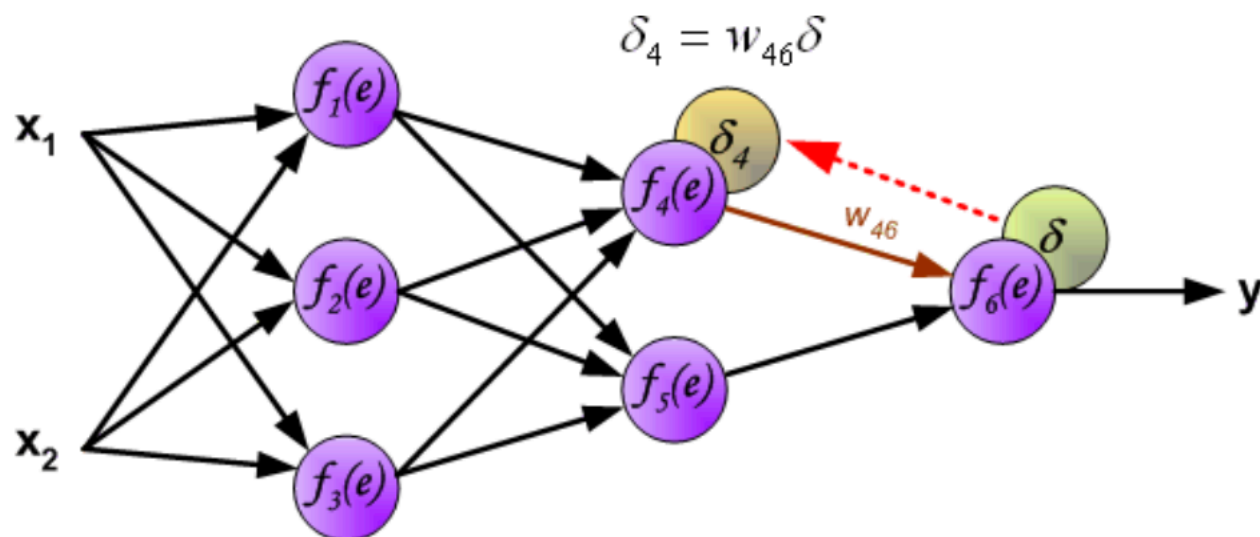


Lesson4 Basic Neural Network



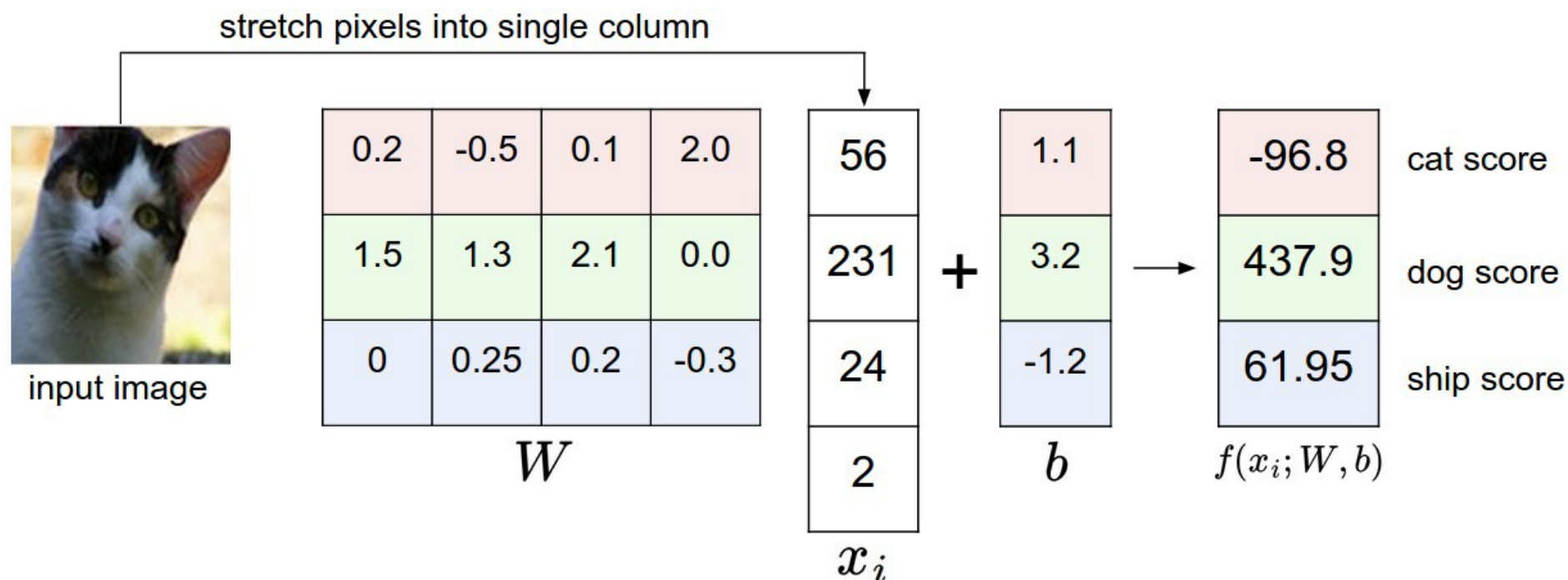
By Daniel
Research Scientist
danielshaving@gmail.com

- Neural Network
- Loss Function
- Gradient descent
- Forward Propagation
- Backward Propagation

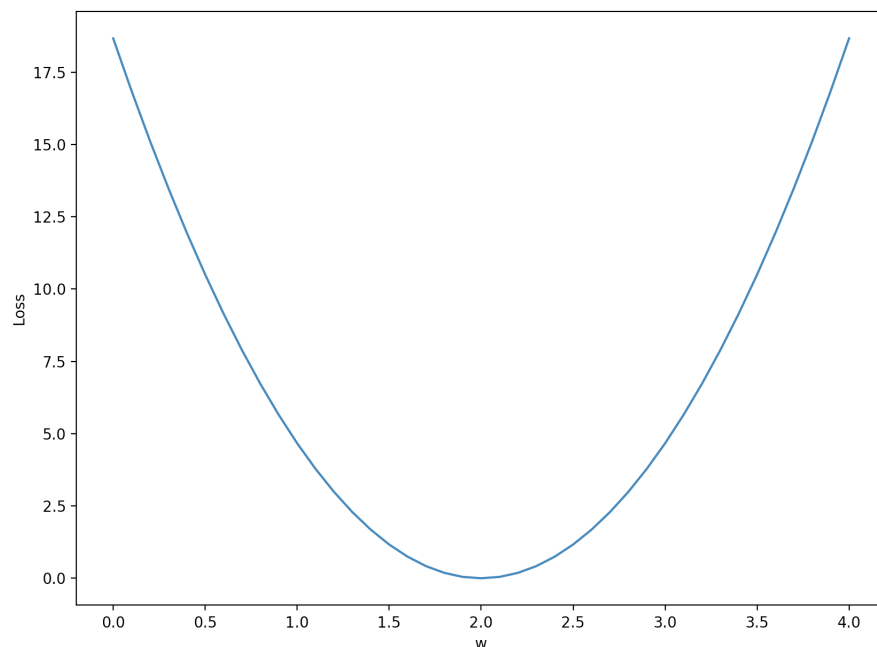
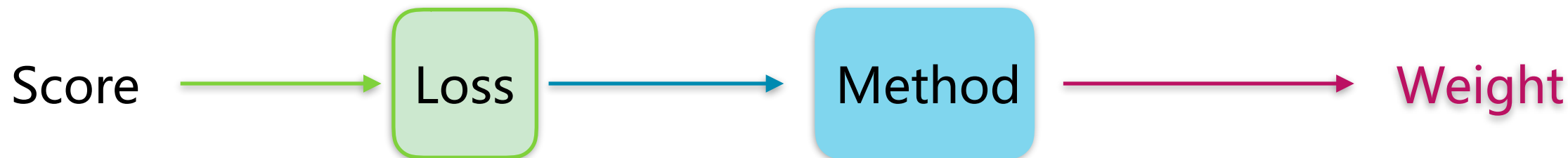
An NN model

Suppose we have a weights matrix $W \in \mathbb{R}^{K \times D}$, a bias vector $b \in \mathbb{R}^{K \times 1}$ where K is number of classes and D is feature dimensionality, then given an input $x_i \in \mathbb{R}^{D \times 1}$, the score function $f : \mathbb{R}^D \rightarrow \mathbb{R}^K$ is:

$$f(x_i, W, b) = Wx_i + b$$



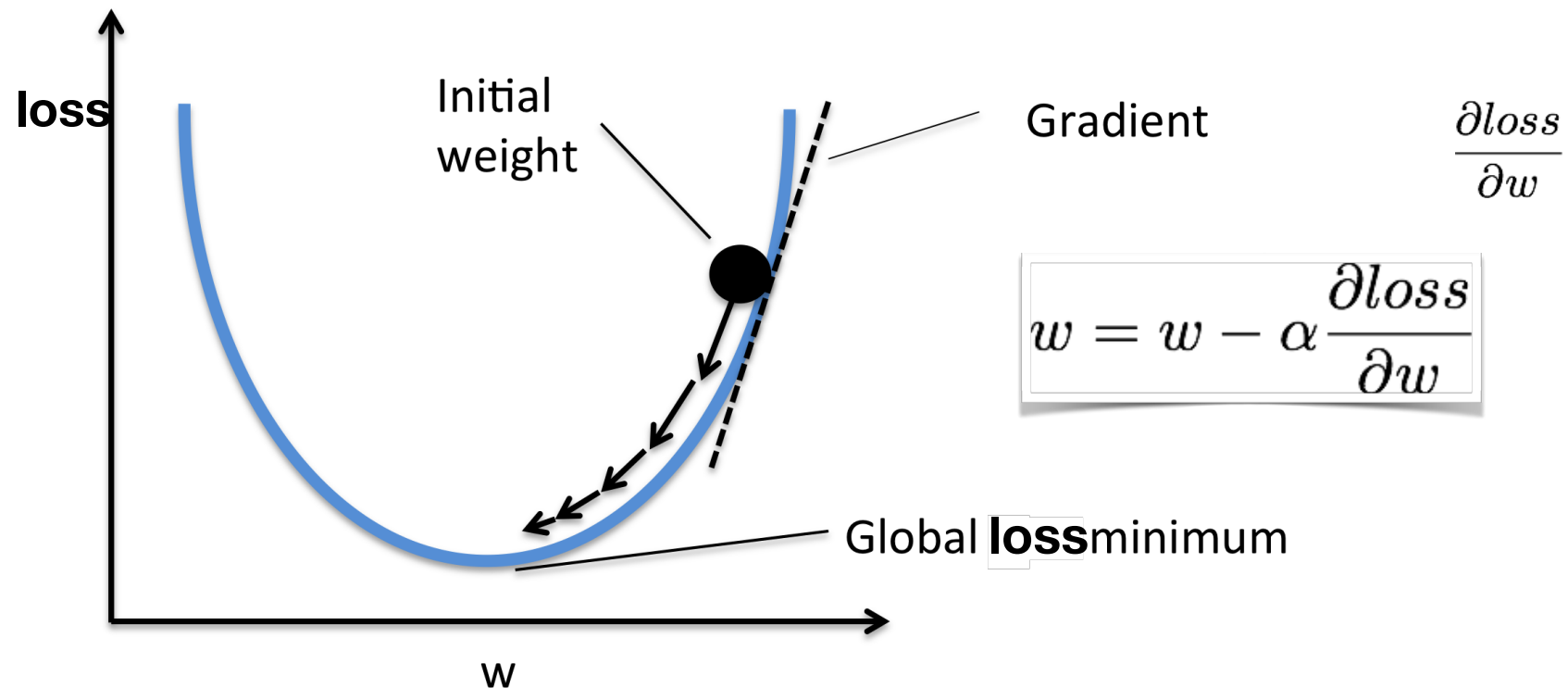
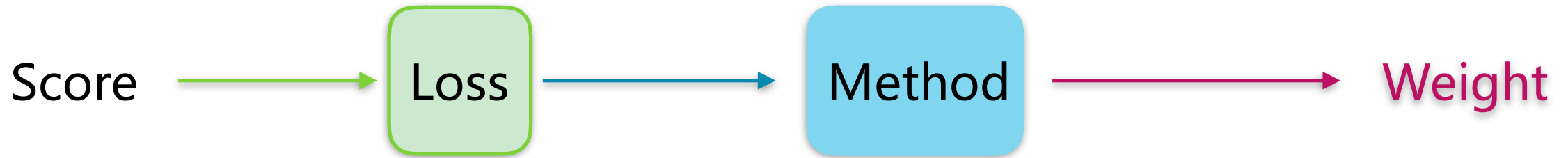
Using score to correct weights (supervising)



$$loss(w) = \frac{1}{N} \sum_{n=1}^N (\hat{y}_n - y_n)^2$$

$$\arg \min_w loss(w)$$

Example: UNet $(x \times w - y)^2$



Derivative $(x \times w - y)^2$

$$loss = (\hat{y} - y)^2 = (x * w - y)^2$$

$$w = w - \alpha \frac{\partial loss}{\partial w}$$

$$\frac{\partial loss}{\partial w} = ?$$

$$w = w - \alpha \frac{\partial loss}{\partial w}$$

$$= w - \alpha * 2x(xw - y)$$

YOUR INPUT:
 $f(w) =$

$$(xw - y)^2$$

Simplify Roots/zeros

FIRST DERIVATIVE:

$$\frac{d}{dw} [f(w)] = f'(w) =$$

The steps of calculation are displayed.

Move the mouse over a derivative $\frac{d}{dw} [\dots]$ or tap it in order to show its calculation.

$$\frac{d}{dw} [(xw - y)^2]$$

$$= 2(xw - y) \cdot \frac{d}{dw} [xw - y]$$

$$= 2 \left(x \cdot \frac{d}{dw} [w] + \frac{d}{dw} [-y] \right) (xw - y)$$

$$= 2(1x + 0)(xw - y)$$

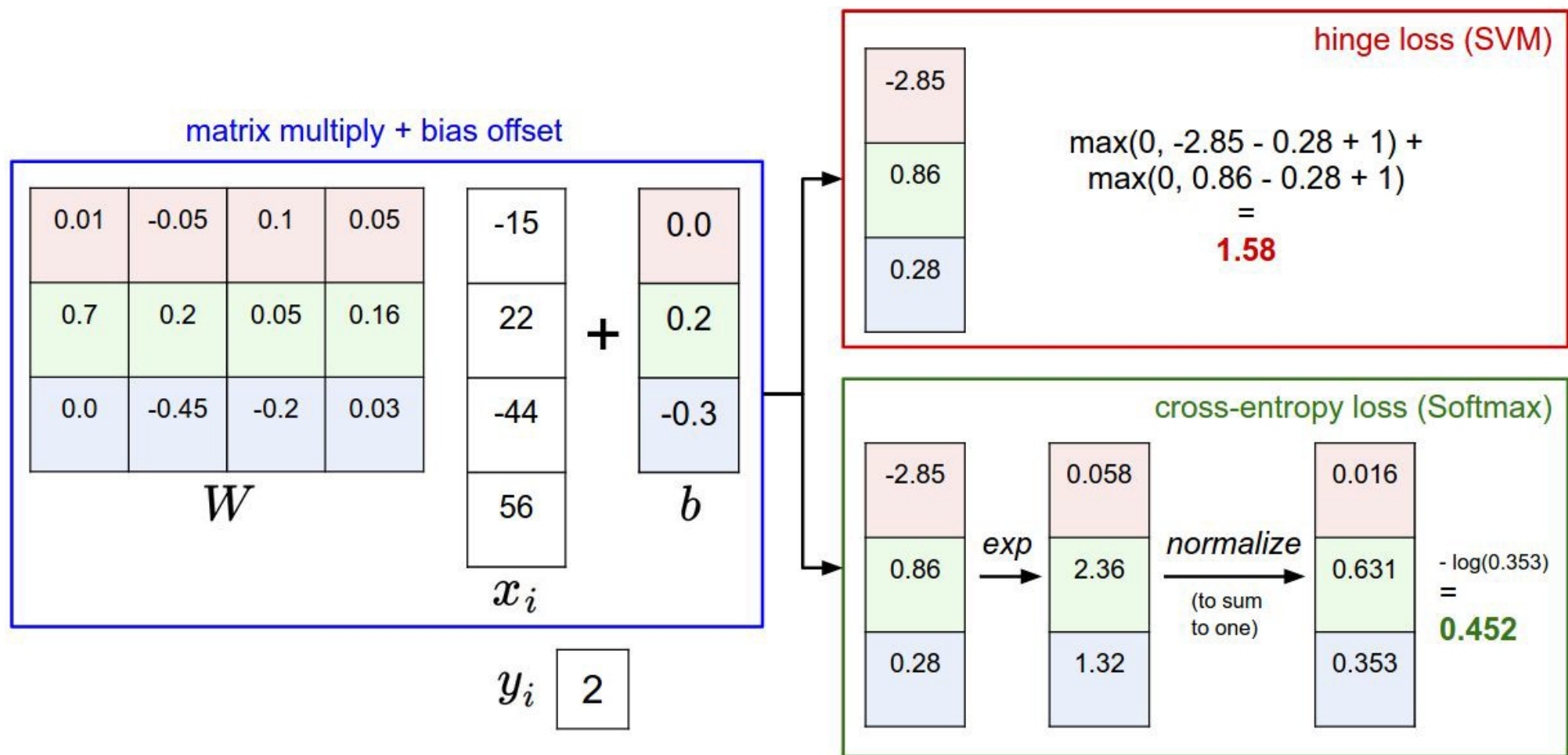
$$= 2x(xw - y)$$

Table of contents



- Neutral Network
- Loss Function
 - Loss function overview
 - Softmax + Cross Entropy loss
 - Hinge Loss
 - Overall Loss
 - Loss summary
- Gradient descent
- Forward Propagation
- Backward Propagation

Loss Function to calculate the errors



Loss Function (Cross-Entropy and Hinge (SVM Loss))



Cross-Entropy-Loss (after Softmax)

Let (x_i, y_i) be a pair of training example, where $x_i \in \mathbb{R}^D$ is training data and $y_i \in \mathbb{R}^K$ is a one-hot vector with all zeros except for its true class index is one. Assume the prediction $\hat{y}_i \in \mathbb{R}^K$ is computed after Softmax, then the cross-entropy loss is:

$$L_i = -y_i \cdot \log(\hat{y})$$

Hinge Loss (SVM Loss)

Note that the exponential computing in cross-entropy loss is expensive and sometimes we do not need a probability. Thus, the hinge loss function is:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + \Delta)$$

- Neutral Network
- Loss Function
 - Loss function overview
 - Softmax + Cross Entropy loss
 - Hinge Loss
 - Overall Loss
 - Loss summary
- Gradient descent
- Forward Propagation
- Backward Propagation

Softmax Classifier (Multinomial Logistic Regression)



cat

3.2

car

5.1

frog

-1.7

scores = unnormalized log probabilities of the classes.

$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{where} \quad s = f(x_i; W)$$

Want to maximize the log likelihood, or (for a loss function) to minimize the negative log likelihood of the correct class:

$$L_i = -\log P(Y = y_i|X = x_i)$$

in summary:
$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$



cat

3.2

car

5.1

frog

-1.7

$$L_i = -\log\left(\frac{e^{sy_i}}{\sum_j e^{s_j}}\right)$$

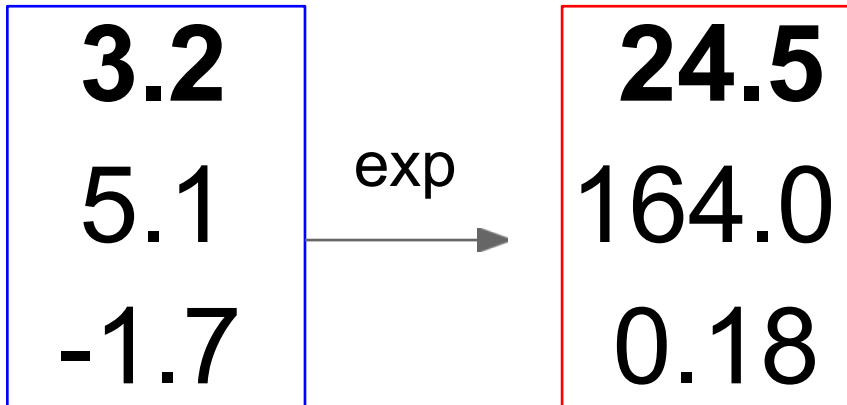
unnormalized log probabilities



$$L_i = -\log\left(\frac{e^{sy_i}}{\sum_j e^{s_j}}\right)$$

unnormalized probabilities

cat
car
frog



unnormalized log probabilities



$$L_i = -\log\left(\frac{e^{sy_i}}{\sum_j e^{s_j}}\right)$$

unnormalized probabilities

cat
car
frog

3.2

5.1

-1.7

exp

24.5

164.0

0.18

normalize

0.13

0.87

0.00

unnormalized log probabilities

probabilities

Softmax + Cross entropy



$$L_i = -y_i \cdot \log(\hat{y})$$

unnormalized probabilities

cat
car
frog

3.2

5.1

-1.7

exp

24.5

164.0

0.18

normalize

0.13

0.87

0.00

$$L_i = -\log(0.13) \\ = 0.89$$

unnormalized log probabilities

probabilities

Table of contents



- Neutral Network
- Loss Function
 - Loss function overview
 - Softmax + Cross Entropy loss
 - Hinge Loss
 - Overall Loss
 - Loss summary
- Gradient descent
- Forward Propagation
- Backward Propagation

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

Multiclass SVM loss:

Given an example (x_i, y_i)
where x_i is the image and
where y_i is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Hinge Loss

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9		

Multiclass SVM loss:

Given an example (x_i, y_i)
where x_i is the image and
where y_i is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$\begin{aligned} L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\ &= \max(0, 5.1 - 3.2 + 1) \\ &\quad + \max(0, -1.7 - 3.2 + 1) \\ &= \max(0, 2.9) + \max(0, -3.9) \\ &= 2.9 + 0 \\ &= 2.9 \end{aligned}$$

Hinge Loss

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	10.9

Multiclass SVM loss:

Given an example (x_i, y_i)
where x_i is the image and
where y_i is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$

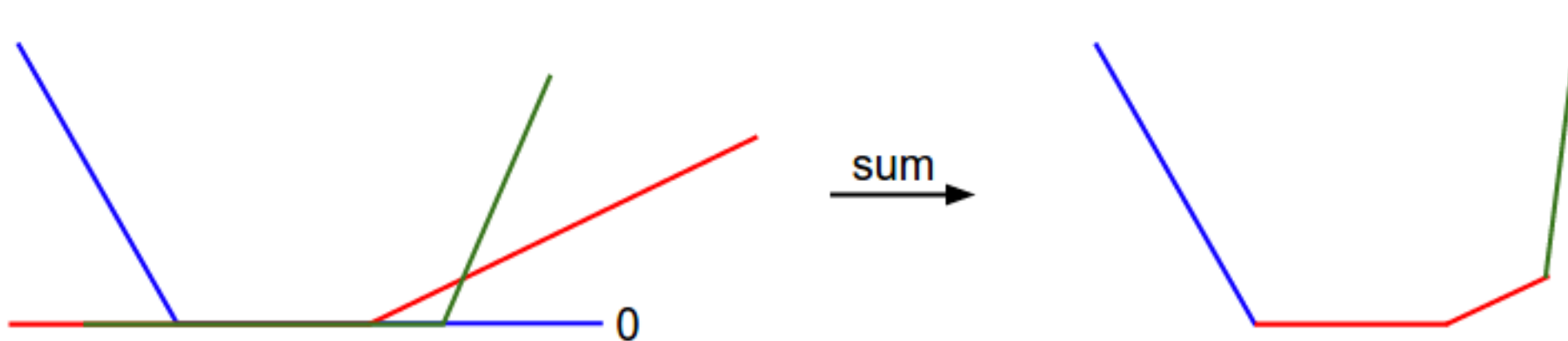
the SVM loss has the form:

$$\begin{aligned} L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\ &= \max(0, 2.2 - (-3.1) + 1) \\ &\quad + \max(0, 2.5 - (-3.1) + 1) \\ &= \max(0, 5.3) + \max(0, 5.6) \\ &= 5.3 + 5.6 \\ &= 10.9 \end{aligned}$$

- Neutral Network
- Loss Function
 - Loss function overview
 - Softmax + Cross Entropy loss
 - Hinge Loss
 - Overall Loss
 - Loss summary
- Gradient descent
- Forward Propagation
- Backward Propagation

The overall loss consists of two items, namely data loss and regularization loss.

$$L = \frac{1}{N} \sum_i L_i + \lambda R(W)$$



- Neutral Network
- Loss Function
 - Loss function overview
 - Softmax + Cross Entropy loss
 - Hinge Loss
 - Overall Loss
 - Loss summary
- Gradient descent
- Forward Propagation
- Backward Propagation

Summary of loss function

- We have some dataset of (x,y)
- We have a **score function**: $s = f(x; W) \stackrel{\text{e.g.}}{=} Wx$
- We have a **loss function**:

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right) \quad \text{Softmax}$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \quad \text{SVM}$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i + R(W) \quad \text{Full loss}$$

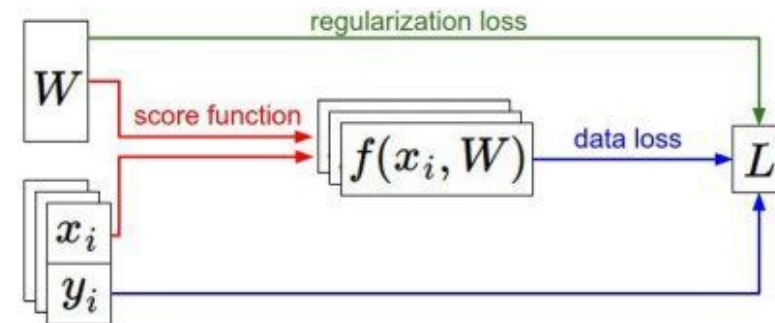


Table of contents



- Neutral Network
- Loss Function
- Gradient descent
- Forward Propagation
- Backward Propagation

In 1-dimension, the derivative of a function:

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

In multiple dimensions, the **gradient** is the vector of (partial derivatives).

Follow the slope

current W:

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25347

W + h (first dim):

[0.34 + **0.0001**,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25322

gradient dW:

[?,
?,
?,
?,
?,
?,
?,
?,
?,...]

Follow the slope

current W:

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25347

W + h (first dim):

[0.34 + **0.0001**,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]

loss 1.25322

gradient dW:

[-2.5,
?,
?,

$$(1.25322 - 1.25347)/0.0001 = -2.5$$

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

?,
?,...]

```
# Vanilla Gradient Descent
```

```
while True:
```

```
    weights_grad = evaluate_gradient(loss_fun, data, weights)
```

```
    weights += - step_size * weights_grad # perform parameter update
```

```
# Vanilla Minibatch Gradient Descent
```

```
while True:
```

```
    data_batch = sample_training_data(data, 256) # sample 256 examples
```

```
    weights_grad = evaluate_gradient(loss_fun, data_batch, weights)
```

```
    weights += - step_size * weights_grad # perform parameter update
```

Gradient Decent

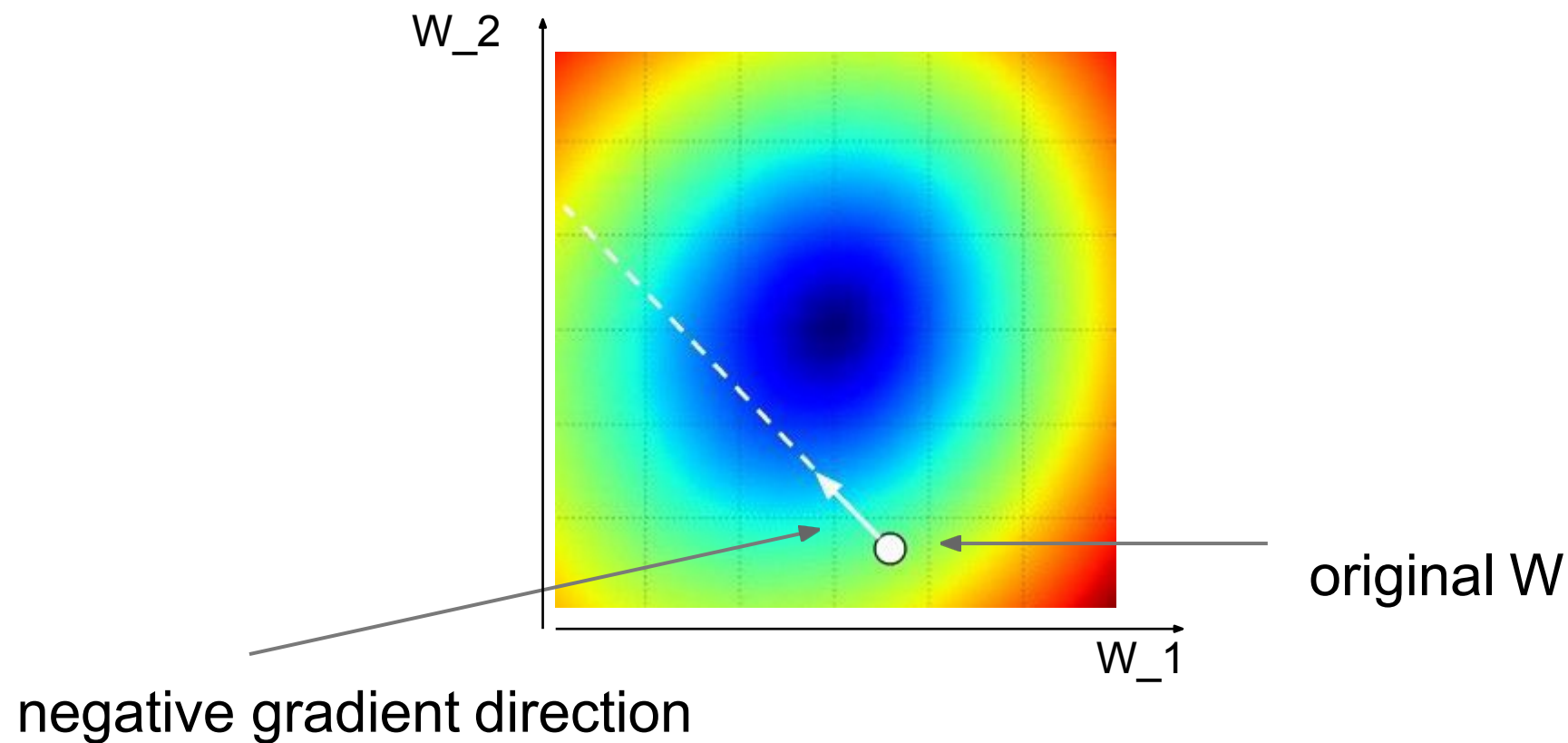
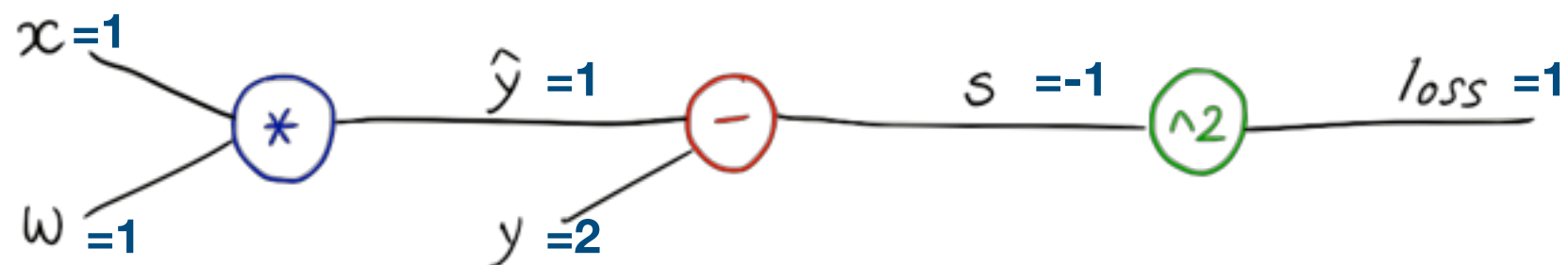


Table of contents



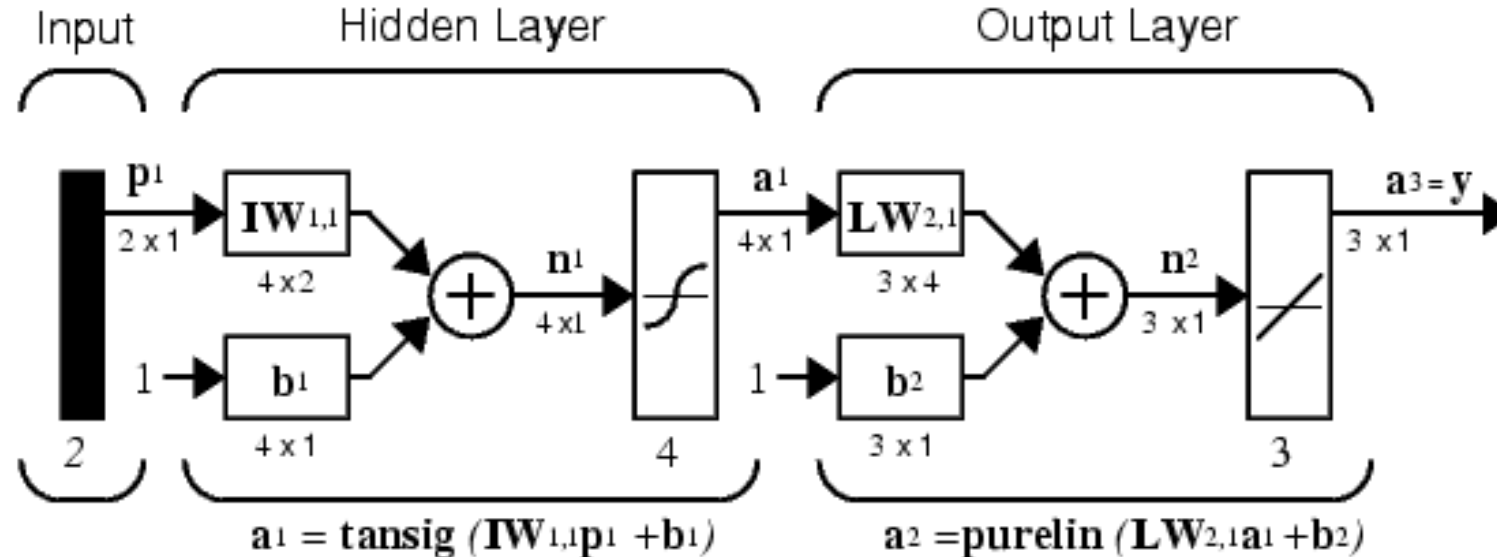
- Neutral Network
- Loss Function
- Gradient descent
- Forward Propagation
- Backward Propagation

Feed forwarding



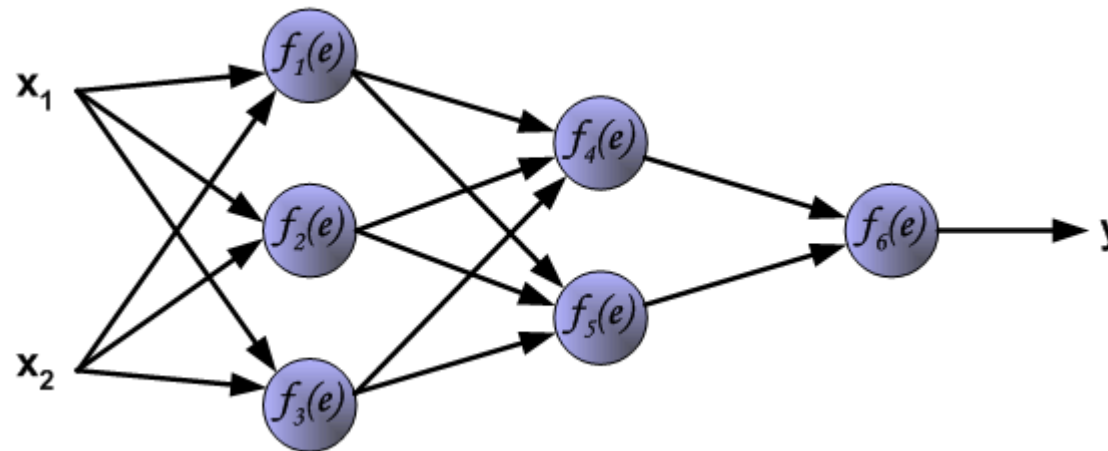
Feed forwarding

Feedforward networks often have one or more hidden layers of sigmoid neurons followed by an output layer of linear neurons. Multiple layers of neurons with nonlinear transfer functions allow the network to learn nonlinear and linear relationships between input and output vectors. The linear output layer lets the network produce values outside the range -1 to +1. On the other hand, if you want to constrain the outputs of a network (such as between 0 and 1), then the output layer should use a sigmoid transfer function (such as logsig).

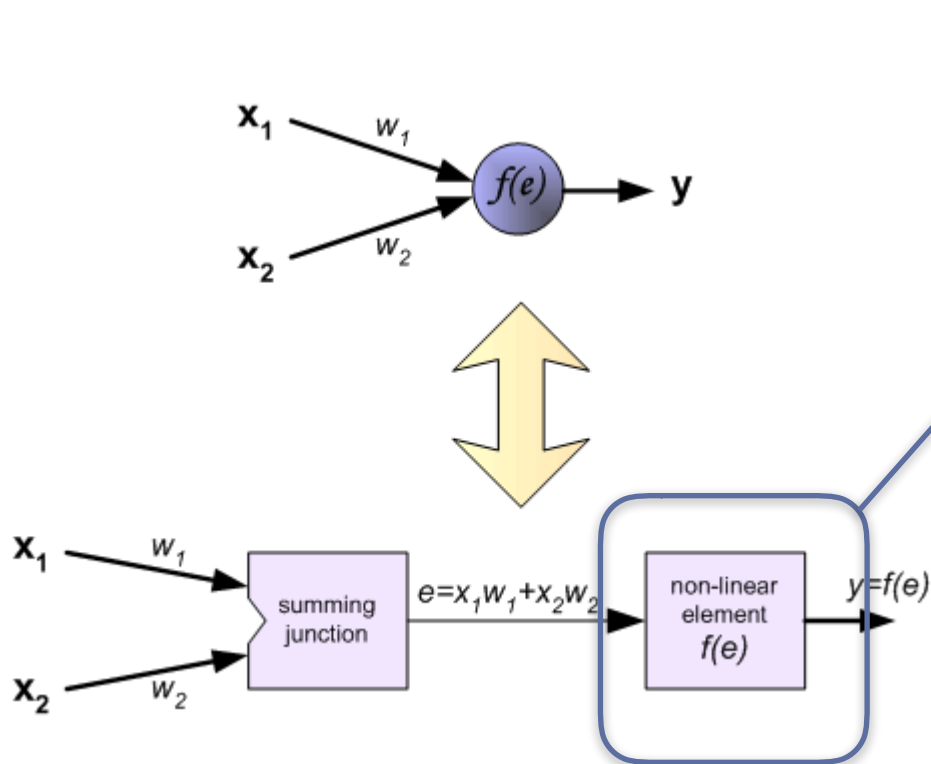


Feed forwarding

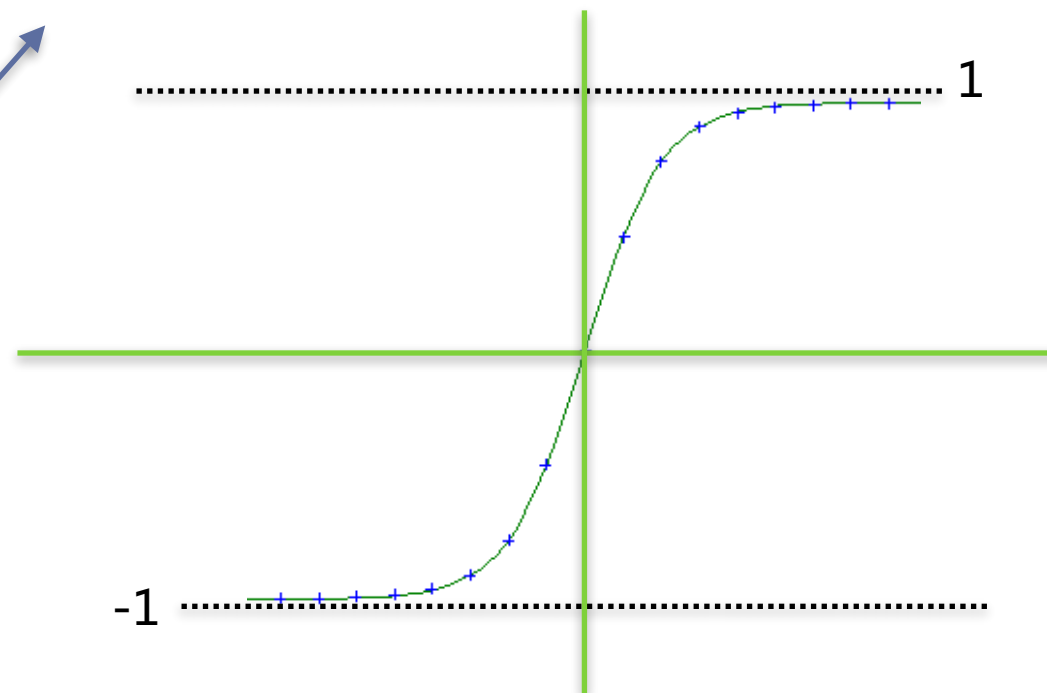
The following slides describes **teaching process** of multi-layer neural network employing **backpropagation** algorithm. To illustrate this process the three layer neural network with two inputs and one output, which is shown in the picture below, is used:



Active Function



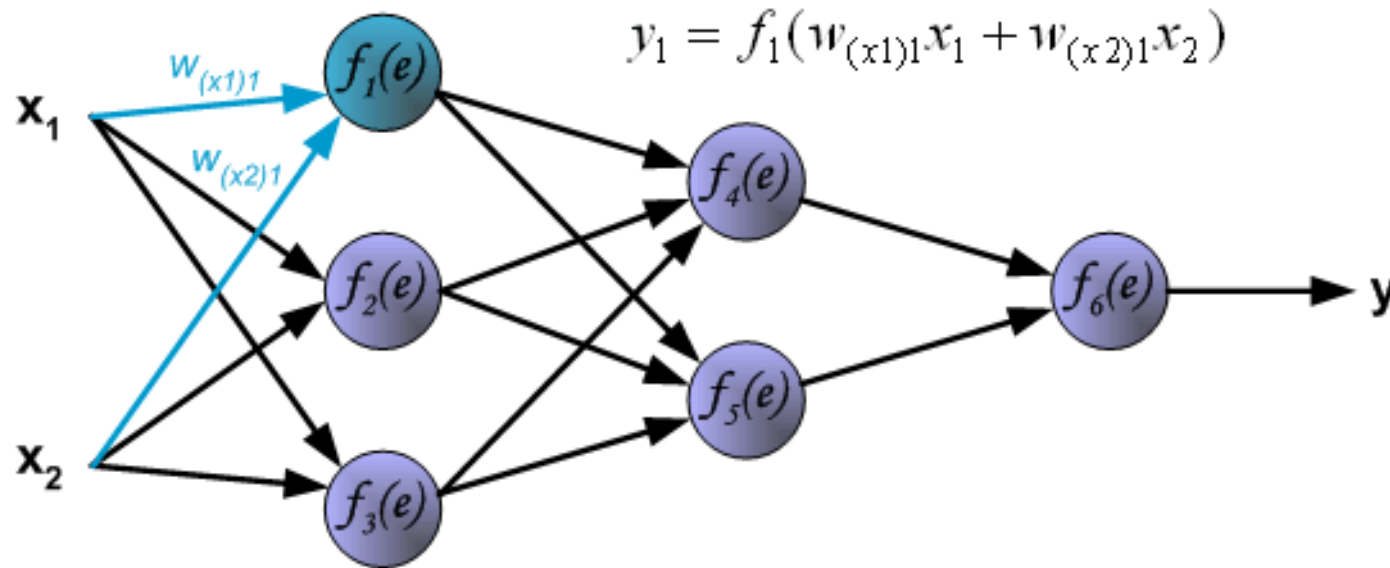
Active Function: Sigmoid



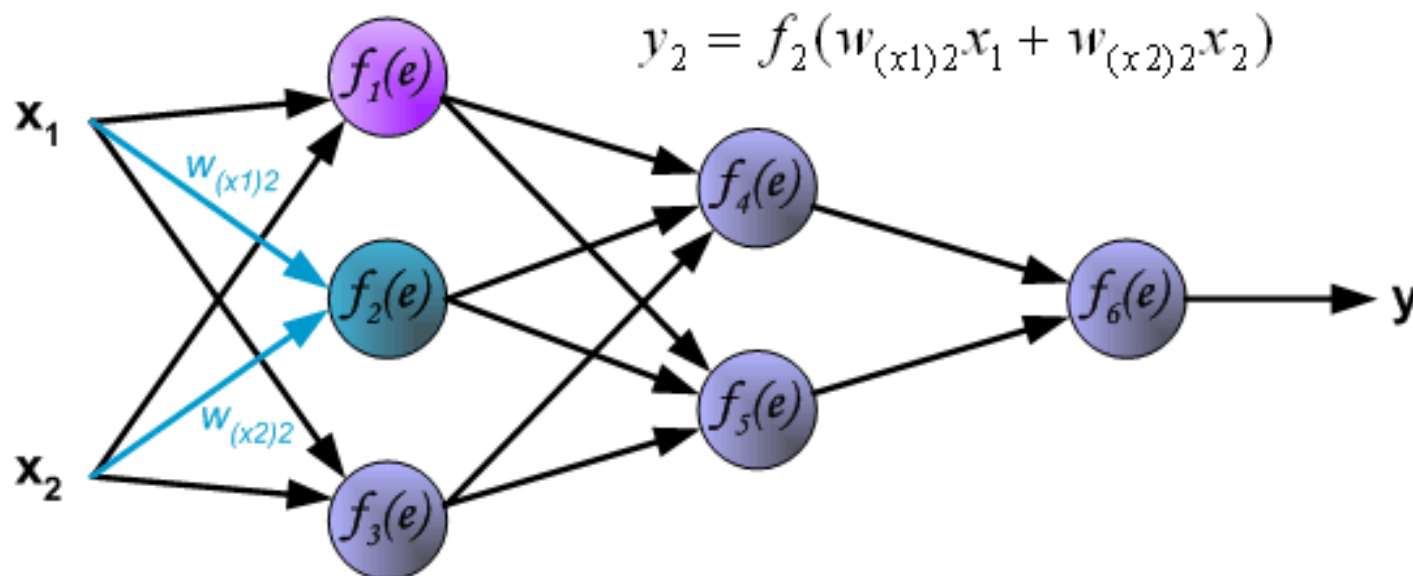
$$f(x) = \frac{1}{1 + e^{-\lambda x}}$$

Feed forwarding

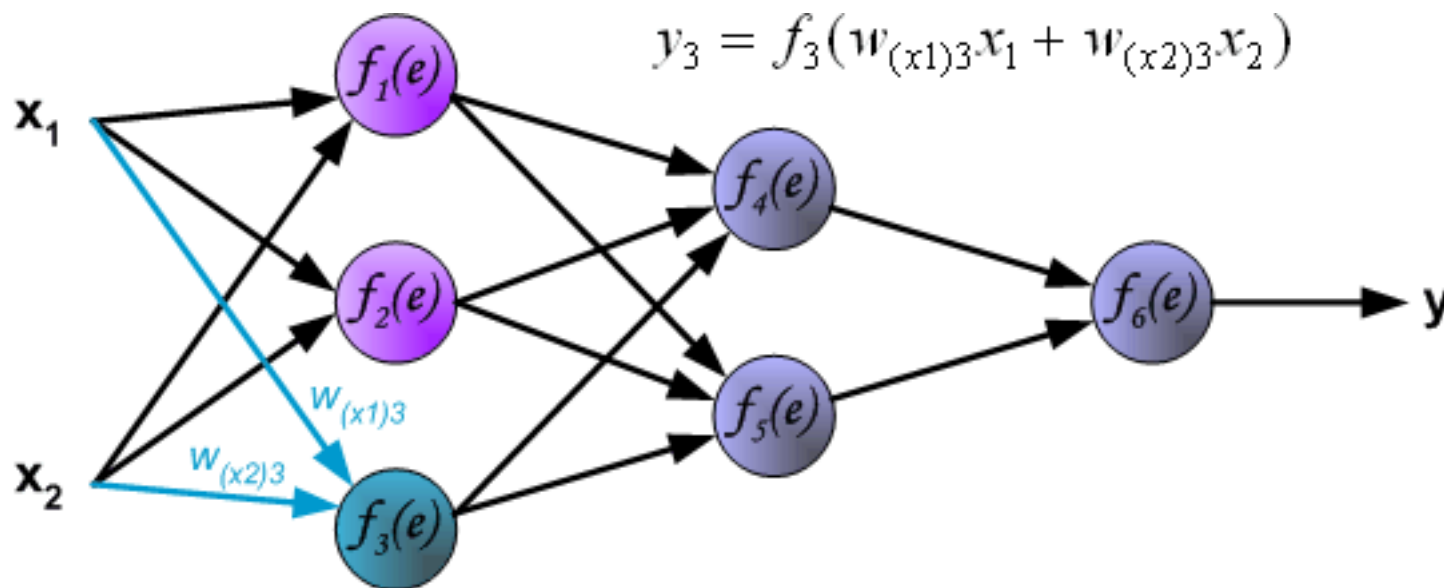
Pictures below illustrate how signal is propagating through the network, Symbols $w_{(xm)n}$ represent weights of connections between network input x_m and neuron n in input layer. Symbols y_n represents output signal of neuron n .



Feed forwarding

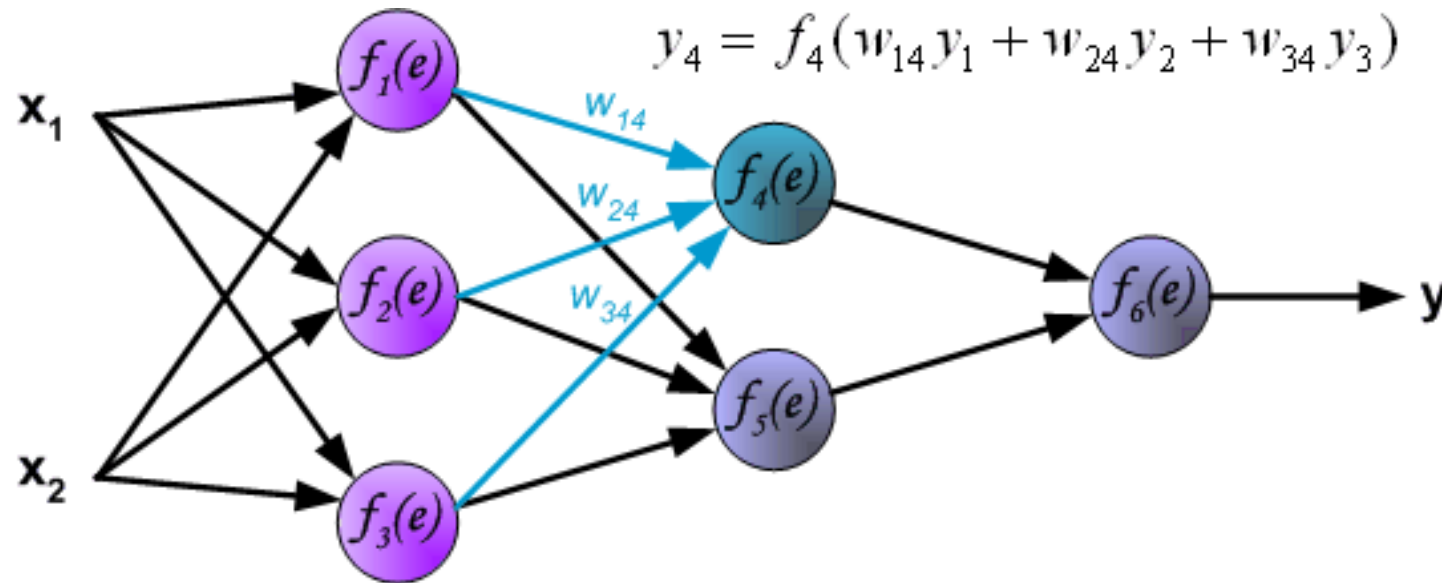


Feed forwarding

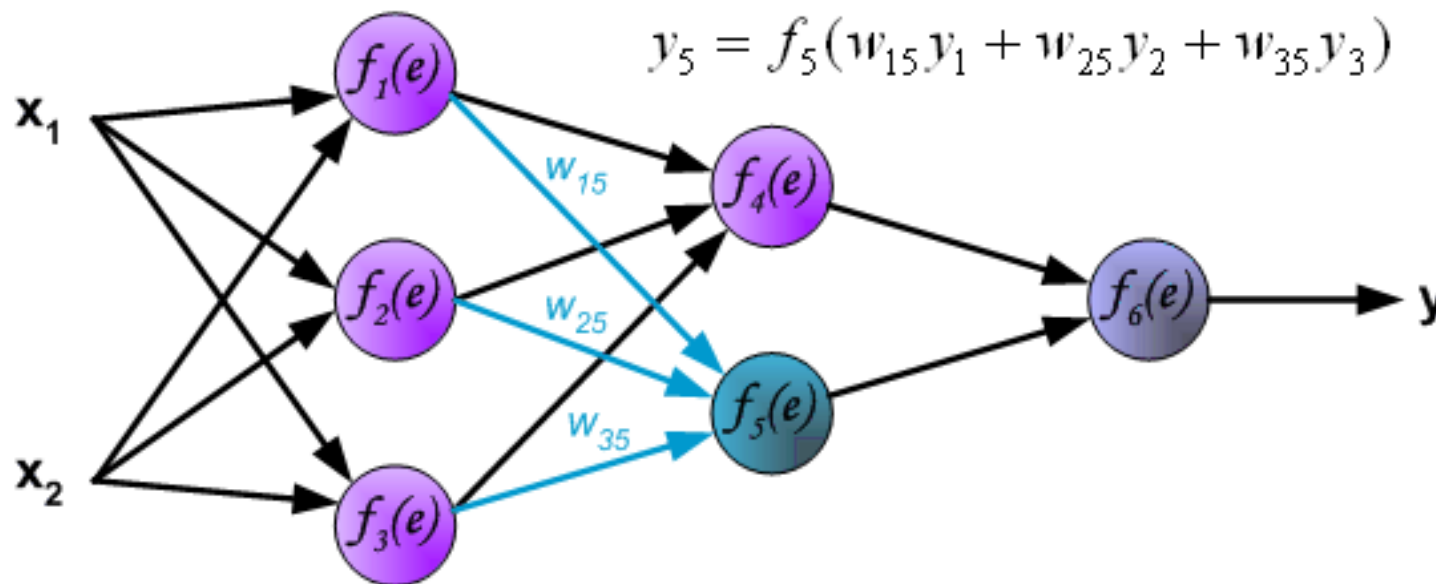


Feed forwarding

Propagation of signals through the hidden layer. Symbols w_{mn} represent weights of connections between output of neuron m and input of neuron n in the next layer.

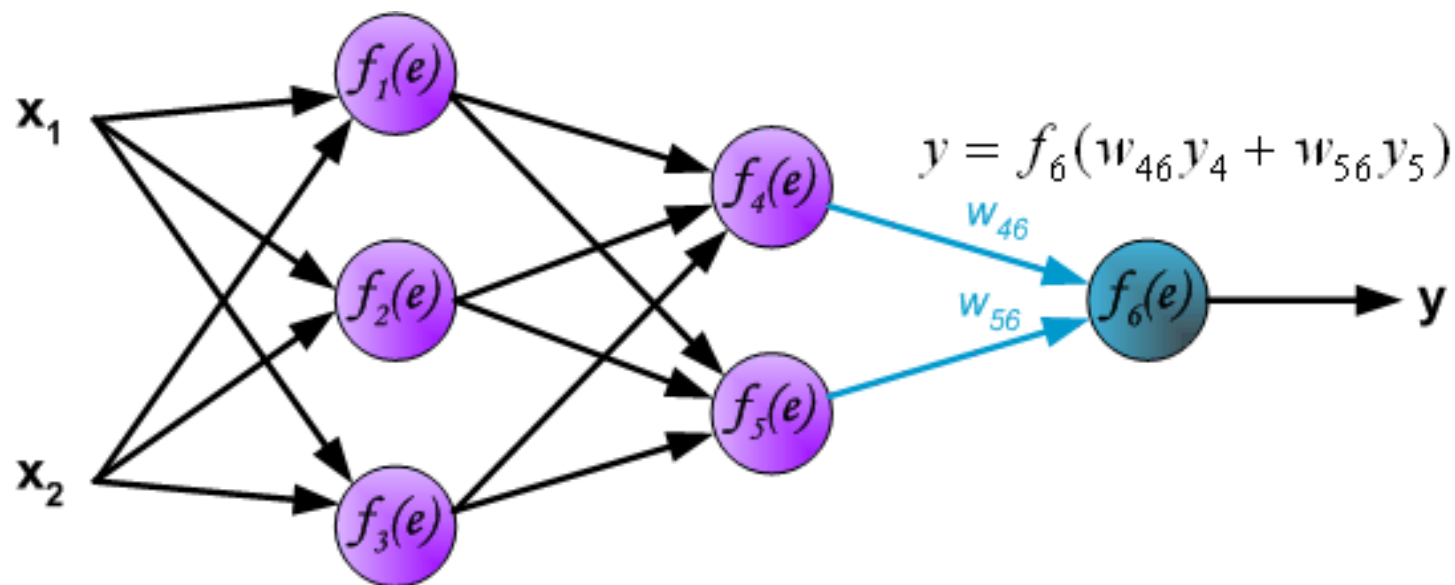


Feed forwarding



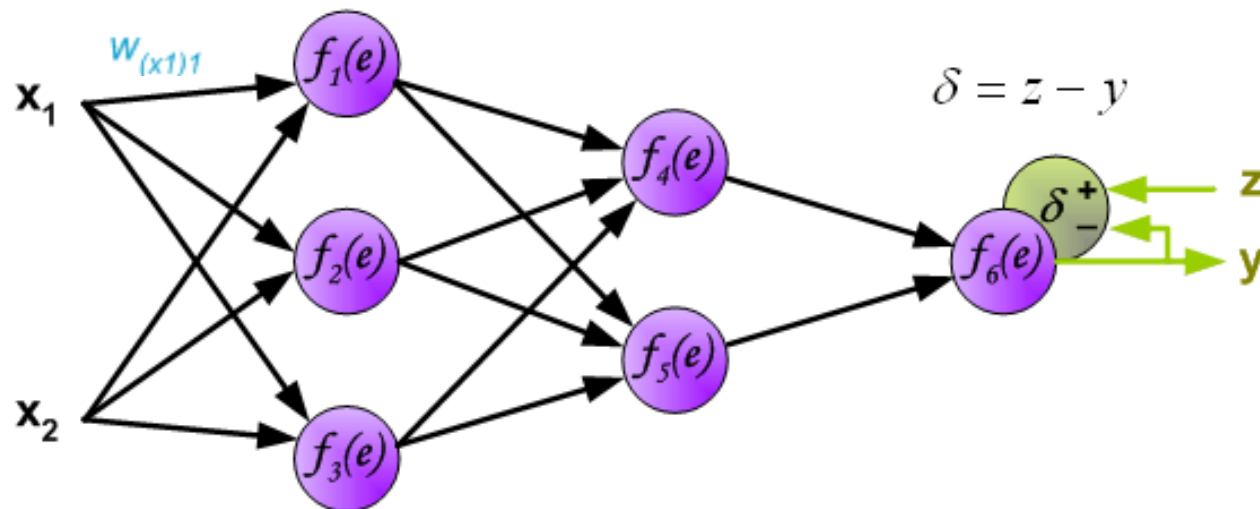
Feed forwarding

Propagation of signals through the output layer.



- Neutral Network
- Loss Function
- Gradient descent
- Forward Propagation
- Backward Propagation
 - Back propagation problem
 - Chain Rule
 - Computational graph
 - Calculation
 - Generalizations

Back Propagation

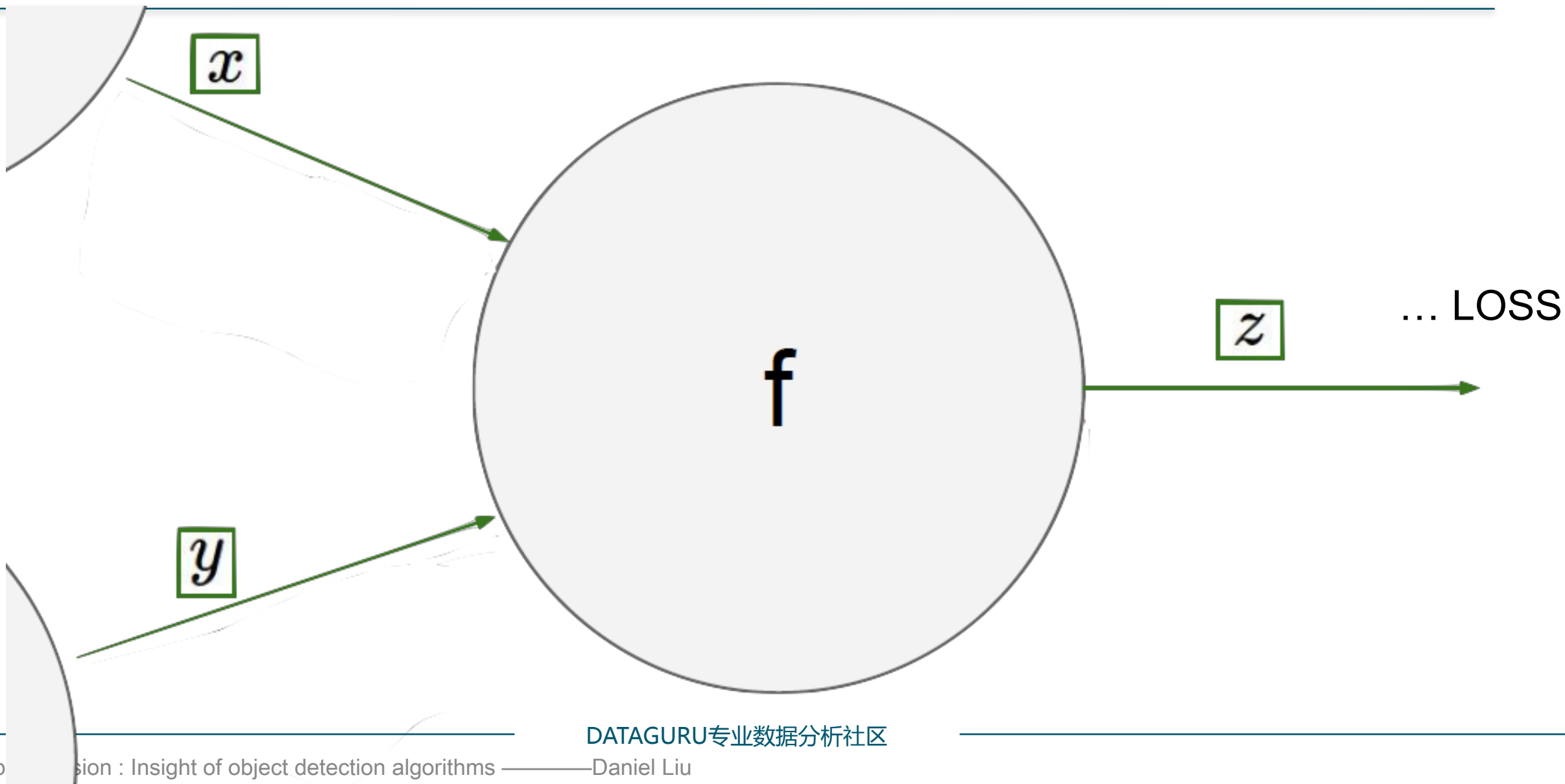


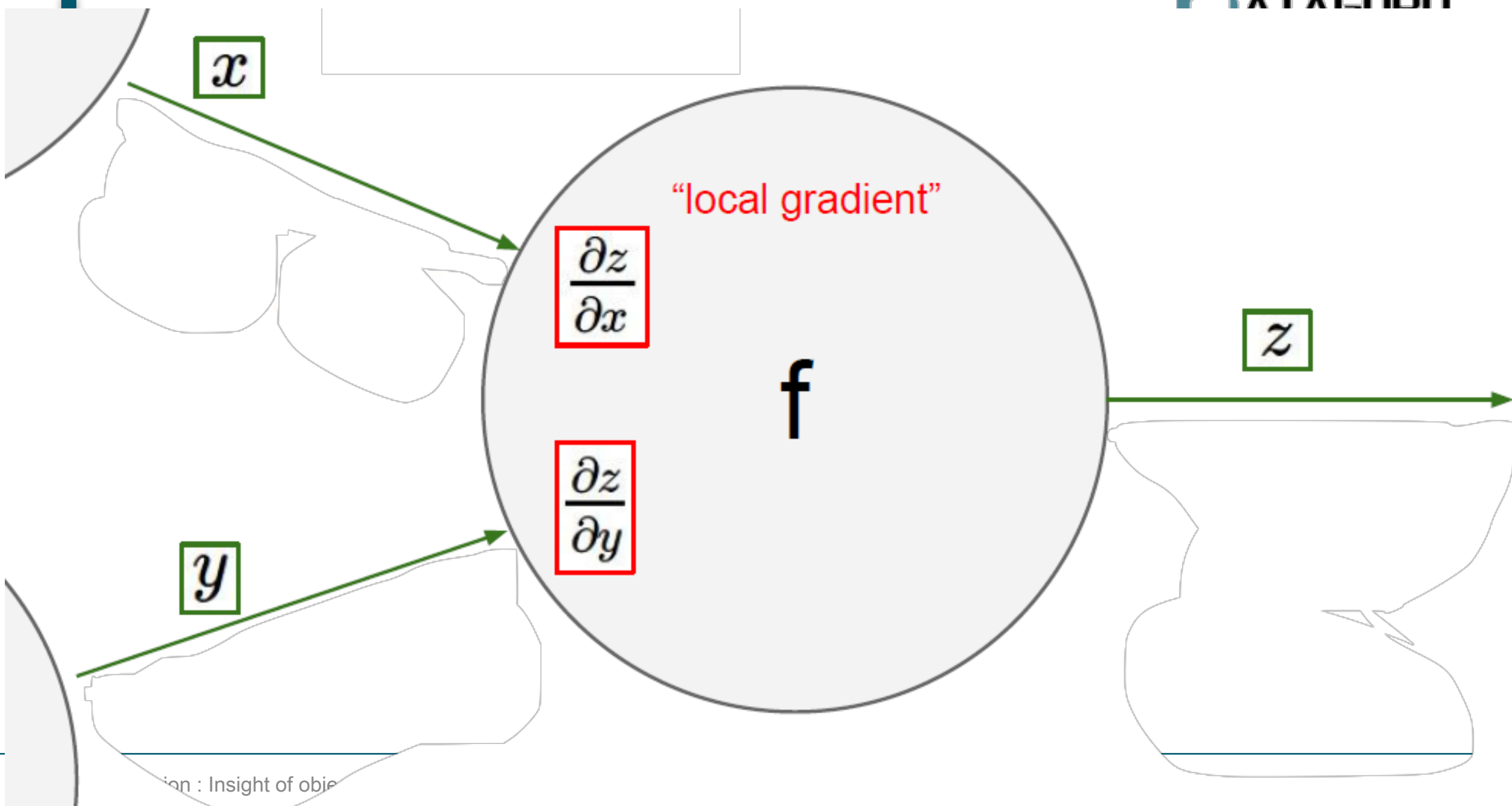
Gradient of **loss**
with respect to w $\frac{\partial \text{loss}}{\partial w} = ?$

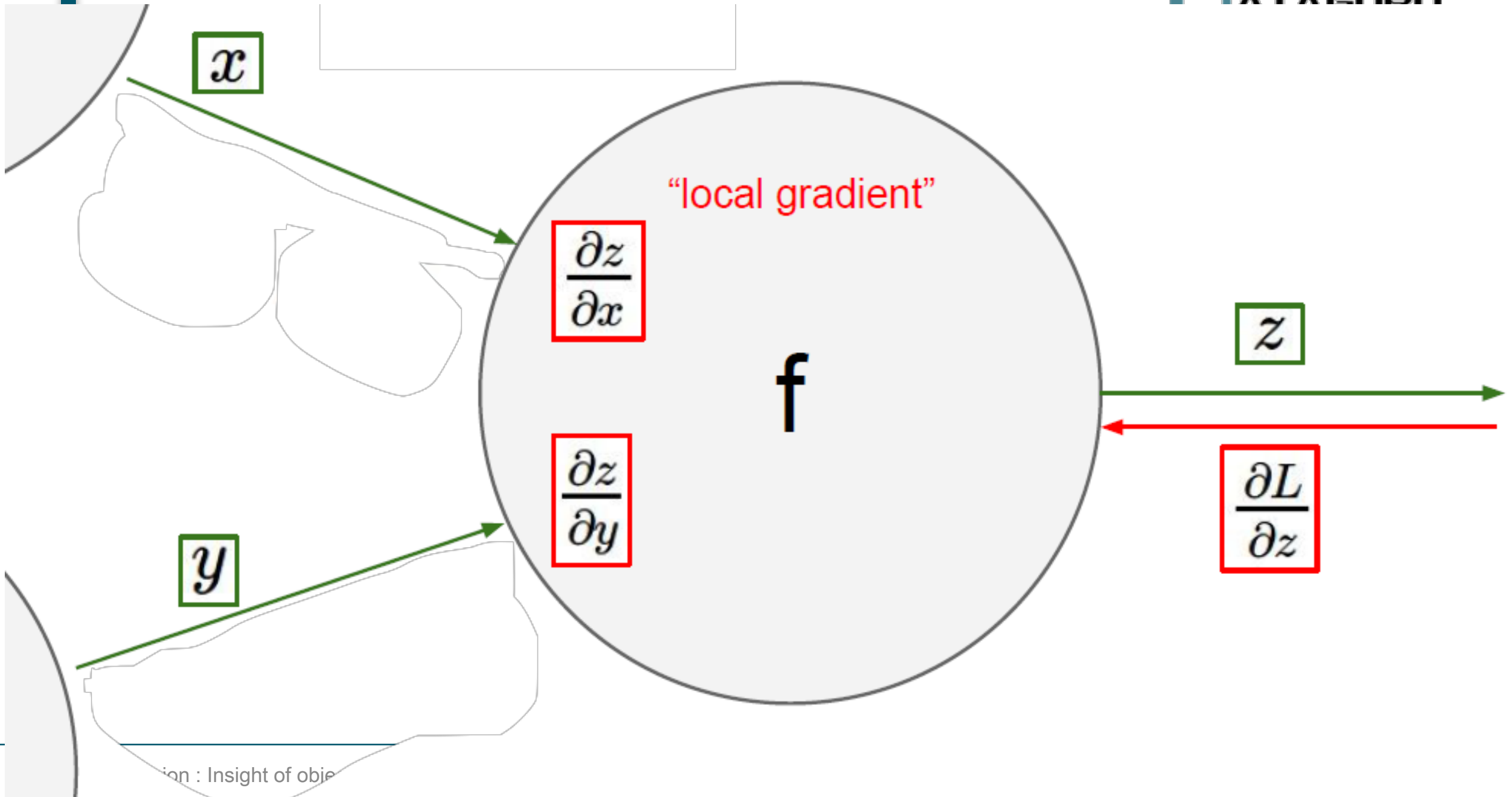
$$\frac{\partial \delta}{\partial w} = ?$$

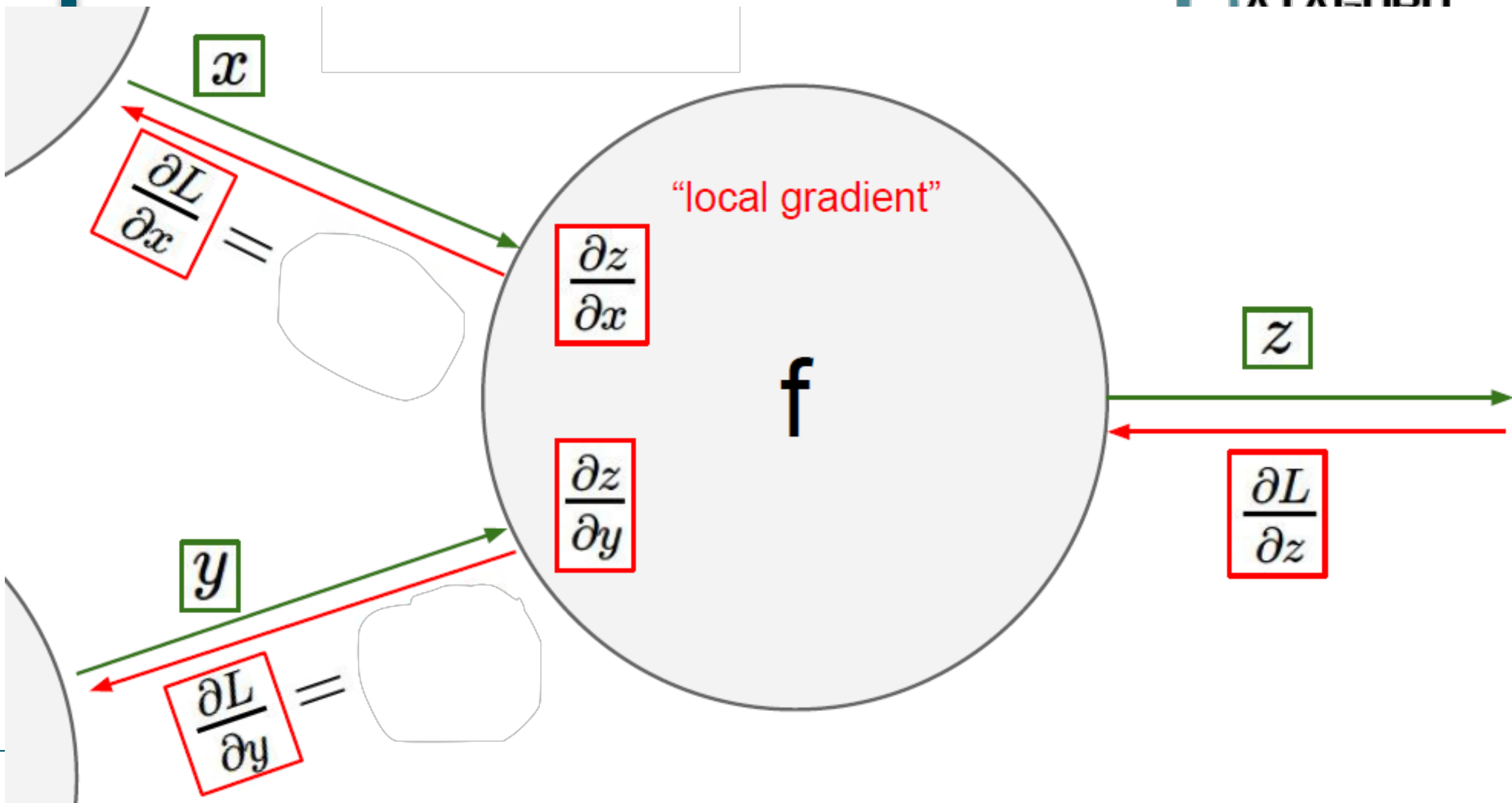
- Neutral Network
- Loss Function
- Gradient descent
- Forward Propagation
- Backward Propagation
 - Back propagation problem
 - Chain Rule
 - Computational graph
 - Calculation
 - Generalizations

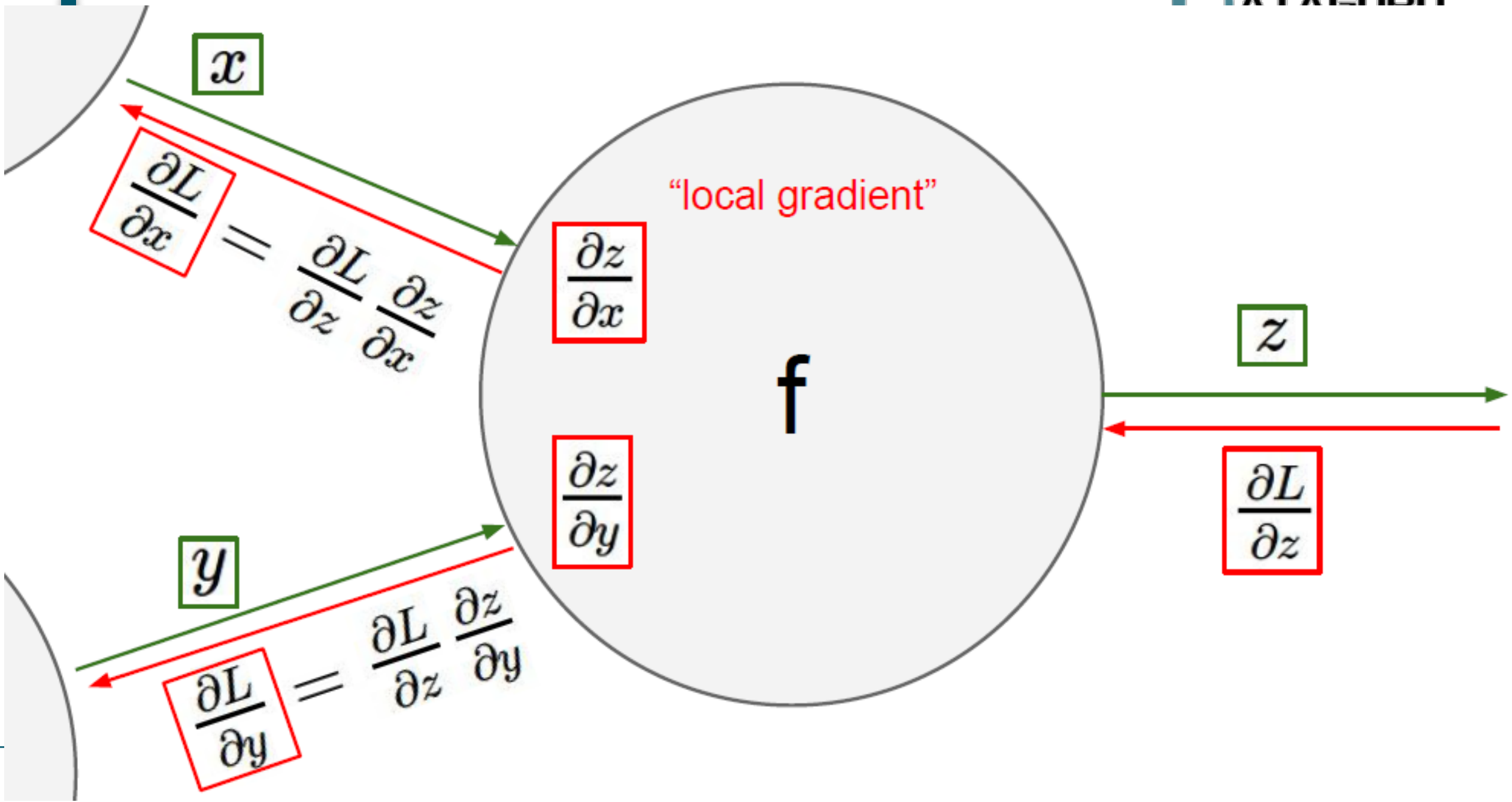
Chain rule

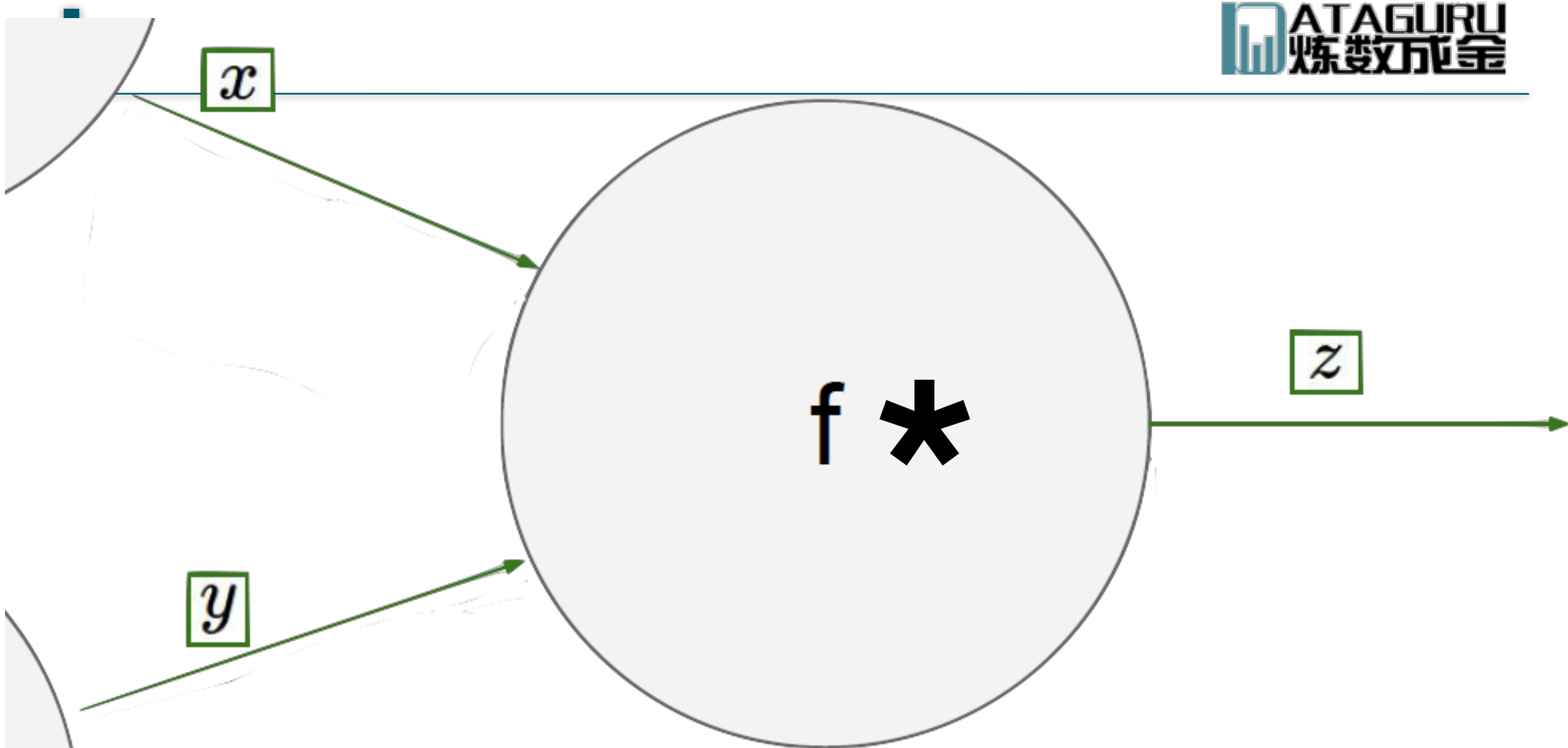












1

Forward pass $x = 2, y = 3$

x

$f *$

z

$$x = 2$$

"local gradient"

$$\frac{\partial z}{\partial x}$$

$$\frac{\partial z}{\partial y}$$

$f *$

$$z = 6$$

$$y = 3$$

2

Backward propagation

$\frac{\partial L}{\partial z} = 5$ is given

$x = 2$

“local gradient”

$\frac{\partial z}{\partial x}$

$\frac{\partial z}{\partial y}$

$f *$

$z = 6$

$\frac{\partial L}{\partial z} = 5$

2

Backward propagation

$\frac{\partial L}{\partial z} = 5$ is given

$x = 2$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x}$$

“local gradient”

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

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

$f *$

$z = 6$

$$\frac{\partial L}{\partial z} = 5$$

$y = 3$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y}$$

2

Backward propagation

$\frac{\partial L}{\partial z} = 5$ is given

$x = 2$

$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x}$
 $= 5 * y = 15$

$y = 3$

$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y}$
 $= 5 * x = 10$

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

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

“local gradient”

$f *$

$z = 6$

$\frac{\partial L}{\partial z} = 5$

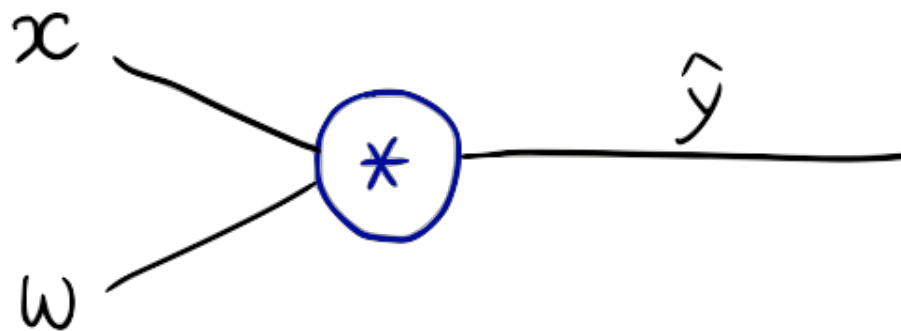
- Neutral Network
- Loss Function
- Gradient descent
- Forward Propagation
- Backward Propagation
 - Back propagation problem
 - Chain Rule
 - Computational graph
 - Calculation
 - Generalizations

Computational graph

$$\hat{y} = x * w$$

Computational graph

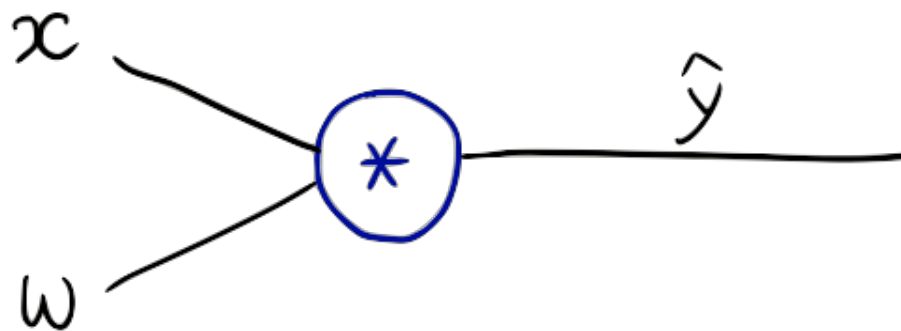
$$\hat{y} = x * w$$



Computational graph

$$\hat{y} = x * w$$

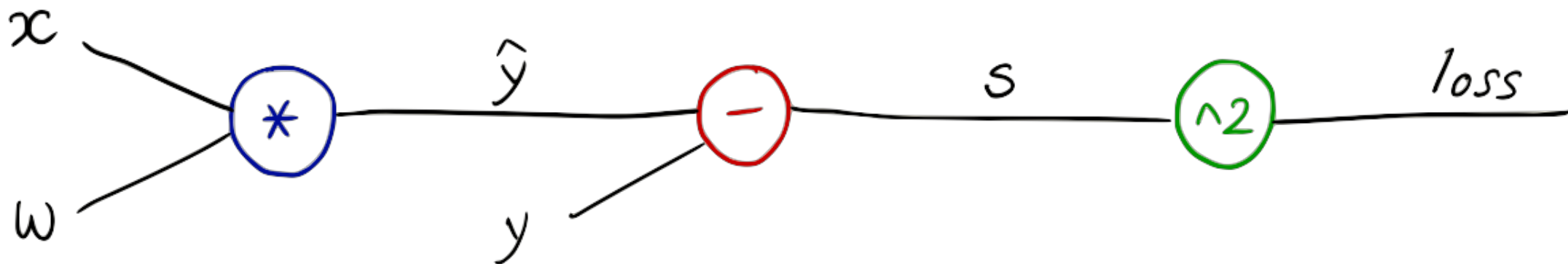
$$loss = (\hat{y} - y)^2 = (x * w - y)^2$$



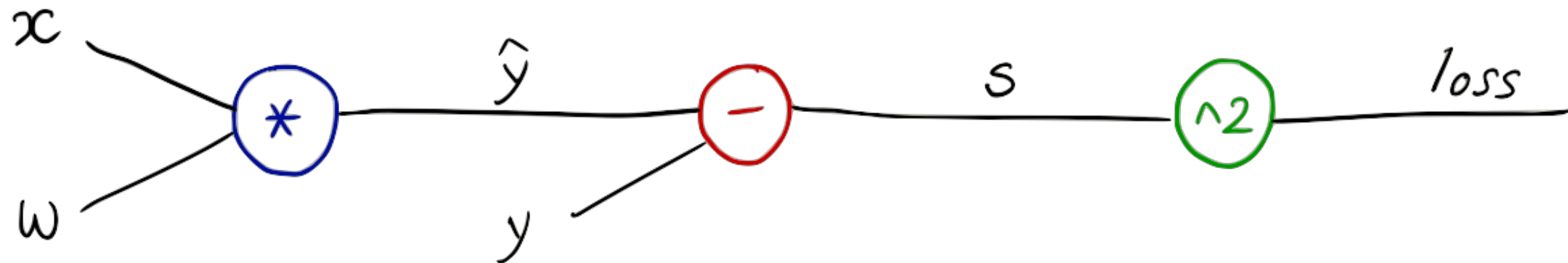
Computational graph

$$\hat{y} = x * w$$

$$loss = (\hat{y} - y)^2 = (x * w - y)^2$$

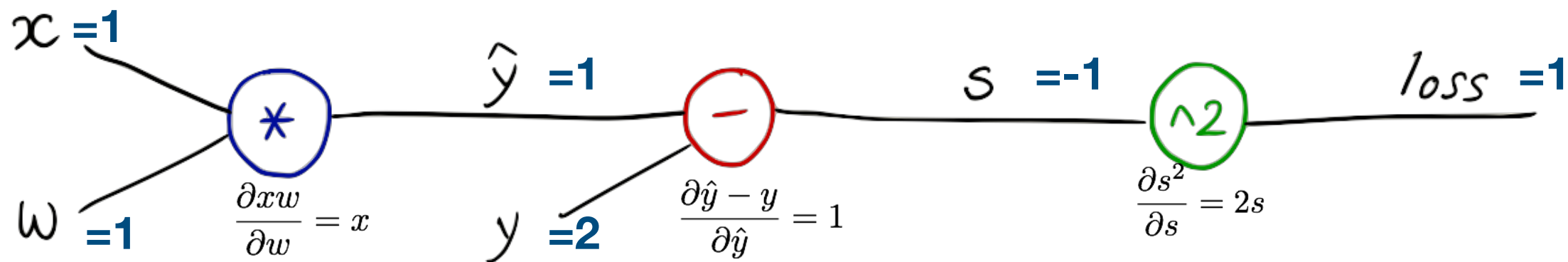


1 Forward pass $x=1, y=2$ where $w=1$



2

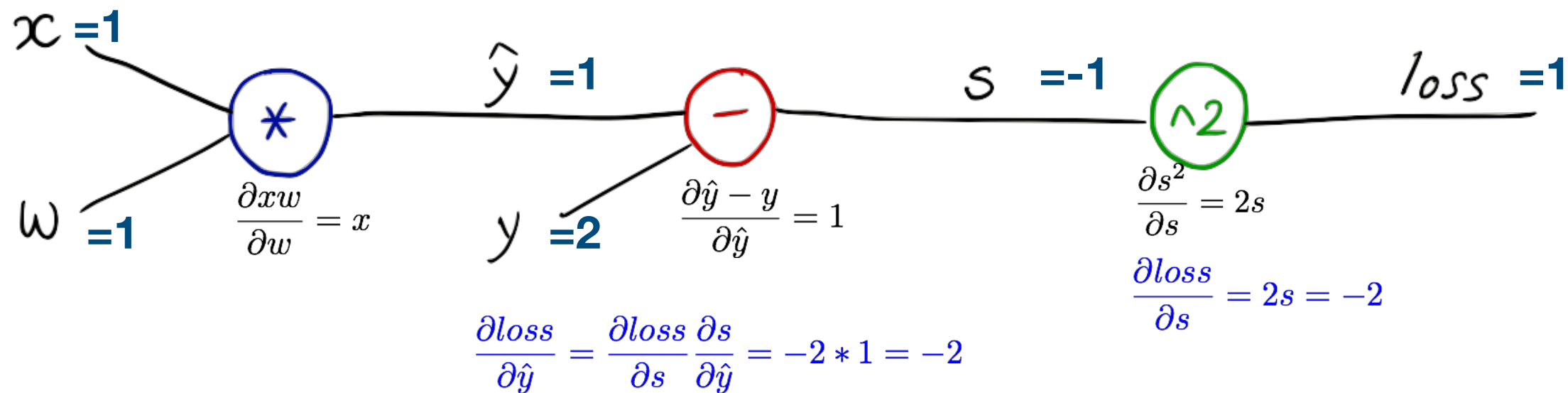
Backward propagation



$$\frac{\partial loss}{\partial w} =$$

2

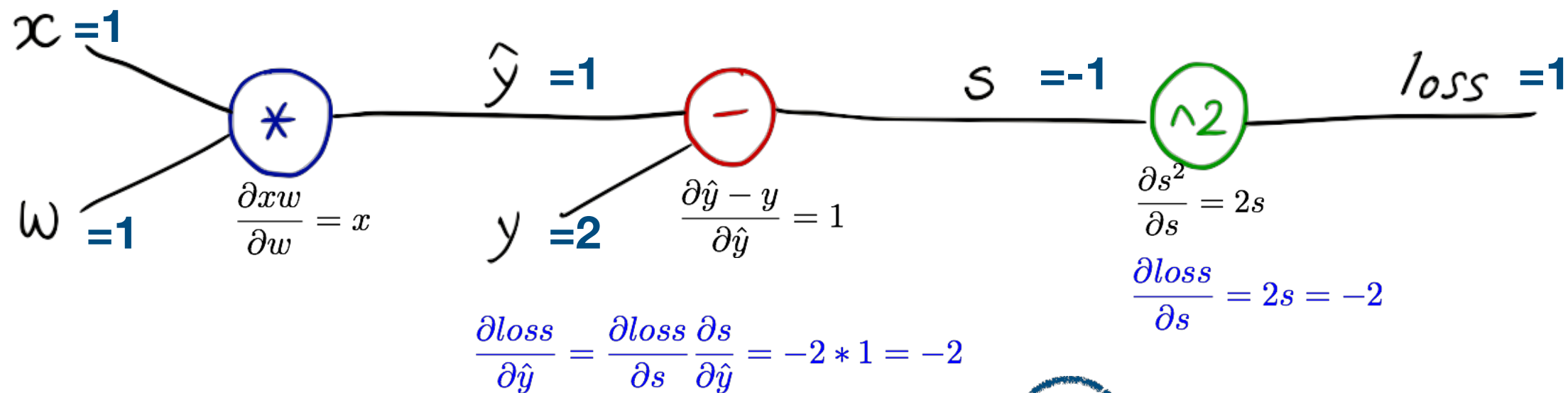
Backward propagation



$$\frac{\partial loss}{\partial w} = \frac{\partial loss}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w} = -2 * x = -2 * 1 = -2$$

2

Backward propagation



$$\frac{\partial loss}{\partial w} = \frac{\partial loss}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w} = -2 * x = -2 * 1 = -2$$

DATA

Daniel

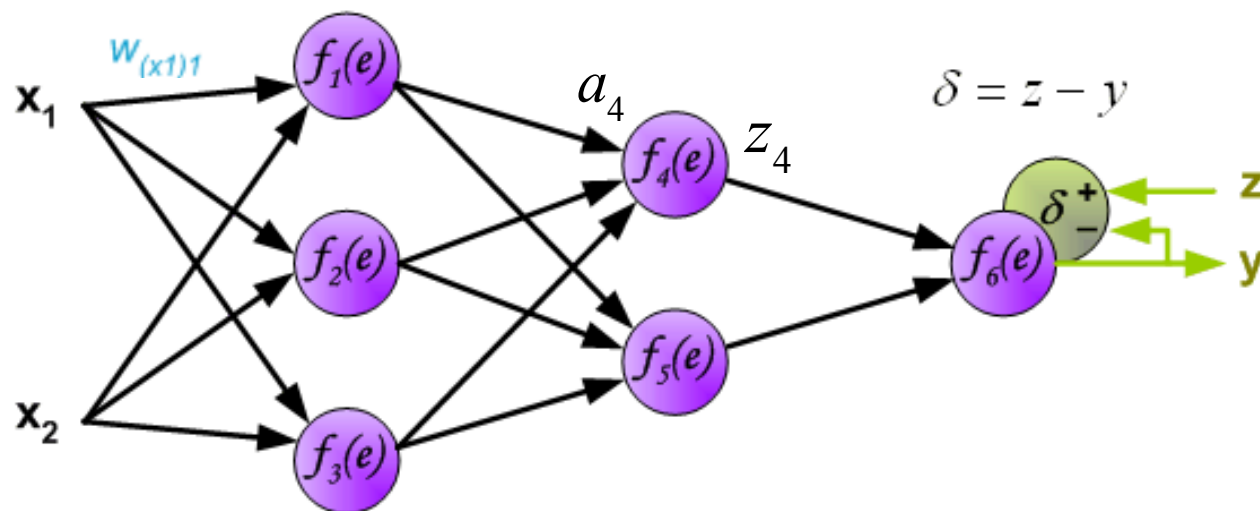
```
grad: 1.0 2.0 -2.0
grad: 2.0 4.0 -7.84
grad: 3.0 6.0 -16.2288
progress: 0 4.919240100095999
```

Table of contents



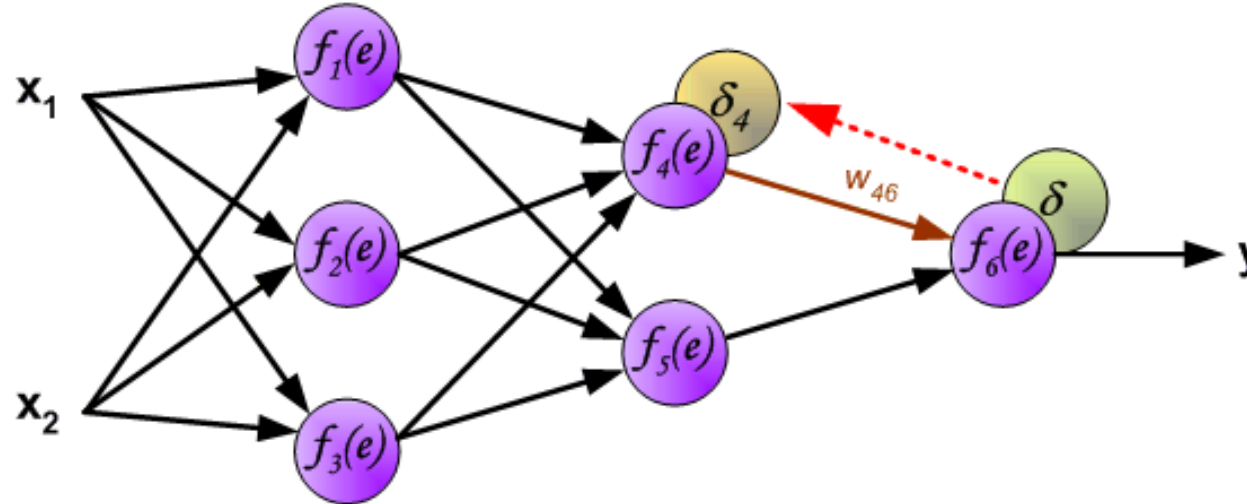
- Neutral Network
- Loss Function
- Gradient descent
- Forward Propagation
- Backward Propagation
 - Back propagation problem
 - Chain Rule
 - Computational graph
 - Calculation
 - Generalizations

Backward propagation



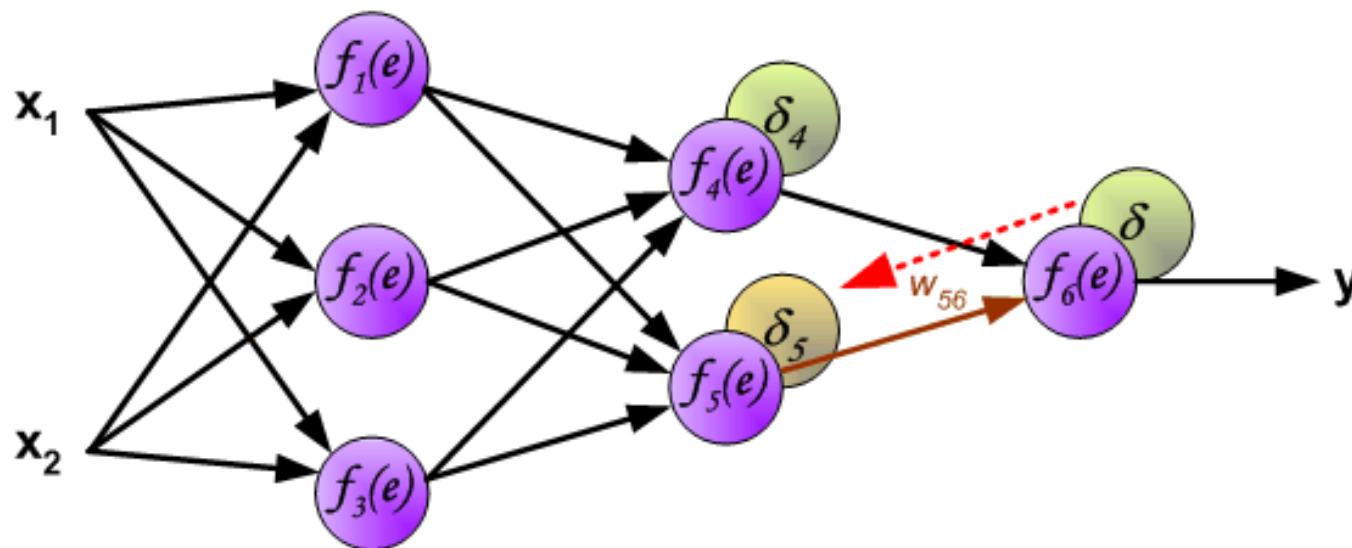
$$\frac{\partial \delta}{\partial z_4} = ?$$

Backward propagation to the 2nd layer



$$\begin{aligned}\delta_4 &= \delta \frac{\partial \delta}{\partial z_4} = \delta \frac{\partial (W_{46} z_4 + b_{46})}{\partial z_4} \cdot f'_6 \\ &= \delta W_{46} f'_6\end{aligned}$$

Backward propagation to the 2nd layer

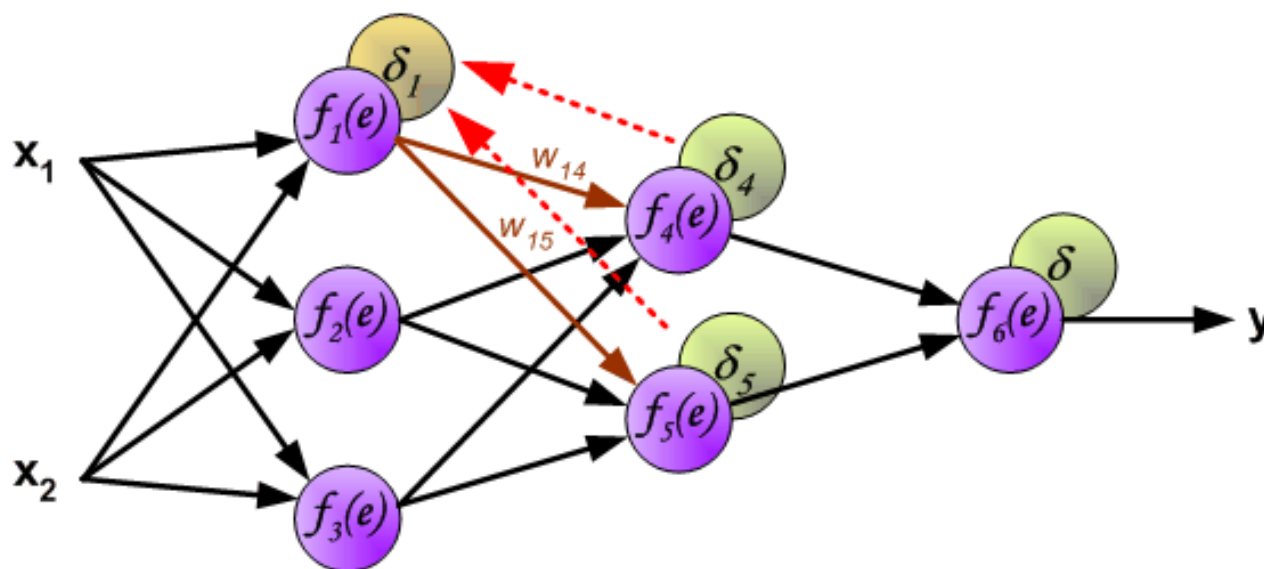


$$\delta_4 = \delta \frac{\partial \delta}{\partial z_4} = \delta \frac{\partial (W_{46} z_4 + b_{46})}{\partial z_4} \cdot f'_6$$

$$= \delta W_{46} f'_6$$

$$\delta_5 = \delta W_{56} f'_6$$

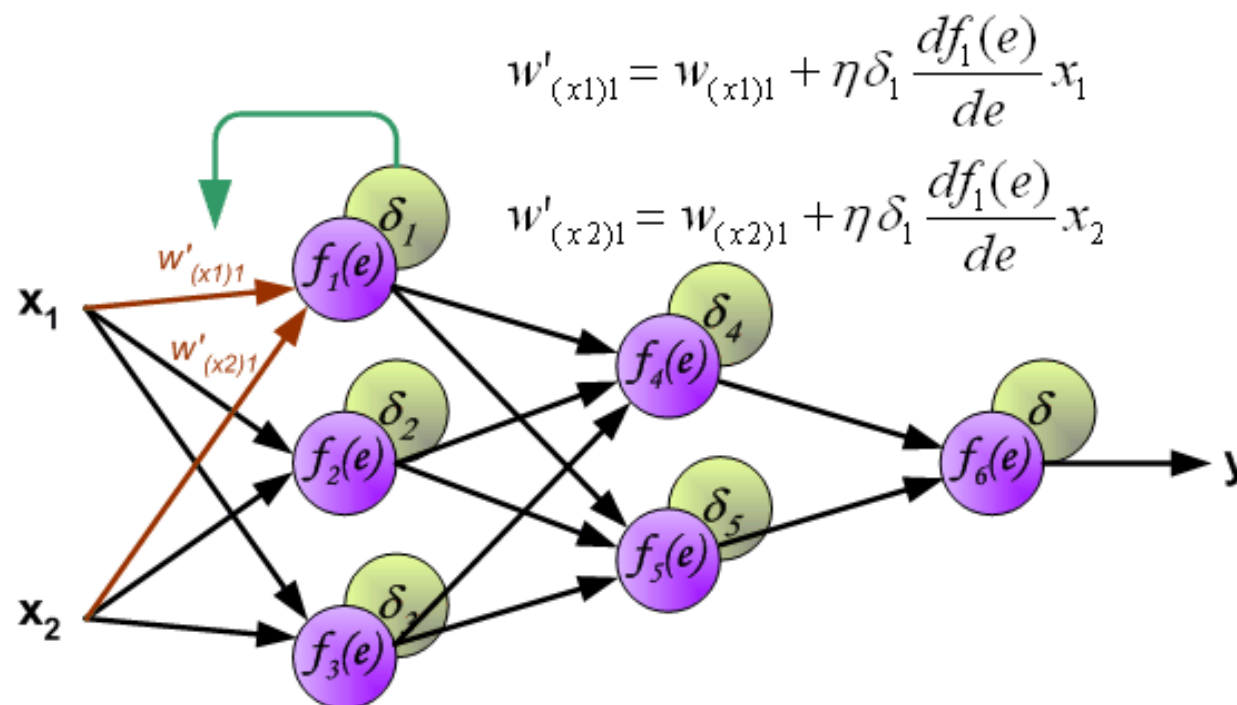
Backward propagation to the 1st layer



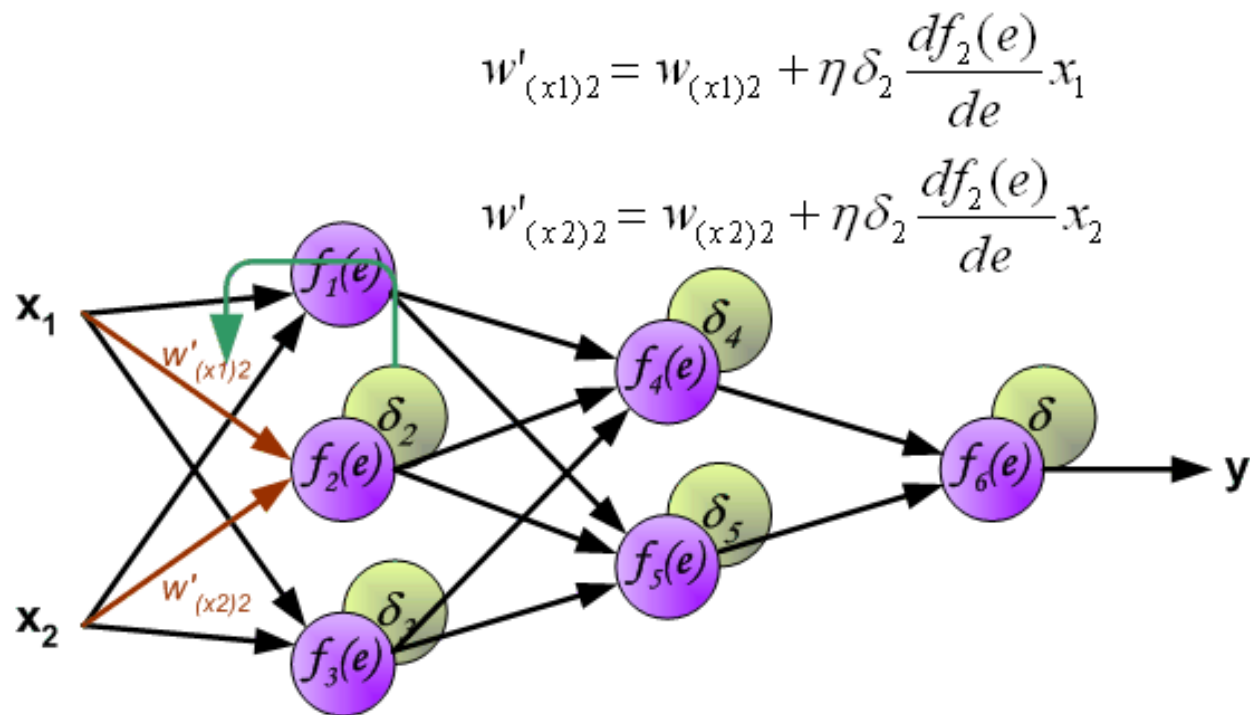
$$\begin{aligned}\delta_1 &= \delta_4 W_{14} f'_4 + \delta_5 W_{15} f'_5 \\ &= \delta f'_6 (W_{14} W_{46} f'_4 + W_{56} W_{15} f'_5)\end{aligned}$$

1st Layer Refresh

When the error signal for each neuron is computed, the weights coefficients of each neuron input node may be modified. In formulas below $df(e)/de$ represents derivative of neuron activation function (which weights are modified).

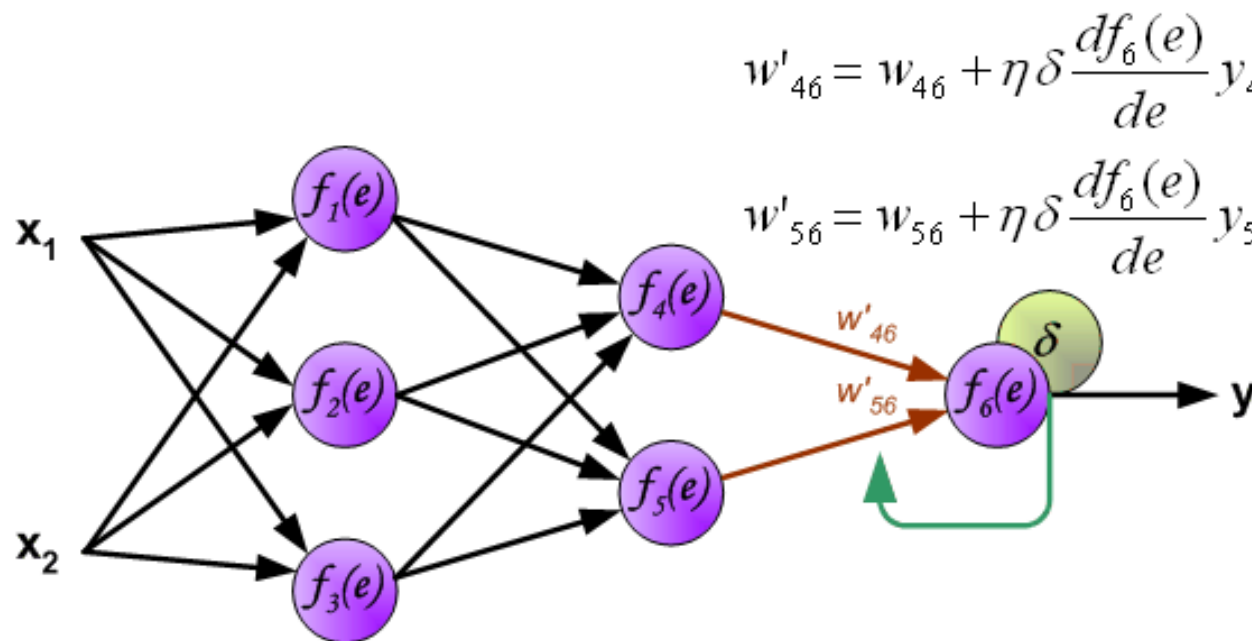


1st Layer Refresh



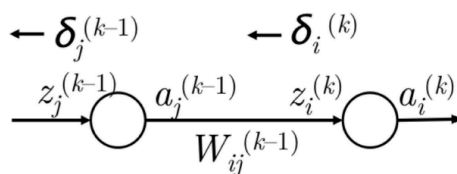
2nd Layer Refresh

When the error signal for each neuron is computed, the weights coefficients of each neuron input node may be modified. In formulas below $df(e)/de$ represents derivative of neuron activation function (which weights are modified).



- Neutral Network
- Loss Function
- Gradient descent
- Forward Propagation
- Backward Propagation
 - Back propagation problem
 - Chain Rule
 - Computational graph
 - Calculation
 - Generalizations

Backpropagation from δ^k to δ^{k-1}



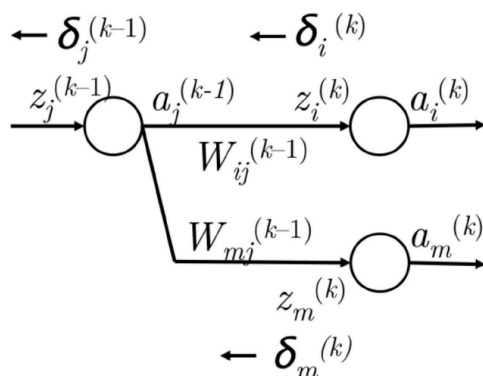
We have error $\delta_i^{(k)}$ propagating backwards from $z_j^{(k)}$, i.e. neuron i at layer k .

We propagate this error backwards to $a_j^{(k-1)}$ by multiplying $\delta_i^{(k)}$ by the path weight $W_{ij}^{(k-1)}$.

Thus, the error received at $a_j^{(k-1)}$ is $\delta_i^{(k)} W_{ij}^{(k-1)}$.

However, $a_j^{(k-1)}$ may have been forwarded to multiple nodes in the next layer (i.e. node m in layer k).

Thus, the total error received at $a_j^{(k-1)}$ is $\delta_i^{(k)} W_{ij}^{(k-1)} + \delta_m^{(k)} W_{mj}^{(k-1)}$.



In fact, we can generalize this to be $\sum_i \delta_i^{(k)} W_{ij}^{(k-1)}$.

Now, we can have the correct error at $a_j^{(k-1)}$, we move it across neuron j at layer $k-1$ by multiplying with the local gradient $\sigma'(z_j^{(k-1)})$.

Finally, the error that reaches at $z_j^{(k-1)}$, called $\delta_j^{(k-1)}$ is $\sigma'(z_j^{(k-1)}) \sum_i \delta_i^{(k)} W_{ij}^{(k-1)}$.

Thanks

FAQ时间