

F1 score

In statistical analysis of binary classification, the **F₁ score** (also **F-score** or **F-measure**) is a measure of a test's accuracy. It considers both the precision *p* and the recall *r* of the test to compute the score: *p* is the number of correct positive results divided by the number of all positive results returned by the classifier, and *r* is the number of correct positive results divided by the number of all relevant samples (all samples that should have been identified as positive). The F₁ score is the harmonic average of the precision and recall, where an F₁ score reaches its best value at 1 (perfect precision and recall) and worst at 0.

Contents

Etymology

Definition

Diagnostic testing

Applications

Criticism

Difference from G-measure

See also

References

Etymology

The name F-measure is believed to be named after a different F function in Van Rijsbergen's book, when introduced to MUC-4.^[1]

Definition

The traditional F-measure or balanced F-score (**F₁ score**) is the harmonic mean of precision and recall:

$$F_1 = \left(\frac{\text{recall}^{-1} + \text{precision}^{-1}}{2} \right)^{-1} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}.$$

The general formula for positive real *β* is:

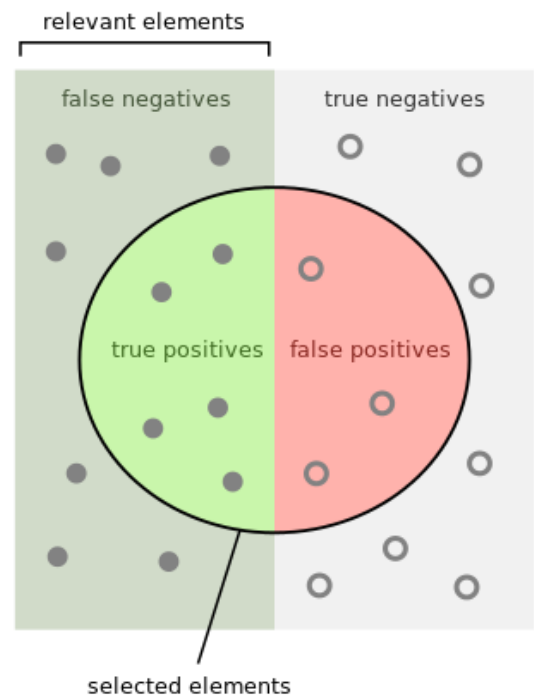
$$F_\beta = (1 + \beta^2) \cdot \frac{\text{precision} \cdot \text{recall}}{(\beta^2 \cdot \text{precision}) + \text{recall}}.$$

The formula in terms of Type I and type II errors

$$F_\beta = \frac{(1 + \beta^2) \cdot \text{true positive}}{(1 + \beta^2) \cdot \text{true positive} + \beta^2 \cdot \text{false negative} + \text{false positive}}.$$

Two other commonly used F measures are the **F₂** measure, which weighs recall higher than precision (by placing more emphasis on false negatives), and the **F_{0.5}** measure, which weighs recall lower than precision (by attenuating the influence of false negatives).

The F-measure was derived so that **F_β** "measures the effectiveness of retrieval with respect to a user who attaches *β* times as much importance to recall as precision".^[2] It is based on Van Rijsbergen's effectiveness measure



How many selected items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

Precision and recall

$$E = 1 - \left(\frac{\alpha}{p} + \frac{1 - \alpha}{r} \right)^{-1}.$$

Their relationship is $F_\beta = 1 - E$ where $\alpha = \frac{1}{1 + \beta^2}$.

The F_1 score is also known as the Sørensen–Dice coefficient or Dice similarity coefficient (DSC).

Diagnostic testing

This is related to the field of binary classification where recall is often termed as Sensitivity. There are several reasons that the F_1 score can be criticized in particular circumstances:^[3]

		True condition				
		Total population	Condition positive	Condition negative	$\text{Prevalence} = \frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	$\text{Accuracy (ACC)} = \frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$
Predicted condition	Predicted condition positive	<u>True positive</u>	<u>False positive</u> <u>Type I error</u>	$\frac{\text{Positive predictive value (PPV), Precision} = \sum \text{True positive}}{\sum \text{Predicted condition positive}}$	$\frac{\text{False discovery rate (FDR)} = \sum \text{False positive}}{\sum \text{Predicted condition positive}}$	
	Predicted condition negative	<u>False negative</u> <u>Type II error</u>	<u>True negative</u>	$\frac{\text{False omission rate (FOR)} = \sum \text{False negative}}{\sum \text{Predicted condition negative}}$	$\frac{\text{Negative predictive value (NPV)} = \sum \text{True negative}}{\sum \text{Predicted condition negative}}$	
		$\frac{\text{True positive rate (TPR), Recall, Sensitivity, probability of detection, Power} = \sum \text{True positive}}{\sum \text{Condition positive}}$	$\frac{\text{False positive rate (FPR), Fall-out, probability of false alarm} = \sum \text{False positive}}{\sum \text{Condition negative}}$	$\frac{\text{Positive likelihood ratio (LR+)} = \frac{\text{TPR}}{\text{FPR}}}$	$\frac{\text{Diagnostic odds ratio (DOR)} = \frac{\text{LR+}}{\text{LR-}}}$	$F_1 \text{ score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$
		$\frac{\text{False negative rate (FNR), Miss rate} = \sum \text{False negative}}{\sum \text{Condition positive}}$	$\frac{\text{Specificity (SPC), Selectivity, True negative rate (TNR)} = \sum \text{True negative}}{\sum \text{Condition negative}}$	$\frac{\text{Negative likelihood ratio (LR-)} = \frac{\text{FNR}}{\text{TNR}}}$		

Applications

The F-score is often used in the field of information retrieval for measuring search, document classification and query classification performance.^[4] Earlier works focused primarily on the F_1 score, but with the proliferation of large scale search engines, performance goals changed to place more emphasis on either precision or recall^[5] and so F_β is seen in wide application.

The F-score is also used in machine learning^[6] Note, however, that the F-measures do not take the true negatives into account, and that measures such as the Matthews correlation coefficient, Informedness or Cohen's kappa may be preferable to assess the performance of a binary classifier^[3]

The F-score has been widely used in the natural language processing literature, such as the evaluation of named entity recognition and word segmentation

Criticism

David Hand and others criticize the widespread use of the F-score since it gives equal importance to precision and recall. In practice, different types of misclassifications incur different costs. In other words, the relative importance of precision and recall is an aspect of the problem.^[7]

Difference from G-measure

While the F-measure is the harmonic mean of recall and precision, the G-measure is the geometric mean^[3]

See also

- BLEU
- Matthews correlation coefficient
- METEOR
- NIST (metric)
- Precision and recall

- Receiver operating characteristic
- ROUGE (metric)
- Sørensen–Dice coefficient
- Uncertainty coefficient, aka Proficiency
- Word error rate (WER)

References

1. Sasaki, Y. (2007). "The truth of the F-measure"(<https://www.toyota-ti.ac.jp/Lab/Denshi/COINpeople/yutaka.sasaki/F-measure-YS-26Oct07.pdf>) (PDF).
2. Van Rijsbergen, C. J. (1979). *Information Retrieval*(<http://www.dcs.gla.ac.uk/Keith/Preface.htm>) (2nd ed.). Butterworth-Heinemann.
3. Powers, David M W (2011). "Evaluation: From Precision, Recall and F-Measure to ROC, Informedness, Markedness & Correlation"(http://www.bioinfopublication.org/files/articles/21_1_JMLT.pdf) (PDF). *Journal of Machine Learning Technologies*. 2 (1): 37–63.
4. Beitzel., Steven M. (2006). *On Understanding and Classifying Web Queries* (Ph.D. thesis). IIT. [CiteSeerX 10.1.1.127.634](https://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.127.634) (<https://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.127.634>)
5. X. Li; Y.-Y. Wang; A. Acero (July 2008). *Learning query intent from regularized click graphs*(<https://pdfs.semanticscholar.org/6718/f8e95461456023196fe6409073151ab0513d.pdf>) (PDF). *Proceedings of the 31st SIGIR Conference*
6. See, e.g., the evaluation of the [1] (<https://dl.acm.org/citation.cfm?id=1119195>)
7. Hand, David. "A note on using the F-measure for evaluating record linkage algorithms - Dimensions"(<https://app.dimensions.ai/details/publication/pub.1084928040>) *app.dimensions.ai* Retrieved 2018-12-08.

Retrieved from 'https://en.wikipedia.org/w/index.php?title=F1_score&oldid=874064435

This page was last edited on 16 December 2018, at 22:23 (UTC).

Text is available under the [Creative Commons Attribution-ShareAlike License](#)additional terms may apply By using this site, you agree to the [Terms of Use](#) and [Privacy Policy](#). Wikipedia® is a registered trademark of the [Wikimedia Foundation, Inc](#), a non-profit organization.