# Loss Functions for Machine Learning

| Function | Usage | Library / Implementation |
|---|---|---|
| **Mean Absolute Error (MAE)**<br><br>**(L1 Norm Loss)**<br><br>$$\mathbf{MAE} = \left(\frac{1}{n}\right) * \sum (y_i - \bar{y})$$ | **REGRESSION**<br><br>Less sensitive to outliers<br><br>Easily interpretable<br><br>Has constant gradient, slower convergence with gradient optimizing methods | Scikit Learn :<br>scikit_learn.metrics.mean_absolute_error<br><br>PyTorch :<br>mae = torch.nn.L1Loss( ) |
| **Mean Squared Error (MSE)**<br><br>**(L2 Norm Loss)**<br><br>$$\mathbf{MSE} = \left(\frac{1}{n}\right) * \sum (y_i - \bar{y})^2$$ | **REGRESSION**<br><br>Sensitive to large errors/outliers due to quadratic nature.<br><br>Has smoother gradient. | Scikit Learn :<br>scikit_learn.metrics.mean_squared_error<br><br>PyTorch :<br>mse_loss = torch.nn.MSELoss() |
| **Smooth Mean Squared Error**<br><br>**(Huber Loss)**<br><br>Large errors :<br>$$L(\delta, y, f(x)) = \delta * |f(x) - y| - \left(\frac{1}{2}\right) * \delta^2$$<br>Small errors:<br>$$L(\delta, y, f(x)) = \left(\frac{1}{2}\right) * (f(x) - y)^2$$ | **REGRESSION**<br><br>Combines advantages of MSE and MAE.<br><br>Handles large and small errors differently based on parameter delta.<br><br>Has medium impact with outliers. | Scikit Learn :<br>from sklearn.linear_model import HuberRegressor<br><br>PyTorch :<br>huber_loss = torch.nn.SmoothL1Loss()<br><br>Tensorflow :<br>huber_loss = tf.keras.losses.Huber() |
| **Binary Cross-Entropy Loss**<br><br>**( Log Loss )**<br><br>$$L(y, f(x)) = -[y * log(f(x)) + (1 - y) * log(1 - f(x))]$$ | **BINARY  CLASSIFICATION**<br><br>Mostly used for classifying elements into two classes.<br><br>Ideal for models which output probabilities. | Scikit Learn :<br>from sklearn.metrics import log_loss<br>PyTorch :<br>bce_loss = nn.BCELoss() OR<br>bce_logits_loss = nn.BCEWithLogitsLoss()<br>Tensorflow :<br>bce = tf.keras.losses.BinaryCrossentropy() |
| **Categorical Cross -Entropy Loss**<br><br>**( Softmax Loss )**<br><br>$$Loss = -\sum_{i=1}^{n} \sum_{c=1}^{C} y_{i,c} \log(p_i, c)$$ | **MULTI CLASS CLASSIFICATION**<br><br>Ideal for models which output probabilities across various categories.<br><br>Used for classifying elements into multiple | Scikit Learn :<br>from sklearn.metrics import log_loss<br><br>PyTorch :<br>cross_entropy_loss = torch.nn.CrossEntropyLoss()<br><br>Tensorflow :<br>cross_entropy_loss = tf.keras.losses.CategoricalCrossentropy() |

| | classes, when labels are one-hot encoded. | |
|---|---|---|
| **Sparse Categorical Cross -Entropy Loss** $$Loss = -\frac{1}{n}\sum_{i=1}^{n} log(p_i, y_i)$$ | **MULTI CLASS CLASSIFICATION** Used for classifying elements into multiple classes, when labels are integers, not one-hot encoded. | PyTorch : cross_entropy_loss = torch.nn.CrossEntropyLoss() Tensorflow : sparse_cross_entropy_loss = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True) |
| **Hinge Loss** $$Loss(y, f(x)) = max(0, 1 - y.f(x))$$ | **CLASSIFICATION (ESPECIALLY SVMs)** Focuses on maximizing the margin between data points (classes) and decision boundary. | Scikit Learn : from sklearn.metrics import hinge_loss PyTorch : hinge_loss = nn.HingeEmbeddingLoss() Tensorflow : loss=tf.keras.losses.Hinge() |
| **K-L Divergence Loss** $$D_{KL}(P||Q) = \sum_{i} P_i.log(\frac{P(i)}{Q(i)})$$ | **FOR PROBABILITY DISTRIBUTIONS** Measure difference between two probability distributions. Used for variational autoencoders, etc. | Scipy : kl_divergence = numpy.sum(scipy.special.rel_entr(P, Q)) PyTorch : kl_divergence = torch.nn.Functional.kl_div(P.log(), Q) Tensorflow : from tensorflow.keras.losses import KLDivergence |
| **Cosine Similarity Loss** $$L(a, b) = -\frac{a.b}{||a||||b||}$$ | **COMPARING SIMILARITY OF VECTORS** Aims to maximize the cosine similarity between predicted and target vectors. Used in natural language processing, recommendation | Scikit Learn : from sklearn.metrics.pairwise import cosine_similarity PyTorch / TensorFlow : Use custom function |
| **Adversarial Loss** | **FOR GENERATIVE ADVERSARIAL NETWORKS** Used for training GANs, where discriminator tries to maximize loss function and generator tries to minimize it. | Defined specifically for the task. |