

# CS 335-337: Assignment 3

Anish Deshpande - 180100013

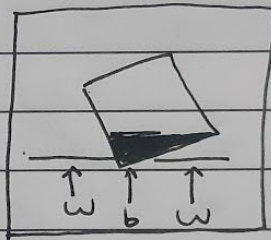
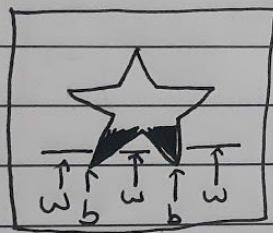
## Assignment 3

Date : \_\_\_\_\_

1.1

b) I have used 5 features to classify the shapes.

i) The first is the number of transitions from white to black or black to white in any given row. This number will be at most 2 for any shape except a star, and for a star, there will be a few rows where there are 4 transitions.

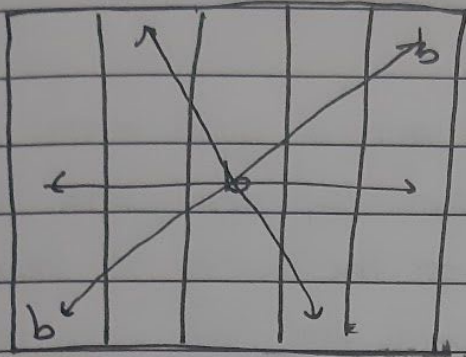


ii) The next 4 features are all extracted from information from  $10 \times 10$  patches surrounding each point.

\* We take patches whose centre is on the shape.

\* We want it to be a corner/edge, so only let 30/100 pixels be black, this excludes the interior of a shape.

\* We distinguish between a corner and edge by taking a smaller  $5 \times 5$  sub-patch, and seeing whether all points on some straight line are black:



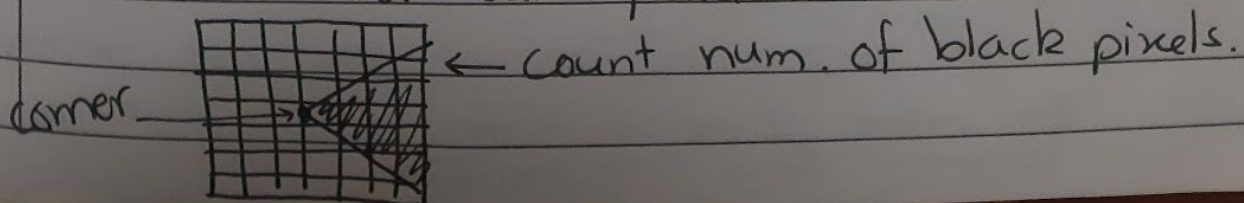
If all points in a straight line are black, on any one line, then increment the 'on-edge' counter.

∴ \* The number of such 'on-edge' pixels is a feature.

After determining that it is not an edge point, we calculate the number of black pixels in the  $10 \times 10$  patch.

The trend:  $\text{num}(\text{black})$ : star < triangle < square. circles are set to 0 because of the patch condition.

\* The min, max, and average of these counts are our features.





2.1

a) Categorical cross-entropy loss:

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

$$(P(Y=k | W_k, \phi(x^{(i)})))^x = \sigma_{W_k}(x^i) = \left( \frac{1}{1 + e^{-W_k^T \phi(x^i)}} \right)$$

In binary logistic regression,  $k = \{1, 2\}$   
(only 2 classes). equivalent to 0,1

$$P(Y=k | W_k, \phi(x^i)) = \frac{e^{W_k^T \phi(x^i)}}{\sum_{k=1}^K e^{W_k^T \phi(x^i)}} \quad \forall k \in 1 \dots K$$

In the special case, we only have 2 categories  
 $P(Y=1)$  and  $P(Y=0)$

Hence, we may write  $P(Y=0)$  as  $1 - P(Y=1)$ .

also, say the class labels are 0 and 1

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

$\therefore y^{(i)}$  is 1 if class label is 1

$(1 - y^{(i)})$  is 1 if class label is 0

and vice versa. ( $y_i^0 = 1 - y_i^1 \quad \forall i$ )

$\therefore$   ~~$E(W) = -\frac{1}{N} \sum$~~  However, taking:

$$P(Y=k | W_k, \phi(x^i)) = \frac{e^{W_k^T \phi(x^i)}}{1 + \sum_{k=1}^{K-1} e^{W_k^T \phi(x^i)}} \quad \forall k \in 1 \dots K-1$$

$$\text{And } P(Y=K | w_K, \phi(x^i)) = \frac{1}{1 + \sum_{k=1}^{K-1} e^{w_k^T \phi(x^i)}} \quad (k=K)$$

This reduces the number of weight vectors/variables by 1, it removes redundancy.

$$\therefore \text{We have } P(Y=0 | w_0, \phi(x^i)) = \frac{e^{w_0^T \phi(x^i)}}{1 + \sum_{k=0}^0 e^{w_k^T \phi(x^i)}}$$

(K runs from 0 to 0, as K=1 (class labels 0,1))

$$\therefore P(Y=0) = \frac{e^{w_0^T \phi(x^i)}}{1 + e^{w_0^T \phi(x^i)}} \quad \text{and } P(Y=1) = \frac{1}{1 + e^{w_0^T \phi(x^i)}}$$

$$\therefore E(W) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=0}^1 y_k^{(i)} \log(P(Y=k | w_K, \phi(x^i)))$$

$$E(W) = -\frac{1}{N} \sum_{i=1}^N \left[ y_0^{(i)} \log\left(\frac{e^{w_0^T \phi(x^i)}}{1 + e^{w_0^T \phi(x^i)}}\right) + y_1^{(i)} \log\left(\frac{1}{1 + e^{w_0^T \phi(x^i)}}\right) \right]$$

$$E(W) = -\frac{1}{N} \sum_{i=1}^N \left[ (1 - y_1^{(i)}) \log\left(\frac{1 + e^{w_0^T \phi(x^i)} - 1}{1 + e^{w_0^T \phi(x^i)}}\right) + y_1^{(i)} \log\left(\frac{1}{1 + e^{w_0^T \phi(x^i)}}\right) \right]$$

$$\therefore \sigma_{w_0}(x^{(i)}) = \frac{1}{1 + e^{w_0^T \phi(x^i)}}$$



$$\therefore E(w) = -\frac{1}{N} \sum_{i=1}^N y_i^{(i)} \log \sigma_{w_0}(x_i) + (1 - y_i^{(i)}) (\log(1 - \sigma_{w_0}(x_i)))$$

Generalising  $y_i$  to  $y$  and  $w_0$  to  $w$

$$E(w) = -\frac{1}{N} \sum_{i=1}^N y_i^{(i)} \log \sigma_w(x_i) + (1 - y_i^{(i)}) \log(1 - \sigma_w(x_i))$$

This is the required expression.

Hence proved.

$$b) E(w) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{ik}^{(i)} \times \log(P(Y=k | W_k, \phi(x^{(i)})))$$

$$\text{to find } \frac{\partial E(w)}{\partial w} = \frac{\partial E(w)}{\partial z} \times \frac{\partial z}{\partial w} \quad (\text{chain rule})$$

$$\text{Let } z = \phi(x) \cdot w \quad (N \times K)$$

$\phi(x)$  is constant w.r.t  $w$

$$\therefore \frac{\partial z}{\partial w} = \phi(x)^T$$

$$E(w) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{ik}^{(i)} \log \left( \frac{e^{z_{ik}}}{\sum_{p=1}^K e^{z_{ip}}} \right)$$

$$E(w) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{ik}^{(i)} \left[ z_{ik} - \log \left( \sum_{p=1}^K e^{z_{ip}} \right) \right]$$

Date : \_\_\_\_\_

$$\therefore \frac{\partial E(w)}{\partial z_{ik}} = -\frac{1}{N} \left( y_k^{(i)} \left[ 1 - \frac{1}{\sum_{p=1}^K e^{z_{ip}}} \times e^{z_{ik}} \right] \right)$$

$$\frac{\partial E(w)}{\partial z_{ik}} = -\frac{1}{N} \left( y_k^{(i)} - \frac{y_k^{(i)} e^{z_{ik}}}{\sum_{p=1}^K e^{z_{ip}}} \right)$$

Say  $Y' = N \times K$  matrix of  $y_k^{(i)}$  s.

$$\frac{\partial E(w)}{\partial z} = -\frac{1}{N} \left( Y' - \text{softmax}(\text{weights}^{(w)}, \phi(x)) \right)$$

$$\frac{\partial E(w)}{\partial w} = \left( \frac{1}{N} \right)^* \phi(x)^T \left( \sigma_w(\phi(x)) - Y \right)$$

Where  $\sigma_w(\phi(x)) = \text{softmax} = \frac{e^{\phi(x)w}}{\sum_{\text{columns}} e^{\phi(x)w}}$



2.2

b) For binary logistic regression,

test accuracy obtained = ~~0.8525~~ 0.863

Model 'M', which predicts '0' all the time,

test accuracy = ~~0.86~~ 0.84

Clearly, accuracy is not a good metric. This is because of the nature of the data. It is extremely unlikely for the song to have made it to the Billboard Hot 100, so any model can cheat and just predict that it won't make it, for high accuracy.

c) For our model, F1 score

= 0.301 (excluding artist name)

= 0.522 (including artist name, preferred).

Consider the model M, which predicts 0 all the time.

$$F1 \text{ score} = \frac{2TP}{2TP + FN + FP}$$

$$\therefore TP = 0 \text{ and } FN > 0$$

$$F1 \text{ score} = 0$$

Clearly, F1 score is a great indicator of the model's performance, as it takes into consideration the extent of misclassification also: (FN, FP) terms, comparing those with the TPs.

An imbalanced class-distribution exists here, so false negatives and false positives are very important. So, F1 score is a good metric for model evaluation.

- e) Perceptron accuracy on D1  $\approx 0.78$   
Logistic regression on D1, accuracy  $\approx 0.84$

Logistic regression outperforms the perceptron.

They are both hyperplanes, but because of the non-linear nature of the logistic curve, it can fit better to real-world situations.

A linear model is more often than not unable to separate and classify actual data, like the fact that the relationship between the properties of a song and whether it makes it to the billboard top 100 is complex.