

## sklearn.decomposition.PCA

```
class sklearn.decomposition.PCA(n_components=None, *, copy=True, whiten=False, svd_solver='auto', tol=0.0,
                                iterated_power='auto', random_state=None)
```

[\[source\]](#)

Principal component analysis (PCA).

Linear dimensionality reduction using Singular Value Decomposition of the data to project it to a lower dimensional space. The input data is centered but not scaled for each feature before applying the SVD.

It uses the LAPACK implementation of the full SVD or a randomized truncated SVD by the method of Halko et al. 2009, depending on the shape of the input data and the number of components to extract.

It can also use the scipy.sparse.linalg ARPACK implementation of the truncated SVD.

Notice that this class does not support sparse input. See [TruncatedSVD](#) for an alternative with sparse data.

Read more in the [User Guide](#).

<b>Parameters:</b>	<b>n_components</b> : <i>int, float, None or str</i> Number of components to keep. If n_components is not set all components are kept: <pre>n_components == min(n_samples, n_features)</pre> If n_components == 'mle' and svd_solver == 'full', Minka's MLE is used to guess the dimension. Use of n_components == 'mle' will interpret svd_solver == 'auto' as svd_solver == 'full'. If 0 < n_components < 1 and svd_solver == 'full', select the number of components such that the amount of variance that needs to be explained is greater than the percentage specified by n_components. If svd_solver == 'arpack', the number of components must be strictly less than the minimum of n_features and n_samples. Hence, the None case results in: <pre>n_components == min(n_samples, n_features) - 1</pre> <b>copy</b> : <i>bool, default=True</i> If False, data passed to fit are overwritten and running fit(X).transform(X) will not yield the expected results, use fit_transform(X) instead. <b>whiten</b> : <i>bool, optional (default False)</i> When True (False by default) the components_ vectors are multiplied by the square root of n_samples and then divided by the singular values to ensure uncorrelated outputs with unit component-wise variances. Whitening will remove some information from the transformed signal (the relative variance scales of the components) but can sometime improve the predictive accuracy of the downstream estimators by making their data respect some hard-wired assumptions. <b>svd_solver</b> : <i>str {'auto', 'full', 'arpack', 'randomized'}</i> <b>if auto</b> : The solver is selected by a default policy based on X.shape and n_components: if the input data is larger than 500x500 and the number of components to extract is lower than 80% of the smallest dimension of the data, then the more efficient 'randomized' method is enabled. Otherwise the exact full SVD is computed and optionally truncated afterwards. <b>if full</b> : run exact full SVD calling the standard LAPACK solver via scipy.linalg.svd and select the components by postprocessing <b>if arpack</b> : run SVD truncated to n_components calling ARPACK solver via scipy.sparse.linalg.svds. It requires strictly 0 < n_components < min(X.shape) <b>if randomized</b> : run randomized SVD by the method of Halko et al. <i>New in version 0.18.0.</i> <b>tol</b> : <i>float &gt;= 0, optional (default .0)</i> Tolerance for singular values computed by svd_solver == 'arpack'. <i>New in version 0.18.0.</i> <b>iterated_power</b> : <i>int &gt;= 0, or 'auto', (default 'auto')</i> Number of iterations for the power method computed by svd_solver == 'randomized'. <i>New in version 0.18.0.</i> <b>random_state</b> : <i>int, RandomState instance, default=None</i> Used when svd_solver == 'arpack' or 'randomized'. Pass an int for reproducible results across multiple function calls. See <a href="#">Glossary</a> . <i>New in version 0.18.0.</i>
<b>Attributes:</b>	<b>components_</b> : <i>array, shape (n_components, n_features)</i> Principal axes in feature space, representing the directions of maximum variance in the data. The components are sorted by explained_variance_. <b>explained_variance_</b> : <i>array, shape (n_components,)</i> The amount of variance explained by each of the selected components. Equal to n_components largest eigenvalues of the covariance matrix of X. <i>New in version 0.18.</i> <b>explained_variance_ratio_</b> : <i>array, shape (n_components,)</i> Percentage of variance explained by each of the selected components. If n_components is not set then all components are stored and the sum of the ratios is equal to 1.0. <b>singular_values_</b> : <i>array, shape (n_components,)</i> The singular values corresponding to each of the selected components. The singular values are equal to the 2-norms of the n_components variables in the lower-dimensional space. <i>New in version 0.19.</i> <b>mean_</b> : <i>array, shape (n_features,)</i> Per-feature empirical mean, estimated from the training set. Equal to X.mean(axis=0). <b>n_components_</b> : <i>int</i> The estimated number of components. When n_components is set to 'mle' or a number between 0 and 1 (with svd_solver == 'full') this number is estimated from input data. Otherwise it equals the parameter n_components, or the lesser value of n_features and n_samples if n_components is None. <b>n_features_</b> : <i>int</i> Number of features in the training data. <b>n_samples_</b> : <i>int</i> Number of samples in the training data. <b>noise_variance_</b> : <i>float</i> The estimated noise covariance following the Probabilistic PCA model from Tipping and Bishop 1999. See "Pattern Recognition and Machine Learning" by C. Bishop, 12.2.1 p. 574 or <a href="http://www.miketipping.com/papers/met-mppca.pdf">http://www.miketipping.com/papers/met-mppca.pdf</a> . It is required to compute the estimated data covariance and score samples. Equal to the average of (min(n_features, n_samples) - n_components) smallest eigenvalues of the covariance matrix of X.

**See also:**

- [KernelPCA](#)  
Kernel Principal Component Analysis.
- [SparsePCA](#)  
Sparse Principal Component Analysis.
- [TruncatedSVD](#)  
Dimensionality reduction using truncated SVD.
- [IncrementalPCA](#)  
Incremental Principal Component Analysis.

**References**

For n\_components == 'mle', this class uses the method of *Minka, T. P. "Automatic choice of dimensionality for PCA". In NIPS, pp. 598-604*

Implements the probabilistic PCA model from: Tipping, M. E., and Bishop, C. M. (1999). "Probabilistic principal component analysis". Journal of the Royal Statistical Society: Series B (Statistical Methodology), 61(3), 611-622. via the score and score\_samples methods. See <http://www.miketipping.com/papers/met-mppca.pdf>

For svd\_solver == 'arpack', refer to [scipy.sparse.linalg.svds](#).

For svd\_solver == 'randomized', see: Halko, N., Martinsson, P. G., and Tropp, J. A. (2011). "Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions". *SIAM review*, 53(2), 217-288. and Martinsson, P. G., Rokhlin, V., and Tygert, M. (2011). "A randomized algorithm for the decomposition of matrices". *Applied and Computational Harmonic Analysis*, 30(1), 47-68.

**Examples**

```
>>> import numpy as np
>>> from sklearn.decomposition import PCA
>>> X = np.array([[1,-1, -1], [-2, -1], [-3, -2], [1, 1], [2, 1], [3, 2]])
>>> pca = PCA(n_components=2)
>>> pca.fit(X)
PCA(n_components=2)
>>> print(pca.explained_variance_ratio_)
[0.9924... 0.0075...]
>>> print(pca.singular_values_)
[6.30861... 0.54988...]
```

```
>>> pca = PCA(n_components=2, svd_solver='full')
>>> pca.fit(X)
PCA(n_components=2, svd_solver='full')
>>> print(pca.explained_variance_ratio_)
[0.9924... 0.0075...]
>>> print(pca.singular_values_)
[6.30861... 0.54988...]
```

```
>>> pca = PCA(n_components=1, svd_solver='arpack')
>>> pca.fit(X)
PCA(n_components=1, svd_solver='arpack')
>>> print(pca.explained_variance_ratio_)
[0.99244...]
```

**Methods**

<b>fit</b> (self, X[, y])	Fit the model with X.
<b>fit_transform</b> (self, X[, y])	Fit the model with X and apply the dimensionality reduction on X.
<b>get_covariance</b> (self)	Compute data covariance with the generative model.
<b>get_params</b> (self, deep)	Get parameters for this estimator.
<b>get_precision</b> (self)	Compute data precision matrix with the generative model.
<b>inverse_transform</b> (self, X)	Transform data back to its original space.
<b>score</b> (self, X[, y])	Return the average log-likelihood of all samples.
<b>score_samples</b> (self, X)	Return the log-likelihood of each sample.
<b>set_params</b> (self, **params)	Set the parameters of this estimator.
<b>transform</b> (self, X)	Apply dimensionality reduction to X.

**\_\_init\_\_**(self, n\_components=None, \*, copy=True, whiten=False, svd\_solver='auto', tol=0.0, iterated\_power='auto', random\_state=None) [\[source\]](#)

Initialize self. See help(type(self)) for accurate signature.

**fit**(self, X, y=None) [\[source\]](#)

Fit the model with X.

<b>Parameters:</b>	<b>X</b> : <i>array-like, shape (n_samples, n_features)</i> Training data, where n_samples is the number of samples and n_features is the number of features. <b>y</b> : <i>None</i> Ignored variable. <b>Returns:</b> <b>self</b> : <i>object</i> Returns the instance itself.
--------------------	--

**fit\_transform**(self, X, y=None) [\[source\]](#)

Fit the model with X and apply the dimensionality reduction on X.

<b>Parameters:</b>	<b>X</b> : <i>array-like, shape (n_samples, n_features)</i> Training data, where n_samples is the number of samples and n_features is the number of features. <b>y</b> : <i>None</i> Ignored variable. <b>Returns:</b> <b>X_new</b> : <i>array-like, shape (n_samples, n_components)</i> Transformed values.
--------------------	---

**Notes**

This method returns a Fortran-ordered array. To convert it to a C-ordered array, use 'np.ascontiguousarray'.

**get\_covariance**(self) [\[source\]](#)

Compute data covariance with the generative model.

cov = components\_.T \* S\*\*2 \* components\_ + sigma2 \* eye(n\_features) where S\*\*2 contains the explained variances, and sigma2 contains the noise variances.

**Returns:** **cov** : *array, shape=(n\_features, n\_features)*  
 Estimated covariance of data.

**get\_params**(self, deep=True) [\[source\]](#)

Get parameters for this estimator.

<b>Parameters:</b>	<b>deep</b> : <i>bool, default=True</i> If True, will return the parameters for this estimator and contained subobjects that are estimators. <b>Returns:</b> <b>params</b> : <i>mapping of string to any</i> Parameter names mapped to their values.
--------------------	---

**get\_precision**(self) [\[source\]](#)

Compute data precision matrix with the generative model.

Equals the inverse of the covariance but computed with the matrix inversion lemma for efficiency.

**Returns:** **precision** : *array, shape=(n\_features, n\_features)*  
 Estimated precision of data.

**inverse\_transform**(self, X) [\[source\]](#)

Transform data back to its original space.

In other words, return an input X\_original whose transform would be X.

<b>Parameters:</b>	<b>X</b> : <i>array-like, shape (n_samples, n_components)</i> New data, where n_samples is the number of samples and n_components is the number of components. <b>Returns:</b> <b>X_original</b> : <i>array-like, shape (n_samples, n_features)</i>
--------------------	---

**Notes**

If whitening is enabled, inverse\_transform will compute the exact inverse operation, which includes reversing whitening.

**score**(self, X, y=None) [\[source\]](#)

Return the average log-likelihood of all samples.

See, "Pattern Recognition and Machine Learning" by C. Bishop, 12.2.1 p. 574 or <http://www.miketipping.com/papers/met-mppca.pdf>

<b>Parameters:</b>	<b>X</b> : <i>array, shape(n_samples, n_features)</i> The data. <b>y</b> : <i>None</i> Ignored variable. <b>Returns:</b> <b>ll</b> : <i>float</i> Average log-likelihood of the samples under the current model.
--------------------	---

**score\_samples**(self, X) [\[source\]](#)

Return the log-likelihood of each sample.

See, "Pattern Recognition and Machine Learning" by C. Bishop, 12.2.1 p. 574 or <http://www.miketipping.com/papers/met-mppca.pdf>

<b>Parameters:</b>	<b>X</b> : <i>array, shape(n_samples, n_features)</i> The data. <b>Returns:</b> <b>ll</b> : <i>array, shape (n_samples,)</i> Log-likelihood of each sample under the current model.
--------------------	--

**set\_params**(self, \*\*params) [\[source\]](#)

Set the parameters of this estimator.

The method works on simple estimators as well as on nested objects (such as pipelines). The latter have parameters of the form <component>\_\_<parameter> so that it's possible to update each component of a nested object.

<b>Parameters:</b>	<b>**params</b> : <i>dict</i> Estimator parameters. <b>Returns:</b> <b>self</b> : <i>object</i> Estimator instance.
--------------------	--

**transform**(self, X) [\[source\]](#)

Apply dimensionality reduction to X.

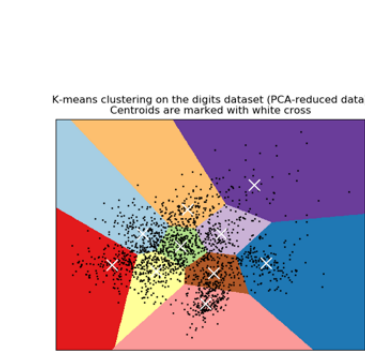
X is projected on the first principal components previously extracted from a training set.

<b>Parameters:</b>	<b>X</b> : <i>array-like, shape (n_samples, n_features)</i> New data, where n_samples is the number of samples and n_features is the number of features. <b>Returns:</b> <b>X_new</b> : <i>array-like, shape (n_samples, n_components)</i>
--------------------	--

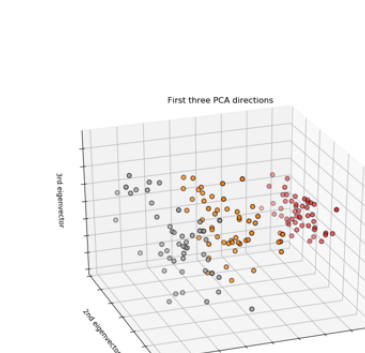
**Examples**

```
>>> import numpy as np
>>> from sklearn.decomposition import IncrementalPCA
>>> X = np.array([[1,-1, -1], [-2, -1], [-3, -2], [1, 1], [2, 1], [3, 2]])
>>> ipca = IncrementalPCA(n_components=2, batch_size=3)
>>> ipca.fit(X)
IncrementalPCA(batch_size=3, n_components=2)
>>> ipca.transform(X) # doctest: +SKIP
```

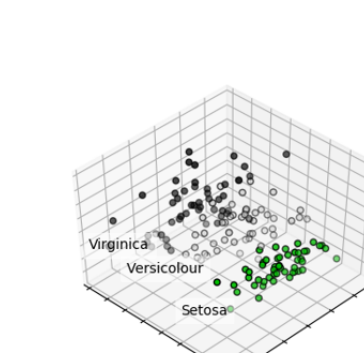
### Examples using sklearn.decomposition.PCA



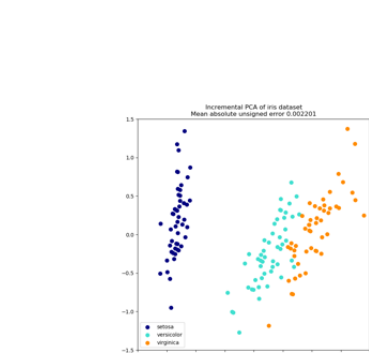
A demo of K-Means clustering on the handwritten digits data




The Iris Dataset



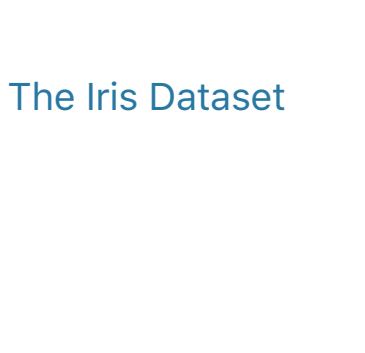
PCA example with Iris data-set



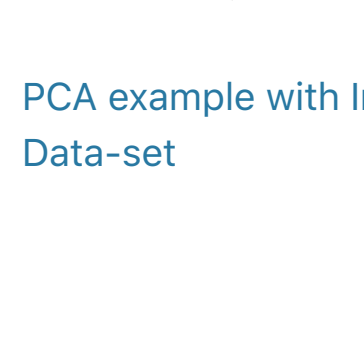
Incremental PCA



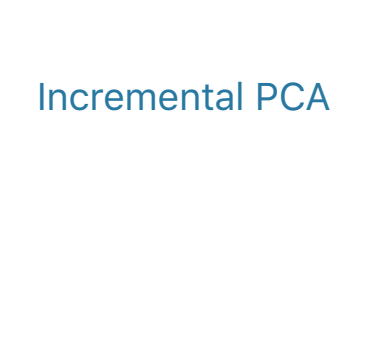
Comparison of LDA and PCA 2D projection of Iris dataset



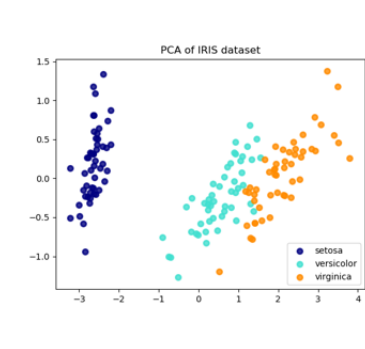
Blind source separation using FastICA



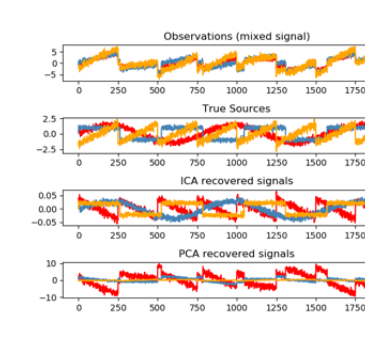
Principal components analysis (PCA)



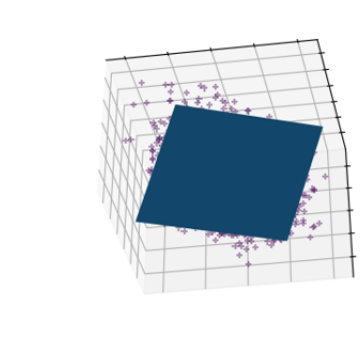
FastICA on 2D point clouds



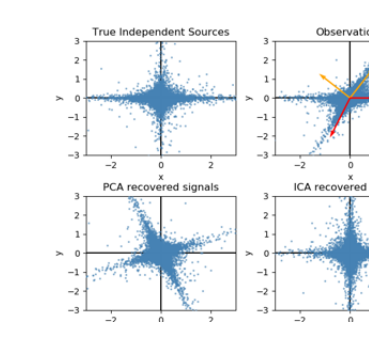
Kernel PCA



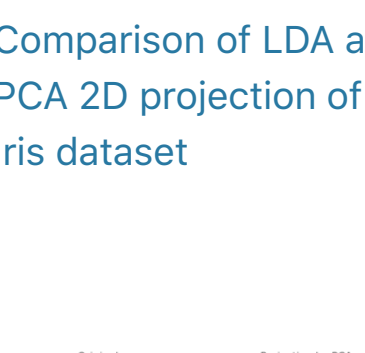
Model selection with Probabilistic PCA and Factor Analysis (FA)




Faces dataset decompositions



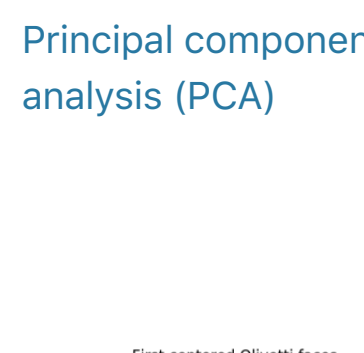
Faces recognition example using eigenfaces and SVMs



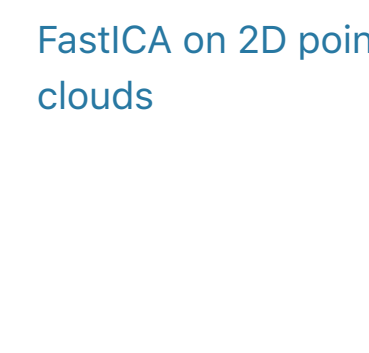
Multi-dimensional scaling



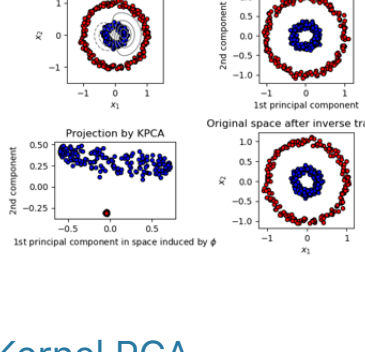
Multilabel classification



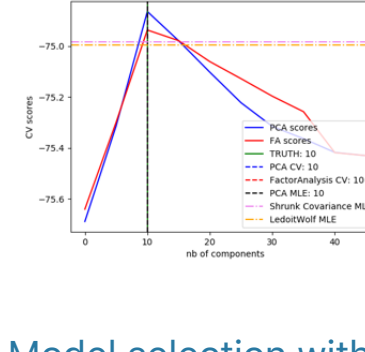
Explicit feature map approximation for RBF kernels



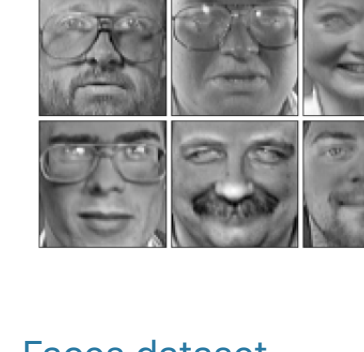
Balance model complexity and cross-validated score



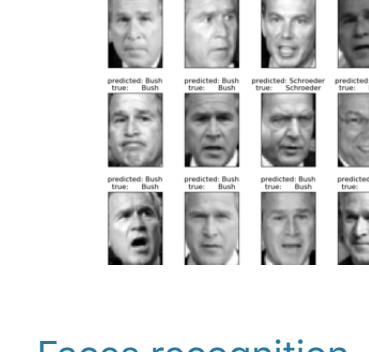
Kernel Density Estimation



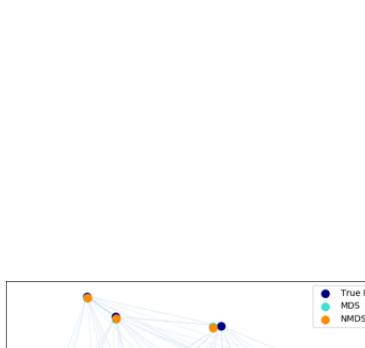
Dimensionality Reduction with Neighborhood Components Analysis



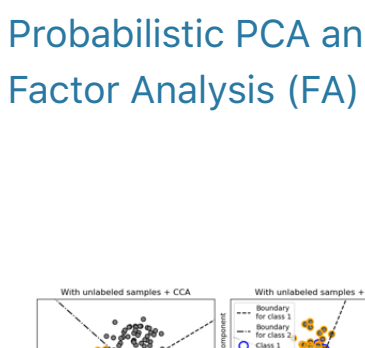
Concatenating multiple feature extraction methods




Pipelining: chaining a PCA and a logistic regression



Selecting dimensionality reduction with Pipeline and GridSearchCV



Using FunctionTransformer to select columns



Importance of Feature Scaling