

CS 335 - CS 337 Assignment - 3.

1.1] ~~Q~~ Features used:-

a) Perimeter.

→ Idea is to use the perimeter of the shape to distinguish between shapes. Stars will have a ~~unusually~~^{very} high perimeter, as compared to others, while triangles will have a lower perimeter.

→ Method.

→ Our input is a reshaped 2500 dimensional vector of the image. This is basically all rows of the image concatenated one after the other.

→ A pixel will be on the perimeter if it has all zeros on one side and all ones on another side.

→ If we convolve a ~~1D~~ vector $[1, 1, 1, 1, 1]$ through the 2500 element input, and chose the elements with value 2, or 3, then we basically get all the vertical edges, ~~pixels~~^{pixels}.

→ Similarly, if we convolve the same mask on $X.reshape(50, 50).T.reshape(-1)$ where X is the input, we will get the number of horizontal edges pixels.

→ If we sum up these two numbers, then we have a rough estimate of the ~~parameter~~ perimeter.

b) Total Area

Idea

Idea for this feature is quite simple: all shapes will occupy different areas and this can be a good estimate to ~~get the~~ ~~total~~ differentiate between classes.

Method

Simply add the input image pixels.

c) Ring Area: (Using just Ring Area, without other features gives 100% Test Accuracy)

Idea

This, according to me, is the best feature that differentiates between classes.

Consider a ring (donut) with ~~at~~ centre (25, 25) and outer radius = max distance of any pixel from centre, ~~inner radius~~ ~~set to outer radius - 6~~, and inner radius = $\max(-4, \text{outer radius} - 6)$

If we consider the white area inside this ring (donut), then it can, by itself, give 100% Test Accuracy.

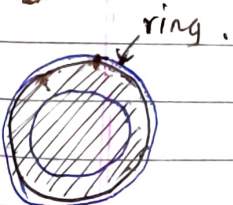
Illustration:-

Circle

Square

Triangle

Star



No white area in ring.



some white area in ring.



Huge ^{white} area in ring.



Lower area in the ring.

Q-2.1.

(a) For a binary logistic regression, our labels were 0, 1. So, let's look at the one hot encoding of these labels.

$$y_{(0)} = [1, 0]$$

$$y_{(1)} = [0, 1]$$

Hence $y_{\text{one hot}} = [(1-y) \ y]$, for any image sample in data set.

So, if we ~~app~~ have the following formula for multi class LR.

$$E(w) = -\frac{1}{N} \sum_{i=1}^N \sum_{\kappa=1}^K y_{\kappa}^{(i)} \log(P(Y=\kappa | w_{\kappa}, \phi(x^{(i)}))).$$

For binary, we have $\kappa=2$, which gives us.

$$E(w) = -\frac{1}{N} \sum_{i=1}^N \left\{ \begin{aligned} &[(1-y^{(i)}) \log(P(Y=0 | w_{00}, \phi(x^{(i)})))] \\ &+ [y^{(i)} \cdot \log(P(Y=1 | w_{01}, \phi(x^{(i)})))] \end{aligned} \right\}$$

In the slides, eqn 3, $P(Y=1 | w_1, \phi(x^{(i)})) = \sigma_w(x^{(i)})$.

Hence, $P(Y=0 | w_0, \phi(x^{(i)})) = 1 - \sigma_w(x^{(i)})$.

Hence we have

$$E(w) = -\frac{1}{N} \sum_{i=1}^N \left(y^{(i)} \sigma_w(x^{(i)}) + (1-y^{(i)}) \right)$$

Hence we have.

$$E(w) = \frac{-1}{N} \left\{ \sum_{i=1}^N \left\{ y^{(i)} \log(\sigma_w(x^{(i)})) + (1 - y^{(i)}) \log(1 - \sigma_w(x^{(i)})) \right\} \right\}$$

Hence proved, binary L.R is a special case of Multiclass L.R.

b) Define $Z_{N \times K} = \phi(x)_{N \times D} \cdot W_{D \times K}$.

Each element of Z is $z_{ik} = \phi(x^{(i)})_{1 \times D} \cdot W_{D \times 1}^k$ (scalar).

We need $\frac{dE}{dW}$

Any ~~element~~ ^{column element} of this matrix,

$$\left(\frac{dE}{dW} \right)_{D \times K} = \frac{dE_{k.}}{dW_{k.}} = \frac{dE_{k.}}{dz} \frac{dz}{dW_{k.}}$$

→ To Find $\frac{dE}{dz}$

$$P(Y=k | \phi(x) \cdot w_k) = \frac{e^{z_{ki}}}{\sum_{k=1}^K e^{z_{ki}}}$$

$$\log(P(Y=k | \phi(x) \cdot w_k)) = z_{ki} - \log\left(\sum_{k=1}^K e^{z_{ki}}\right)$$

b). Define $Z_{N \times K} = \phi(x)_{N \times D} \cdot W_{D \times K}$

$Z_{i,k}$ = each element of $Z_{N \times K} = \phi(x^{(i)})_{1 \times D} \cdot W_{D \times K}^k$

We have

$$E = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_k^{(i)} \log(P(y=k | w_k, \phi(x^{(i)})))$$

$$= -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_k^{(i)} \left\{ Z_{ik} - \log \left(\sum_k e^{Z_{ik}} \right) \right\}$$

$$= -\frac{1}{N} \left\{ \sum_{i=1}^N \sum_{k=1}^K y_k^{(i)} Z_{ik} - \sum_i \left(\sum_k y_k^{(i)} \right) \log \left(\sum_k e^{Z_{ik}} \right) \right\}$$

Since $\sum_k y_{ik} = 1$
as it is one-hot
encoding
class

$$= -\frac{1}{N} \left\{ \sum_i \sum_k y_{ik} Z_{ik} - \sum_i \log \left(\sum_k e^{Z_{ik}} \right) \right\} \quad \text{--- (1)}$$

So, we want $\frac{dE}{dw_{d,k}}$ to the

The $(d,k)^{th}$ element of this Matrix will be $\frac{dE}{dw_{d,k}}$

$$\frac{dE}{dw_{d,k}} = \left(\sum_i \frac{dE}{dz_{ik}} \right) \cdot \frac{dz_{ik}}{dw_{d,k}}$$

$$\rightarrow \frac{dz_{ik}}{dw_{d,k}} = \sum_i \phi(x^{(i)})_{d,k}$$

→ For ~~dzik~~ $\frac{dz_{ik}}{dw_{dk}}$,

$$z_{ik} = \sum_d w_{dk} \cdot \phi(x^{(i)})_d$$

$$\frac{dz_{ik}}{dw_{dk}} = \phi(x^{(i)})_d = \phi_{i,d} \quad \text{--- (2)}$$

→ For $\frac{dE}{dz_{ik}}$,

Using eqⁿ (1)

$$\frac{dE}{dz_{ik}} = -\frac{1}{N} \left(y_{ik} - \frac{e^{z_{ik}}}{\sum_k e^{z_{ik}}} \right)$$

~~$\Rightarrow \frac{dE}{dw_{dk}} = \frac{1}{N} \phi_{i,d} (y_{ik} - \frac{e^{z_{ik}}}{\sum_k e^{z_{ik}}})$~~

$$\Rightarrow \frac{dE}{dw_{dk}} = -\frac{1}{N} \sum_i \phi_{i,d} \left(y_{ik} - \frac{e^{z_{ik}}}{\sum_k e^{z_{ik}}} \right)$$

Vectorized:-

$$\vec{\nabla}_{w_{dk}} E = -\frac{1}{N} \phi(x)^T_{D \times N} \left(Y_{N \times K} - \text{Softmax}(Z, \text{dim}=1) \right)$$

Here $\text{Softmax}(Z, \text{dim} = 1)$
 $= \frac{\text{np.exp}(Z)}{\text{np.sum}(Z.\text{exp}(), 1)}$

\Rightarrow Vectorized:

$$\vec{\nabla}_W E_{D \times K} = -\frac{1}{N} \Phi(X)_{D \times N}^T \left(Y_{N \times K} - \text{Softmax} \left(\Phi(X)_{N \times D} \cdot W_{D \times K}, \text{dim} = 1 \right) \right)$$

Q-2.2.

b). Test Accuracy = 86.33 %

For Model M (all zeros)

\hookrightarrow test accuracy = 84.18 %

Here Accuracy is not a good metric to gauge the model performance. This is because the data is heavily imbalanced. The number of negative samples are much more than the positive samples, and hence accuracy will give a high value even for such a bad model.

c). F_1 score = 0.301

F_1 score for Model M = 0.

I believe F_1 score is a good metric as compared to Accuracy. In cases where data is non uniform, with low numbers of positives, like Cancer detector / Fraud Detector, accuracy will give high ~~accuracy~~ value even for bad models.

But F_1 score, on the other hand, considers False Negatives, False positives as well, and hence will give a better estimate in such cases, ~~whether~~
of whether our model is good or not.

Also, most real life datasets are not uniform, and hence, F_1 score is better than Accuracy.

2.2) c) Test Accuracy for multiclass LR = 84.43 %
Test Accuracy for multiclass perceptron = 78.28 %

Logistic Regression achieves a higher test accuracy. The reason for this is, although both the models try to exploit the linear separability in the data, ~~there actually is~~ linear regression (since it uses sigmoid/softmax) has a continuous range of outputs, whereas, the perceptron gives a binary Yes/NO output.

Also, while training a perceptron, the gradient from an image is not considered even if it is 'marginally' correctly classified, whereas in case of Logistic Regression, the gradient tries to adjust the ~~para~~ model even for correctly classified images, in order to maximise their classification score (probability). This leads to a better separating plane (~~more~~ closer to the ideal optima) in case of Logistic Regression as compared to perceptron.

Hence, the accuracy of Logistic Regression is better than that of a perceptron.