

# Understanding the F1 Score: A Comprehensive Guide for Evaluating Classification Models

Authored by  
**Mohammed Iooti**

November 2, 2025

## RECOMMENDED CITATION

Mohammed Iooti (2025). *Understanding the F1 Score: A Comprehensive Guide for Evaluating Classification Models*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=8678>

When engineering sophisticated systems in [Machine Learning](#) (ML), particularly those focused on [classification](#) tasks, the need for a rigorous and reliable metric to assess model performance is paramount. While simple metrics such as overall accuracy might seem intuitive, they often fail dramatically when applied to real-world scenarios, especially those involving skewed or imbalanced datasets.

A fundamental metric designed to overcome the shortcomings of simple accuracy is the **F1 Score**. This metric provides a crucial balance, offering a single numerical representation of a model's effectiveness. By combining two critical performance measures--[Precision](#) and **Recall**--through a harmonic mean, the F1 Score becomes an indispensable tool for evaluating models where the correct identification of minority classes is essential for success.

Understanding the F1 Score is not just about memorizing a formula; it is about grasping the delicate trade-offs inherent in building predictive models. A high F1 Score indicates that the model is minimizing both false positives and false negatives simultaneously, thereby demonstrating a robust and reliable predictive capability across all classes.

## The Fundamental Role of Precision and Recall

To fully appreciate the utility of the [F1 Score](#), we must first dissect the two constituent metrics it harmonizes. Both Precision and Recall are derived directly from the outcomes recorded in the [Confusion Matrix](#), which systematically tabulates the four possible outcomes of a binary classifier: True Positives (TP), True Negatives (TN), False Positives (FP), and False Negatives (FN).

These metrics address two fundamentally different, yet equally important, questions regarding the classifier's performance. The balance achieved between them is critical, as efforts to improve one metric often result in the degradation of the other. For instance, tuning a model to be extremely sensitive (high Recall) might increase the number of incorrect positive predictions (low Precision), and vice versa. The F1 Score serves as the single optimization target for navigating this trade-off.

**Precision:** This metric quantifies the quality of the positive predictions made by the model. It is the ratio of correctly identified positive cases (True Positives) to the total number of cases the model labeled as positive (True Positives + False Positives). Precision answers the question: "When the model says something is positive, how often is it actually correct?" High precision is crucial in scenarios where minimizing false alarms or unnecessary interventions (False Positives) is the priority, such as in spam detection or criminal sentencing prediction.

**Recall** (or Sensitivity): This metric quantifies the completeness of the positive predictions. It is the ratio of correctly identified positive cases (True Positives) to the total number of actual positive cases in the dataset (True Positives + False Negatives). Recall addresses the question: "Of all the truly positive instances in the dataset, how many did the model successfully find?" High recall is vital in applications where failing to detect an actual positive case (False Negative) carries a high

cost, such as diagnosing a serious medical condition or identifying fraudulent transactions.

## Calculating the F1 Score: The Harmonic Mean

The **F1 Score** is formally defined as the weighted average, or more specifically, the harmonic mean, of Precision and Recall. The use of the harmonic mean, rather than a simple arithmetic average, is intentional and crucial to its function. The harmonic mean is mathematically structured to heavily penalize scores that are extremely low in either Precision or Recall. If a model achieves a perfect score of 1.0 in one metric but a dismal score of 0.1 in the other, the resulting harmonic mean will be significantly lower than the arithmetic mean, forcing practitioners to build models that perform moderately well across both dimensions.

The mathematical formulation for the **F1 Score** is:

$$\mathbf{F1\ Score} = 2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$$

This single metric effectively encapsulates the model's overall performance in handling positive class predictions. A high score suggests that the model is maintaining a strong balance, avoiding excessive False Positives (a sign of high Precision) while simultaneously avoiding excessive False Negatives (a sign of high Recall). It is the preferred single-figure metric for classification tasks where the class distribution is known to be imbalanced.

## Practical Application: Analyzing NBA Draft Predictions (Case Study)

To demonstrate the calculation and interpretation of the F1 Score, let us examine a practical classification scenario. Suppose we are utilizing a [Logistic Regression](#) model designed to predict whether 400 college basketball players will ultimately be selected in the NBA draft (a positive outcome) or not (a negative outcome). Since only a minority of college players are drafted, this task naturally involves class imbalance.

After running the initial model on our test data, the results are summarized in the following [Confusion Matrix](#), illustrating the distribution of the model's predictions versus the actual outcomes:

		Predicted	
		Drafted = Yes	Drafted = No
Actual	Drafted = Yes	120 (True Positive)	40 (False Negative)
	Drafted = No	70 (False positive)	170 (True Negative)

Using these results, we can methodically calculate the F1 score for this initial iteration of the model:

### Calculate Precision:

Precision is calculated as True Positive / (True Positive + False Positive). This shows that out of the 190 players predicted to be drafted, 120 actually were drafted.

$$\text{Precision} = 120 / (120 + 70) = 120 / 190 = \mathbf{0.63157}$$

### Calculate Recall:

Recall is calculated as True Positive / (True Positive + False Negative). This shows that out of the 160 players who were actually drafted, the model successfully identified 120 of them.

$$\text{Recall} = 120 / (120 + 40) = 120 / 160 = \mathbf{0.75}$$

### Calculate F1 Score:

The F1 Score combines these two metrics using the harmonic mean formula, emphasizing the lower of the two scores.

$$\text{F1 Score} = 2 * (0.63157 * 0.75) / (0.63157 + 0.75) = 2 * (0.4736775) / (1.38157) = \mathbf{0.6857}$$

An F1 score of 0.6857 provides a concrete measure of performance. While this number is now calculated, the critical question remains unanswered: Is a score of 0.6857 deemed acceptable, or is it merely mediocre? Answering this requires contextualizing the score against established benchmarks.

## Defining "Good": Contextualizing the F1 Score Threshold

Data science professionals frequently struggle with the subjective nature of evaluation metrics, specifically asking what numerical threshold defines a "good" F1 score. The fundamental truth is that the F1 Score, like most performance metrics, operates on a simple principle: **higher F1**

**scores are unequivocally better.**

The F1 Score is mathematically bounded, ranging from 0.0 to 1.0. A score of 1.0 signifies a model that has achieved perfect classification, meaning it has correctly identified every positive and negative instance without error (100% Precision and 100% Recall). Conversely, a score of 0.0 indicates a model that is completely useless, failing to classify any observation correctly. Achieving a score near 1.0 is the goal, but the acceptable level of deviation from perfection is entirely dependent on the specific context and inherent difficulty of the problem.

What constitutes a "good" score is heavily influenced by three primary factors: the domain of application, the complexity of the data, and the established industry standards. For instance, in complex, high-stakes domains such as clinical diagnostics or seismic activity prediction, where the data is noisy and the consequences of error are severe, an F1 score of 0.85 might be lauded as excellent. However, in a standardized, low-variance task like optical character recognition (OCR) or high-volume quality control, where near-perfect results are expected, a score of 0.85 might be considered inadequate and signals the need for significant model refinement.

To further illustrate the concept of perfection, consider the hypothetical scenario where our NBA draft model achieves flawless prediction across all 400 players:

		Predicted	
		Drafted = Yes	Drafted = No
Actual	Drafted = Yes	240 (True Positive)	0 (False Negative)
	Drafted = No	0 (False positive)	160 (True Negative)

In this mathematically ideal scenario, both False Positives and False Negatives are zero. The resulting metrics confirm the upper limit of the scale:

$$\text{Precision} = 240 / (240 + 0) = 1$$

$$\text{Recall} = 240 / (240 + 0) = 1$$

$$\text{F1 Score} = 2 * (1 * 1) / (1 + 1) = 1$$

This confirms that the maximum possible F1 score is indeed 1.0, achievable only when both Precision and Recall are perfect.

## Establishing the Benchmark: The Importance of a Baseline Model

Because the definition of "good" is subjective, the most crucial step in evaluating any performance

metric, including the F1 Score, is comparing the model's output against a meaningful **baseline model**. A baseline model represents the minimal acceptable performance level. It is typically the simplest, most naive model possible--often one that ignores all features and simply makes the same prediction for every observation (e.g., always predicting the majority class or making a random guess).

If our complex, feature-rich model cannot significantly outperform this simple baseline, then the effort and resources invested in building the complex model are not justified. The baseline establishes the floor; any useful model must demonstrate a substantial lift above this threshold.

Let's apply this concept to our NBA draft case study. We define a naive baseline model that assumes every single one of the 400 college players will be drafted (i.e., always predicting the positive class). The resulting Confusion Matrix for this naive approach is dramatically different:

		Predicted	
		Drafted = Yes	Drafted = No
Actual	Drafted = Yes	160 (True Positive)	0 (False Negative)
	Drafted = No	240 (False positive)	0 (True Negative)

Calculating the metrics for this naive **baseline model**:

Precision =  $160 / (160 + 240) = 160 / 400 = \mathbf{0.4}$  (Only 40% of its positive predictions were correct.)

Recall =  $160 / (160 + 0) = 160 / 160 = \mathbf{1}$  (It captured all actual drafted players, but only by predicting everyone would be drafted.)

F1 Score =  $2 * (0.4 * 1) / (0.4 + 1) = 0.8 / 1.4 = \mathbf{0.5714}$

This score of 0.5714 is the minimal acceptable standard for our draft prediction task. Now, returning to our initial Logistic Regression model, which achieved an F1 score of **0.6857**, we can make an informed judgment. While 0.6857 is indeed higher than the baseline of 0.5714, the difference is only approximately 11 percentage points. This marginal improvement suggests that while the model is technically superior to a naive guess, it still has significant deficiencies and requires further tuning, feature engineering, or the application of more advanced algorithms before it can be considered truly effective or "good" in an industrial context.

## Strategic Model Selection Using F1 Score

In practical machine learning workflows, determining a "good" F1 score is ultimately a comparative exercise. Instead of searching for an absolute target number, practitioners focus on maximizing the

F1 score relative to competing models and the established baseline. This structured, comparative approach ensures that the chosen solution is the optimal available model for the specific constraints of the problem.

The standard process for model selection based on the F1 Score involves the following disciplined steps:

### **Step 1: Fit a Naive Baseline Model.**

The initial and most critical step is to establish the floor performance. Calculate the F1 score for the simplest possible model, often one that merely predicts the majority class. This score represents the benchmark that all subsequent, complex models must demonstrably exceed.

### **Step 2: Develop and Evaluate Candidate Models.**

Fit and train a variety of candidate classification models, potentially including sophisticated algorithms such as Support Vector Machines, Random Forests, Gradient Boosting Machines, or Neural Networks. Calculate the F1 score for each model on the held-out test data, ensuring a fair comparison across all architectures.

### **Step 3: Choose the Optimal Model.**

Select the model that produces the highest F1 score among all the candidates. The decision is straightforward: the model with the highest balanced performance (best F1 score) is selected. Furthermore, rigorously verify that this optimal score represents a substantial and meaningful performance gain when compared against the F1 score generated by the initial baseline model. If the gain is minimal, the added complexity of the optimal model may not be worth the marginal benefit.

By using the F1 Score within this comparative framework, data scientists ensure that the chosen model delivers the best possible balance between Precision and Recall, making the definition of "good" directly relative to the difficulty of the problem and the best available algorithmic solutions.

## **Further Reading and Resources**

To deepen your expertise in evaluation metrics within machine learning and understand how the **F1 Score** fits into the broader landscape of model assessment, we recommend exploring these related concepts:

The relationship between the F-beta score, which is a generalized version that allows for explicitly weighting the importance of Recall versus Precision, where the F1 Score is the special case when beta equals 1.

The specific applications of [Precision](#) and [Recall](#) in highly specialized fields such as information retrieval and document ranking systems.

The utility of other graphical evaluation methods, including ROC curves and AUC (Area Under the Curve), which offer a holistic view of classifier performance across various threshold settings.