

Original Article

# Loss Functions in Artificial Intelligence and Machine Learning: A Comprehensive Overview

Rishi Mohan

<sup>1</sup>Trading Tech, Charles Schwab, Texas, United States.

Corresponding Author : [Rishi.Mohan1@schwab.com](mailto:Rishi.Mohan1@schwab.com)

Received: 30 May 2025

Revised: 22 June 2025

Accepted: 14 July 2025

Published: 28 July 2025

**Abstract** - Loss functions play a central role in machine learning by quantifying how far predicted outcomes deviate from actual values, directly influencing model optimization. This paper presents a structured and comparative overview of key loss functions used across regression and classification tasks. It evaluates regression-based losses—such as Mean Squared Error (MSE), Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), Huber Loss, and Log-Cosh—focusing on their sensitivity to outliers and interpretability. Classification loss functions—such as Cross-Entropy, Hinge Loss, Kullback-Leibler Divergence, and Focal Loss—are assessed for their effectiveness in probabilistic modeling and handling imbalanced datasets. Each loss function is presented with its mathematical formulation, practical examples, and trade-offs, followed by a comparative analysis to guide selection based on task requirements. A comparative table outlines their strengths, limitations, and ideal use cases. The paper not only demystifies the mathematical underpinnings of loss functions but also provides practical insights for selecting the appropriate loss mechanism across various machine learning contexts, with the goal of improving training efficiency and model accuracy.

**Keywords** - Loss functions, Model optimization, Regression metrics, Classification loss, Robust learning, Deep learning, Prediction error.

## 1. Introduction

Loss functions can be considered the heart and soul of Machine Learning (ML) and Artificial Intelligence (AI) systems, providing a measurable objective that guides the training of models. A loss function provides a numerical assessment of how far off a model's predictions are from the actual target values, thereby directing learning algorithms such as gradient descent in adjusting model parameters throughout training. Therefore, choosing the right loss function is crucial to ensure effective training and reliable results.

Although loss functions are fundamental to almost all supervised learning tasks, practitioners frequently face difficulties in choosing the most appropriate one for their specific problem. Most studies focus only on regression or classification losses, or they look at specific areas like image recognition or language processing. Consequently, comprehensive guidance for selecting loss functions across a wide range of tasks is still lacking.

To bridge this knowledge gap, the present study offers a comparative and task-specific analysis of widely adopted loss functions in both artificial intelligence and machine learning. It evaluates both regression and classification loss functions in

terms of mathematical formulation, sensitivity to outliers, application contexts, and performance trade-offs. In addition, the paper highlights current trends in adaptive and domain-specific loss functions, aiming to support informed decision-making in model development.

## 2. Related Work/ Literature Review

The role of loss functions in machine learning has been extensively studied in both theoretical and applied contexts. Early foundational work by Huber [5] introduced robust estimation techniques that inspired hybrid loss formulations, such as the Huber Loss, which balances sensitivity and robustness to outliers. Bishop's comprehensive text [6] explored loss functions in pattern recognition, highlighting their statistical aspects.

In recent years, advances have targeted specialized loss functions to handle challenges like class imbalance and model calibration. For instance, Lin et al. [4] proposed Focal Loss to mitigate the impact of the overwhelming majority of classes in object detection tasks, thereby improving performance on rare categories. Similarly, LeCun et al. [1] and Kingma & Ba [2] discussed optimization strategies in deep neural networks, although they provided limited guidance on selecting among various loss functions.



Despite these contributions, existing literature often isolates regression or classification losses and tends to focus on domain-specific applications such as computer vision or natural language processing. There remains a lack of comprehensive studies that systematically compare loss functions across both task types, considering practical implications and trade-offs.

This paper seeks to fill this gap by presenting an integrated review of common regression and classification loss functions, along with illustrative examples and a comparative analysis.

### 3. Types of Loss Functions

Loss functions in machine learning are usually grouped by the type of task—regression or classification. Each group includes specific loss functions tailored to improve model performance in different situations.

#### 3.1. Regression Loss Functions

Regression tasks involve predicting continuous outcomes. The aim is to reduce the discrepancy between model predictions and actual target values across all samples in a dataset. The formulas below in this section would use predicted values as  $\hat{y}_i$  and true values as  $y_i$  for a dataset of size  $n$ .

##### 3.1.1. Mean Squared Error (MSE)

MSE provides the fluctuation of a model from reality by taking the difference between the model's predicted and actual values, then squaring those differences and finding the average. Because it squares the prediction errors, MSE places greater emphasis on large discrepancies, encouraging the model to focus on reducing substantial mistakes during training.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Example:

Consider a scenario where a model is used to predict a student's final exam score based on factors like study hours, attendance, and past performance. If actual scores come as 92 and the model's prediction was 82, the error is 10. The squared error would be  $10^2 = 100$ . Squaring the error amplifies the impact of larger differences, which is why Mean Squared Error is more sensitive to large prediction mistakes.

##### 3.1.2. Mean Absolute Error (MAE)

MAE computes the average absolute difference between predicted and true values, treating all errors equally. MAE is less influenced by extreme values since it calculates the average of absolute errors without squaring them, which limits the impact of large deviations.

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Example:

Imagine you're building a model to predict the delivery time for online orders. If the actual delivery time is 30 minutes and the model predicts 28 minutes, the error is 2 minutes. Mean Absolute Error (MAE) adds up these absolute differences—just 2 minutes in this case—making it less affected by unusually large errors compared to Mean Squared Error (MSE).

##### 3.1.3. Root Mean Squared Error (RMSE)

Root Mean Squared Error takes the square root of the MSE, offering an error value that matches the original units of the target variable, which improves its readability and practical relevance. It penalizes large errors similar to MSE.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Example:

Consider a model predicting the price of different stocks over a month and the original stock price units (say \$100,000). If we have  $MSE = 250000$ , then RMSE is \$500; it means the model's average error is approximately \$500, which is more interpretable when compared to the original stock price.

##### 3.1.4. Huber Loss

Huber Loss combines the benefits of MSE and MAE, reducing sensitivity to outliers. It behaves like MSE for small errors and like MAE for large errors by introducing a threshold hyperparameter.

$$L_{\delta}(y, \hat{y}) = \begin{cases} \frac{1}{2}(y_i - \hat{y}_i)^2 & \text{for } |y_i - \hat{y}_i| \leq \delta \\ \delta|y_i - \hat{y}_i| - \frac{1}{2}\delta^2 & \text{otherwise} \end{cases}$$

Example:

Consider a scenario in which a model is trained to predict sales numbers for a business. A few days of data have massive fluctuations in sales due to seasonal effects. Huber Loss handles deviations by applying a quadratic approach to small errors and a linear one to larger errors, unlike MSE, which heavily penalizes large deviations. This makes Huber Loss less sensitive to extreme fluctuations.

##### 3.1.5. Log-Cosh Loss

It can also be visualized as Huber loss, but it makes optimization easier by handling extreme errors in a differentiable manner.

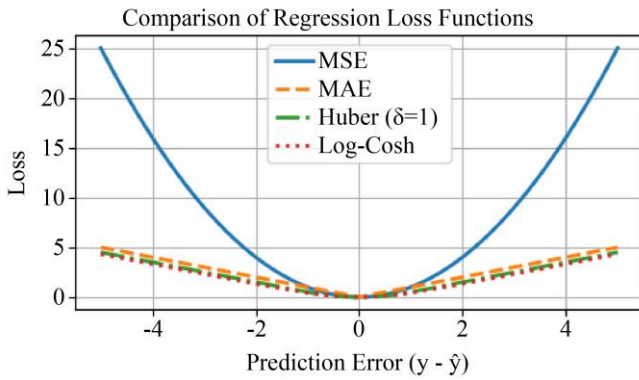
$$\text{Log-Cosh Loss} = (y, \hat{y}) = \sum_{i=1}^n \log(\cosh(y_i - \hat{y}_i))$$

Example:

Consider a model predicting daily temperatures. Suppose the actual temperature is 30°C. If the model predicts 29°C, the error is small (1°C), and Log-Cosh Loss treats it similarly to MSE. If the predicted temperature is 20 degrees and the error is 10 degrees. Unlike MSE, which heavily penalizes large errors by squaring them ( $10^2 = 100$ ), Log-Cosh Loss handles these extreme errors more smoothly and in a fully differentiable way, similar to Huber Loss but with easier optimization. This results in smoother optimization and helps prevent extreme weight updates caused by outliers.

### 3.1.6. Comparative Behaviour

The comparative behaviour of common regression loss functions is summarized to highlight their sensitivity to prediction errors and outliers.



## 3.2. Classification Loss Functions

Classification tasks involve predicting discrete labels. Common loss functions for classification include:

### 3.2.1. Cross-Entropy Loss (Log Loss)

It is widely used in neural networks for handling multi-class classification problems. It works by comparing the predicted probability distribution to the actual class labels, helping the model learn to assign higher confidence to correct predictions.

$$\text{Log Loss} = - \sum_{i=1}^n y_i \log(\hat{y}_i)$$

Here,  $y_i$  is true probabilities and  $\hat{y}_i$  is the predicted probability.

Example:

Imagine a model classifying whether an email is spam or not. If the model predicts a 90% probability that an email is spam, but it is actually not spam, the loss would be high, since cross-entropy penalizes such incorrect high-probability predictions. On the other hand, if the model predicted 50% for both classes, the penalty would be smaller, but the classification would still be wrong.

### 3.2.2. Hinge Loss

Hinge loss encourages maximizing the margin between classes by penalizing incorrect classifications. Primarily used in Support Vector Machines (SVMs).

$$\text{Hinge loss} = \max(0, 1 - y_i \hat{y}_i)$$

Here,  $y_i$  is -1 or 1 and  $\hat{y}_i$  is predicted value.

Example:

In a spam detection system using an SVM model, suppose an email that is actually "spam" is classified with a score of 0.6 for "not spam" and 0.4 for "spam." Since the model assigns a higher score to the wrong class, hinge loss applies a penalty based on the margin between the correct and incorrect classifications. If the model correctly assigns a higher score to "spam," the margin would be positive, and no loss would be applied. Misclassifications result in a loss that increases as the predicted score moves further away from the correct label.

### 3.2.3. Kullback-Leibler (KL) Divergence

The Kullback-Leibler (KL) divergence quantifies how one probability distribution diverges from another, comparing distributions PPP and QQQ. It is commonly applied in models that rely on probability distributions, such as those used in variational inference techniques.

$$D_{KL}(P \parallel Q) = \sum_{i=1}^n P(i) \log \frac{P(i)}{Q(i)}$$

Example:

For a language model predicting the next word in a sentence, KL Divergence can measure how much the predicted word distribution differs from the true distribution. When the actual distribution strongly favors a particular word (e.g., 80%), but the model assigns it a much lower probability (e.g., 50%), KL Divergence captures this gap, signaling the model to adjust its predictions toward the more accurate distribution.

### 3.2.4. Focal Loss

Created to handle imbalanced datasets, it reduces the influence of correctly predicted examples and places greater focus on harder, misclassified cases by adjusting the traditional cross-entropy formula.

$$FL(p_t) = -\alpha_t (1 - p_t)^\gamma \log(p_t)$$

Where  $p_t$  is the model's estimated probability for the true class,  $\alpha_t$  balances class importance, and  $\gamma$  adjusts focusing strength.

Example:

Imagine developing a facial recognition system that often fails to match known individuals, resulting in frequent missed identifications.

Using focal loss, which down-weights easy-to-classify examples, you focus the model's learning on difficult-to-classify instances, helping the model improve its performance in identifying people correctly, despite the imbalance in the dataset.

#### 4. Considerations for Selecting Loss Functions

Choosing an appropriate loss function is essential for effectively directing model training and evaluating its performance. Choosing the right loss function is influenced by various considerations, including:

- **Nature of the Problem:** The choice of loss function depends on whether the model is solving a regression task (predicting continuous values) or a classification task (predicting categories). Applying an inappropriate loss type can hinder training and reduce prediction reliability.
- **Computational Efficiency:** Some functions are computationally expensive or require additional hyperparameter tuning (e.g., Focal Loss, Huber Loss), which can impact training time and scalability.
- **Uneven Class Distribution:** When dealing with classification problems where certain classes are underrepresented, using specialized loss functions such as Focal Loss or class-weighted Cross-Entropy can help the

model focus more on correctly predicting the minority class.

- **Context-Specific Needs:** Different fields, like healthcare or finance, often require tailored loss functions that prioritize factors like accuracy, clarity of results, or resilience to prediction errors, depending on the critical nature of decisions made from the model's output.
- **Result Interpretability:** Loss metrics like RMSE express errors using the same scale as the predicted variable, which enhances their readability and makes them more accessible when communicating findings to non-technical audiences.
- **Resistance to Outliers:** Certain loss functions—such as MAE, Huber, and Log-Cosh—offer greater stability when working with noisy or imperfect data, as they limit the influence of extreme deviations. On the other hand, metrics like MSE and RMSE can amplify the impact of large errors, making them less ideal for outlier-prone datasets.

#### 5. Comparative Analysis of Loss Functions

The following table summarizes key properties of commonly used loss functions:

Table 1. Comparative analysis of key loss functions

Loss Function	Type	Sensitivity To Outliers	Common Applications
MSE	Regression	High	General Predictive Models
MAE	Regression	Low	Robust Regression Models
RMSE	Regression	High	Interpretable Error Metrics
Huber	Regression	Moderate	Noise Resistant Models
Log-Cosh	Regression	Low	Smooth Optimization
Cross Entropy	Classification	High	Neural Networks, Deep Learning
Hinge	Classification	High	SVMs
KL Divergence	Classification	High	Probabilistic Models
Focal Loss	Classification	Low	Class-Imbalanced Learning

#### 6. Future Directions

As machine learning systems continue to evolve, so does the need for more adaptive, robust, and context-aware loss functions. Ongoing research is investigating loss functions that can adjust themselves in real time based on the evolving patterns in the data during training. These dynamic approaches reduce the dependence on manual tuning and help improve the model's performance on new, unseen inputs.

In addition, there is growing interest in domain-specific loss functions, particularly in fields such as healthcare, autonomous systems, and finance, where conventional loss functions may not reflect real-world costs or constraints. Hybrid loss functions—combining features of multiple traditional losses—are also being explored to balance robustness, sensitivity, and interpretability. Future studies could aim to systematically evaluate different loss functions using diverse, real-world datasets and applications, helping to create more concrete recommendations for their optimal use.

Advances in explainable AI (XAI) may integrate better interpretability into loss function design, promoting transparency and trust in predictive systems.

#### 7. Conclusion

This paper presented a comprehensive overview of both regression and classification loss functions, highlighting their mathematical formulations, application contexts, and sensitivity to data characteristics. Unlike prior studies that focus on narrow tasks or individual loss functions, this work offers a unified comparison framework supported by practical examples and a comparative table. It helps both researchers and professionals gain deeper insight into how different loss function selections influence model training results and overall effectiveness.

Although this review offers conceptual and practical insights, it does not include empirical evaluations across

datasets. Future work may extend this study through experimental benchmarks to further guide the selection of loss functions in applied machine learning.

This study links core theoretical principles with hands-on implementation, helping to guide the development of more effective and streamlined models in artificial intelligence and machine learning.

## References

- [1] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton, "Deep Learning," *Nature*, vol. 521, pp. 436-444, 2015. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [2] D. P. Kingma, and J. Ba, "Adam: A Method for Stochastic Optimization," *arXiv:1412.6980*, pp. 1-15, 2014. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [3] Sebastian Ruder, "An Overview of Gradient Descent Optimization Algorithms," *arXiv:1609.04747*, 2016. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [4] Tsung-Yi Lin et al., "Focal Loss for Dense Object Detection," *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2980-2988, 2017. [[Google Scholar](#)] [[Publisher Link](#)]
- [5] Peter J. Huber, "Robust Estimation of a Location Parameter," *The Annals of Mathematical Statistics*, vol. 35, no. 1, pp. 73-101, 1964. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [6] Christopher M. Bishop, *Pattern Recognition and Machine Learning*, New York, NY, USA: Springer, 2006. [[Google Scholar](#)] [[Publisher Link](#)]