

Model evaluation metrics and methods		
Method Name	Description	Code Syntax
classification_report	Generates a report with precision, recall, F1 score, and support for each class in classification problems. Useful for model evaluation. <b>Hyperparameters:</b> target_names: List of labels to include in the report. <b>Pros:</b> Provides a comprehensive evaluation of classification models. <b>Limitations:</b> May not provide enough insight for imbalanced datasets.	<pre>from sklearn.metrics import classification_report # y_true: True labels # y_pred: Predicted labels # target_names: List of target class names report = classification_report(y_true, y_pred, target_names=["class1", "class2"])</pre>
confusion_matrix	Computes a confusion matrix to evaluate the classification performance, showing counts of true positives, false positives, true negatives, and false negatives. <b>Hyperparameters:</b> labels: List of class labels to include. <b>Pros:</b> Essential for understanding classification errors. <b>Limitations:</b> Doesn't give insights into prediction probabilities.	<pre>from sklearn.metrics import confusion_matrix # y_true: True labels # y_pred: Predicted labels conf_matrix = confusion_matrix(y_true, y_pred)</pre>
mean_squared_error	Calculates the mean squared error (MSE), a common metric for regression models. Lower values indicate better performance. <b>Hyperparameters:</b> sample_weight: Weights to apply to each sample. <b>Pros:</b> Simple and widely used metric. <b>Limitations:</b> Sensitive to outliers, as large errors are squared.	<pre>from sklearn.metrics import mean_squared_error # y_true: True values # y_pred: Predicted values # sample_weight: Optional, array of sample weights mse = mean_squared_error(y_true, y_pred)</pre>
root_mean_squared_error	Calculates the root mean squared error (RMSE), which is the square root of the MSE. RMSE gives more interpretable results as it is in the same units as the target. <b>Hyperparameters:</b> sample_weight: Weights to apply to each sample. <b>Pros:</b> More interpretable than MSE. <b>Limitations:</b> Like MSE, it can be sensitive to large errors and outliers.	<pre>from sklearn.metrics import root_mean_squared_error # y_true: True values # y_pred: Predicted values # sample_weight: Optional, array of sample weights rmse = root_mean_squared_error(y_true, y_pred)</pre>
mean_absolute_error	Measures the average magnitude of errors in predictions, without considering their direction. Useful for understanding the average error size. <b>Hyperparameters:</b> sample_weight: Optional sample weights. <b>Pros:</b> Less sensitive to outliers compared to MSE. <b>Limitations:</b> Does not penalize large errors as much as MSE or RMSE.	<pre>from sklearn.metrics import mean_absolute_error # y_true: True values # y_pred: Predicted values mae = mean_absolute_error(y_true, y_pred)</pre>
r2_score	Computes the coefficient of determination (R <sup>2</sup> ), which represents the proportion of variance explained by the model. A higher value indicates a better fit. <b>Pros:</b> Provides a clear indication of model performance. <b>Limitations:</b> Doesn't always represent model quality, especially for non-linear models.	<pre>from sklearn.metrics import r2_score # y_true: True values # y_pred: Predicted values r2 = r2_score(y_true, y_pred)</pre>
silhouette_score	Measures the quality of clustering by assessing the cohesion within clusters and separation between clusters. Higher scores indicate better clustering. <b>Hyperparameters:</b> metric: Distance metric to use. <b>Pros:</b> Useful for validating clustering performance. <b>Limitations:</b> Sensitive to outliers and choice of distance metric.	<pre>from sklearn.metrics import silhouette_score # X: Data used in clustering # labels: Cluster labels for each sample score = silhouette_score(X, labels, metric='euclidean')</pre>
silhouette_samples	Provides silhouette scores for each individual sample, indicating how well it fits in assigned cluster. <b>Hyperparameters:</b> metric: Distance metric to use. <b>Pros:</b> Offers granular insight into each sample's clustering quality. <b>Limitations:</b> Same as silhouette_score, sensitive to outliers and distance metric.	<pre>from sklearn.metrics import silhouette_samples # X: Data used in clustering # labels: Cluster labels for each sample samples = silhouette_samples(X, labels, metric='euclidean')</pre>
devies_bouldin_score	Measures the average similarity ratio of each cluster with the most similar cluster. Lower values indicate better clustering. <b>Pros:</b> Provides a simple, effective clustering evaluation. <b>Limitations:</b> May not work well with highly imbalanced clusters.	<pre>from sklearn.metrics import devies_bouldin_score # X: Data used in clustering # labels: Cluster labels for each sample db_score = devies_bouldin_score(X, labels)</pre>
Voronoi	Computes the Voronoi diagram, which partitions space based on the nearest neighbor. <b>Pros:</b> Useful for spatial analysis and clustering. <b>Limitations:</b> Limited to use cases that involve spatial partitioning of data.	<pre>from scipy.spatial import Voronoi # points: Coordinates for Voronoi diagram vor = Voronoi(points)</pre>
voronoi_plot_2d	Plots the Voronoi diagram in 2D for visualizing clustering results. <b>Hyperparameters:</b> show_vertices: Whether to display the vertices. <b>Pros:</b> Great for visualizing spatial clustering. <b>Limitations:</b> Limited to 2D spaces and large datasets may cause performance issues.	<pre>from scipy.spatial import voronoi_plot_2d # vor: Voronoi diagram object voronoi_plot_2d(vor, show_vertices=True)</pre>
matplotlib patches Patch	Creates custom shapes such as rectangles, circles, or ellipses for adding to plots. <b>Hyperparameters:</b> color: Fill color of the shape. <b>Pros:</b> Versatile for visual customization. <b>Limitations:</b> May not support all shapes or complex customizations.	<pre>import matplotlib.patches as patches # Create a rectangle with specified width, height, and position rectangle = patches.Rectangle((0, 0), 1, 1, color='blue')</pre>
explained_variance_score	Measures the proportion of variance explained by the model's predictions. A higher score indicates better performance. <b>Pros:</b> Helps in assessing the fit of regression models. <b>Limitations:</b> Not suitable for classification tasks.	<pre>from sklearn.metrics import explained_variance_score # y_true: True values # y_pred: Predicted values ev_score = explained_variance_score(y_true, y_pred)</pre>
Ridge regression	Performs ridge regression (L2 regularization) to avoid overfitting by penalizing large coefficients. <b>Hyperparameters:</b> alpha: Regularization strength. <b>Pros:</b> Helps reduce overfitting in regression models. <b>Limitations:</b> May not work well with sparse data.	<pre>from sklearn.linear_model import Ridge # alpha: Regularization strength (Larger values indicate stronger regularization) ridge = Ridge(alpha=0.01)</pre>
Lasso regression	Performs lasso regression (L1 regularization), which encourages sparsity by penalizing the absolute value of coefficients. <b>Hyperparameters:</b> alpha: Regularization strength. <b>Pros:</b> Encourages sparse solutions, useful for feature selection. <b>Limitations:</b> May struggle with multicollinearity.	<pre>from sklearn.linear_model import Lasso # alpha: Regularization strength (Larger values indicate stronger regularization) lasso = Lasso(alpha=0.1)</pre>
Pipeline	Chains multiple steps of preprocessing and modeling into a single object, ensuring efficient workflow. <b>Pros:</b> Simplifies code, ensures reproducibility. <b>Limitations:</b> May not work well with complex pipelines requiring dynamic configurations.	<pre>from sklearn.pipeline import Pipeline # steps: List of tuples with name and estimator/transformer pipeline = Pipeline(steps=[('scalar', StandardScaler()), ('model', Ridge(alpha=0.01))])</pre>
GridSearchCV	Performs exhaustive search over a specified parameter grid to find the best model configuration. <b>Hyperparameters:</b> param_grid: Dictionary of parameter grids. <b>Pros:</b> Ensures optimal model parameters. <b>Limitations:</b> Computationally expensive for large grids.	<pre>from sklearn.model_selection import GridSearchCV # estimator: Model to be tuned # param_grid: Dictionary with parameters to search over grid_search = GridSearchCV(estimator=model, param_grid={'alpha': (0.1, 1.0, 10.0)})</pre>

Visualization strategies for k-means evaluation

Process Name	Brief Description	Code Snippet
Multiple runs of k-means	Executes KMeans clustering multiple times with different random initializations to assess variability in cluster assignments. <b>Advantage:</b> Helps visualize consistency. <b>Limitation:</b> Computationally costly for large datasets.	<pre># Number of runs for KMeans with different random states n_runs = 4 inertia_values = [] plt.figure(figsize=(12, 12)) # Run KMeans multiple times with different random states for i in range(n_runs):     means = KMeans(n_clusters=4, random_state=None) # Use the default 'n_init'     inertia_values.append(means.inertia_) # Plot the clustering result plt.subplot(2, 2, 1 + i) plt.scatter(X[:, 0], X[:, 1], c=means.labels_, cmap='tab10', alpha=0.6, edgecolor='k') plt.title(f'KMeans Clustering Run {i + 1}') plt.xlabel('Feature 1') plt.ylabel('Feature 2') plt.legend() # Print inertia values plt.show() # Print inertia values for i, inertia in enumerate(inertia_values, start=1):     print(f'Run {i}: Inertia={inertia:.2f}')</pre>

		<pre>about:blank</pre>
Elbow method	Evaluates the optimal number of clusters by plotting inertia (within-cluster sum of squares) for different k values. <b>Advantage:</b> Easy to interpret. <b>Limitation:</b> Subjective elbow point.	<pre># Range of k values to test k_values = range(2, 15) # Store performance metrics inertia_values = [] for k in k_values:     means = KMeans(n_clusters=k, random_state=0)     y_means = means.fit_predict(X)     # Calculate and store metrics     inertia_values.append(means.inertia_)     plt.figure(figsize=(10, 8))     plt.subplot(2, 1, 1)     # Plot the inertia values (Elbow Method)     plt.plot(k_values, inertia_values, marker='o')     plt.title('Elbow Method: Inertia vs. k')     plt.xlabel('Number of Clusters (k)')     plt.ylabel('Inertia')</pre>
Silhouette method	Determines the optimal number of clusters by evaluating Silhouette Scores for different k values. <b>Advantage:</b> Considers both cohesion and separation. <b>Limitation:</b> High computation for large datasets.	<pre># Range of k values to test k_values = range(2, 15) # Store performance metrics silhouette_scores = [] for k in k_values:     means = KMeans(n_clusters=k, random_state=0)     y_means = means.fit_predict(X)     silhouette_scores.append(silhouette_score(X, y_means)) # Plot the Silhouette Scores plt.figure(figsize=(10, 8)) plt.subplot(2, 1, 1) # Plot the inertia values, silhouette_scores, marker='o') plt.plot(k_values, silhouette_scores, marker='o') plt.title('Silhouette Score vs. k') plt.xlabel('Number of Clusters (k)') plt.ylabel('Silhouette Score')</pre>
Devies-Bouldin Index	Evaluates clustering performance by calculating DBI for different k values. <b>Advantage:</b> Quantifies compactness and separation. <b>Limitation:</b> Sensitive to cluster shapes and density.	<pre># Range of k values to test k_values = range(2, 15) # Store performance metrics devies_bouldin_indices = [] for k in k_values:     means = KMeans(n_clusters=k, random_state=0)     y_means = means.fit_predict(X)     inertia_values.append(means.inertia_)     devies_bouldin_indices.append(devies_bouldin_score(X, y_means)) # Plot the Devies-Bouldin Index plt.figure(figsize=(10, 8)) plt.subplot(2, 1, 1) # Plot the inertia values, devies_bouldin_indices, marker='o') plt.plot(k_values, devies_bouldin_indices, marker='o') plt.title('Devies-Bouldin Index vs. k') plt.xlabel('Number of Clusters (k)') plt.ylabel('Devies-Bouldin Index')</pre>

Authors

[Anil Choudhary](#)  
[Abhishek Gargua](#)



Skills Network

