

Answer to Q1:

a) Insufficient training data -

Ans.1. Linear Discriminant Analysis

Because Logistic regression can become unstable when there are few examples from which to estimate the parameters.

b) Class imbalance -

Ans. 3. Either of the two. We need to use additional sampling in case of class imbalance to improve the accuracy of model

c) Different co-variance matrices for the classes with Gaussian distribution.

Ans. 2. Logistic Regression works well when co-variance matrices differ for the classes, because LDA expects Homogeneity of variance /covariance among classes.

d) Uniform distribution instead of Gaussian distribution.

Ans. 2. Logistic Regression because LDA for simplification assumes that data is gaussian.

Answer to Q2:

Initial Entropy:

$$H(x) = -(((4/10)*\log_2(4/10))+((6/10)*\log_2(6/10))) \\ = 0.971$$

Information Gain (Weather):

Entropy(Weather) =

$$\begin{aligned} & (4/10) * (-(3/4) * \log_2(3/4) + (1/4) * \log_2(1/4)) \text{ (Weather=Fine)} \\ & + (4/10) * (-(3/4) * \log_2(3/4) + (1/4) * \log_2(1/4)) \text{ (Weather=Rain)} \\ & + (2/10) * 0 \text{ (Weather=Cloudy)} \\ & = 0.646 \end{aligned}$$

$$\text{Information Gain(Weather)} = \text{Initial Entropy} - \text{Entropy(Weather)} = 0.971 - 0.646 = 0.325$$

Information Gain (Humidity):

Entropy(Humidity) =

$$\begin{aligned} & (5/10) * (-(2/5) * \log_2(2/5) + (3/5) * \log_2(3/5)) \text{ (Humidity = High)} \\ & + (5/10) * (-(2/5) * \log_2(2/5) + (3/5) * \log_2(3/5)) \text{ (Humidity = Medium)} \\ & = 0.843 \end{aligned}$$

$$\begin{aligned} \text{Information Gain (Humidity)} &= \text{Initial Entropy} - \text{Entropy(Humidity)} \\ &= 0.971 - 0.843 = 0.128 \end{aligned}$$

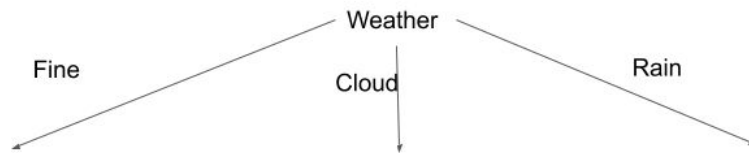
Information Gain (Wind):

Entropy(Wind) =

$$\begin{aligned} & (7/10) * (-(5/7) * \log_2(5/7) + (2/7) * \log_2(2/7)) \text{ (Wind = None)} \\ & + (3/10) * (-(1/3) * \log_2(1/3) + (2/3) * \log_2(2/3)) \text{ (Wind = Breezy)} \\ & = 0.88 \end{aligned}$$

$$\begin{aligned} \text{Information Gain (Wind)} &= \text{Initial Entropy} - \text{Entropy(Wind)} \\ &= 0.971 - 0.88 = 0.091 \end{aligned}$$

Since, Information gain is high if we split data set on attribute weather, we will pick Weather as an attribute to split initially.



RID	Humidity	Wind	Play Golf
1	high	none	no
3	high	none	no
4	medium	none	yes
5	high	breezy	no

RID	Humidity	Wind	Play Golf
2	medium	breezy	yes
6	high	none	yes

RID	Humidity	Wind	Play Golf
7	high	none	yes
8	medium	none	yes
9	medium	breezy	no
10	medium	none	yes

Now new computing information gain for each individual sub data sets

1. Sub Data Set Weather=Fine

Initial Entropy:

$$H(x) = -\left(\left(\frac{3}{4}\right) \cdot \log_2\left(\frac{1}{4}\right)\right) + \left(\left(\frac{1}{4}\right) \cdot \log_2\left(\frac{3}{4}\right)\right)$$

$$= 0.811$$

Information Gain (Humidity):

Entropy(Humidity) =

$$\left(\frac{3}{4}\right) \cdot -\left(\left(\frac{3}{3}\right) \cdot \log_2\left(\frac{3}{3}\right)\right) + \left(\frac{0}{3}\right) \cdot \log_2\left(\frac{0}{3}\right) \quad (\text{Humidity} = \text{High})$$

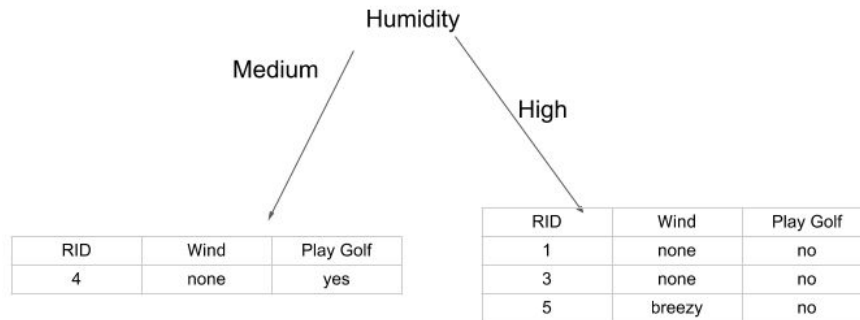
$$+ \left(\frac{1}{4}\right) \cdot -\left(\left(\frac{1}{1}\right) \cdot \log_2\left(\frac{1}{1}\right)\right) + \left(\frac{0}{1}\right) \cdot \log_2\left(\frac{0}{1}\right) \quad (\text{Humidity} = \text{Medium})$$

$$= 0$$

$$\text{Information Gain (Humidity)} = \text{Initial Entropy} - \text{Entropy(Humidity)}$$

$$= 0.811 - 0 = 0.811$$

Since here information gain is maximum here for this split, we can skip other attributes and split on Humidity



2. Sub Data Set Weather=Cloud

We can stop examining further attributes for this split since information gain is maximum here.

3. Sub Data Set Weather=Rain

Initial Entropy:

$$H(x) = -\left(\left(\frac{1}{4}\right) \log_2\left(\frac{1}{4}\right) + \left(\frac{3}{4}\right) \log_2\left(\frac{3}{4}\right)\right) \\ = 0.811$$

Information Gain (Humidity):

$$\text{Entropy(Humidity)} = \\ \left(\frac{3}{4}\right) * \left(-\left(\left(\frac{2}{3}\right) \log_2\left(\frac{2}{3}\right) + \left(\frac{1}{3}\right) \log_2\left(\frac{1}{3}\right)\right)\right) \text{ (Humidity = Medium)} \\ + \left(\frac{1}{4}\right) * \left(-\left(\left(\frac{1}{1}\right) \log_2\left(\frac{1}{1}\right) + \left(\frac{0}{1}\right) \log_2\left(\frac{0}{1}\right)\right)\right) \text{ (Humidity = High)} \\ = 0.68875$$

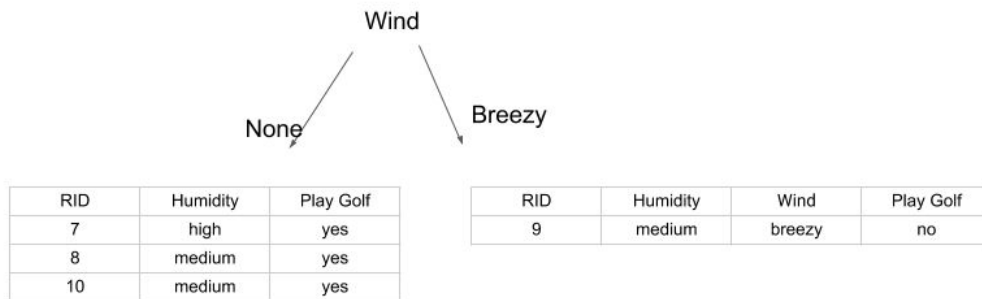
$$\text{Information Gain (Humidity)} = \text{Initial Entropy} - \text{Entropy(Humidity)} \\ = 0.811 - 0.68875 = 0.12225$$

Information Gain (Wind):

$$\text{Entropy(Wind)} = \\ \left(\frac{3}{4}\right) * \left(-\left(\left(\frac{3}{3}\right) \log_2\left(\frac{3}{3}\right) + \left(\frac{0}{3}\right) \log_2\left(\frac{0}{3}\right)\right)\right) \text{ (Wind = none)} \\ + \left(\frac{1}{4}\right) * \left(-\left(\left(\frac{1}{1}\right) \log_2\left(\frac{1}{1}\right) + \left(\frac{0}{1}\right) \log_2\left(\frac{0}{1}\right)\right)\right) \text{ (Wind = breeze)} \\ = 0$$

$$\text{Information Gain (Wind)} = \text{Initial Entropy} - \text{Entropy(Wind)} \\ = 0.811 - 0 = 0.811$$

So for this split, if we split of Wind we get maximum gain.



Final Decision Tree will look like this.

