Practical Examples: Supervised & Unsupervised Learning, Dimensionality Reduction



Lida Kanari

Mathematical Institute University of Oxford

Introduction to Machine Learning, November 2025





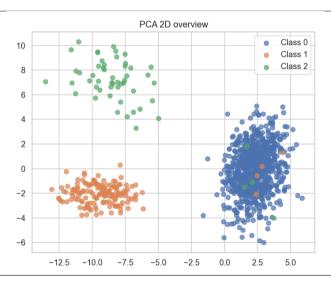


Table of Contents



- ► Example Datasets
- ▶ Sklearn
- Classification
- Clustering
- ► Dimensionality Reduction



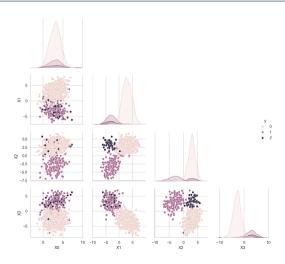




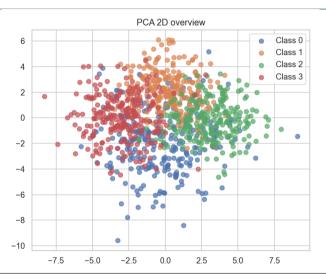


Dataset A





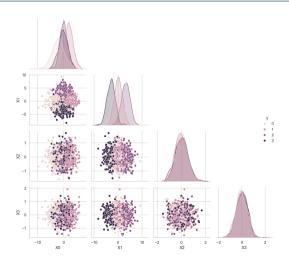














https://scikit-learn.org/stable/index.html



Packages you will need for the practical lecture



- numpy
- pandas
- ► seaborn
- sklearn
- matplotlib
- umap



Dataset A: simple three class dataset

```
31: # Generate dataset
    n \text{ samples} = 1000
    n features = 10
    weights = [0.8, 0.15, 0.05]
    X, y = make classification(n samples=n samples,
                                n features=n features.
                                n informative=8.
                                n_redundant=0,
                                n repeated=0.
                                n classes=3,
                                n_clusters_per_class=1,
                                weights=weights,
                                class sep=3.0.
                                flip y=0.01,
                                 random state=RND)
    feature_names = [f'X{i}' for i in range(X.shape[1])]
    df = pd.DataFrame(X. columns=feature names)
    df['v'] = v
```

Part 3: Dimensionality Reduction Algorithms





Summary Table

Algorithm	Туре	Linear / Non- linear	Preserves Local Structure	Preserves Global Structure	Key Parameters	Main Use Case	Link
PCA	Principal components	Linear	Partial	✓ Strong	n_components	Noise reduction, variance analysis	Docs
ICA	Independent components	Linear	Partial	Moderate	n_components	Source separation, signal analysis	Docs
t-SNE	Probabilistic	X Non-linear	Excellent	Weak	<pre>n_components , perplexity , learning_rate</pre>	Visualizing high- dimensional clusters	Docs
UMAP	Graph-based	X Non-linear	Excellent	☑ Good	<pre>n_neighbors , min_dist , n_components</pre>	Fast, interpretable 2D embeddings	Docs
Isomap	Geodesic	X Non-linear	✓ Good	✓ Good	<pre>n_neighbors , n_components , metric</pre>	Unfolding curved manifolds	Docs
MDS	Distance-based	X Non-linear	Partial	✓ Strong	<pre>n_components, metric, max_iter, dissimilarity</pre>	Preserving pairwise distances	Docs

Oxford Practical Examples October 2025 13/62
Mathematics



▼ 3.1 PCA (Principal Component Analysis)

scikit-learn docs: https://scikit-learn.org/stable/modules/decomposition.html#pca

```
[153]: pca2 = PCA(n_components=2, random_state=RND).fit(X)
X_pca2 = pca2.transform(X)
explained = pca2.explained_variance_ratio_
```

Oxford Practical Examples October 2025 14 / 62

Principal Component Analysis (PCA)



Type: Linear Projection

Concept: Finds orthogonal directions (principal components) that capture the maximum variance in the data.

Strengths:

Fast, interpretable, widely used.

Provides insight into feature variance and structure.

Limitations:

Captures only linear relationships.

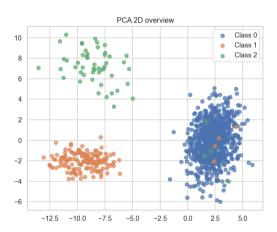
Sensitive to feature scaling.

Key Parameters: n_components, svd_solver, whiten.

Reference: scikit-learn - PCA

Principal Component Analysis (PCA)





Independent Component Analysis (ICA)



Type: Linear (Statistical Independence)

Concept: Decomposes data into statistically independent components.

Strengths:

Identifies independent hidden factors.

Limitations:

Sensitive to noise and scaling.

Results depend on initialization.

Key Parameters: n_components, max_iter, tol.

Reference: scikit-learn - ICA

t-Distributed Stochastic Neighbor Embedding (t-SNE)



Type: Non-linear / Probabilistic Manifold Learning

Concept: Preserves local structure by modeling pairwise similarities in high- and low-dimensional space.

Strengths:

► Good for visualizing non-linear clusters.

Limitations:

Computationally heavy.

► Non-deterministic results.

Key Parameters: n_components, perplexity.

Reference: scikit-learn - tSNE

Uniform Manifold Approximation and Projection (UMAP)



Type: Non-linear / Graph-based

Concept: Builds a weighted graph and learns a low-dimensional embedding that tried

to preserve local and global structures.

Strengths:

► Faster than t-SNE.

Limitations:

► Non-deterministic results.

Key Parameters: n_neighbors, min_dist, metric.

Reference: umap-learn

Isomap (Isometric Mapping)



Type: Non-linear / Geodesic

Concept: Preserves geodesic distances computed along a neighborhood graph.

Strengths:

► Captures curved manifolds like the Swiss roll.

Limitations:

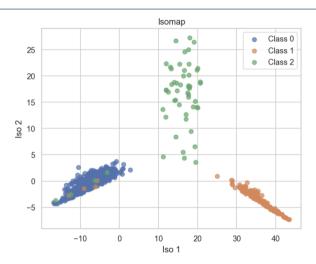
Sensitive to noise.

Key Parameters: n_neighbors, n_components, metric.

Reference: scikit-learn - Isomap

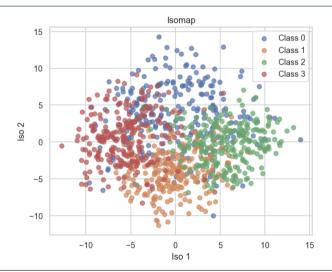
Isomap





Isomap





Multidimensional Scaling (MDS)



Type: Non-linear / Distance-preserving

Concept: Embeds points to preserve pairwise distances as faithfully as possible.

Strengths:

▶ Works with arbitrary distance metrics.

Limitations:

Computationally expensive; can get stuck in local minima.

Key Parameters: n_components, metric, max_iter, dissimilarity.

Reference: scikit-learn - MDS





- ► Linear methods (PCA, ICA) are fast and interpretable.
- Non-linear methods (t-SNE, UMAP, Isomap, MDS) uncover complex manifolds.
- ▶ It may be useful to scale features before projection.

Oxford Practical Examples October 2025 24 / 62





Summary Table

Algorithm	Туре	Linear / Non- linear	Preserves Local Structure	Preserves Global Structure	Key Parameters	Main Use Case	Link
PCA	Principal components	✓ Linear	Partial	✓ Strong	n_components	Noise reduction, variance analysis	Docs
ICA	Independent components	Linear	Partial	Moderate	n_components	Source separation, signal analysis	Docs
t-SNE	Probabilistic	X Non-linear	Excellent	Weak	<pre>n_components , perplexity , learning_rate</pre>	Visualizing high- dimensional clusters	Docs
UMAP	Graph-based	X Non-linear	Excellent	☑ Good	<pre>n_neighbors , min_dist , n_components</pre>	Fast, interpretable 2D embeddings	Docs
Isomap	Geodesic	X Non-linear	☑ Good	☑ Good	<pre>n_neighbors , n_components , metric</pre>	Unfolding curved manifolds	Docs
MDS	Distance-based	X Non-linear	Partial	☑ Strong	<pre>n_components , metric , max_iter , dissimilarity</pre>	Preserving pairwise distances	Docs

Oxford Practical Examples October 2025 25/62 Mathematics

Part 4: Classification Algorithms

Oxford Practical Examples October 2025 26/62 Mathematics

Part 4: Classification Algorithms



Summary

Algorithm	Type	Linear	Probabilistic	Handles Non-linearity	Key Strength	Link
Gaussian Naive Bayes	Probabilistic	✓ Yes	✓ Yes	X No	Very fast baseline	Docs
LDA	Linear Generative	✓ Yes	✓ Yes	X No	Simple & interpretable	Docs
QDA	Quadratic Generative	X No	✓ Yes	✓ Yes	Captures curved boundaries	Docs
Logistic Regression	Linear Discriminative	✓ Yes	✓ Yes	X No	Robust baseline classifier	Docs
SVM	Margin-based	(depends on kernel)	X No	Yes (with kernel)	Strong generalization	Docs
k-NN	Instance-based	X No	X No	✓ Yes	Simple and non-parametric	Docs
Decision Trees	Tree-based	X No	(optional)	✓ Yes	Interpretable structure	Docs
Random Forests	Ensemble	X No	(probabilities via votes)	✓ Yes	Robust ensemble learner	Docs

Oxford
MathematicsPractical ExamplesOctober 202527 / 62





Gaussian Naive Bayes (GaussianNB)

sklearn: https://scikit-learn.org/stable/modules/naive_bayes.html

```
from sklearn.naive bayes import GaussianNB
clf = GaussianNB()
try:
    clf.fit(X_train_s, y_train)
   v pred = clf.predict(X test s)
    res_gnb = classification_scores(y_test, y_pred)
    print('GaussianNB metrics:')
    for k in ['accuracy', 'balanced_accuracy', 'precision_macro', 'recall_macro', 'f1_macro']:
        v = res gnb.get(k, np.nan)
        trv:
            print(f' {k}: {v:.4f}')
        except Exception:
            print(f' {k}: {v}')
    plot preds pca(v pred. 'GaussianNB predictions', X test pca)
    if res_lda['confusion_matrix'] is not None:
        plot confusion(res gnb['confusion matrix'])
except Exception as e:
    print('GaussianNB failed:'. e)
```

Gaussian Naive Bayes



Type: Probabilistic (Generative)

Concept: Assumes features are conditionally independent given the class and normally distributed

Strengths:

- Very fast and simple; good baseline for many problems.
- ▶ Works well with small training sets and high-dimensional data.

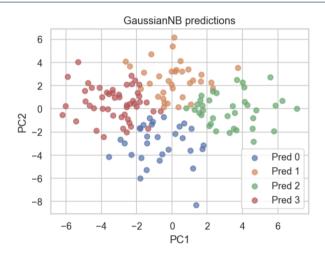
Limitations:

- ▶ Independence assumption is often violated in practice.
- Performance suffers when features are strongly correlated.

Reference: scikit-learn - Naive Bayes

Gaussian Naive Bayes (example)





Linear Discriminant Analysis (LDA)



Type: Linear (Generative)

Concept: Models each class with a Gaussian distribution and assumes identical covariance matrices; finds a linear boundary maximizing class separability.

Strengths:

► Interpretable; effective for linearly separable problems.

Provides probabilistic outputs.

Limitations:

► Assumes normally distributed features and equal covariances across classes.

Reference: scikit-learn - LDA/QDA

Linear Discriminant Analysis (LDA)



LDA metrics:

accuracy: 0.9933

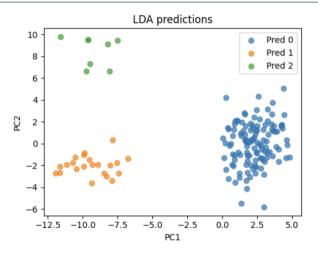
balanced accuracy: 0.9855 precision macro: 0.9972

recall macro: 0.9855

f1 macro: 0.9912



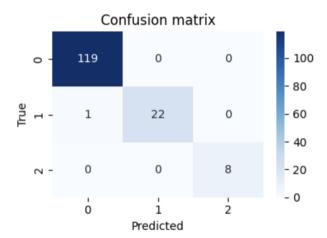




Practical Examples







Oxford Mathematics Practical Examples

Quadratic Discriminant Analysis (QDA)



Type: Quadratic (Generative)

Concept: Like LDA but allows each class to have its own covariance — yields quadratic decision boundaries.

Strengths:

► Can model non-linear boundaries between classes.

Limitations:

Estimates many parameters (covariances); prone to overfitting on small datasets.

Reference: scikit-learn – LDA/QDA

Logistic Regression



Type: Linear (Discriminative)

Concept: Models class probability using the logistic function applied to a linear combination of features.

Strengths:

Interpretable coefficients, outputs probabilities, robust baseline.

Limitations:

Only linear decision boundary unless features are transformed or interactions added.

Reference: scikit-learn — Logistic Regression

Support Vector Machine (SVM)



37/62

Type: Margin-based (Kernel)

Concept: Finds a hyperplane that maximizes margin between classes; kernels allow non-linear boundaries.

Strengths:

Effective in high-dimensional spaces; flexible via kernels.

Limitations:

Sensitive to kernel choice and regularization; can be slow on large datasets.

Key Parameters: kernel (linear, rbf, poly).

Reference: scikit-learn - SVM

k-Nearest Neighbors (k-NN)



Type: Instance-based / Non-parametric

Concept: Classifies a sample by majority vote among its k nearest neighbors.

Strengths:

Simple, non-parametric; adapts to complex boundaries.

Limitations:

▶ Prediction cost scales with dataset size; requires appropriate scaling and distance metric.

Reference: scikit-learn - Neighbors

Decision Trees



Type: Tree-based

Concept: Recursively splits the feature space to maximize classification accuracy.

Strengths:

Interpretable and easy to visualize.

Limitations:

▶ Prone to overfitting if not pruned.

Reference: scikit-learn - Decision Trees

Random Forests



Type: Ensemble of Trees

Concept: Aggregates many decision trees trained on bootstrapped samples and random feature subsets.

Strengths:

► Robust, reduces overfitting compared to single trees; often strong predictive performance.

Limitations:

Less interpretable than random trees; can be slower and memory-intensive.

Reference: scikit-learn - Random Forests

Classification — Summary



- Generative methods (Naive Bayes, LDA, QDA) model class-conditional distributions.
- ▶ **Discriminative linear methods** (Logistic Regression) model boundaries directly.
- Kernel methods (SVM) handle non-linearities more efficiently.
- ► Tree-based methods (Decision Tree, Random Forest) handle non-linearities more efficiently.
- ► **Evaluation** with confusion matrices, precision/recall/F1, and unbalanced accuracy.



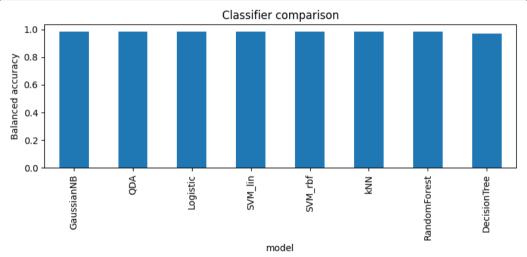


Summary

Algorithm	Type	Linear	Probabilistic	Handles Non-linearity	Key Strength	Link
Gaussian Naive Bayes	Probabilistic	✓ Yes	✓ Yes	X No	Very fast baseline	Docs
LDA	Linear Generative	✓ Yes	✓ Yes	X No	Simple & interpretable	Docs
QDA	Quadratic Generative	X No	✓ Yes	✓ Yes	Captures curved boundaries	Docs
Logistic Regression	Linear Discriminative	✓ Yes	✓ Yes	X No	Robust baseline classifier	Docs
SVM	Margin-based	(depends on kernel)	X No	Yes (with kernel)	Strong generalization	Docs
k-NN	Instance-based	X No	X No	✓ Yes	Simple and non-parametric	Docs
Decision Trees	Tree-based	X No	(optional)	✓ Yes	Interpretable structure	Docs
Random Forests	Ensemble	X No	(probabilities via votes)	✓ Yes	Robust ensemble learner	Docs

Classification — Summary





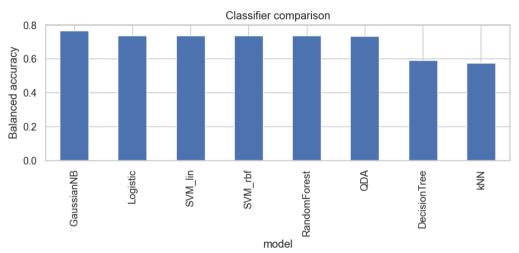
Oxford Mathematics

Practical Examples

October 2025

Classification — Summary





Oxford Mathematics

Practical Examples

Part 5: Clustering Algorithms



Part 5: Clustering Algorithms



Summary

Algorithm	Туре	Needs k	Handles Noise	Suitable For	Link
K-Means	Centroid-based	Yes	X No	Spherical clusters	Docs
Agglomerative	Hierarchical	Yes	× No	Hierarchical structure	Docs
DBSCAN	Density-based	× No	Yes	Arbitrary shapes	Docs
Spectral	Spectral	Yes	× No	Non-convex shapes	Docs
BIRCH	Hierarchical (tree-based)	Yes	X Limited	Large datasets / streaming	Docs
Affinity Propagation	Message-passing	× No	X Limited	Automatically finding exemplars	Docs
Gaussian Mixture Models	Probabilistic / Model-based	Yes	X Limited	Overlapping / elliptical clusters	Docs

Part 5: Clustering Algorithms



K-Means

sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html

```
from sklearn.cluster import KMeans

clf = KMeans(n_clusters=len(np.unique(y_train)), random_state=42)

try:
    clf.fit(X)
    y_pred = clf.predict(X)
    res_kmeans = unsupervised_clustering_metrics(X, y_pred, y_true=y)
```

K-Means



Type: Centroid-based

Concept: Iteratively assigns points to nearest centroid and updates centroids to minimize within-cluster variance.

Strengths:

► Fast and efficient for spherical, well-separated clusters.

Limitations:

 \blacktriangleright Requires specifying number of clusters k; sensitive to initialization and outliers.

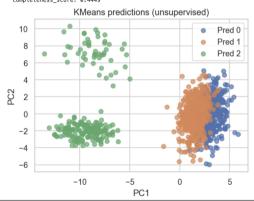
Key Parameters: n_clusters, n_init.

Reference: scikit-learn - KMeans

K-Means



KMeans unsupervised metrics: silhouette_score: 0.1988 calinski_harabasz_score: 535.5923 davies_bouldin_score: 1.6674 mutual_info_score: 0.4683 rand_score: 0.6651 homogeneity_score: 0.7428 fowlkes_mallows_score: 0.6996 completeness_score: 0.4449



Agglomerative Clustering



Type: Hierarchical (Bottom-up)

Concept: Start with each point as a cluster and iteratively merge the closest clusters per chosen linkage.

Strengths:

Produces a dendrogram; flexible linkage criteria.

Limitations:

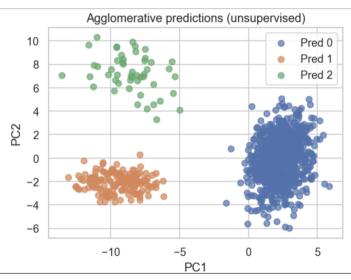
 $ightharpoonup n^3$ - worst-case complexity; not ideal for very large datasets.

Key Parameters: n_clusters, linkage (ward, complete, average), affinity.

Reference: scikit-learn – AgglomerativeClustering

Agglomerative Clustering





DBSCAN



Type: Density-based

Concept: Forms clusters as high-density regions separated by low-density areas; labels low-density points as outliers.

Strengths:

▶ Detects arbitrary-shaped clusters and outliers; does not require specifying k.

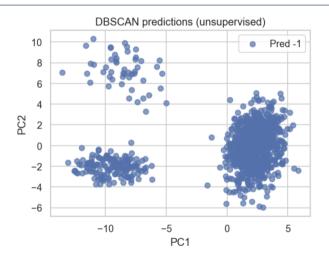
Limitations:

Sensitive to parameter selection, struggles with varying densities.

Key Parameters: eps, min_samples. **Reference:** scikit-learn – DBSCAN

DBSCAN





Spectral Clustering



Type: Graph / Spectral

Concept: Uses eigenvectors of a graph Laplacian built from data graph, then clusters

in the spectral embedding.

Strengths:

Effective for non-convex clusters; leverages graph structure.

Limitations:

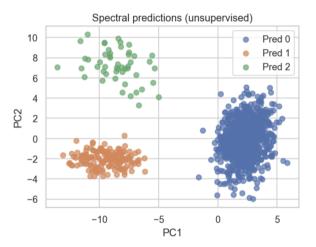
Computationally expensive; sensitive to affinity construction.

Key Parameters: n_clusters, affinity (rbf, nearest_neighbors).

Reference: scikit-learn – SpectralClustering

Spectral Clustering





BIRCH



Type: Hierarchical, Tree-based

Concept: Incrementally builds a CF-tree that summarizes data and clusters using

limited memory.

Strengths:

Scales well to large datasets; useful for streaming data.

Limitations:

Sensitive to threshold; not ideal for arbitrary-shaped clusters.

Key Parameters: threshold, branching_factor, n_clusters.

Reference: scikit-learn - Birch

Affinity Propagation



Type: Message-passing / Graph-based

Concept: Exchanges messages between points to identify exemplars; clusters form around exemplars.

Strengths:

Automatically selects number of clusters and exemplar points.

Limitations:

► High memory/time costs; sensitive to preference parameter.

Key Parameters: damping, preference, max_iter.

Reference: scikit-learn – AffinityPropagation

Gaussian Mixture Models (GMM)



Type: Probabilistic / Model-based

Concept: Models data as a set of Gaussian distributions, each representing a cluster.

Assigns soft probabilities rather than hard labels.

Strengths:

- ► Handles overlapping clusters; captures elliptical shapes.
- Provides probabilistic cluster membership.

Limitations:

ightharpoonup Requires specifying the number of components k; assumes Gaussian-shaped clusters; may converge to local minimum.

Key Parameters: n_components, covariance_type.

Reference: scikit-learn – GaussianMixture

Clustering — Summary



- ► Centroid-based: K-Means: fast, simple, assumes spherical clusters.
- Hierarchical: Agglomerative, BIRCH: reveals structure, scales to large datasets.
- ▶ **Density-based:** DBSCAN: identifies arbitrary shapes and noise.
- ► **Spectral/Graph-based:** Spectral Clustering, Affinity Propagation: captures complex relationships.
- ▶ **Probabilistic:** Gaussian Mixture Models: elliptical clusters.

Clustering — Summary

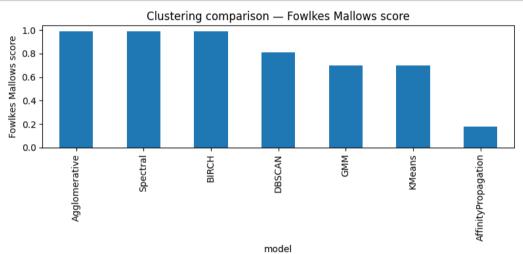


Summary

Algorithm	Туре	Needs k	Handles Noise	Suitable For	Link
K-Means	Centroid-based	Yes	X No	Spherical clusters	Docs
Agglomerative	Hierarchical	Yes	× No	Hierarchical structure	Docs
DBSCAN	Density-based	× No	Yes	Arbitrary shapes	Docs
Spectral	Spectral	Yes	× No	Non-convex shapes	Docs
BIRCH	Hierarchical (tree-based)	Yes	X Limited	Large datasets / streaming	Docs
Affinity Propagation	Message-passing	× No	X Limited	Automatically finding exemplars	Docs
Gaussian Mixture Models	Probabilistic / Model-based	Yes	X Limited	Overlapping / elliptical clusters	Docs

Clustering — Summary





Thank you! Questions?