

## ② Softmax Reg.

$$w = [w^{(1)} \dots w^{(c)}]$$

$$\hat{y}_k = \frac{\exp z_k}{\sum_{k'=1}^c \exp z_{k'}}$$

$$z_k = n^T w^k + b_k$$

$$f_{CE}(w, b) = -\frac{1}{n} \sum_{q=1}^n \sum_{k=1}^c y_{1k}^{(q)} \log \hat{y}_{1k}^{(q)}$$

$$\begin{aligned} \nabla_{w^{(l)}} f_{CE}(w, b) &= -\frac{1}{n} \sum_{q=1}^n \sum_{k=1}^c y_{1k}^{(q)} \nabla_{w^l} \log \hat{y}_{1k}^{(q)} \\ &= -\frac{1}{n} \sum_{q=1}^n \sum_{k=1}^c y_{1k}^{(q)} \left( \frac{\nabla_{w^l} \hat{y}_{1k}^{(q)}}{\hat{y}_{1k}^{(q)}} \right) \end{aligned}$$

Case 1 for  $l=k$

$$\nabla_{w(u)} \hat{y}_k = \nabla_{w(u)} \frac{\exp z_k}{\sum_{k'=1}^c \exp z_{k'}}$$

$$\frac{\partial}{\partial u} \left[ \frac{f(u)}{g(u)} \right] = \frac{g(u) \cdot f'(u) - f(u) g'(u)}{[g(u)]^2}$$

$$\frac{\partial}{\partial u} \sum_{k'=1}^c \exp z_{k'} \cdot \nabla_{w(u)} \exp z_k -$$

$$\exp z_k \cdot \nabla_{w(u)} \frac{\sum_{k'=1}^c \exp z_{k'}}{\left( \sum_{k'=1}^c \exp z_{k'} \right)^2}$$

$$\left( \sum_{k'=1}^c \exp z_{k'} \right)^2$$

$$\sum_{k'=1}^c \exp^{2k'} \cdot \exp^{2k} \cdot u^{(p)} - \exp^{2k} \cdot u^{(q)} \cdot \sum_{k'=1}^c \exp^{2k'}$$

$$\sum_{k'=1}^c \exp^{2k'}$$

$$\left( \sum_{k'=1}^c \exp^{2k'} \right)^2$$

$$= \frac{n^{(q)} \exp^{z_k}}{\sum_{k'=1}^c \exp^{z_{k'}}} \left[ 1 - \frac{\exp^{z_k}}{\sum_{k'=1}^c \exp^{z_{k'}}} \right]$$

$$= n^{(q)} \hat{y}_k^q (1 - \hat{y}_k^q) \quad \therefore (1)$$

Now for case (2)  $k \neq k'$  we get

$$= \sum_{k'=1}^c \exp^{z_{k'}} \cdot (0) - e^{z_k} \cdot n^{(q)} \cdot e^{z_k}$$

$$= n^{(q)} \cdot \hat{y}_k^q \cdot \hat{y}_{k'}^q \quad \therefore (2)$$

Now to find total gradient of fce.

$$\nabla_{W^{(q)}} fce(l, b) = -\frac{1}{n} \sum_{q=1}^l \sum_{k=1}^c y_k^q \nabla_{W^{(q)}} \log \hat{y}_k^q$$

$$= -\frac{1}{n} \sum_{q=1}^l \sum_{k=1}^c y_k^q \frac{1}{\hat{y}_k^q} \cdot \nabla_{W^{(q)}} \hat{y}_k^q$$

$$= -\frac{1}{n} \sum_{q=1}^l \sum_{k=1}^c y_k^q \cdot \frac{1}{\hat{y}_k^q} \left[ n^q \hat{y}_k^q (1 - \hat{y}_k^q) - \frac{n^q \hat{y}_k^q}{\sum_{k'=1}^c \hat{y}_{k'}^q} \right]$$

$\sum_k y_k = 1$  given &  $\sum_k \hat{y}_k = 1$  (softmax layer)  
softmax 2

$$= -\frac{1}{n} \sum_{q=1}^l \sum_{k=1}^c y_k^q \cdot \frac{1}{\hat{y}_k^q} \cdot n^q \hat{y}_k^q [1 - \hat{y}_k^q - \hat{y}_{k'}^q]$$

$$\text{Also } 1 - \hat{y}_k^q = \hat{y}_{k'}^q$$

$$= -\frac{1}{n} \sum_{q=1}^l \sum_{k=1}^c \frac{1}{2} \cdot n^q (1 - \hat{y}_k^q - 1 + \hat{y}_k^q)$$

$$= -\frac{1}{n} \sum_{q=1}^n n^p (y_q^p - \hat{y}_q^p) \quad \text{... (3)}$$

Now since  $(l, b)$  will be

$$\nabla_b = \left( -\frac{1}{n} \sum_{q=1}^n \sum_{k=1}^c y_k^p \log \hat{y}_k^p \right)$$

$$= -\frac{1}{n} \sum_{q=1}^n \sum_{k=1}^c y_k^p \nabla_b \log \left( \frac{\exp^{-2k}}{\sum_{k'=1}^c \exp^{-2k'}} \right)$$

$$\nabla_b \log \left( \frac{\exp^{-2k}}{\sum_{k'=1}^c \exp^{-2k'}} \right) = \nabla_b \log (\exp^{-2k}) - \nabla_b \log \left( \sum_{k'=1}^c \exp^{-2k'} \right)$$

$$= 1 - \frac{\exp^{-2k}}{\sum_{k'=1}^c \exp^{-2k'}}$$

$$\therefore \nabla_b = -\frac{1}{n} \sum_{q=1}^n \sum_{k=1}^c y_k^p \left[ 1 - \frac{\exp^{-2k}}{\sum_{k'=1}^c \exp^{-2k'}} \right]$$

$$= -\frac{1}{n} \sum_{q=1}^n \left[ \sum_{k=1}^c y_k^p - \sum_{k=1}^c y_k^p \cdot \frac{\exp^{-2k}}{\sum_{k'=1}^c \exp^{-2k'}} \right]$$

$$= -\frac{1}{n} \sum_{q=1}^n (y^p - \hat{y}^p) \quad \text{... (4)}$$

②

Cross-Entropy as Neg. log-likelihood

$$\hat{y}_k = P(y_k = 1 | x_i, w, b) \cdot \text{AIC} \in (0, 1)$$

Now  $P(y_i = 1 | x_i, w, b) = P(y_i = 1 | x_i, w, b)^{y_i} \times \dots$

$$\therefore P(y_c = 1 | x_i, w, b)^{y_c}$$

$$P(y | x, w, b) = \prod_{k=1}^c \hat{y}_k^{y_k}$$

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$$

$$-\log P(D | w, b) = -\log P(y | x, w, b)$$

$$= -\log \left( \prod_{k=1}^c \hat{y}_k^{y_k} \right) = -\sum_{k=1}^c \log (\hat{y}_k^{y_k})$$

$$\text{Now } -\log a^b = b \cdot \log a$$

$\therefore$  we get  $-\sum_{k=1}^c y_k^p \log \hat{y}_k^p$   $\leftarrow$  to sum over all samples we get

$$-\sum_{i=1}^n \sum_{k=1}^c y_k^p \log \hat{y}_k^p = \text{f}_{CE}(D; w, b)$$

$$\textcircled{3} \quad f\left(\begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix}\right) = \begin{bmatrix} u_1 - 2u_2 + u_3 \\ u_1 \end{bmatrix}$$

$$g\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1 + 2x_2 \\ x_2 \\ -x_1 + 3 \end{bmatrix}$$

$\textcircled{4}$

$$g = \sigma(w_1 u + b_1) \quad \text{where } \sigma = \text{linear or transformation}$$

$$f = \sigma(w_2 u + b_2) \quad \text{non-linear}$$

$$\nabla g(u) = \begin{bmatrix} \nabla_{w_1} (u_1 + 2u_2) & \nabla_{w_2} (u_1 + 2u_2) \\ \nabla_{w_2} (u_2) & \nabla_{w_2} u_2 \\ \nabla_{u_1} (-u_1 + 3) & \nabla_{u_2} (-u_1 + 3) \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 2 \\ 0 & 1 \\ -1 & 0 \end{bmatrix} \quad B_{PA_0} = \begin{bmatrix} 0 \\ 0 \\ 3 \end{bmatrix}$$

$$\therefore g([u_1 \ u_2]) = \begin{bmatrix} 1 & 2 \\ 0 & 1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 3 \end{bmatrix}$$

$$= L_1 x + b_1$$

Now for f using similar approach we get

$$= \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix} \left[ \begin{bmatrix} 1 & -2 & \frac{1}{4} \end{bmatrix} \right] + \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\Rightarrow f([u_1 \ u_2 \ u_3]) = \begin{bmatrix} 1 & -2 & \frac{1}{4} \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$= L_2 x + b_2$$

$$g = \sigma(L_1 x + b_1) \quad \text{or linear or non-linear}$$

$$f = \sigma(L_2 x + b_2) \quad \text{transformation } f^n$$

$$\text{i.e. } f = \sigma[L_2(\sigma[L_1 x + b_1]) + b_2]$$

(3e)

$$z_1 = f(z_1) \quad \frac{\partial z_1}{\partial b_2}(z_1) = ?$$

$$z_2 = L_2(z_1) + b_2 \quad \frac{\partial z_2}{\partial b_2} = (0+1) = 1$$

(3f)

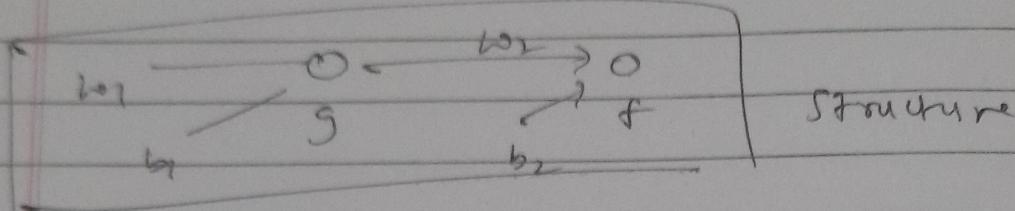
$$\text{given } g = L_1 x + b_1 \quad f = L_2 x + b_2$$

$$\frac{\partial(f \circ g)(x)}{\partial b_2} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial b_2} = 2$$

$$\text{i.e. } f(L_1 x + b_1) = h(x).$$

$$2\textcircled{e} \quad \frac{\partial(f \circ g)}{\partial b_i} x = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial b_i}$$

$$\frac{\partial (w_1x + b_1)}{\partial b_1} = 1 \quad \text{so } \frac{\partial (f \circ g)(u)}{\partial b_1} \approx 0.$$



$$3\textcircled{f} \quad \text{unrec} \left[ \frac{\partial(f \circ g)}{\partial \text{vec}(w_2)} u \right]$$

$$\partial \text{vec}(w_2) = ? \quad w \Rightarrow w_{n \times m} \text{ matrix}$$

$$\text{vec}(w) = [w_{11} \dots w_{1m} \dots w_{n1} \dots w_{nm}]$$

$$\frac{\partial(f \circ g)}{\partial \text{vec}(w_2)} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial \text{vec}(w_2)} \quad \frac{\partial f}{\partial g} = p^T$$

$$\frac{\partial g}{\partial \text{vec}(w_2)} = \begin{bmatrix} \frac{\partial g_1}{\partial w_{11}} & \dots & \frac{\partial g_1}{\partial w_{1m}} & \dots & \frac{\partial g_1}{\partial w_{n1}} & \dots & \frac{\partial g_1}{\partial w_{nm}} \\ \vdots & & & & \vdots & & & \end{bmatrix}$$

$$\text{unrec} \left[ \frac{\partial(f \circ g)}{\partial \text{vec}(w_2)} \right] \approx p^T \odot$$

$$3\textcircled{g} \quad \text{unrec} \left[ \frac{\partial(f \circ g)}{\partial \text{vec}(w_1)} u \right] = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial \text{vec}(w_1)}$$

$$\text{Suppose } \frac{\partial f}{\partial g} = p^T \quad \frac{\partial g}{\partial \text{vec}(w_1)} = \begin{bmatrix} \frac{\partial g_1}{\partial w_{11}} & \dots & \frac{\partial g_1}{\partial w_{1m}} \\ \vdots & & \vdots \\ \frac{\partial g_n}{\partial w_{11}} & \dots & \frac{\partial g_n}{\partial w_{1m}} \end{bmatrix}$$

$$\frac{\partial g}{\partial \text{vec}(w_1)} = \begin{bmatrix} n^T & \dots & 0^T \\ \vdots & & \vdots \\ 0^T & \dots & x^T \end{bmatrix} = x^T \quad \frac{\partial g_1}{\partial w_{11}} \dots \frac{\partial g_1}{\partial w_{1m}}$$

$$\therefore \text{unrec} \left[ \frac{\partial(f \circ g)}{\partial \text{vec}(w_1)} u \right] = p^T \cdot \text{unrec} [p_1 n^T \dots p_n x^T] \quad (= p x^T)$$

```

1 # hw_3 problem 3b
2 import numpy as np
3 def AffineTransformation(W: np.ndarray,x: np.ndarray,b: np.ndarray) -> np.ndarray:
4     return np.dot(W,x)+b
5 W1: np.ndarray = np.array([[1, 2], [0, 1],[-1, 0]])
6 b1: np.ndarray = np.array([0, 0, 3])
7 W2: np.ndarray = np.array([[1, -2, 1/4]])
8 b2: np.ndarray = np.array([0])
9 x1: np.ndarray = np.array([1, 2])
10 x2: np.ndarray = np.array([1, 2, 3])
11 result_g: np.ndarray = AffineTransformation(W1,x1,b1)
12 result_f: np.ndarray = AffineTransformation(W2,x2,b2)
13 print(f"Function g result: {result_g} & Function f result: {result_f}")

```

Function g result: result\_g = array([5, 2, 2]) & Function f result: result\_f = array([-2.25])

```

1 # hw_3 problem 3c
2 def Composition(L: list, x: np.ndarray) -> list:
3     z_values: list = list()
4     for W,b in L:
5         # get pre-activation
6         z = AffineTransformation(W,x,b)
7         z_values.append(z)
8         # use for the next iteration
9         x: np.ndarray = z
10    return z_values
11
12 L: list = [
13     (np.array([[1,2,1]]), np.array([3,4,5])),
14     (np.array([[1,-2,1/4]]), np.array([10,20,30])),
15 ]
16
17 x: np.ndarray = np.array([100,200,300])
18 z_values: list = Composition(L,x)
19 print(f"Composition result: {z_values}")

```

Composition result: z\_values = [array([803, 804, 805]), array([-593.75, -583.75, -573.75])]

```

1 # hw_3 problem 3d
2 final_output = Composition(L,x)[-1]
3 # Print the final output
4 print(f"Final Output (f o g)(x): {final_output}")

```

Final Output (f o g)(x): final\_output = array([-593.75, -583.75, -573.75])

```

1 # hw_3 problem 3h
2 def ComputeJacobians(L: np.ndarray, x: np.ndarray, z: np.ndarray) -> tuple:
3     num_layers = len(L)
4     dz_db_list: list = list()
5     dz_dW_list: list = list()
6     # Compute the Jacobians for each layer in reverse order
7     for i in range(num_layers - 1, -1, -1):
8         W, b = L[i]
9         z_i: np.ndarray = z[i]
10        dz_dx: np.ndarray = W
11        dz_db: np.ndarray = np.identity(len(b))
12        dz_dW: np.ndarray = np.outer(np.ones_like(z_i), x)
13        dz_db_list.append(np.dot(dz_dx, dz_db))
14        dz_dW_list.append(np.dot(dz_dx, dz_dW))
15        x: np.ndarray = AffineTransformation(W, x, b)
16        dz_db_list.reverse()
17        dz_dW_list.reverse()
18    return dz_dW_list, dz_db_list
19 dz_dW, dz_db = ComputeJacobians(L, x, z_values)
20 print(f"Jacobians : \nw -> {dz_dW = } \nb -> {dz_db = }")

```

Jacobians :

```
w -> dz_dW = [array([[-860., -820., -780.]]), array([[ -75., -150., -225.]])]
b -> dz_db = [array([[1., 2., 1.]]), array([[ 1. , -2. ,  0.25]])]
```

```

# hw_3 problem 4
import numpy as np
X_tr = np.load("fashion_mnist_train_images.npy")
y_tr = np.load("fashion_mnist_train_labels.npy")
X_te = np.load("fashion_mnist_test_images.npy")
y_te = np.load("fashion_mnist_test_labels.npy")
print(X_tr.shape,y_tr.shape,X_te.shape,y_te.shape)

split_80_20: int = int(X_tr.shape[0]*0.8)
print('80 / 20 Split : ',split_80_20)
X_tr_80, y_tr_80 = X_tr[:split_80_20,:], y_tr[:split_80_20]
X_vl_20, y_vl_20 = X_tr[split_80_20:,:], y_tr[split_80_20:]
del X_tr,y_tr
print(X_tr_80.shape,y_tr_80.shape,X_te.shape,y_te.shape,X_vl_20.shape,y_vl_20.shape)

grid_search_mini_batch_size : list = [100,120,140,160,180,200]
grid_search_learning_rate : list = [3e-7,3e-5,3e-3,3e-1]
grid_search_no_of_epochs : list = [8,16,20,24,32,48,64,80]
grid_search_regularization_val : list = [1e-5,1e-3,1e-1]

# custom helper functions
def get_initialize_weights_and_bias_liner_reg(input_dims:int,output_dims:int) -> tuple:
    return np.random.uniform(low=1e-5,high=1e-2,size=(input_dims,output_dims)),np.random.random(size=output_dims)

def get_y_hat_pred_func(x: np.ndarray, w: np.ndarray, b: np.ndarray) -> np.ndarray:
    return np.dot(x, w) + b

def get_softmax(y_pred: np.ndarray) -> np.ndarray:
    exp_y_pred = np.exp(y_pred - np.max(y_pred, axis=1, keepdims=True))
    return exp_y_pred / np.sum(exp_y_pred, axis=1, keepdims=True)

def get_cross_entropy_regularized_loss(y_true: np.ndarray,y_pred: np.ndarray,W: np.ndarray,alpha_reg: float) -> float:
    n = len(y_true)
    # addition of 1e-10 is to silence - RuntimeWarning: divide by zero encountered in log
    cross_entropy = -1.0 * np.sum(y_pred * np.log(y_pred + 1e-10)) + (0.5 * alpha_reg) * np.sum(W.T @ W)
    cross_entropy = cross_entropy / n
    return cross_entropy

```

```
def get_cross_entropy_unregularized_loss(y_true: np.ndarray,y_pred: np.ndarray) -> float:
    n = len(y_true)
    cross_entropy = -1.0 * np.sum(y_pred * np.log(y_pred + 1e-10))
    cross_entropy = cross_entropy / n
    return cross_entropy

# hyperparameter tuning
best_params: dict = {
    'batch_size': 0,
    'learning_rate': 0,
    'num_epochs': 0,
    'regularization_strength': 0,
    'reg_loss': 0,
    'unreg_loss': 0,
    'accuracy': 0,
}

best_loss: float = 10_000.0
best_accuracy:float = 0.70

# Shuffle training data
shuffling_indices: np.ndarray = np.arange(X_tr_80.shape[0])
np.random.shuffle(shuffling_indices)
X_train_shuffled: np.ndarray = X_tr_80[shuffling_indices]
y_train_shuffled: np.ndarray = y_tr_80[shuffling_indices]

input_dims1,output_dims1 = X_train_shuffled.shape[1],10
```

```
for initial_mini_batch_size in grid_search_mini_batch_size:
    for initial_no_of_epochs in grid_search_no_of_epochs:
        for initial_learning_rate in grid_search_learning_rate:
            for initial_regularization_strength in grid_search_regularization_val:

                # Initialize weights and bias
                w, b = get_initialize_weights_and_bias_liner_reg(input_dims1, output_dims1) Loading...

                # Training Loop
                for epoch in range(initial_no_of_epochs):

                    for batch_start in range(0, X_tr_80.shape[0], initial_mini_batch_size):

                        # taking mini batches from training data
                        X_mini_batch: np.ndarray = X_train_shuffled[batch_start:batch_start+initial_mini_batch_size,:]
                        # X_mini_batch: np.ndarray = X_mini_batch / 255.0
                        len_X_mini_batch = len(X_mini_batch)
                        y_mini_batch: np.ndarray = y_train_shuffled[batch_start:batch_start+initial_mini_batch_size]

                        # we compute the y hat predict
                        y_pred_mini: np.ndarray = get_softmax(get_y_hat_func(X_mini_batch, w, b))

                        # Backpropagation
                        # when pred is correct for the class remove it so that false case gets more weight
                        y_pred_mini[np.arange(len(y_mini_batch)), y_mini_batch] -= 1

                        # we compute gradient of w and b
                        grad_w: np.ndarray = np.dot(X_mini_batch.T, y_pred_mini)
                        grad_b: np.ndarray = np.sum(y_pred_mini, axis=0)

                        # update learning weights and biases
                        w -= initial_learning_rate * (grad_w + initial_regularization_strength * w)
                        b -= initial_learning_rate * grad_b

                # Validation
                valid_probabilities: np.ndarray = get_softmax(get_y_hat_func(X_vl_20,w,b))
                valid_reg_loss: float = get_cross_entropy_regularizerized_loss(y_vl_20,valid_probabilities,w,initial_regularization_strength)
                valid_unreg_loss: float = get_cross_entropy_unregularized_loss(y_vl_20,valid_probabilities)
                valid_accuracy: float = np.mean(np.argmax(valid_probabilities, axis=1)==y_vl_20)
```

```

grad_b: np.ndarray = np.sum(y_pred_mini, axis=0)

# update learning weights and biases
w -= initial_learning_rate * (grad_w + initial_regularization_strength * w)
b -= initial_learning_rate * grad_b

# Validation
valid_probabilities: np.ndarray = get_softmax(get_y_hat_pred_func(X_vl_20,w,b))
valid_reg_loss: float = get_cross_entropy_regularizerized_loss(y_vl_20,valid_probabilities,w,initial_regularization_strength)
valid_unreg_loss: float = get_cross_entropy_unregularized_loss(y_vl_20,valid_probabilities)
valid_accuracy: float = np.mean(np.argmax(valid_probabilities, axis=1)==y_vl_20)

if valid_unreg_loss > 0 and valid_reg_loss > 0 and best_accuracy > 0.70:
    if valid_unreg_loss < best_loss and valid_accuracy < best_accuracy:
        best_loss: float = valid_unreg_loss
        best_accuracy: float = valid_accuracy
        best_params = {
            'batch_size': initial_mini_batch_size,
            'learning_rate': initial_learning_rate,
            'num_epochs': initial_no_of_epochs,
            'regularization_strength': initial_regularization_strength,
            'valid_reg_loss': valid_reg_loss,
            'valid_unreg_loss': valid_unreg_loss,
            'valid_accuracy': valid_accuracy,
        }
    else:
        continue
else:
    break

# testing dataset
test_probabilities: np.ndarray = get_softmax(get_y_hat_pred_func(X_te,w,b))
test_reg_loss: float = get_cross_entropy_regularizerized_loss(y_te,test_probabilities,w,initial_regularization_strength)
test_unreg_loss: float = get_cross_entropy_unregularized_loss(y_te,test_probabilities)
test_accuracy: float = np.mean(np.argmax(test_probabilities, axis=1)==y_te)
best_params['test_reg_loss'] = test_reg_loss
best_params['test_unreg_loss'] = test_unreg_loss
best_params['test_accuracy'] = test_accuracy
print(f"Computing Loss on Test Data Set ->\nBest Test Reg Loss: {test_reg_loss},\nBest Test UnReg Loss: {test_unreg_loss},\nBest Accuracy: {test_accuracy}")

```

```

288     valid_reg_loss: float = get_cross_entropy_regularized_loss(y_vl_20,valid_probabilities,w,initial_regularization_strength)
289     valid_unreg_loss: float = get_cross_entropy_unregularized_loss(y_vl_20,valid_probabilities)
290     valid_accuracy: float = np.mean(np.argmax(valid_probabilities, axis=1)==y_vl_20)
291
292     if valid_unreg_loss > 0 and valid_reg_loss > 0 and best_accuracy > 0.70:
293         if valid_unreg_loss < best_loss and valid_accuracy < best_accuracy:
294             best_loss: float = valid_unreg_loss
295             best_accuracy: float = valid_accuracy
296             best_params = {
297                 'batch_size': initial_mini_batch_size,
298                 'learning_rate': initial_learning_rate,
299                 'num_epochs': initial_no_of_epochs,
300                 'regularization_strength': initial_regularization_strength,
301                 'valid_reg_loss': valid_reg_loss,
302                 'valid_unreg_loss': valid_unreg_loss,
303                 'valid_accuracy': valid_accuracy,
304             }
305         else:
306             continue
307         else:
308             break
309
310 # testing dataset
311 test_probabilities: np.ndarray = get_y_hat_pred_func(X_te,w,b)
312 test_reg_loss: float = get_cross_entropy_regularized_loss(y_te,test_probabilities,w,initial_regularization_strength)
313 test_unreg_loss: float = get_cross_entropy_unregularized_loss(y_te,test_probabilities)
314 test_accuracy: float = np.mean(np.argmax(test_probabilities, axis=1)==y_te)
315 best_params['test_reg_loss'] = test_reg_loss
316 best_params['test_unreg_loss'] = test_unreg_loss
317 best_params['test_accuracy'] = test_accuracy
318 print(f"Computing Loss on Test Data Set ->\nBest Test Reg Loss: {test_reg_loss},\nBest Test UnReg Loss: {test_unreg_loss},\nBest Accuracy: {test_accuracy}")

```

```

↳ (60000, 784) (60000,) (10000, 784) (10000,)
80 / 20 Split : 48000
(48000, 784) (48000,) (10000, 784) (10000,) (12000, 784) (12000,)

Best Hyperparameters: {'batch_size': 160, 'learning_rate': 3e-05, 'num_epochs': 32, 'regularization_strength': 0.001, 'valid_reg_loss': 0.0005008797915, 'valid_unreg_loss': 0.0005008761, 'valid_accuracy': 0.8
Best Validation Reg Loss: 0.0005008797915,
Best Validation UnReg Loss: 0.0005008761,
Best Accuracy: 0.8065001898765323

Computing Loss on Test Data Set ->
Best Test Reg Loss: 0.0006779493,
Best Test UnReg Loss: 0.000677274824,
Best Accuracy: 0.80671435352

```