# NN Lecture - 3
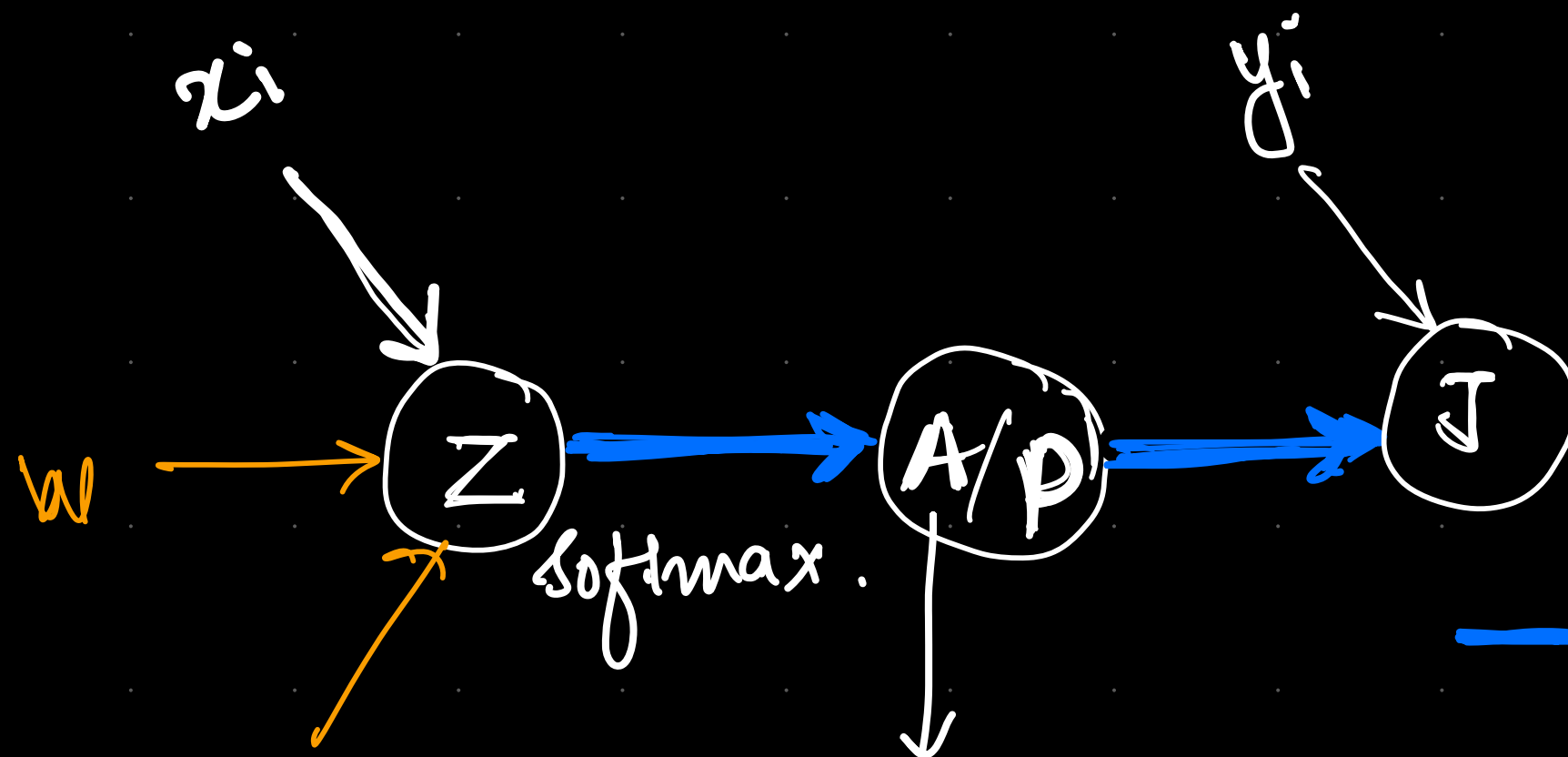
Backward Prop- Softmax classifier.

N-layer Neural Network.

# Computational graph for softmax classifier
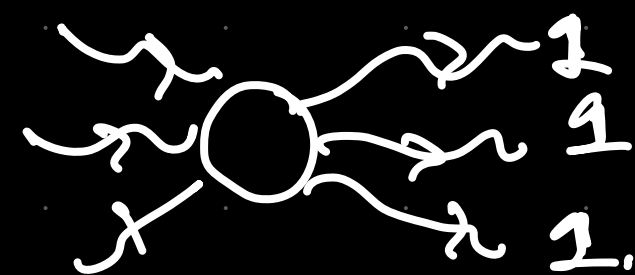
$x_i$

$y_i$

$K$ = #classes.



W

b

Z   softmax.

A/p

J

CG loss = $\sum_{1}^{K} y_k \log p_k$

— Forward Propogation.

$xw+b$ .
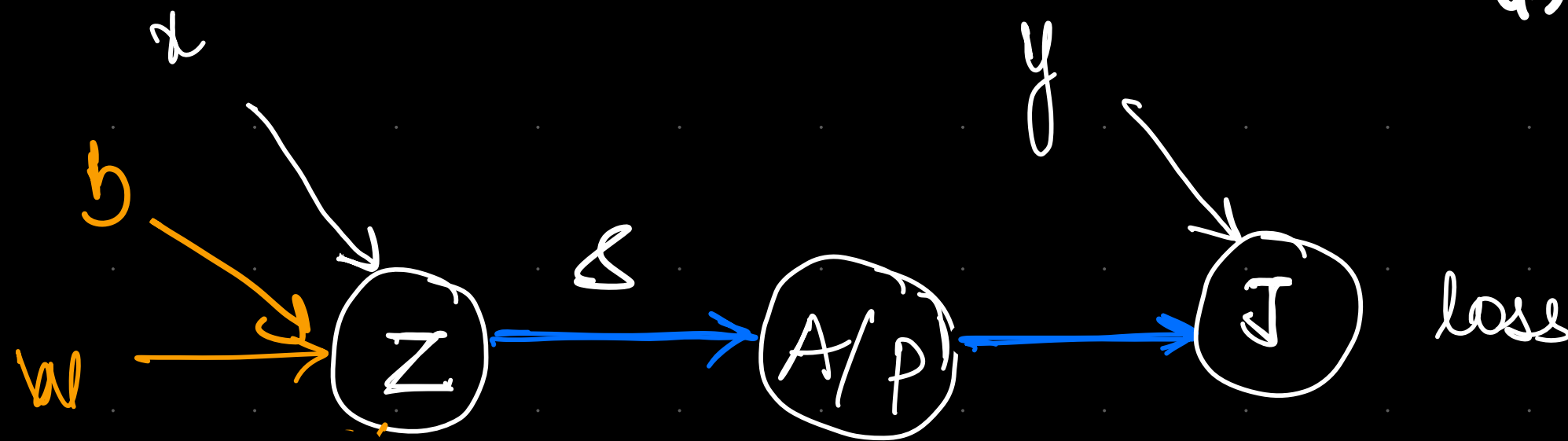
A- Activation.

P- probability

Computational graph.

1
1
1.
Activated.

# Backward Propogation $\left(\frac{\partial J}{\partial w}, \frac{\partial J}{\partial b}\right)$
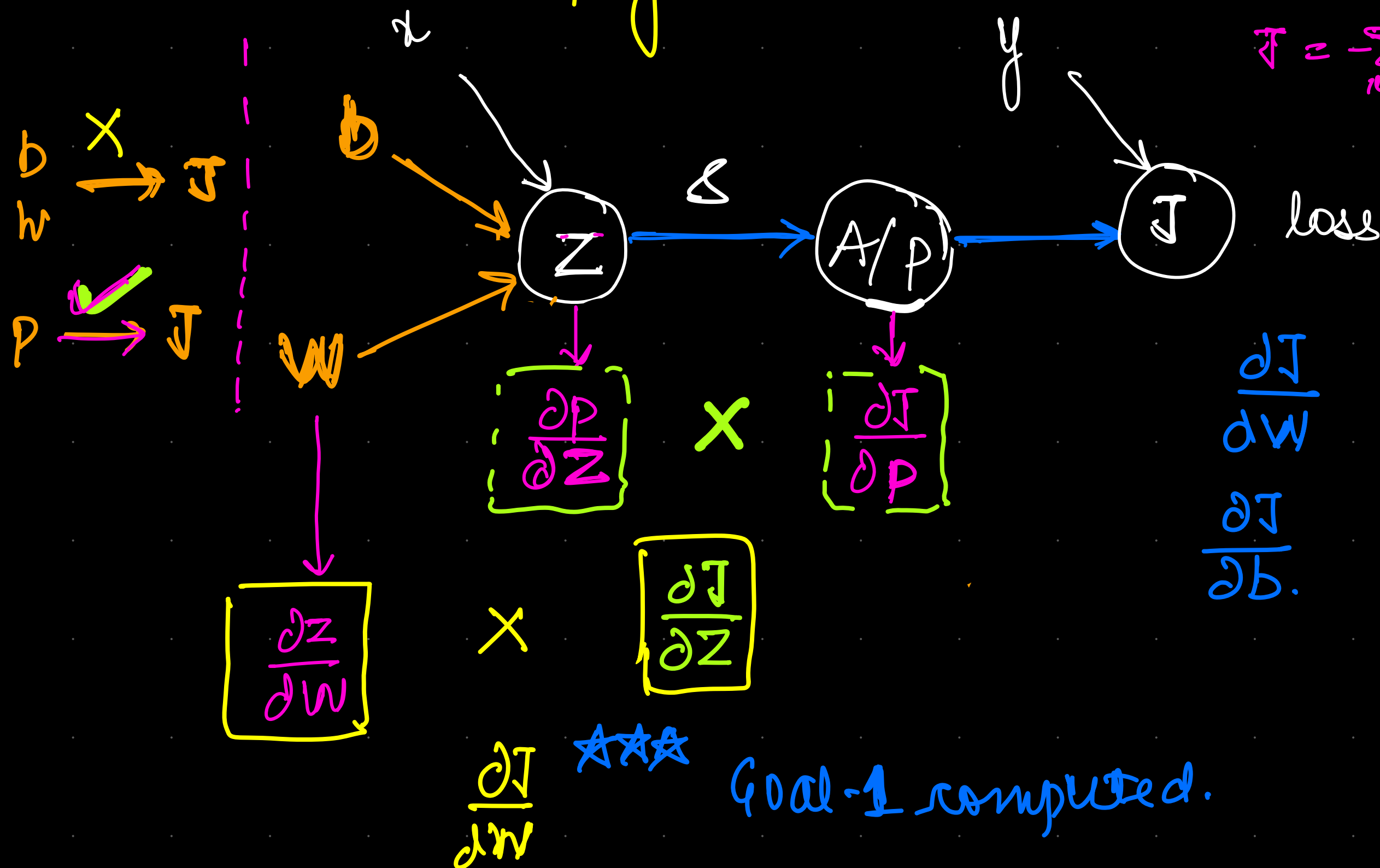
Gradients.



Goal : Calculate $\frac{\partial J}{\partial w}, \frac{\partial J}{\partial b}$

# Backward Propogation

$$Z = WX + b$$
$$P = e^Z / \Sigma e^Z$$
$$J = -\sum_k y \log P_k$$

x

b
w    X
  →  J

P  →  J

b  →  Z  →  A/P  →  J   loss

y

$$\boxed{\frac{\partial P}{\partial Z}} \times \boxed{\frac{\partial J}{\partial P}}$$

$$\boxed{\frac{\partial Z}{\partial W}} \times \boxed{\frac{\partial J}{\partial Z}}$$

$$\frac{\partial J}{\partial W}$$

$$\frac{\partial J}{\partial W}$$

$$\frac{\partial J}{\partial b}$$

★★★  Goal-1 computed.

# Backward Propogation chain rule

Chain rule for $\partial J / \partial W$:

$$\frac{\partial J}{\partial W} = \frac{\partial J}{\partial p} \frac{\partial p}{\partial z} \frac{\partial z}{\partial W}$$

Chain rule for $\partial J / \partial b$:

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial p} \frac{\partial p}{\partial z} \frac{\partial z}{\partial b}$$

# Backward Propogation - shorthands

chain rule for $\partial J/\partial w$

$$\frac{\partial J}{\partial w} = \underset{1}{\frac{\partial J}{\partial p}} \; \underset{2}{\frac{\partial p}{\partial z}} \; \underset{3}{\frac{\partial z}{\partial w}}$$

$\partial J/\partial \square$

Rule: All derivatives for $\boxed{\partial J/\partial \square}$ can be written as $\boxed{\partial \square}$

$dw$     $\partial p$     $\partial z$    $dw$

✓     ✓     X    X

# Backward Propogation - shorthands

chain rule for $\partial J / \partial w$

$$\frac{\partial J}{\partial w} = \boxed{\frac{\partial J}{\partial p} \frac{\partial p}{\partial z}} \frac{\partial z}{\partial w}$$

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial z} \frac{\partial z}{\partial w}$$

✦✦✦

$$\boxed{\partial w = \partial z \; \partial z / \partial w}$$

That is only for coding.

↓

V·V·Nice solution.

→ Very easy.

$$\frac{\partial z}{\partial w} = \frac{\partial (wx + b)}{\partial w}$$
$$= x$$

# Backward Propogation - $\partial z$, $\partial z/\partial w$  ✗✗

① $\dfrac{\partial z}{\partial w} = \dfrac{\partial(wx+b)}{\partial w} = \boxed{x}$

$$dw = \boxed{(p-y)x}$$

Near solution.

② $\partial z = P_k - \boxed{I(k=y)}$ ← Indicator function.

Probab.f class(k)

$k \neq y$ → 0

$k = y$ → 1

How? Derivation in post-read (Difficult)
Not expected.

code :

$$\boxed{\partial z = p - y}$$

Prob vector     Ground. T. vector.

$\begin{bmatrix} P_1 \\ P_2 \\ P_3 \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$

✗✗ Residual (like)

But why are we $\boxed{\text{subtracting 1}}$ ?

Ground Truth $\quad y = [\ 0 \ 1 \ 0\ ]$    Intuitively

$$\overset{P_1 \quad\quad P_2 \quad\quad\quad P_3}{\text{Predicted. probs} = [\ 0.2,\ 0.3\ ,\ 0.5\ ]}$$

$$\frac{\partial J}{\partial z} = \boxed{dz} = [\ 0.2\ ,\ -0.7,\ 0.5\ ]$$

★★★
$$\boxed{dz = p - y}$$

error.

① $P_2 \uparrow \longrightarrow \partial z_2 \downarrow$, $P_2 \downarrow \longrightarrow \partial z_2 \uparrow$

② $P_1 \uparrow \longrightarrow \partial z_1 \uparrow$, $P_1 \downarrow \longrightarrow \partial z_2 \downarrow$

(or $P_3$)

# lets calculate $dW$

Input Matrix.

$$dW = \partial Z \cdot X$$

$\partial Z \longrightarrow p-y = \boxed{m \times n} - \boxed{m \times n} = (m, n) = (300, 3)$

$X \longrightarrow (m, d) \Longrightarrow (300, 2)$

$\boxed{\partial W}$ shape $\longrightarrow (d, n) \longrightarrow (2, 3)$ #Classes

$\boxed{W = W - \alpha \, \partial W}$

$\parallel$ features   Neuron   "   same shape.

$W$

$$(2, 300)(300, 3) \longrightarrow (2, 3)$$

$$X^T \cdot \partial Z = \partial W$$

# lets calculate $\partial b$

Nice

$$\partial b = \partial z \cdot \partial z / \partial b$$

$$= (p - y) \cdot 1$$

$$\boxed{\partial b = \partial z}$$

$$\frac{\partial z}{\partial b} = \frac{\partial (x \cdot w + b)}{\partial b}$$

$$= 1$$

[0.1 0.3 0.7]

argmax → Index

0.1 ②

$\partial z \longrightarrow P - y = (300, 3) = (m, n)$

$\partial b \longrightarrow (1, n)$

shape

$b (1, n)$

$\partial b = np.sum(\partial z, axis = 0, keepdims = True)$

Matrix

For m samples

$(1, n) = (m, n)$

Mean is the solution.

⭐⭐⭐

$m$ -YD array.

# Output of softmax classifier



Task-I

✓ : Adapted LRVs to work for multi-class classification

✗ : learning non-linear decision boundary