

Sistemas de Información y Telemedicina. *

Marta Girones Sanguesa

Silvia Marset Gomis

Ignacio Amat Hernández

Sofía Gutiérrez Santamaría

December 2, 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
9.2 Biplot	13
10 Modelos de Clasificación	14
10.1 Clasificación Lineal	14
10.2 Clasificación Cuadrática	15

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

List of Figures

1	Histogramas para datos con y sin anomalias.	4
2	Kernel Density para datos con y sin anomalias.	5
3	Boxplots para datos con y sin anomalias.	6
4	QQplots para datos con y sin anomalias.	7
5	Corrplot para datos con anomalias.	8
6	Corrplot para datos sin anomalias.	9
7	Diagrama de Pareto.	12
8	Biplot.	13

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	Varianza explicada por componente y suma acumulada.	12
13	Código generador del diagrama de Pareto	12
14	Código generador del Biplot.	13
15	Validación del modelo lineal.	14
16	Validación según distintos métodos.	14
17	Validación del modelo cuadrático.	15
18	Validación según distintos métodos.	15

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.

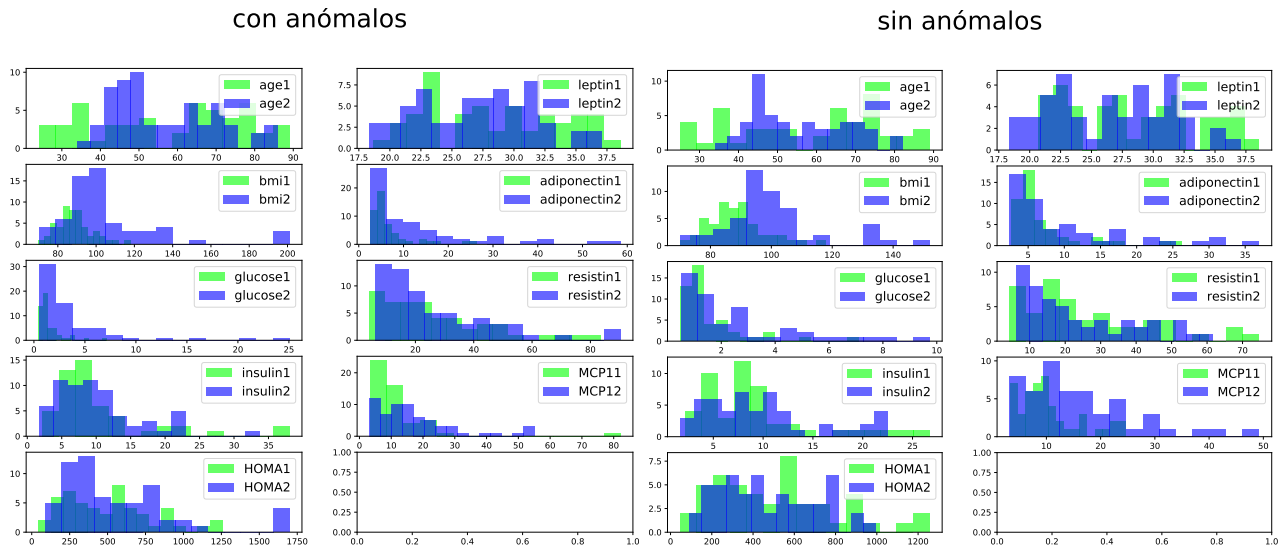


Fig. 1: Histogramas para datos con y sin anomalías.

```

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 anómalos', 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

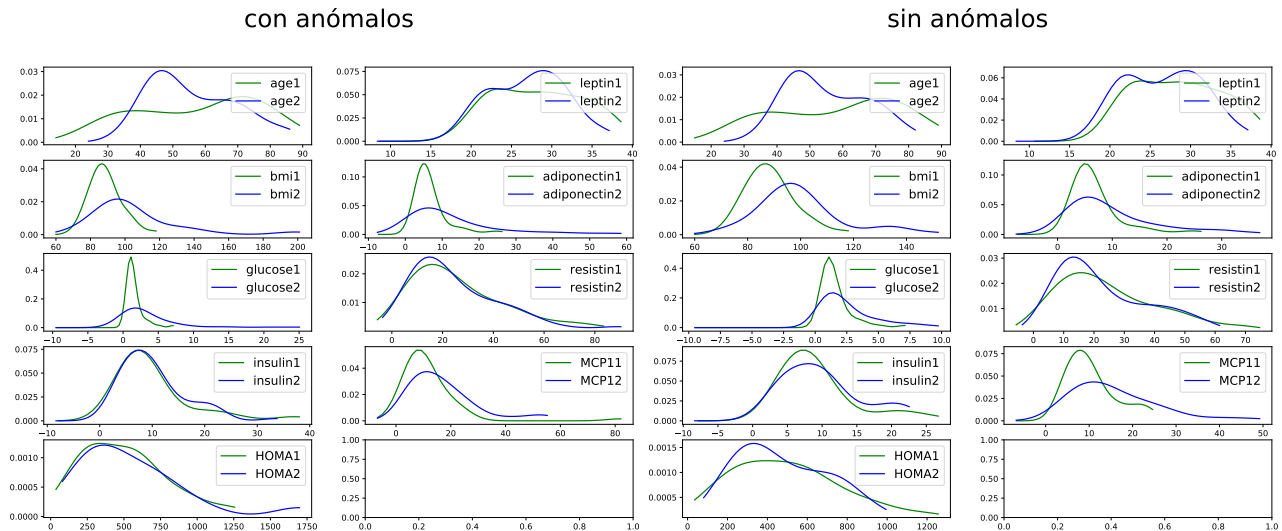


Fig. 2: Kernel Density para datos con y sin anomalías.

```

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

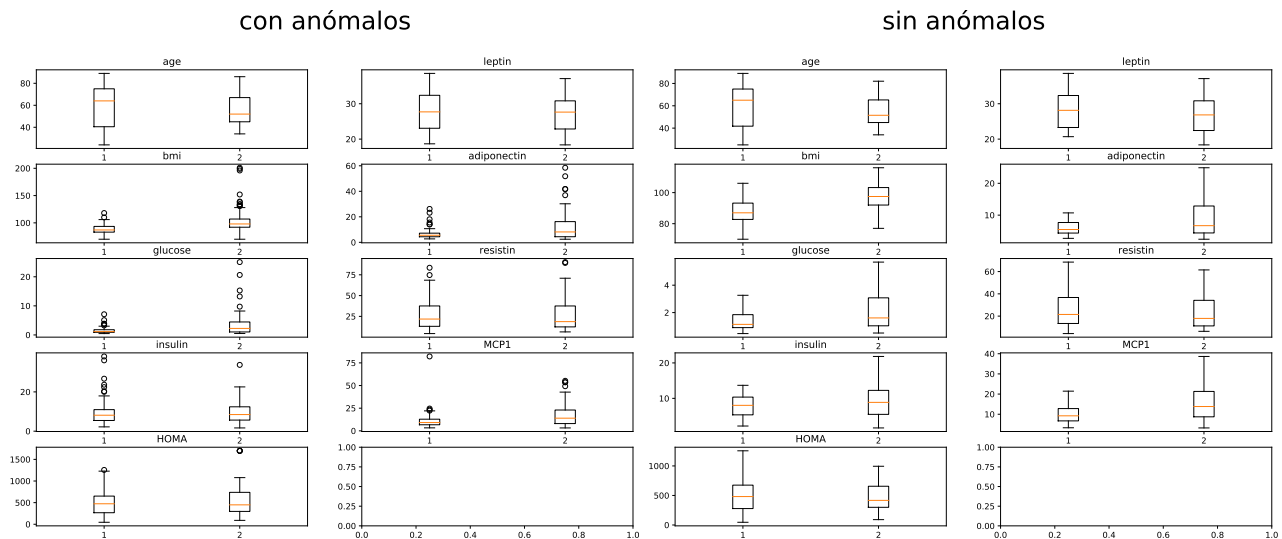


Fig. 3: Boxplots para datos con y sin anomalias.

```

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

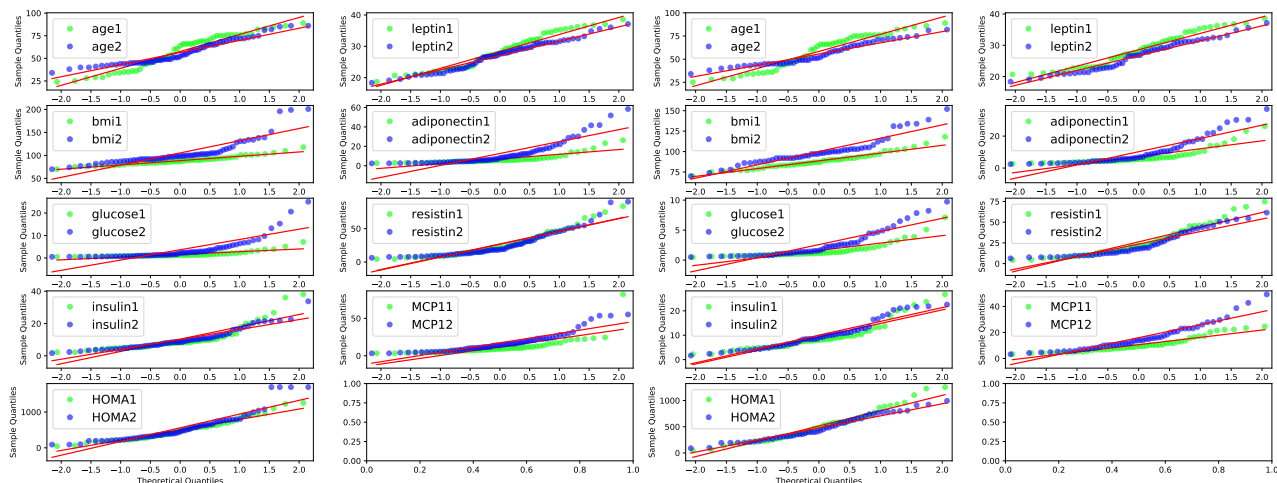


Fig. 4: QQplots para datos con y sin anomalías.

```

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

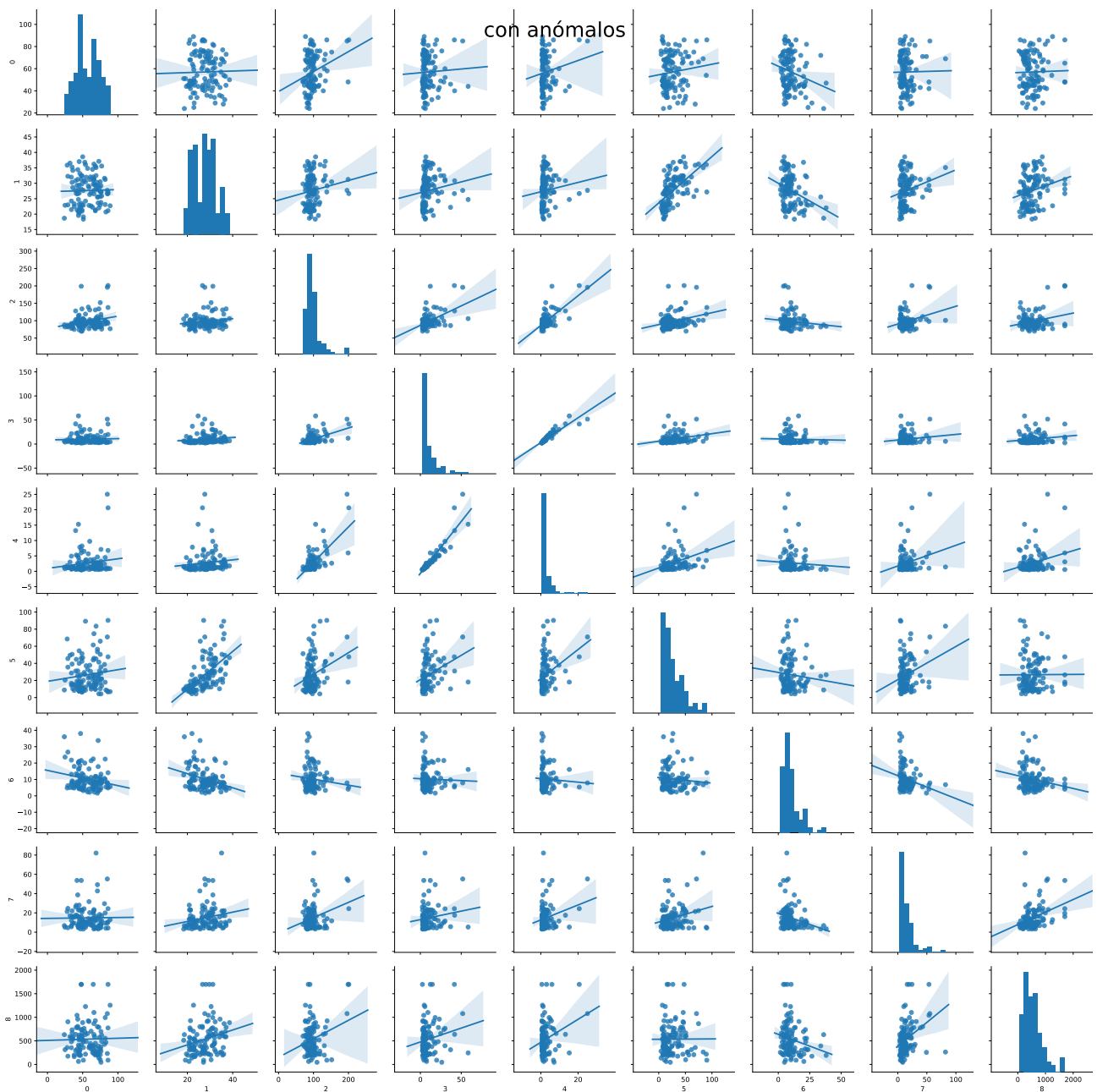


Fig. 5: Corrplot para datos con anomalías.

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

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

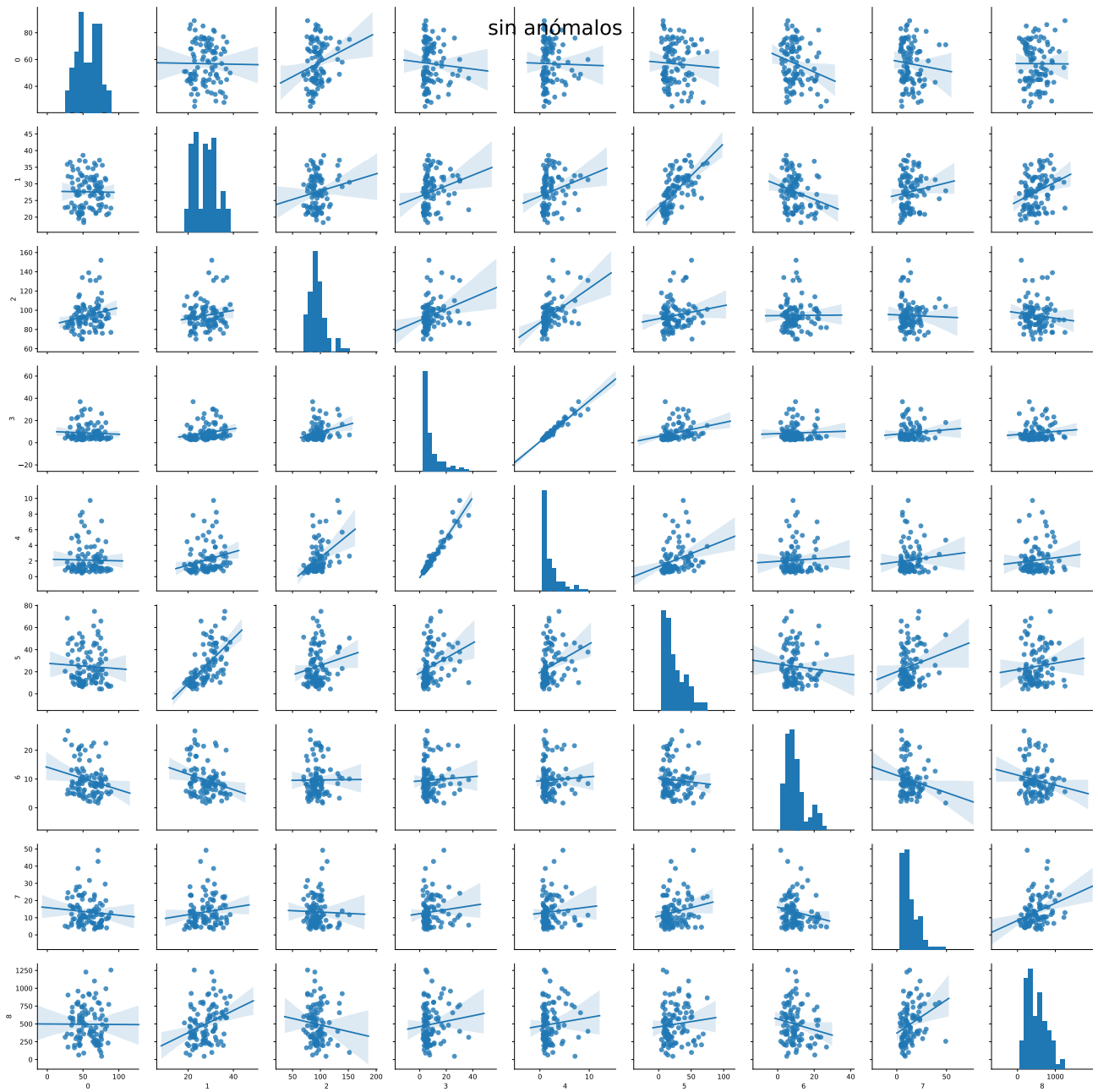


Fig. 6: Corrplot para datos sin anomalías.

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

Listing 11: *Principal Component Analysis*

```
1 [0.29146865 0.18490568 0.14125105 0.11727276 0.08486126 0.07999359
2  0.06636991 0.03254865 0.00132847]
3 [0.29146865 0.47637432 0.61762537 0.73489813 0.81975939 0.89975298
4  0.96612289 0.99867153 1.          ]
```

Listing 12: Varianza explicada por componente y suma acumulada.

9.1 Pareto

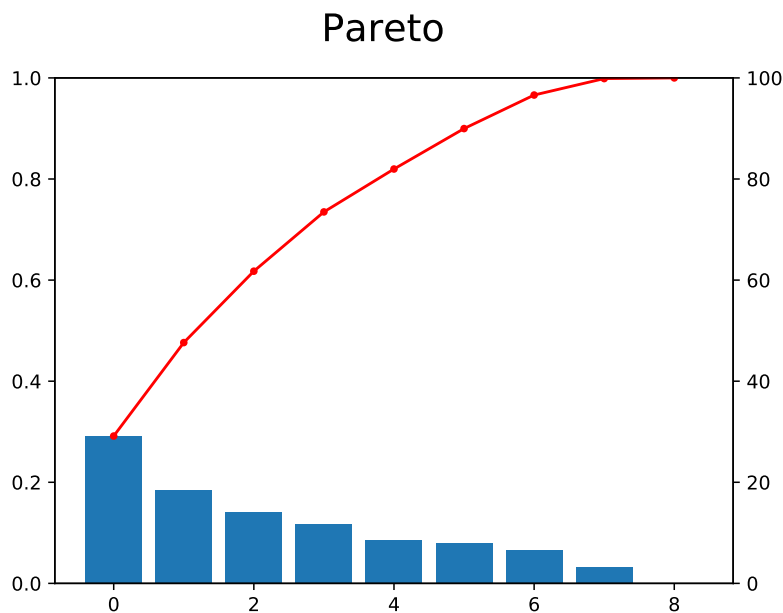


Fig. 7: Diagrama de Pareto.

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

Listing 13: Código generador del diagrama de Pareto

9.2 Biplot

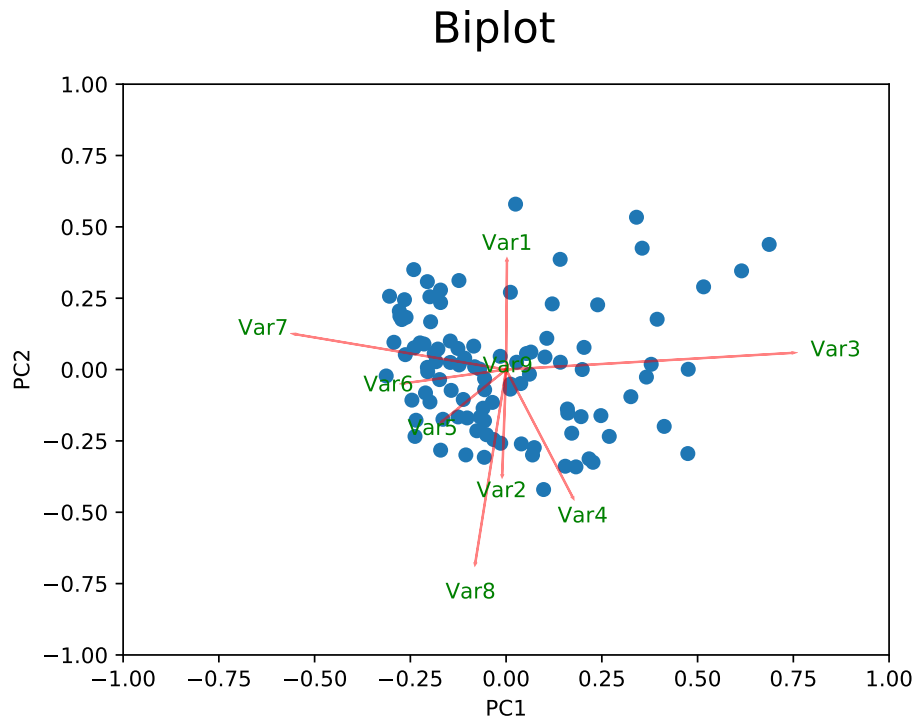


Fig. 8: Biplot.

```
1 def biplot(score,coeff,pcax,pcay,labels=None):
2     pca1=pcax-1; pca2=pcay-1
3     xs = score[:,pca1]; ys = score[:,pca2]
4     n=score.shape[1]
5     scalex = 1.0/(xs.max()- xs.min()); scaley = 1.0/(ys.max()- ys.min())
6     plt.scatter(xs*scalex,ys*scaley)
7     for i in range(n):
8         plt.arrow(0, 0, coeff[i,pca1], coeff[i,pca2],color='r',alpha=0.5)
9         if labels is None:
10             plt.text(coeff[i,pca1]* 1.15, coeff[i,pca2] * 1.15, "Var"+str(i+1), \
11                      color='g', ha='center', va='center')
12         else:
13             plt.text(coeff[i,pca1]* 1.15, coeff[i,pca2] * 1.15, labels[i], color='g', \
14                      , ha='center', va='center')
15     plt.xlim(-1,1); plt.ylim(-1,1)
16     plt.xlabel("PC{}".format(pcax)); plt.ylabel("PC{}".format(pcay))
17     return plt
18 bp = biplot (pca.fit_transform (x_no), pca.components_,1,2)
19 bp.suptitle ('Biplot', fontsize = 20)
20 bp.savefig ('../images/biplotpca.pdf', bbox_inches = 'tight', pad_inches = 0)
```

Listing 14: Código generador del Biplot.

10 Modelos de Clasificación

10.1 Clasificación Lineal

```
1 from sklearn.model_selection import cross_val_score, KFold, ShuffleSplit, \
    LeaveOneOut
2 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA
3
4 lda = LDA ()
5 error = 1 - cross_val_score (lda, x, y, cv = 10)
6 print ('Linear error CV media: %.2f std: %.2f'
7       %(np.mean (error), np.std (error)))
8
9 error = 1 - cross_val_score (lda, x, y, cv = KFold (n_splits = 10, shuffle = True))
10 print ('Linear error KF media: %.2f std: %.2f'
11       %(np.mean (error), np.std (error)))
12
13 error = 1 - cross_val_score (lda, x, y, cv = ShuffleSplit (n_splits = 10))
14 print ('Linear error SS media: %.2f std: %.2f'
15       %(np.mean (error), np.std (error)))
16
17 error = 1 - cross_val_score (lda, x, y, cv = LeaveOneOut ())
18 print ('Linear error LO media: %.2f std: %.2f'
19       %(np.mean (error), np.std (error)))
```

Listing 15: Validación del modelo lineal.

```
1 Linear error CV media: 0.25 std: 0.13
2 Linear error KF media: 0.23 std: 0.15
3 Linear error SS media: 0.24 std: 0.10
4 Linear error LO media: 0.24 std: 0.43
```

Listing 16: Validación según distintos métodos.

10.2 Clasificación Cuadrática

```
1 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis as QDA
2
3 qda = QDA ()
4 error = 1 - cross_val_score (qda, x, y, cv = 10)
5 print ('Quadratic error CV media: %.2f std: %.2f'
6        %(np.mean (error), np.std (error)))
7
8 error = 1 - cross_val_score (qda, x, y, cv = KFold (n_splits = 10, shuffle = True))
9 print ('Quadratic error KF media: %.2f std: %.2f'
10        %(np.mean (error), np.std (error)))
11
12 error = 1 - cross_val_score (qda, x, y, cv = ShuffleSplit (n_splits = 10))
13 print ('Quadratic error SS media: %.2f std: %.2f'
14        %(np.mean (error), np.std (error)))
15
16 error = (1 - cross_val_score (qda, x, y, cv = LeaveOneOut ()))
17 print ('Quadratic error LO media: %.2f std: %.2f'
18        %(np.mean (error), np.std (error)))
```

Listing 17: Validación del modelo cuadrático.

```
1 Quadratic error CV media: 0.34 std: 0.19
2 Quadratic error KF media: 0.28 std: 0.13
3 Quadratic error SS media: 0.22 std: 0.15
4 Quadratic error LO media: 0.27 std: 0.44
```

Listing 18: Validación según distintos métodos.