

Amrita Vishwa Vidyapeetham
College Of Engineering,
Bengaluru

Name: Raman Kumar
Program: Mtech DS
Course: Machine Learning
Registration number: BL.EN.P2DSC22009

Contents:-

Assignment 02

(1) Confusion Matrix for binary-Class Scenario.

- Binary class produces outputs with only 2-label values. example (Yes or No), (1/0) (+ve / -ve) for the given data.
- We use the observed classes to compare with predicted classes for the different evaluation performances after classification.
- Confusion matrix table describe the performance metrics on a set of test data, for which the positive values or true values are known.
- Confusion matrix label will be exactly the same if the performance of classifier will be perfect.
- It has four outcomes
 - True Positive : Correct +ve prediction
 - True Negative : Incorrect +ve prediction
 - False Positive : correct -ve prediction
 - False Negative : Incorrect -ve prediction

		-ve	+ve
Actual	-ve	True Negative	False Positive
	+ve	False Negative	True Positive

for observed "Predicted"

"Type-I" Error

"Type-II" error

- Accuracy is calculated as the number of all correct predictions divided by total number of dataset

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

- Recall: it is sensitivity (or) true positive rate.

When it's actually positive, how often does it predict positive?

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\frac{\text{TP}}{\text{TP} + \text{FN}}$$

- Precision: is calculated as number of correct positive predictions divided by total no. of positive predictions.

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

Confusion matrix for multiclass classification:

		Predicted			Actual
		TN	FP	FN	
Actual	TN	FP	FN	TN	
	FP	TN	FN	FP	

Q.2). Beta-variation in F-measure:

The F-score is a way of combining the precision and recall of the model, and it is defined as harmonic mean of the model's precision and recall.

- The F-score balances the precision and recall.
- F-measure measures two types of errors that are made for positive classes.

- Minimizing FP, maximizing Precision
- Minimizing FN, maximizing Recall.

F-score provide a way to combine both precision & recall into single measure.

$$\text{F-measure} = \frac{(2 * \text{Precision} * \text{Recall})}{\text{Precision} + \text{Recall}}$$

- In F- β measure, the coefficient β controls the balance of precision & recalled.

$$F - \beta = \frac{(1 + \beta^2) * \text{Precision} * \text{Recall}}{\beta * \text{Precision} + 1}$$

Three common values of beta-parameter are -

- 1) F-0.5 score ($\beta=0.5$): More weight on precision, less weight on recall.

- 2) F-1 score ($\beta=1$): Balance the weight on precision & recall.
- 3) F-2 score ($\beta=2$): Less weight on precision, more weight on recall.

(Q.3). Overfitting in Decision Trees:

Overfitting is a modeling error which occurs when a function is too closely fit to a limited set of data points.

and h : hypothesis error
error on training data: $\text{error}_{\text{train}}(h)$
error over the entire distribution D of data: $\text{error}_D(h)$.

When hypothesis h overfits the training data,
if there is an alternative hypothesis h'

$\text{error}_{\text{train}}(h) < \text{error}_{\text{train}}(h')$

$\text{error}(h) < \text{error}_D(h)$

Causes of Overfitting:

- Overfitting due to presence of noise
- Overfitting due to lack of representative instances
- Overfitting & multiple comparison procedure

Underfitting in Decision Trees:

Learning algorithm had the opportunity to learn more from training data, but didn't, that comes under underfitting.

It's when model is too simple for the data.

Causes: • Less no. of features, or small dataset.

(Q4) E.S. OR gate, for following question is [] = ~~stab~~ []

[(-2, 0.75), (0, 0.63), (0, 1), (1, 1)]

Input 1	Input 2	Output
0	0	0
0	1	1
1	0	1
1	1	1

- Initial weights are $(-2, 0.75, 0.63)$

- First row of truth table, the inputs are $(0, 0)$

$$\text{output} = 1 \text{ if } [(-2 * 0) + (0.75 * 0) + (0.63 * 1)] \geq 0, \text{ else } 0$$

$$\text{output} = 0$$

- Target output and observed output are difficult, update the weights:

$$\Delta W_0 = \alpha * |\text{predicted value} - \text{observed value}| * \text{input}$$

$$= 0.1 * |0 - 0| * 0 = 0$$

$$W_0 \leftarrow W_0 + \Delta W_0 = -2 + 0 = -2$$

- Change in weight of W_1 is:

$$\Delta W_1 = 0.1 * |0 - 0| * 0 = 0$$

$$W_1 \leftarrow 0.75 + 0 = 0.75$$

- Change in weight of W_2 is:

$$\Delta W_2 = 0.1 * |0 - 0| * 1 = 0$$

$$W_2 \leftarrow W_2 + \Delta W_2 = 0.63 + 0 = 0.63$$

- For the 2nd row of truth table, inputs: $(0, 1)$, outputs: 1, weights $(-2, 0.75, 0.63)$

- output = 1, if $(-2 * 0) + (0.75 * 0) + (0.63 * 1) \geq 0$

$$\text{output} = 0$$

\therefore Since the target value & observed output are different, update the weights.

- change in weight for W_0 is:

$$\Delta W_0 = 0.1 * |1 - 0| * 0 = 0$$

$$W_0 \leftarrow W_0 + \Delta W_0 = -2 + 0 = -2$$

- The change in weight w_1 is:
 $\Delta w_1 = 0.1 * |1 - 0| * 0.75 = 0.75$
 $w_1 \leftarrow w_1 + \Delta w_1 = 0.75 + 0 = 0.75$
- change in weight w_2 is:
 $\Delta w_2 = 0.1 * |1 - 0| * 0 = 0.1$
 $w_2 \leftarrow w_2 + \Delta w_2 = 0.63 + 0.1 = 0.73$

Third row of truth table
 inputs = $(1, 1, 0)$, target output = 1, weights =
 $(-2, 0.75, 0.73)$

- output = 1, if $(-2 * 1) + (0.75 * 0) + (0.73 * 1) \geq 0$; else 0

$$\text{output} = 1$$

Since the target value & observed output are same.

- now no need to update the weights.

- For 4th row input $(1, 1)$, output = 1.



(Q.5). Trace a decision tree

- Information gain: The goal is to maximize the information gain at each split, which means we want to select the feature that results in the most homogeneous subgroups.

- Entropy of label (Sex) at root node:

$$\text{Entropy} = -(P(M) * \log_2(P(M)) + P(F) * \log_2(P(F)))$$

kind of trees see $P(M)$: proportions of males

and $P(F)$: proportions of females

where, 4 males & 3 females

$$\text{Entropy} = -(4/7 * \log_2(4/7) + 3/7 * \log_2(3/7))$$

Result = 0.985

- IG on each feature i.e. entropy of subgroups.

- height feature:

$$\text{three subgroups for height}$$

High	Medium	Low
$= -\left(\frac{1}{7} * \log_2\left(\frac{1}{7}\right)\right)$	$= -\left(\frac{2}{7} * \log_2\left(\frac{2}{7}\right)\right)$	$= -\left(\frac{4}{7} * \log_2\left(\frac{4}{7}\right)\right)$

IG = Entropy (root) - (weighted average of entropies of subgroups)

$$\geq 0.985 - \left(\frac{1}{7} * 0 + \frac{2}{7} * 0 + \frac{4}{7} * 1\right)$$

$$= 0.985 - \frac{4}{7} = 0.169$$

- We can repeat the process for the other features to calculate their IG.

* weight : $IG = 0.985 - \left(\frac{1}{7} * 0 + \frac{3}{7} * 0.918 + \frac{3}{7} * 0\right)$

$$= 0.151$$

$$\text{Hair length} = IG = 0.985 - \left(\frac{2}{7} * 0.918 + \frac{5}{7} * 0\right)$$

$$= \underline{\underline{0.043}}$$

- Based on the calculations, highest IG feature is height.

• Gini Index: at root node,

$$\text{Gini} = 1 - (P(M)^2 + P(F)^2)$$

$$= 1 - \left(\frac{4}{7}\right)^2 + \left(\frac{3}{7}\right)^2 = 0.48$$

- For height feature:

$$\text{High : Gini} = 1 - \left(\frac{1}{1}\right)^2 - \left(\frac{0}{1}\right)^2 = 0$$

$$\text{Med : Gini} = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5$$

$$\text{Low : Gini} = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5$$

$$\text{Gini} = \left(\frac{1}{7}\right) * 0 + \left(\frac{2}{7} * 0.5\right) + \left(\frac{4}{7} * 0.5\right) = 0.429$$

- Weight feature:

$$\text{Gini} = \left(\frac{1}{7} * 0\right) + \left(\frac{3}{7} * 0.5\right) + \left(\frac{3}{7} * 0\right) = 0.286$$

$$\text{Hair length} = \left(\frac{2}{7} * 0.5\right) + \left(\frac{5}{7} * 0\right) = 0.143$$

Gini index of ~~highest~~ 'height' is 'lowest'.

- (6). To calculate the parameters 'm' & 'c' for a line of the form $Y = mX + c$ using linear regression, you can use the following steps:
1. Calculate the mean of the X values & the mean of the Y values.
 2. Calculate the difference between each X value and the mean of the X values, and the difference between each Y value and the mean of the Y values.
 3. Calculate the sum of the products of the differences calculated in step 2.
 4. ~~5.~~ Calculate the sum of the squares of the differences in the X values.
 5. Divide the sum from Step 3 by the sum from Step 4. The result is the value of m.
 6. Calculate the value of c by substituting the mean of the X values and the mean of Y values into the equation $Y = mX + c$ & solving for c.

. data = [(1.0, 1.75), (2, 2.14), (2.5, 2.47), (3, 2.38),
(3.5, 3.55), (4, 2.72), (5, 2.91), (4.5, 4.2)]

Step 1:
x-mean = sum(x for x, y in data) / len(data)

y-mean = sum(y for x, y in data) / len(data)

Step 2:
xy-diff = [(x - x-mean) * (y - y-mean) for
(x, y) in data]

x_sq_diff = [(x - x-mean) ** 2 for x, y in data]

Step 3:
xy-diff_sum = sum(xy-diff)

Step 4:
x_sq_diff_sum = sum(x_sq_diff)

Step 5:
m = xy-diff_sum / x_sq_diff_sum

Step 6:
c = y-mean - m * x-mean

printf("m = %f")

printf("c = %f")

$\Rightarrow m = 0.4$ $c = 1.3$

for i in range(10):

 o = 4 * (o - 0) * 1.0 + 1.3

 print(o)

for i in range(10):

 o = (o + (o * 2.0) + (o * 2.0)) / 3.0

 print(o)



(8). In SVM, the kernel function plays a crucial role in mapping the input data into a higher-dimensional space, where it becomes possible to perform linear classification. This is important because many datasets are not linearly separable in their original feature space, but they might be linearly separable in a higher-dimensional space.

- e.g., suppose we have a dataset with two classes, A and B, and we want to build an SVM classifier to separate the two classes.

If the data is not linearly separable in the original feature space, we can apply a kernel function to map the data into a higher dimensional space where it becomes linearly separable.

- Types of Kernal with SVM:
 - linear kernel
 - polynomial kernel
 - radial basis kernel (RBF)

- For multi-class classification, there are several ways in which SVM's can be used. One approach is to build a one-versus-all (OVA) classifier, in which a separate SVM is trained for each class, with that class as the positive class and all other classes as the negative class.

- Another approach is to use a one-versus-one (OVO) classifier, in which a separate SVM is trained for each pair of classes.

(7.)

Support Vectors	$y_i * \alpha_i$	Class Values
[1.7 0.5]	-0.4485	-1
[1.9 0.2]	-1	-1
[1.9 0.4]	-1	-1
[3.3 1]	0.436651	1
[5.1 1.6]	0.160646	1
[3 1.1]	1	1

$$\text{Intercept } (b) = -6.8$$

$$\text{Vector, } A = [7 \ 3]$$

Class Value for test vector X

$$= \left[\text{Sign}(\sum (y_i * \alpha_i * \text{Dot_Product}(SupportVector, } X)) + b \right]$$

$[\text{Support vectors, } X] \Rightarrow \text{Dot product of support vector and } X.$

$$= \begin{bmatrix} 1.7 & 0.5 \\ 1.9 & 0.2 \\ 1.9 & 0.4 \\ 3.3 & 1 \\ 5.1 & 1.6 \\ 3 & 1.1 \end{bmatrix}_{6 \times 2} \cdot \begin{bmatrix} 7 \\ 3 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} 13.4 \\ 13.9 \\ 14.5 \\ 26.1 \\ 40.5 \\ 24.3 \end{bmatrix}_{6 \times 1}$$

Applying the formula, $\sum (y_i * \alpha_i * \text{Dot_Product}(SupportVector, } X)$

$$= [-0.4485 \ -1 \ -1 \ 0.436651 \ 0.160646 \ 1] \begin{bmatrix} 13.4 \\ 13.9 \\ 14.5 \\ 26.1 \\ 40.5 \\ 24.3 \end{bmatrix}_{6 \times 1}$$

$$= (0.4485 \times 13.4) + (-1 \times 13.9) + (-1 \times 14.5) + (0.43665 \times 26.1) + (0.160646 \times 40.5) + (1 \times 24.3)$$

$$= \underline{7.792828}$$

$$X = 7.792828 + \text{intercept} = 7.792828 - 6.8 = \underline{0.992828}$$

Sign is +ve.

$\therefore \underline{\text{Class value is +1}}$

