

Mean Squared Error (MSE)

Formula

Mean Squared Error (MSE) can be written formally as:

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

where y_i represents the ground-truth or observed values, and \hat{y}_i denotes the predictions generated by our model for each of the n data points. We use J_{MSE} as our loss (or cost) function, though any notation like J or C would be acceptable in principle.

When we compute this average of squared differences, we obtain a measure of how far our predictions are from the true targets on average. It is a straightforward metric that sums up the squared errors and then normalizes by the number of samples, n . Although this might seem like a simple arithmetic formula, it has several important implications:

1. **Squared errors penalize large deviations more strongly than smaller ones.** This means that if a single prediction is far off from its target value, its squared difference will be disproportionately large and can dominate the overall cost.
2. **The function is continuous and differentiable with respect to \hat{y}_i .** As a result, MSE is very convenient for gradient-based optimization methods in machine learning. In practice, we simply compute partial derivatives and backpropagate them.
3. **It naturally corresponds to a model assumption that errors are drawn from a Gaussian (Normal) distribution.** From a probabilistic standpoint, minimizing MSE is equivalent to maximizing the likelihood of the data if errors are Gaussian with zero mean.

When MSE Is Useful

MSE is frequently chosen in regression tasks, especially when we expect data to be free of significant outliers or when we value penalizing large errors in a more pronounced way than smaller ones. Since the squaring process magnifies larger deviations, MSE encourages models to focus on correcting substantial mistakes.

Characteristics in More Detail

While MSE is often described with bullet points, we can delve into each characteristic in a more narrative style:

Differentiability and ease of computation: MSE can be computed by summing and squaring differences, which is computationally inexpensive. Because it is a smooth function, it fits neatly into optimization pipelines that rely on calculating gradients, such as stochastic gradient descent.

Sensitivity to outliers: By squaring the residuals, MSE magnifies the impact of points that are far from the model's predictions. This quality can be advantageous if we specifically want to ensure that large errors are minimized. However, it can become problematic in datasets containing many extreme or anomalous values, as these outliers can dominate the training process.

Encouraging predictions to cluster around the mean: When using MSE, the optimal single-point estimate (for a unimodal, symmetric distribution) is often the mean of the observed data. This results from the fact that the partial derivatives push the average prediction toward the average target.

Typical usage in continuous regression scenarios: Because MSE deals with squared differences, it is best suited to numeric, continuous outcomes. If you are modeling categories or ordinal variables, there might be other metrics (like cross-entropy or MAE) that are more appropriate.

By recognizing these details, practitioners can decide whether MSE is the right metric for their task, or whether alternative losses like Mean Absolute Error (MAE), Huber Loss, or others might better address their requirements.

Interpretations

Gaussian Distribution

The probability density function (PDF) of a Gaussian (Normal) distribution with mean μ and variance σ^2 for a variable x is:

$$\mathcal{L}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right),$$

where we use \mathcal{L} to indicate the likelihood function.

Probabilistic Perspective (Derivation)

Assuming the prediction errors follow a Gaussian distribution with mean 0 and variance σ^2 , the likelihood (\mathcal{L}) for a single observation y_i given the predicted value \hat{y}_i is:

$$\mathcal{L}(y_i | \hat{y}_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \hat{y}_i)^2}{2\sigma^2}\right).$$

1. Write out the PDF:

$$\mathcal{L}(y_i | \hat{y}_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \hat{y}_i)^2}{2\sigma^2}\right).$$

2. Take the logarithm:

$$\log \mathcal{L}(y_i | \hat{y}_i) = \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) + \log\left(\exp\left(-\frac{(y_i - \hat{y}_i)^2}{2\sigma^2}\right)\right).$$

3. Use log rules:

$$\log \mathcal{L}(y_i | \hat{y}_i) = -\frac{1}{2} \log(2\pi\sigma^2) - \frac{(y_i - \hat{y}_i)^2}{2\sigma^2}.$$

4. Multiply by -1 to define the negative log-likelihood (our cost J):

$$J_i = -\log \mathcal{L}(y_i | \hat{y}_i) = \frac{(y_i - \hat{y}_i)^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2).$$

Single-Point vs. Full Likelihood It's common to discuss a single data point's probability density as a likelihood in a derivation. Formally, the word likelihood typically refers to the joint function of all data points, but each individual term in that product (or each individual summand in the log-likelihood) can also be referred to as the likelihood contribution of a single observation. Thus, we can show the negative log-likelihood derivation on a per-data-point basis, then multiply (or sum in log space) to account for the entire dataset.

Negative Log-Likelihood for the Entire Dataset

For n independent observations, the joint likelihood is the product of individual likelihoods. Consequently, the negative log-likelihood for all observations becomes the sum:

$$J_{\text{NLL}} = -\log \mathcal{L}(\{y_i\} \mid \{\hat{y}_i\}) = \sum_{i=1}^n \left[\frac{(y_i - \hat{y}_i)^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2) \right].$$

Note on Curly Braces ($\{\}$): When we write $\{y_i\}$ or $\{\hat{y}_i\}$, we're indicating the entire set (or collection) of y_i or \hat{y}_i values for $i = 1, \dots, n$. This is standard mathematical notation for describing a group of elements, rather than a single element.

Ignoring the constant term $\frac{1}{2}n \log(2\pi\sigma^2)$ gives:

$$J_{\text{NLL}} = \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + (\text{const.}).$$

A Note on the $\frac{1}{2\sigma^2}$ Factor

From a purely optimization perspective (where σ^2 is fixed and we only optimize \hat{y} or model parameters), this constant factor does not affect the location of the minimum. So, many treatments omit it when they only care about the solution for \hat{y} . However, when the variance σ^2 is also being optimized (as in a full maximum-likelihood approach for both mean and variance), we need to keep that factor to reflect the exact Gaussian log-likelihood.

Minimizing this with respect to \hat{y}_i is equivalent to minimizing the sum of squared errors. Hence, **minimizing the negative log-likelihood** under these Gaussian assumptions is equivalent to **minimizing MSE**.

Geometric Perspective (Derivation)

The MSE can also be interpreted as the squared Euclidean (L2) distance between the predicted and true values:

$$\|y - \hat{y}\|_2^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$

Minimizing this distance is equivalent to finding the point \hat{y} closest (in an L2 sense) to the actual target y .

Connection to Linear Regression

MSE is the standard cost function in Ordinary Least Squares (OLS) regression. In linear regression, the model is often expressed as:

$$y = X\beta + \varepsilon,$$

where X is the design matrix of input features, β is the parameter vector, and ε represents normally distributed noise. Minimizing the MSE cost:

$$J_{\text{MSE}}(\beta) = \sum_{i=1}^n (y_i - X_i\beta)^2$$

yields the OLS estimates. If the columns of X are linearly independent, the closed-form solution is:

$$\hat{\beta} = (X^T X)^{-1} X^T y.$$

Detailed Derivation of the OLS Solution

Here is a step-by-step outline of how we arrive at the formula $\hat{\beta} = (X^T X)^{-1} X^T y$ in linear algebra terms:

1. Set up the cost function:

$$J_{\text{MSE}}(\beta) = \sum_{i=1}^n (y_i - X_i \beta)^2.$$

In matrix form, if y is the $n \times 1$ vector of targets and X is the $n \times d$ matrix of features (with each row corresponding to one observation), then:

$$J_{\text{MSE}}(\beta) = (y - X\beta)^T (y - X\beta).$$

2. Take the gradient w.r.t. β : We can expand the above cost:

$$(y - X\beta)^T (y - X\beta) = y^T y - y^T X\beta - (X\beta)^T y + (X\beta)^T (X\beta).$$

Notice that $y^T X\beta$ and $(X\beta)^T y$ are scalars (single numbers), so they are equal. Thus:

$$(y - X\beta)^T (y - X\beta) = y^T y - 2y^T X\beta + \beta^T X^T X\beta.$$

Now take the derivative (gradient) w.r.t. β :

$$\nabla_{\beta} J_{\text{MSE}}(\beta) = -2X^T y + 2X^T X\beta.$$

(We use basic rules of matrix calculus here, noting that the derivative of $\beta^T A\beta$ w.r.t. β is $2A\beta$ if A is symmetric. Here, $X^T X$ is symmetric.)

3. Set the gradient to zero (the Normal Equations): To minimize $J_{\text{MSE}}(\beta)$, we set its gradient to 0:

$$-2X^T y + 2X^T X\hat{\beta} = 0 \quad \Rightarrow \quad X^T X\hat{\beta} = X^T y.$$

This equation is known as the *Normal Equation*.

4. Solve for $\hat{\beta}$: Provided $X^T X$ is invertible (which requires that X has full column rank), we can multiply both sides by $(X^T X)^{-1}$:

$$\hat{\beta} = (X^T X)^{-1} X^T y.$$

5. Interpretation:

- $X^T X$ captures the correlations among features.
- $X^T y$ captures how each feature relates to the target vector.
- $(X^T X)^{-1} X^T y$ thus gives the coefficients that minimize the overall sum of squared errors.

This derivation relies on understanding how to take matrix derivatives and the condition that $X^T X$ be invertible. When these assumptions hold, the formula neatly expresses the solution that best fits the data in the least squares sense.

Summation-Based Derivation (Without Matrix Algebra)

While matrix notation is concise, we can also derive the same result using summations explicitly. Let's assume our model is:

$$\hat{y}_i = \beta_0 + \sum_{j=1}^d \beta_j x_{ij},$$

where x_{ij} is the value of the j -th feature for the i -th data point, β_0 is an intercept term, and β_j are the feature coefficients.

Our goal is to minimize the MSE cost function:

$$J_{\text{MSE}}(\beta_0, \beta_1, \dots, \beta_d) = \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$

Substituting \hat{y}_i :

$$J_{\text{MSE}} = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij} \right)^2.$$

1. Take partial derivatives w.r.t. each β_k For β_0 :

$$\frac{\partial J_{\text{MSE}}}{\partial \beta_0} = \sum_{i=1}^n -2 \left(y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij} \right) \cdot 1 = 0$$

(when set to 0 for the optimum). Simplify to get:

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij} \right) = 0.$$

For each β_k with $k \geq 1$:

$$\frac{\partial J_{\text{MSE}}}{\partial \beta_k} = \sum_{i=1}^n -2 \left(y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij} \right) x_{ik} = 0.$$

2. System of $d+1$ equations These partial derivatives give us $d+1$ simultaneous equations (one for β_0 , and one for each β_k , $k = 1, \dots, d$). In compact form:

$$\begin{cases} \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij}) = 0, \\ \sum_{i=1}^n x_{i1} (y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij}) = 0, \\ \vdots \\ \sum_{i=1}^n x_{id} (y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij}) = 0. \end{cases}$$

Solving this system yields the same result as applying the matrix-based Normal Equation: each equation here corresponds to a row in $X^T X \hat{\beta} = X^T y$.

Hence, whether we choose the summation approach (manually solving these $d+1$ equations) or the matrix approach (using linear algebra), we arrive at the same solution. The matrix form is simply a more compact and elegant representation of these summations.

Gradient-Based Optimization

Because J_{MSE} is smooth and differentiable, it fits perfectly in gradient-based methods. The partial derivative of J_{MSE} with respect to the prediction \hat{y}_i is:

$$\frac{\partial J_{\text{MSE}}}{\partial \hat{y}_i} = \frac{2}{n} (\hat{y}_i - y_i).$$

When using a parametric model, this gradient is propagated back through the model parameters (e.g., weights in neural networks) by the chain rule.

Relationship to R^2

Minimizing J_{MSE} is closely related to maximizing the coefficient of determination R^2 , which measures how well the model explains the variance in the observed data. The R^2 value is defined (in one of its forms) as:

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

where \bar{y} is the mean of the true targets y_i . Reducing $\sum_{i=1}^n (y_i - \hat{y}_i)^2$ (i.e., the SSE) will increase R^2 . A higher R^2 indicates a tighter fit of the model to the data.