

Introduction to Machine Learning

Ranga **Raju** Vatsavai, Ph.D.

Chancellors Faculty Excellence Associate Professor in Geospatial Analytics
Department of Computer Science, North Carolina State University (NCSU)

Feb. 25-27, 2019

How do you find accuracy of a model

- General concepts
 - Training data (to fit a model)
 - Test data (to validate a model)

2/26/19

© Raju Vatsavai

CSC-591.2

Confusion Matrix

- Confusion Matrix:

	PREDICTED CLASS	
	Class=Yes	Class=No
	Class=Yes	Class=No
ACTUAL CLASS	a	b
	c	d

	PREDICTED CLASS	
	Class=Yes	Class=No
	Class=Yes	Class=No
ACTUAL CLASS	a (TP)	b (FN) (Type II)
	c (FP) (Type I)	d (TN)

- a: TP (true positive): Actual "Y" samples that are also classified as "Y"
- b: FN (false negative): Actual "Y" that were incorrectly classified as "N"
- c: FP (false positive): Actual "N" that were incorrectly classified as "Y"
- d: TN (true negative): Actual "N" that are correctly classified as "N"

- Most widely-used metric:

$$\text{Accuracy} = \frac{a + d}{a + b + c + d} = \frac{TP + TN}{TP + TN + FP + FN}$$

		True condition		Prevalence $= \frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	Accuracy (ACC) = $\frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$
		Condition positive	Condition negative		
Predicted condition	Predicted condition positive	True positive, Power $= \frac{\sum \text{True positive}}{\sum \text{Predicted condition positive}}$	False positive, Type I error $= \frac{\sum \text{False positive}}{\sum \text{Predicted condition positive}}$	Positive predictive value (PPV), Precision = $\frac{\sum \text{True positive}}{\sum \text{Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\sum \text{False positive}}{\sum \text{Predicted condition positive}}$
	Predicted condition negative	False negative, Type II error $= \frac{\sum \text{False negative}}{\sum \text{Predicted condition negative}}$	True negative $= \frac{\sum \text{True negative}}{\sum \text{Predicted condition negative}}$	False omission rate (FOR) = $\frac{\sum \text{False negative}}{\sum \text{Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\sum \text{True negative}}{\sum \text{Predicted condition negative}}$
		True positive rate (TPR), Recall, Sensitivity, probability of detection $= \frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False positive rate (FPR), Fail-out, probability of false alarm = $= \frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	Positive likelihood ratio (LR+) $= \frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) = $\frac{\text{LR+}}{\text{LR-}}$
		False negative rate (FNR), Miss rate $= \frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	True negative rate (TNR), Specificity (SPC) $= \frac{\sum \text{True negative}}{\sum \text{Condition negative}}$	Negative likelihood ratio (LR-) $= \frac{\text{FNR}}{\text{TNR}}$	F ₁ score = $\frac{2}{\text{Recall} + \text{Precision}}$

Source: Wikipedia

Problem with Accuracy

- Consider a 2-class problem
 - Number of Class 0 examples = 990
 - Number of Class 1 examples = 10

Problem with Accuracy

- Consider a 2-class problem
 - Number of Class NO examples = 990
 - Number of Class YES examples = 10
- If a model predicts everything to be class NO, accuracy is $990/1000 = 99\%$
 - This is misleading because the model does not detect any class YES example
 - Detecting the rare class is usually more interesting

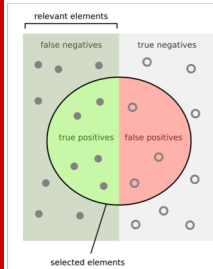
Class Imbalance Problem

- Lots of classification problems where the classes are skewed (more records from one class than another)
 - Credit card fraud
 - Intrusion detection
 - Defective products in manufacturing assembly line

Challenges

- Evaluation measures such as accuracy is not well-suited for imbalanced class
- Detecting the rare class is like finding needle in a haystack

Classification Accuracy Measures



	PREDICTED CLASS	
	Class=Yes	Class=No
	Class=Yes	Class=No
ACTUAL CLASS	a (TP)	b (FN)
	c (FP)	d (TN)

$$\text{Precision } (p) = \frac{a}{a+c} \quad (\text{PPV})$$

$$\text{Recall } (r) = \frac{a}{a+b} \quad (\text{TPR or Sensitivity})$$

$$F\text{-measure } (F) = \frac{2rp}{r+p} = \frac{2a}{2a+b+c}$$

Alternative Measures

	PREDICTED CLASS	
	Class=Yes	Class=No
	Class=Yes	Class=No
ACTUAL CLASS	10	0
	10	980

$$\text{Precision } (p) = \frac{10}{10+0} = 0.5$$

$$\text{Recall } (r) = \frac{10}{10+0} = 1$$

$$F\text{-measure } (F) = \frac{2 * 1 * 0.5}{1 + 0.5} = 0.62$$

$$\text{Accuracy} = \frac{990}{1000} = 0.99$$

Alternative Measures

	PREDICTED CLASS	
	Class=Yes	Class=No
	Class=Yes	Class=No
ACTUAL CLASS	10	0
	10	980

$$\text{Precision } (p) = \frac{10}{10+0} = 0.5$$

$$\text{Recall } (r) = \frac{10}{10+0} = 1$$

$$F\text{-measure } (F) = \frac{2 * 1 * 0.5}{1 + 0.5} = 0.62$$

$$\text{Accuracy} = \frac{990}{1000} = 0.99$$

	PREDICTED CLASS	
	Class=Yes	Class=No
	Class=Yes	Class=No
ACTUAL CLASS	1	9
	0	990

$$\text{Precision } (p) = \frac{1}{1+0} = 1$$

$$\text{Recall } (r) = \frac{1}{1+9} = 0.1$$

$$F\text{-measure } (F) = \frac{2 * 0.1 * 1}{1 + 0.1} = 0.18$$

$$\text{Accuracy} = \frac{991}{1000} = 0.991$$

Alternative Measures

		PREDICTED CLASS	
		Class=Yes	Class=No
ACTUAL CLASS	Class=Yes	40	10
	Class=No	10	40

Precision (p) = 0.8
 Recall (r) = 0.8
 F-measure (F) = 0.8
 Accuracy = 0.8

Alternative Measures

		PREDICTED CLASS	
		Class=Yes	Class=No
ACTUAL CLASS	Class=Yes	40	10
	Class=No	10	40

Precision (p) = 0.8
 Recall (r) = 0.8
 F-measure (F) = 0.8
 Accuracy = 0.8

		PREDICTED CLASS	
		Class=Yes	Class=No
ACTUAL CLASS	Class=Yes	40	10
	Class=No	1000	4000

Precision (p) = ~ 0.04
 Recall (r) = 0.8
 F-measure (F) = ~ 0.08
 Accuracy = 0.8

Measures of Classification Performance

		PREDICTED CLASS	
		Yes	No
ACTUAL CLASS	Yes	TP	FN
	No	FP	TN

α is the probability that we reject the null hypothesis when it is true. This is a Type I error or a false positive (FP).
 β is the probability that we accept the null hypothesis when it is false. This is a Type II error or a false negative (FN).

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

$$ErrorRate = 1 - accuracy$$

$$Precision = Positive Predictive Value = \frac{TP}{TP + FP}$$

$$Recall = Sensitivity = TP Rate = \frac{TP}{TP + FN}$$

$$Specificity = TN Rate = \frac{TN}{TN + FP}$$

$$FP Rate = \alpha = \frac{FP}{TN + FP} = 1 - specificity$$

$$FN Rate = \beta = \frac{FN}{FN + TP} = 1 - sensitivity$$

$$Power = sensitivity = 1 - \beta$$

Alternative Measures

		PREDICTED CLASS	
		Class=Yes	Class=No
ACTUAL CLASS	Class=Yes	10	40
	Class=No	10	40

Precision (p) = 0.5
 TPR = Recall (r) = 0.2
 FPR = 0.2

		PREDICTED CLASS	
		Class=Yes	Class=No
ACTUAL CLASS	Class=Yes	25	25
	Class=No	25	25

Precision (p) = 0.5
 TPR = Recall (r) = 0.5
 FPR = 0.5

		PREDICTED CLASS	
		Class=Yes	Class=No
ACTUAL CLASS	Class=Yes	40	10
	Class=No	40	10

Precision (p) = 0.5
 TPR = Recall (r) = 0.8
 FPR = 0.8

ROC (Receiver Operating Characteristic)

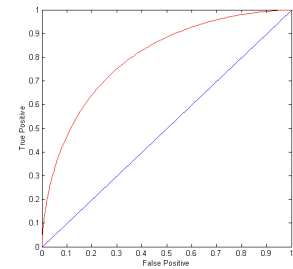
- A graphical approach for displaying trade-off between detection rate and false alarm rate
- Developed in 1950s for signal detection theory to analyze noisy signals
- ROC curve plots TPR against FPR
 - Performance of a model represented as a point in an ROC curve

ROC Curve

(TPR, FPR):

- (0,0): declare everything to be negative class
- (1,1): declare everything to be positive class
- (0,1): ideal
- Diagonal line:
 - Random guessing
 - Below diagonal line:
 - prediction is opposite of the true class

	PREDICTED CLASS		
ACTUAL CLASS		Class=Yes	Class=No
	Class=Yes	a (TP)	b (FN)
	Class=No	c (FP)	d (TN)



Handling Class Imbalance Problem

- Cost-sensitive classification
 - Misclassifying rare class as majority class is more expensive than misclassifying majority as rare class

Cost Matrix

		PREDICTED CLASS	
ACTUAL CLASS		Class=Yes	Class=No
	Class=Yes	f(Yes, Yes)	f(Yes, No)
	Class=No	f(No, Yes)	f(No, No)

$C(i,j)$: Cost of misclassifying class i example as class j

Cost Matrix	PREDICTED CLASS		
	C(i, j)	Class=Yes	Class=No
	Class=Yes	C(Yes, Yes)	C(Yes, No)
	Class=No	C(No, Yes)	C(No, No)

$$\text{Cost} = \sum C(i, j) \times f(i, j)$$

NC STATE UNIVERSITY		Computing Cost of Classification		
Cost Matrix		PREDICTED CLASS		
		C(i,j)	+	-
ACTUAL CLASS	+	-1	100	
	-	1	0	

Model M_1	PREDICTED CLASS		
		+	-
ACTUAL CLASS	+	150	40
	-	60	250

Accuracy = 80%
Cost = 3910

Model M_2	PREDICTED CLASS		
		+	-
ACTUAL CLASS	+	250	45
	-	5	200

Accuracy = 90%
Cost = 4255