

| Data I/O & External API Reference | | | |
|-----------------------------------|---|---|---|
| Module | Item | What it does | Example |
| sklearn.datasets | load_iris, load_wine, fetch_california_housing | Load toy or real datasets from the library or internet. | <pre>X, y = load_iris(return_X_y=True)</pre> |
| | make_classification, make_regression | Programmatically create synthetic classification/regression data. | <pre>X, y = make_classification(n_samples=200)</pre> |
| | make_blobs, make_moons | Create simple toy clusters / decision boundaries. | <pre>X, y = make_blobs(n_samples=300)</pre> |
| | make_circles | Generates a circular decision boundary. | <pre>X, y = make_circles(n_samples=400)</pre> |
| Dimensionality Reduction | | | |
| Module | Item | What it does | Example |
| sklearn.preprocessing | StandardScaler, MinMaxScaler, RobustScaler | Scale/normalize numeric features. | <pre>X_scaled = StandardScaler().fit_transform(X)</pre> |
| | Normalizer | Scale each sample to unit norm. | <pre>X_norm = Normalizer().fit_transform(X)</pre> |
| | LabelEncoder, OneHotEncoder | Encode categorical labels ordinal or one-hot. | <pre>X_dummies = OneHotEncoder().fit_transform(X)</pre> |
| | PolynomialFeatures | Generate interaction / polynomial terms. | <pre>X_poly = PolynomialFeatures(degree=2).fit_transform(X)</pre> |
| | Binarizer | Threshold values to 0/1. | <pre>X_binarized = Binarizer(threshold=0.5).fit_transform(X)</pre> |
| | PowerTransformer (e.g. Johnson_Betancourt) | Stabilize variance / normalize distributions. | <pre>X_sqrt = PowerTransformer().fit_transform(X)</pre> |
| | FeatureUnion / ColumnTransformer | Combine several transforms on different columns. | <pre>processor = ColumnTransformer([('num', StandardScaler(), num_cols), ('cat', OneHotEncoder(), cat_cols)])</pre> |
| | PCA | Non-linear PCA via kernels | <pre>X_pca = KernelPCA(kernel='rbf').fit_transform(X)</pre> |
| Feature Selection | | | |
| Module | Item | What it does | Example |
| sklearn.feature_selection | SelectKBest, RFE, RecursiveFeatureEliminationCV | Select top-n or recursively eliminate features based on a score. | <pre>selector = SelectKBest(score_func=chi2, k=10).fit(X, y)</pre> |
| | VarianceThreshold | Remove features with low variance. | <pre>v1 = VarianceThreshold(threshold=0.8).fit(X)</pre> |
| | SelectFromModel | Select features based on importance from a model (e.g., tree). | <pre>v2 = SelectFromModel(RandomForestClassifier(erl)).fit(X, y)</pre> |
| | PCA_IncrementalPCA (n_decomposition) | Principal Component Analysis for dimensionality reduction. | <pre>pca = PCA(n_components=5).fit_transform(X)</pre> |
| | TruncatedSVD | SVD for sparse matrices (e.g., TF-IDF). | <pre>svd = TruncatedSVD(n_components=100).fit_transform(X)</pre> |
| | KernelPCA | Non-linear PCA via kernels. | <pre>k pca = KernelPCA(kernel='rbf').fit_transform(X)</pre> |
| | cross_val_score, cross_validate | Compute cross-validated scores for a model. | <pre>cross_val_score(X, y, cv=5)</pre> |
| | GridSearchCV, RandomizedSearchCV | Exhaustive / random search over hyperparameter grid. | <pre>grid_search = GridSearchCV(LogisticRegression()), param_grid, cv=5).fit(X, y)</pre> |
| sklearn.model_selection | ShuffleSplit, StratifiedKFold | Custom split strategies for CV. | <pre>cv = StratifiedKFold(n_splits=5)</pre> |
| | ParameterSampler (used by RandomizedSearchCV) | Sample parameters randomly from distributions. | <pre>parameters = ParameterSampler(..., n_iter=10)</pre> |
| Validation Curves | | | |
| Module | Item | What it does | Example |
| sklearn.pipeline | Pipeline | Chain transformers + estimator into a single object. | <pre>pipe = Pipeline([('scale', StandardScaler()), ('clf', LogisticRegression(max_iter=200))])</pre> |
| | FeatureUnion | Combine multiple feature extraction pipelines in parallel. | <pre>union = FeatureUnion([('pca', PCA(1)), ('poly', PolynomialFeatures(2))])</pre> |
| | TransformedTargetRegressor | Apply a transform to the target variable (e.g., log-transform). | <pre>reg = TransformedTargetRegressor(regressor=LinearRegression(), transformer=np.log1p)</pre> |
| | make_pipeline, make_union | Short-hand constructors for pipelines / unions. | <pre>pipe = make_pipeline(StandardScaler(), LogisticRegression())</pre> |
| | make_union | | |
| | make_pipeline | | |
| | make_regressor | | |
| | make_classifier | | |
| Classification | | | |
| Module | Item | What it does | Example |
| sklearn.linear_model | LogisticRegression, SGDClassifier | Linear models for classification (logreg, adaboost). | <pre>clf = LogisticRegression(max_iter=200).fit(X, y)</pre> |
| | LinearDiscriminantAnalysis | LDA for dimensionality reduction + classification. | |
| | QuadraticDiscriminantAnalysis | QDA (quadratic decision boundaries). | |
| | SVC, NuSVC, LinearSVC | Support Vector Machines (linearized or linear). | <pre>svc = SVC(kernel='rbf', C=1.0).fit(X, y)</pre> |
| | DecisionTreeClassifier | CART decision tree. | |
| | RandomForestClassifier, GradientBoostingClassifier, AdaBoostClassifier, BaggingClassifier | Ensemble of trees / boosting. | |
| | VotingClassifier, StackingClassifier | Meta-learning by combining models. | |
| | KNeighborsClassifier | Instance-based nearest neighbor. | |
| sklearn.naive_bayes | GaussianNB, MultinomialNB, BernoulliNB | Naive Bayes variants. | |
| | MLPClassifier | Feed-forward neural network. | |
| | Perceptron | Perceptron classifier with GPUs. | |
| | SGDClassifier | Stochastic gradient classifier. | |
| Clustering | | | |
| Module | Item | What it does | Example |
| sklearn.cluster | KMeans, MiniBatchKMeans | Partition data into k clusters. | <pre>kmeans = KMeans(n_clusters=3).fit(X)</pre> |
| | AgglomerativeClustering | Hierarchical clustering. | |
| | BIRCH | Density-based clustering. | |
| | DBSCAN | Other clustering algorithms. | |
| | KNeighborsClassifier | Estimate density and classify via kernels. | |
| | KNeighborsRegressor | Instance-based regression. | |
| | NeuralNetwork | Neural network regression. | |
| | PosteriorProcessRegressor | Posterior regression with GPs. | |
| Density Estimation | | | |
| Module | Item | What it does | Example |
| sklearn.neighbors | KNeighborsClassifier, KNeighborsRegressor | K-nearest neighbor queries (used in DBSCAN, etc.) | |
| | RadiusNeighborsClassifier, RadiusNeighborsRegressor | | |
| | EM | | |
| | KernelDensity | | |
| | ParzenWindows | | |
| | KernelDensityEstimator | | |
| | KernelDensity | | |
| | KernelDensity | | |
| Evaluation Metrics | | | |
| Module | Item | What it does | Example |
| sklearn.metrics | accuracy_score, precision_score, recall_score, f1_score, roc_auc_score, confusion_matrix, classification_report | Basic or n-ary classification metrics. | <pre>accuracy = accuracy_score(y_true, y_pred)</pre> |
| | mean_absolute_error, mean_squared_error, r2_score | Common regression metrics. | <pre>regression = mean_squared_error(y_true, y_pred)</pre> |
| | adjusted_rand_score, silhouette_score | Cluster evaluation metrics. | <pre>clustering = adjusted_rand_score(y_true, y_pred)</pre> |
| | classification_report | Model selection. | <pre>report = classification_report(y_true, y_pred)</pre> |
| | precision_recall_fscore_support | Compute precision, recall, f-score and support. | <pre>precision, recall, f1, support = precision_recall_fscore_support(y_true, y_pred)</pre> |
| | confusion_matrix | Produce arrays for plots. | <pre>confusion = confusion_matrix(y_true, y_pred)</pre> |
| | cross_val_predict, validation_curve, learning_curve | Generate predictions / plots across CV folds. | <pre>learning = learning_curve(...)</pre> |
| | roc_auc_score | | |
| Text Processing | | | |
| Module | Class / Function | What it does (short explanation) | Typical usage pattern |
| sklearn.feature_extraction.text | CountVectorizer | Converts a collection of text documents to a new term-frequency matrix (one row per document, one column per word or n-gram). | <pre>CountVectorizer(ngram_range=(1,2), stop_words='english')</pre> |
| | TfidfVectorizer | Same as CountVectorizer but applies TF-IDF weighting (term frequency-inverse document frequency). Used for most ML pipelines. | <pre>TfidfVectorizer(max_features=5000)</pre> |
| | HashingVectorizer | Same idea as CountVectorizer but uses a hash table instead of building a sparse matrix - great for very large corpora. | <pre>HashingVectorizer(n_features=2**20)</pre> |
| | stop_words | Stop words for English language. | |
| | Normalizer | Scales feature vectors to unit norm (used for cosine-similarity based models). | <pre>Normalizer(norm='l2')</pre> |
| | CharNgramVectorizer | Character n-gram vectorizer, commonly used for character-n-gram models. | |
| | PolynomialFeatures (mainly used in NLP) | Creates polynomial features (commonly used for character-n-gram models). | <pre>Pipeline([('vec', CharNgramVectorizer()), ('tfidf', TfidfVectorizer()), ('clf', LogisticRegression())])</pre> |
| | Pipeline | Creates a pipeline of several feature extraction/pipelines (e.g., word-level + character-level). | <pre>featureunion = FeatureUnion(...)</pre> |
| sklearn.decomposition | FeatureUnion / make_union | Performs latent semantic analysis on sparse TF-IDF matrices; reduces dimensionality while preserving semantics. | <pre>lsa = TruncatedSVD(n_components=100)</pre> |
| | TruncatedSVD (aka LSA) | Non-negative matrix factorization - useful for visualizing high-dimensional text embeddings. | <pre>NMF(n_components=20)</pre> |
| | TR | Non-linear dimensionality reduction for visualizing high-dimensional text embeddings. | <pre>T-SNE(n_components=2)</pre> |
| | classification_report | Standard evaluation metrics for text classifiers. | |
| | confusion_matrix | Produce arrays for plots. | |
| | cross_val_score | Compute cross-validation scores. | |
| | cross_val_predict | Generate predictions / plots across CV folds. | |
| | cross_val_score | | |
| sklearn.linear_model | LogisticRegression, SGDClassifier (with hinge or log loss) | Fast linear classifiers for text (when the baseline is linear). | <pre>LogisticRegression(max_iter=1000)</pre> |
| | PassiveAggressiveClassifier | Online, large-scale linear model that works well for sparse texts. | |
| | sklearn.naive_bayes | Classic Naive Bayes models for word-count or character-n-gram. | <pre>MultinomialNB(alpha=0.1)</pre> |
| | sklearn.svm | Linear SVC (kernel='linear') | <pre>LinearSVC(C=1)</pre> |
| | sklearn.tree | Decision trees - often outperform linear models on medium-size text corpora. | |
| | sklearn.ensemble | GradientBoostingClassifier, AdaBoostClassifier. | |
| | KNeighborsClassifier (rarely used) | Instance-based classification; works with multi-class classification. | |
| | MLPClassifier | Combine multiple vectorizers - e.g., word-level + character-n-gram. | |
| sklearn.pipeline | FeatureUnion / make_union | Combine multiple vectorizers - e.g., word-level + character-n-gram. | <pre>make_unionword_vec, char_vec</pre> |
| | make_union | Will run custom preprocessing functions inside a pipeline. | |
| | FunctionTransformer | Will run a Python function as a custom model. | |
| | text_dataset (user-defined) | Will return X, y for the above pipeline. | |
| | Utilities | | |