B1) (a)

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \phi(y_i(w^T x_i + b)) + \lambda ||w||^2$$

$$\phi(t) = \log(1 + e^t)$$

$$J = \min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \log\left(1 + e^{-y_i(w^T x_i + b)}\right) + \lambda ||w||^2$$

$$\frac{\partial J}{\partial w} \Rightarrow \frac{1}{n} \sum_{i=1}^{n} \underbrace{\frac{1}{1 + e^{-(y_i(w^T x_i + b))}} \left[-y_i x_i e^{-y_i(w^T x_i + b)}\right]}_{} + 2\lambda w$$

dividing the numerator and denominator

by $e^{-y_i(w^T x_i + b)}$

We get,

$$\boxed{\frac{\partial J}{\partial w} = \frac{1}{n} \sum_{i=1}^{n} \frac{-y_i x_i}{1 + e^{y_i(w^T x_i + b)}} + 2\lambda w}$$

Similarly,

$$\frac{\partial J}{\partial b} = \frac{1}{n} \sum_{i=1}^{n} \underbrace{\frac{1}{1 + e^{-y_i(w^T x_i + b)}} \left[-y_i e^{-y_i(w^T x_i + b)}\right]}_{} + 0$$

dividing numerator and denominator by $e^{-y_i(w^T x_i + b)}$

We get,

$$\boxed{\frac{\partial J}{\partial b} = \frac{1}{n} \sum_{i=1}^{n} \frac{-y_i}{1 + e^{y_i(w^T x_i + b)}}}$$

Gradient descent update rule

$$w^{k+1} = w^k - \eta \frac{\partial J}{\partial w^k}$$

$$b^{k+1} = b^k - \eta \frac{\partial J}{\partial b^k}$$

where $w^k, b^k$ indicate hyperparameters at $k$th iteration.

$$w^{k+1} = w^k - \eta \left[ \frac{1}{n} \sum_{i=1}^{n} \frac{-y_i x_i}{1 + e^{y_i(w^k x_i + b^k)}} + 2\lambda w^k \right]$$

→ This is $(w^k)^T$

→ Update rules
=

$$b^{k+1} = b^k - \eta \left[ \frac{1}{n} \sum_{i=1}^{n} \frac{-y_i}{1 + e^{y_i(w^{kT} x_i + b^k)}} \right]$$

**81 (b)**
=

In 81(a) we have got

$$\frac{\partial J}{\partial w} = \frac{1}{n} \sum_{i=1}^{n} \frac{-y_i x_i}{1 + e^{y_i(w^T x_i + b)}} + 2\lambda w$$

As the objective $J$ is given to be convex, a unique global minimizer exists., $(w^*, b^*)$. And @ global minimizer

$$\frac{\partial J}{\partial w^*} = 0 \qquad \frac{\partial J}{\partial b^*} = 0$$

We will just take $\frac{\partial J}{\partial w^*} = 0$ and examine it.

$$\frac{\partial J}{\partial w^*} = 0 \implies \frac{1}{n} \sum_{i=1}^{n} \frac{-y_i x_i}{1 + e^{y_i(w^{*T} x_i + b^*)}} + 2\lambda w^* = 0$$

$$\Rightarrow \quad 2\lambda w^* = +\frac{1}{5} \sum_{i=1}^{n} \frac{y_i z_i}{1+e^{y_i(w^*\vec{5}x_i+b^*)}}$$

$$\Rightarrow \quad 2\lambda w^* = \frac{1}{5} \; x^T \begin{bmatrix} \dfrac{y_i}{1+e^{y_i(w^{*T}x_i+b^*)}} \\ \vdots \\ \dfrac{y_n}{1+e^{y_n(w^{*T}x_n+b)}} \end{bmatrix}$$

where $\quad x^T = \begin{bmatrix} 1 & 1 & & 1 \\ x_1 & x_2 & \cdots\cdots & x_n \\ 1 & 1 & & 1 \end{bmatrix}$

$$\Rightarrow \quad w^* = x^T \underbrace{\begin{bmatrix} \dfrac{1}{2\lambda n} \begin{bmatrix} \dfrac{y_i}{1+e^{y_i(2w^{*T}x_i+b^*)}} \\ \vdots \\ \dfrac{y_n}{1+e^{y_n(w^{*T}x_n+b)}} \end{bmatrix} \end{bmatrix}}_{\alpha^*}$$

$$\therefore \quad \boxed{w^* = x^T\alpha^*} \qquad \text{for some } \alpha^*$$

Q1

(C) (a) $w^{old} = x^T \alpha^{old}$

We know from Q1 (a) update formula

$$w^{new} = w^{old} - \eta\left[\frac{1}{n}\sum_{i=1}^{n}\frac{-y_i x_i}{1+e^{y_i(w^{old T}x_i+b^{old})}} + 2\lambda w^{old}\right] \qquad - ①$$

This is $(w^{old})^T$

Let us substitute $w^{old} = X^T \alpha^{old}$ in eqn ①

$$w^{new} = X^T \alpha^{old} - \eta \left[ \frac{1}{n} \sum \frac{-y_i x_i}{1 + e^{y_i (X^T \alpha^{old})^T x_i + b^{old}}} + 2\lambda X^T \alpha^{old} \right] \qquad -②$$

Let us try to analyse this term

$$\Rightarrow \frac{1}{n} \sum_{i=1}^{n} \frac{-y_i x_i}{1 + e^{y_i (\alpha^{old T} X x_i + b^{old})}} \qquad \Rightarrow \qquad \frac{1}{n} X^T \begin{bmatrix} \dfrac{-y_1}{1 + e^{y_1 (\alpha^{old T} X x_1 + b^{old})}} \\ \vdots \\ \dfrac{-y_n}{1 + e^{y_n (\alpha^{old T} X x_n + b^{old})}} \end{bmatrix} \qquad -③$$

Substituting eqn ③ back in ②

$$w^{new} = X^T \alpha^{old} - \eta \left[ \frac{1}{n} X^T \begin{bmatrix} \dfrac{-y_1}{1 + e^{y_1 (\alpha^{old T} X x_1 + b^{old})}} \\ \vdots \\ \dfrac{-y_n}{1 + e^{y_n (\alpha^{old T} X x_n + b^{old})}} \end{bmatrix} + 2\lambda X^T \alpha^{old} \right]$$

$$\Rightarrow w^{new} = X^T \left[ \alpha^{old} - \eta \left[ \frac{1}{n} \begin{bmatrix} \dfrac{-y_1}{1 + e^{y_1 (\alpha^{old T} X x_1 + b^{old})}} \\ \vdots \\ \dfrac{-y_n}{1 + e^{y_n (\alpha^{old T} X x_n + b^{old})}} \end{bmatrix} + 2\lambda \alpha^{old} \right] \right]$$

$$\underbrace{\phantom{XXXXXXXXXXXXXXXXXXXXXXXXXXXXXX}}_{\alpha^{new}}$$

$$\Rightarrow \boxed{w^{new} = X^T \alpha^{new}}$$

where

$$\boxed{\alpha^{new} = \alpha^{old} - \eta \left[ \frac{1}{n} \begin{bmatrix} \dfrac{-y_1}{1 + e^{y_1 (\alpha^{old T} X x_1 + b^{old})}} \\ \vdots \\ \dfrac{-y_n}{1 + e^{y_n (\alpha^{old T} X x_n + b^{old})}} \end{bmatrix} + 2\lambda \alpha^{old} \right]} \quad \rightarrow \text{Update}$$

B1

(d)     From   B1(C)   we know   that

$$\alpha^{new} = \alpha^{old} - \eta \left[ \frac{1}{n} \left[ \begin{array}{c} \dfrac{-y_1}{1+e^{y_1(\alpha^{old T}Xx_1 + b^{old})}} \\ \vdots \\ \dfrac{-y_n}{1+e^{y_n(\alpha^{old T}Xx_n + b^{old})}} \end{array} \right] + 2\lambda \alpha^{old} \right] \qquad —①$$

Let us just examine a general term in this vector

$$\Rightarrow \quad * \quad \frac{-y_i}{1+e^{y_i(\alpha^{old T}Xx_i + b^{old})}}$$

→ we see that $Xx_i$ is an inner product

$$\text{(ie)} \quad Xx_i = \left[ \begin{array}{c} x_1^T x_i \\ x_2^T x_i \\ \vdots \\ x_n^T x_i \end{array} \right] \qquad —②$$

This inner product can be replaced by using a kernel.

$$Xx_i = \left[ \begin{array}{c} k(x_1, x_i) \\ k(x_2, x_i) \\ \vdots \\ k(x_n, x_i) \end{array} \right] = K(x, x_i) \qquad —③$$

I am denoting this column vector as $K(x, x_i)$

Now substituting the same, back in eqn ①

We get

$$\alpha^{new} = \alpha^{old} - \eta \left[ \frac{1}{n} \left[ \frac{y_i}{1 + e^{y_i (\alpha^{old T} K(X, x_i) + b^{old})}} \right]_{i=1....n} + 2\lambda \alpha^{old} \right]$$

also

This is $b^{old}$

$$b^{new} = b^{old} - \eta \sum_{i=1}^{n} \frac{-y_i}{1 + e^{y_i ((\alpha^{old})^T K(X, x_i) + b^{old})}}$$

Thus the update rule is kernelized.

We can set $\alpha^{old}$ to zero vector to start with and keep updating $\alpha^{new}, b^{new}$ till we end up with $(\alpha^*, b^*)$, the we know $w^* = X^T \alpha^*$

Final classifier,

$$\mathbb{P}(y=1 \mid x) = \frac{1}{1 + e^{-(\alpha^{*T} K(X, x_i) + b^*)}}$$

This is the probability that Class = 1,

and

$$\begin{cases} \mathbb{P}(y=1 \mid x) \geq \frac{1}{2} & \text{predict class} = 1 \\ \\ \mathbb{P}(y=1 \mid x) < \frac{1}{2} & \text{predict class} = -1 \end{cases}$$

This is $-1$

**Q 1(e)** Kernel Logistic Regression

Accuracy = 0.796

Error => 1 - 0.796 => 0.204

**Q 2(a)** Statistics

$$X\_mean = \begin{bmatrix} 0.50161345 & , & 0.45612671 & , & 0.3824407 \end{bmatrix}$$

$$X\_std = \begin{bmatrix} 0.24617303 & , & 0.23615181 & , & 0.23905821 \end{bmatrix}$$

**(1)**

- **Resizing:**

  ✓ Benefit — resizing is useful to have images with the same pixels in the data-set. The raw dataset might contain different images with different dimensions for example 1000×800 , 200×400 , etc, but resizing helps us to work with uniform dimensions through out data set.

  ✓ Draw back: By doing resizing, some important features are interpolated, and also, when ever we resize a smaller image, we would create a padding around, it which migh mislead our classifier, if there are several images with padding.

- **Normalizing:**

  ✓ Benefit: Normalizing the pixel intensities, will help the convergence of the optimizer, there wontbe much of zig-zagging during optimization. The centered data will keep the gradients in control when we use back propagation

Normalizing data

**Drawback** - Normalizing data can sometimes lead to loss of information. For example, Normalizing just preserves relative variations in data. If we want to classify day and night images; normalizing our data will affect the classification performance as it only captures relative information

## Q2(a)

(2) per-channel mean and standard deviation are use to scale our training data. Validation set is not a part of our training data, when we scale our training data using training per channel mean and per-channel standard deviation, it helps us center the data for training, but during validation and testing phase we should use this same scaling as the neural network that is trained is learnt for the information scaled this particular way. Hence, it doesnot make sense to use a different scaling metric for validation images.

## Q2(b)

(1) **Layer - 1**

Filter = $5 \times 5$

Input channel = 3

Output channel = 16

No. of parameters = $(5 \times 5 \times 3 + 1) \times 16$

This is for bias

$$= 1216 \quad \text{— Convolution layer-1}$$

**Layer - 2**

Filter = $5 \times 5$

Input channel = 16

Output channels = 32

No. of parameters = $(5 \times 5 \times 16 + 1) \times 32$

Bias

$$= 12832 \quad \text{— Convolution layer-2}$$

## Layer-3

Filter = 5×5

Input channel = 32

Output channel = 64

No. of parameters $= (5 \times 5 \times 32 + 1) \times 64$ — Bias

$= \boxed{51264}$ — Convolution Layer-3

## Layer-4

Filter = 5×5

Input channel = 64

Output Channels = 128

No. of parameters $= (5 \times 5 \times 64 + 1) \times 128$ — Bias

$= \boxed{204,928}$ — Convolution Layer-4

## Layer-5 (Fully connected)

Inputs = 128×2×2

Output = 64

No. of parameters $= (128 \times 2 \times 2 + 1) \times 64$ — Bias

$= \boxed{32,832}$ — fully connected Layer-1

## Layer-6 (Fully connected)

Inputs = 64

Outputs = 5

No. of parameters $= (64 + 1) \times 5$ — Bias

$= \boxed{325}$ — fully connected Layer-2

Total parameters $=$
$$\begin{array}{r} 1216 \; + \\ 12832 \; + \\ 51264 \; + \\ 204928 \; + \\ 32832 \; + \\ 325 \\ \hline 303397 \end{array}$$

Q2(b) Q. One reason as to why we don't want to initialize

our neutral network to zero is because, when intialized with o's

then all the layers perform the same calulation, there won't be

any symmetry breaking, the gradients computed take the learning

no, where, This is the reason why we should initialize the

weights
network, randomly.

Q2 (d)    Final accuracies and Losses
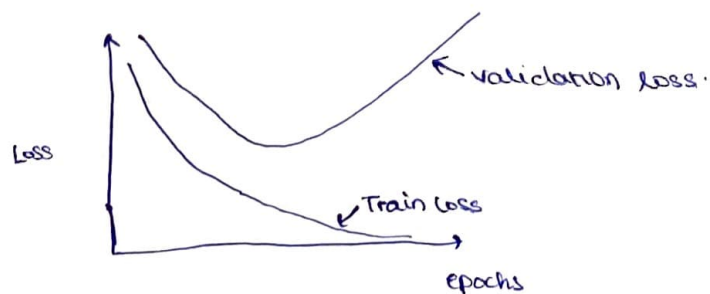
        Validation Loss:- 2.3131

        Validation accuracy :- 0.4872

        Train loss :- 0.007707

        Train accuracy :- 1

(e)    1) The training loss keeps decreasing, and the validation

loss drops at first and then starts increasing after a point.



As we keep on training the model

our training loss would be going to

zero, while our validation loss

will shoot up to higher values.

This means we have done an overfit to

our data.

Q2 (e)   a.   We should always try to minimize Validation loss no the training loss.

Based on the plot, we should stop training our model after 2 epochs as after 2 epochs the validation loss start increasing.

We should not try to maximize training accuracy because, it leads to overfitting ∮for the data, and the model learnt, iwon't be a good generalization of the true nature of data.


Q2 (f)   Q2(f) Final accuracies and losses:

Validation Loss:- 0.18
Validation accuracy = 0.935897
train loss : 0.11598
train accuray: 0.956


Q2 (g)   We can't use accuracy during an imbalance dataset case because, it is not fully reflective of the correct classification for each class. Por-class accuracy metric makes more sense in this scenario to asses how good our model is.

**8.2(h)** | **(Non-Weighted) Loss**

Imbalance dataset (5 epochs) | Weighted Cross entropy Loss

Imbalance dataset (5 epochs)

| | |
|---|---|
| Validation loss :- 0.1246 | Validation loss :- 0.12327 |
| Validation accuracy :- 0.9545 | Validation accuracy :- 0.963636 |
| train loss : 0.03430 | train loss :- 0.047591 |
| train accuracy : 0.990746 | train accuracy :- 0.9857142 |

Per class accuracy :-
$$\begin{bmatrix} Cat & dog \ (+ve \ class) \\ 1 & 0.5 \end{bmatrix}$$

$\left(\dfrac{dog+ve}{class}\right)$ $\begin{cases} precision = 1 \\ recall = 0.5 \\ F1\text{-}Score = 0.66666 \end{cases}$

Per class accuracy :-
$$\begin{bmatrix} Cat & dog \\ 0.98 & 0.8 \end{bmatrix}$$

$\begin{cases} Precision = 0.8 \\ Recall = 0.8 \\ F1 \ Scale = 0.8 \end{cases}$ $\left(\dfrac{dog \ is}{+ve \ class}\right)$

⌐ The un-weighted model have more training accuracy and train loss, but as the dataset is imbalanced we can see it performs poorly on other metrics like F-1 score, recall, per-class accuracy. What was happening here is that the model was being overfit for the features of Cat as the examples are more in number.

In the weighted case, it can be observed that weighting dog class by a factor improved the performance of the model, this can be seen from metrics like F1 score, precision, recall. Here due to weight, the model was able to distinguish features of dog from that of cat, though the examples are pretty less in number.