

Creado por:

Isabel Maniega

Análisis de Componentes Principales (PCA)



El análisis de componentes principales (PCA) es una técnica de reducción de dimensionalidad lineal que se puede utilizar para extraer información de un espacio de alta dimensión proyectándola en un subespacio de menor dimensión. Intenta preservar las partes esenciales que tienen más variación de los datos y eliminar las partes no esenciales con menos variación.

Una cosa importante a tener en cuenta sobre PCA es que es una técnica de reducción de dimensionalidad no supervisada, puede agrupar los puntos de datos similares en función de la correlación de características entre ellos sin supervisión (o etiquetas).

es un procedimiento estadístico que utiliza una transformación ortogonal para convertir un conjunto de observaciones de variables posiblemente correlacionadas (entidades cada una de las cuales toma varios valores numéricos) en un conjunto de valores de variables linealmente no correlacionadas llamadas componentes principales.

```
In [1]: from IPython import display  
display.Image("pca_df.png")
```

Out[1]:

Features / Attributes / Variables					
					
Samples					
	sepal-length	sepal-width	petal-length	petal-width	
	145	6.7	3.0	5.2	2.3
	146	6.3	2.5	5.0	1.9
	147	6.5	3.0	5.2	2.0
	148	6.2	3.4	5.4	2.3
149	5.9	3.0	5.1	1.8	

Pero, ¿dónde se puede aplicar PCA?

Visualización de datos: cuando se trabaja en cualquier problema relacionado con datos, el desafío en el mundo actual es el gran volumen de datos y las variables/características que definen esos datos. Para resolver un problema donde los

datos son la clave, necesita una amplia exploración de datos, como descubrir cómo se correlacionan las variables o comprender la distribución de algunas variables. Teniendo en cuenta que hay una gran cantidad de variables o dimensiones a lo largo de las cuales se distribuyen los datos, la visualización puede ser un desafío y casi imposible.

Por lo tanto, PCA permite visualizar los datos en un espacio 2D o 3D a simple vista.

Aceleración del algoritmo de aprendizaje automático (ML): dado que la idea principal de PCA es la reducción de la dimensionalidad, puede aprovecharla para acelerar el entrenamiento y el tiempo de prueba de su algoritmo de aprendizaje automático, teniendo en cuenta que sus datos tienen muchas características y que el aprendizaje del algoritmo ML es demasiado lento.

En un nivel abstracto, toma un conjunto de datos que tiene muchas características y simplifica ese conjunto de datos seleccionando algunas Principales Componentes de las características originales.

¿Qué es un componente principal?

Los componentes principales son la clave de PCA. En términos sencillos, cuando los datos se proyectan en una dimensión más baja (suponga tres dimensiones) desde un espacio más alto, las tres dimensiones no son más que los tres componentes principales que capturan (o contienen) la mayor parte de la variación (información) de sus datos .

Los componentes principales tienen dirección y magnitud. La dirección representa a través de qué ejes principales se distribuyen principalmente los datos o tienen la mayor variación y la magnitud indica la cantidad de variación que el Componente principal captura de los datos cuando se proyecta en ese eje. Los componentes principales son una línea recta y el primer componente principal tiene la mayor variación en los datos. Cada componente principal posterior es ortogonal al último y tiene una varianza menor. De esta forma, dado un conjunto de "x" variables correlacionadas sobre "y" muestras, se obtiene un conjunto de "u" componentes principales no correlacionados sobre las mismas "y" muestras.

La razón por la que obtiene componentes principales no correlacionados de las características originales es que las características correlacionadas contribuyen al mismo componente principal, reduciendo así las características de datos originales a componentes principales no correlacionados; cada uno representa un conjunto diferente de características correlacionadas con diferentes cantidades de variación.

Cada componente principal representa un porcentaje de la variación total capturada de los datos.

```
In [1]: from sklearn.datasets import load_breast_cancer  
  
breast = load_breast_cancer()
```

```
In [2]: breast_data = breast.data
```

```
In [3]: breast_data.shape
```

```
Out[3]: (569, 30)
```

```
In [4]: breast_labels = breast.target
```

```
In [5]: breast_labels.shape
```

```
Out[5]: (569,)
```

```
In [6]: import numpy as np
```

```
In [7]: labels = np.reshape(breast_labels,(569,1))
```

```
In [8]: final_breast_data = np.concatenate([breast_data,labels],axis=1)
```

```
In [9]: final_breast_data.shape
```

```
Out[9]: (569, 31)
```

```
In [10]: import pandas as pd
```

```
In [11]: breast_dataset = pd.DataFrame(final_breast_data)
breast_dataset
```

```
Out[11]:
```

	0	1	2	3	4	5	6	7	8	9
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.30010	0.14710	0.2419	0.07871
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.08690	0.07017	0.1812	0.05667
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.19740	0.12790	0.2069	0.05995
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.24140	0.10520	0.2597	0.09744
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.19800	0.10430	0.1809	0.05883
...
564	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.24390	0.13890	0.1726	0.05623
565	20.13	28.25	131.20	1261.0	0.09780	0.10340	0.14400	0.09791	0.1752	0.05533
566	16.60	28.08	108.30	858.1	0.08455	0.10230	0.09251	0.05302	0.1590	0.05646
567	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0.15200	0.2397	0.07016
568	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0.00000	0.1587	0.05884

569 rows × 31 columns



```
In [12]: features = breast.feature_names
```

```
In [13]: features
```

```
Out[13]: array(['mean radius', 'mean texture', 'mean perimeter', 'mean area',
               'mean smoothness', 'mean compactness', 'mean concavity',
               'mean concave points', 'mean symmetry', 'mean fractal dimension',
               'radius error', 'texture error', 'perimeter error', 'area error',
               'smoothness error', 'compactness error', 'concavity error',
               'concave points error', 'symmetry error',
               'fractal dimension error', 'worst radius', 'worst texture',
               'worst perimeter', 'worst area', 'worst smoothness',
               'worst compactness', 'worst concavity', 'worst concave points',
               'worst symmetry', 'worst fractal dimension'], dtype='<U23')
```

```
In [14]: features_labels = np.append(features, 'label')
```

```
In [15]: breast_dataset.columns = features_labels
```

```
In [16]: breast_dataset.head()
```

```
Out[16]:
```

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430

5 rows × 31 columns



```
In [17]: breast_dataset['label'].replace(0, 'Benign', inplace=True)
breast_dataset['label'].replace(1, 'Malignant', inplace=True)
```

/tmp/ipykernel_11724/2839743579.py:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.

The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always has as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
breast_dataset['label'].replace(0, 'Benign', inplace=True)
```

```
In [18]: breast_dataset.tail()
```

Out[18]:

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points
564	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.24390	0.13890
565	20.13	28.25	131.20	1261.0	0.09780	0.10340	0.14400	0.09791
566	16.60	28.08	108.30	858.1	0.08455	0.10230	0.09251	0.05302
567	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0.15200
568	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0.00000

5 rows × 31 columns



```
In [19]: from sklearn.preprocessing import StandardScaler
x = breast_dataset.loc[:, features].values
x = StandardScaler().fit_transform(x) # normalizing the features
```

In [20]: x.shape

Out[20]: (569, 30)

In [21]: np.mean(x), np.std(x)

Out[21]: (np.float64(-6.53516165871914e-17), np.float64(1.0))

In [22]: feat_cols = ['feature'+str(i) for i in range(x.shape[1])]

In [23]: normalised_breast = pd.DataFrame(x, columns=feat_cols)

In [24]: normalised_breast.tail()

	feature0	feature1	feature2	feature3	feature4	feature5	feature6	feature7
564	2.110995	0.721473	2.060786	2.343856	1.041842	0.219060	1.947285	2.3209
565	1.704854	2.085134	1.615931	1.723842	0.102458	-0.017833	0.693043	1.2636
566	0.702284	2.045574	0.672676	0.577953	-0.840484	-0.038680	0.046588	0.1057
567	1.838341	2.336457	1.982524	1.735218	1.525767	3.272144	3.296944	2.6588
568	-1.808401	1.221792	-1.814389	-1.347789	-3.112085	-1.150752	-1.114873	-1.2618

5 rows × 30 columns



```
In [25]: from sklearn.decomposition import PCA
pca_breast = PCA(n_components=2)
principalComponents_breast = pca_breast.fit_transform(x)
```

```
In [26]: principal_breast_Df = pd.DataFrame(data = principalComponents_breast
, columns = ['principal component 1', 'principal component 2'])
```

In [27]: principal_breast_Df.tail()

Out[27]:

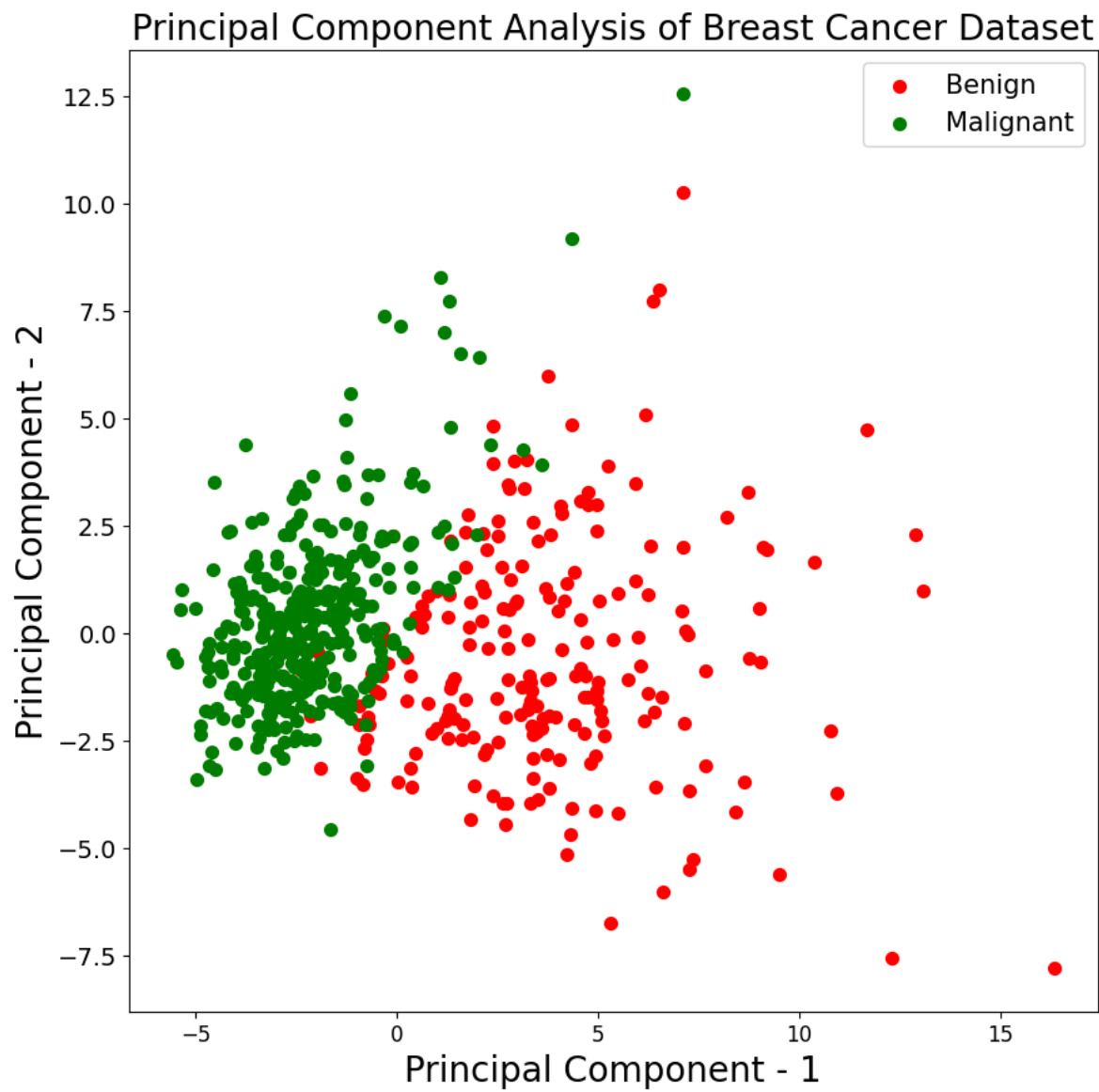
	principal component 1	principal component 2
564	6.439315	-3.576817
565	3.793382	-3.584048
566	1.256179	-1.902297
567	10.374794	1.672010
568	-5.475243	-0.670637

In [28]: `print('Explained variation per principal component: {}'.format(pca_breast`
 Explained variation per principal component: [0.44272026 0.18971182]

In [29]: `import matplotlib.pyplot as plt`

In [30]: `plt.figure()
 plt.figure(figsize=(10,10))
 plt.xticks(fontsize=12)
 plt.yticks(fontsize=14)
 plt.xlabel('Principal Component - 1',fontsize=20)
 plt.ylabel('Principal Component - 2',fontsize=20)
 plt.title("Principal Component Analysis of Breast Cancer Dataset",fontsize=20)
 targets = ['Benign', 'Malignant']
 colors = ['r', 'g']
 for target, color in zip(targets,colors):
 indicesToKeep = breast_dataset['label'] == target
 plt.scatter(principal_breast_Df.loc[indicesToKeep, 'principal component 1'],
 principal_breast_Df.loc[indicesToKeep, 'principal component 2'],
 color=color,
 s=150)
 plt.legend(targets,prop={'size': 15})`

Out[30]: <matplotlib.legend.Legend at 0x7137f61ae8d0>
 <Figure size 640x480 with 0 Axes>



Creado por:

Isabel Maniega