

Model Evaluation Metrics in Machine Learning

Model Evaluation:

- Important component in Data Science lifecycle which occurs post Model Training.
- Assess performance of a model on a dataset.
- Determine how well the model is generalize to unseen data.
- To select best performing model.

Model Evaluation Metrics:

1. Classification:

- **Accuracy**
- **Precision** (Positive prediction Value)
- **Recall** (True Positive Rate or Sensitivity)
- **F1-Score**
- **Confusion Matrix**
- **AUC-ROC**

2. Regression:

- **Mean Squared Error**
- **Root Mean Squared Error**
- **Mean Absolute Error**
- **R-Squared Score (R²)**
- **Adjusted R-Squared Score**

3. Clustering:

- **Silhouette Coefficient**

Classification:

Confusion Matrix: A Table that summarizes model predictions against Actual values (ground truth)

Predictions	Actual Values	
	Positive	Negative
Positive	True Positive (TP)	False Positive (FP)
Negative	False Negative (FN)	True Negative (TN)

Accuracy: Measures overall Correctness of the Model.

$$Accuracy = \frac{(TP + TN)}{(TP + TN + FP + FN)}$$

Note: Useful for only Balanced Dataset.

Precision: Measures ability of Model to identify True Positives out of all positive predictions.

$$Precision = \frac{TP}{(TP + FP)}$$

1. It should be used when False Positives are costly.
2. It ensures our positive predictions are correct.

Example: Spam Detection

Recall: Measures ability of Model to identify True Positives out of all Actual Positives.

$$\text{Recall} = \frac{TP}{(TP + FN)}$$

1. It should be used when False Negatives are costly.
2. It ensures to capture all the positive instances.

Example: Cancer Detection

F1-Score: Harmonic mean of Precision and Recall.

$$F1 - Score = \frac{2 \times precision \times recall}{(precision + recall)}$$

1. Useful for imbalanced data.
2. Should use in applications where Precision and Recall both are important.

Example: Fraud Detection

AUC – ROC (Area Under Receiver Operating Characteristic Curve):

AUC - Degree or Measure of Separability.

ROC - Probability Curve.

1. Measure's ability of a model to distinguish between positive and Negative instances.
2. It is a plot between True Positive rate (TPR) and False Positive Rate (FPR)
3. Useful for imbalance dataset.

4. Values ranges from 0 to 1. Where 0 is worse and 1 is best.
5. 0.5 means Model has no class separation capability.

Example: Medical Diagnosis.

Regression:

Mean Squared Error (MSE): Average of squared differences between prediction and actual values.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y - y_{pred})^2$$

1. It is useful when we want to penalize large errors more heavily.
2. Smaller the MSE, Better the fit.

Example: Stock Price prediction

Root Mean Squared Error (RMSE): Square root of Average of squared differences between prediction and actual values.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y - y_{pred})^2}$$

1. It is useful when we want to evaluate model performance in terms of units of Target variable.
2. Easy to interpret.

Example: Plant Height (It tells avg difference in cms between prediction and Actual heights.)

Mean Absolute Error (MAE): Average of absolute differences between prediction and actual values.

$$MAE = \frac{1}{n} \sum_{i=1}^n |y - y_{pred}|$$

1. It is useful when we want to know average error in terms of Target variable units.
2. Less sensitive to Outliers.

Example: House Price prediction.

R-Squared Score (R2): It measures the proportion of variance in the target variable that can be explained by the model.

$$R2 = 1 - \frac{ss_{res}}{ss_{tot}}$$

$$ss_{res} = \sum_{i=1}^n (y - y_{pred})^2 \quad ss_{tot} = \sum_{i=1}^n (y - y_{mean})^2$$

1. It determines overall performance of the model which ranges from 0 to 1.
2. Near to 1, Better the Model.

Problems with R2:

1. R2 Score doesn't have to do with correlation between independent and dependent features. It simply increases whenever we add new feature to the Model.
2. Because, The Linear Regression tries to assign coefficients in such a way that ss_res always decreases.

Adjusted R2:

1. It is same as R2 Score, but it also considers number of independent variables are used to predict target variable.
2. By this we can determine whether adding new variables to the model actually increases the model fit or not.

$$\text{Adjusted } R2 = 1 - \frac{(1 - R2)(N - 1)}{(N - P - 1)}$$

- N - Total number of datapoints.
- P - Number of independent features.
- R2 - R2_Score determined by the model.

It penalizes when features are not correlated to dependent variable.

$$\text{Adjusted_}R2 \leq R2$$

Clustering:

Silhouette Coefficient: It measures the similarity of an object to its own cluster compared to other clusters.

$$\textit{Silhouette Score} = \frac{(b - a)}{\max(a, b)}$$

- a – average intra-cluster distance (avg distance between each point within a cluster)
 - b - average inter-cluster distance (avg distance between all the clusters)
1. It ranges from -1 to 1.
 - 1 - Clusters are well apart from each other.
 - 0 - Distance between clusters are insignificant.
 - -1 - Clusters assigned in wrong way.
 2. It helps to evaluate quality of clusters and determine optimal number of clusters.

***** Thank You *****