

# ML Assignment 2

Anshul Mendiratta  
2018219

## Question 1

(a) Models built on datasets with high dimensionality(having a large number of features) tend to overfit. Further, the lower the dimensions the faster we are able to fit the training data. Hence, in datasets like these we want to reduce the number of features. PCA helps us in doing so. PCA is a linear unsupervised learning algorithm which helps us construct a set of uncorrelated features from a set of correlated features. We define a set of axis using a linear combination of the features and rank them according to the amount of variance.

We want the highest ranked axis to have maximum variability and hence the ability to capture the most information. All other subsequent axis should capture the remaining variability while not being correlated with the other axis. We will extract the most important axis from the ones obtained and treat them as our new set of features.

(b) As seen from above, datasets with large number of features aren't desirable. SVD is a way of reducing these features using matrix decomposition.

Suppose we have a matrix A, which represents the set of features and their values for each datapoint. Using SVD, we can decompose A into,

$$A = U * \Sigma * V^T$$

Where A is our  $m \times n$  original matrix, U is an  $m \times m$  matrix of Left singular vectors,  $\Sigma$  is an  $m \times n$  diagonal matrix and V is an  $n \times n$  matrix of right singular vectors. Using these matrices we can represent the original matrix as a linear combination of low rank matrices. We can then pick the first few singular vectors and truncate the matrices accordingly. Hence, we have reduced our features.

(c) t-SNE (t-Distributed Stochastic Neighbor Embedding) is a non-linear algorithm which helps in reducing the dimension of a dataset. Contrary to PCA, which places emphasis on placing dissimilar data points far apart, t-SNE also focuses on representing similar ones closer. Furthermore, unlike PCA, t-SNE is based on probability distributions to find relationships between features. There are mainly 3 steps -

1. We measure similarities between points in the high dimensional space.
2. Then we try to measure similarities between points in the low dimensional space.
3. Finally, we want to map both of these structures and optimize the similarity measures using a cost function.

(d) Stratified sampling is a way of sampling the dataset such that the training and test data after splitting have approximately the same percentages of each target class. As expected, the class frequency of both test and train are similar. (Code in python file)

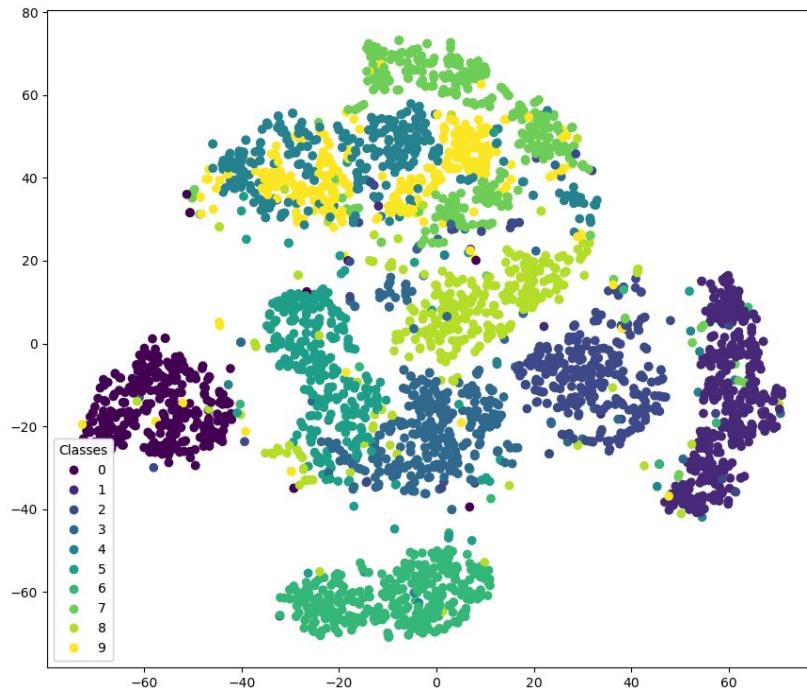
Class frequencies (Class vs frequency percentage) of train and test are -

```
{0: 9.523809523809524, 1: 11.755952380952381, 2: 9.345238095238095,  
3: 10.089285714285715, 4: 9.910714285714285, 5: 9.464285714285714,  
6: 10.505952380952381, 7: 10.267857142857142, 8: 9.761904761904763,  
9: 9.375}
```

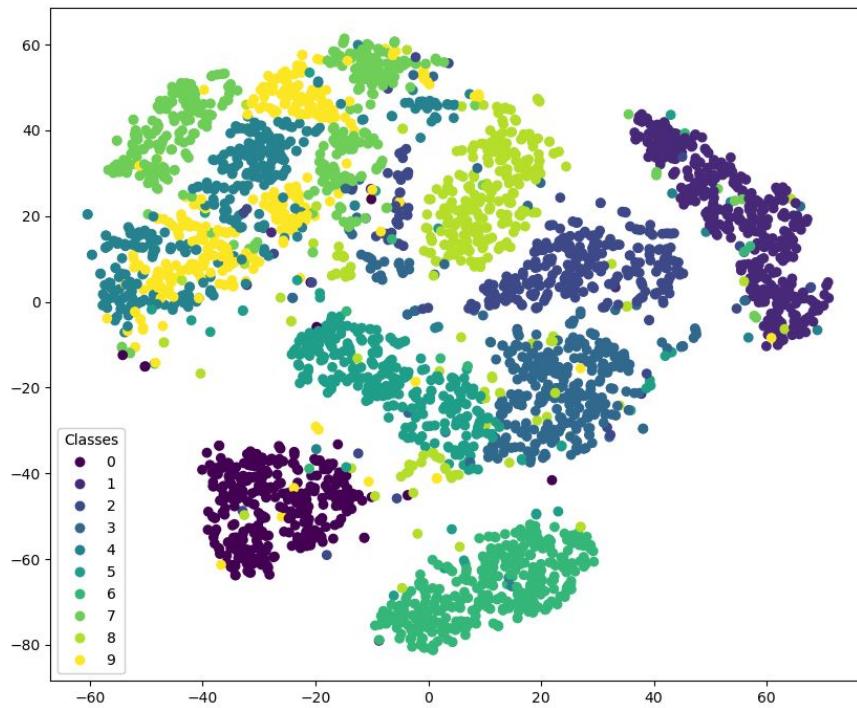
```
{0: 9.523809523809524, 1: 11.785714285714285, 2: 9.404761904761905,  
3: 10.119047619047619, 4: 9.880952380952381, 5: 9.523809523809524,  
6: 10.476190476190476, 7: 10.238095238095237, 8:  
9.761904761904763, 9: 9.285714285714286}
```

(e) Using PCA and training Logistic Regression Model we get accuracy of 0.883. We then apply t-SNE on the training data obtained after applying PCA and plot a scatter plot, where x and y axis are the features obtained by t-SNE and each point is colored according to the class it belongs to.

We get the following plot -



(f) Using SCD and training Logistic Regression Model we get accuracy of 0.878. We then apply t-SNE and make a graph following same steps as above. We get the following plot -



(g) We get very similar accuracies for SVD and PCA, even though PCA performs negligibly better. We do see from the graphs obtained that the features extracted from SVD and PCA are different as the clusters of classes are not at the same place in both the graphs. Hence even though different features are being extracted in SVD and PCA, they perform well in both cases.

## Question 2

We apply bootstrapping in the following way -

- 1) First we create a holdout from X. In this case we kept 20% of X as holdout. The rest is used for creating samples and we call this testX.
- 2) Next we create B samples each of size B.
- 3) We train our model on each sample  $S_b$  and compute values for each  $x$  in holdout. Hence we have B predicted values for each  $x$  in holdout.
- 4) Now we calculate bootstrap estimate, bias and variance for every point  $x$  in testX
- 5) We then average these values over testX and report them
- 6) Similarly to find MSE, we find MSE at each  $x$  in testX and calculate  $MSE - bias^2 - variance$  at every point  $x$  and report its mean.

(a) The bias and variance of the model came out to be - 0.169 and - 0.53 respectively.

(b)  $MSE - bias^2 - variance$  came out to be = -0.23

We know  $MSE = bias^2 + variance + irreducible\ uncertainty$

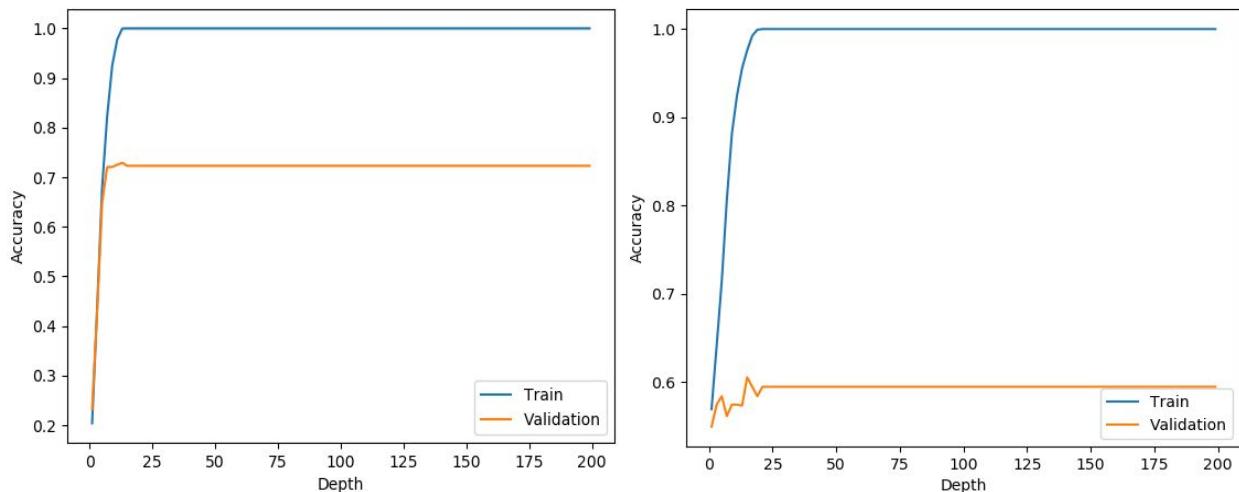
Hence we can see that our bootstrap model satisfies the equation as -0.23 is a comparatively smaller value and can be attributed to irreducible uncertainty

### Question 3

(a)

1. Dataset A
  - a. Decision Tree - Applying grid search and K fold we found the average accuracy (over all the folds) for depths from 1 to 25. We got the best accuracy at depth 16.
  - b. Gaussian Naive Bayes - Using K fold we found the best accuracy of 0.619 was achieved at fold 4. Further average K fold accuracy was 0.59
2. Dataset B
  - c. Decision Tree - Applying grid search and K fold we found the average accuracy (over all the folds) for depths from 1 to 25. We got the best accuracy at depth 20.
  - d. Gaussian Naive Bayes - Using K fold we found the best accuracy of 0.595 was achieved at fold 2. Further average K fold accuracy was 0.569

(b) Plots for training and validation accuracy vs tree depth for dataset A and B are (Left is A and right is B) -



(c) Implemented in code and saved in weights folder

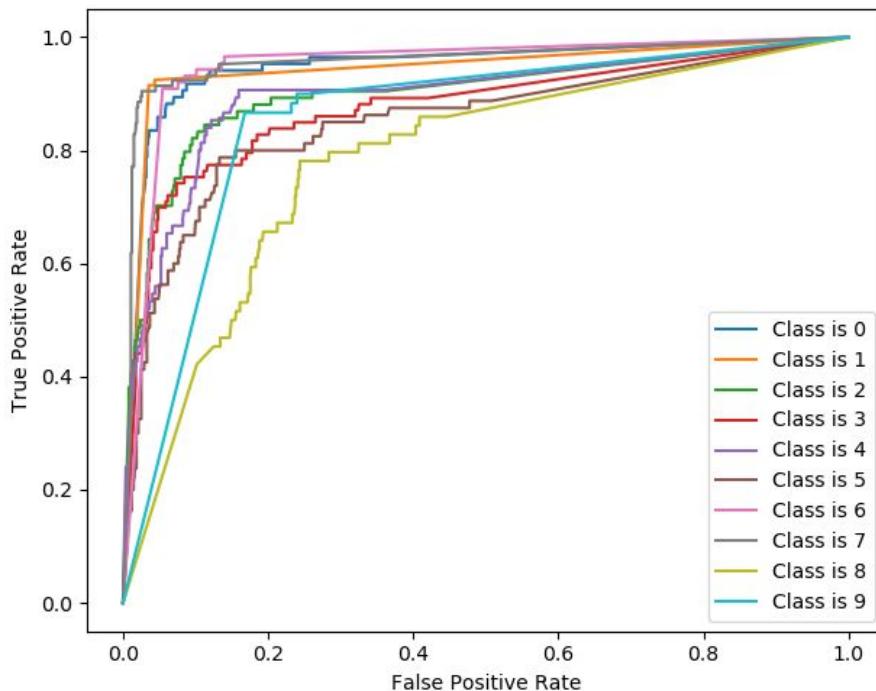
(d)

1. For dataset A

a. Gaussian Naive Bayes

```
Confusion Matrix -  
[[63, 0, 7, 1, 4, 6, 2, 1, 1, 1],  
 [0, 97, 1, 6, 0, 6, 3, 3, 6, 1],  
 [2, 0, 31, 1, 1, 0, 1, 0, 1, 0],  
 [0, 2, 6, 40, 0, 4, 0, 0, 1, 0],  
 [0, 1, 1, 18, 1, 0, 2, 0, 0, 0],  
 [1, 2, 1, 0, 0, 13, 1, 0, 2, 0],  
 [4, 2, 18, 10, 6, 5, 80, 0, 2, 0],  
 [0, 0, 0, 1, 1, 0, 39, 2, 3],  
 [9, 0, 17, 20, 10, 36, 1, 1, 29, 3],  
 [6, 2, 2, 13, 35, 8, 0, 59, 20, 52]]  
Accuracy 0.55  
Macro precision 0.6467557459590406  
Precision for every class {0: 0.7325581395348837, 1: 0.7886178861788617, 2: 0.8378378378378378, 3: 0.7547169811320755, 4: 0.75, 5: 0.65, 6: 0.6299212598425197, 7: 0.8297872340425532, 8: 0.23015873015873015, 9: 0.2639593908629442}  
Macro Recall 0.5458237103326363  
Recall for every class {0: 0.7411764705882353, 1: 0.9150943396226415, 2: 0.36904761904761907, 3: 0.43010752688172044, 4: 0.24, 5: 0.1625, 6: 0.9090909090909091, 7: 0.37142857142857144, 8: 0.453125, 9: 0.8666666666666667}  
Macro F1 score 0.5235258300489847  
F1 score for every class {0: 0.736842105263158, 1: 0.8471615720524018, 2: 0.5123966942148761, 3: 0.5479452054794521, 4: 0.3636363636363636, 5: 0.26, 6: 0.7441860465116279, 7: 0.5131578947368421, 8: 0.30526315789473685, 9: 0.40466926070038917}  
Micro Precision 0.55  
Micro Recall 0.55  
Micro F1Score 0.55
```

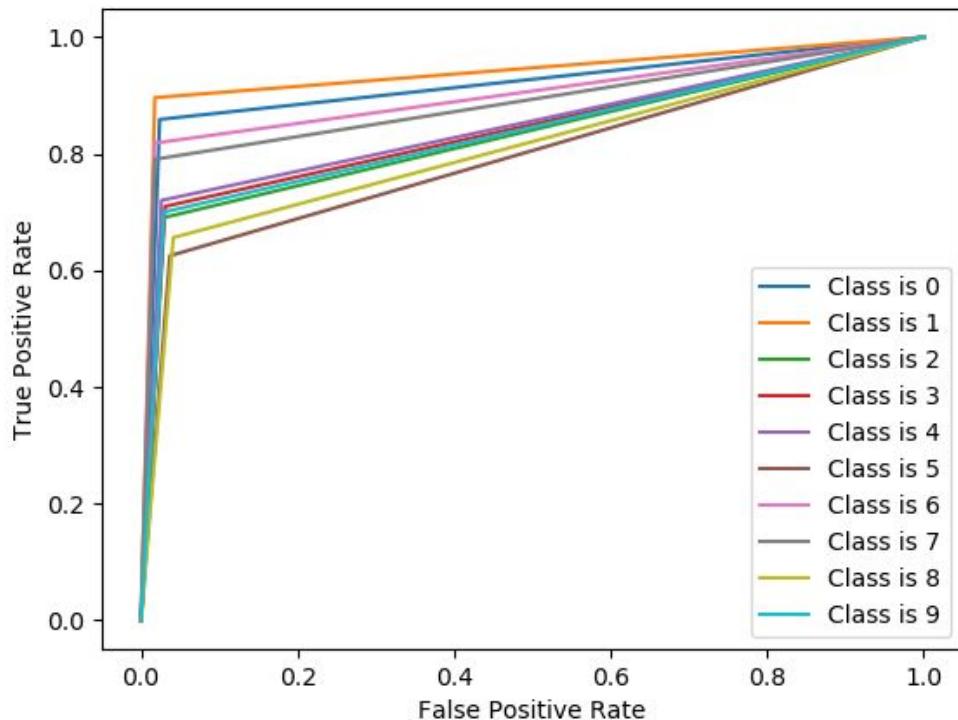
ROC Plot -



## b. Decision Tree

```
Confusion Matrix -  
[[73, 0, 1, 2, 0, 7, 1, 3, 3, 1],  
 [0, 95, 3, 3, 1, 2, 1, 2, 0, 1],  
 [0, 2, 58, 4, 3, 2, 5, 2, 4, 0],  
 [3, 0, 5, 66, 1, 3, 0, 3, 6, 2],  
 [1, 3, 1, 1, 54, 1, 1, 2, 2, 8],  
 [3, 0, 5, 8, 2, 50, 4, 2, 3, 1],  
 [0, 0, 2, 0, 1, 9, 72, 0, 1, 0],  
 [0, 1, 5, 1, 1, 1, 0, 83, 0, 4],  
 [2, 5, 4, 6, 5, 3, 2, 4, 42, 1],  
 [3, 0, 0, 2, 7, 2, 2, 4, 3, 42]]  
Accuracy 0.7559523809523809  
Macro precision 0.7444519406874827  
Precision for every class {0: 0.8021978021978022, 1: 0.8796296296296297, 2: 0.  
725, 3: 0.7415730337078652, 4: 0.7297297297297297, 5: 0.6410256410256411, 6: 0.  
8470588235294118, 7: 0.8645833333333334, 8: 0.5675675675675675, 9: 0.6461538461  
538462}  
Macro Recall 0.7465111562995143  
Recall for every class {0: 0.8588235294117647, 1: 0.8962264150943396, 2: 0.6904  
761904761905, 3: 0.7096774193548387, 4: 0.72, 5: 0.625, 6: 0.81818181818182,  
7: 0.7904761904761904, 8: 0.65625, 9: 0.7}  
Macro F1 score 0.744666756858741  
F1 score for every class {0: 0.8295454545454546, 1: 0.8878504672897196, 2: 0.70  
73170731707318, 3: 0.7252747252747254, 4: 0.7248322147651007, 5: 0.632911392405  
0633, 6: 0.8323699421965318, 7: 0.8258706467661692, 8: 0.6086956521739131, 9: 0  
.6719999999999999}  
Micro Precision 0.7559523809523809  
Micro Recall 0.7559523809523809  
Micro F1Score 0.7559523809523809
```

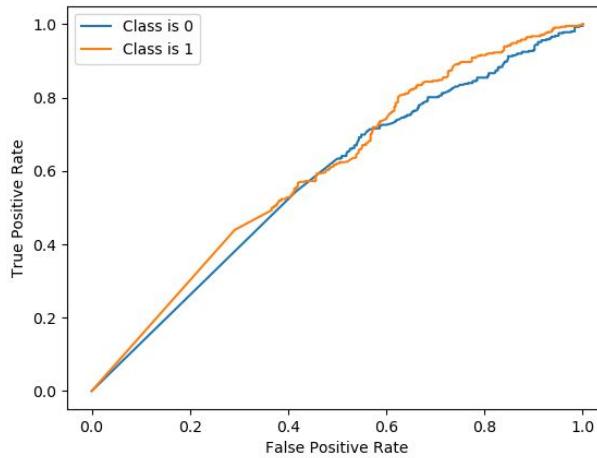
ROC Curve -



## 2. For Dataset B

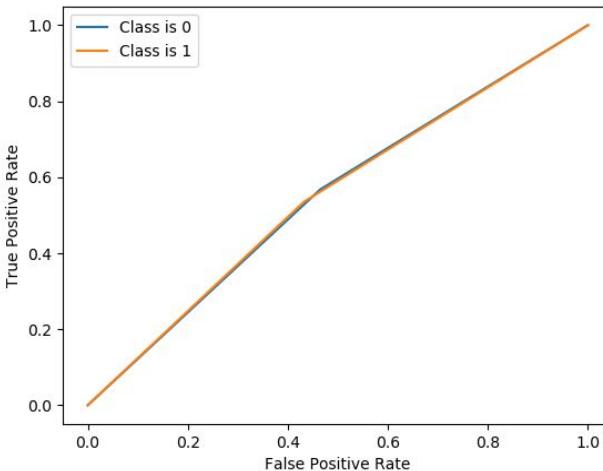
### a. Gaussian

```
Confusion Matrix -  
[261, 215]  
[151, 213]  
Accuracy 0.5642857142857143  
Macro precision 0.5667420814479638  
Precision for every class {0: 0.5483193277310925, 1: 0.5851648351648352}  
Macro Recall 0.5655793485164686  
Recall for every class {0: 0.633495145631068, 1: 0.4976635514018692}  
Macro F1 score 0.5628583128583129  
F1 score for every class {0: 0.5878378378378378, 1: 0.537878787878788}
```



### b. Decision Tree

```
Confusion Matrix -  
[234, 199]  
[178, 229]  
Accuracy 0.5511904761904762  
Macro precision 0.5515346335207767  
Precision for every class {0: 0.5404157043879908, 1: 0.5626535626535627}  
Macro Recall 0.5515039470102532  
Recall for every class {0: 0.5679611650485437, 1: 0.5350467289719626}  
Macro F1 score 0.551174573929065  
F1 score for every class {0: 0.5538461538461538, 1: 0.5485029940119761}
```



#### **Question 4**

##### **1. Dataset 1**

Accuracy of my implementation - 0.53

Accuracy of skLearn's implementation - 0.55

##### **2. Dataset 2**

Accuracy of my implementation - 0.55

Accuracy of skLearn's implementation - 0.56

**Rest of the Questions have been attached below**

Q.5

a)

Initial Entropy,

$$H(Y) = - \sum_{i=1}^n P(Y=y_i) \log_2 P(Y=y_i)$$

$$= - \left( \frac{5}{14} \log_2 \frac{5}{14} + \frac{9}{14} \log_2 \frac{9}{14} \right)$$

$$= -(-0.9402) \approx 0.94$$

Now we find information gain for each parameter,

$$IG(x) = H(Y) - H(Y|X)$$

$$\begin{aligned} \therefore IG(\text{Outlook}) &= 0.94 - \left[ - \left( \frac{5}{14} \left( \frac{3}{5} \log \frac{3}{5} + \frac{2}{5} \log \frac{2}{5} \right) + \right. \right. \\ &\quad \left. \left. \frac{4}{14} \left( \frac{4}{6} \log \frac{4}{6} + \frac{2}{6} \log \frac{2}{6} \right) + \frac{5}{14} \left( \frac{3}{5} \log \frac{3}{5} + \frac{2}{5} \log \frac{2}{5} \right) \right) \right] \\ &= 0.94 + \left[ \frac{10}{14} \left( \frac{3}{5} \log \frac{3}{5} + \frac{2}{5} \log \frac{2}{5} \right) \right] \\ &= 0.94 - 0.6935 \approx 0.24 \end{aligned}$$

$$IG(\text{Temperature}) = 0.94 + \left[ \frac{4}{14} \left( \frac{2}{6} \log \frac{2}{6} + \frac{4}{6} \log \frac{4}{6} \right) + \right.$$

$$\left. \frac{6}{14} \left( \frac{4}{6} \log \frac{4}{6} + \frac{2}{6} \log \frac{2}{6} \right) + \frac{4}{14} \left( \frac{3}{5} \log \frac{3}{5} + \frac{1}{5} \log \frac{1}{5} \right) \right]$$

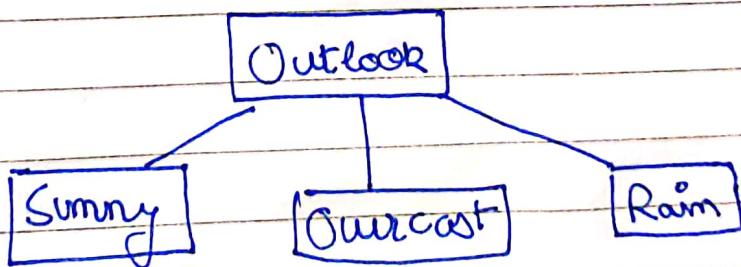
$$= 0.94 + \left( -\frac{4}{14} - 0.39 - 0.23 \right)$$

$$\begin{aligned} &= 0.94 + (-0.96) \\ &= 0.04 \end{aligned}$$

$$\begin{aligned}
 I_G(\text{Humidity}) &= 0.94 + \left[ \frac{7}{14} \left( \frac{4}{7} \log \frac{4}{7} + \frac{3}{7} \log \frac{3}{7} \right) + \frac{7}{14} \left( \frac{6}{7} \cdot \right. \right. \\
 &\quad \left. \left. \log \frac{6}{7} + \frac{1}{7} \log \frac{1}{7} \right) \right] \\
 &= 0.94 + [-0.49 + (-0.29)] \\
 &= 0.94 - 0.78 \\
 &= 0.16
 \end{aligned}$$

$$\begin{aligned}
 I_G(\text{Wind}) &= 0.94 + \left[ \frac{8}{14} \left( \frac{2}{8} \log \frac{2}{8} + \frac{6}{8} \log \frac{6}{8} \right) + \right. \\
 &\quad \left. \frac{6}{14} \left( \frac{3}{6} \log \frac{3}{6} + \frac{3}{6} \log \frac{3}{6} \right) \right] \\
 &= 0.94 + (-0.463 + (-0.423)) \\
 &= 0.05
 \end{aligned}$$

We get max information gain at Outlook, so  
first split is at outlook



~~Now we do~~ Now we do ~~outlook~~  $\Rightarrow$   $\emptyset$

Now we find Entropy (Play Match | sunny)

$$= -\left(\frac{3}{5} \log \frac{3}{5} + \frac{2}{5} \log \frac{2}{5}\right)$$

$$= 0.97$$

We then find corresponding  $I_G$ .

$$I_G(\text{Temp, Sunny}) = 0.97 + \left[ \frac{2}{5} \left( \frac{2}{2} \log \frac{2}{2} \right) + \right.$$

$$\left. \frac{2}{5} \left[ \frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} \right] + \right. \\ \left. \frac{1}{5} \left( + \log \frac{1}{1} \right) \right]$$

$$= 0.97 + (-0.4) = 0.57$$

$$I_G(\text{Humidity, Sunny}) = 0.97 + \left[ \frac{3}{5} \left( \frac{3}{3} \log \frac{3}{3} \right) + \right.$$

$$\left. \frac{2}{5} \left( \frac{2}{2} \log \frac{2}{2} \right) \right]$$

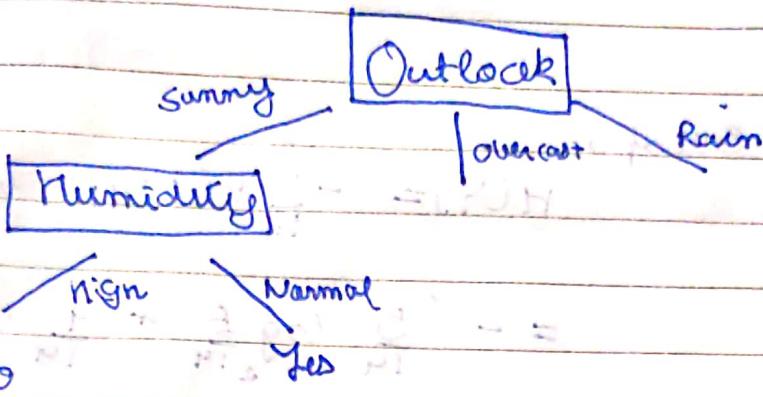
$$= 0.97$$

$$I_G(\text{Wind, Sunny}) = 0.97 + \left[ \frac{3}{5} \left( \frac{2}{3} \log \frac{2}{3} + \frac{1}{3} \log \frac{1}{3} \right) + \right.$$

$$\left. \frac{2}{5} \left( \frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} \right) \right]$$

$$= 0.02$$

Now since there are 3 factors, we split at humidity for sunny

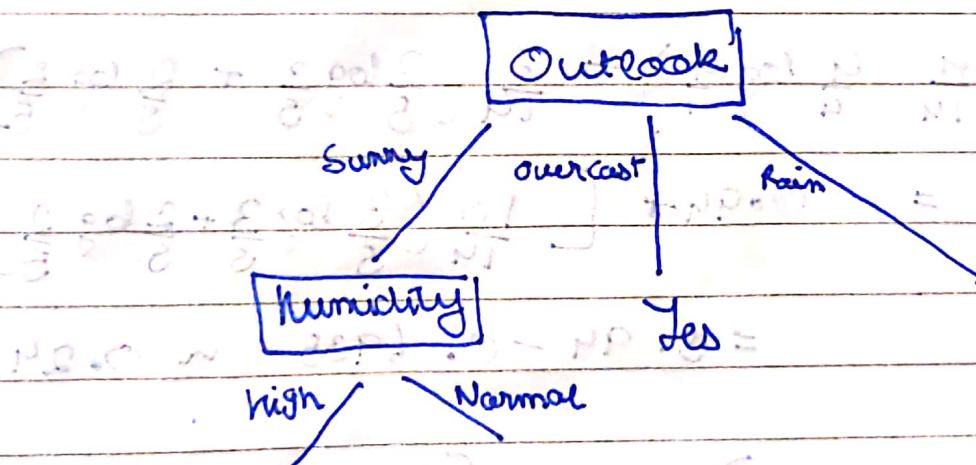


Now, for overcast,

Entropy,  $H(\text{PlayMatch}, \text{Overcast})$

$$H.P.O = -\left(\frac{4}{9} \log \frac{4}{9}\right) = 0$$

No splitting occurs



Now for rain

$$H(\text{PlayMatch}, \text{Rain}) = -\left(\frac{3}{5} \log \frac{3}{5} + \frac{2}{5} \log \frac{2}{5}\right)$$

$$= +0.97$$

Entropy + P.P.O

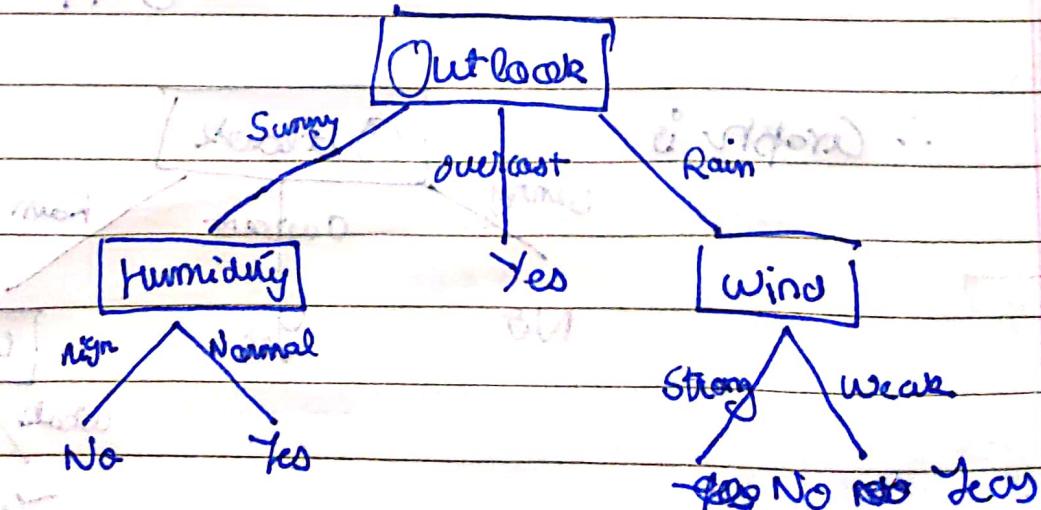
Now, I<sub>G</sub> corresponding to Rain

$$I_{G_r}(\text{Temp, Rainy}) = 0.97 + \left[ \frac{3}{5} \left( \frac{2}{3} \log \frac{2}{3} + \frac{1}{3} \log \frac{1}{3} \right) + \frac{2}{5} \left( \log \frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} \right) \right] = 0.02$$

$$I_{G_r}(\text{Humidity, Rain}) = 0.97 + \left[ \frac{3}{5} \left( \frac{1}{3} \log \frac{1}{3} + \frac{2}{3} \log \frac{2}{3} \right) + \frac{2}{5} \left( \frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} \right) \right] = 0.02$$

$$I_{G_r}(\text{Wind, Rain}) = 0.97 + \left[ \frac{3}{5} \cdot \left( \frac{2}{3} \log \frac{2}{3} \right) + \frac{2}{5} \cdot \left( \frac{2}{2} \log \frac{2}{2} \right) \right] = 0.97$$

$\therefore$  Max  $I_{G_r}$  for Wind  $\therefore$  Finally we have



Ans Day

D<sub>1</sub>D<sub>2</sub>D<sub>3</sub>D<sub>4</sub>D<sub>5</sub>D<sub>6</sub>D<sub>7</sub>D<sub>8</sub>D<sub>9</sub>D<sub>10</sub>D<sub>11</sub>D<sub>12</sub>D<sub>13</sub>

Pred

Actual

No No

No No

Yes Yes

Yes Yes

Yes Yes

No No

Yes Yes

No No

Yes Yes

Yes Yes

Yes Yes

Yes Yes

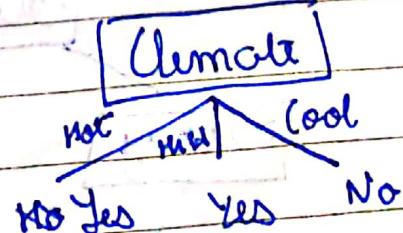
Yes Yes

No No

As seen from above all predicted values match  
Hence accuracy is 100%

- b) Yes it is possible to consider the dataset {D<sub>3</sub>, D<sub>13</sub>, D<sub>1</sub>, D<sub>10</sub>, D<sub>11</sub>, D<sub>12</sub>, D<sub>13</sub>}

In such a dataset Climate on its own is able to accurately predict the play Match.  
Decision Tree would look like



c) Entropy,  $H(Y) = -\left(\frac{4}{7} \log \frac{4}{7} + \frac{3}{7} \log \frac{3}{7}\right)$

$= 0.985$

$I_{a1}(\text{Outlook}) = 0.985 + \left(\frac{2}{7} \cdot \left(\frac{2}{2} \log \frac{2}{2}\right) + \frac{3}{7} \left(\frac{2}{3} \log \frac{2}{3}\right)\right)$

$+ \frac{3}{7} \left(\frac{2}{3} \log \frac{2}{3} + \frac{1}{3} \log \frac{1}{3}\right)$

$= 0.985 - 0.39 = 0.595$

$I_{a1}(\text{Climate}) = 0.985 + \left(\frac{3}{7} \left(\frac{2}{3} \log \frac{2}{3} + \frac{1}{3} \log \frac{1}{3}\right)\right)$

$+ 0 + \frac{3}{7} \left(\frac{2}{3} \log \frac{2}{3} + \frac{2}{3} \log \frac{2}{3}\right)$

$= 0.985 - 0.78 = 0.205$

$I_{a1}(\text{Humidity}) = 0.985 + \left[\frac{4}{7} \left(\frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4}\right) + \frac{3}{7} \left(\frac{2}{3} \log \frac{2}{3} + \frac{1}{3} \log \frac{1}{3}\right)\right]$

$= 0.985 + [-0.39 - 0.51]$   
 $= 0.045$

$I_{a1}(\text{Wind}) = 0.985 + \left[\frac{4}{7} \left(\frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4}\right) + \frac{3}{7} \log \left(\frac{1}{3} \cdot \frac{1}{3} + \frac{2}{3} \log \frac{2}{3}\right)\right]$

$= 0.985 + [-0.39 + (-0.46)]$   
 $= 0.135$

$\therefore$  We split at Outlook



$$H(Y | \text{Outlook} = \text{Sunny}) = 0$$

$$H(Y | \text{Outlook} = \text{Overcast}) = 0$$

$$H(Y | \text{Rain}) = -\left(\frac{2}{3}\log\frac{2}{3} + \frac{1}{3}\log\frac{1}{3}\right) = 0.91$$

Finding  $I_{AB}$

$$I_{AB}(\text{Climate}, \text{Rain}) = 0.91 + [0 + \frac{2}{3}(\log\frac{1}{2} \cdot \frac{1}{2} + \frac{1}{3}\log\frac{1}{3})]$$

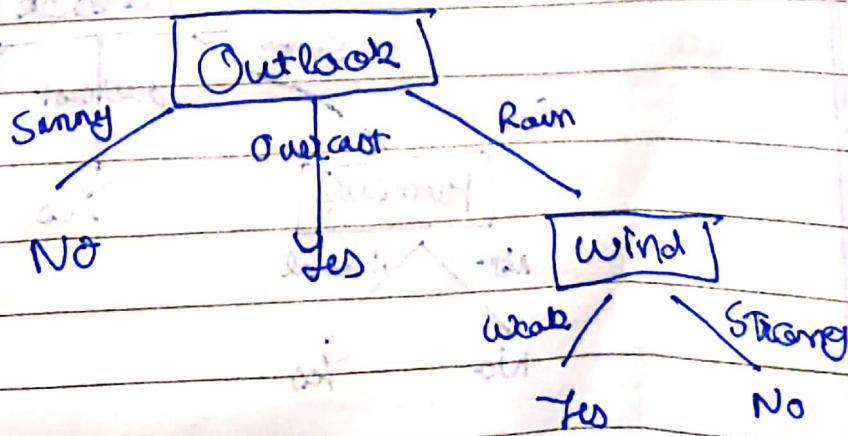
$$\begin{aligned} &= 0.91 - 0.66 \\ &= 0.25 \end{aligned}$$

$$I_{AB}(\text{Humidity}, \text{Rain}) = 0.91 + \left(\frac{2}{3}(\log\frac{1}{2} + \frac{1}{2}\log\frac{1}{2})\right)$$

$$= 0.25$$

$$\therefore I_{AB}(\text{Wind}, \text{Rain}) = 0.91 + (0+0) \\ = 0.91$$

$\therefore$  Graph is



## Finding Accuracy

<u>Day</u>	<u>Predicted</u>	<u>Actual</u>
D <sub>3</sub>	No	No
D <sub>9</sub>	No	Yes
D <sub>10</sub>	Yes	Yes
D <sub>11</sub>	No	Yes
D <sub>12</sub>	Yes	Yes
D <sub>13</sub>	Yes	Yes
+ D <sub>14</sub>	No	No

$$\therefore \text{Accuracy} = \frac{5}{7} \times 100 = 71.4\%$$

There are a few reasons why this might be happening.

One could be that the data being trained on isn't diverse enough. This is evident as sunny always predicts no, however training testing set varies and has data where sunny can take different values depending on different factors. The model is hence overfitting on the given training data leading to large errors on the unseen testing data.

d) Overfitting usually occurs when a model is too complex. To make our tree simpler we want to reduce the number of leaves. We can use the following approach -

- ① Split data into training and testing
- ② Use training to build up a tree.
- ③ Now start from the bottom and try replacing head of a leaf node with a leaf node (This leaf node will be labelled with the majority class)
- ④ If error in testing decreases then we keep the leaf node else we revert back
- ⑤ Repeat until no longer possible.

The tree we get as a result will be much simpler and hence shouldn't overfit

Q6 In a first Order Markov Model, probability of next state only depends on the current state  
 $P(x_{n+1} = x_1 | x_1 = x_1, x_2 = x_2, \dots, x_n = x_n) = P(x_{n+1} = x_{n+1} | x_n = x_n)$

Now, we need to find the conditional probability given the following information  $P(w_3 | w_1, w_2, w_n)$

From markov model's properties,

$$P(w_3 | w_1, w_2, w_n) = P(w_3 | w_2, w_n)$$

Now using Bayes Rule for multiple Variables,

$$(P(R|I, S) = \frac{P(I|R, S)P(R|S)}{P(I|S)}) \text{, we have,}$$

$$P(w_3 | w_2, w_n) = \frac{P(w_2 | w_3, w_n) \cdot P(w_3 | w_2)}{P(w_n | w_2)}$$

Again from Markov models' properties,

$$P(w_n | w_3) = \frac{P(w_n | w_3) \cdot P(w_3 | w_2)}{P(w_n | w_2)}$$

$$\text{Now } P(w_n | w_2) = P(w_3 = \text{tough} | w_2) \cdot P(w_4, w_3 = \text{tough})$$

$$+ P(w_3 = \text{cause} | w_2) \cdot P(w_4, w_3 = \text{cause})$$

$$= [(0.5 \times 0.3)] + [(0.5 \times 0.5)] = 0.4$$

$$\therefore P(W_3 = \text{tough} | W_2, W_1) = \frac{P(\text{aux} | \text{Tough}) \cdot P(\text{Tough} | \text{aux})}{0.4}$$

$$= \frac{0.3 \times 0.5}{0.4} = \frac{15}{40} = \frac{3}{8}$$

$$\text{and } P(W_3 = \text{cause} | W_2, W_1) = \frac{P(\text{cause} | \text{cause}) \cdot P(\text{cause} | \text{cause})}{0.4}$$
$$= \frac{0.25}{0.40} = \frac{25}{40} = \frac{5}{8}$$

Q.7

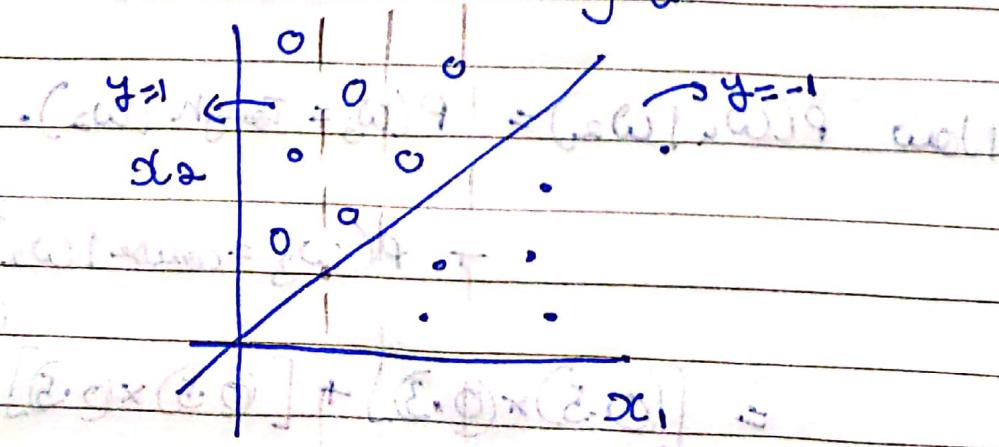
a) In the case of logistic regression we only have one decision boundary (line or plane dividing data into different classes). So for logistic regression to work properly, we'll need data to be linearly separable, whereas decision trees can make "mopes" by using boundaries parallel to the axes and hence isn't bounded by this constraint.

b) Due to the above mentioned fact, decision tree also tend to overfit in case of large number of features. In these cases, logistic regression tends to perform better.

c) We have been given  $n$  2-dimensional points which can be classified using a regression classifier into  $y=1$  or  $y=-1$ , as

$$y = \begin{cases} 1, & \omega^T x + b > 0 \\ -1, & \omega^T x + b \leq 0 \end{cases}$$

This looks something like



Now, we can split basis according to  $x_1$ . Now, as values are given as being linearly separable, there exists  $x_2$  for every  $x_1$ , above which points belong to different class ( $(x_1, x_2)$  lies on the original boundary). In worst case every point has unique  $x_2$ , hence  $\log n$  nodes to decide  $x_1$  and one more to decide  $x_2$ , to represent norm to check if point belongs to  $y=1$  or  $y=-1$ . Hence we need  $\log(n)$  height tree (2n nodes).

d) In case they aren't linearly separable, we can still divide using values of  $x_1$ , however, we don't have  $x_2$  corresponding to  $x_1$ . Instead we'll have to check for  $x_2$  over n points as well. Therefore we will have total of  $n^2$  nodes, to represent this we need  $\log(n^2)$   
 $= 2\log(n)$  height tree

Q.8

We are given  $x$ , set of  $n$  Boolean vectors,  $x_1, x_2, \dots, x_n$   
 and  $y = y_0$  where  $y_i = 0$  or  $1$   $\therefore$   $y$  binary as well

$$\text{Now, } P(Y=1|x) = \frac{P(x_1|y) \times P(y)}{P(x)}$$

$$= \prod_{i=1}^n \frac{P(x_i|y) \times P(y)}{P(x_1, x_2, \dots, x_n)}$$

Now as  $y$  is binary,

$$P(y) = P(y=0) + P(y=1)$$

$$\text{Also, } P(x \cap y = y_0) = P(x|y=y_0) \cdot P(y=y_0)$$

$$\therefore \text{We have, } P(Y=1|x) = \frac{\prod_{i=1}^n P(x_i|y) \times P(y)}{P(y=0) \cdot P(x) + P(y=1) \cdot P(x|y)}$$

$$= \frac{\prod_{i=1}^n P(x_i|y) \times P(y)}{P(y=0) \cdot \prod_{i=1}^n P(x_i|y=0) + P(y=1) \cdot \prod_{i=1}^n P(x_i|y=1)}$$

Now, dividing numerator and denominator by  $P(y=0) \cdot \prod_{i=1}^n P(x_i|y=0) \cdot P(y=1)$

We have

$$1 + \frac{\prod_{i=1}^n P(x_i|y=0)}{\prod_{i=1}^n P(x_i|y=1)}$$

$$= \frac{1}{1 + e^{(\ln(\frac{\prod_{i=1}^n P(x_i|y=0)}{\prod_{i=1}^n P(x_i|y=1)}) - \ln(P(y=0)))}}$$

$$P(Y=1 | X) = \frac{1}{1 + e^{[\ln(\frac{P(Y=1)}{P(Y=0)}) + \sum_{i=1}^n \ln(\frac{P(x_i | Y=1)}{P(x_i | Y=0)})]}}$$

Now,  $P(Y=1) = \pi \quad \therefore P(Y=0) = 1 - \pi$

Also,  $P(x_i | Y=1) = \theta_{i1} \quad P(x_i | Y=0) = \theta_{i0}$

Also,  $P(x_i = j | Y=k) = \theta_{jk} \circ (1 - \theta_{jk})^{1-j}$

from given defn in hint

$\therefore$  We have

$$P(Y=1 | X) = \frac{1}{1 + e^{[\ln(\frac{1-\pi}{\pi}) + \sum_{i=1}^n \ln(\frac{\theta_{i0}^{x_i} (1-\theta_{i0})^{1-x_i}}{\theta_{i1}^{x_i} (1-\theta_{i1})^{1-x_i}})]}}$$

Defined Likelihood at given observation

at  $i^{th}$  observation

$$L_i = \frac{1}{1 + e^{[\ln(\frac{1-\pi}{\pi}) + \sum_{i=1}^n \ln(\frac{\theta_{i0}^{x_i} (1-\theta_{i0})^{1-x_i}}{\theta_{i1}^{x_i} (1-\theta_{i1})^{1-x_i}})]}}$$

at  $i^{th}$  observation

$$= \frac{1}{1 + e^{[\ln(\frac{1-\pi}{\pi}) + \sum_{i=1}^n x_i \ln(\frac{\theta_{i0}}{\theta_{i1}}) + \sum_{i=1}^n (1-x_i) \ln(\frac{(1-\theta_{i0})}{1-\theta_{i1}})]}}$$

$$= \frac{1}{1 + e^{[\ln(\frac{1-\pi}{\pi}) + \sum_{i=1}^n x_i \ln\left(\frac{\theta_{i0}(1-\theta_{i1})}{\theta_{i1}(1-\theta_{i0})}\right) + \sum_{i=1}^n \ln\left(\frac{1-\theta_{i0}}{1-\theta_{i1}}\right)]}}$$

Taking  $w_0 = \ln(\frac{1-\pi}{\pi}) + \sum_{i=1}^n \ln\left(\frac{1-\theta_{i0}}{1-\theta_{i1}}\right)$

$$w_i = \ln\left(\frac{\theta_{i0}(1-\theta_{i1})}{\theta_{i1}(1-\theta_{i0})}\right)$$

$$\therefore P(Y=1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i x_i)}$$

Hence we have shown what was required

$$\frac{(1)^9 * (0)^{10}}{10!} = \frac{1}{10!}$$

$$\frac{(1)^9 * (0)^{10}}{10!} =$$