

# Agenda:

- ① Multi-class classfn: (softmax)
- ② CE-loss (generalization of log-loss)
- ③ Code → Python  
3-class classifier  
(backprop)  
→ TF2 & Keras  
= = = ↑ ↗  
GPU / TPU ...

④

Dropout

Binary - case

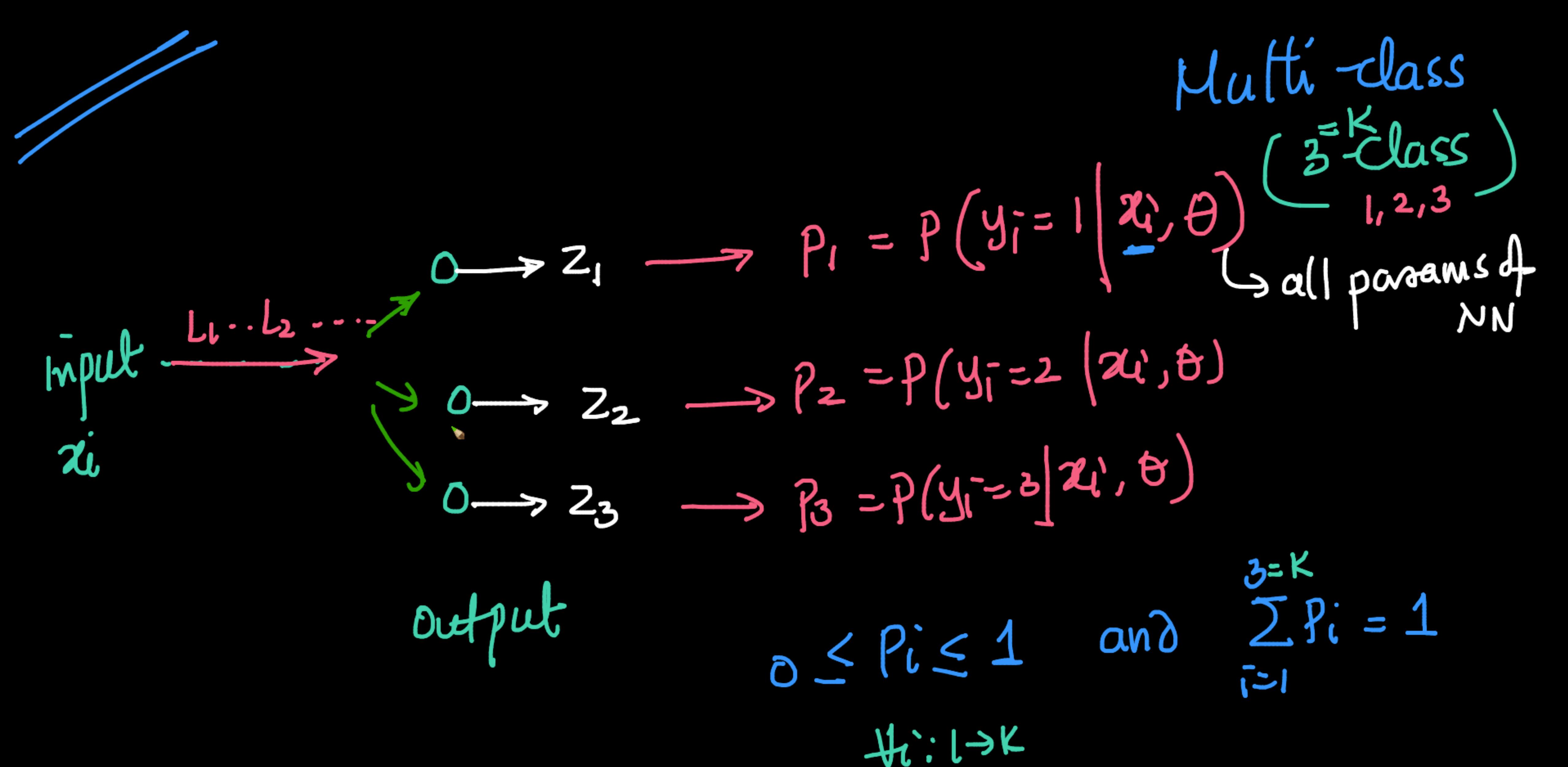
$$y_i \in \{0, 1\}$$

$$\left\{ \begin{array}{l} \omega^T x_i + b = z_i \\ \text{Sigmoid}(z_i) = P(y_i=1 | x_i, \omega, b) = p_i \\ \frac{1}{1+e^{-z_i}} = \frac{e^{z_i}}{1+e^{z_i}} \end{array} \right. \quad P(y_i=0 | x_i, \omega, b) = 1 - p_i$$

Log-loss: (Binary CE)

$$\{ = - \left[ y_i \log \hat{y}_i + (1-y_i) \log(1-\hat{y}_i) \right]$$

$\downarrow$   
 $p_i$



[Softmax:] ~

$$p_i = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

$$0 \leq p_i \leq 1$$

$$p_2 = \frac{e^{z_2}}{\sum_{i=1}^k e^{z_i}}$$

$$p_3 = \frac{e^{z_3}}{\sum_{i=1}^k e^{z_i}}$$

if  $k=3$ :

$$\underline{p_1 + p_2 + p_3 = 1}$$



softmax  $\approx$  sigmoid-like fn for  
multi-class Classfn

Q

$$P_i = \frac{z_i}{\sum_{i=1}^n z_i}$$

$\infty \leq z_i \leq \infty$



$$z_i : 2 \quad 6 \quad 8$$

$$e^{2i} : e^2 \quad e^6 \quad e^8$$

$$\rightarrow p_i : 0.662578$$

keisan.casio.com/exec/system/15168444286206

### SOFTMAX FUNCTION Calculator

Home / Special Function / Activation function

Calculates the softmax function.

The softmax function is used in the activation function of the neural network.

$\checkmark$   $P_i = \frac{z_i}{\sum z_j} \Rightarrow 1:3:6$

$a = \begin{matrix} \text{vector } a \\ \hline a_1 & 2 \\ a_2 & 6 \\ a_3 & 12 \end{matrix}$

$P_i = \frac{e^{z_i}}{\sum e^{z_i}} \Rightarrow ?$

Softmax

$0$

$0$

$1$

Max

$b$  value

b	value
b <sub>1</sub>	4.5285621964412E-5
b <sub>2</sub>	0.0024725111823572
b <sub>3</sub>	0.99748220319568

$\sigma(z)_j$

$$\sigma(z)_j = \frac{e^{z_j}}{\sum e^{z_k}} \text{ for } j = 1, \dots, K.$$

Not registered.

+ Bookmark

History

Softmax function

Related Calculator

- Sigmoid function
- Derivative Sigmoid function
- Second Derivative Sigmoid function
- Sigmoid function (chart)
- Softsign function
- Derivative Softsign function
- Softsign function (chart)
- Softplus function
- Derivative Softplus function
- Softplus function (chart)
- Softmax function**
- ReLU
- ReLU (chart)
- Leaky ReLU
- Leaky ReLU (chart)
- tanh(x) function
- Derivative tanh(x)

17 / 17

(Q) why  $e^z$  are common

✓ {  $\frac{\partial e^z}{\partial z} = e^z$

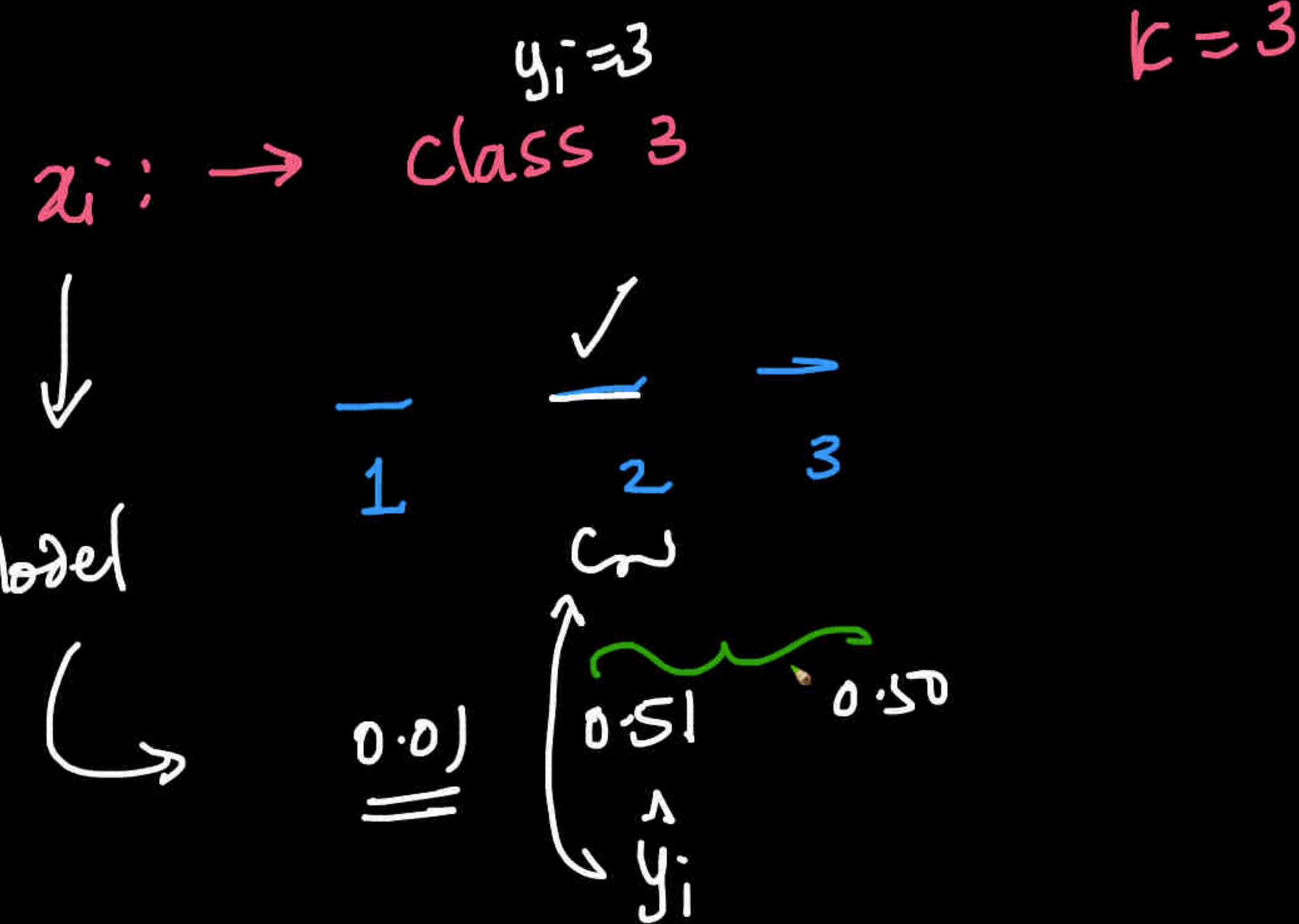
Q  
||

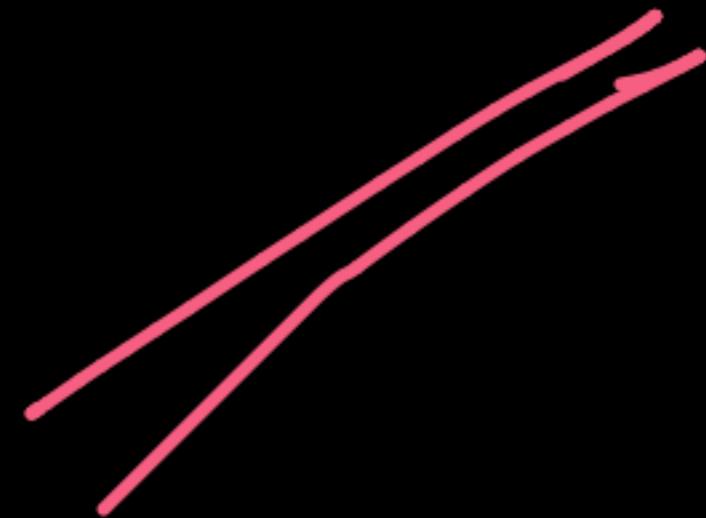
Binary:

$$\textcircled{O} \rightarrow z_i \rightarrow P(y_i=1) = p_1 \\ P_0 = 1 - p_1$$

k-ary

$$\textcircled{O} \rightarrow 1, 2, 3, \dots, k \quad (\text{integer?})$$





Log-loss in Binary classfn.

Multi-class classfn

← (Cross-Entropy)

•  $p_{i1} \ p_{i2} \ p_{ij} \dots p_{ik}$   
 $x_i \rightarrow 1 \ 2 \ \dots j \ \dots k$  classes  
↳  $\sum p_{ij} \ \forall j=1 \rightarrow k$

$x_i, y_i$	$\rightarrow$	$y_{i1} \quad y_{i2} \quad y_{i3} \quad y_{ij} \quad \dots \quad y_{ik}$	$\sigma$
0	1	0	-
1	2	$\dots j$	$\dots k$
$p_{i1} \quad p_{i2} \quad \dots p_{ij} \quad \dots p_{ik}$			

↑ Prob of  $x_i \in \text{class } j$

$$CE_i = - \sum_{j=1}^K y_{ij} \log(p_{ij})$$

$$CE_i = - \sum_{j=1}^k y_{ij} \log p_{ij} \quad y_i \in \text{class 1}$$

let  $k=2$

$$= - [y_{i1} \log p_{i1} + y_{i2} \log p_{i2}]$$

$y_{i1}$        $y_{i2}$

$p_{i1}$        $p_{i2}$

$(1-y_{i1})$        $(1-y_{i2})$

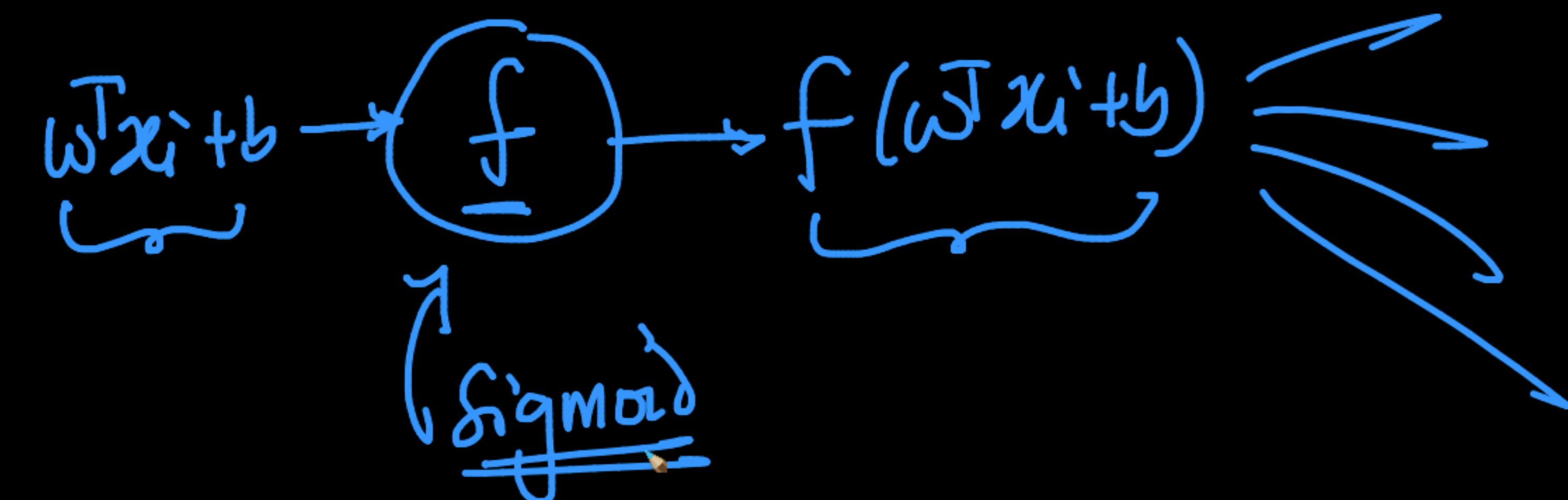
log-loss = Binary CE

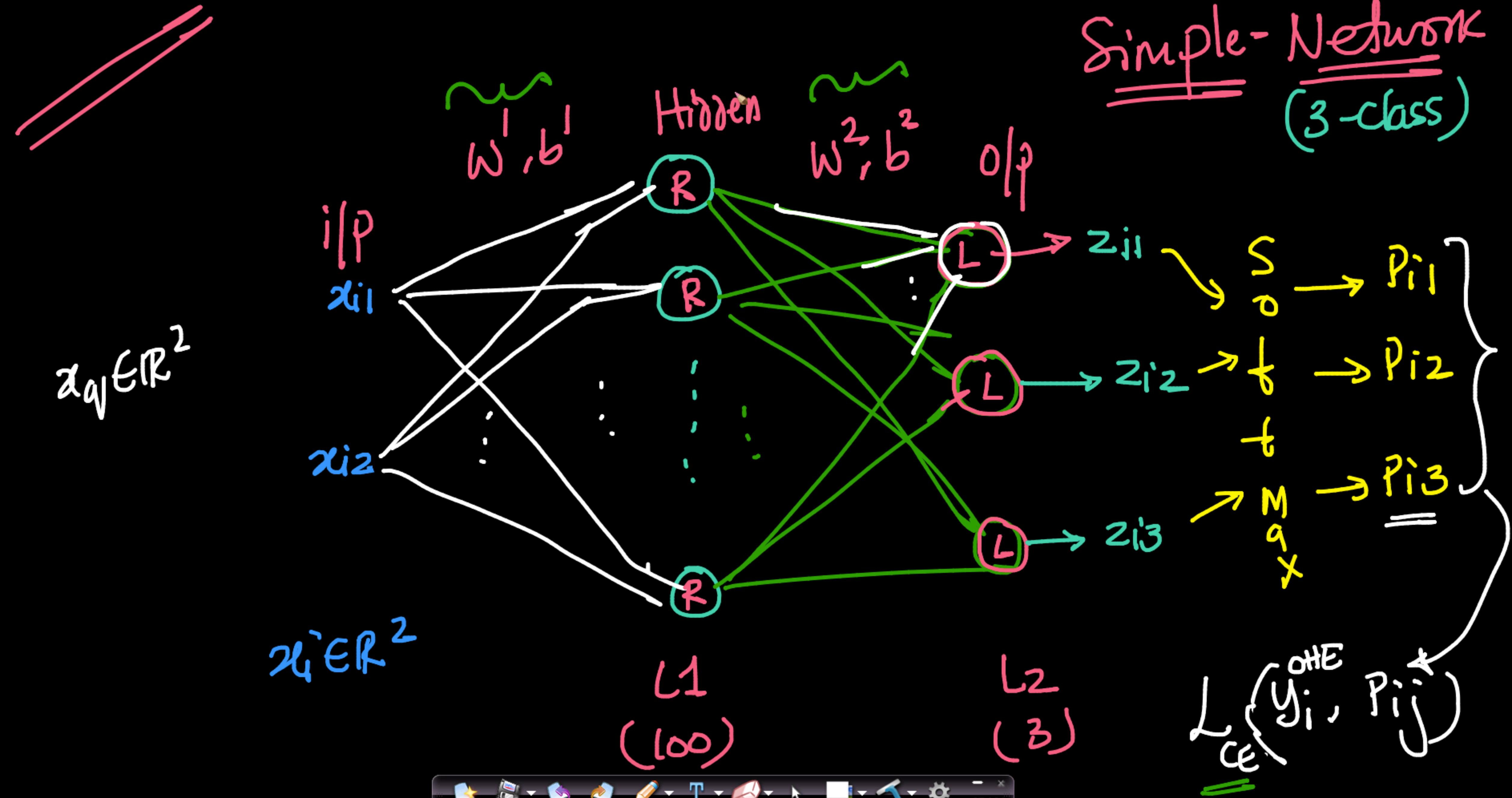
Binary  
Sigmoid  $\rightarrow$  k-ary  
Softmax  
(Pois  $\rightarrow$ )

loss  
Bin. CE /  
log-loss  
 $\overline{CE}$

~~Q~~

$\{ z_1 \rightarrow$  softmax as actvn-fn  $\{ p_1$   
 $z_2 \rightarrow p_2$   
 $z_3 \rightarrow p_3 \}$





colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=y9Kkza74VRoH

+ Code + Text

-rw-r--r-- 1 root root 12867 Sep 5 16:14 spiral.csv

RAM Disk

# Very simple case of a back-prop

```
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
df = pd.read_csv("spiral.csv")
plt.scatter(df["x1"], df["x2"], c=df["y"], s=40, cmap=plt.cm.Spectral)
plt.show()
x = df.iloc[:, :-1].to_numpy()
y = df.iloc[:, -1].to_numpy()
```

100  
0.75  
0.50  
0.25  
0.00  
-0.25  
-0.50  
-0.75

-1.00 -0.75 -0.50 -0.25 0.00 0.25 0.50 0.75 1.00

28 / 28

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=y9Kkza74VRoH

+ Code + Text

-rw-r--r-- 1 root root 12867 Sep 5 16:14 spiral.csv

RAM Disk

# Very simple case of a back-prop

```
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
df = pd.read_csv("spiral.csv")
plt.scatter(df["x1"], df["x2"], c=df["y"], s=40, cmap=plt.cm.Spectral)
plt.show()
x = df.iloc[:, :-1].to_numpy()
y = df.iloc[:, -1].to_numpy()
```

100  
75  
50  
25  
0.00  
-25  
-50  
-75

-1.00 -0.75 -0.50 -0.25 0.00 0.25 0.50 0.75 1.00

29 / 29

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=y9Kkza74VRoH

+ Code + Text

-rw-r--r-- 1 root root 12867 Sep 5 16:14 spiral.csv

RAM Disk

RAM Disk

# Very simple case of a back-prop

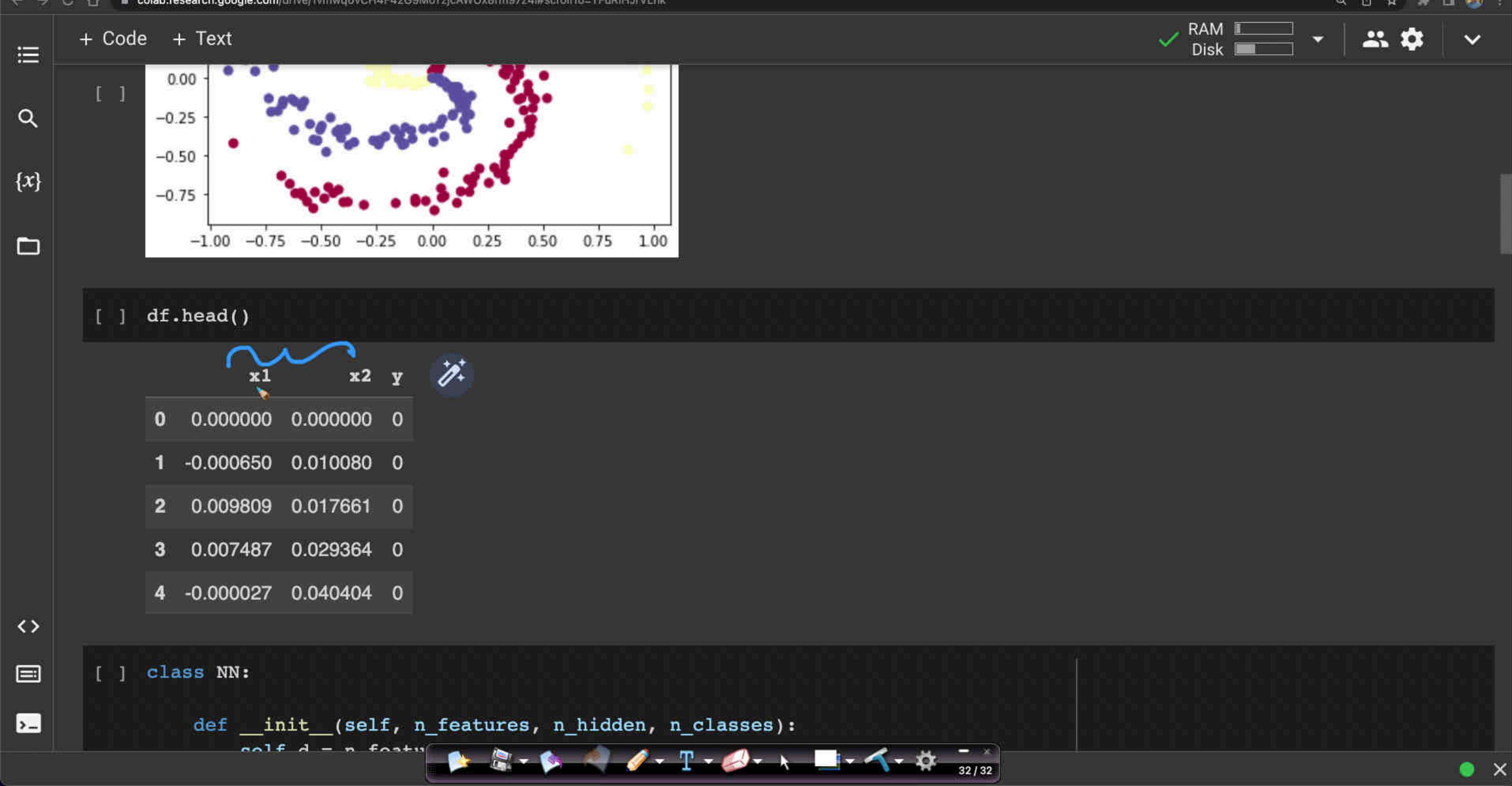
```
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
df = pd.read_csv("spiral.csv")
plt.scatter(df["x1"], df["x2"], c=df["y"], s=40, cmap=plt.cm.Spectral)
plt.show()
x = df.iloc[:, :-1].to_numpy()
y = df.iloc[:, -1].to_numpy()
```

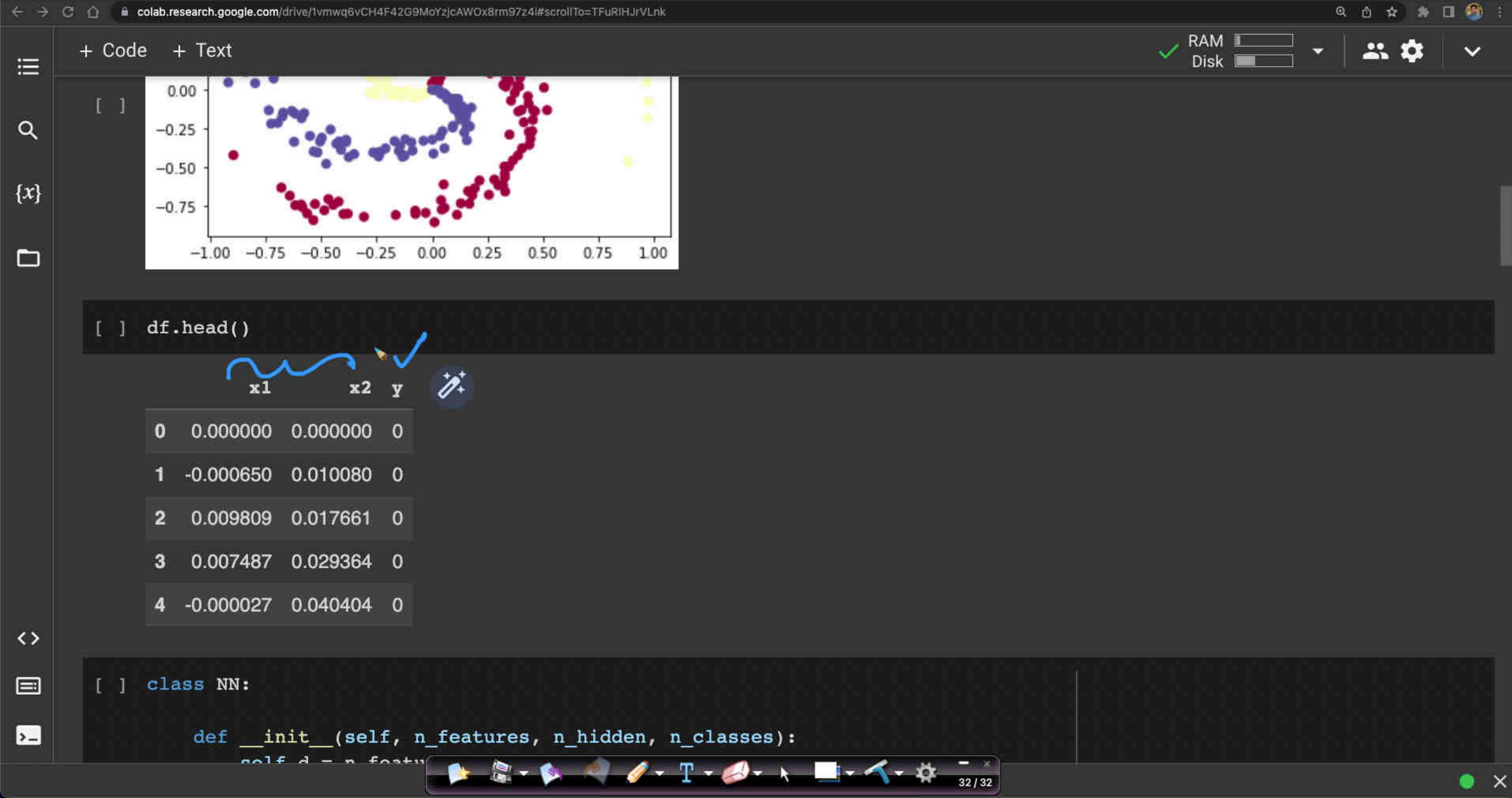
100  
0.75  
0.50  
0.25  
0.00  
-0.25  
-0.50  
-0.75

-1.00 -0.75 -0.50 -0.25 0.00 0.25 0.50 0.75 1.00

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100

30 / 30





colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=TFuRIHJrVLnk

+ Code + Text

RAM Disk

```
[ ] class NN:
    def __init__(self, n_features, n_hidden, n_classes):
        self.d = n_features
        self.h = n_hidden
        self.n = n_classes
        self.W1 = 0.01 * np.random.randn(self.d, self.h)
        self.b1 = np.zeros((1, self.h))
        self.W2 = 0.01 * np.random.randn(self.h, self.n)
        self.b2 = np.zeros((1, self.n))

    def fwd_prop(self, X):
        Z1 = np.dot(X, self.W1) + self.b1
        A1 = np.maximum(0, Z1)
        Z2 = np.dot(A1, self.W2) + self.b2
        Z2 = np.exp(Z2)
        A2 = Z2 / np.sum(Z2, axis=1, keepdims=True)
        return A1, A2

    def cce_loss(self, y, probs):
        num_examples = y.shape[0]
        correct_logprobs = -np.log(probs[range(num_examples), y])
        loss = np.sum(correct_logprobs)/num_examples
        return loss
```

**Hidden**

i/p

$x_{i1}$

$x_{i2}$

$w_{11}, w_{12}, w_{13}, w_{14}$

$w_{21}, w_{22}, w_{23}, w_{24}$

$b_1$

$b_2$

$(100)$

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=TFuRIHJrVLnk

+ Code + Text RAM Disk

```
[ ] class NN:  
  
{x}     def __init__(self, n_features, n_hidden, n_classes):  
        self.d = n_features  
        self.h = n_hidden  
        self.n = n_classes  
        self.W1 = 0.01 * np.random.randn(self.d, self.h)  
        self.b1 = np.zeros((1, self.h))  
        self.W2 = 0.01 * np.random.randn(self.h, self.n)  
        self.b2 = np.zeros((1, self.n))  
  
    def fwd_prop(self, X):  
        Z1 = np.dot(X, self.W1) + self.b1  
        A1 = np.maximum(0, Z1)  
        Z2 = np.dot(A1, self.W2) + self.b2  
        Z2 = np.exp(Z2)  
        A2 = Z2 / np.sum(Z2, axis=1, keepdims=True)  
        return A1, A2  
  
<>    def cce_loss(self, y, probs):  
        num_examples = y.shape[0]  
        correct_logprobs = -np.log(probs[range(num_examples), y])  
        loss = np.sum(correct_logprobs)/num_examples  
        return loss
```

+ Code + Text

self.n = n\_classes  
self.W1 = 0.01 \* np.random.randn(self.d, self.h)  
self.b1 = np.zeros((1, self.h))  
self.W2 = 0.01 \* np.random.randn(self.h, self.n)  
self.b2 = np.zeros((1, self.n))

def fwd\_prop(self, X):  
 Z1 = np.dot(X, self.W1) + self.b1  
 A1 = np.maximum(0, Z1)  
 Z2 = np.dot(A1, self.W2) + self.b2  
 Z2 = np.exp(Z2)  
 A2 = Z2 / np.sum(Z2, axis=1, keepdims=True)  
 return A1, A2

def cce\_loss(self, y, probs):  
 num\_examples = y.shape[0]  
 correct\_logprobs = -np.log(probs[range(num\_examples), y])  
 loss = np.sum(correct\_logprobs)/num\_examples  
 return loss

def back\_prop(self, X, A1, A2, y):  
 # compute the gradient on scores  
 num\_examples = y.shape[0]  
 dZ2 = A2  
 dZ2[range(num\_examples), y] -= 1  
 dZ2 /= num\_examples  
 # first backprop

RAM Disk

x 

35 / 37

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=uPrKn3MnjCeY

+ Code + Text

RAM Disk

class NN:

```

    def __init__(self, n_features, n_hidden, n_classes):
        self.d = n_features
        self.h = n_hidden
        self.n = n_classes
        self.W1 = 0.01 * np.random.randn(self.d, self.h)
        self.b1 = np.zeros((1, self.h))
        self.W2 = 0.01 * np.random.randn(self.h, self.n)
        self.b2 = np.zeros((1, self.n))

    def fwd_prop(self, X):
        Z1 = np.dot(X, self.W1) + self.b1
        A1 = np.maximum(0, Z1) → ReLU
        Z2 = np.dot(A1, self.W2) + self.b2
        Z2 = np.exp(Z2)
        A2 = Z2 / np.sum(Z2, axis=1, keepdims=True)
        return A1, A2

    def cce_loss(self, y, probs):
        num_examples = y.shape[0]
        correct_logprobs = -np.log(probs[range(num_examples), y])
        loss = np.sum(correct_logprobs)/num_examples
        return loss

    def back_prop(self, X, A1, A2, y):

```

$X_{mxd}$

$d=2$

$A_1$

$h=10$

$A_2$

$n=3$

$Z_2 = X \cdot W_1 + b_1$

RAM Disk

36 / 38

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=uPrKn3MnjCeY

+ Code + Text

RAM Disk ✓

Code Editor

class NN:

def \_\_init\_\_(self, n\_features, n\_hidden, n\_classes):

    self.d = n\_features

    self.h = n\_hidden

    self.n = n\_classes

    self.W1 = 0.01 \* np.random.randn(self.d, self.h)

    self.b1 = np.zeros((1, self.h))

    self.W2 = 0.01 \* np.random.randn(self.h, self.n)

    self.b2 = np.zeros((1, self.n))

def fwd\_prop(self, X):

    Z1 = np.dot(X, self.W1) + self.b1

    A1 = np.maximum(0, Z1)

    Z2 = np.dot(A1, self.W2) + self.b2

    Z2 = np.exp(Z2)

    A2 = Z2 / np.sum(Z2, axis=1, keepdims=True)

    return A1, A2

def cce\_loss(self, y, probs):

    num\_examples = y.shape[0]

    correct\_logprobs = -np.log(probs[range(num\_examples), y])

    loss = np.sum(correct\_logprobs)/num\_examples

    return loss

def back\_prop(self, X, A1, A2, y):

    # Implementation of back propagation

    # A1: Activation of hidden layer

    # A2: Activation of output layer

    # y: Ground truth labels

    # X: Input features

    # W1, b1, W2, b2: Weights and biases from forward pass

    # Calculate gradients for weights and biases

    # Update weights and biases using gradients

    # Return gradients for input features

    # Implementation of back propagation

    # A1: Activation of hidden layer

    # A2: Activation of output layer

    # y: Ground truth labels

    # X: Input features

    # W1, b1, W2, b2: Weights and biases from forward pass

    # Calculate gradients for weights and biases

    # Update weights and biases using gradients

    # Return gradients for input features

    # Implementation of back propagation

    # A1: Activation of hidden layer

    # A2: Activation of output layer

    # y: Ground truth labels

    # X: Input features

    # W1, b1, W2, b2: Weights and biases from forward pass

    # Calculate gradients for weights and biases

    # Update weights and biases using gradients

    # Return gradients for input features

+ Code + TextRAM Disk

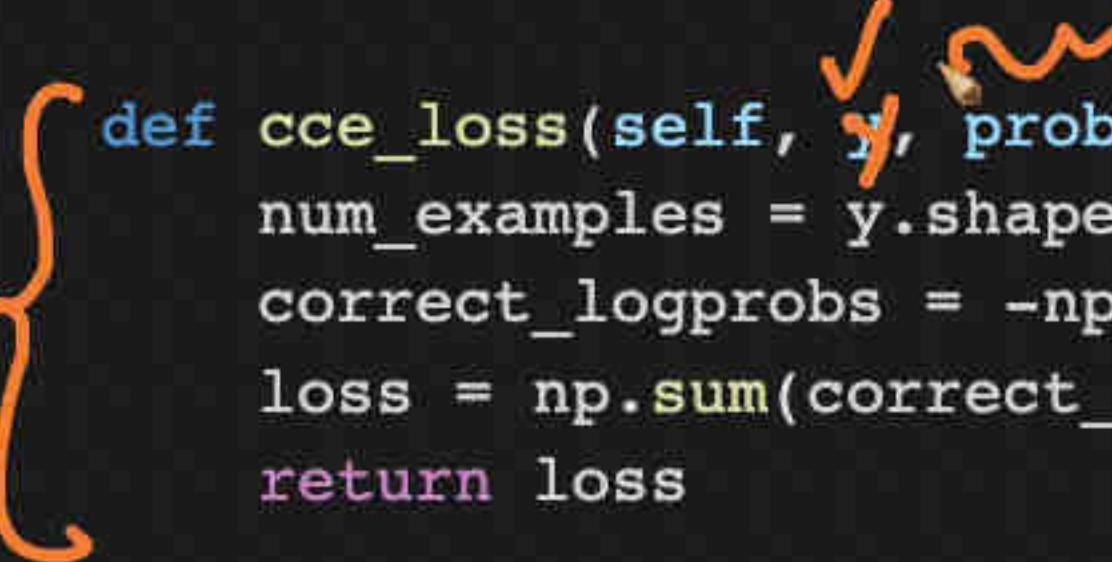
PlayUp Down Left Right Refresh Settings Copy Delete More

```
def fwd_prop(self, X):
    Z1 = np.dot(X, self.W1) + self.b1
    A1 = np.maximum(0, Z1)
    Z2 = np.dot(A1, self.W2) + self.b2
    Z2 = np.exp(Z2)
    A2 = Z2 / np.sum(Z2, axis=1, keepdims=True)
    return A1, A2
```

{  

```
def cce_loss(self, y, probs):
    num_examples = y.shape[0]
    correct_logprobs = -np.log(probs[range(num_examples),y])
    loss = np.sum(correct_logprobs)/num_examples
    return loss
```

```
def back_prop(self, X, A1, A2, y):
    # compute the gradient on scores
    num_examples = y.shape[0]
    dZ2 = A2
    dZ2[range(num_examples),y] -= 1
    dZ2 /= num_examples
    # first backprop into parameters W2 and b2
    dW2 = np.dot(A1.T, dZ2)
    db2 = np.sum(dZ2, axis=0, keepdims=True)
    # next backprop into hidden layer, A1
    dA1 = np.dot(dZ2, self.W2.T)
    # backprop the R
```



38 / 40

X

```
+ Code + Text
```



```
def fwd_prop(self, X):
    Z1 = np.dot(X, self.W1) + self.b1
    A1 = np.maximum(0, Z1)
    Z2 = np.dot(A1, self.W2) + self.b2
    Z2 = np.exp(Z2)
    A2 = Z2 / np.sum(Z2, axis=1, keepdims=True)
    return A1, A2

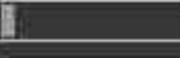
def cce_loss(self, y, probs):
    num_examples = y.shape[0]
    correct_logprobs = -np.log(probs[range(num_examples),y])
    loss = np.sum(correct_logprobs)/num_examples
    return loss

def back_prop(self, X, A1, A2, y):
    # compute the gradient on scores
    num_examples = y.shape[0]
    dZ2 = A2
    dZ2[range(num_examples),y] -= 1
    dZ2 /= num_examples
    # first backprop into parameters W2 and b2
    dW2 = np.dot(A1.T, dZ2)
    db2 = np.sum(dZ2, axis=0, keepdims=True)
    # next backprop into hidden layer, A1
    dA1 = np.dot(dZ2, self.W2.T)
    # backprop the R
```



{x}



✓ RAM  Disk 



colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=uPrKn3MnjCeY

+ Code + Text

RAM Disk

$\text{CE} = \frac{1}{m} \sum_{i=1}^m \text{CE}_i$

$\text{CE}_i = -\sum_{j=1}^K y_{ij} \log p_{ij}$

1

let  $y_i \in \text{class } 2$

$\begin{cases} y_{i2} = 1 \\ y_{ij} \neq 2 \end{cases}$

```

def fwd_prop(self, X):
    Z1 = np.dot(X, self.W1) + self.b1
    A1 = np.maximum(0, Z1)
    Z2 = np.dot(A1, self.W2) + self.b2
    Z2 = np.exp(Z2)
    A2 = Z2 / np.sum(Z2, axis=1, keepdims=True)
    return A1, A2

def cce_loss(self, y, probs):
    num_examples = y.shape[0]
    correct_logprobs = -np.log(probs[range(num_examples), y])
    loss = np.sum(correct_logprobs)/num_examples
    return loss

def back_prop(self, X, A1, A2, y):
    # compute the gradient on scores
    num_examples = y.shape[0]
    dZ2 = A2
    dZ2[range(num_examples), y] -= 1
    dZ2 /= num_examples
    # first backprop into parameters W2 and b2
    dW2 = np.dot(A1.T, dZ2)
    db2 = np.sum(dZ2, axis=0, keepdims=True)
    # next backprop into hidden layer, A1
    dA1 = np.dot(dZ2, self.W2.T)
    # backprop the R

```

+ Code + Text

✓ RAM Disk



```
def fwd_prop(self, X):
    Z1 = np.dot(X, self.W1) + self.b1
    A1 = np.maximum(0, Z1)
    Z2 = np.dot(A1, self.W2) + self.b2
    Z2 = np.exp(Z2)
    A2 = Z2 / np.sum(Z2, axis=1, keepdims=True)
    return A1, A2

def cce_loss(self, y, probs):
    num_examples = y.shape[0]
    correct_logprobs = -np.log(probs[range(num_examples),y])
    loss = np.sum(correct_logprobs)/num_examples
    return loss

def back_prop(self, X, A1, A2, y):
    # compute the gradient on scores
    num_examples = y.shape[0]
    dZ2 = A2
    dZ2[range(num_examples),y] -= 1
    dZ2 /= num_examples
    # first backprop into parameters W2 and b2
    dW2 = np.dot(A1.T, dZ2)
    db2 = np.sum(dZ2, axis=0, keepdims=True)
    # next backprop into hidden layer, A1
    dA1 = np.dot(dZ2, self.W2.T)
    # backprop the R
```

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=uPrKn3MnjCeY

+ Code + Text

RAM Disk

```

num_examples = y.shape[0]
correct_logprobs = -np.log(probs[range(num_examples),y])
loss = np.sum(correct_logprobs)/num_examples
return loss

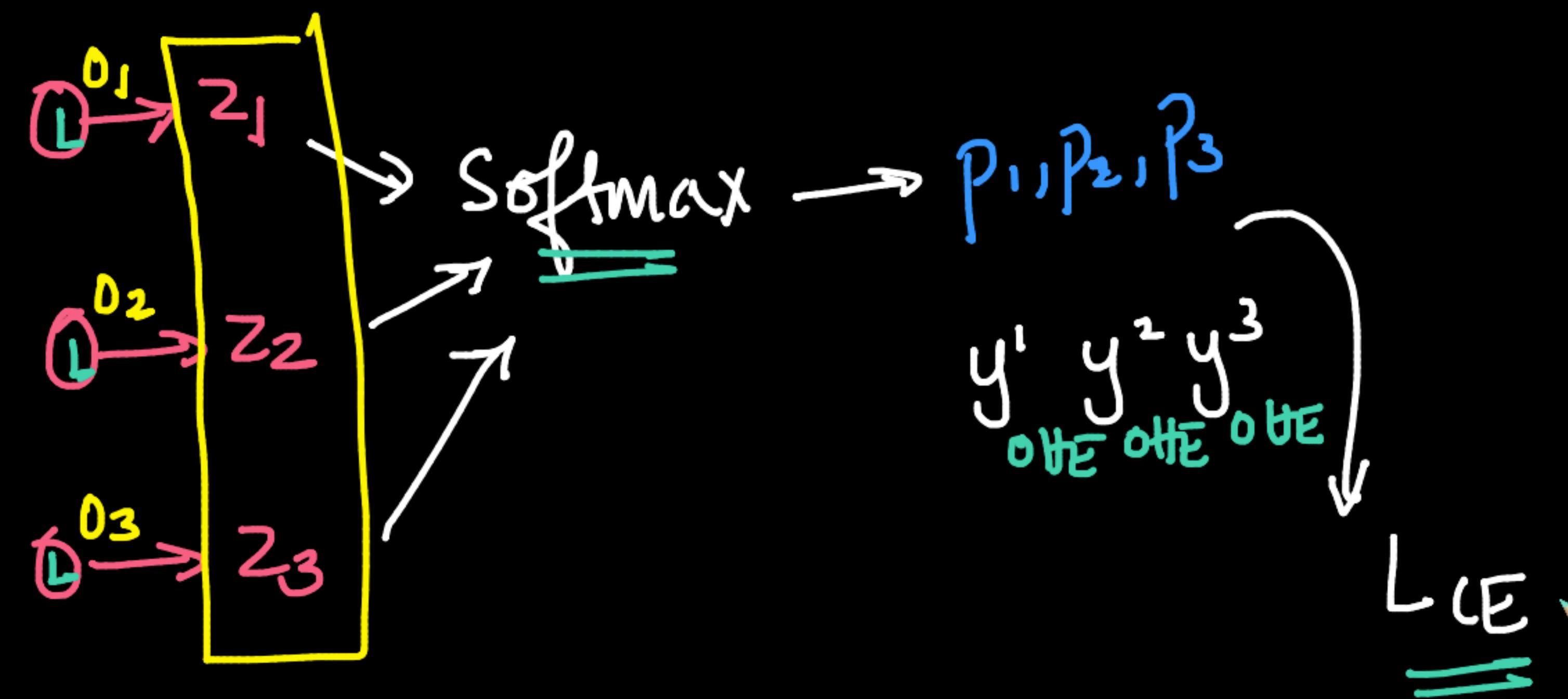
{x}

def back_prop(self, X, A1, A2, y):
    # compute the gradient on scores
    num_examples = y.shape[0]
    dZ2 = A2
    dZ2[range(num_examples),y] -= 1
    dZ2 /= num_examples
    # first backprop into parameters W2 and b2
    dW2 = np.dot(A1.T, dZ2)
    db2 = np.sum(dZ2, axis=0, keepdims=True)
    # next backprop into hidden layer, A1
    dA1 = np.dot(dZ2, self.W2.T)
    # backprop the ReLU non-linearity
    dA1[A1 <= 0] = 0
    # finally into w,b
    dZ1 = dA1
    dW1 = np.dot(X.T, dZ1)
    db1 = np.sum(dZ1, axis=0, keepdims=True)
    return dW1, db1, dW2, db2

def fit(self, X, lr, reg, max_iters):
    num_examples = X.shape[0]

```

The diagram illustrates a neural network architecture. It starts with an input vector  $X$  (represented by a yellow arrow) which is multiplied by weights  $W_1$  and bias  $b_1$  (also in yellow) to produce hidden layer activations  $Z_1$ .  $Z_1$  then passes through a ReLU activation function (indicated by the handwritten text "ReLU" above the arrow) to produce  $A_1$ .  $A_1$  is then multiplied by weights  $W_2$  and bias  $b_2$  (both in orange) to produce output layer activations  $Z_2$ .  $Z_2$  passes through a softmax function (indicated by the handwritten text "softmax" below the arrow) to produce  $A_2$ . Finally,  $A_2$  is compared against the target  $y$  to calculate the cross-entropy loss (indicated by the handwritten text "CE" at the bottom).



$$\frac{\partial L_{CE}}{\partial z_1}, \frac{\partial L_{CE}}{\partial z_2} = \frac{\partial L_{CE}}{\partial z_3}$$

3-class  
classfn's

$$\frac{\partial L_{CE}}{\partial z_2}$$

(e.g)

$x \rightarrow$  class 2

Model

$$\cancel{p_1}$$

$$\cancel{p_2}$$

$$\cancel{p_3}$$

$$\frac{\partial L_{CE}}{\partial z_j} = p_j - y_j^{OHE}$$

↳ looks that  $x \in$  class j

$y_1^{OHE}$	$y_2^{OHE}$	$y_3^{OHE}$
0	1	0
1	2	3

$$p_2 - 1 = \frac{\partial L_{CE}}{\partial z_2}$$



Thomas Kurpiel

212 Followers

Advanced Computer Vision & AI Research Engineer  
at APTIV Germany

Follow



## More from Medium

Daniel Reiff in Towards Data Science

**Generalizing Your Model: An Example With EfficientNetV2 and Cats & Dogs**

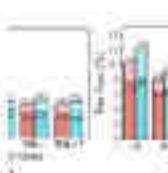
iva vrtaric in Artificialis

**Essential skill of Hypertuning with KerasTuner**

Shenyanyan

**A Machine Learning Template based on Pytorch Lightning**

Sik-Ho Tsang

**Review—There Are Many Consistent Explanations of Unlabeled Data: Why You Should...**

$$\frac{\partial z_i}{\partial z_j} = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases}$$

which can be concisely written using the indicator function  $1\{\cdot\}$ . The indicator function takes on a value of 1 if its argument is true, and 0 otherwise.

The second term on the right-hand side can be evaluated by applying the chain rule:

$$\frac{\partial}{\partial z_j} \log s_i = 1\{i = j\} \cdot \frac{1}{\sum_{l=1}^n e^{z_l}} \cdot \left( \frac{\partial}{\partial z_j} \sum_{l=1}^n e^{z_l} \right)$$

In the step above we used the derivative of the natural logarithm:

$$\frac{d}{dx} \log(x) = \frac{1}{x}$$

Obtaining the partial derivative of the sum is trivial:

$$\frac{\partial}{\partial z_j} \sum_{l=1}^n e^{z_l} = \frac{\partial}{\partial z_j} [e^{z_1} + e^{z_2} + \dots + e^{z_j} + \dots + e^{z_n}] = \frac{\partial}{\partial z_j} [e^{z_j}] = e^{z_j}$$

Respond

Plugging the result into the form:



[towardsdatascience.com/derivative-of-the-softmax-function-and-the-categorical-cross-entropy-loss-ffceefc081d1](#)

Get started Sign In

Search

Thomas Kurbiel  
212 Followers  
Advanced Computer Vision & AI Research Engineer at APTIV Germany

Follow

More from Medium

Daniel Reiff in Towards Data Science  
**Generalizing Your Model: An Example With EfficientNetV2 and Cats & Dogs**

iva vrtaric in Artificialis  
**Essential skill of Hypertuning with KerasTuner**

Shenyan  
**A Machine Learning Template based on Pytorch Lightning**

Sik-Ho Tsang  
**Review—There Are Many**

363 6

$\frac{\partial \mathcal{L}}{\partial z_j} = \sum_{i=1}^C y_i \cdot s_j - y_j$

Next, we pull  $s$  out of the sum, since it does not depend on index  $i$ :

$$\frac{\partial \mathcal{L}}{\partial z_j} = s_j \cdot \sum_{i=1}^C y_i - y_j = s_j - y_j$$

$\frac{\partial \mathcal{L}}{\partial z} = s - y$

More readings

[Eli Bendersky's website](#)

[Raúl Gómez blog](#)

47 / 47

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=uPrKn3MnjCeY

+ Code + Text

```

def back_prop(self, X, A1, A2, y):
    # compute the gradient on scores
    num_examples = y.shape[0]
    dZ2 = A2
    dZ2[range(num_examples), y] -= 1
    dZ2 /= num_examples
    # first backprop into parameters w2 and b2
    dW2 = np.dot(A1.T, dZ2)
    db2 = np.sum(dZ2, axis=0, keepdims=True)
    # next backprop into hidden layer, A1
    dA1 = np.dot(dZ2, self.W2.T)
    # backprop the ReLU non-linearity
    dA1[A1 <= 0] = 0
    # finally into w,b
    dZ1 = dA1
    dW1 = np.dot(X.T, dZ1)
    db1 = np.sum(dZ1, axis=0, keepdims=True)
    return dW1, db1, dW2, db2

def fit(self, X, lr, reg, max_iters):
    num_examples = X.shape[0]
    for i in range(max_iters):
        # forward prop
        A1, A2 = self.fwd_prop(X)
        # calculate loss

```

$\frac{\partial L}{\partial A_1} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial A_1}$

$X \xrightarrow{w_1, b_1} Z_1 \xrightarrow{\text{ReLU}} A_1 \xrightarrow{w_2, b_2} Z_2 \xrightarrow{\text{Softmax}} A_2 \xrightarrow{\text{prob}} L_{CE}$

$\frac{\partial z_2}{\partial w_2} = A_1$

$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2}$

RAM Disk

48 / 49

[colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=uPrKn3MnjCeY](#)

+ Code + Text  



```
def back_prop(X, A1, A2, y):
    dw1 = np.sum(dZ1, axis=0, keepdims=True)
    return dw1, db1, dw2, db2

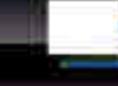
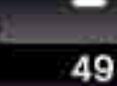
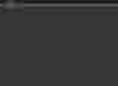
def fit(self, X, lr, reg, max_iters):
    num_examples = X.shape[0]
    for i in range(max_iters):
        #forward prop
        A1, A2 = self.fwd_prop(X)
        # calculate loss
        data_loss = self.cce_loss(y, A2)
        reg_loss = 0.5*reg*np.sum(self.W1*self.W1) + 0.5*reg*np.sum(self.W2*self.W2)
        loss = data_loss + reg_loss

        if i % 1000 == 0:
            print("iteration %d: loss %f" % (i, loss))

        dw1, db1, dw2, db2 = self.back_prop(X, A1, A2, y)

        # add regularization gradient contribution
        dw2 += reg * self.W2
        dw1 += reg * self.W1

        # perform a parameter update
        self.W1 += -lr * dw1
        self.b1 += -lr * db1
        self.W2 += -lr * dw2
```

49 / 50

```
+ Code + Text
```

dW1 = np.dot(X.T, dZ1)  
db1 = np.sum(dZ1, axis=0, keepdims=True)  
return dW1, db1, dW2, db2

{x}

def fit(self, X, lr, reg, max\_iters):  
 num\_examples = X.shape[0]  
 ✓ for i in range(max\_iters): → 10,000  
 ✓ { # forward prop  
 A1, A2 = self.fwd\_prop(X)  
 # calculate loss  
 data\_loss = self.cce\_loss(y, A2)  
 reg\_loss = 0.5\*reg\*np.sum(self.W1\*self.W1) + 0.5\*reg\*np.sum(self.W2\*self.W2)  
 loss = data\_loss + reg\_loss  
 if i % 1000 == 0:  
 print("iteration %d: loss %f" % (i, loss))  
  
 dW1, db1, dW2, db2 = self.back\_prop(X, A1, A2, y)  
  
 # add regularization gradient contribution  
 dW2 += reg \* self.W2  
 dW1 += reg \* self.W1  
  
 # perform a parameter update  
 self.W1 += -lr \* dW1  
 self.b1 += -lr \* db1

<>

50 / 51

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=uPrKn3MnjCeY

+ Code + Text

RAM Disk

$L_2 \cdot \text{reg}$

$\frac{1}{2} \omega^T \omega$

$\frac{1}{2} \omega \cdot \omega$

```

dW1 = np.dot(X.T, dZ1)
db1 = np.sum(dZ1, axis=0, keepdims=True)
return dw1, db1, dw2, db2

{x}
def fit(self, X, lr, reg, max_iters):
    num_examples = X.shape[0]
    for i in range(max_iters):
        # forward prop
        A1, A2 = self.fwd_prop(X)
        # calculate loss
        data_loss = self.cce_loss(y, A2)
        reg_loss = 0.5 * reg * np.sum(self.W1 * self.W1) + 0.5 * reg * np.sum(self.W2 * self.W2)
        loss = data_loss + reg_loss
        if i % 1000 == 0:
            print("iteration %d: loss %f" % (i, loss))

        dw1, db1, dw2, db2 = self.back_prop(X, A1, A2, y)

        # add regularization gradient contribution
        dw2 += reg * self.W2
        dw1 += reg * self.W1

        # perform a parameter update
        self.W1 += -lr * dw1
        self.b1 += -lr * db1
    
```

51 / 52

```
dW1 = np.dot(X.T, dZ1)
db1 = np.sum(dZ1, axis=0, keepdims=True)
return dW1, db1, dW2, db2

def fit(self, X, lr, reg, max_iters):
    num_examples = X.shape[0]
    for i in range(max_iters):
        # forward prop
        A1, A2 = self.fwd_prop(X)
        # calculate loss
        data_loss = self.cce_loss(y, A2)
        reg_loss = 0.5*reg*np.sum(self.W1*self.W1) + 0.5*reg*np.sum(self.W2*self.W2)
        loss = data_loss + reg_loss
        if i % 1000 == 0:
            print("iteration %d: loss %f" % (i, loss))
        dW1, db1, dW2, db2 = self.back_prop(X, A1, A2, y)
        # add regularization gradient contribution
        dW2 += reg * self.W2
        dW1 += reg * self.W1
        # perform a parameter update
        self.W1 += -lr * dW1
        self.b1 += -lr * db1
```

```
+ Code + Text
```

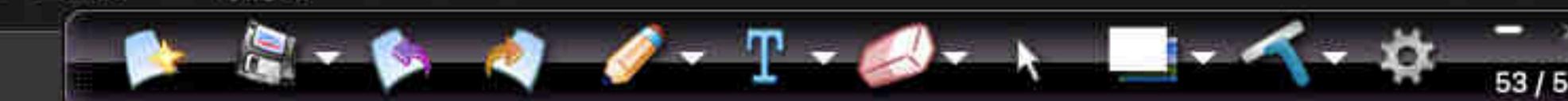
```
dW1 = np.dot(X.T, dZ1)
db1 = np.sum(dZ1, axis=0, keepdims=True)
return dw1, db1, dw2, db2

{x}
def fit(self, X, lr, reg, max_iters):
    num_examples = X.shape[0]
    for i in range(max_iters):
        # forward prop
        A1, A2 = self.fwd_prop(X)
        # calculate loss
        data_loss = self.cce_loss(y, A2)
        reg_loss = 0.5*reg*np.sum(self.W1*self.W1) + 0.5*reg*np.sum(self.W2*self.W2)
        loss = data_loss + reg_loss
        if i % 1000 == 0:
            print("iteration %d: loss %f" % (i, loss))

        dw1, db1, dw2, db2 = self.back_prop(X, A1, A2, y)

        # add regularization gradient contribution
        dw2 += reg * self.W2
        dw1 += reg * self.W1

        # perform a parameter update
        self.W1 += -lr * dw1
        self.b1 += -lr * db1
```



+ Code + Text

```
dW1 = np.dot(X.T, dZ1)
db1 = np.sum(dZ1, axis=0, keepdims=True)
return dW1, db1, dW2, db2

def fit(self, X, lr, reg, max_iters):
    num_examples = X.shape[0]
    for i in range(max_iters):
        #forward prop
        A1, A2 = self.fwd_prop(X)
        # calculate loss
        data_loss = self.cce_loss(y, A2)
        reg_loss = 0.5*reg*np.sum(self.W1*self.W1) + 0.5*reg*np.sum(self.W2*self.W2)
        loss = data_loss + reg_loss

        if i % 1000 == 0:
            print("iteration %d: loss %f" % (i, loss))

        dW1, db1, dW2, db2 = self.back_prop(X, A1, A2, y)

        # add regularization gradient contribution
        dW2 += reg * self.W2
        dW1 += reg * self.W1

        # perform a parameter update
        self.W1 += -lr * dW1
        self.b1 += -lr * db1
```

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=uPrKn3MnjCeY

+ Code + Text

RAM Disk

iteration 0: loss 0.0

```

loss = data_loss + reg_loss

if i % 1000 == 0:
    print("iteration %d: loss %f" % (i, loss))

    dw1, db1, dw2, db2 = self.back_prop(X, A1, A2, y)

    # add regularization gradient contribution
    dw2 += reg * self.W2
    dw1 += reg * self.W1

    # perform a parameter update
    self.W1 += -lr * dw1
    self.b1 += -lr * db1
    self.W2 += -lr * dw2
    self.b2 += -lr * db2

def predict(self, X):
    A1 = np.maximum(0, np.dot(X, self.W1) + self.b1)
    Z2 = np.dot(A1, self.W2) + self.b2
    y_hat = np.argmax(Z2, axis=1)
    return y_hat

nn_model = NN(n_features=2, n_hidden=100, n_classes=3)
nn_model.fit(X, lr=1, reg=1e-3, max_iters=10000)
print('training accuracy: %.2f' % (np.mean(nn_model.predict(X) == y)))

```

$\frac{\partial L}{\partial w_2}$

loss

$X \cdot w_2^T$

$+ \lambda \frac{1}{2} \cdot w_2^T w_2$

$\lambda \cdot \underline{w_2}$

RAM Disk

55 / 56

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=uPrKn3MnjCeY

+ Code + Text

RAM Disk

# finally into W,b  
dZ1 = dA1  
dW1 = np.dot(X.T, dZ1)  
db1 = np.sum(dZ1, axis=0, keepdims=True)  
return dW1, db1, dW2, db2

def fit(self, X, lr, reg, max\_iters):  
 num\_examples = X.shape[0]  
 for i in range(max\_iters):  
 #forward prop  
 A1, A2 = self.fwd\_prop(X)  
 # calculate loss  
 data\_loss = self.cce\_loss(y, A2)  
 reg\_loss = 0.5\*reg\*np.sum(self.W1\*self.W1) + 0.5\*reg\*np.sum(self.W2\*self.W2)  
 loss = data\_loss + reg\_loss  
 if i % 1000 == 0:  
 print("iteration %d: loss %f" % (i, loss))  
 dw1, db1, dw2, db2 = self.back\_prop(X, A1, A2, y)  
 # add regularization gradient contribution  
 dw2 += reg \* self.W2  
 dw1 += reg \* self.W1

# perform a

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=uPrKn3MnjCeY

+ Code + Text

RAM Disk

# add regularization gradient contribution  
dw2 += reg \* self.W2  
dw1 += reg \* self.W1

{x}

# perform a parameter update  
self.W1 += -lr \* dw1  
self.b1 += -lr \* db1  
self.W2 += -lr \* dw2  
self.b2 += -lr \* db2

def predict(self, X):  
 A1 = np.maximum(0, np.dot(X, self.W1) + self.b1)  
 Z2 = np.dot(A1, self.W2) + self.b2  
 y\_hat = np.argmax(Z2, axis=1)  
 return y\_hat

nn\_model = NN(n\_features=2, n\_hidden=100, n\_classes=3)  
nn\_model.fit(X, lr=1, reg=1e-3, max\_iters=10000)  
print('training accuracy: %.2f' % (np.mean(nn\_model.predict(X) == y)))

iteration 0: loss 1.098614  
iteration 1000: loss 0.305437  
iteration 2000: loss 0.260879  
iteration 3000: loss 0.255363  
iteration 4000: loss 0.254227  
iteration 5000: loss 0.253527

57 / 58

+ Code + Text

✓ RAM Disk

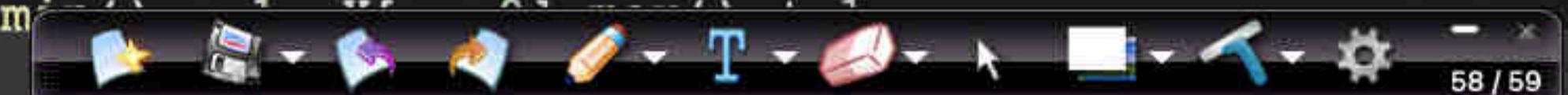
```
self.b2 += -lr * db2

def predict(self, X):
    A1 = np.maximum(0, np.dot(X, self.W1) + self.b1)
    Z2 = np.dot(A1, self.W2) + self.b2
    y_hat = np.argmax(Z2, axis=1)
    return y_hat

nn_model = NN(n_features=2, n_hidden=100, n_classes=3)
nn_model.fit(X, lr=1, reg=1e-3, max_iters=10000)
print('training accuracy: %.2f' % (np.mean(nn_model.predict(X) == y)))
```

```
iteration 0: loss 1.098614
iteration 1000: loss 0.305437
iteration 2000: loss 0.260879
iteration 3000: loss 0.255363
iteration 4000: loss 0.254227
iteration 5000: loss 0.253764
iteration 6000: loss 0.264896
iteration 7000: loss 0.253992
iteration 8000: loss 0.253224
iteration 9000: loss 0.251573
training accuracy: 0.96
```

```
[ ] # create a 2D grid
step = 0.02
x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
```



# Code # Text

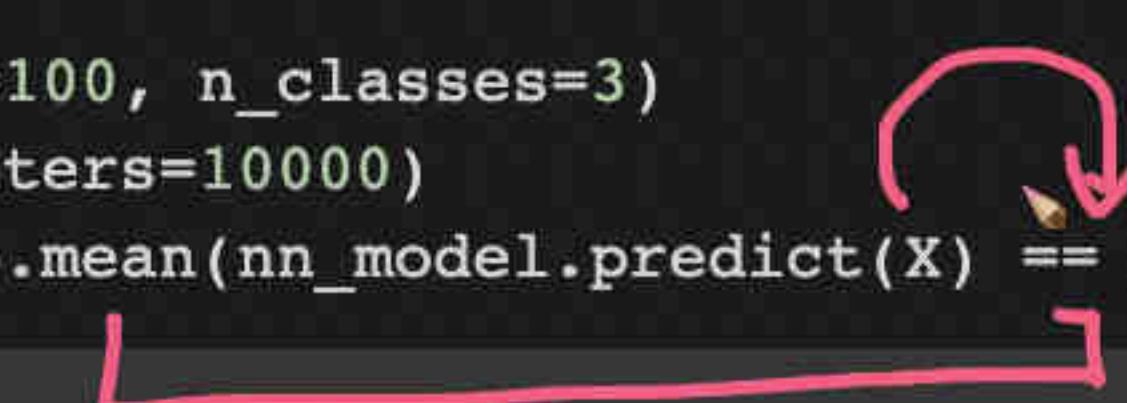
✓ RAM Disk



```
self.b2 += -lr * db2

def predict(self, X):
    A1 = np.maximum(0, np.dot(X, self.W1) + self.b1)
    Z2 = np.dot(A1, self.W2) + self.b2
    y_hat = np.argmax(Z2, axis=1)
    return y_hat

model = NN(n_features=2, n_hidden=100, n_classes=3)
model.fit(X, lr=1, reg=1e-3, max_iters=10000)
print('training accuracy: %.2f' % (np.mean(nn_model.predict(X) == y)))
```



```
iteration 0: loss 1.098614
iteration 1000: loss 0.305437
iteration 2000: loss 0.260879
iteration 3000: loss 0.255363
iteration 4000: loss 0.254227
iteration 5000: loss 0.253764
iteration 6000: loss 0.264896
iteration 7000: loss 0.253992
iteration 8000: loss 0.253224
iteration 9000: loss 0.251573
training accuracy: 0.96
```

```
[ ] # create a 2D grid  
step = 0.02  
x_min, x_max = X[:,
```

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=uPrKn3MnjCeY

+ Code + Text

self.b2 += -lr \* db2

def predict(self, X):  
 A1 = np.maximum(0, np.dot(X, self.W1) + self.b1)  
 Z2 = np.dot(A1, self.W2) + self.b2  
 y\_hat = np.argmax(Z2, axis=1)  
 return y\_hat

nn\_model = NN(n\_features=2, n\_hidden=100, n\_classes=3)  
nn\_model.fit(X, lr=1, reg=1e-3, max\_iters=10000)  
print('training accuracy: %.2f' % (np.mean(nn\_model.predict(X) == y)))

iteration 0: loss 1.098614  
iteration 1000: loss 0.305437  
iteration 2000: loss 0.260879  
iteration 3000: loss 0.255363  
iteration 4000: loss 0.254227  
iteration 5000: loss 0.253764  
iteration 6000: loss 0.264896  
iteration 7000: loss 0.253992  
iteration 8000: loss 0.253224  
iteration 9000: loss 0.251573  
training accuracy: 0.96

LCE

low loss

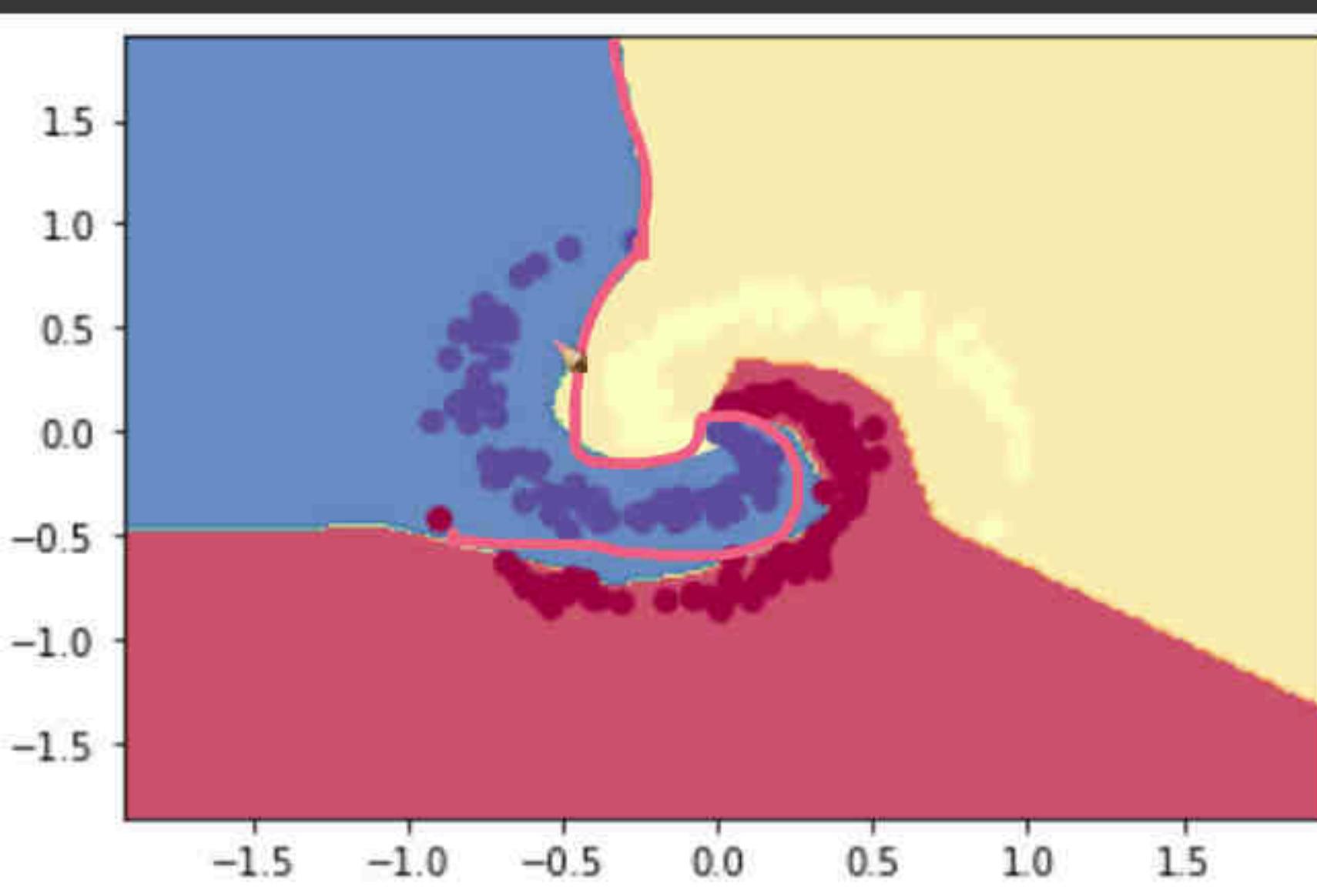
iter

[ ] # create a 2D grid  
step = 0.02  
x\_min, x\_max = X[:, 0].min() - 1, X[:, 0].max() + 1  
y\_min, y\_max = X[:, 1].min() - 1, X[:, 1].max() + 1

+ Code + Text

```
[ ] y_hat = nn_model.predict(np.c_[xx.ravel(), yy.ravel()])
[ ] y_hat = y_hat.reshape(xx.shape)

# plot
fig = plt.figure()
plt.contourf(xx, yy, y_hat, cmap=plt.cm.Spectral, alpha=0.8)
plt.scatter(X[:, 0], X[:, 1], c=y, s=40, cmap=plt.cm.Spectral)
plt.xlim(xx.min(), xx.max())
plt.ylim(yy.min(), yy.max())
plt.show()
```



let  
 $f(x) = x^2$

$Q$

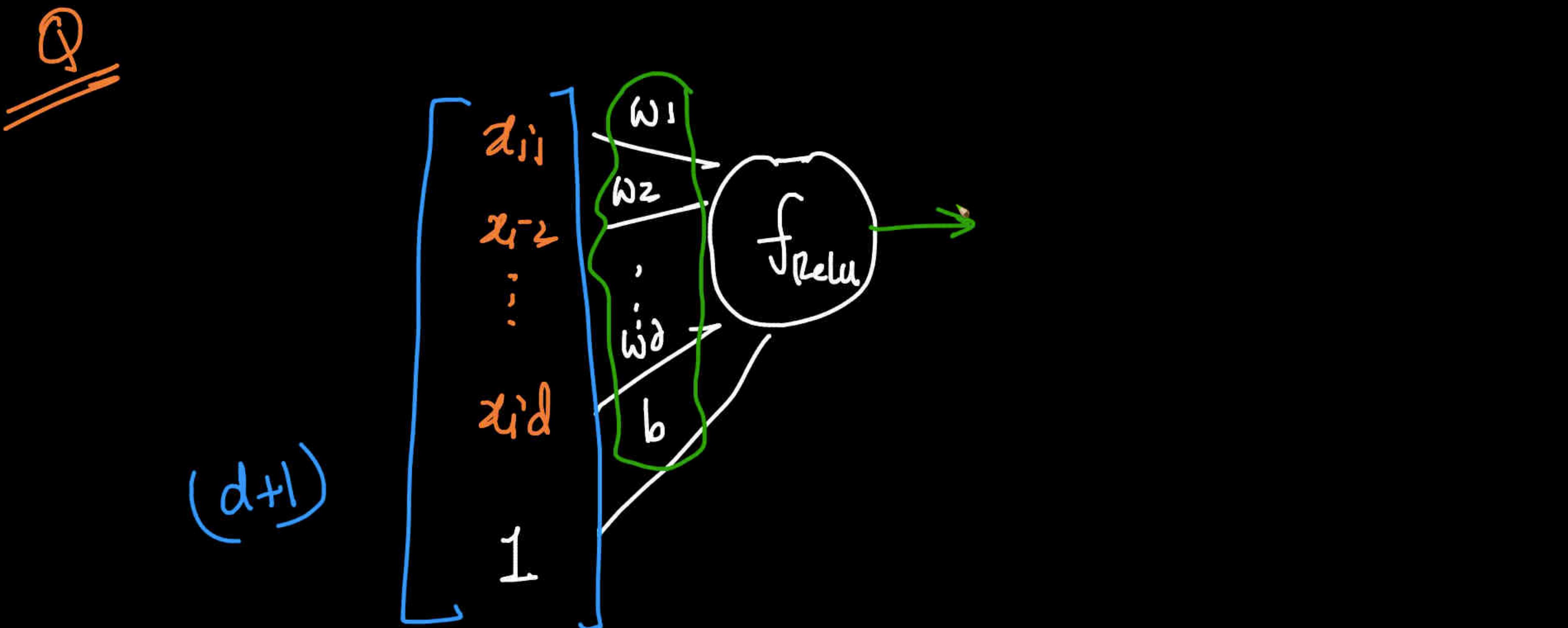
$$z_1 = 2x + 3$$

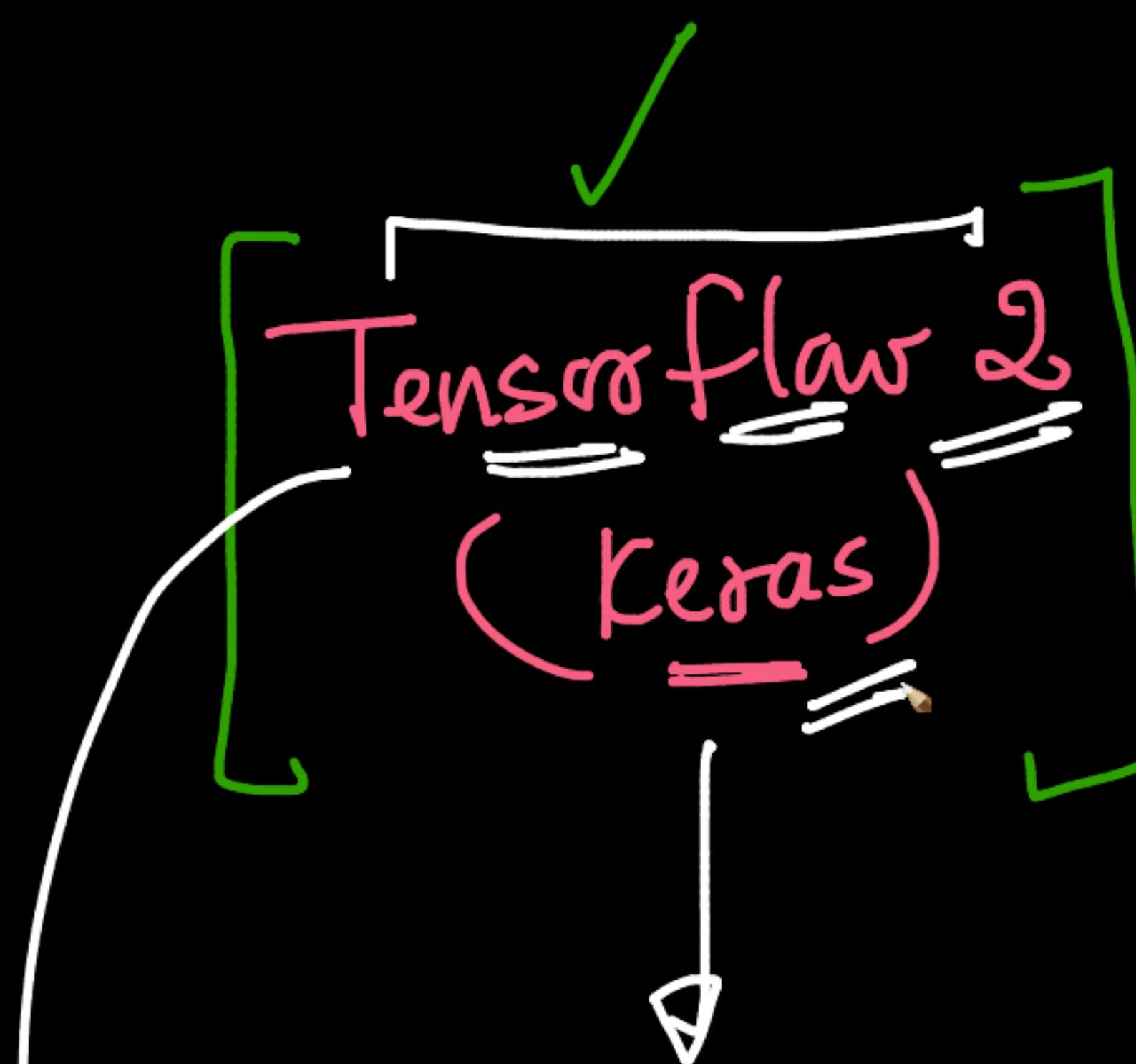
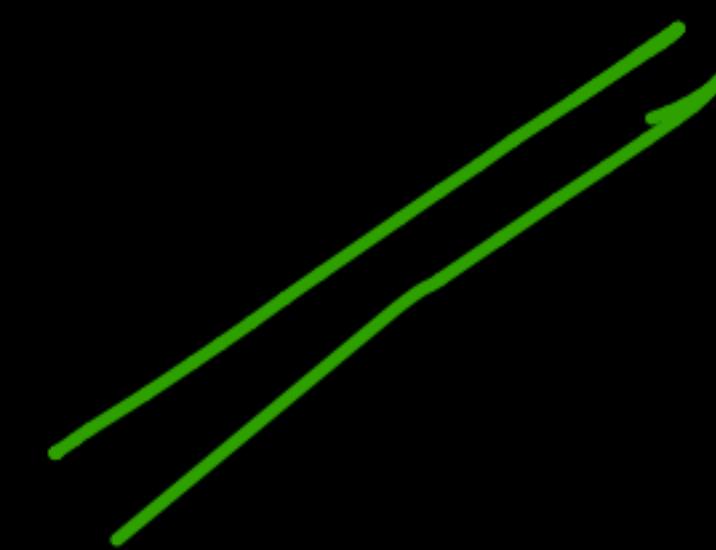
$$x \xrightarrow{w_1, b_1} z_1 = \underline{\underline{w_1^T x + b_1}} \rightarrow f^0 \rightarrow \underline{\underline{z_2}}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_2} \cdot \boxed{\frac{\partial z_2}{\partial w_1}}$$

Loss

$\hookrightarrow \frac{\partial (2x+3)}{\partial x} =$





A hand-drawn diagram on a black background. It consists of a white rectangular box with a green border. Inside the box, the words "Py-Torch" are written in pink, with a green bracket underneath. A green arrow points from the bottom of the box towards the text "→ Keras: Very simple" below it.

→ Keras: Very simple

→ TF2

Desmos | Graphing x BackpropFromScratch x K Why choose Keras x K Getting started x K Introduction to Keras x K Introduction to Keras x srivastava14a.pdf x Calculus on Comp x Softmax function x Derivative of the Sigmoid x +

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=YMVEIVUyWGrt

+ Code + Text RAM Disk

## TF 2 and Keras: Sequential API

```
import tensorflow as tf
from tensorflow import keras
from tensorflow.keras import layers

print(tf.__version__)

2.8.2
```

```
[ ] #Keras Sequential API
model = keras.Sequential()
model.add(keras.Input(shape=(2,)))
model.add(layers.Dense(100, activation="relu"))
model.add(layers.Dense(3, activation="softmax"))

[ ] model.summary()

Model: "sequential_1"
```

Desmos | Graphing x BackpropFromScratch x K Why choose Keras x K Getting started x K Introduction to Keras x K Introduction to Keras x srivastava14a.pdf x Calculus on Comp x Softmax function x Derivative of the Sigmoid x +

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=YMVEIVUyWGrt

+ Code + Text RAM Disk

TF 2 and Keras: Sequential API

{x}

import tensorflow as tf  
from tensorflow import keras  
from tensorflow.keras import layers

print(tf.\_\_version\_\_)

2.8.2

[ ] #Keras Sequential API  
model = keras.Sequential()  
model.add(keras.Input(shape=(2,)))  
model.add(layers.Dense(100, activation="relu"))  
model.add(layers.Dense(3, activation="softmax"))

[ ] model.summary()

Model: "sequential\_1"

Page Number: 66 / 66

Desmos | Graphing x BackpropFromScratch x K Why choose Keras x K Getting started x K Introduction to Keras x K Introduction to Keras x srivastava14a.pdf x Calculus on Comp x Softmax function x Derivative of the Sigmoid x +

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=YMVEIVUyWGrt

+ Code + Text RAM Disk

TF 2 and Keras: Sequential API

{x}

input → activ → o/p loss

```
import tensorflow as tf
from tensorflow import keras
from tensorflow.keras import layers

print(tf.__version__)

2.8.2
```

```
[ ] #Keras Sequential API
model = keras.Sequential()
model.add(keras.Input(shape=(2,)))
model.add(layers.Dense(100, activation="relu"))
model.add(layers.Dense(3, activation="softmax"))

[ ] model.summary()

Model: "sequential_1"
```

67 / 67

+ Code + Text

✓ RAM Disk

-1.5 -1.0 -0.5 0.0 0.5 1.0 1.5



## {x} ▾ TF 2 and Keras: Sequential API

```
[ ] import tensorflow as tf  
from tensorflow import keras  
from tensorflow.keras import layers
```

$x_i$  Dense

```
print(tf.__version__)
```

2.8.2

```
[ ] #Keras Sequential API  
model = keras.Sequential()  
model.add(keras.Input(shape=(2,)))  
model.add(layers.Dense(100, activation="relu"))  
model.add(layers.Dense(3, activation="softmax"))
```

```
[ ] model.summary()
```

Model: "sequential\_1"



Desmos | Graphing x BackpropFromScratch x K Why choose Keras x K Getting started x K Introduction to Keras x K Introduction to Keras x srivastava14a.pdf x Calculus on Comp x Softmax function x Derivative of the Sigmoid x +

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=ZBvZDZajWHHk

+ Code + Text

RAM Disk

Code:

```
from tensorflow.keras import layers
[ ]
print(tf.__version__)
{x}
2.8.2
```

[ ] #Keras Sequential API

```
model = keras.Sequential()
model.add(keras.Input(shape=(2,)))
model.add(layers.Dense(100, activation="relu"))
model.add(layers.Dense(3, activation="softmax"))
```

[ ] model.summary()

Model: "sequential\_1"

Layer (type)	Output Shape	Param #
dense_2 (Dense)	(None, 100)	300
dense_3 (Dense)	(None, 3)	303

Total params: 603  
Trainable params: 603

File Edit View Insert Cell Kernel Help

99 / 69

+ Code + Text

```
from tensorflow.keras import layers
[ ]
```

print(tf.\_\_version\_\_)

2.8.2

```
[ ] #Keras Sequential API
model = keras.Sequential()
model.add(keras.Input(shape=(2,)))
model.add(layers.Dense(100, activation="relu"))
model.add(layers.Dense(3, activation="softmax"))
```

[ ] model.summary()

Model: "sequential\_1"

Layer (type)	Output Shape	Param #
dense_2 (Dense)	(None, 100)	300
dense_3 (Dense)	(None, 3)	303

Total params: 603  
Trainable params: 603

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=ZBvZDZajWHHk

+ Code + Text

RAM Disk

Code:

```
from tensorflow.keras import layers
[ ]
print(tf.__version__)
{x}
2.8.2
```

Model:

```
[ ] #Keras Sequential API
model = keras.Sequential()
model.add(keras.Input(shape=(2,)))
model.add(layers.Dense(100, activation="relu"))
model.add(layers.Dense(3, activation="softmax"))
[ ] model.summary()
```

Diagram:

Model: "sequential\_1"

Layer (type)	Output Shape	Param #
dense_2 (Dense)	(None, 100)	300
dense_3 (Dense)	(None, 3)	303

Total params: 603  
Trainable params: 603

Toolbar:

Page Number: 71/72

Desmos | Graphing x BackpropFromScratch x K Why choose Keras x K Getting started x K Introduction to Keras x K Introduction to Keras x srivastava14a.pdf x Calculus on Comp x Softmax function x Derivative of the Sigmoid x + colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=RI24jRUzgdGz

+ Code + Text RAM Disk

```
[ ] #Keras Sequential API
model = keras.Sequential()
model.add(keras.Input(shape=(2,)))
model.add(layers.Dense(100, activation="relu"))
model.add(layers.Dense(3, activation="softmax"))

model.summary()
```

Model: "sequential\_1"

Layer (type)	Output Shape	Param #
dense_2 (Dense)	(None, 100)	300
<b>dense_3 (Dense)</b>	(None, 3)	303
Total params: 603		
Trainable params: 603		
Non-trainable params: 0		

```
[ ] model.compile(optimizer=keras.optimizers.Adam(), # suggest in colab : Option+Esc in Mac
loss=keras.losses.CategoricalCrossentropy(),
```

*Relu*

ReLU

$(2, 2) \xrightarrow{\text{FC}} \text{dense} \xrightarrow{\text{FC}} (3)$

100 neurons

3

Desmos | Graphing x BackpropFromScratch x K Why choose Keras x K Getting started x K Introduction to Keras x K Introduction to Keras x srivastava14a.pdf x Calculus on Comp x Softmax function x Derivative of the Sigmoid x +

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=RI24jRUzgdGz

+ Code + Text RAM Disk

```
[ ] #Keras Sequential API
model = keras.Sequential()
model.add(keras.Input(shape=(2,)))
model.add(layers.Dense(100, activation="relu"))
model.add(layers.Dense(3, activation="softmax"))

{ } model.summary()

Model: "sequential_1"
Layer (type)          Output Shape         Param #
=====            (None, 100)           300
dense_2 (Dense)      (None, 3)            303
=====
Total params: 603
Trainable params: 603
Non-trainable params: 0

[ ] model.compile(optimizer=keras.optimizers.Adam(), # suggest in colab : Option+Esc in Mac
                  loss=keras.losses.CategoricalCrossentropy(),
```

73 / 74

+ Code + Text

✓ RAM Disk

dense_2 (Dense)	(None, 100)	300
dense_3 (Dense)	(None, 3)	303

Total params: 60

Trainable params: 603

Non-trainable params: 0

```
[ ] model.compile(optimizer=keras.optimizers.Adam(), # suggest in colab : Option+Esc in Mac  
[ ]         loss=keras.losses.CategoricalCrossentropy(),  
[ ]         metrics=[keras.metrics.BinaryAccuracy()])
```

```
[ ] print(X.shape)  
      print(y.shape)
```

(300, 2)  
(300, )

```
[ ] print(y)
```

+ Code + Text

✓ RAM Disk

dense_2 (Dense)	(None, 100)	300
dense_3 (Dense)	(None, 3)	303

```
Total params: 603  
Trainable params: 603  
Non-trainable params: 0
```

```
[ ] model.compile(optimizer=keras.optimizers.Adam(), # suggest in colab : Option+Esc in Mac  
✓ loss=keras.losses.CategoricalCrossentropy(),  
metrics=[keras.metrics.BinaryAccuracy()])
```

```
[ ] print(X.shape)  
      print(y.shape)
```

(300, 2)  
(300, )

```
[ ] print(y)
```

+ Code + Text

✓ RAM Disk

A set of small, light-gray navigation icons located at the bottom right of the page. From left to right, they include: a downward arrow, a circular arrow, a speech bubble, a gear, a square with a diagonal line, a trash can, and three vertical dots.

dense_2 (Dense)	(None, 100)	300
dense_3 (Dense)	(None, 3)	303

100

```
Total params: 603  
Trainable params: 603  
Non-trainable params: 0
```

```
[ ] model.compile(optimizer=keras.optimizers.Adam(), # suggest in colab : Option+Esc in Mac  
    loss=keras.losses.CategoricalCrossentropy(),  
    metrics=[keras.metrics.BinaryAccuracy()])
```

```
[ ] print(X.shape)  
print(y.shape)
```

(300, 2)  
(300, )

```
[ ] print(y)
```

+ Code + Text

✓ RAM Disk



dense 3 (Dense) (None, 3) 303

```
Total params: 603  
Trainable params: 603  
Non-trainable params: 0
```

```
[ ] model.compile(optimizer=keras.optimizers.Adam(), # suggest in colab : Option+Esc in Mac
                  loss=keras.losses.CategoricalCrossentropy(),
                  metrics=[keras.metrics.BinaryAccuracy()])
```

```
[ ] print(x.shape)  
      print(y.shape)
```

✓ (300, 2)  
(300, )

[ ] print(y)

+ Code + Text

[ ] print(y)

```
[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2  
 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2  
 2 2 2 2]
```

```
history = model.fit(X, y, epochs=10, batch_size=32)
```

Epoch 1/10

```
/usr/local/lib/python3.7/dist-packages/tensorflow/python/data/ops/structured_function.py:265: UserWarning: Even thou  
"Even though the `tf.config.experimental_run_functions_eagerly`"
```

```
ValueError Traceback (most recent call last)
```

```
<ipython-input-39-2de9726b75a6> in <module>  
----> 1 history = model.fit(X, y, epochs=10, batch_size=32)
```

1 frames

```
/usr/local/lib/python3.7/dist-packages/keras/backend.py in categorical_crossentropy(target, output, from_logits,  
axis)
```

```
5081 target = tf.convert_to_tensor(target)
```

Desmos | Graphing x BackpropFromScratch x K Why choose Keras x K Getting started x K Introduction to Keras x K Introduction to Keras x srivastava14a.pdf x Calculus on Comp x Softmax function x Derivative of the Sigmoid x +

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=RI24jRUzgdGz

+ Code + Text RAM Disk

[ ] history = model.fit(X, y, epochs=10, batch\_size=32)

{ } Epoch 1/10  
/usr/local/lib/python3.7/dist-packages/tensorflow/python/data/ops/structured\_function.py:265: UserWarning: Even though the `tf.config.experimental\_run\_functions\_eagerly`

ValueError Traceback (most recent call last)  
<ipython-input-39-2de9726b75a6> in <module>  
----> 1 history = model.fit(X, y, epochs=10, batch\_size=32)

1 frames  
/usr/local/lib/python3.7/dist-packages/keras/backend.py in categorical\_crossentropy(target, output, from\_logits, axis)  
5081 target = tf.convert\_to\_tensor(target)  
5082 output = tf.convert\_to\_tensor(output)  
-> 5083 target.shape.assert\_is\_compatible\_with(output.shape)  
5084  
5085 # Use logits whenever they are available. `softmax` and `sigmoid`

ValueError: Shapes (32, 1) and (32, 3) are incompatible

SEARCH STACK OVERFLOW

[ ] from tensorflow.keras.utils import to\_categorical

+ Code + Text

✓ RAM Disk

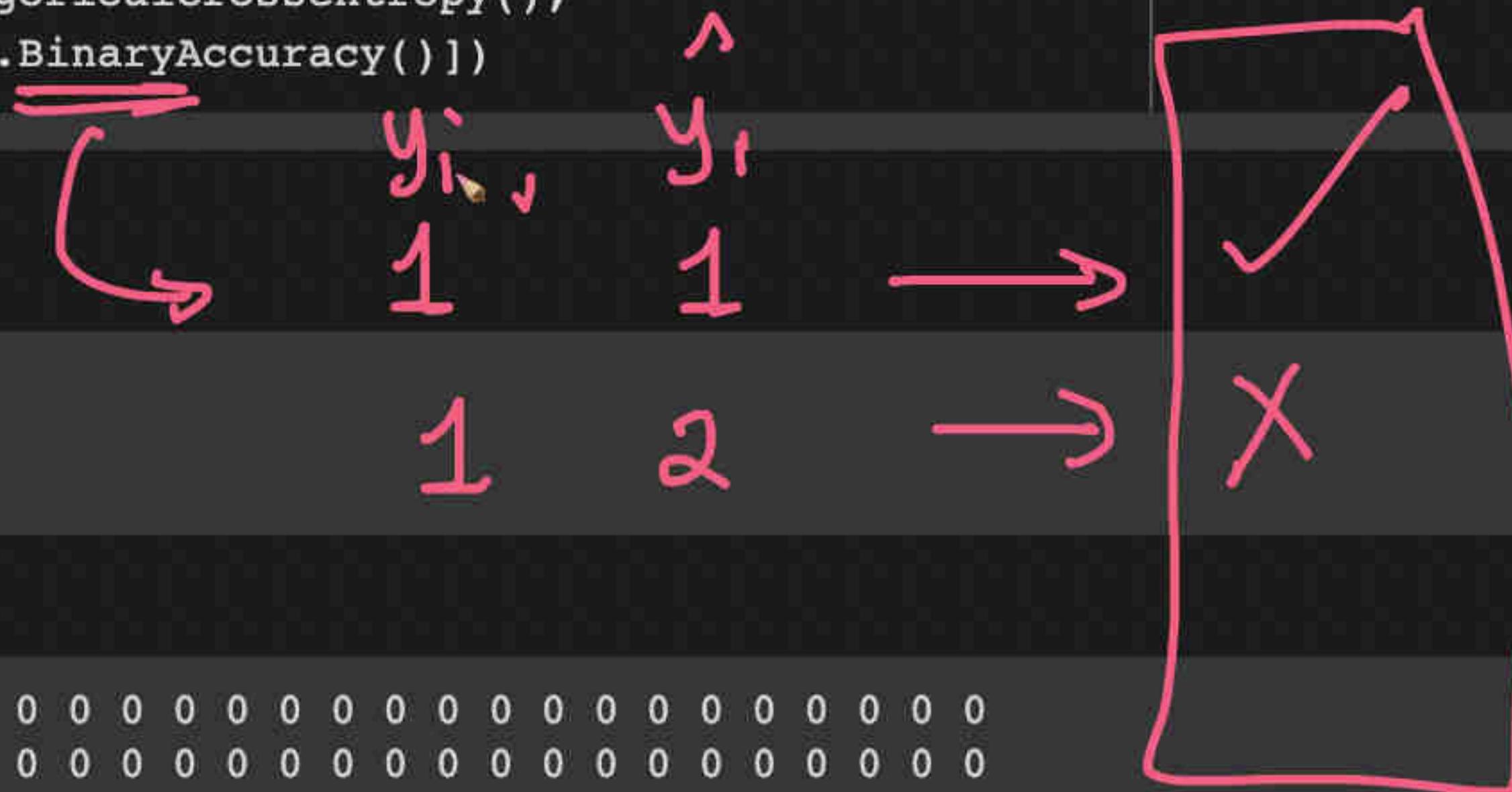
Total params: 603  
Trainable params: 603  
Non-trainable params: 0

```
[ ] model.compile(optimizer=keras.optimizers.Adam(), # suggest in colab : Option+Esc in Mac  
    loss=keras.losses.CategoricalCrossentropy(),  
    metrics=[keras.metrics.BinaryAccuracy()])
```

```
[ ] print(X.shape)  
print(y.shape)
```

(300, 2)  
(300, )

```
[ ] print(y)
```



Two pink parallel lines.

$$m-pls = \underline{\underline{300}} \text{ (let)}$$

GD

pass all m pls



Loss<sub>all m  
pls</sub>

$$\frac{\partial L_i}{\partial w_j}$$

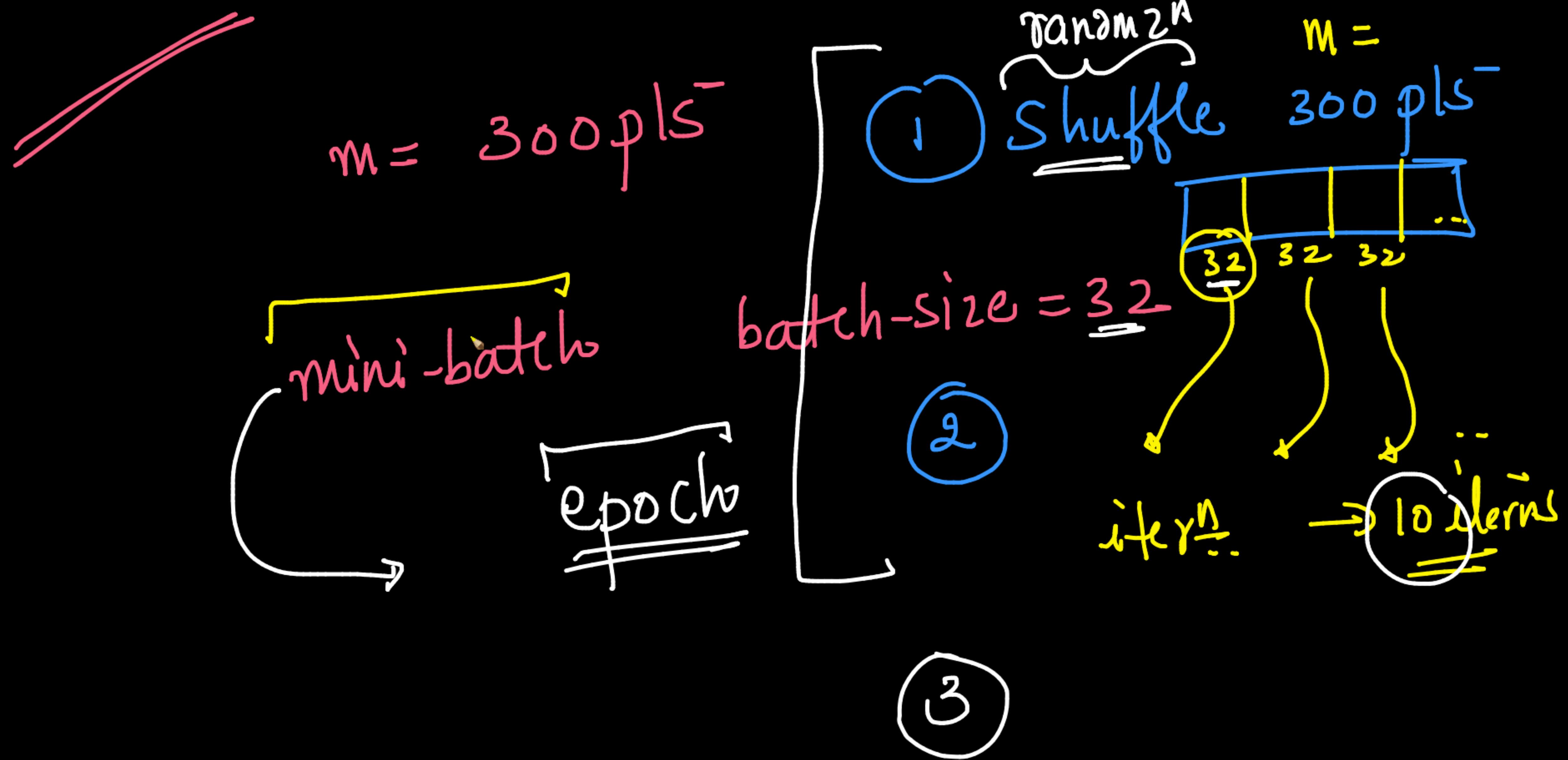
using all pls

mini batch GD:

batch-size = 32  
out of  
 $300 = M$

$$\frac{\partial L}{\partial w_j} \text{ using minibatch}$$

SGD: use only 1 pt in each iter



+ Code + Text

✓ RAM Disk

▼

(300, 2)  
(300, )

```
[ ] print(y)
```

```
[ ] history = model.fit(X, y, epochs=10, batch_size=32)
```

Epoch 1/1

```
/usr/local/lib/python3.7/dist-packages/tensorflow/python/data/ops/structured_function.py:265: UserWarning: Even thou  
"Even though the `tf.config.experimental_run_functions_eagerly` "
```

## ValueError

## Traceback (most recent call last)

[<ipython-input-39-2de9726b75a6>](#) in [<module](#)

```
--> 1 history = model.fit(X, y, epochs=10, batch_size=32)
```

1 frames

colab.research.google.com/drive/1vmwg6vCH4F42G9MoYzicAWOx8rm97z4i#scrollTo=R124jRUzqdGz

+ Code + Text

```
[ ]      5085 # Use logits whenever they are available. sof  
ValueError: Shapes (32, 1) and (32, 3) are incompatible
```

## SEARCH STACK OVERFLOW

```
[ ] from tensorflow.keras.utils import to_categorical  
y_OHE = to_categorical(y)  
print(y_OHE)
```

Desmos | Graphing x BackpropFromScratch x K Why choose Keras x K Getting started x K Introduction to Keras x K Introduction to Keras x srivastava14a.pdf x Calculus on Comp x Softmax function x Derivative of the Sigmoid x +

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=RI24jRUzgdGz

+ Code + Text

RAM Disk

{x}

```
[ ] [0. 0. 1.]  
[0. 0. 1.]  
[0. 0. 1.]  
[0. 0. 1.]  
[0. 0. 1.]
```

[ ] `tf.config.run_functions_eagerly(True)`

history = model.fit(X, y\_OHE, epochs=10, batch\_size=256)

Epoch 1/10  
2/2 [=====] - 0s 32ms/step - loss: 1.1073 - binary\_accuracy: 0.6667  
Epoch 2/10  
2/2 [=====] - 0s 29ms/step - loss: 1.0991 - binary\_accuracy: 0.6667  
Epoch 3/10  
2/2 [=====] - 0s 35ms/step - loss: 1.0913 - binary\_accuracy: 0.6667  
Epoch 4/10  
2/2 [=====] - 0s 26ms/step - loss: 1.0835 - binary\_accuracy: 0.6667  
Epoch 5/10  
2/2 [=====] - 0s 63ms/step - loss: 1.0758 - binary\_accuracy: 0.6667  
Epoch 6/10  
2/2 [=====] - 0s 26ms/step - loss: 1.0682 - binary\_accuracy: 0.6667  
Epoch 7/10  
2/2 [=====] - 0s 28ms/step - loss: 1.0606 - binary\_accuracy: 0.6667  
Epoch 8/10  
2/2 [=====] - 0s 27ms/step - loss: 1.0531 - binary\_accuracy: 0.6667  
Epoch 9/10  
2/2 [=====1 - 0s 67ms/step - loss: 1.0458 - binary\_accuracy: 0.6667

87/87

```
+ Code + Text ✓ RAM Disk
```

[ ] tf.config.run\_functions\_eagerly(True)  
history = model.fit(X, y\_OHE, epochs=10, batch\_size=256)

{x} Epoch 1/10  
2/2 [=====] - 0s 32ms/step - loss: 1.1073 - binary\_accuracy: 0.6667

Epoch 2/10  
2/2 [=====] - 0s 29ms/step - loss: 1.0991 - binary\_accuracy: 0.6667

Epoch 3/10  
2/2 [=====] - 0s 35ms/step - loss: 1.0913 - binary\_accuracy: 0.6667

Epoch 4/10  
2/2 [=====] - 0s 26ms/step - loss: 1.0835 - binary\_accuracy: 0.6667

Epoch 5/10  
2/2 [=====] - 0s 63ms/step - loss: 1.0758 - binary\_accuracy: 0.6667

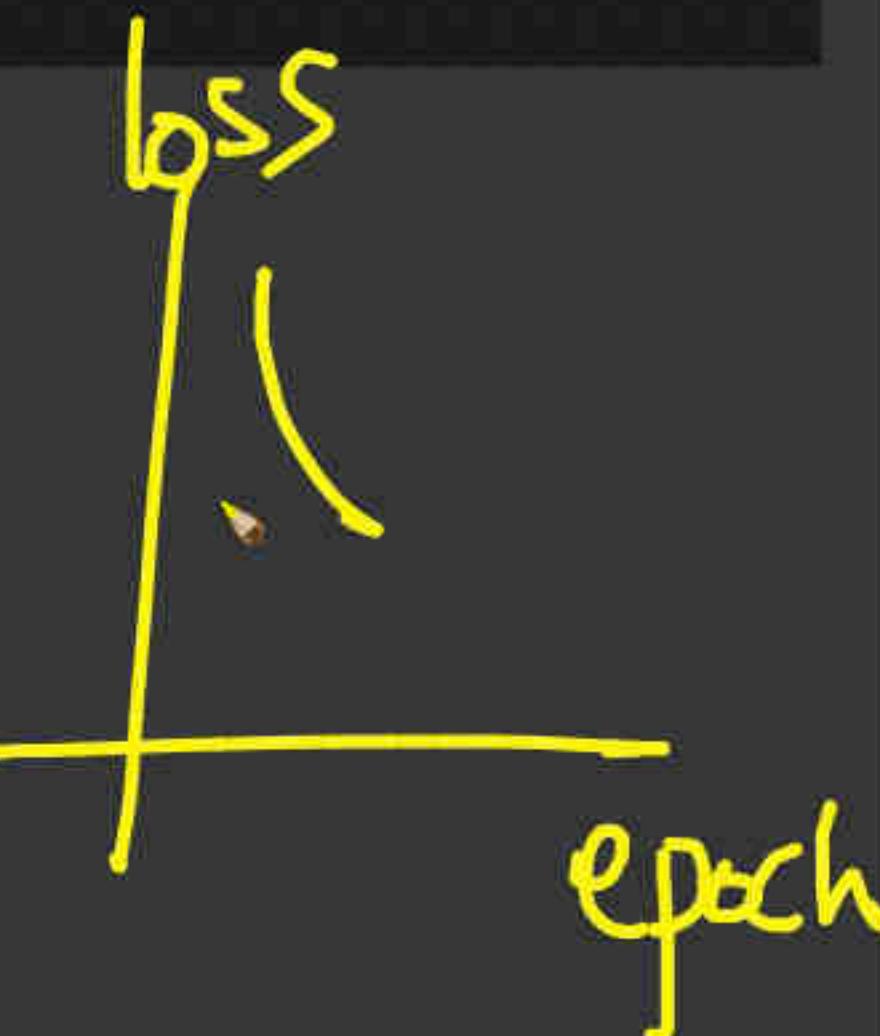
Epoch 6/10  
2/2 [=====] - 0s 26ms/step - loss: 1.0682 - binary\_accuracy: 0.6667

Epoch 7/10  
2/2 [=====] - 0s 28ms/step - loss: 1.0606 - binary\_accuracy: 0.6667

Epoch 8/10  
2/2 [=====] - 0s 27ms/step - loss: 1.0531 - binary\_accuracy: 0.6667

Epoch 9/10  
2/2 [=====] - 0s 67ms/step - loss: 1.0458 - binary\_accuracy: 0.6667

Epoch 10/10  
2/2 [=====] - 0s 46ms/step - loss: 1.0385 - binary\_accuracy: 0.6667



```
[ ] tf.config.run_functions_eagerly(True)  
history = model.fit(X, y_OHE, epochs=100, batch_size=256)
```

Desmos | Graphing x BackpropFromScratch x K Why choose Keras x K Getting started x K Introduction to Keras x K Introduction to Keras x srivastava14a.pdf x Calculus on Comp x Softmax function x Derivative of the Sigmoid x +

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=RI24jRUzgdGz

+ Code + Text ✓ RAM Disk

```
[ ] 2/2 [=====] - 0s 28ms/step - loss: 0.1136 - binary_accuracy: 0.9800
[ ] Epoch 998/1000
2/2 [=====] - 0s 35ms/step - loss: 0.1134 - binary_accuracy: 0.9811
Epoch 999/1000
2/2 [=====] - 0s 27ms/step - loss: 0.1132 - binary_accuracy: 0.9811
Epoch 1000/1000
2/2 [=====] - 0s 24ms/step - loss: 0.1131 - binary_accuracy: 0.9811
```

Q {x} D ▶ TF & Keras: Functional API

[ ] ↳ 5 cells hidden

◀ Computational Graph

[ ] ↳ 3 cells hidden

<>

☰

▶

Page-Footer

colab.research.google.com/drive/1vmwq6vCH4F42G9MoYzjcAWOx8rm97z4i#scrollTo=cYHEvbE-4zXy

+ Code + Text RAM Disk ✓

## Simple NN (in Python)

{x}

[1] !gdown 1dLOPwh01o3k8p\_hK633ixhD1ehz6nNWk

Downloading...  
From: [https://drive.google.com/uc?id=1dLOPwh01o3k8p\\_hK633ixhD1ehz6nNWk](https://drive.google.com/uc?id=1dLOPwh01o3k8p_hK633ixhD1ehz6nNWk)  
To: /content/spiral.csv  
100% 12.9k/12.9k [00:00<00:00, 17.1MB/s]

[2] !ls -lrt

total 20  
drwxr-xr-x 1 root root 4096 Aug 31 13:47 sample\_data  
-rw-r--r-- 1 root root 12867 Sep 5 16:14 spiral.csv

# Very simple case of a back-prop  
import pandas as pd  
import matplotlib.pyplot as plt  
import numpy as np  
df = pd.read\_csv("spiral. spiral.csv")  
plt.scatter(df["x1"], df["x2"], c=df["y"], s=40, cmap=plt.cm.Spectral)

90/90 ✓ 1s completed at 21:44