

DATA-DRIVEN PROBLEM SOLVING IN MECHANICAL ENGINEERING

Classification

MASOUD MASOUMI

ME 364 - Spring 2022

Department of Mechanical Engineering
The Cooper Union for the Advancement of Science and Art

April, 2022



Classification is the problem of predicting the right label for a given input record. The task differs from regression in that labels are discrete entities not continuous function values.

Classification is a supervised learning approach and the target variable is a categorical variable.

Here are some classification algorithms:

- (a) Decision Tree (ID3, C4.5, C5.0)
- (b) Naive Bayes
- (c) Linear Discriminant Analysis
- (d) K-Nearest Neighbors
- (e) Logistic Regression
- (f) Support Vector Machine (SVM)
- (g) Neural Networks



- (a) **Jaccard Index:** defined as the size of the intersection divided by the size of the union of two label sets. It is used to compare set of predicted labels for a sample to the corresponding set of true labels.

$$J(y, \hat{y}) = \frac{|y \cap \hat{y}|}{|y \cup \hat{y}|}$$

- (b) **F-Score:** It is calculated based the Precision and Recall.

First let's define some terms corresponding to all the possible outcomes of a classifier:

True Positive (TP): Classifier labels a positive item as positive

True Negative (TN): Classifier labels a negative item as negative

False Positive (FP): Classifier labels a negative item as positive (Type I error)

False Negative (FN): Classifier labels a positive item as negative (Type II error)

Evaluation Metrics for Classification



We can have a *Confusion Matrix*

		Predicted	
		Positive	Negative
Actual	Positive	True Positive (TP)	False Negative (FN)
	Negative	False Positive (FP)	True Negative (TN)

Let's take a look at the definition of Precision

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

Precision can be thought of as a measure of exactness (what percentage of target variables labeled as positive are actually positive)



A perfect precision score of 1.0 for a class positive means that every target variable that the classifier labeled as belonging to positive does indeed belong to class positive. However, it does not tell us anything about the number of class positive that the classifier mislabeled.

looking at the definition of Recall

$$recall = \frac{TP}{TP + FN}$$

Recall is a measure of completeness (what percentage of positive target variables are labeled as such).

A perfect recall score of 1.0 for class positive means that every item from class positive was labeled as such, but it does not tell us how many other target variables were incorrectly labeled as belonging to class positive.



There tends to be an inverse relationship between precision and recall, where it is possible to increase one at the cost of reducing the other.

When a search engine returns 30 pages, only 20 of which are relevant (TP), 10 are irrelevant (FP) while failing to return 40 additional relevant pages (FN), its precision is $20/30 = 2/3$, which tells us how valid the results are, while its recall is $20/60 = 1/3$, which tells us how complete the results are.

Finally, *F-score* is a combination of precision and recall, a harmonic mean of the two:

$$F = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

(c) **Accuracy:** The ratio of the number of correct predictions over total predictions.

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