

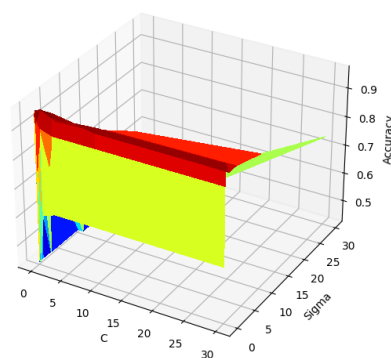
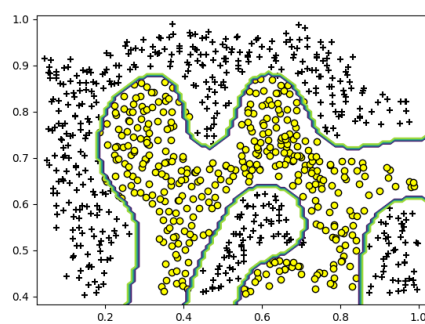
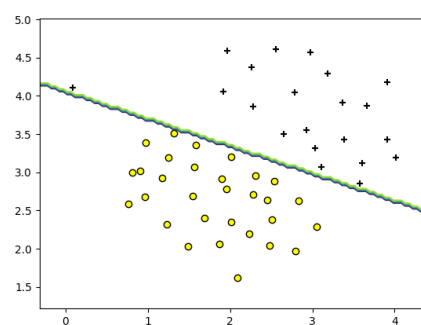
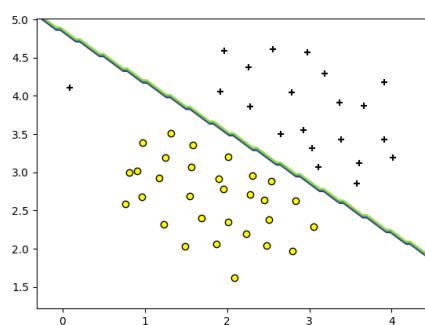
# Práctica 5: Regresión lineal regularizada: sesgo y varianza

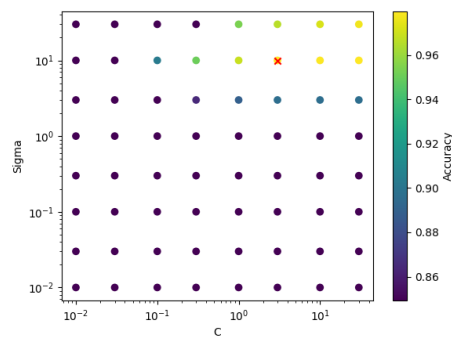
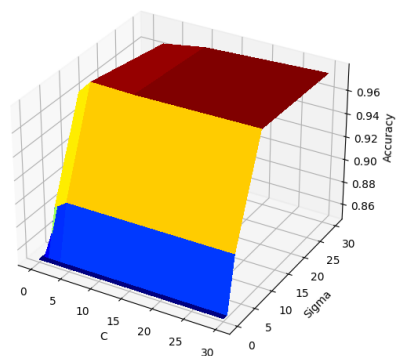
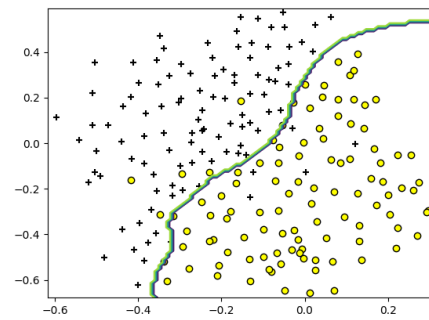
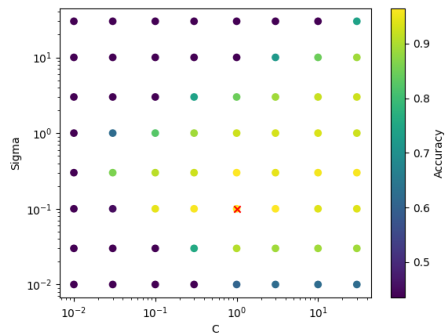
Ana Martín Sánchez, Nicolás Pastore Burgos

21/09/2021

## 1 Solución propuesta

### 1.1 Resultados obtenidos





## 1.2 Implementación

```

1000 import numpy as np
1001 import matplotlib.pyplot as plt
1002 import sklearn.svm as svm
1003 import codecs
1004 import os, os.path
1005
1006 from matplotlib import cm
1007 from mpl_toolkits.mplot3d import Axes3D
1008 from scipy.io import loadmat
1009 from sklearn.metrics import accuracy_score
1010
1011 from Data.Parte2.process_email import email2TokenList
1012 from Data.Parte2.get_vocab_dict import getVocabDict
1013
1014 def scatter_mat(x,y):
1015     pos = (y == 1).ravel()
1016     neg = (y == 0).ravel()
1017     plt.figure()
1018     plt.scatter(x[pos, 0], x[pos, 1], color='black', marker='+')
1019     plt.scatter(

```

```

1020     x[neg, 0], x[neg, 1], color='yellow', edgecolors='black', marker='o')
1021
1022 def visualize_boundary(x, y, svm, zoom):
1023     x1 = np.linspace(x[:, 0].min() - zoom, x[:, 0].max() + zoom, 100)
1024     x2 = np.linspace(x[:, 1].min() - zoom, x[:, 1].max() + zoom, 100)
1025
1026     x1, x2 = np.meshgrid(x1, x2)
1027     yp = svm.predict(np.array([x1.ravel(), x2.ravel()]).T).reshape(x1.shape)
1028
1029     scatter_mat(x,y)
1030
1031     plt.contour(x1, x2, yp)
1032     plt.show()
1033
1034 def partel(x,y):
1035     s1 = svm.SVC(kernel='linear', C=1.0)
1036     s1.fit(x,y)
1037
1038     s2 = svm.SVC(kernel='linear', C=100.0)
1039     s2.fit(x,y)
1040
1041     zoom = 0.4
1042     visualize_boundary(x,y,s1,zoom)
1043     visualize_boundary(x,y,s2,zoom)
1044
1045 def partel_2(x, y, C, sigma):
1046     s = svm.SVC(kernel='rbf', C=C, gamma= 1 / (2 * sigma ** 2))
1047     s.fit(x,y)
1048
1049     visualize_boundary(x,y,s, 0.02)
1050     plt.show()
1051
1052 def searchBestCandSigma(x, xVal, y, yVal):
1053     cs = np.array([0.01, 0.03, 0.1, 0.3, 1, 3, 10, 30])
1054     sigmas = np.array([0.01, 0.03, 0.1, 0.3, 1, 3, 10, 30])
1055
1056     numCs = cs.shape[0]
1057     numSigmas = cs.shape[0]
1058
1059     res = np.zeros(numCs * numSigmas).reshape(numCs, numSigmas)
1060
1061     bestAcc = 0
1062     bestC = -1
1063     bestSigma = -1
1064     bestSVM = {}
1065
1066     for i in np.arange(numCs):
1067         for j in np.arange(numSigmas):
1068             s = svm.SVC(kernel='rbf', C=cs[i], gamma= 1 / (2 * sigmas[j] **
1069                 2))
1070             s.fit(x,y)

```

```

1071         res[i,j] = accuracy_score(yVal, s.predict(xVal))
1072         if bestAcc < res[i,j]:
1073             bestAcc = res[i,j]
1074             bestC = cs[i]
1075             bestSigma = sigmas[j]
1076             bestSVM = s
1077         print('Tested sigma ' + str(sigmas[j]))
1078
1079     print('Tested c ' + str(cs[i]))
1080
1081     print("Best accuracy: " + str(bestAcc))
1082     print("C: " + str(bestC))
1083     print("Sigma: " + str(bestSigma))
1084
1085     fig = plt.figure()
1086
1087     cc, ss = np.meshgrid(sigmas, cs)
1088
1089     ax = Axes3D(fig, auto_add_to_figure=False)
1090     ax.set_xlabel('C')
1091     ax.set_ylabel('Sigma')
1092     ax.set_zlabel('Accuracy')
1093
1094     fig.add_axes(ax)
1095
1096     surf = ax.plot_surface(ss, cc, res, cmap=cm.jet, linewidth=0, antialiased
1097                             =False)
1098     plt.show()
1099
1100     plt.scatter(ss, cc, c=res)
1101     plt.xscale('log')
1102     plt.yscale('log')
1103     plt.xlabel('C')
1104     plt.ylabel('Sigma')
1105     plt.clim(np.min(res), np.max(res))
1106     plt.colorbar().set_label('Accuracy')
1107     plt.scatter(bestC, bestSigma, marker='x', color='r')
1108     plt.show()
1109
1110     return bestSVM, bestAcc, bestC, bestSigma
1111
1112 def partel_3(x, xVal, y, yVal):
1113     s, acc, c, sigma = searchBestCandSigma(x, xVal, y, yVal)
1114
1115     visualize_boundary(x,y,s, 0.02)
1116     plt.show()
1117
1118     """
1119
1120
1121
1122     """

```

```

1123
1124 def filterEmailWithDictionary(dict, email):
1125     emailDict = np.zeros([len(dict)])
1126     for word in email:
1127         if word in dict:
1128             emailDict[dict[word]-1] = 1
1129
1130     return emailDict
1131
1132 def readEmails(dict, route, folder, format, size):
1133     emails = np.empty((size, len(dict)))
1134     for i in np.arange(size):
1135         email_contents = codecs.open(route + folder + '{0:04d}'.format(i+1)
1136         + format, 'r', encoding='utf 8', errors='ignore').read()
1137         emails[i] = filterEmailWithDictionary(dict, email2TokenList(
1138         email_contents))
1139     print('Done reading and preparing ' + folder)
1140     return emails
1141
1142 def createXandY(spam, easyHam, hardHam, ySpam, yEasy, yHard):
1143     return np.append(np.append(spam, easyHam, axis=0), hardHam, axis=0), np
1144     .append(np.append(ySpam, yEasy, axis=0), yHard, axis=0)
1145
1146 def parte2():
1147     dict = getVocabDict()
1148
1149     route = 'Data/Parte2/'
1150     folders = ['spam/', 'easy_ham/', 'hard_ham/']
1151     format = '.txt'
1152
1153     print('Comencing to read data')
1154
1155     spam = readEmails(dict, route, folders[0], format, len([name for name
1156     in os.listdir('./' + route + folders[0])]))
1157     spamSize = len(spam)
1158     ySpam = np.ones(spamSize)
1159
1160     easy = readEmails(dict, route, folders[1], format, len([name for name
1161     in os.listdir('./' + route + folders[1])]))
1162     easySize = len(easy)
1163     yEasy = np.zeros(easySize)
1164
1165     hard = readEmails(dict, route, folders[2], format, len([name for name
1166     in os.listdir('./' + route + folders[2])]))
1167     hardSize = len(hard)
1168     yHard = np.zeros(hardSize)
1169
1170     trainPerc = 0.6
1171     valPerc = 0.3 + trainPerc
1172     # rest for test
1173
1174     trainSpam = int(trainPerc * spamSize)
1175     trainEasy = int(trainPerc * easySize)

```

```

1170     trainHard = int(trainPerc * hardSize)
1171
1172     valSpam = int(valPerc * spamSize)
1173     valEasy = int(valPerc * easySize)
1174     valHard = int(valPerc * hardSize)
1175
1176     x, y = createXandY(spam[:trainSpam],
1177                        easy[:trainEasy],
1178                        hard[:trainHard],
1179                        ySpam[:trainSpam],
1180                        yEasy[:trainEasy],
1181                        yHard[:trainHard])
1182
1183     xVal, yVal = createXandY(spam[trainSpam:valSpam],
1184                              easy[trainEasy:valEasy],
1185                              hard[trainHard:valHard],
1186                              ySpam[trainSpam:valSpam],
1187                              yEasy[trainEasy:valEasy],
1188                              yHard[trainHard:valHard])
1189
1190     xTest, yTest = createXandY(spam[valSpam:],
1191                                easy[valEasy:],
1192                                hard[valHard:],
1193                                ySpam[valSpam:],
1194                                yEasy[valEasy:],
1195                                yHard[valHard:])
1196
1197     s, acc, c, sigma = searchBestCandSigma(x, xVal, y, yVal)
1198
1199     print('Accuracy over Test sample: ' + str(accuracy_score(yTest, s.
1200 predict(xTest)) * 100) + '%')
1201
1202 def main():
1203     data1 = loadmat("Data/Partel/ex6data1.mat")
1204
1205     x1 = data1['X']
1206     y1 = data1['y']
1207     y1R = np.ravel(y1)
1208
1209     data2 = loadmat("Data/Partel/ex6data2.mat")
1210
1211     x2 = data2['X']
1212     y2 = data2['y']
1213     y2R = np.ravel(y2)
1214
1215     data3 = loadmat("Data/Partel/ex6data3.mat")
1216
1217     x3 = data3['X']
1218     y3 = data3['y']
1219     y3R = np.ravel(y3)
1220
1221     x3Val = data3['Xval']
1222     y3Val = data3['yval']

```

```
1222     y3ValR = np.ravel(y3Val)
1223
1224     parte1(x1, y1R)
1225     parte1_2(x2, y2R, 1.0, 0.1)
1226     parte1_3(x3, x3Val, y3R, y3ValR)
1227     parte2()
1228
1229 if __name__ == "__main__":
1230     main()
```

main.py