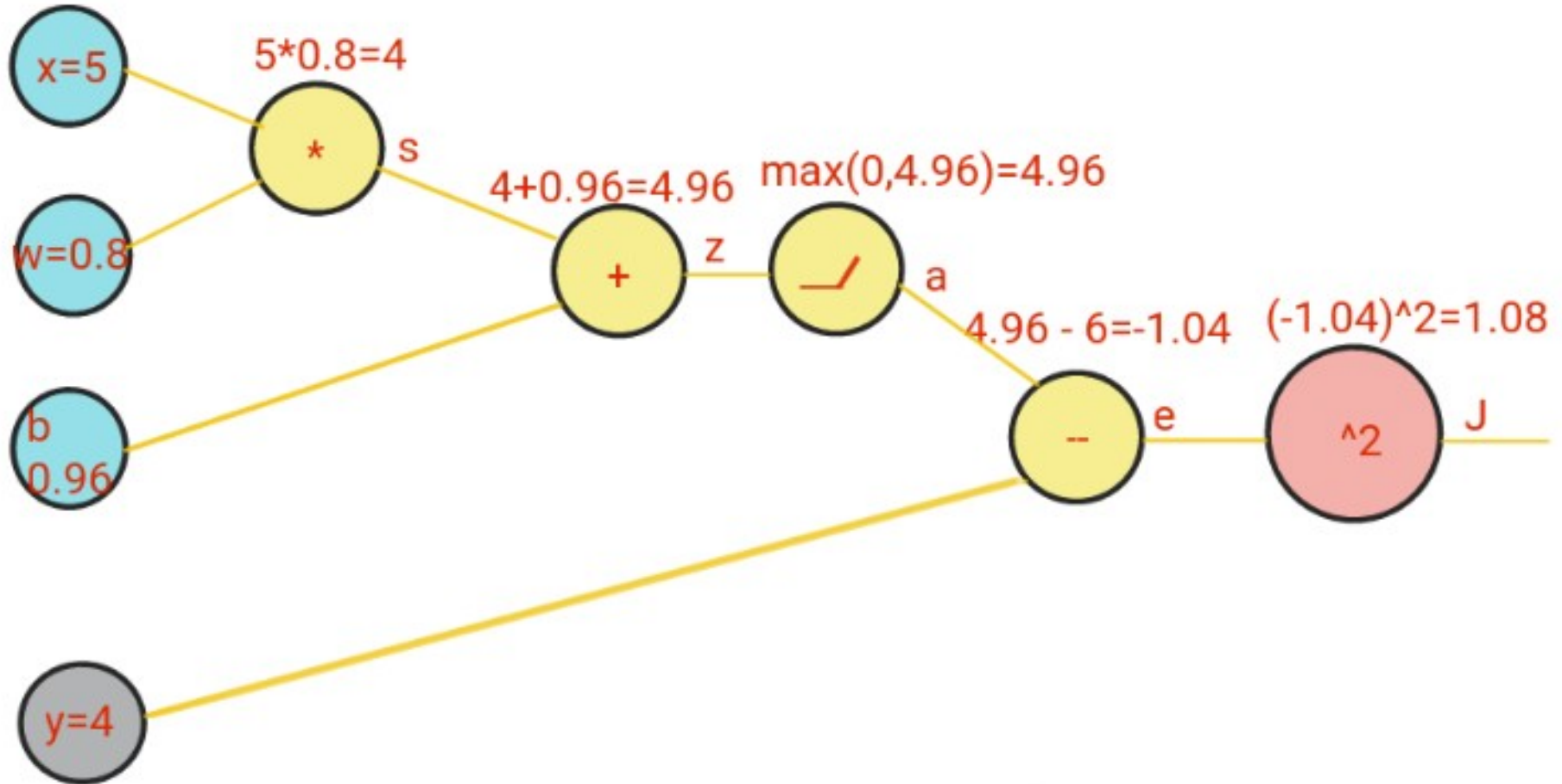
$$w_{new} = w - \alpha \frac{\partial J}{\partial w}$$

$$b_{new} = b - \alpha \frac{\partial J}{\partial b}$$

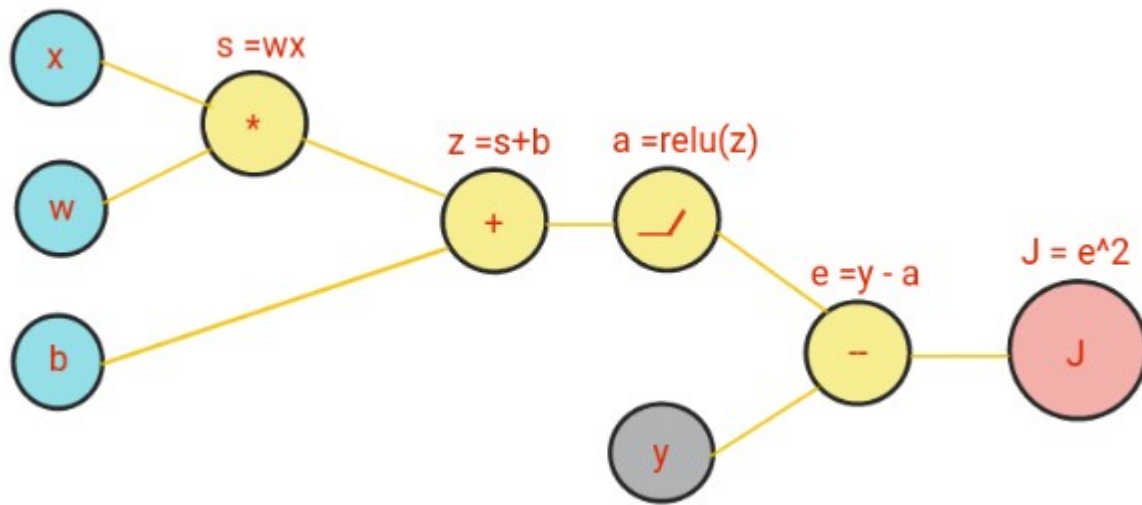
$$w_{new} = 1 - (0.01) 20 = 0.8$$

$$b_{new} = 1 - (0.01) 4 = 0.96$$

Gradient descent



Operations with chain rule

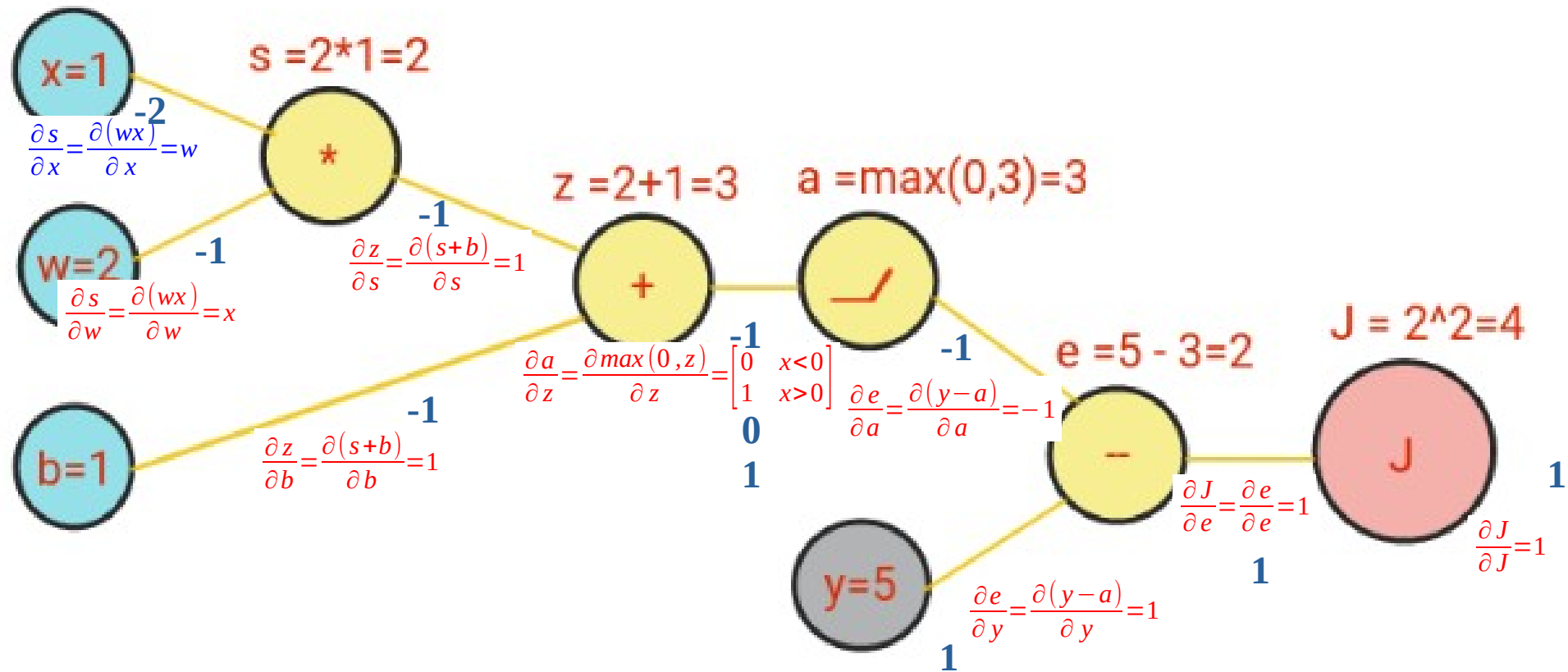


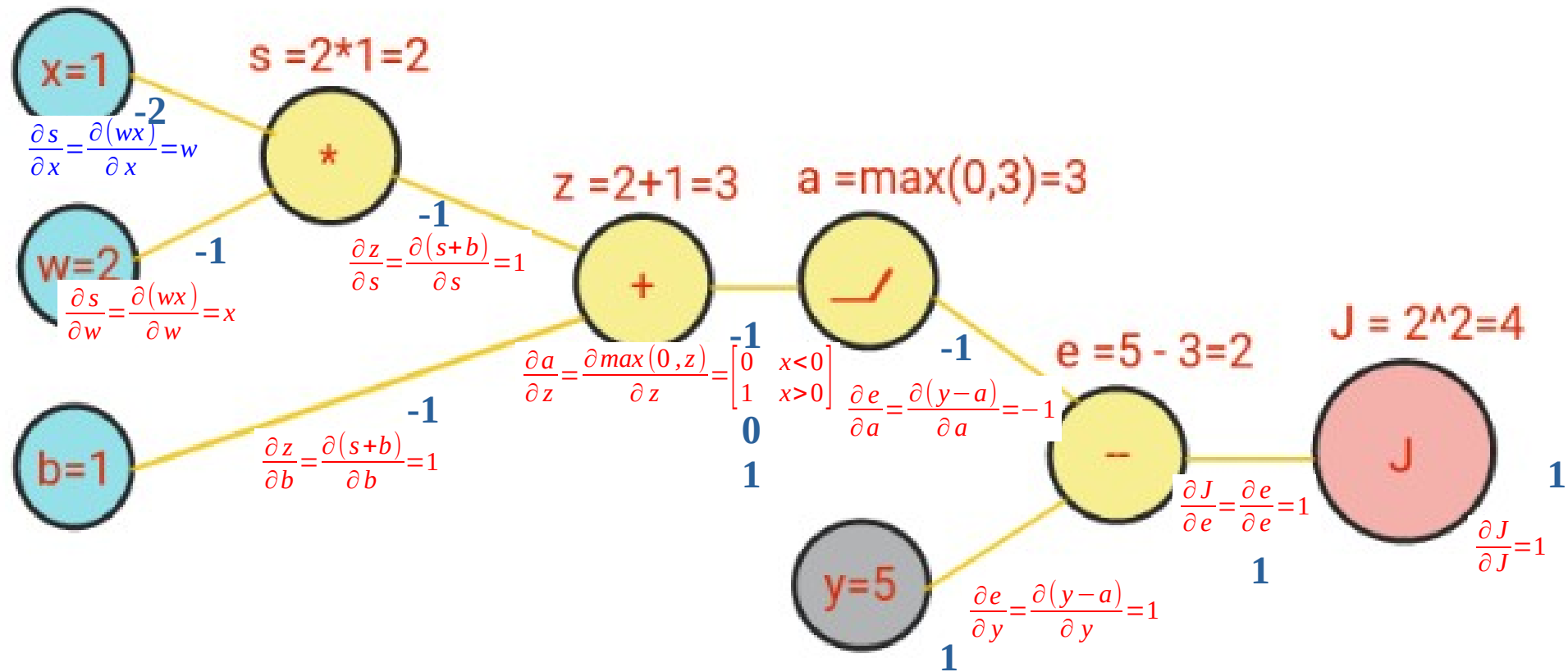
$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial e} \frac{\partial e}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial s} \frac{\partial s}{\partial w} = 2 * (-1) * \begin{bmatrix} 0 & x < 0 \\ 1 & x > 0 \end{bmatrix} * 1 * x$$

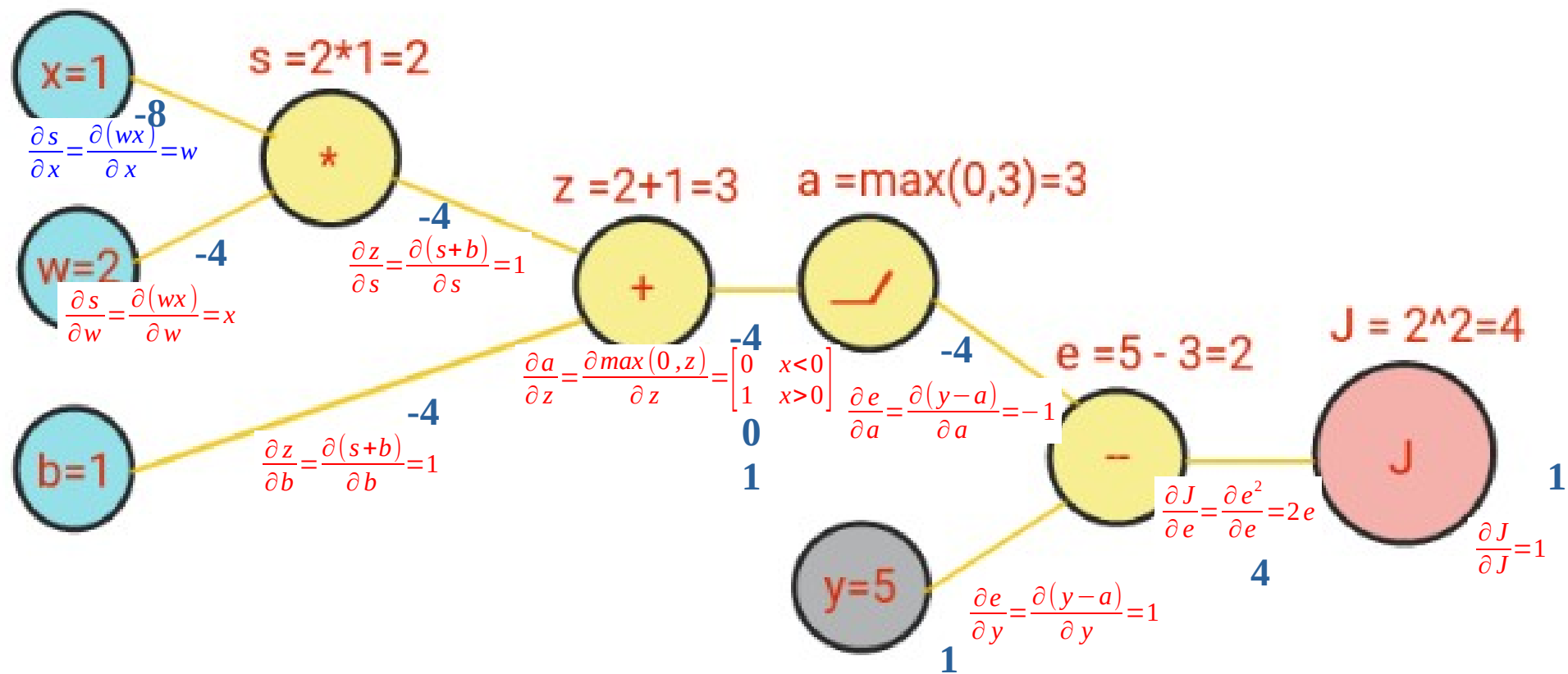
$\frac{\partial s}{\partial w} = \frac{\partial (wx)}{\partial w} = x$	$\frac{\partial s}{\partial x} = \frac{\partial (wx)}{\partial x} = w$
$\frac{\partial z}{\partial s} = \frac{\partial (s+b)}{\partial s} = 1$	$\frac{\delta w_6}{\delta w_3} = \frac{\delta w_5 + w_3}{\delta w_3} = 1$
$\frac{\partial a}{\partial z} = \frac{\partial \max(0, z)}{\partial z} = \begin{cases} 0 & x < 0 \\ 1 & x > 0 \end{cases}$	$\frac{\partial e}{\partial a} = \frac{\partial (y-a)}{\partial a} = -1$
$\frac{\delta w_8}{\delta w_4} = \frac{\delta w_4 - w_7}{\delta w_4} = 1$	$\frac{\partial J}{\partial e} = \frac{\partial e^2}{\partial e} = 2$

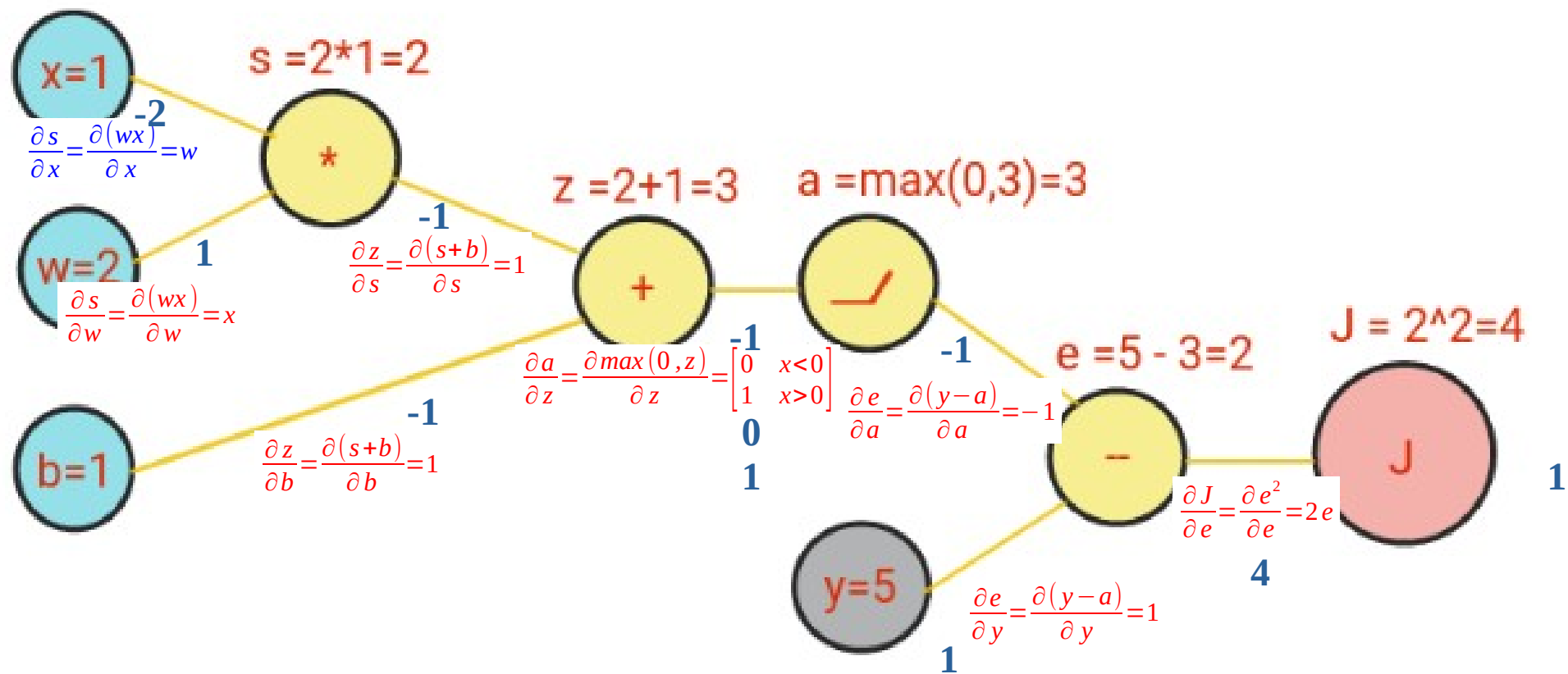
$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial e} \frac{\partial e}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial s} \frac{\partial s}{\partial w}$$

reference: <https://cs231n.github.io/optimization-2/>

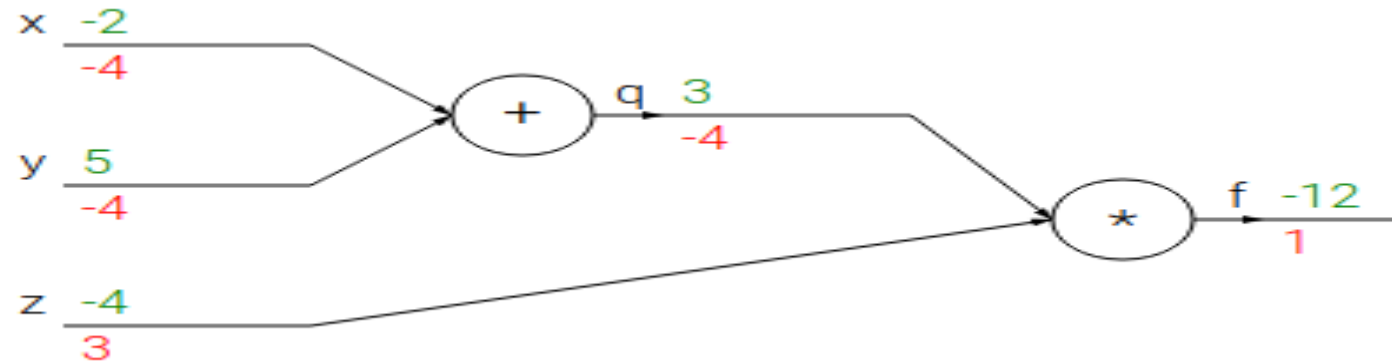






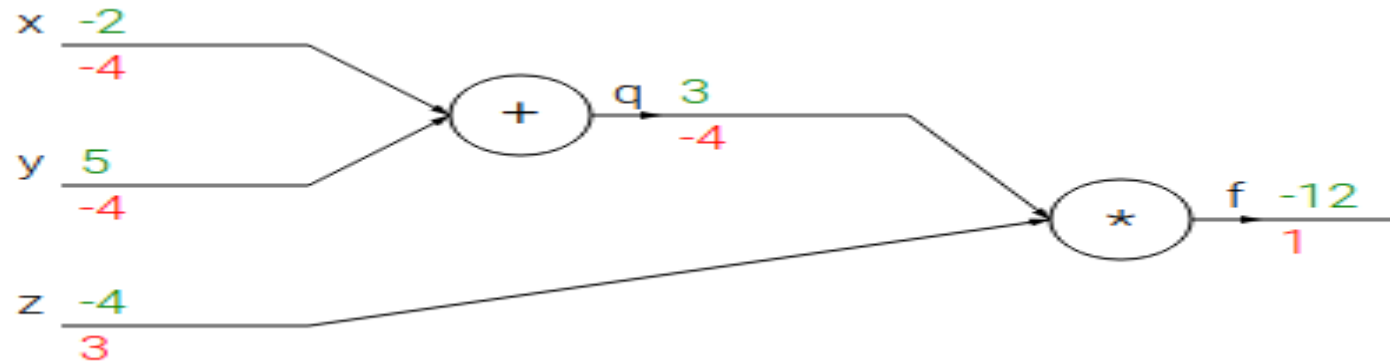


Backprop-Compound expressions with chain rule



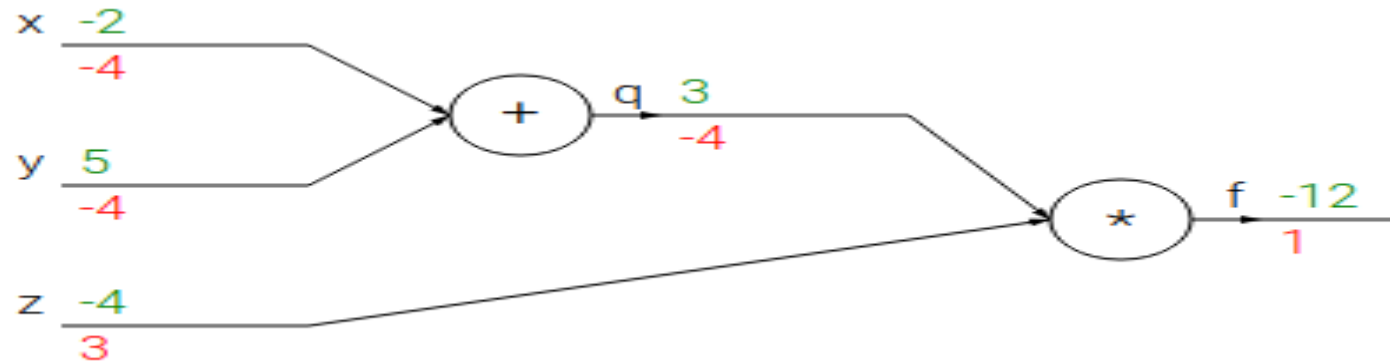
- The add operation received inputs $[-2, 5]$ and computed output 3 .
- Since the gate is computing the addition operation, its local gradient for both of its inputs is $+1$.
- The rest of the graph computed the final value, which is -12 .

Backprop-Compound expressions with chain rule



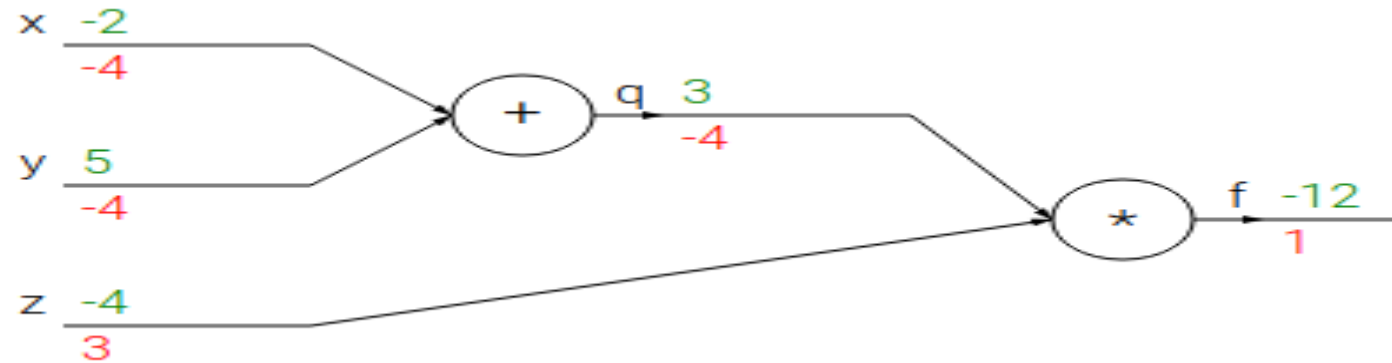
- During the backward pass in which the chain rule is applied recursively backwards through the computational graph
- The add operation (which is an input to the multiply operation) learns that the gradient for its output was -4.

Backprop-Compound expressions with chain rule



- Assume we want to maximize the value of graph then the output of the add operation to be lower (due to negative sign), and with a force of 4.
- To continue the recurrence and to chain the gradient, the add operation takes that gradient and multiplies it to all of the local gradients for its inputs (making the gradient on both x and y $1 * -4 = -4$).

Backprop-Compound expressions with chain rule



- Notice that this has the desired effect
- If x,y were to decrease (responding to their negative gradient) then the add gate's output would decrease, which in turn makes the multiply gate's output increase.

Gradient descent

- Because $\|\nabla J\|^2 \geq 0$ this guarantees that $\Delta J \leq 0$
- Means J will always decrease, never increase, if we change θ according to the $\Delta \theta = -\alpha \nabla J$

$$\Delta w = -\alpha \frac{\partial J}{\partial w}$$

$$\Delta b = -\alpha \frac{\partial J}{\partial b}$$

$$w_{new} = w + \Delta w = w - \alpha \frac{\partial J}{\partial w}$$

$$b_{new} = b + \Delta b = b - \alpha \frac{\partial J}{\partial b}$$

Gradient-Based Optimization

- Lets take simple example $f(x, y) = 2x^2 + y^2$
- The partial derivative $\frac{\partial f(x, y)}{\partial x}$ measures how f changes as only the variable x increases.
- The gradient of f is the vector containing all of the partial derivatives
- The vector of partial derivatives

$$\nabla f = \begin{bmatrix} \frac{\partial f(x, y)}{\partial x} \\ \frac{\partial f(x, y)}{\partial y} \end{bmatrix} = \begin{bmatrix} 4x \\ 2y \end{bmatrix}$$

The derivatives tell you:

They indicate the rate of change of a function with respect to that variable

surrounding an infinitesimally small region near a particular point:

$$\Delta z = 4x \cdot \Delta x + 2y \cdot \Delta y$$

Gradient-Based Optimization

- Here $f(x, y) = 2x^2 + y^2$ $f(x, y) = z$

$$\nabla f = \begin{bmatrix} 4x \\ 2y \end{bmatrix} \quad \Delta z = 4x \cdot \Delta x + 2y \cdot \Delta y$$

$$(0.5, 0.5) = 0.75$$

$$(0.6, 0.5) = 0.97 \quad \Delta z = 4(0.5) \cdot (0.1) = 0.2$$

$$(0.6, 0.6) = 1.08 \quad \Delta z = 2(0.5) \cdot (0.1) = 0.1$$

$$\Delta z = 0.2 + 0.1 = 0.3$$

$$z = z + \Delta z = 0.75 + 0.3 = 1.05$$

Gradient-Based Optimization

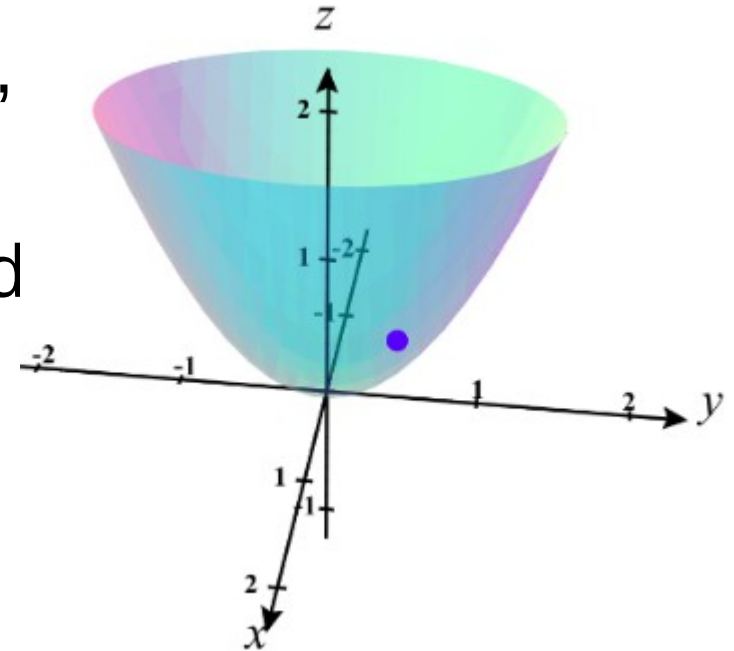
- Consider a simple addition function of two numbers

$$f(x, y) = x + y \qquad \frac{\partial f(x, y)}{\partial x} = \frac{\partial f}{\partial x} = 1 \qquad \frac{\partial f(x, y)}{\partial y} = \frac{\partial f}{\partial y} = 1$$

- That is, the derivative on both x, y is 1 regardless of what the values of x, y are.
- This makes sense, since increasing either x, y would increase the output of f , and the rate of that increase would be independent of what the actual values of x, y are (unlike the case of multiplication above).

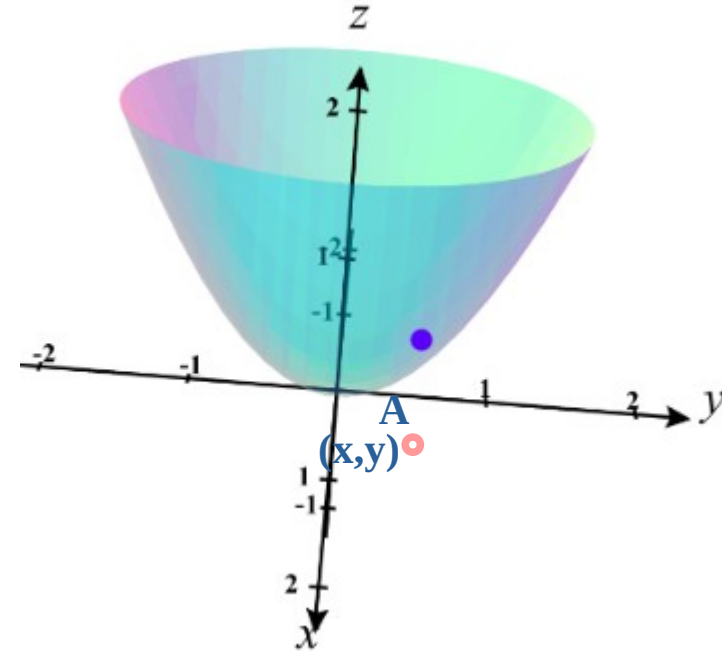
Gradient-Multivariate function

- Let $z = f(x, y)$ is multivariate function
- function $z = f(x, y)$ of two variables, whose graph will be a surface.
- Imagine we change x to $x + \Delta x$ and y to $y + \Delta y$ with Δx and Δy very small.
- What is the corresponding change in z ?



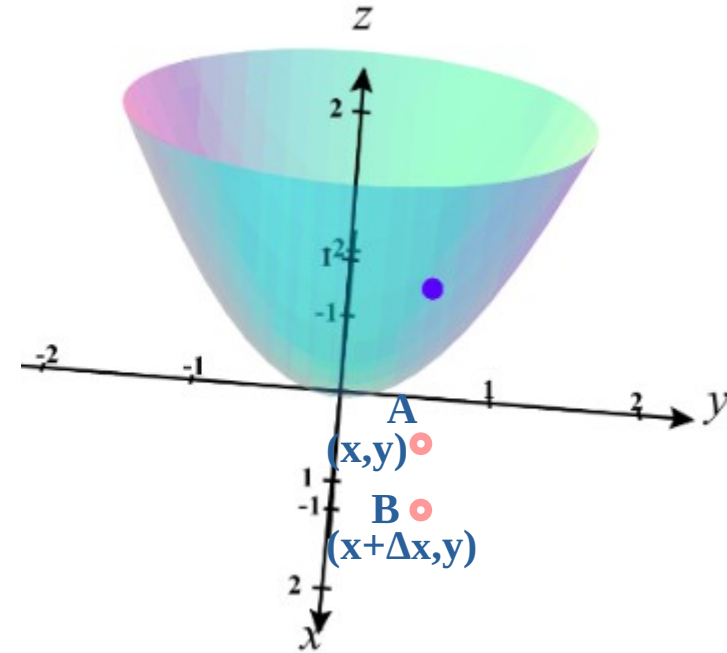
Gradient-Multivariate function

- In the (x, y) plane let A be the point with coordinates (x, y)



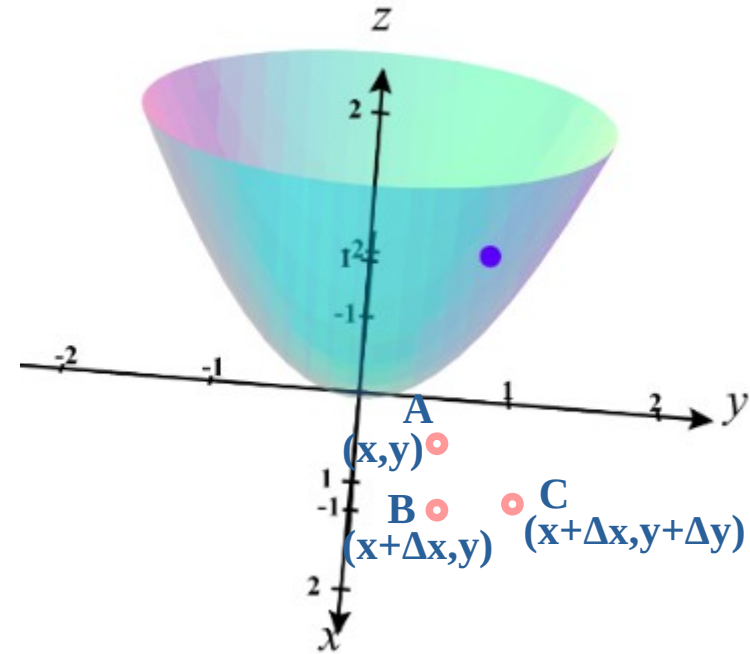
Gradient-Multivariate function

- In the (x, y) plane let A be the point with coordinates (x, y)
- let B be the point with coordinates $(x + \Delta x, y)$



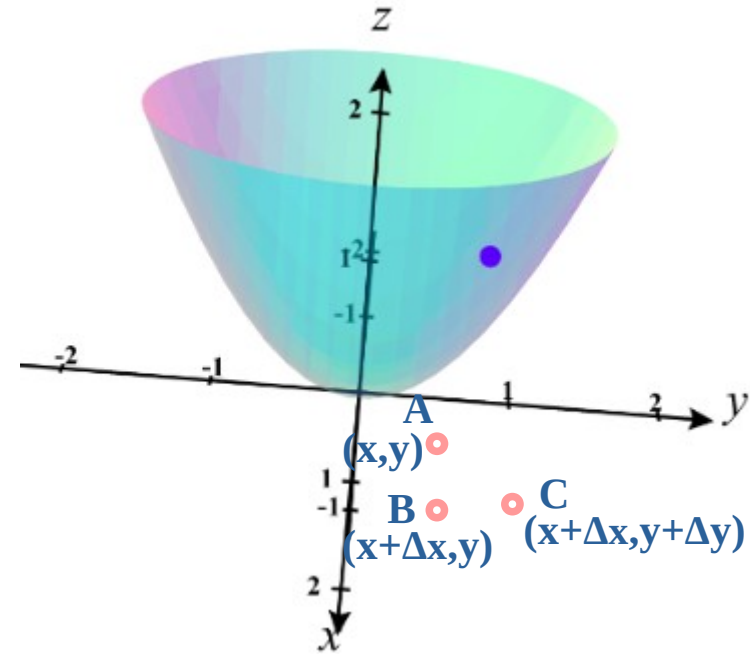
Gradient-Multivariate function

- In the (x, y) plane let A be the point with coordinates (x, y)
- let B be the point with coordinates $(x + \Delta x, y)$
- and C the point with coordinates $(x + \Delta x, y + \Delta y)$.



Gradient-Multivariate function

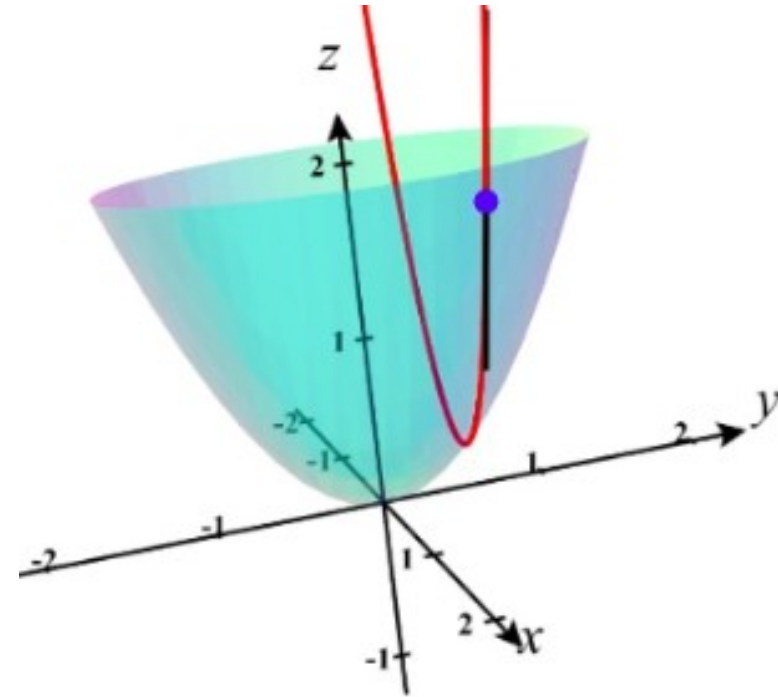
- The overall change in height, Δz , from A to C is given by
- $\Delta z = (\text{change in height A to B}) + (\text{change in height B to C})$



Gradient-Multivariate function

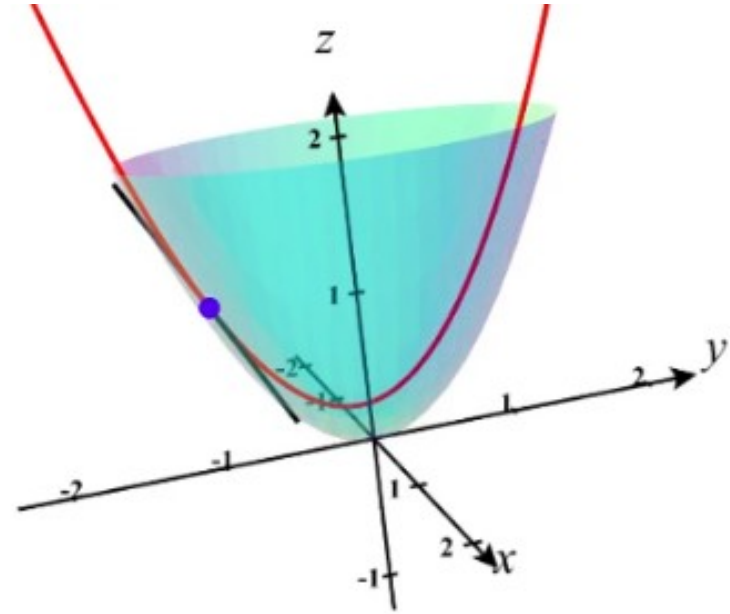
- In calculating the change in height from A to B we are travelling across the surface from A to B along a curve in which y is held fixed
- change in height A to B $\approx \frac{\partial z}{\partial x} \Delta x$

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



Gradient-Multivariate function

- In calculating the change in height from B to C we are travelling across the surface from B to C along a curve in which x is held fixed
- change in height B to C $\approx \frac{\partial z}{\partial y} \Delta y$



Gradient-Multivariate function

- The overall change in height, Δz , from A to C is given by
- $\Delta z = (\text{change in height A to B}) + (\text{change in height B to C})$
- Therefore

$$\Delta z = \frac{\partial z}{\partial x} \Delta x + \frac{\partial z}{\partial y} \Delta y$$

