



# PCA in Machine Learning



AMAN KHARWAL / JULY 8, 2020 / MACHINE LEARNING

In this article, you will explore what is perhaps one of the most broadly used of unsupervised algorithms, principal component analysis (PCA). PCA is fundamentally a dimensionality reduction algorithm, but it can also be useful as a tool for visualization, for noise filtering, for feature extraction and engineering, and much more.

## D3 short cuts for exploration

A higher level library by the creators of D3 meets pre-built UI

Observable Plot

## Also, read – Best Data Science Books

After a brief conceptual discussion of the PCA algorithm, we will see a couple examples of these further applications. I will start with the standard imports:

```
%matplotlib inline  
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns; sns.set()
```

# Introducing Principal Component Analysis (PCA)

Principal component analysis (PCA) is a fast and flexible unsupervised method for dimensionality reduction in data. Its

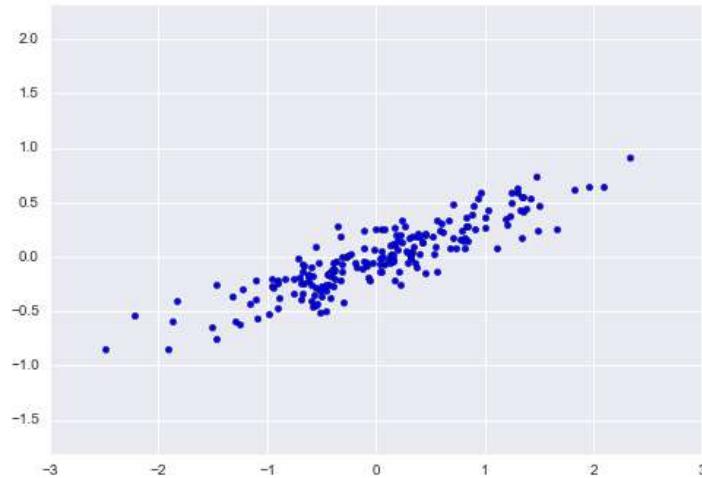
behavior is easiest to visualize by looking at a two-dimensional dataset. Consider the following 200 points:

### D3 short cuts for exploration

A higher level library by the creators of D3 meets pre-built UI

Observable Plot

```
rng = np.random.RandomState(1)
X = np.dot(rng.rand(2, 2), rng.randn(2, 200)).T
plt.scatter(X[:, 0], X[:, 1])
plt.axis('equal')
```



By eye, it is clear that there is a nearly linear relationship between the x and y variables. But the problem here states: rather than attempting to predict the y values from the x values, the unsupervised learning problem attempts to learn about the relationship between the x and y values.

In principal component analysis (PCA), this relationship is quantified by finding a list of the principal axes in the data, and using those axes to describe the dataset. Using Scikit-Learn's `PCA` estimator, we can compute this as follows:

```
from sklearn.decomposition import PCA
pca = PCA(n_components=2)
pca.fit(X)
```

```
PCA(copy=True, n_components=2, whiten=False)
```

The fit learns some quantities from the data, most importantly the “components” and “explained variance”:

## Pay only after placement

AlmaBetter

```
print(pca.components_)
```

```
[[ 0.94446029  0.32862557]
 [ 0.32862557 -0.94446029]]
```

```
print(pca.explained_variance_)
```

```
[ 0.75871884  0.01838551]
```

To see what these numbers mean, let’s visualize them as vectors over the input data, using the “components” to define the direction of the vector, and the “explained variance” to define the squared-length of the vector:

```
def draw_vector(v0, v1, ax=None):
    ax = ax or plt.gca()
    arrowprops=dict(arrowstyle='->',
                    linewidth=2,
                    shrinkA=0, shrinkB=0)
    ax.annotate('', v1, v0, arrowprops=arrowprops)

# plot data
plt.scatter(X[:, 0], X[:, 1], alpha=0.2)
for length, vector in zip(pca.explained_variance_, pca.components_):
    v = vector * 3 * np.sqrt(length)
    draw_vector(pca.mean_, pca.mean_ + v)
plt.axis('equal')
```



These vectors represent the principal axes of the data, and the length of the vector is an indication of how “important” that axis is in describing the distribution of the data—more precisely, it is a measure of the variance of the data when projected onto that axis. The projection of each data point onto the principal axes are the “principal components” of the data.

## PCA as Dimensionality Reduction

Using PCA for dimensionality reduction involves zeroing out one or more of the smallest principal components, resulting in a lower-dimensional projection of the data that preserves the maximal data variance.

**Pay only after placement**

AlmaBetter

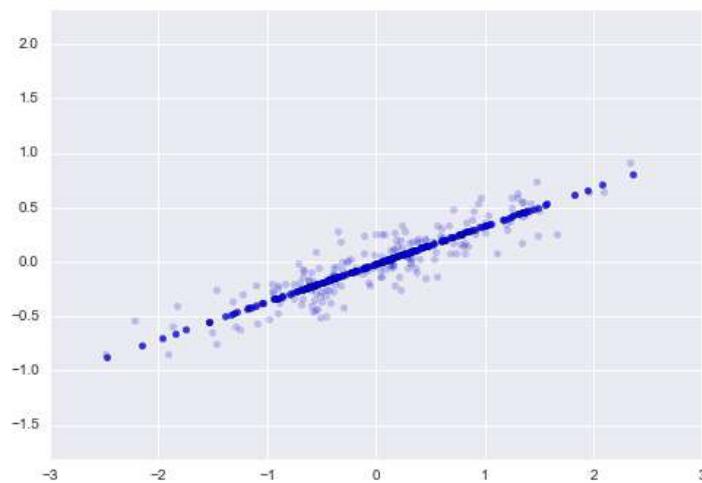
Here is an example of using PCA as a dimensionality reduction transform:

```
pca = PCA(n_components=1)
pca.fit(X)
X_pca = pca.transform(X)
print("original shape: ", X.shape)
print("transformed shape:", X_pca.shape)
```

```
original shape: (200, 2)
transformed shape: (200, 1)
```

The transformed data has been reduced to a single dimension. To understand the effect of this dimensionality reduction, we can perform the inverse transform of this reduced data and plot it along with the original data:

```
X_new = pca.inverse_transform(X_pca)
plt.scatter(X[:, 0], X[:, 1], alpha=0.2)
plt.scatter(X_new[:, 0], X_new[:, 1], alpha=0.8)
plt.axis('equal')
```



The light points are the original data, while the dark points are the projected version. This makes clear what a PCA dimensionality reduction means: the information along the least important principal axis or axes is removed, leaving only the component(s) of the data with the highest variance.

The fraction of variance that is cut out (proportional to the spread of points about the line formed in this figure) is roughly a measure of how much “information” is discarded in this reduction of dimensionality.

## PCA for visualization: Hand-written digits

The usefulness of the dimensionality reduction may not be entirely apparent in only two dimensions, but becomes much more clear when looking at high-dimensional data. To see this,

let's take a quick look at the application of PCA to the digits data.

I will start by loading the data:

### D3 short cuts for exploration

A higher level library by the creators of D3 meets pre-built UI

Observable Plot

```
from sklearn.datasets import load_digits
digits = load_digits()
digits.data.shape
```

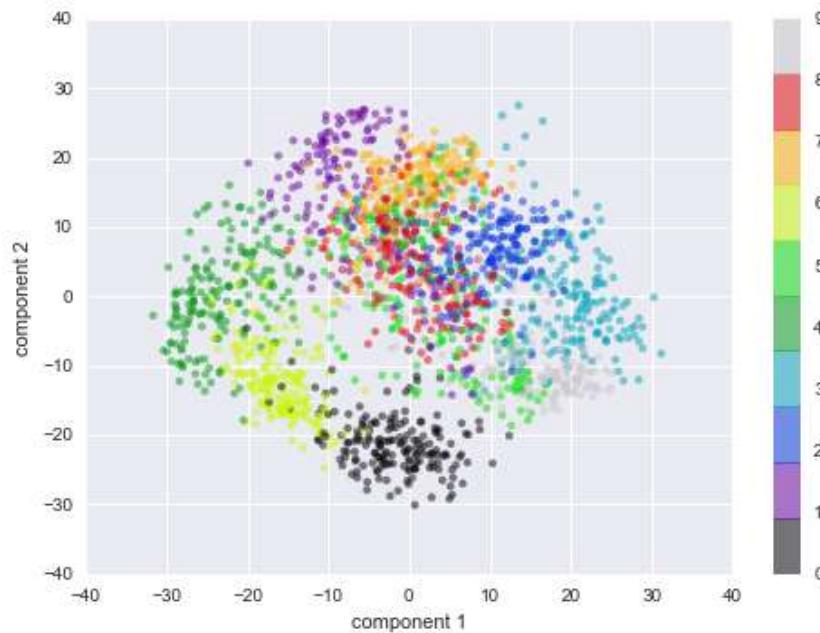
To gain some intuition into the relationships between these points, we can use PCA to project them to a more manageable number of dimensions, say two:

```
pca = PCA(2) # project from 64 to 2 dimensions
projected = pca.fit_transform(digits.data)
print(digits.data.shape)
print(projected.shape)
```

```
(1797, 64)
(1797, 2)
```

We can now plot the first two principal components of each point to learn about the data:

```
plt.scatter(projected[:, 0], projected[:, 1],
            c=digits.target, edgecolor='none', alpha=0.5,
            cmap=plt.cm.get_cmap('spectral', 10))
plt.xlabel('component 1')
plt.ylabel('component 2')
plt.colorbar()
```

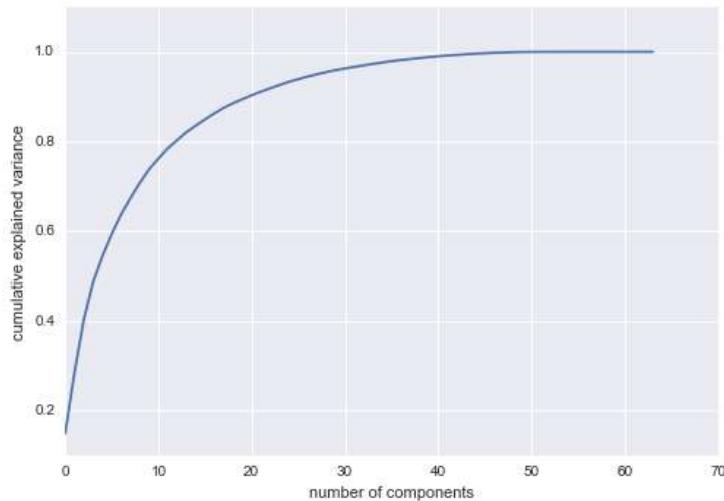


the full data is a 64-dimensional point cloud, and these points are the projection of each data point along the directions with the largest variance. Essentially, we have found the optimal stretch and rotation in 64-dimensional space that allows us to see the layout of the digits in two dimensions, and have done this in an unsupervised manner—that is, without reference to the labels.

## Choosing the Number of Components

A vital part of using PCA in practice is the ability to estimate how many components are needed to describe the data. This can be determined by looking at the cumulative explained variance ratio as a function of the number of components:

```
pca = PCA().fit(digits.data)
plt.plot(np.cumsum(pca.explained_variance_ratio_))
plt.xlabel('number of components')
plt.ylabel('cumulative explained variance')
```



## Principal Component Analysis Example: Eigenfaces

Here we will take a look and explore the Labeled Faces in the Wild dataset made available through Scikit-Learn:

**Pay only after  
placement**

AlmaBetter

```
from sklearn.datasets import fetch_lfw_people
faces = fetch_lfw_people(min_faces_per_person=60)
print(faces.target_names)
print(faces.images.shape)

['Ariel Sharon' 'Colin Powell' 'Donald Rumsfeld' 'George W Bush'
 'Gerhard Schroeder' 'Hugo Chavez' 'Junichiro Koizumi' 'Tony
 Blair']
(1348, 62, 47)
```

Let's take a look at the principal axes that span this dataset. Because this is a large dataset, we will use `RandomizedPCA` —it contains a randomized method to approximate the first N principal components much more quickly than the standard `PCA` estimator, and thus is very useful for high-

dimensional data (here, a dimensionality of nearly 3,000). We will take a look at the first 150 components:

```
from sklearn.decomposition import RandomizedPCA
pca = RandomizedPCA(150)
pca.fit(faces.data)

RandomizedPCA(copy=True, iterated_power=3, n_components=150,
               random_state=None, whiten=False)
```

In this case, it can be interesting to visualize the images associated with the first several principal components:

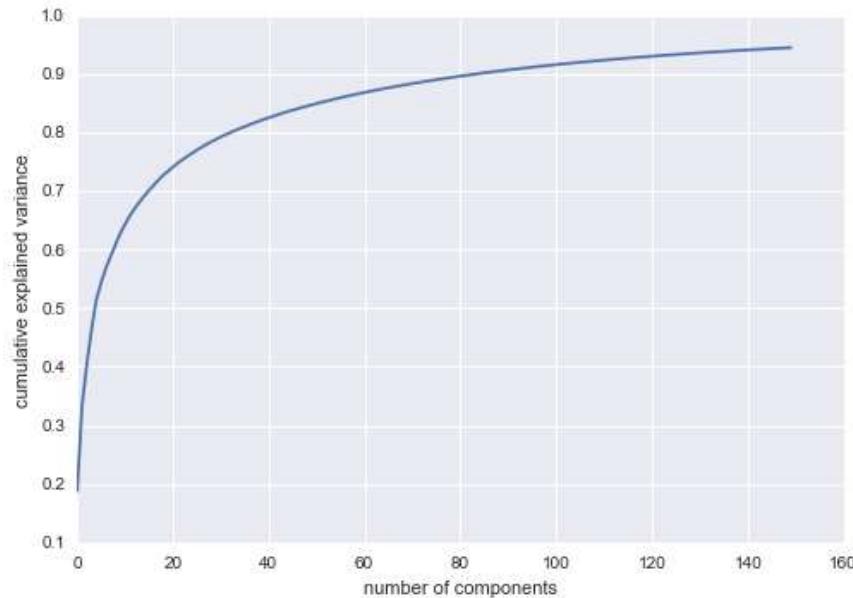
```
fig, axes = plt.subplots(3, 8, figsize=(9, 4),
                       subplot_kw={'xticks':[], 'yticks':[]},
                       gridspec_kw=dict(hspace=0.1, wspace=0.1))
for i, ax in enumerate(axes.flat):
    ax.imshow(pca.components_[i].reshape(62, 47), cmap='bone')
```



The results are very interesting, and give us insight into how the images vary: for example, the first few eigenfaces seem to be associated with the angle of lighting on the face, and later principal vectors seem to be picking out certain features, such as eyes, noses, and lips.

Let's take a look at the cumulative variance of these components to see how much of the data information the projection is preserving:

```
plt.plot(np.cumsum(pca.explained_variance_ratio_))
plt.xlabel('number of components')
plt.ylabel('cumulative explained variance')
```



We see that these 150 components account for just over 90% of the variance. That would lead us to believe that using these 150 components, we would recover most of the essential characteristics of the data. To make this more concrete, we can compare the input images with the images reconstructed from these 150 components:

**Pay only after placement**

AlmaBetter

```
pca = RandomizedPCA(150).fit(faces.data)
components = pca.transform(faces.data)
projected = pca.inverse_transform(components)

fig, ax = plt.subplots(2, 10, figsize=(10, 2.5),
                      subplot_kw={'xticks':[], 'yticks':[]},
                      gridspec_kw=dict(hspace=0.1, wspace=0.1))
for i in range(10):
    ax[0, i].imshow(faces.data[i].reshape(62, 47), cmap='binary_r')
    ax[1, i].imshow(projected[i].reshape(62, 47), cmap='binary_r')

ax[0, 0].set_ylabel('full-dim\\ninput')
ax[1, 0].set_ylabel('150-dim\\nreconstruction')
```



The top row here shows the input images, while the bottom row shows the reconstruction of the images from just 150 of the ~3,000 initial features. This visualization makes clear why the PCA feature selection used in was so successful: although it reduces the dimensionality of the data by nearly a factor of 20, the projected images contain enough information that we might, by eye, recognize the individuals in the image.

### [Also, read – 10 Machine Learning Projects to Boost your Portfolio](#)

What this means is that our classification algorithm needs to be trained on 150-dimensional data rather than 3,000-dimensional data, which depending on the particular algorithm we choose, can lead to a much more efficient classification.

### Follow Us:



**Aman Kharwal**

Coder with the ❤ of a Writer || Data Scientist | Solopreneur |  
Founder

ARTICLES: 1202





PREVIOUS



NEXT



## Recommended For You

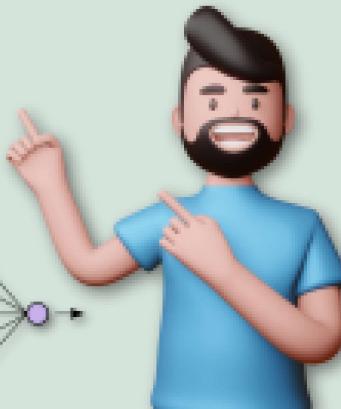
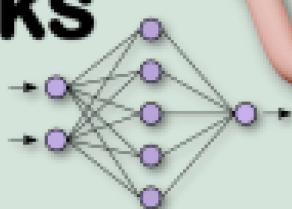
# The Best Approaches for Time Series Analysis



[Best Approaches for Time Series Analysis](#)

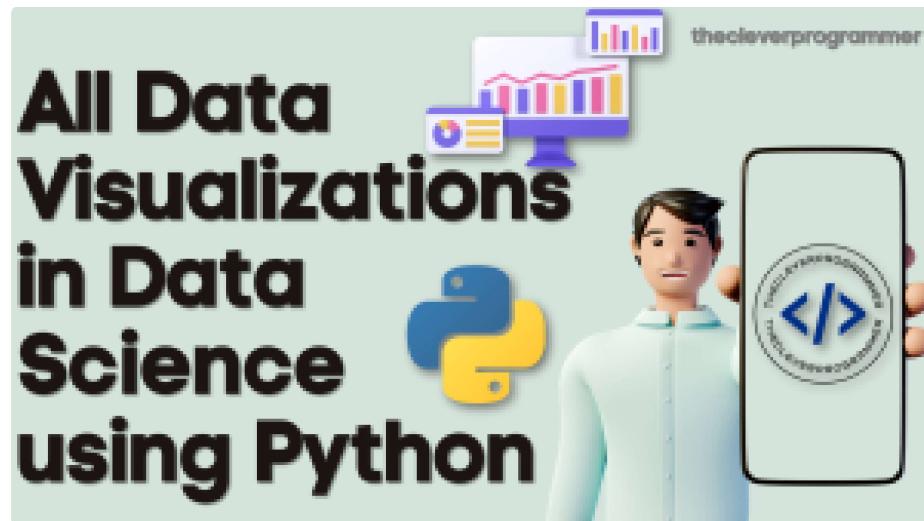
January 8, 2022

# Understand How Neural Networks Work



[How Neural Network Works](#)

January 6, 2022



## All Data Visualizations in Data Science using Python

January 5, 2022



## Visualize a Machine Learning Algorithm using Python

December 29, 2021

### Leave a Reply

Enter your comment here...