

▼ Algoritmo para analisis de componentes principales

```
1 from cmath import sqrt
2 import numpy as np
3 import numpy.linalg as alg
4
5 from sklearn.decomposition import PCA
6
7 import seaborn as sns
8 from matplotlib import pyplot as plt
9
10 import pandas as pd
11
12 np.set_printoptions(precision=6, suppress=True)
```

▼ La siguiente clase implementa el algoritmo de analisis de componentes principales

Para la implementación del algoritmo se siguen los siguientes pasos:

- Centrar y reducir los datos
- Calcular la matriz de correlación
- Calcular los valores y vectores propios de la matriz de correlación
- Ordenar los vectores propios basado en la magnitud de sus valores propios asociados
- Selección de las n columnas utilizadas como componente principal

A Este punto se pueden calcular metricas para la eficiencia del algoritmo, como el vector de inercia que representa el porcentaje de importancia para cada columna de la transformación obtenida.

```
1 class My_PCA:
2     def __init__(self, n_components = -1):
3         self.n_components = n_components
4         self.matrix = []
5         self.inertia = []
6         self.points = []
7         self.eigenvalues = []
8         self.eigenvectors = []
9     pass
10
11     def fit(self, X):
12         redux = self.center_and_scale(X)
```

```
13     corr = self.corr_matrix(redux)
14     self.eigenvalues, self.eigenvectors = self.eigensomething(corr)
15     self.matrix = self.pca_matrix(self.eigenvectors)
16     self.remove_components()
17     self.inertia = self.calculate_inertia(self.eigenvalues)
18     self.points = self.calculate_points(self.eigenvalues)
19     pass
20
21
22     def transform(self, X):
23         return np.matmul(X, self.matrix)
24
25
26     def center_and_scale(self, X):
27         _X = X
28
29         mmean = np.mean(_X, axis=0)
30         mstd = np.std(_X, axis=0)
31
32         _X = (_X - mmean) / mstd
33         return _X
34
35
36     def corr_matrix(self, X):
37         n, m = X.shape
38         return (1/n) * np.matmul(X.transpose(), X)
39
40
41     def eigensomething(self, R):
42         w, v = alg.eigh(R)
43         sort_index = np.argsort(abs(w))[:, -1]
44         sorted_eigenvals = w[sort_index]
45         sorted_eigenvecs = v[:, sort_index]
46
47         return (sorted_eigenvals, sorted_eigenvecs)
48
49
50     def pca_matrix(self, vectors):
51         return np.array(vectors)
52
53
54     def remove_components(self):
55         self.matrix = np.delete(self.matrix, np.s_[self.n_components - 1 : -1], axis = 1)
56         self.eigenvalues = self.eigenvalues[0:self.n_components]
57         pass
58
59     def calculate_inertia(self, eigenvalues):
60         n, m = self.matrix.shape
61         return eigenvalues / m
62
63
```

```

64     def calculate_points(self, eigenvalues):
65         w0 = eigenvalues[0]
66         w1 = eigenvalues[1]
67         v0 = self.matrix[:,0]
68         v1 = self.matrix[:,1]
69
70         return ((v0 * sqrt(w0)).real, (v1 * sqrt(w1)).real)

```

▼ Poniendo a prueba la imlementación

Se importará el dataset del Titanic para analizar cuales de sus características pueden ser de interes

```

1 # Importación del dataset
2 base_dataset = pd.read_csv("/content/titanic.csv")
3 print(base_dataset.head())
4
5 # Búsqueda de valores ausentes
6 percent_missing = base_dataset.isnull().sum() * 100 / len(base_dataset)
7 missing_value_df = pd.DataFrame({'column_name': base_dataset.columns,
8                                  'percent_missing': percent_missing})
9
10 print(missing_value_df)

```

	PassengerId	Survived	Pclass	\
0	1	0	3	
1	2	1	1	
2	3	1	3	
3	4	1	1	
4	5	0	3	

	Name	Sex	Age	SibSp	\
0	Braund, Mr. Owen Harris	male	22.0	1	
1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	
2	Heikkinen, Miss. Laina	female	26.0	0	
3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	
4	Allen, Mr. William Henry	male	35.0	0	

	Parch	Ticket	Fare	Cabin	Embarked
0	0	A/5 21171	7.2500	NaN	S
1	0	PC 17599	71.2833	C85	C
2	0	STON/O2. 3101282	7.9250	NaN	S
3	0	113803	53.1000	C123	S
4	0	373450	8.0500	NaN	S

	column_name	percent_missing
PassengerId	PassengerId	0.000000
Survived	Survived	0.000000
Pclass	Pclass	0.000000
Name	Name	0.000000

Sex	Sex	0.000000
Age	Age	20.091673
SibSp	SibSp	0.000000
Parch	Parch	0.000000
Ticket	Ticket	0.000000
Fare	Fare	0.076394
Cabin	Cabin	77.463713
Embarked	Embarked	0.152788

▼ Información básica

Podemos observar que hay datos útiles para ubicar a cada uno de los pasajeros, sin embargo son datos que no serán útiles para la clasificación de si este pasajero sobrevive o no, por lo que las columnas "PassengerId", "Name", "Ticket" ser completamente eliminadas, también podemos notar que la columna de Cabina tiene un 77% de información faltante por lo que podemos eliminarla completamente.

Codificación de información nominal

Las clases "Sex", "Embarked" y "Survived" se pueden codificar con la técnica de One Hot Encoding para facilitar el proceso de clasificación a futuro, la variable "Pclass" puede ser más compleja de tratar ya que es una variable ordinal, sin embargo no es una variable continua por lo que será mejor evitar la complejidad de esta variable y codificarla igualmente.

```

1 # Eliminación de variables innecesarias
2 titanic_data = base_dataset.drop(["PassengerId", "Name", "Cabin", "Ticket"], axis=1)
3 titanic_data.dropna(axis=0, how="any", inplace=True)
4
5 # One hot encoding
6 cnames = ["Survived", "Pclass", "Sex", "Embarked"]
7 for cname in cnames:
8     dummies = pd.get_dummies(titanic_data[cname], prefix=cname)
9     titanic_data = titanic_data.drop(cname, axis=1)
10    titanic_data = titanic_data.join(dummies)

```

▼ Preparación de los datos para utilizar PCA

Para poder procesar los datos con el algoritmo de PCA es necesario transformar los datos de un Pandas dataset a una matriz numérica de NumPy.

En la implementación de PCA propia se considera que los datos serán centrados y reducidos dentro de la función de 'fit' sin embargo al utilizar PCA implementado por sklearn los datos deben ser centrados y reducidos antes de ser introducidos en el algoritmo de PCA.

```

1 # Conversion de datos a numpy matrix
2 titanic_as_np = titanic_data.to_numpy()
3
4 # Uso del metodo PCA implementado
5 mypca = My_PCA(n_components=4)
6 mypca.fit(titanic_as_np)
7
8 # Uso del metodo PCA de la biblioteca sklearn
9 pca = PCA(n_components=4)
10
11 # Centrado y reducido previo al uso de sklearn
12 titanic_prep = mypca.center_and_scale(titanic_as_np)
13 pca.fit(titanic_prep)

PCA(n_components=4)

```

▼ Uso del algoritmo de PCA sobre los datos originales

Para modificar la matriz de datos originales solo es necesario centrar y reducir los datos, luego multiplicar la matriz resultante por la matriz obtenida con el algoritmo de PCA

```

1 # titanic_prep se encuentra centrada y reducida del bloque de codigo anterior
2
3 # Transformación de los datos a partir del PCA implementado
4 my_matrix = mypca.transform(titanic_prep)
5
6 # Transformacion de los daatos a partir del PCA sklearn
7 sk_matrix = pca.transform(titanic_prep)

```

▼ Graficación de los resultados

Podemos ver como los datos se agrupan en 2 categorías principales, sin embargo existe un tercer grupo donde los datos se encuentran mezclados, aún así este grupo es menor a los 2 grupos principales, esto se logra al transformar las coordenadas de cada punto aumentando al maximo la varianza entre las clases.

```

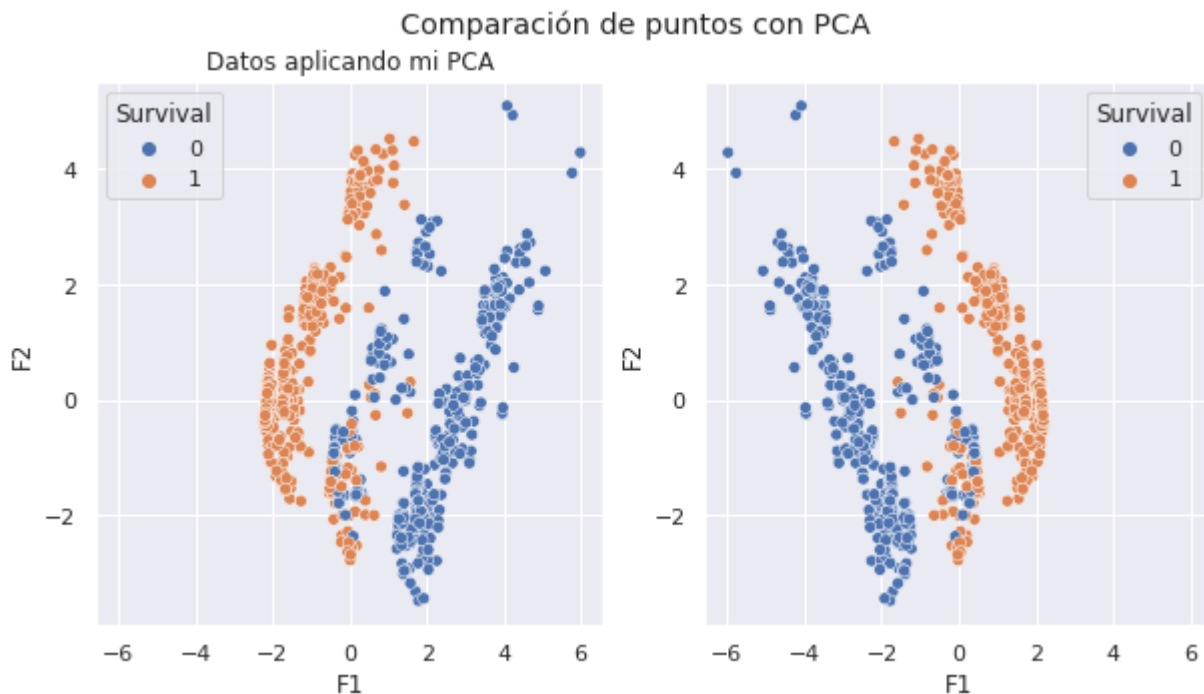
1 transformed_sk_data = pd.DataFrame({'F1' : sk_matrix[:, 0], 'F2': sk_matrix[:, 1], 'Surviv
2 transformed_my_data = pd.DataFrame({'F1' : my_matrix[:, 0], 'F2': my_matrix[:, 1], 'Surviv
3 sns.set()
4
5 fig, axes = plt.subplots(1, 2, sharex=True, figsize=(10,5))
6 fig.suptitle('Comparación de puntos con PCA')
7
8 sns.scatterplot(ax=axes[0], x="F1", y="F2", hue='Survival', data=transformed_sk_data)

```

```

9 axes[0].set_title('Datos aplicando PCA de sklearn')
10
11 sns.scatterplot(ax=axes[1], x="F1", y="F2", hue='Survival', data=transformed_my_data)
12 axes[0].set_title('Datos aplicando mi PCA')
13
14 plt.show()

```



▼ Inercia e importancia de características

A partir de la inercia podemos saber que tan importante será cada columna de datos final, pero también podemos obtener esta información de manera gráfica dentro de un círculo de correlación, esto nos dirá cuánta correlación o "importancia" tendrá cada columna para obtener un resultado.

Para este caso podemos ver que las variables con una correlación más fuerte son:

- Supervivencia
- Sexo
- Clase
- Donde ha embarcado
- Tarifa
- Edad

```

1 x, y = mypca.points
2
3 (fig, ax) = plt.subplots(figsize=(10, 10))
4 for i in range(0, len(titanic_data.columns)):
5     ax.arrow(0, 0,

```

```

6         x[i],y[i],
7         head_width=0.1,head_length=0.1)
8     plt.text(x[i],y[i],titanic_data.columns[i])
9
10 an = np.linspace(0, 2 * np.pi, 100)
11 plt.plot(np.cos(an), np.sin(an))
12 plt.axis('equal')
13 ax.set_title('Circulo de correlación')
14 plt.show()

```

