

Task Category	Task	Target Variable Type	Methods	Rationale	Evaluation Measures	R Library:Function
Supervised Learning	Regression – predict continuous values	Continuous (Quantity)	Linear Regression	Best fit line approximates linear relationship between independent variable and dependent (target) variable; line equation can be used to predict target for new data	Prediction Errors (MAE, MSE, RMSE)	stats:lm()
	Classification – predict discrete categories	Categorical (Class)	K-Nearest Neighbors (kNN) (can be used for regression)	The class label of a data point can be estimated by a majority vote of its closest k neighbors	Classification Accuracy, Precision, Recall, F1	class:knn
			Decision Tree (can be used for regression)	A tree of conditions can be induced from a dataset where each condition checks the value of a variable and traversing the tree (resolving conditions) arrives at a class label	Classification Accuracy, Precision, Recall, F1	rpart:rpart()
			Random Forest (can be used for regression)	Growing many decision trees, each based on a random sample and taking a majority vote of the predictions of all trees to classify new data	Classification Accuracy, Precision, Recall, F1	randomForest:randomForest()
			Support Vector Machines (SVM) (can be used for regression)	A maximum margin decision boundary between data points of 2 classes based on a number of “support vectors” (data points) can be used to classify	Classification Accuracy, Precision, Recall, F1	e1071:svm()
Unsupervised Learning	Clustering – find natural structure (clusters) in data	None	k-means	To find k means (centroids) that are centers of k clusters of points where points in a cluster are closer to each other than to points in other clusters	Total Within Cluster Sum of Squares, Average Silhouette Width, Calinski–Harabasz index (CHI)	stats:kmeans()
			Partitioning Around Medoids (PAM) / k-medoids	To find k points (medoids) that are centers of k clusters of points where points in a cluster are closer to each other than to points in other clusters	Total Within Cluster Distances, Average Silhouette Width, Calinski–Harabasz index (CHI)	cluster:pam()
			Hierarchical Clustering (Agglomerative)	Begin with each data point in a cluster and iteratively merge close points into clusters until all points are in one cluster	Dendrogram	stats:hclust()
	Dimensionality Reduction – reduce dimension of dataset while preserving patterns	Principal Component Analysis (PCA)	Find new axes/features called “Principal Components” that are linear combinations of the original features such that the new axes maximize the variance explained		stats:princomp()	