

Sistemas de Información y Telemedicina. *

Marta Girones Sanguesa

Silvia Marset Gomis

Ignacio Amat Hernández

Sofía Gutiérrez Santamaría

November 28, 2019

Contents

Sección	Página
1 Preámbulo	3
2 Histogramas	4
3 Kernel Density	5
4 Boxplot	6
5 QQplot	7
6 Corrplot	8
7 Filter Methods	10
8 Wrapper Methods	11
9 PCA	12
9.1 Pareto	12

*Grado en Ingeniería Biomédica, Escuela Técnica Superior de Ingenieros Industriales, Valencia, España.

List of Figures

1	Histogramas para datos con y sin anomalías.	4
2	Kernel Density para datos con y sin anomalías.	5
3	Boxplots para datos con y sin anomalías.	6
4	QQplots para datos con y sin anomalías.	7
5	Corrplot para datos con anomalías.	8
6	Corrplot para datos sin anomalías.	9
7	Diagrama de Pareto.	12

Listings

1	Importaciones iniciales y preparacion de datos.	3
2	Código generador de los histogramas con datos anómalos.	4
3	Código generador de los kernel density plots con datos anómalos.	5
4	Código generador de los boxplots con datos anómalos.	6
5	Código generador de los QQplots con datos anómalos.	7
6	Código generador de los corrplots con datos anómalos.	8
7	Aplicación métodos <i>filter</i> de selección características.	10
8	Ranking de variables según los métodos filter.	10
9	Aplicación métodos <i>wrapper</i> de selección características.	11
10	Resultados del filtrado mediante wrappers.	11
11	<i>Principal Component Analysis</i>	12
12	Código generador del diagrama de Pareto	12

1 Preámbulo

```
1 import numpy as np
2 from scipy import stats
3
4 # names of variables
5 labels = ['age', 'leptin', 'bmi', 'adiponectin', 'glucose',
6           'resistin', 'insulin', 'MCP1', 'HOMA']
7
8 # loads data
9 data = np.loadtxt (open (r'../..data.csv', 'rb'), delimiter = ',')
10
11 # rewrites data as all the rows of data w/out nan cells
12 data = data [~np.isnan (data).any (axis=1)]
13
14 # separates parameters into matrix x
15 x = np.array ([list (data [x][: -1]) for x in range (len (data))])
16
17 # and class (1, 2) into vector y
18 y = np.array ([int (data [x][ -1]) for x in range (len (data))])
19
20 # removes outliers
21 data_no = data [(np.abs (stats.zscore (data)) < 3).all (axis = 1)]
22 # ↑ = No Outliers
23
24 x_no = np.array ([list (data_no [x][: -1]) for x in range (len (data_no))])
25 y_no = np.array ([int (data_no [x][ -1]) for x in range (len (data_no))])
```

Listing 1: Importaciones iniciales y preparacion de datos.

2 Histogramas

En este apartado dibujamos los histogramas comparativos.

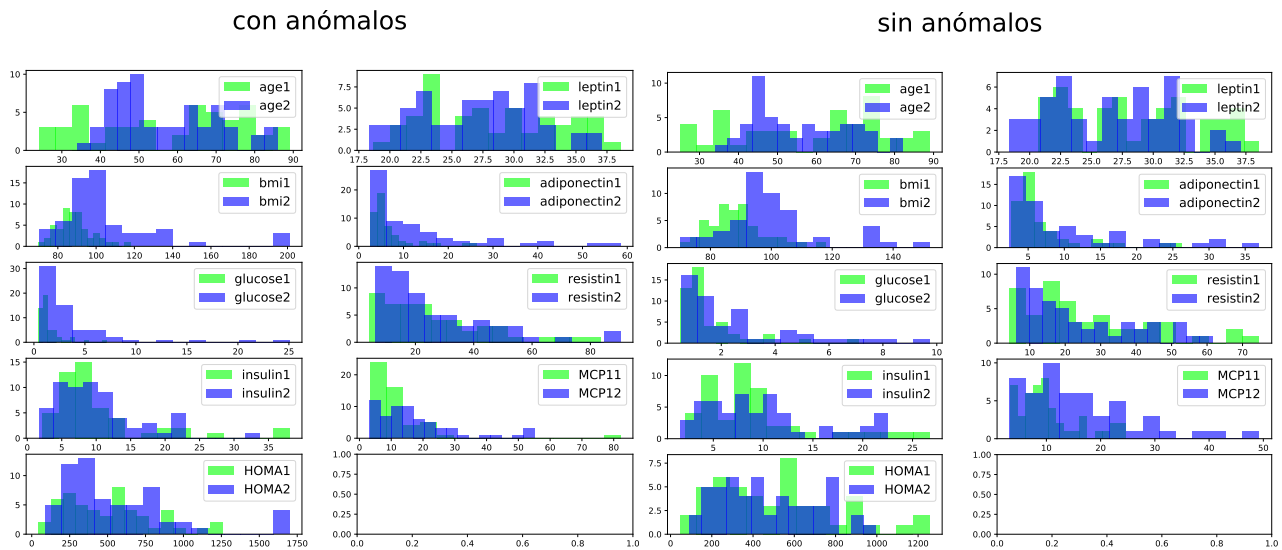


Fig. 1: Histogramas para datos con y sin anomalías.

```

1  import matplotlib as mpl
2  import matplotlib.pyplot as plt
3
4  # load preprocessed data, x and y are raw, x_no and y_no contain no outliers
5  from preprocessing import x, y, x_no, y_no, labels
6
7  # colours for the histograms
8  fc = [(0, 1, 0, 0.6), (0, 0, 1, 0.6)]
9  #      (R, G, B,  α )← transparency
10
11 fig, ax = plt.subplots (nrows = 5, ncols = 2, figsize = (13, 10))
12 ax = ax.flatten ()
13
14 # draws each of the histograms, two for each variable
15 for i in range (0, 9):
16     for j in [1, 2]:
17         ax[i].hist (x [y == j, i], bins = 15, fc = fc [j], label = labels [i] + str \
18                     (j))
19         ax[i].legend (loc = 1, prop={'size': 15})
20
21 fig.suptitle ('con anomalías', fontsize = 30)
22 fig.savefig ('../images/hist.pdf', bbox_inches = 'tight', pad_inches = 0)

```

Listing 2: Código generador de los histogramas con datos anómalos.

3 Kernel Density

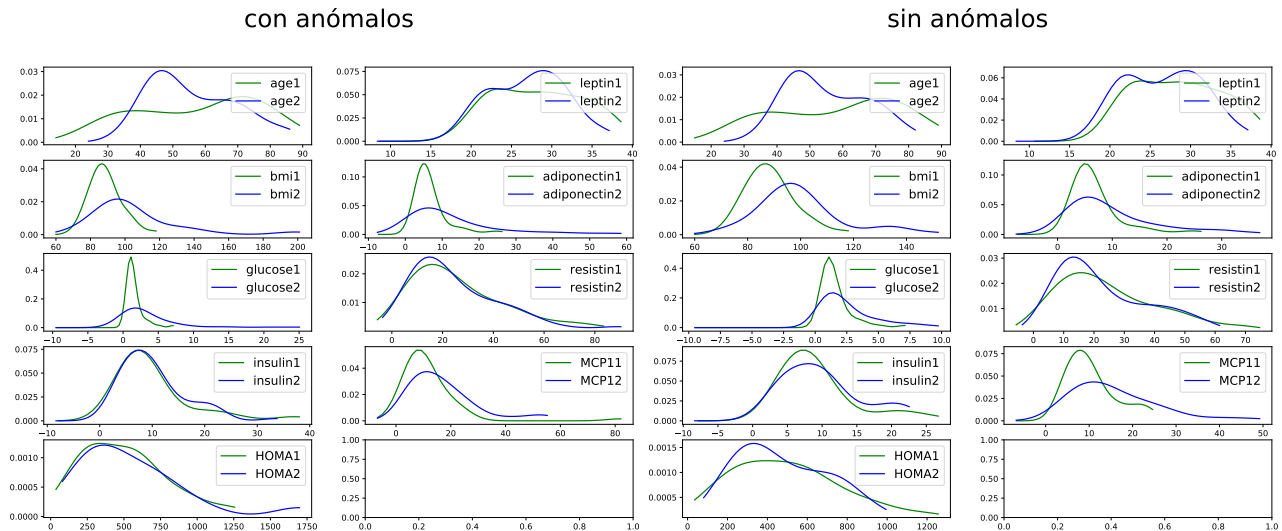


Fig. 2: Kernel Density para datos con y sin anomalías.

```

1  import matplotlib as mpl
2  import matplotlib.pyplot as plt
3  import numpy as np
4  from scipy.stats import gaussian_kde
5
6  # load preprocessed data, x and y are raw, x_no and y_no contain no outliers
7  from preprocessing import x, y, x_no, y_no, labels
8
9  # colours
10 fc = ['', 'green', 'blue']
11
12 fig, ax = plt.subplots (nrows = 5, ncols = 2, figsize = (13, 10))
13 ax = ax.flatten ()
14
15 # same loop in principle as before
16 for i in range (0, 9):
17     for j in [1, 2]:
18         kde = gaussian_kde (x_ := x [y == j, i])
19         xs = np.linspace(np.min (x_) - 10, np.max (x_), num=len (x_))
20         ax[i].plot (xs, kde(xs), c = fc[j], label = labels [i] + str (j))
21         ax[i].legend (loc = 1, prop={'size': 15})
22
23 fig.suptitle ('con anomalos', fontsize = 30)
24 fig.savefig ('../images/kden.pdf', bbox_inches = 'tight', pad_inches = 0)

```

Listing 3: Código generador de los kernel density plots con datos anómalos.

4 Boxplot

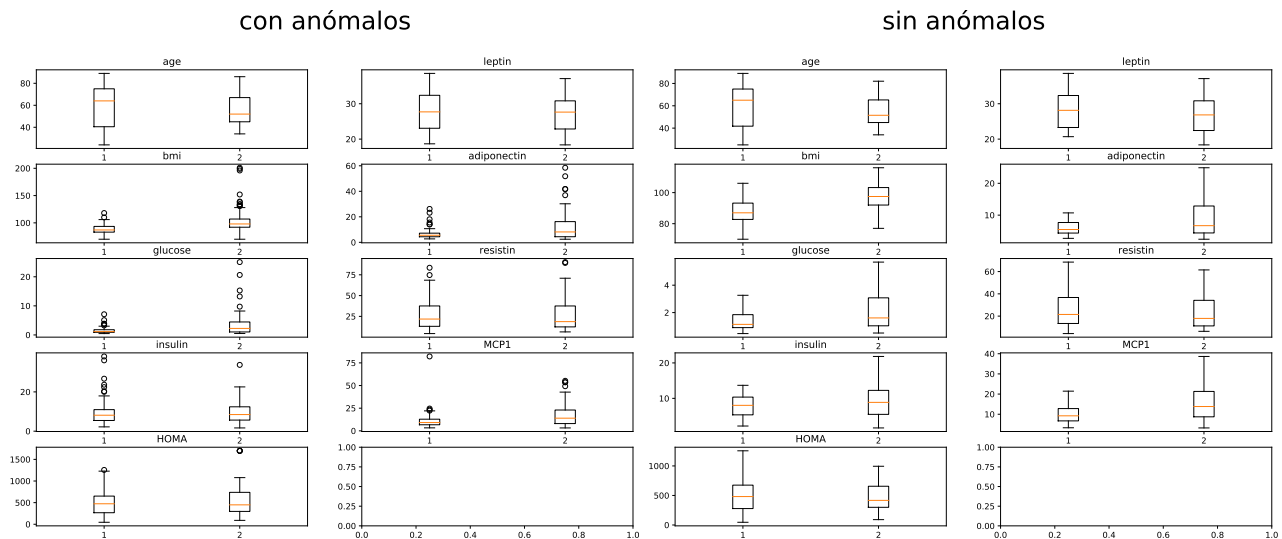


Fig. 3: Boxplots para datos con y sin anomalías.

```

1 import matplotlib as mpl
2 import matplotlib.pyplot as plt
3
4 # load preprocessed data, x and y are raw, x_no and y_no contain no outliers
5 from preprocessing import x, y, x_no, y_no, labels
6
7 fig, ax = plt.subplots (nrows = 5, ncols = 2, figsize = (13, 10))
8 ax = ax.flatten ()
9
10 for i in range (0, 9):
11     ax[i].boxplot ([x [y == 1, i], x [y == 2, i]])
12     ax[i].title.set_text (labels [i])
13
14 fig.suptitle ('con anomalos', fontsize = 30)
15 fig.savefig ('../images/boxp.pdf', bbox_inches = 'tight', pad_inches = 0)

```

Listing 4: Código generador de los boxplots con datos anómalos.

5 QQplot

con anomalías

sin anomalías

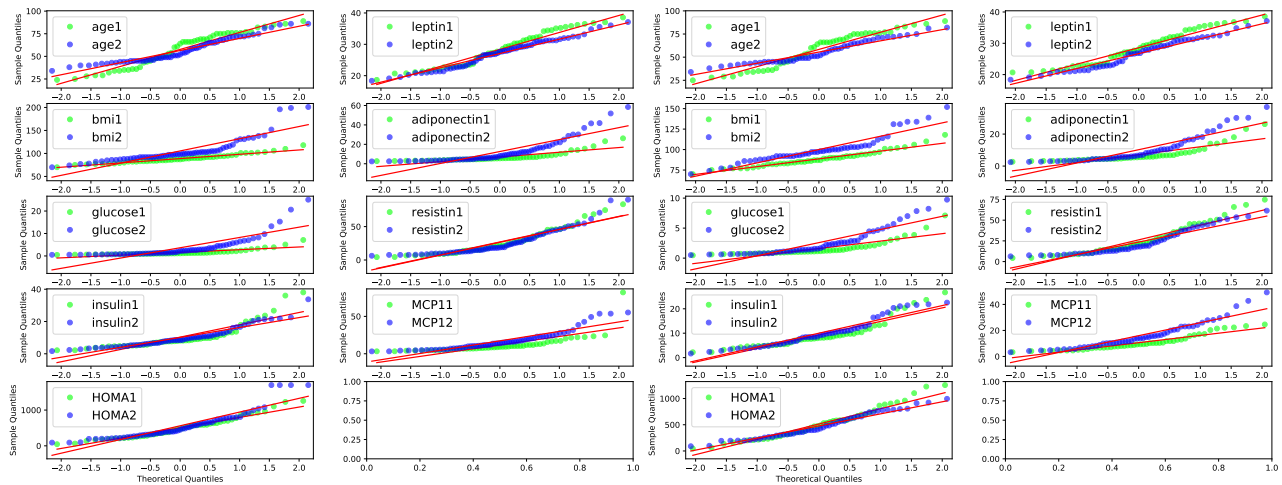


Fig. 4: QQplots para datos con y sin anomalías.

```

1 import matplotlib as mpl
2 import matplotlib.pyplot as plt
3
4 # load preprocessed data, x and y are raw, x_no and y_no contain no outliers
5 from preprocessing import x, y, x_no, y_no, labels
6 import statsmodels.api as sm
7
8 fc = [(0, (0, 1, 0, 0.6), (0, 0, 1, 0.6))]
9
10 fig, ax = plt.subplots (nrows = 5, ncols = 2, figsize = (13, 10))
11 ax = ax.flatten ()
12
13 for i in range (0, 9):
14     for j in [1, 2]:
15         sm.qqplot (x [y == j, i], ax = ax[i], c = fc[j],
16                     line = 's', label = labels [i] + str (j))
17         ax[i].legend (loc = 2, prop={'size': 15})
18
19 fig.suptitle ('con anomalías', fontsize = 30)
20 fig.savefig ('../images/qqp.pdf', bbox_inches = 'tight', pad_inches = 0)

```

Listing 5: Código generador de los QQplots con datos anómalos.

6 Corrplot

```
1 dataframe = pd.DataFrame.from_records(x)
2 sns.pairplot(dataframe, kind = 'reg')
3 plt.suptitle ('con anomalos', fontsize = 30)
4 plt.savefig ('../images/corrp.pdf', bbox_inches = 'tight', pad_inches = 0)
```

Listing 6: Código generador de los corrplots con datos anómalos.

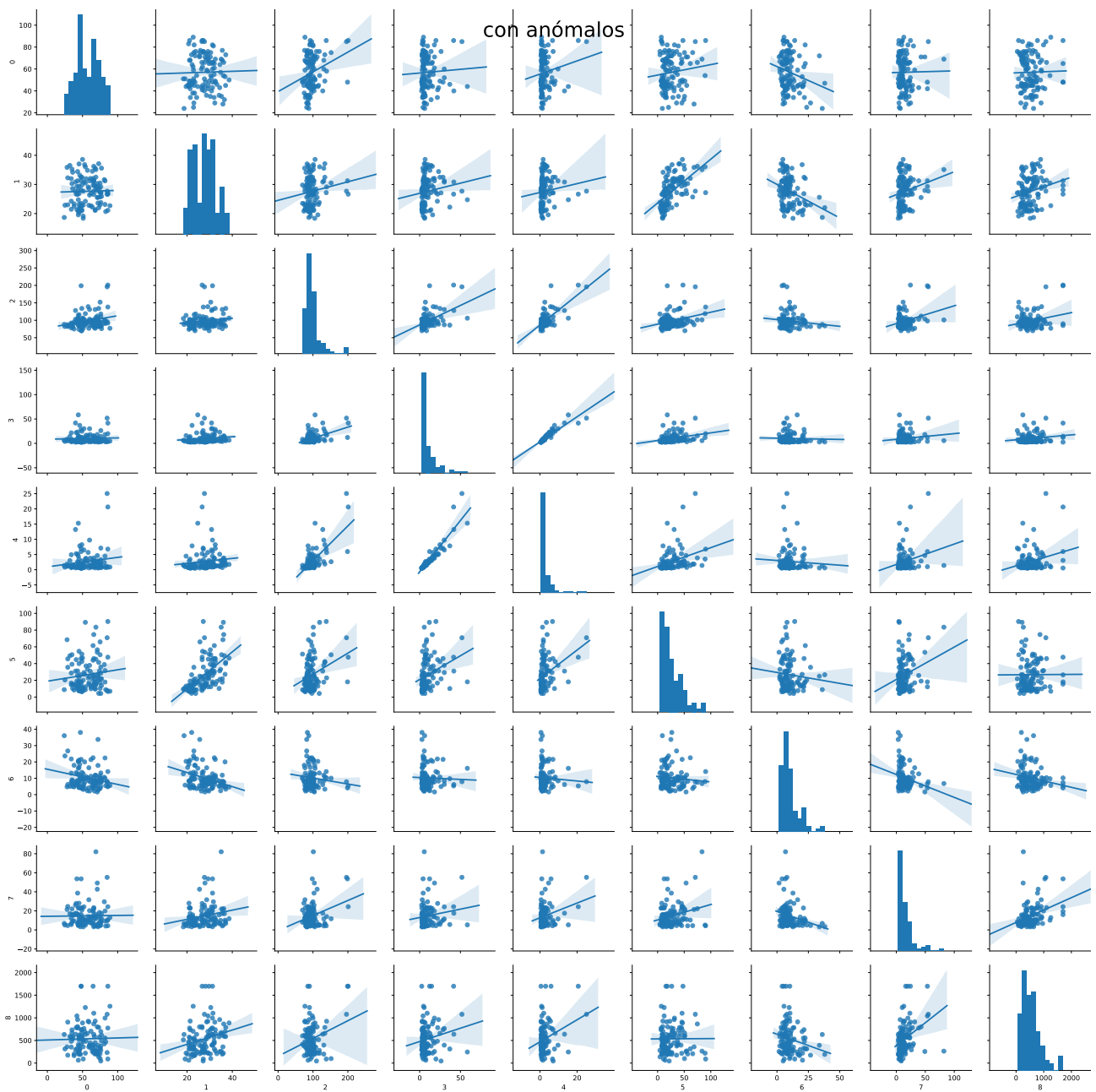


Fig. 5: Corrplot para datos con anomalias.

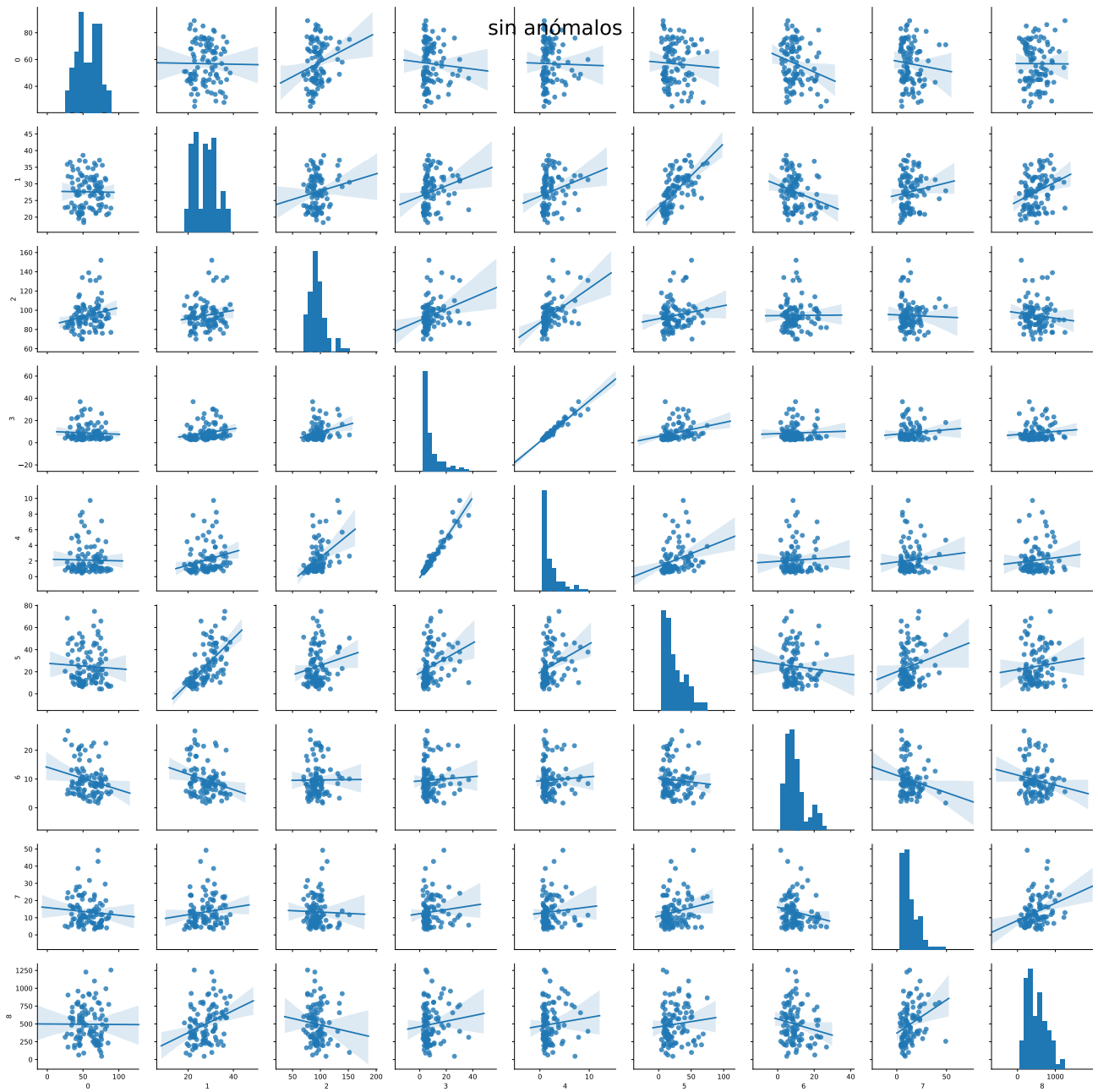


Fig. 6: Corplot para datos sin anomalías.

7 Filter Methods

```
1 # Filter Methods
2 import sklearn.feature_selection as sk
3
4 Fscore, pval = sk.f_classif (x_no, y_no)
5 r1 = Fscore.argsort().argsort() # fscore rank
6 print (r1+1)
7
8 import ReliefF as rl
9
10 r2 = rl.ReliefF (n_neighbors = 1) # relieff rank
11 r2.fit(x_no, y_no)
12 r2 = r2.top_features
13 print (r2+1)
14
15 diferencias = abs (r1-r2)
16 media = np.mean (diferencias)
```

Listing 7: Aplicación métodos *filter* de selección características.

```
1 [4 5 9 6 7 3 1 8 2] -> fscore
2 [1 9 8 7 6 5 4 2 3] -> relieff
3 [3 4 1 1 1 2 3 6 1] -> diferencias
4 2.4444444444444446 -> media
```

Listing 8: Ranking de variables según los métodos filter.

8 Wrapper Methods

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.neighbors import KNeighborsClassifier
3 from mlxtend.feature_selection import SequentialFeatureSelector
4
5 knn = KNeighborsClassifier (n_neighbors = 50)
6
7 sfs = SequentialFeatureSelector (knn,
8                                 k_features = 4,
9                                 forward = True,
10                                scoring = 'accuracy',
11                                cv = 10)
12
13 sfs.fit (x_no, y_no, custom_feature_names = labels)
14 print (sfs.k_score_)
15 print ('Sequential Forward Selection', sfs.k_feature_names_, end = '\n\n')
16
17 sfs.forward = False
18
19 sfs.fit (x_no, y_no, custom_feature_names = labels)
20 print (sfs.k_score_)
21 print ('Sequential Backward Selection', sfs.k_feature_names_, end = '\n\n')
```

Listing 9: Aplicación métodos *wrapper* de selección características.

```
1 0.7054545454545454
2 Sequential Forward Selection ('leptin', 'bmi', 'glucose', 'MCP1')
3
4 0.7094949494949495
5 Sequential Backward Selection ('leptin', 'bmi', 'glucose', 'insulin')
```

Listing 10: Resultados del filtrado mediante wrappers.

9 PCA

```
1 from sklearn.preprocessing import StandardScaler
2 x_no = StandardScaler ().fit_transform (x_no) # typify
3
4 from sklearn.decomposition import PCA
5 pca = PCA (n_components = 9)
6 principalComponents = pca.fit_transform(x_no)
7 evr = pca.explained_variance_ratio_
```

Listing 11: *Principal Component Analysis*

9.1 Pareto

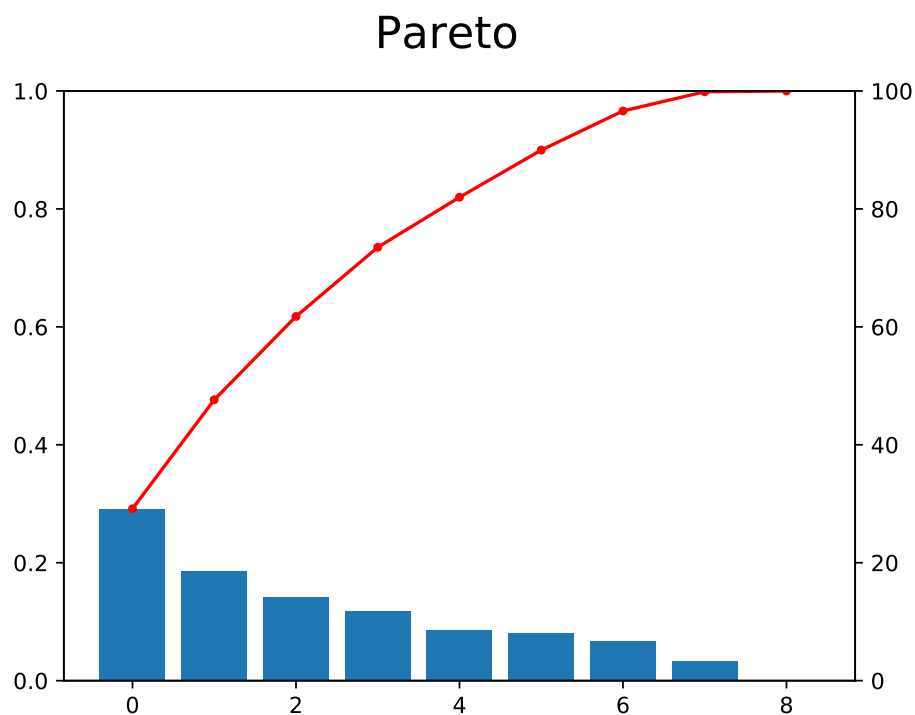


Fig. 7: Diagrama de Pareto.

```
1 fig, ax = plt.subplots ()
2 ax.bar (range (len (evr)), evr)
3 ax.set_ylim (top=1)
4
5 ax1 = ax.twinx ()
6 ax1.set_ylim (top=100)
7 ax1.plot (range (len (evr)), np.cumsum (evr)*100, marker = '.', color = 'red')
8
9 fig.suptitle ('Pareto', fontsize = 20)
10 fig.savefig ('../images/pareto.pdf', bbox_inches = 'tight', pad_inches = 0)
```

Listing 12: Código generador del diagrama de Pareto