

scikit-learn

The scikit-learn workflow, from splitting data and building pipelines to tuning, scoring, and saving a model.

Data: Load & Split

Get arrays in, hold out a test set you never touch while modeling.

PURPOSE	COMMAND	RESULT
Load a toy dataset as a DataFrame.	<code>X, y =load_iris(return_X_y=True, as_frame=True)</code>	
Fetch a real dataset by name.	<code>X, y =fetch_openml("titanic", version=1, return_X_y=True, as_frame=True)</code>	
Make synthetic data for demos.	<code>X, y =make_classification(n_samples=1000, n_classes=3, n_informative=3)</code>	
Hold out a test set (stratified).	<code>train_test_split(X, y, test_size=0.2, stratify=y, random_state=0)</code>	

X is (n_samples, n_features); split off the test set and do not look at it until final evaluation

Preprocess: Scale, Encode, Impute






Transformers learn statistics with .fit, then reshape any data with .transform.

PURPOSE	COMMAND	RESULT
Standardize numeric columns.	<code>StandardScaler().fit_transform(X_tr)</code>	
One-hot encode categoricals.	<code>OneHotEncoder(handle_unknown="ignore", sparse_output=False)</code>	
Fill in missing values.	<code>SimpleImputer(strategy="median")</code>	
Get output column names back.	<code>enc.get_feature_names_out()</code>	
Return a DataFrame, not an array.	<code>StandardScaler().set_output("pandas")</code>	

fit learns statistics from training data only; that is what keeps test data from leaking in

Estimator: Fit, Predict, Score



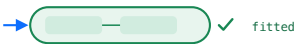

.fit learns, .predict applies, .score reports. Same three calls for every algorithm.

PURPOSE	COMMAND	RESULT
Pick & train a classifier.	<code>clf = RandomForestClassifier(n_estimators=300).fit(X_tr, y_tr)</code>	
Predict labels on new data.	<code>y_pred = clf.predict(X_te)</code>	
Predict class probabilities.	<code>proba = clf.predict_proba(X_te)</code>	
Default score on a test set.	<code>clf.score(X_te, y_te)</code>	
Swap in another algorithm.	<code>from sklearn.svm import SVC</code>	

the interface is uniform, so swapping algorithms is a one-line change

Pipeline: Chain Steps Safely





Wrap preprocessing and a model in one estimator that travels together.

PURPOSE	COMMAND	RESULT
Chain a scaler + model.	<code>pipe = make_pipeline(StandardScaler(), LogisticRegression(max_iter=1000))</code>	
Route numeric vs categorical cols.	<code>ColumnTransformer([("num", StandardScaler(), num_cols), ("cat", OneHotEncoder(), cat_cols)])</code>	
Auto-select columns by dtype.	<code>make_column_selector(dtype_include="number")</code>	
Fit the whole chain at once.	<code>pipe.fit(X_tr, y_tr)</code>	
Reach into a named step.	<code>pipe.named_steps["logisticregression"].coef_</code>	

transforms are fit only on training folds, which prevents the most common data leakage

Validate & Tune

Rotate the held-out fold for honest scores, then search and refit the best config.

PURPOSE	COMMAND	RESULT									
K-fold cross-validated score.	<code>scores = cross_val_score(pipe, X, y, cv=5)</code>										
Multiple metrics + train scores.	<code>cross_validate(pipe, X, y, cv=5, scoring=["accuracy", "f1_macro"])</code>	<table border="1" data-bbox="1128 367 1396 420"> <thead> <tr> <th>test_acc</th> <th>test_f1</th> <th>fit_time</th> </tr> </thead> <tbody> <tr> <td>0.94</td> <td>0.93</td> <td>0.08s</td> </tr> <tr> <td>0.96</td> <td>0.95</td> <td>0.07s</td> </tr> </tbody> </table>	test_acc	test_f1	fit_time	0.94	0.93	0.08s	0.96	0.95	0.07s
test_acc	test_f1	fit_time									
0.94	0.93	0.08s									
0.96	0.95	0.07s									
Grid search over a pipeline.	<code>gs = GridSearchCV(pipe, {"logisticregression__C": [0.1, 1, 10]}, cv=5)</code>										
Read the winner.	<code>gs.best_params_</code> and <code>gs.best_score_</code>										
Randomized search (big spaces).	<code>RandomizedSearchCV(pipe, ..., n_iter=20)</code>										

pass the whole pipeline so preprocessing is re-fit inside every fold; tune with `step_param` naming

Metrics & Evaluation

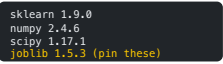
Pick the metric for the task; accuracy hides class imbalance.

PURPOSE	COMMAND	RESULT
Per-class precision/recall/F1.	<code>print(classification_report(...))</code>	<pre>prec rec f1 sup class 0 0.93 0.95 0.94 40 class 1 0.91 0.88 0.89 35 macro avg 0.92 0.91 0.92 75</pre>
Confusion matrix (plotted).	<code>ConfusionMatrixDisplay.from_estimator(...)</code>	
Ranking quality (AUC).	<code>roc_auc_score(y_te, proba[:, 1])</code>	
Regression error.	<code>root_mean_squared_error</code>	
Custom scorer for search.	<code>make_scorer(f1_score, average="macro")</code>	

`classification_report`, confusion matrix, and ROC AUC for classes; RMSE / R^2 for regression

Persist: Save & Load


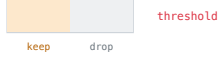
Train once, serve many times. Persist the whole fitted pipeline with joblib.

PURPOSE	COMMAND	RESULT
Save a fitted pipeline.	<code>joblib.dump(pipe, "model.joblib")</code>	
Load it back later.	<code>pipe = joblib.load("model.joblib")</code>	
Predict in a fresh process.	<code>joblib.load("model.joblib").predict(X_new)</code>	
Check installed versions.	<code>sklearn.show_versions()</code>	

the serialized format is not guaranteed across versions: pin and record them; only unpickle models you trust

Inspect & Select

Understand what the model uses and trim the inputs.

PURPOSE	COMMAND	RESULT
Tree feature importances.	<code>clf.feature_importances_</code>	
Model-agnostic importance.	<code>permutation_importance(clf, X_te, y_te)</code>	
Keep the best K features.	<code>SelectKBest(f_classif, k=10).fit_transform(X, y)</code>	
Select via a model's weights.	<code>SelectFromModel(estimator, threshold="median")</code>	
Compress with PCA.	<code>PCA(n_components=2).fit_transform(X)</code>	
Visualize the estimator.	<code>set_config(display="diagram")</code>	

feature_importances_ is fast; permutation_importance works for any fitted estimator