

Machine Learning EE514 - CS535

Analysis and Evaluation of Classifier's Performance and Multi-class Classification



School of Science and Engineering Lahore University of Management Sciences

https://www.zubairkhalid.org/ee514 2023.html





Outline

- Classification Accuracy (0/1 Loss)
- TP, TN, FP and FN
- Confusion Matrix
- Sensitivity, Specificity, Precision Trade-offs, ROC, AUC
- F1-Score and Matthew's Correlation Coefficient
- Multi-class Classification, Evaluation, Micro, Macro Averaging



Classification Accuracy, Misclassification Rate (0/1 Loss):

$$\mathcal{L}_{0/1}(h) = \frac{1}{n} \sum_{i=1}^{n} 1 - \delta_{h(\mathbf{x}_i) - y_i} \qquad \qquad \delta_k = \begin{cases} 1, & k = 0\\ 0 & \text{otherwise} \end{cases}$$

- For each test-point, the loss is either 0 or 1; whether the prediction is correct or incorrect.
- Averaged over n data-points, this loss is a 'Misclassification Rate'.

Interpretation:

- Misclassification Rate: Estimate of the probability that a point is incorrectly classified.
- Accuracy = 1 Misclassification rate

Issue:

- Not meaningful when the classes are imbalanced or skewed.



Classification Accuracy (0/1 Loss):

Example:

- Predict if a bowler will not bowl a **no-ball**?
 - Assuming 15 no-balls in an inning, a **model that says 'Yes' all the time** will have **95%** accuracy.
 - Using accuracy as performance metric, we can say that a model is very accurate, but it is not useful or valuable in fact.

<u>Why?</u>

- Total points: 315 (assuming other balls are legal 🙂)
- No-ball label: Class O (4.76% are from this class)
- Not a no-ball label: Class 1 (95.24% are from this class)

Imbalanced Classes



TP, TN, FP and FN:

- Consider a binary classification problem.

 $D = \{ (\mathbf{x_1}, y_1), (\mathbf{x_2}, y_2), \dots, (\mathbf{x_n}, y_n) \} \subseteq \mathcal{X}^d \times \mathcal{Y}$

 $\mathcal{Y} = \{0, 1\}$ (Referring 0 as Negative, 1 as Positive)

 \boldsymbol{y} - Actual labels, Ground truth, Gold labels or Standards

We have a classifier (hypothesis function) $h(\mathbf{x}) = \hat{y}$.

$$y, \hat{y}$$
 - Positive (1) or Negative (0)
 \hat{y} - True if $\hat{y} = y$, False if $\hat{y} \neq y$



TP, TN, FP and FN:

- TP True Positive Number of points with y = 1 and are classified as $\hat{y} = 1$
- TN True Negative Number of points with y = 0 and are classified as $\hat{y} = 0$
- FP False Positive Number of points with y = 0 and are classified as $\hat{y} = 1$
- FN False Negative Number of points with y = 1 and are classified as $\hat{y} = 0$



TP, TN, FP and FN:

Example:

- Predict if a bowler will not bowl a **no-ball**?
 - 15 no-balls in an inning (Total balls: 315)
 - Bowl no-ball (Class O), Bowl regular ball (Class 1)
 - Model(*) predicted 10 no-balls (8 correct predictions, 2 incorrect)
 - TP True Positive TP 298
 - TN True Negative TN 8
 - FP False Positive FP 7
 - FN False Negative FN 2



* Assume you have a model that has been observing the bowlers for the last 15 years and used these observations for learning.

Confusion Matrix (Contingency Table):

- (TP; TN; FP; FN); usefully summarized in a table, referred to as confusion matrix:
 - the rows correspond to predicted class (\hat{y})
 - and the columns to true class (y)

	ŀ	Actual Labels		
		1 (Positive)	0 (Negative)	Total
Predicted Labels	1 (Positive)	ТР	FP	Predicted Total Positives
	0 (Negative)	FN	TN	Predicted Total Negatives
	Total	P= TP+FN Actual Total Positives	N= P+TN Actual Total Negatives	



Confusion Matrix:

Example:

- Disease Detection :
- Given pathology reports and scans, predict heart disease
- Yes: 1, No: O

	A			
		1 (Positive)	0 (Negative)	Total
Predicted	1 (Positive)	TP = 100	FP = 10	110
Labels	0 (Negative)	FN = 5	TN = 50	55
	Total	P = 105	N = 60	

Interpretation:

Out of 165 cases

- Predicted: "Yes" 110 times, and "No" 55 times
- In reality: "Yes" 105 times, and "No" 60 times



Confusion Matrix:

Example:

 Predict if a bowler will not bowl a no-ball?

Interpretation:

Out of 315 balls, we had 15 no-balls.

Model predicted 305 regular balls and 10 no-balls (8 correct predictions, 2 incorrect).



	А					
		1 (Positive)	0 (Negative)	Total		
Predicted Labels	1 (Positive)	TP = 298	FP = 7	305		
	0 (Negative)	FN = 2	TN = 8	10		
	Total	P = 300	N = 15			

Confusion Matrix:

Metrics using Confusion Matrix:

- Accuracy: Overall, how frequently is the classifier correct?

$$Accuracy = \frac{TP + TN}{Total} = \frac{TP + TN}{P + N}$$

- **Actual Labels** 1 (Positive) 0 (Negative) Total Predicted 1 (Positive) Predicted TP FP Total Positives Labels Predicted 0 (Negative) FN Total Negatives TN N= P+TN Total P= TP+FN Actual Total Actual Total Positives Negatives
- Misclassification or Error Rate: Overall, how frequently is it wrong?

$$1 - Accuracy = \frac{FP + FN}{Total} = \frac{FP + FN}{P + N}$$

- Sensitivity or Recall or True Positive Rate (TPR): How often does it predict Positive when it is actually Positive?

$$TPR = S_e = \frac{TP}{TP + FN} = \frac{TP}{P}$$



Confusion Matrix:

Metrics using Confusion Matrix:

- False Positive Rate: Actual Negative, how often does it predict Positive?

$$FPR = \frac{FP}{TN + FP} = \frac{FP}{N}$$

			1 (Positive)	0 (Negative)	Total
	Predicted Labels	1 (Positive)	ТР	FP	Predicted Total Positives
		0 (Negative)	FN	TN	Predicted Total Negatives
•		Total	P= TP+FN Actual Total Positives	N= P+TN Actual Total Negatives	

- **Specificity or True Negative Rate (TNR)**: When it's actually Negative, how often does it predict Negative?

$$TNR = S_p = \frac{\mathrm{TN}}{\mathrm{TN} + \mathrm{FP}} = \frac{\mathrm{TN}}{\mathrm{N}} = 1 - FPR$$

- Precision: When it predicts Positive, how often is it Positive?

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

Confusion Matrix Metrics:



$$TPR = S_e = \frac{11}{\text{TP} + \text{FN}} = \frac{11}{\text{P}} \quad TNR = S_p = \frac{11}{\text{TN} + \text{FP}} = \frac{11}{\text{N}}$$



Actual Labels

1 (Positive) 0 (Negative) Total

Confusion Matrix:

Metrics using Confusion Matrix (Example: Disease Prediction):

- Accuracy: Disease/Healthy prediction accuracy $\frac{Predicted}{Labels} = \frac{1 (Positive)}{Predicted} = \frac{100}{P + 10} = 100$ $\frac{Predicted}{Predicted} = \frac{1}{P + 10} = 100$ $\frac{1}{P = 100} = 100$ $\frac{1}{P = 100} = 100$ $\frac{1}{P = 100} = 100$
- Misclassification or Error Rate: Disease/Healthy prediction accuracy

 $1 - Accuracy = \frac{FP + FN}{Total} = \frac{FP + FN}{P + N} = (10+5)/165 = 0.09$

- Sensitivity or Recall or True Positive Rate (TPR): When it's positive, how often does the model detected disease?

$$TPR = S_e = rac{TP}{TP + FN} = rac{TP}{P}$$
 = 100/105 = 0.95



Confusion Matrix:

Metrics using Confusion Matrix (Example: Disease Prediction):

- False Positive Rate: Actually heathy, how often does it predict yes?

$$FPR = \frac{FP}{TN + FP} = \frac{FP}{N} = 10/60 = 0.17$$



- Specificity or True Negative Rate (TNR): When it's actually health, how often does it predict healthy? $TNR = S_p = \frac{TN}{TN + FP} = \frac{TN}{N} = 50/60 = 0.83$

- Precision: When it predicts disease, how often is it correct?

$$\frac{\text{TP}}{\text{TP} + \text{FP}} = 100/110 = 0.91$$



Confusion Matrix:

Metrics using Confusion Matrix:

- When to use which?
- Disease Detection: We do not want FN

$$TPR = S_e = \frac{TP}{TP + FN} = \frac{TP}{P}$$

- Fraud Detection: We do not want FP

$$TNR = S_p = \frac{TN}{TN + FP} = \frac{TN}{N}$$
 Precision $= \frac{TP}{TP + FP}$



	Actual Labels			
		1 (Positive)	0 (Negative)	
Predicted	1 (Positive)	ТР	FP	
Labels	0 (Negative)	FN	TN	

Outline

- Classification Accuracy (0/1 Loss)
- TP, TN, FP and FN
- Confusion Matrix
- Sensitivity, Specificity, Precision Trade-offs, ROC, AUC
- F1-Score and Matthew's Correlation Coefficient
- Multi-class Classification, Evaluation, Micro, Macro Averaging



Confusion Matrix:

Precision and Sensitivity (Recall) Trade-off:

- Disease Detection:

Sensitivity or Recall $TPR = S_e = \frac{TP}{TP + FN} = \frac{TP}{P}$



- Recall or Sensitivity (S_e) ; how good we are at detecting diseased people.
- Precision: How many have been correctly diagnosed as unhealthy.
- If we have diagnosed everyone unhealthy, S_e=1 (diagnose all unhealthy people correctly) but Precision may be low (because TN=0 that increases the value of FP).

	Actual Labels			
		1 (Positive)	0 (Negative)	
Predicted	1 (Positive)	ТР	FP	
Labels	0 (Negative)	FN	TN	

- We want high **Precision** and high S_e (=1, Ideally).
- We should combine precision and sensitivity to evaluate the performance of classifier.
 - F1-Score



Confusion Matrix:

Sensitivity and Specificity Trade-off:



- S_p and S_e ; how good we are at detecting healthy and diseased people, respectively.
- If we have diagnosed everyone healthy, $S_p=1$ (diagnose all healthy people correctly) but $S_e=0$ (diagnose all unhealthy people incorrectly)

- Ideally: we want $S_p = S_e = 1$ (perfect sensitivity and specificity) but unrealistic.



1.0

Confusion Matrix:



- Trade-off is better explained by ROC curve and AUC.



Confusion Matrix:

ROC (Receiver Operating Characteristic) Curve:

Plot of TPR (Sensitivity) against FPR (1 – Specificity)
 for different values of threshold.

- Also referred to as Sensitivity-(1-Specificity) plot.
- Threshold of 0.0, every case is diagnosed as positive.
 S_e= TPR = 1
 FPR = 1
 S_p= 0
- Threshold of 1.0, every case is diagnosed as negative.
 S_e= TPR = 0
 FPR = 0
 S_p= 1





Confusion Matrix:

ROC Curve and AUC:

ROC Curve



- The best possible prediction method - $S_e = S_p = 1$ (Upper left corner of ROC space)
- Random guess; a point along a diagonal line (the so-called line of no-discrimination), No Power!

- Area Under the ROC Curve, abbreviated as (AUC) quantifies the power of the classifier.





Outline

- Classification Accuracy (0/1 Loss)
- TP, TN, FP and FN
- Confusion Matrix
- Sensitivity, Specificity, Precision Trade-offs, ROC, AUC
- F1-Score and Matthew's Correlation Coefficient
- Multi-class Classification, Evaluation, Micro, Macro Averaging



F1-Score:

- We observed trade-off between recall and precision.

$$TPR = S_e = \frac{TP}{TP + FN} = \frac{TP}{P}$$
 Precision = $\frac{TP}{TP + FP}$

- Higher levels of recall may be obtained at the price of lower values of precision.
- We need to define a single measure that combines recall and precision or other metrics to evaluate the performance of a classifier.
- Some combined measures:
 - F1 Score
 - Matthew's Correlation Coefficient
 - 11-point average precision
 - The Breakeven point



F1 Score:

- One measure that assesses recall and precision trade-off is weighted harmonic mean (HM) of recall and precision, that is,

$$F_{\beta} = \frac{1+\beta^2}{\frac{1}{\text{Precision}} + \frac{\beta^2}{\text{Recall}}}, \quad \beta \ge 0$$

For $\beta = 1$, we have harmonic mean of precision and recall, that is,

$$F_1 = \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}} = \frac{2(\text{Precision})(\text{Recall})}{(\text{Precision}) + (\text{Recall})} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$$



F1 Score:

Why harmonic mean?

- We could also use arithmetic mean (AM) or geometric mean (GM).
- HM is preferred as it penalizes model the most; a conservative average, that is, for two real positive numbers, we have

$\rm HM \leq \rm GM \leq \rm AM$

 Improvement in HM implies improvement in AM or GM.



precision. Recall=70% is fixed.



Matthew's Correlation Coefficient (MCC):

- Precision, Recall and F1-score are asymmetric. Get a different result if the classes are switched.
- Matthew's correlation coefficient determines the correlation between true class and predicted class. The higher the correlation between true and predicted values, the better the prediction.

 $(\mathbf{m}\mathbf{n})$

$$\frac{(TP)(TN) - (FP)(FN)}{(TP + FN)(TP + FP)(TN + FN)(TN + FP)}$$

 $(\mathbf{D}\mathbf{D})(\mathbf{D}\mathbf{N}\mathbf{I})$

 $|\mathrm{MCC}| \le 1$

MCC=1 when FP = FN = O (Perfect classification)

- MCC=-1 when TP = TN = O (Perfect misclassification)
- MCC=0; Performance of classifier is not better than a random classifier (flip coin)
- MCC is symmetric by design



<u>11-point Average Precision:</u>

- Adjust threshold of the classifier such that the recall takes the following 11 values 0.0, 0.1., ..., 0.9, 1.0.
- For each value of the recall, determine the precision and find the average value of precision, referred to as average precision (AP).
- This is just uniformly-spaced sampling of Precision-Recall curve and taking average value. **The Breakeven Point:**
- Compute precision as a function of recall for different values of thresholds.
- When Precision = Recall, we have a breakeven.



Outline

- Classification Accuracy (0/1 Loss)
- TP, TN, FP and FN
- Confusion Matrix
- Sensitivity, Specificity, Precision Trade-offs, ROC, AUC
- F1-Score and Matthew's Correlation Coefficient
- Multi-class Classification, Evaluation, Micro, Macro Averaging



Multi-Class Classification

Formulation:

• We assume we have training data D given by

$$D = \{ (\mathbf{x_1}, y_1), (\mathbf{x_2}, y_2), \dots, (\mathbf{x_n}, y_n) \} \subseteq \mathcal{X}^d \times \mathcal{Y}$$

• $\mathcal{Y} = \{1, 2, \dots, M\}$ (M-class classification)

Examples:

- Emotion Detection.
- Vehicle Type, Make, model, color of the vehicle from the images streamed by safe city camera.
- Speaker Identification from Speech Signal.
- State (rest, ramp-up, normal, ramp-down) of the process machine in the plant.
- Sentiment Analysis (Categories: Positive, Negative, Neutral), Text Analysis.
- Take an image of the sky and determine the pollution level (healthy, moderate, hazard).
- Record Home WiFi signals and identify the type of appliance being operated.







Multi-Class Classification

Implementation (Possible options using binary classifiers):

Option 1: Build a one-vs-all (OvA) one-vs-rest (OvR) classifier:

Train M different binary classifiers $h_1(\mathbf{x}), h_2(\mathbf{x}), \ldots, h_M(\mathbf{x})$.

Classifier $h_i(\mathbf{x})$ is trained to classify if \mathbf{x} belongs to *i*-th class or not.

For a new test point \mathbf{z} , get scores for each classifier, that is, $s_i = h_i(\mathbf{z})$. For example, s_i can be assigned the probability that \mathbf{z} belongs to class i.

Predict the label as $\hat{y} = \max_{i=1,2,...,M} s_i$

Option 2: Build an all-vs-all classifier:

Train $\binom{M}{2} = \frac{(M)(M-1)}{2}$ different binary classifiers $h_{i,j}(\mathbf{x}), i, j = 1, 2, ..., k$

Classifier $h_{i,j}(\mathbf{x})$ is trained to classify if \mathbf{x} belongs to *i*-th class or *j*-th class.

For a new test point \mathbf{z} , get scores for each classifier, that is, $s_{i,j} = h_{i,j}(\mathbf{z})$. For example, $s_{i,j} = 1$ if \mathbf{z} is likely class i or 0 if \mathbf{z} is likely class j.

Predict the label \hat{y} that has been predicted multiple times.

<u>There can be other options...</u>



Multiclass Classification:

- How do we define the measures for the evaluation of the performance of multi-class classifier?
- Macro-averaging: We compute performance for each class and then average.
- Micro-averaging: Compute confusion matrix after collecting decisions for all classes and then evaluate.



Multiclass Classification:

Confusion Matrix

- Predict if a bowler will bowl a no-ball, wide bowl, regular bowl?
 - 15 no-balls, 20 wide-balls in an inning (Total balls: 335)
 - Model Predictions:

		No-ball	ACTUAI Wide-ball	Regular ball	Precision
	No-ball	8	5	20	$\frac{8}{8+5+20}$
Classifier Output	Wide-ball	2	10	10	$\frac{10}{2+10+10}$
·	Regular ball	5	5	270	$\frac{270}{5+5+270}$
R	ecall	$\frac{8}{8+2+5}$	$\frac{10}{5+10+5}$	$\frac{270}{20+10+270}$	

A . I . I



Multiclass Classification:

Confusion Matrix – Recall and Precision:

 $C_{i,j}$ represents the entry of the confusion matrix at *i*-th row and *j*-th column. **Recall**

- For i-th class, recall represents the fraction of data-points classified correctly, that is,

 $\sum_{i=1}^{m} C_{i,j}$

Precision

- For i-th class, precision represents the fraction of data-points predicted to be in class i are actually in the i-th class, that is,

$$Precision_i = \frac{C_{i,i}}{\sum\limits_{j=1}^{M} C_{i,j}}$$

 $\operatorname{Recall}_i = \frac{C_{i,i}}{M}$

Accuracy

- Fraction of data points classified correctly, that is,

Accuracy =
$$\frac{\sum_{i=1}^{M} C_{i,i}}{\sum_{i=1}^{M} \sum_{j=1}^{M} C_{i,j}}$$

	No-ball	Wide-ball	Regular ball
No-ball	8	5	20
Wide-ball	2	10	10
Regular ball	5	5	270



Multiclass Classification:

Confusion Matrix – Macro-Averaging:

We compute performance for each class and then average.

Confusion Matrix – Each Class:



Actual

 $\frac{270}{300} = 0.90$

30

Not Regular

10

25





Multiclass Classification:

Confusion Matrix – Each Class:

No-ball

Not a no-

ball

Classifier

Output

A Not-for-Profit University

Confusion Matrix – Micro-Averaging:

- Compute confusion matrix after collecting decisions for all classes and then evaluate.

No-ball

8

7

Actual Not a

No-ball

25

295



Multiclass Classification:

Micro-Averaging vs Macro Averaging:

- Note Micro-average recall= Micro-average precision = F1 Score = Accuracy (computed from confusion matrix)
 - Micro-average is termed as a global metric.
 - Consequently, it is not a good measure when classes are not balanced.
- Macro-average is relatively a better as we can see a zoomed-in picture before averaging.
- Note Macro-averaging does not take class imbalance into account.

Weighted-average Recall:

- Weighted-averaging; Similar to Macro averaging but takes a weighted mean instead where weight for each class is the total number of data-points of that class.

 $\frac{(15 \times 0.53) + (20 \times 0.50) + (300 \times 0.90)}{(20 \times 0.50) + (300 \times 0.90)} = 0.86$



References:

- KM 5.7.2

