

Práctica 1: Regresión Lineal

Ana Martín Sánchez, Nicolás Pastore Burgos

21/09/2021

1 Descripción de la práctica

En esta práctica, se pedía aplicar el método de regresión lineal sobre dos conjuntos de datos.

En primer lugar, se nos da un conjunto de datos que relaciona dos variables: la población de una ciudad con los beneficios de una compañía de distribución de comida en dicha ciudad. Después, se nos presenta un conjunto de datos que relaciona el precio de casas vendidas en Portland con el tamaño en pies cuadrados, el número de habitaciones y el precio de dicha casa.

Para el primer conjunto de datos, se puede aplicar el método de regresión lineal con una variable; es evidente que, en el segundo caso, no es así.

Para agilizar esta segunda parte, es necesario normalizar los atributos, antes de aplicar un descenso de gradiente.

Por último, aplicaremos el método de la ecuación normal sobre el segundo conjunto de datos, con el fin de comprobar que nuestros cálculos anteriores son correctos.

2 Solución propuesta

2.1 Resultados obtenidos

2.1.1 Parte 1

En esta parte hemos conseguido los resultados esperados; hemos podido implementar una función de descenso de gradiente que es capaz, dados unos datos, de calcular una ecuación lineal que minimice los costes. Para comprobarlo, hemos utilizado el ejemplo propuesto en el enunciado (para una población de 70.000 habitantes, los beneficios estimados son de 45.282\$).

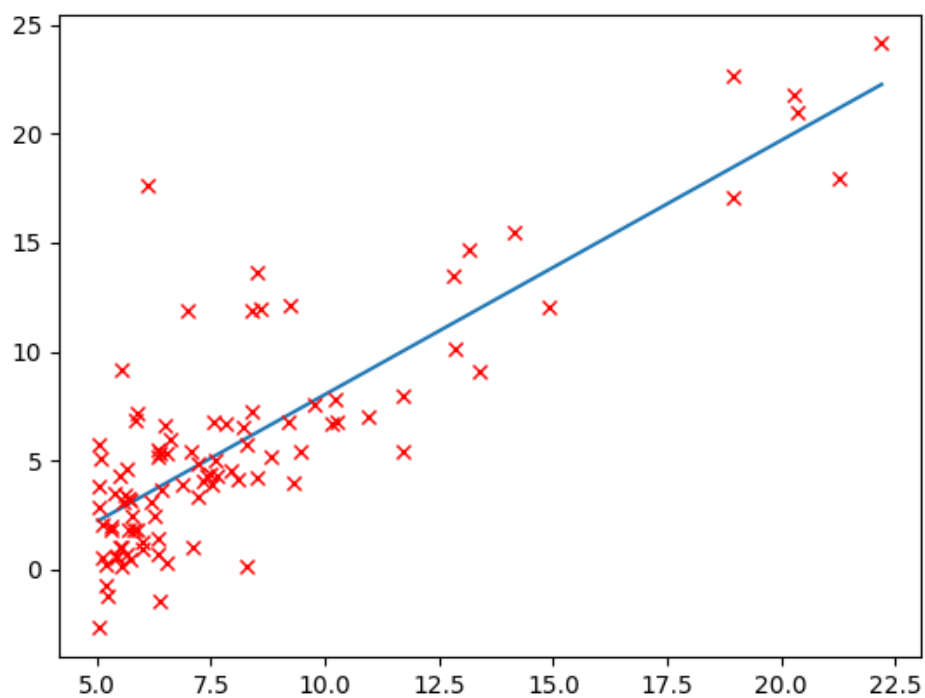


Figure 1: Gráfica que muestra la recta calculada para los datos proporcionados por el enunciado.

Además, hemos generado las gráficas que muestran los costes, como se muestran a continuación.

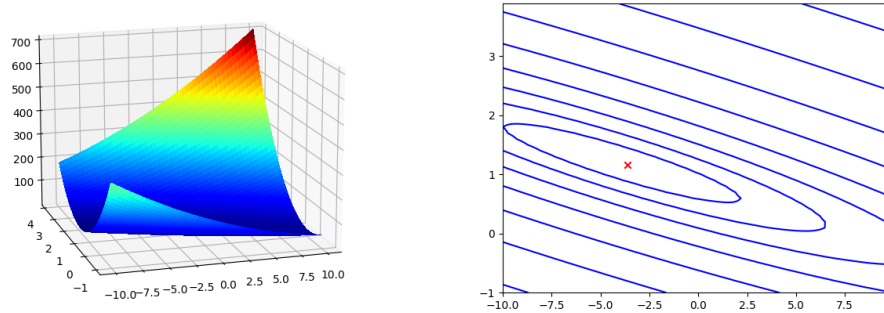


Figure 2: Representación de una sección de la ecuación que siguen las thetas.

2.1.2 Parte 2

En esta segunda parte, también hemos implementado un descenso de gradiente, pero con un número n de atributos. Además hemos implementado la ecuación normal.

Para comprobar que ambas implementaciones eran correctas, hemos utilizado el ejemplo propuesto en el enunciado (la entrada $[1650, 3]$), y ambas nos daban el mismo resultado (293098.46), con una diferencia de menos del 0.01% entre ambos resultados.

Como comprobación adicional, hemos observado que la función de costes devuelve un resultado cada vez menor, como se puede observar a continuación.

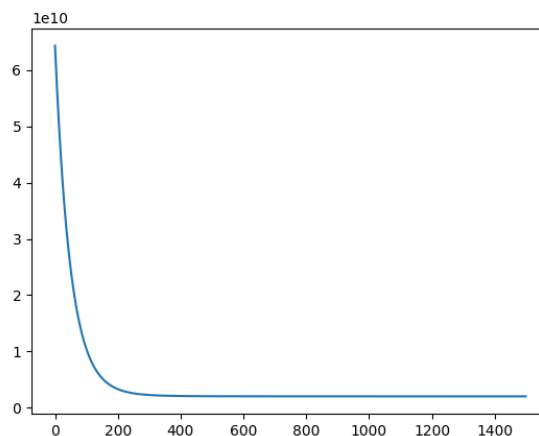


Figure 3: La función de costes es cada vez menor.

2.2 Implementación

```
1000 import numpy as np
1001 import matplotlib.pyplot as plt
1002 from numpy.core.fromnumeric import shape
1003 from numpy.core.numeric import ones_like, zeros_like
1004 from pandas.io.parsers import read_csv
1005 from mpl_toolkits.mplot3d import Axes3D, axes3d
1006 from matplotlib import cm, colors
1007 from matplotlib.ticker import LinearLocator, FormatStrFormatter
1008
1009 def loadCSV(fileName):
1010     return read_csv(fileName, header=None).to_numpy().astype(float)
1011
1012 def make_data(t0_range, t1_range, x, y):
1013     step = 0.1
1014     Theta0 = np.arange(t0_range[0], t0_range[1], step)
1015     Theta1 = np.arange(t1_range[0], t1_range[1], step)
1016
1017     Theta0, Theta1 = np.meshgrid(Theta0, Theta1)
1018
1019     Coste = np.empty_like(Theta0)
1020     for ix, iy in np.ndindex(Theta0.shape):
1021         Coste[ix, iy] = coste(Theta0[ix, iy], Theta1[ix, iy], x, y, len(x))
1022
1023     return [Theta0, Theta1, Coste]
1024
1025 def paint(x, y, t0, t1, t0Mat, t1Mat, costeMat):
1026     fig = plt.figure()
1027
1028     # to avoid a warning we do...
1029     ax = Axes3D(fig, auto_add_to_figure=False)
1030     fig.add_axes(ax)
1031
1032     surf = ax.plot_surface(t0Mat, t1Mat, costeMat, cmap=cm.jet, linewidth=0,
1033                             antialiased=False)
1034     plt.show()
1035
1036     plt.scatter(t0, t1, marker='x', color='red')
1037     plt.contour(t0Mat, t1Mat, costeMat, np.logspace(-2, 3, 20), colors='blue')
1038     plt.show()
1039
1040     minX = min(x)
1041     maxX = max(x)
1042     minY = t0 + t1 * minX
1043     maxY = t0 + t1 * maxX
1044     plt.plot([minX, maxX], [minY, maxY])
1045     plt.plot(x, y, "x", color='red')
1046     plt.show()
1047
1048     print([t0, t1])
1049
1050 def minimizeCost(x, y):
```

```

1050
1051 m = len(x)
1052 alpha = 0.01
1053
1054 theta0 = theta1 = 0
1055
1056 for _ in range(1500):
1057     """
1058     Esto visualiza la progresion de la pendiente gradualmente en
1059     funcion de los errores
1060     """
1061     # minX = min(x)
1062     # maxX = max(x)
1063     # minY = theta0 + theta1 * minX
1064     # maxY = theta0 + theta1 * maxX
1065     # plt.plot([minX, maxX], [minY, maxY], "--", linewidth=0.5)
1066
1067     """
1068     Ajusta una vez la theta0 y theta1 segun el error calculado
1069     """
1070     sum0 = np.sum((x * theta1 + theta0) - y)
1071     sum1 = np.sum(((x * theta1 + theta0) - y) * x)
1072
1073     theta0 = theta0 - (alpha/m) * sum0
1074     theta1 = theta1 - (alpha/m) * sum1
1075
1076     # 70.000 habs = 4.53...
1077     # print(theta0 + theta1*7)
1078
1079     return [theta0, theta1]
1080
1081 def coste (theta0, theta1, x, y, m):
1082     return np.sum(((theta0 + theta1 * x) - y) ** 2) / (2 * m)
1083
1084 def parte1():
1085     data = loadCSV("Data/ex1data1.csv")
1086
1087     x = data[:, 0]
1088     y = data[:, 1]
1089
1090     t0, t1 = minimizeCost(x, y)
1091
1092     t0Mat, t1Mat, costeMat = make_data([-10,10], [-1,4],x, y)
1093     paint(x, y, t0, t1, t0Mat, t1Mat, costeMat)
1094
1095     """
1096     =====
1097     =====
1098     =====
1099     """
1100
1101 def normalizaMat(mat):

```

```

1102 mu = np.array(np.mean(mat, axis=0))
1103 sigma = np.array(np.std(mat, axis=0))
1104
1105 # res = (mat - mu) / sigma
1106 # np.nan_to_num(res, False, 1.0)
1107
1108 # [1, 1541.3, 4]
1109 matNorm = ones_like(mat)
1110 # [0, 1]
1111 for i in np.arange(np.shape(mat)[1]-1):
1112     # [1, 2]
1113     matNorm[:, i+1] = (mat[:, i+1] - mu[i+1]) / sigma[i+1]
1114
1115     return matNorm, mu, sigma
1116
1117 def costeVec(x, y, thetas):
1118     xTh = np.dot(x, thetas)
1119     return np.sum((xTh - y)**2) / (2*len(x))
1120
1121 def descensoGradiente(x, y, alpha):
1122     m = np.shape(x)[0]
1123     n = np.shape(x)[1]
1124
1125     #thetas = np.zeros(n)
1126     thetas2 = np.zeros(n)
1127
1128     costes = np.zeros(1500)
1129
1130     for i in range(len(costes)):
1131         xTh = np.dot(x, thetas2)
1132
1133         NuevaTheta = thetas2
1134         Aux = xTh - y
1135         for j in range(n):
1136             Aux_j = Aux * x[:, j]
1137             NuevaTheta[j] -= (alpha / m) * Aux_j.sum()
1138         costes[i] = costeVec(x, y, thetas2)
1139
1140         thetas2 = NuevaTheta
1141
1142     # for i in range(len(costes)):
1143     #     xTh = np.dot(x, thetas)
1144
1145     #     temp = np.dot(np.transpose(x), (xTh - y))
1146     #     thetas = thetas - (alpha/m) * temp
1147     #     costes[i] = costeVec(x, y, thetas)
1148
1149     plt.plot(np.arange(len(costes)), costes)
1150
1151     plt.show()
1152
1153     return thetas2
1154

```

```

1155 def ecuacionNormal(x,y):
1156     return np.dot(np.linalg.pinv(np.dot(np.transpose(x), x)), np.dot(np.
transpose(x), y))
1157
1158 def parte2():
1159     data = loadCSV("Data/ex1data2.csv")
1160
1161     x = data[:, :-1]
1162     y = data[:, -1]
1163
1164     m = np.shape(x)[0]
1165     n = np.shape(x)[1]
1166
1167     xNew = np.hstack([np.ones([m, 1]), x])
1168
1169     xNorm, mu, sigma = normalizaMat(xNew)
1170
1171     alpha = 0.01
1172     thetasDG = descensoGradiente(xNorm, y, alpha)
1173     thetasEN = ecuacionNormal(xNew,y)
1174
1175     example = [1.0, 1650.0, 3.0]
1176     exampleNorm = ones_like(example)
1177
1178     # for i in np.arange(len(example)-1):
1179     #     exampleNorm[i+1] = (example[i+1] - mu[i+1]) / sigma[i+1]
1180
1181     exampleNorm[1:] = (example[1:] - mu[1:]) / sigma[1:]
1182
1183
1184 if __name__ == "__main__":
1185     parte1()
1186     parte2()

```

main.py