

Principal Components Analysis

Descripción del Dataset

El dataset Breast Cancer Wisconsin (Diagnostic) pertenece al área médica y se utiliza para la detección de cáncer de mama.

El objetivo del dataset es clasificar tumores como:

- 0 = Maligno
 - 1 = Benigno
- Los datos fueron obtenidos a partir de imágenes digitalizadas de muestras tomadas mediante biopsia (Fine Needle Aspirate).
- El dataset contiene 30 variables predictoras numéricas. Estas variables describen características de las células, como:
- Radio
 - Textura
 - Perímetro
 - Área
 - Suavidad
 - Compacidad
 - Concavidad
 - Simetría
 - Dimensión fractal
- Cada una de estas características se calcula en tres versiones:
- Promedio (mean)
 - Error estándar (se)
 - Peor valor observado (worst)
- En total, el dataset contiene 569 observaciones.

id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave_points_mean	symmetry_mean	radius_worst	texture_worst	perimeter_worst	area_worst	smoothness_worst	compactness_worst	concavity_worst
0	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	...	17.33	184.60	2019.0	0.1622	0.6656	0.71	
1	M	20.57	17.77	132.90	1326.0	0.09474	0.07964	0.0869	0.07017	...	23.41	158.80	1956.0	0.1238	0.1866	0.24	
2	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	...	25.53	152.50	1709.0	0.1444	0.4245	0.45	
3	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	...	26.50	98.87	567.7	0.2098	0.8663	0.68	
4	M	20.29	14.34	135.10	1297.0	0.10300	0.13280	0.1980	0.10430	...	16.67	152.20	1575.0	0.1374	0.2056	0.40	

5 rows × 31 columns

Limpieza de datos

El conjunto de datos contiene una columna de identificación ("id") que no proporciona información predictiva. Además, el conjunto de datos incluye una columna vacía ("Sin nombre: 32") que solo contiene valores faltantes.

Estas columnas se eliminarán, ya que no contribuyen a la tarea de clasificación.

id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave_points_mean	symmetry_mean	radius_worst	texture_worst	perimeter_worst	area_worst	smoothness_worst	compactness_worst
0	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	...	17.33	184.60	2019.0	0.1622	0.6656	0.71
1	M	20.57	17.77	132.90	1326.0	0.09474	0.07964	0.0869	0.07017	...	23.41	158.80	1956.0	0.1238	0.1866	0.24
2	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	...	25.53	152.50	1709.0	0.1444	0.4245	0.45
3	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	...	26.50	98.87	567.7	0.2098	0.8663	0.68
4	M	20.29	14.34	135.10	1297.0	0.10300	0.13280	0.1980	0.10430	...	16.67	152.20	1575.0	0.1374	0.2056	0.40

5 rows × 31 columns

Ahora verificamos si el conjunto de datos contiene valores faltantes.

id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave_points_mean	symmetry_mean	radius_worst	texture_worst	perimeter_worst	area_worst	smoothness_worst	compactness_worst
0	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	...	17.33	184.60	2019.0	0.1622	0.6656	0.71
1	M	20.57	17.77	132.90	1326.0	0.09474	0.07964	0.0869	0.07017	...	23.41	158.80	1956.0	0.1238	0.1866	0.24
2	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	...	25.53	152.50	1709.0	0.1444	0.4245	0.45
3	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	...	26.50	98.87	567.7	0.2098	0.8663	0.68
4	M	20.29	14.34	135.10	1297.0	0.10300	0.13280	0.1980	0.10430	...	16.67	152.20	1575.0	0.1374	0.2056	0.40

5 rows × 31 columns

Ahora verificamos si el conjunto de datos contiene valores faltantes.

id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave_points_mean	symmetry_mean	radius_worst	texture_worst	perimeter_worst	area_worst	smoothness_worst	compactness_worst
0	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	...	17.33	184.60	2019.0	0.1622	0.6656	0.71
1	M	20.57	17.77	132.90	1326.0	0.09474	0.07964	0.0869	0.07017	...	23.41	158.80	1956.0	0.1238	0.1866	0.24
2	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	...	25.53	152.50	1709.0	0.1444	0.4245	0.45
3	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	...	26.50	98.87	567.7	0.2098	0.8663	0.68
4	M	20.29	14.34	135.10	1297.0	0.10300	0.13280	0.1980	0.10430	...	16.67	152.20	1575.0	0.1374	0.2056	0.40

5 rows × 31 columns

Análisis de datos

Es importante conocer cuántos casos son benignos y cuántos malignos, ya que esto nos permite entender si el dataset está balanceado.

id	diagnosis	count
0	0	357
1	1	212

dtype: int64

La variable objetivo "diagnóstico" es categórica (M = Maligno, 0 = Benigno). Se codificará como:

M → 0 → 1

Esta codificación convierte el problema en una tarea de clasificación binaria.

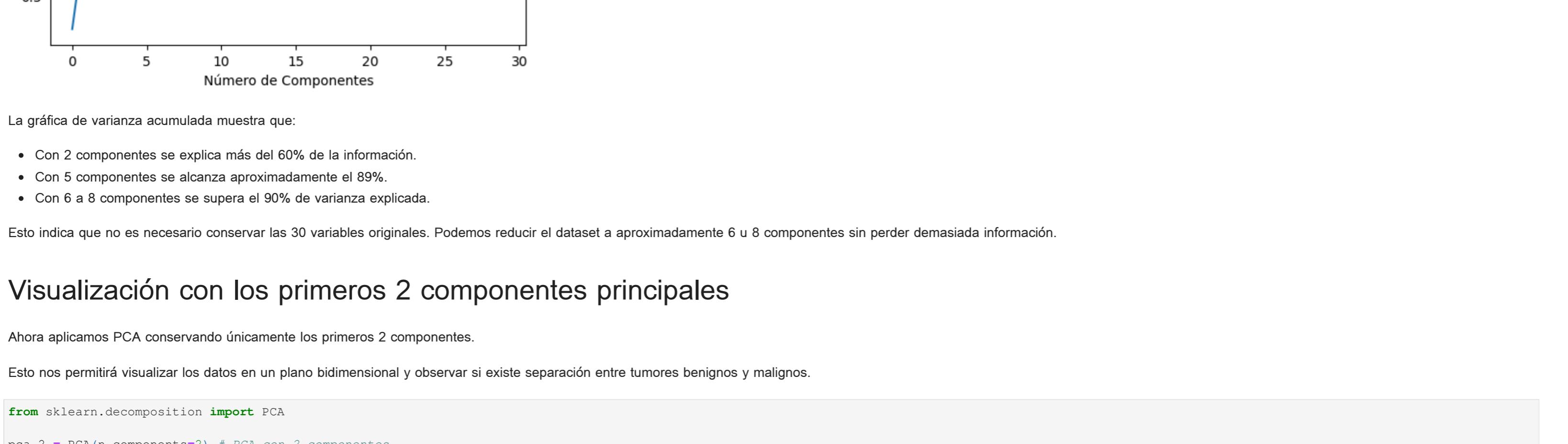
id	diagnosis	map	df
0	M	0	df["diagnosis"] = df["diagnosis"].map({0: "M", 1: "B"})
1	B	1	df["diagnosis"] = df["diagnosis"].map({0: "B", 1: "M"})
2	M	0	df["diagnosis"] = df["diagnosis"].map({0: "M", 1: "B"})
3	B	1	df["diagnosis"] = df["diagnosis"].map({0: "B", 1: "M"})
4	M	0	df["diagnosis"] = df["diagnosis"].map({0: "M", 1: "B"})

5 rows × 31 columns

Se observa que el dataset contiene más casos benignos que malignos.

Sin embargo, la diferencia no es extremadamente grande, por lo que el dataset puede considerarse razonablemente balanceado.

Ahora visualizamos la distribución de la variable objetivo mediante una gráfica.



A continuación observaremos estadísticas descriptivas de las variables predictoras.

Este nos permitirá entender la escala y la variabilidad de los datos.

id	diagnosis	