

Aprendizaje Automático y Minería de Datos: Proyecto Final

Ana Martín Sánchez, Nicolás Pastore Burgos

[Repositorio en GitHub](#)

21/09/2021

1 Propuesta de proyecto

Nuestro proyecto consiste en clasificar setas como comestibles o venenosas, dependiendo de un total de 20 atributos, como los siguientes:

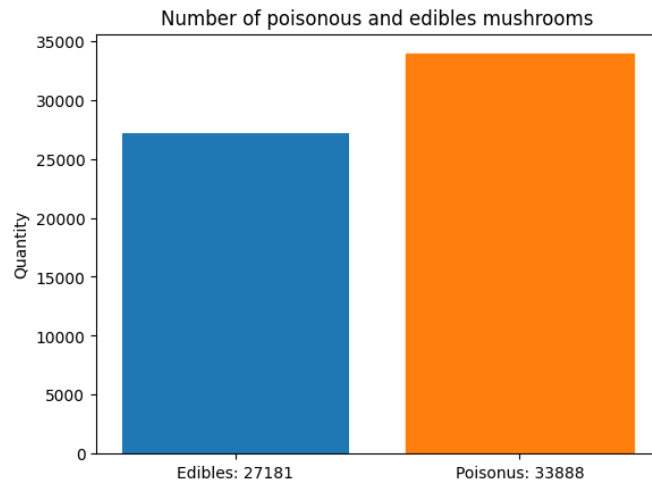
- Diámetro del sombrero: un float que representa el diámetro, en cm.
- Forma del sombrero: un char que representa una de las posibles formas.
- Superficie del sombrero: un char que representa el adjetivo que mejor describe la superficie del sombrero.

Para poder analizar los datos, desplegamos algunos de estos atributos (los que eran de tipo enumerado) a matrices de booleanos. Después de esta operación, cada entrada de la base de datos tenía 122 atributos.

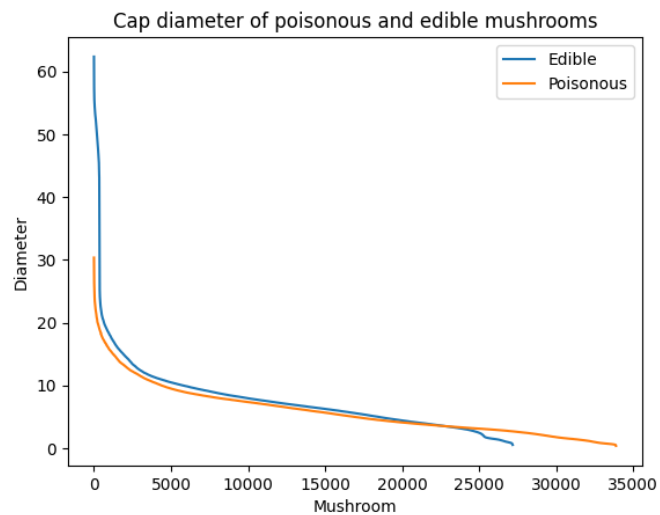
La base de datos tiene un total de 61069 entradas. Por este motivo, no hemos utilizado todos los datos para el proyecto; para las pruebas, de manera general, hemos escogido un 20% de los datos para entrenar los sistemas, otro 20% para validarlos y otro 20% para hacer la prueba final.

La base de datos original, extraída de la plataforma [Kaggle](#), tenía los datos organizados de manera que escoger una muestra en el orden establecido no era útil. Por esta razón, decidimos utilizar los datos según estaban organizados en [este repositorio](#), ya que estaban mezclados.

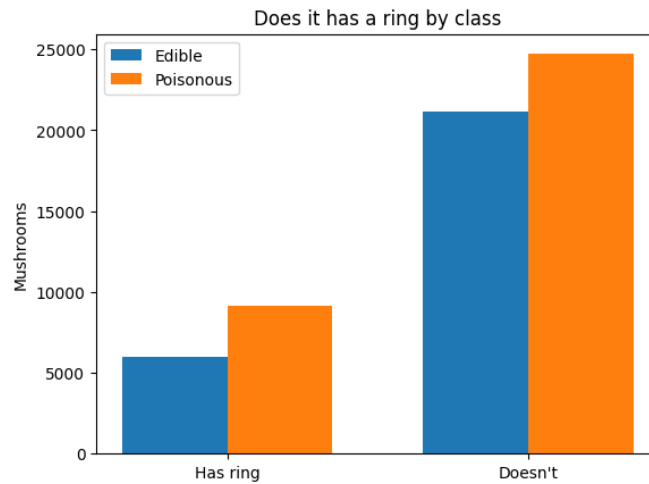
Para entender mejor los datos, presentamos una serie de gráficos que permiten hacer comparaciones entre las setas comestibles y las venenosas:



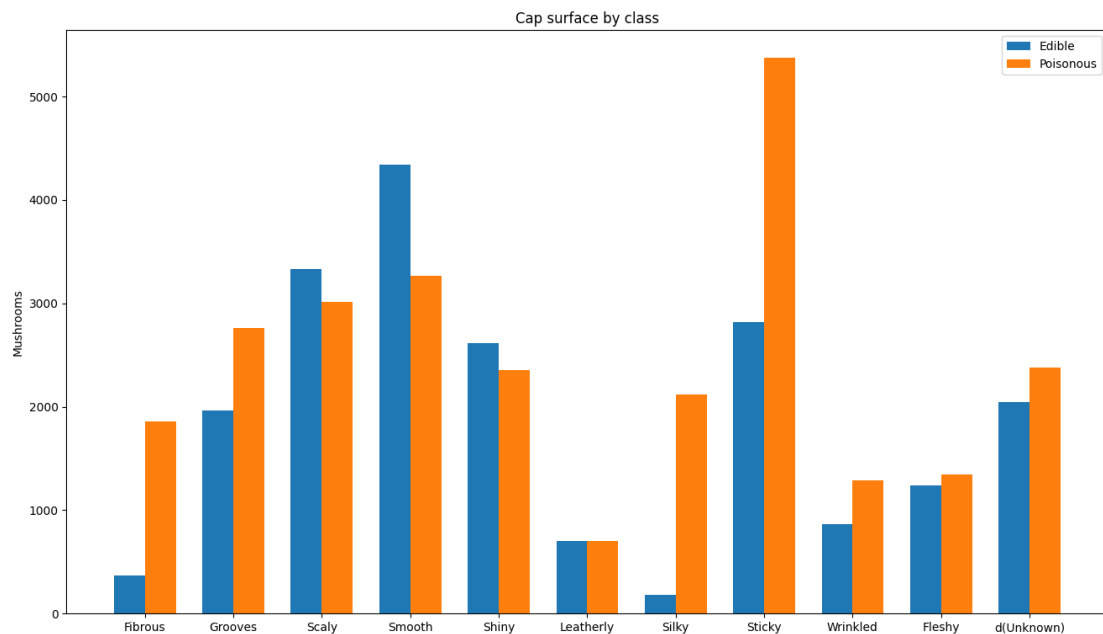
Como se puede observar, la cantidad de muestras de setas venenosas supera en número a la de setas comestibles, en más de 6000 entradas.



Nos pareció interesante mostrar cómo, cuando el diámetro del sombrero de las setas es muy grande, es muy probable que las setas sean comestibles. Sin embargo, en tamaños intermedios, resulta más difícil diferenciar las setas comestibles de las venenosas.



En este caso, vemos que la mayoría de las setas no tienen anillo. Hemos escogido esta gráfica como muestra de algunos de los valores binarios que se van a analizar más adelante.



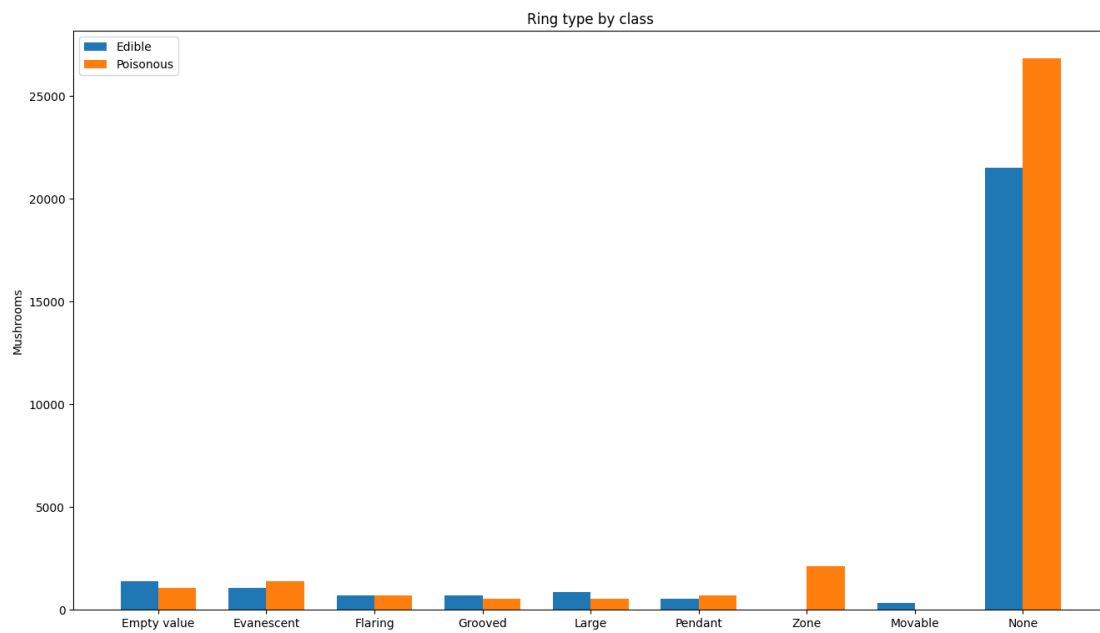
Por último, se muestran una serie de adjetivos que describen más fielmente el tipo de superficie del sombrero, y cuáles son venenosas y comestibles. Como se puede ver, y como comentamos a continuación, muchas de las entradas tenían un valor "d" para esta característica (un valor desconocido).

1.1 Problemas encontrados

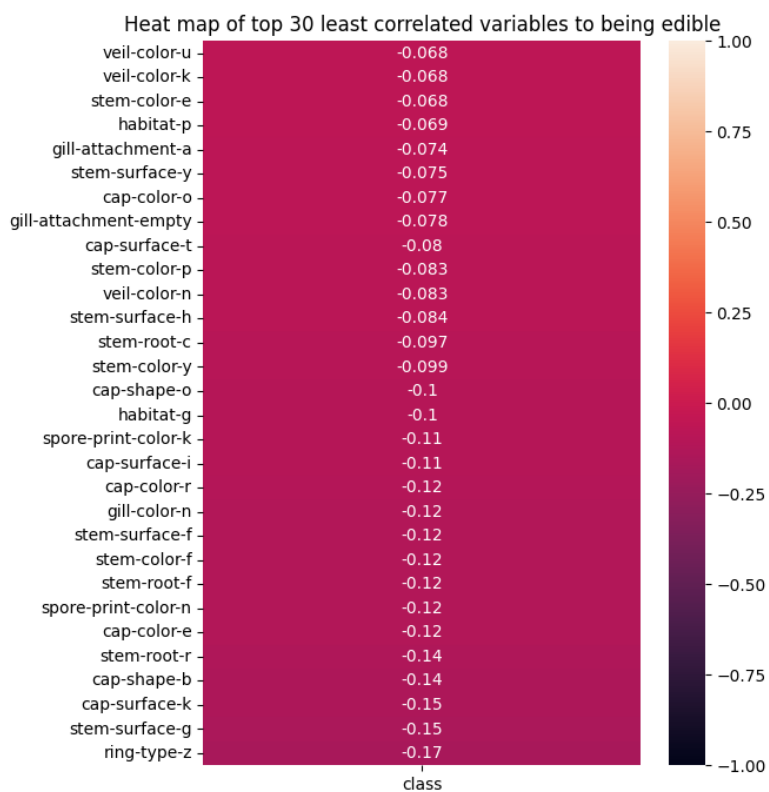
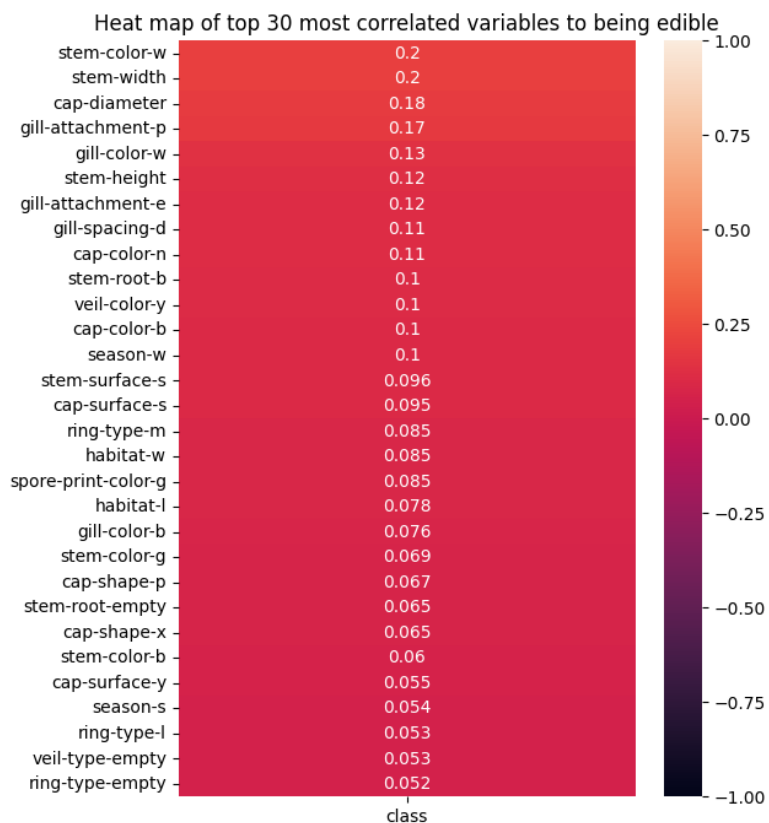
Al analizar los datos, nos dimos cuenta de que la base de datos nos supondría algunos problemas para realizar el proyecto:

En primer lugar, hay columnas en las que todos los valores son 0. También se dio que, en la gran mayoría de filas, alguna de las columnas no tenía un valor.

Por otra parte, algunos de los datos eran contradictorios. En concreto, algunas de las entradas tenían marcado como "t" la columna de has-ring (lo que significa que, efectivamente, tienen un anillo); pero, posteriormente, en la columna de ring-type, tenían marcado "f" (que se corresponde con el tipo "ninguno"). Esto no supone un problema a la hora de implementar los sistemas de aprendizaje automático, pero sí son un problema desde el punto de vista semántico, y nos hacen dudar de la validez de los datos.



Por último, encontramos un problema con la correlación entre los datos y la salida. Como se puede observar en las siguientes gráficas, ninguna de las entradas tiene una correlación que supere el 0.2, ni en el eje positivo ni en el negativo.



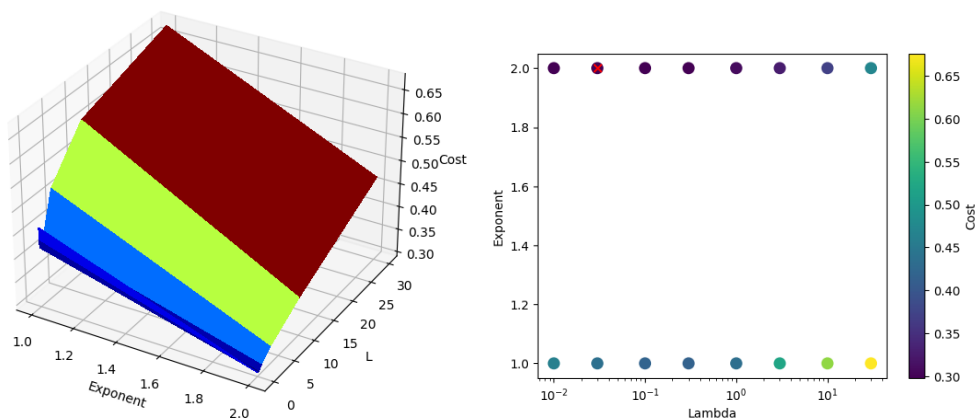
2 Resultados obtenidos

Para obtener los mejores resultados posibles, hemos utilizado varios métodos de los estudiados en clase, a fin de poder compararlos.

2.1 Regresión Logística

La regresión logística se puede entender como un caso especial de la regresión lineal, en la que la variable dependiente puede tomar dos valores: 0 o 1. Se emplea para calcular probabilidades (ya que los valores obtenidos estarán entre 0 y 1), o para clasificar eventos en dos categorías; en este caso, se puede utilizar para calcular la probabilidad de que una seta sea comestible (o, dicho de otra manera, para clasificar una seta como comestible o venenosa).

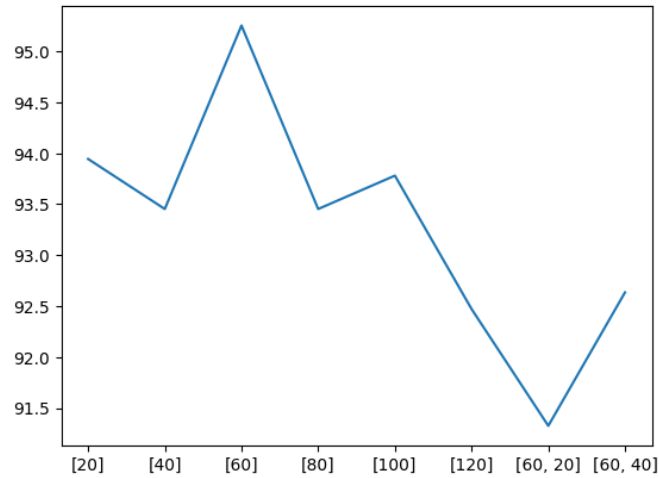
Para utilizar este método, implementamos la función de gradiente, de coste y la función sigmoide como hemos visto en clase, y obtuvimos una tasa de aciertos del 97.14262%, con un coste de 0.12467. Para llegar a estos resultados, utilizamos un valor para $\lambda = 3.0$ y un exponente = 2.



2.2 Redes Neuronales

Se entiende como red neuronal a un conjunto de capas ocultas, una capa de entrada y una capa de salida que intentan imitar el funcionamiento del cerebro humano. Cada nodo de las capas (el que equivaldría a una "neurona artificial"), se conecta a otros, y tiene un peso y umbral asociados. Si el valor de un nodo supera su umbral, la "neurona" se activa y envía los datos a la siguiente capa de la red.

Tras probar con varias configuraciones para la red neuronal, llegamos a un porcentaje de aciertos del 99.828%, con un coste de 0.0025. Para ello, utilizamos una red neuronal con dos capas ocultas (la primera, con 60 nodos; y la segunda, con 40), y una constante de regularización de 0.3. El algoritmo dio 1000 vueltas para llegar a este resultado.



2.3 SVM

Una Máquina de Vectores de Soporte (o SVM, por sus siglas en inglés), es un conjunto de algoritmos relacionados con problemas de clasificación y regresión. Una SVM intenta construir un hiperplano de un espacio N-dimensional que consiga clasificar los datos correctamente (donde N es el número de características que se proporcionan).

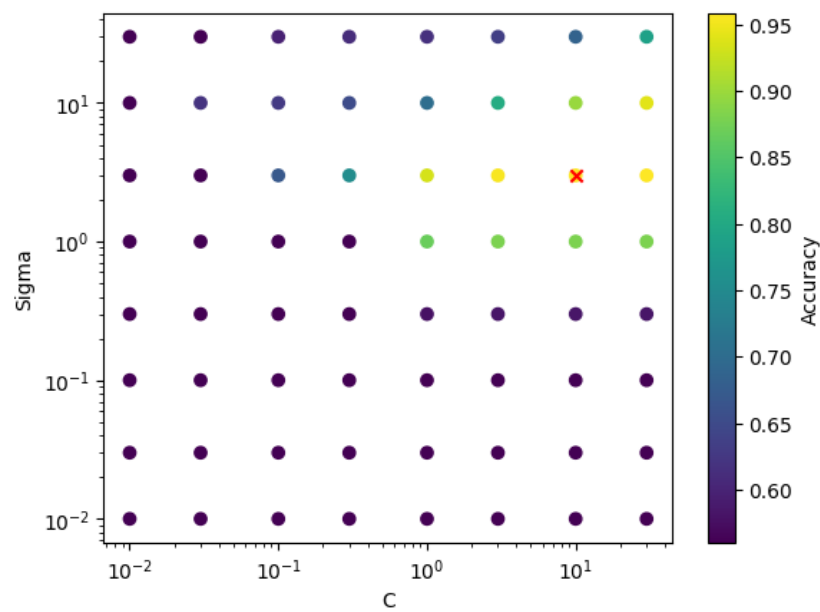
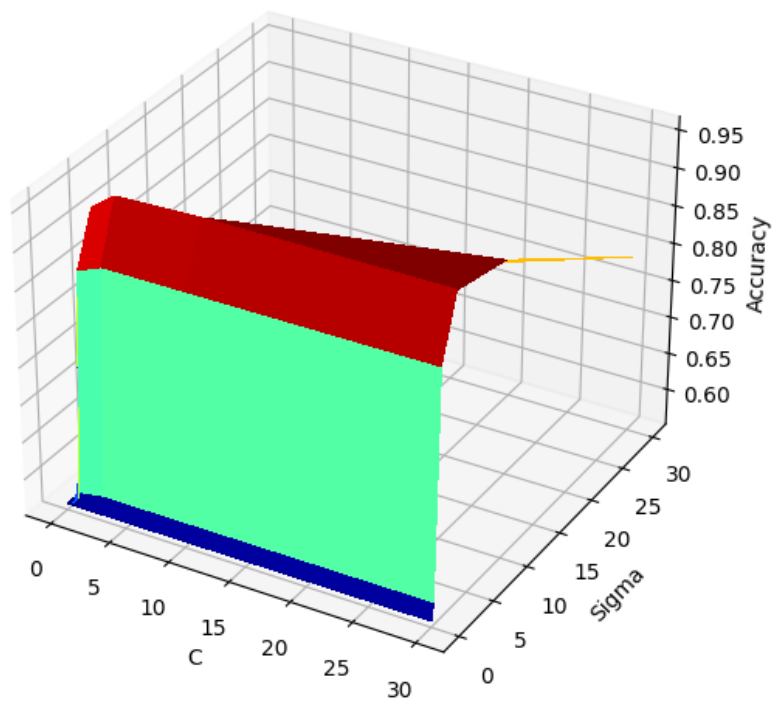
Entre las pruebas que realizamos, comparamos las que nos dieron los siguientes resultados:

Resultados obtenidos al aplicar SVM a nuestra base de datos.						
	C	Sigma	Porcentaje de datos usados (Training - Validation - Test)	Mejor probabilidad de acierto	Porcentaje de aciertos (%)	Tiempo de análisis (s)
1	3.0	3.0	0.20 - 0.20 - 0.20	0.999	99.984	5572
2	3.0	3.0	0.20 - 0.20 - 0.20	0.999	99.984	3856
3	10.0	3.0	0.05 - 0.05 - 0.05	0.996	99.705	104
4	10.0	3.0	0.01 - 0.01 - 0.01	0.959	95.908	1.95

Entre las dos primeras pruebas que se muestran, es notable la diferencia de tiempos. Esto se debe a que, después de varias pruebas, decidimos utilizar hebras para agilizar este proceso.

Entre los siguientes dos sets, el tiempo también se reduce drásticamente, pero por otro motivo: en estos casos, escogimos muestras más pequeñas para hacer las pruebas. Es interesante observar que, aunque el proceso es más rápido, el porcentaje de aciertos también decae.

Nuestro mejor resultado se corresponde, entonces, con la segunda entrada de la tabla: un 99.984% de aciertos con un valor de $C = 3.0$ y un valor de $\sigma = 3.0$, y reduciendo en casi 30 minutos el tiempo necesario para analizar los datos.

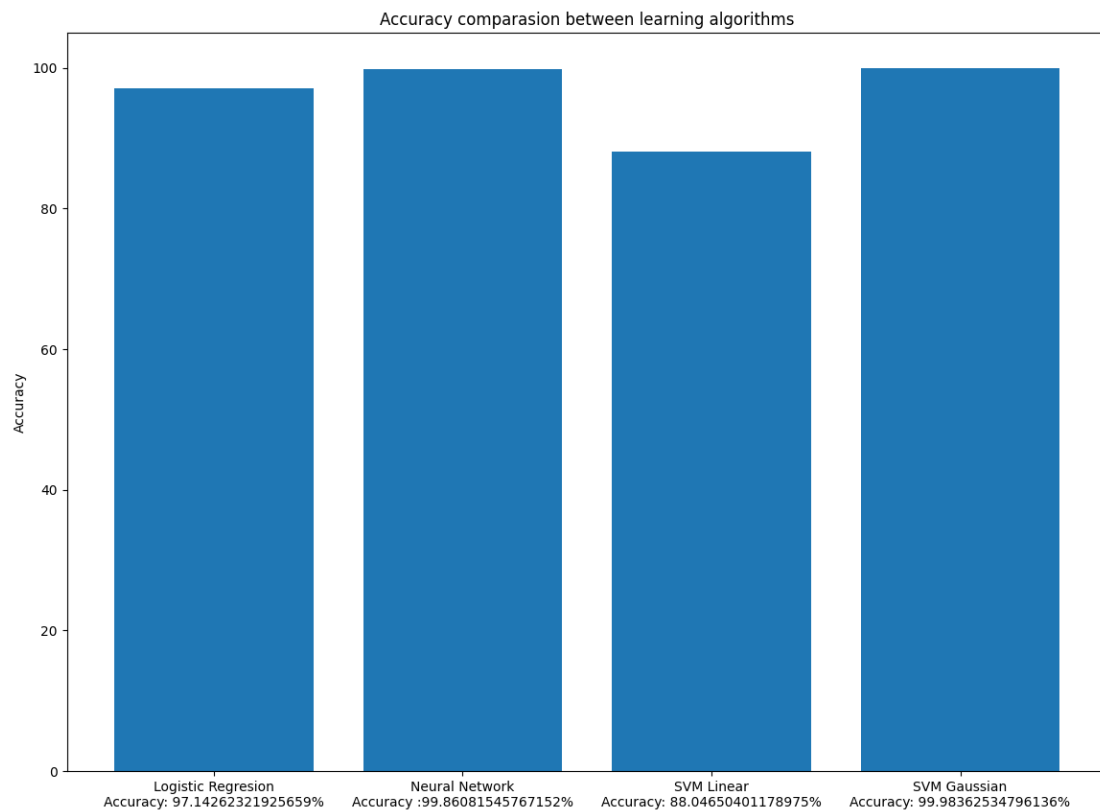


3 Conclusiones

Para calcular los algoritmos, utilizamos en todos los casos unos datasets con las siguientes características:

- Tamaño de la muestra de entrenamiento: 12213 entradas
- Tamaño de la muestra de validación: 12213 entradas
- Tamaño de la muestra de testeo: 12213 entradas
- En todos los casos, las muestras son diferentes entre sí.

Tras realizar todas las pruebas, obtuvimos los siguientes resultados:



3.1 Implementación

Para conseguir los resultados anteriores, implementamos las siguientes funciones, repartidas en varios archivos:

```
1000 import numpy as np
1001
1002 import sklearn.preprocessing as sk
1003 import scipy.optimize as opt
1004
1005 import matplotlib.pyplot as plt
1006 from matplotlib import cm
1007 from mpl_toolkits.mplot3d import Axes3D
1008
1009 import time
1010
1011 from threadRetVal import ThreadWithReturnValue
1012
1013 def sigmoide(z):
1014     return (1 / (1 + np.exp(-z)))
1015
1016 def coste(thetas, x, y):
1017     h = sigmoide(np.dot(x, thetas))
1018     return -1/len(x) * (np.dot(np.log(h), y) + np.dot(np.log(1-h), 1-y))
1019
1020 def costeReg(thetas, x, y, l):
1021     h = sigmoide(np.dot(x, thetas))
1022     return -((np.dot(np.log(h), y) + np.dot(np.log(1 - h), 1 - y)) / len(x)
1023              + (1/(2*len(x))) * l * np.sum(thetas[1:] ** 2))
1024
1025 def gradienteReg(thetas, x, y, l):
1026     h = sigmoide(np.dot(x, thetas))
1027     return np.dot(x.T, h-y) / len(y) + (thetas * l) / len(y)
1028
1029 def threadMethod(xPolTrain, xPolVal, yTrain, yVal, n, exp, l):
1030     print('Testing for exp: ' + str(exp) + ' and lambda: ' + str(l))
1031     thetas = np.zeros(n)
1032     thetas = opt.fmin_tnc(func=costeReg, x0=thetas, disp=False, fprime=
1033                          gradienteReg, args=(xPolTrain, yTrain, l))[0]
1034     cost = coste(thetas, xPolVal, yVal)
1035     print('Completed test for exp: ' + str(exp) + ' and lambda: ' + str(l))
1036     return thetas, cost
1037
1038 def evalLogisticReg(xTrain, xVal, yTrain, yVal):
1039     ls = np.array([0.01, 0.03, 0.1, 0.3, 1.0, 3.0, 10.0, 30.0])
1040     exps = np.array([1, 2]) # 3 asks for too much memory
1041
1042     numLs = ls.shape[0]
1043     numExps = exps.shape[0]
1044
1045     pol = np.empty(numExps, dtype=object)
1046
1047     for i in np.arange(numExps):
1048         pol[i] = sk.PolynomialFeatures(exps[i])
```

```

1047
1048     startTime = time.time()
1049
1050     resCost = np.zeros(numExps * numLs).reshape(numExps, numLs)
1051     resThet = np.empty_like(resCost, dtype=object)
1052     threads = np.empty_like(resCost, dtype=object) # can't be used, for
lack of storage
1053
1054     for i in np.arange(numExps):
1055         xPolTrain = pol[i].fit_transform(xTrain)
1056         xPolVal = pol[i].fit_transform(xVal)
1057         n = np.shape(xPolTrain)[1]
1058         for j in np.arange(numLs):
1059             threads[i,j] = ThreadWithReturnValue(target=threadMethod, args
=(xPolTrain, xPolVal, yTrain, yVal, n, exps[i], ls[j],))
1060             threads[i,j].start()
1061
1062     for i in np.arange(numExps):
1063         for j in np.arange(numLs):
1064             resThet[i,j], resCost[i,j] = threads[i,j].join()
1065
1066     bestCost = np.min(resCost)
1067     w = np.where(resCost == bestCost)
1068     bestExpIndex = w[0][0]
1069     bestLIndex = w[1][0]
1070     bestTheta = resThet[bestExpIndex, bestLIndex]
1071
1072     print()
1073     print('Best cost: ' + str(bestCost))
1074     print('Best lambda: ' + str(ls[bestLIndex]))
1075     print('Best exponent: ' + str(exps[bestExpIndex]))
1076
1077     endTime = time.time()
1078     print('Seconds elapsed of test: ' + str(endTime - startTime))
1079     print()
1080
1081     fig = plt.figure()
1082
1083     expexp, ll = np.meshgrid(ls, exps)
1084
1085     ax = Axes3D(fig, auto_add_to_figure=False)
1086     ax.set_xlabel('Exponent')
1087     ax.set_ylabel('L')
1088     ax.set_zlabel('Cost')
1089
1090     fig.add_axes(ax)
1091
1092     ax.plot_surface(ll, expexp, resCost, cmap=cm.jet, linewidth=0,
antialiased=False)
1093     plt.savefig('../Results/LogReg/LOG.1.png', bbox_inches='tight')
1094     plt.close()
1095
1096     plt.scatter(expexp, ll, s=80, c=resCost)

```

```

1097     plt.xscale('log')
1098     plt.xlabel('Lambda')
1099     plt.ylabel('Exponent')
1100     plt.clim(np.min(resCost), np.max(resCost))
1101     plt.colorbar().set_label('Cost')
1102     plt.scatter(ls[bestLIndex], exps[bestExpIndex], s=40, marker='x', color
='r')
1103     plt.savefig('../Results/LogReg/LOG.2.png', bbox_inches='tight')
1104     plt.close()
1105
1106     xPolVal = pol[bestExpIndex].fit_transform(xVal)
1107
1108     res = sigmoide(np.dot(bestTheta, xPolVal.T))
1109     acertados = np.sum((res >= 0.5) == yVal)
1110     accuracy = acertados*100/np.shape(res)[0]
1111     print("Accuracy of train: " + str(accuracy) + "%")
1112
1113     return bestTheta, pol[bestExpIndex], accuracy, bestCost, ls[bestLIndex
], exps[bestExpIndex]
1114
1115 def getNumAcertadosLog(theta, x, y):
1116     res = sigmoide(np.dot(theta, x.T))
1117     return np.sum((res >= 0.5) == y)

```

src/evaluateLogistic.py

```

1000 import numpy as np
1001
1002 from scipy.optimize import minimize
1003
1004 def sigmoide(z):
1005     return 1 / (1 + np.exp(-z)) + 1e-9
1006
1007 def forwardProp(x, num_capas, thetas):
1008     a = np.empty(num_capas + 1, dtype="object")
1009     a[0] = x
1010     for i in range(num_capas):
1011         aNew = np.hstack([np.ones([x.shape[0], 1]), a[i]])
1012         a[i] = aNew
1013         a[i+1] = sigmoide(np.dot(aNew, thetas[i].T))
1014     return a
1015
1016 def coste(x, y_ones, num_capas, thetas):
1017     res = forwardProp(x, num_capas, thetas)[num_capas]
1018
1019     return np.sum((-y_ones * np.log(res)) - ((1 - y_ones) * np.log(1-res)
)) / x.shape[0]
1020
1021 def costeRegul(x, y_ones, num_capas, thetas, reg):
1022     cost = coste(x, y_ones, num_capas, thetas)
1023
1024     val = 0
1025
1026     for i in range(num_capas):

```

```

1027         val += np.sum(np.power(thetas[i][1:], 2))
1028
1029     regul = val * (reg / (2*x.shape[0]))
1030
1031     return cost + regul
1032
1033 def grad_Delta(m, Delta, theta, reg):
1034     gradient = Delta
1035
1036     col0 = gradient[0]
1037     gradient = gradient + (reg/m)*theta
1038     gradient[0] = col0
1039
1040     return gradient
1041
1042 def lineal_back_prop(x, y, thetas, reg):
1043     m = x.shape[0]
1044     Delta1 = np.zeros_like(thetas[0])
1045     Delta2 = np.zeros_like(thetas[1])
1046
1047     hThetaTot = forwardProp(x, thetas.shape[0], thetas)
1048
1049     dlts = np.empty_like(thetas)
1050     Deltas = np.empty_like(thetas)
1051
1052     dlts[-1] = (hThetaTot[-1].T - y).T
1053
1054     for t in range(m):
1055         a1t = hThetaTot[0][t, :] # (401,)
1056         a2t = hThetaTot[1][t, :] # (26,)
1057         ht = hThetaTot[2][t, :] # (10,)
1058         yt = np.reshape(y[t], (1,)) # (10,)
1059
1060         d3t = ht - yt # (10,)
1061         d2t = np.dot(thetas[1].T, d3t) * (a2t * (1 - a2t)) # (26,)
1062
1063         Delta1 = Delta1 + np.dot(d2t[1:, np.newaxis], a1t[np.newaxis, :])
1064         Delta2 = Delta2 + np.dot(d3t[:, np.newaxis], a2t[np.newaxis, :])
1065
1066     Delta1 = Delta1 / m
1067     Delta2 = Delta2 / m
1068
1069     gradient1 = grad_Delta(m, Delta1, thetas[0], reg)
1070     gradient2 = grad_Delta(m, Delta2, thetas[1], reg)
1071
1072     return np.append(gradient1, gradient2).reshape(-1)
1073
1074 def vect_back_prop(x, y, thetas, reg):
1075     m = x.shape[0]
1076
1077     hThetaTot = forwardProp(x, thetas.shape[0], thetas)
1078
1079     dlts = np.empty_like(thetas)

```

```

1080     Deltas = np.empty_like(thetas)
1081
1082     #dlts[-1] = (hThetaTot[-1].T - y).T
1083     dlts[-1] = hThetaTot[-1] - y
1084
1085     for i in range(1, thetas.shape[0]):
1086         a = hThetaTot[-(i+1)]
1087
1088         delta = np.dot(thetas[-i].T, dlts[-i].T).T
1089
1090         delta = delta * a * (1-a)
1091
1092         delta = delta[:,1:]
1093
1094         dlts[-(i+1)] = delta
1095
1096     res = []
1097
1098     for i in range(thetas.shape[0]):
1099         Deltas[i] = np.dot(dlts[i].T, hThetaTot[i]) / m
1100         Deltas[i] = np.append(Deltas[i][0], Deltas[i][1:] + (reg/m) *
thetas[i][1:])
1101         res = np.append(res, Deltas[i]).reshape(-1)
1102
1103     return res
1104
1105 def backprop(params_rn, num_entradas, num_ocultas, num_etiquetas, x, y, reg
):
1106     # backprop devuelve una tupla (coste, gradiente) con el coste y el
gradiente de
1107     # una red neuronal de tres capas o mas, con num_entradas, num_ocultas
nodos en las capas
1108     # ocultas y num_etiquetas nodos en la capa de salida. Si m es el numero
de ejemplos
1109     # de entrenamiento, la dimension de 'X' es (m, num_entradas) y la de 'y
' es
1110     # (m, num_etiquetas)
1111
1112     if(num_ocultas.shape[0] + 2 < 3):
1113         print("ERROR: num_capas incorrect, must have an input, at least one
hidden and an output layer")
1114         return (0,0)
1115
1116     # calculo de thetas
1117     thetas = np.empty(num_ocultas.shape[0] + 1, dtype='object')
1118     pointer = num_ocultas[0] * (num_entradas + 1)
1119
1120     thetas[0] = np.reshape(params_rn[: pointer], (num_ocultas[0], (
num_entradas + 1)))
1121
1122     for i in range(1, num_ocultas.shape[0]):
1123         thetas[i] = np.reshape(params_rn[pointer : pointer + num_ocultas[i]
* (num_ocultas[i-1] + 1)], (num_ocultas[i], (num_ocultas[i-1] + 1)))

```

```

1124         pointer += num_ocultas[i] * (num_ocultas[i-1] + 1)
1125
1126         thetas[num_ocultas.shape[0]] = np.reshape(params_rn[pointer :], (
num_etiquetas, (num_ocultas[-1] + 1)))
1127
1128         #return costeRegul(x, y, thetas.shape[0], thetas, reg),
lineal_back_prop(x, y, thetas, reg)
1129         return costeRegul(x, y, thetas.shape[0], thetas, reg), vect_back_prop(x
, y, thetas, reg)
1130
1131 def getThetas(num_entradas, num_ocultas, num_etiquetas, out):
1132     thetas = np.empty(shape=[num_ocultas.shape[0] + 1], dtype='object')
1133     pointer = num_ocultas[0] * (num_entradas + 1)
1134
1135     thetas[0] = np.reshape(out.x[: pointer] , (num_ocultas[0], (
num_entradas + 1)))
1136
1137     for i in range(1, num_ocultas.shape[0]):
1138         thetas[i] = np.reshape(out.x[pointer : pointer + num_ocultas[i] * (
num_ocultas[i-1] + 1)], (num_ocultas[i], (num_ocultas[i-1] + 1)))
1139         pointer += num_ocultas[i] * (num_ocultas[i-1] + 1)
1140
1141     thetas[-1] = np.array(np.reshape(out.x[pointer :] , (num_etiquetas, (
num_ocultas[-1] + 1))), dtype='float')
1142
1143     return thetas
1144
1145 def evalNN(num_entradas, num_ocultas, num_etiquetas, xTrain, xVal, yTrain,
yVal, tagsTrain, tagsVal):
1146     eIni = np.sqrt(6) / np.sqrt(num_etiquetas + num_entradas) # = sqrt(6) /
sqrt(Lin + Lout)
1147     regs = np.array([0.01, 0.03, 0.1, 0.3, 0.1, 0.3, 1.0, 3.0, 10.0, 30.0])
1148     laps = np.array([25, 50, 75, 100, 250, 500, 750, 1000])
1149
1150     numRegs = regs.shape[0]
1151     numLaps = laps.shape[0]
1152
1153     resCost = np.zeros(numRegs * numLaps).reshape(numRegs, numLaps)
1154     resThet = np.empty_like(resCost, dtype=object)
1155
1156     pesosSize = (num_entradas + 1) * num_ocultas[0] + (num_ocultas[-1] + 1)
* num_etiquetas
1157
1158     for i in range(1, num_ocultas.shape[0]):
1159         pesosSize = pesosSize + ((num_ocultas[i-1] + 1) * num_ocultas[i])
1160
1161     randomTries = 1
1162
1163     bestCost = 10000
1164
1165     for i in np.arange(numRegs):
1166         for j in np.arange(numLaps):
1167             print(str(num_ocultas) + ' Testing for reg: ' + str(regs[i]) +

```

```

1168         pesos = np.random.uniform(-eIni, eIni, pesosSize)
1169         # because its random, finds the best out of a number of trials
1170         for x in np.arange(randomTries):
1171             out = minimize(fun = backprop, x0= pesos,
1172                 args = (num_entradas, num_ocultas, num_etiquetas,
1173 xTrain, tagsTrain, regs[i]),
1174                 method='TNC', jac = True, options = {'maxiter': laps[j]
1175 ]})
1176             thetas = getThetas(num_entradas, num_ocultas, num_etiquetas
1177 , out)
1178             cost = coste(xTrain, tagsTrain, thetas.shape[0], thetas)
1179             if(cost < bestCost):
1180                 bestCost = cost
1181                 bestTheta = thetas
1182                 bestRegIndex = i
1183                 bestLapsIndex = j
1184
1185 print()
1186
1187 print(str(num_ocultas) + " Best cost: " + str(bestCost))
1188 print(str(num_ocultas) + " Best Reg: " + str(regs[bestRegIndex]))
1189 print(str(num_ocultas) + " Best Laps: " + str(laps[bestLapsIndex]))
1190
1191 res = forwardProp(xVal, bestTheta.shape[0], bestTheta)[-1]
1192 maxIndices = np.argmax(res, axis=1)
1193 acertados = np.sum(maxIndices == yVal)
1194 accuracy = acertados*100/np.shape(res)[0]
1195 print(str(num_ocultas) + " Accuracy of train: " + str(accuracy) + "%")
1196 print()
1197
1198 return bestTheta, accuracy, bestCost, regs[bestRegIndex], laps[
1199 bestLapsIndex]
1200
1201 defgetNumAcertadosNN(x, y, thetas):
1202     res = forwardProp(x, thetas.shape[0], thetas)[-1]
1203     maxIndices = np.argmax(res, axis=1)
1204     return np.sum(maxIndices == y)

```

src/evaluateNeuronal.py

```

1000 import numpy as np
1001
1002 import sklearn.svm as svm
1003
1004 import matplotlib.pyplot as plt
1005 from matplotlib import cm
1006 from mpl_toolkits.mplot3d import Axes3D
1007
1008 from sklearn.metrics import accuracy_score
1009
1010 import time
1011
1012 from threadRetVal import ThreadWithReturnValue

```



```

1013
1014 def threadMethod(x, xVal, y, yVal, c, sigma):
1015     print('Testing c ' + str(c) + ' sigma ' + str(sigma))
1016     s = svm.SVC(kernel='rbf', C=c, gamma= 1 / (2 * sigma ** 2))
1017     s.fit(x,y)
1018
1019     print('Completed c ' + str(c) + ' sigma ' + str(sigma))
1020     return s, accuracy_score(yVal, s.predict(xVal))
1021
1022 def evaluateSVM(x, xVal, y, yVal):
1023     startTime = time.time()
1024
1025     cs = np.array([0.01, 0.03, 0.1, 0.3, 1, 3, 10, 30])
1026     sigmas = np.array([0.01, 0.03, 0.1, 0.3, 1, 3, 10, 30])
1027
1028     numCs = cs.shape[0]
1029     numSigmas = cs.shape[0]
1030
1031     resAcc = np.zeros(numCs * numSigmas).reshape(numCs, numSigmas)
1032     resSVM = np.empty_like(resAcc, dtype=object)
1033     threads = np.empty_like(resAcc, dtype=object)
1034
1035     for i in np.arange(numCs):
1036         for j in np.arange(numSigmas):
1037             threads[i,j] = ThreadWithReturnValue(target=threadMethod, args
=(x, xVal, y, yVal, cs[i], sigmas[j],))
1038             threads[i,j].start()
1039
1040     for i in np.arange(numCs):
1041         for j in np.arange(numSigmas):
1042             resSVM[i,j], resAcc[i,j] = threads[i,j].join()
1043
1044     bestAcc = np.max(resAcc)
1045     w = np.where(resAcc == bestAcc)
1046     bestC = cs[w[0][0]]
1047     bestSigma = sigmas[w[1][0]]
1048     bestSVM = resSVM[w[0][0],w[1][0]]
1049
1050     print()
1051     print("Best accuracy: " + str(bestAcc))
1052     print("Best C: " + str(bestC))
1053     print("Best Sigma: " + str(bestSigma))
1054
1055     endTime = time.time()
1056     print('Seconds elapsed of test: ' + str(endTime - startTime))
1057
1058     fig = plt.figure()
1059
1060     cc, ss = np.meshgrid(sigmas, cs)
1061
1062     ax = Axes3D(fig, auto_add_to_figure=False)
1063     ax.set_xlabel('C')
1064     ax.set_ylabel('Sigma')

```

```

1065     ax.set_zlabel('Accuracy')
1066
1067     fig.add_axes(ax)
1068
1069     ax.plot_surface(ss, cc, resAcc, cmap=cm.jet, linewidth=0, antialiased=
False)
1070     plt.savefig('../Results/SVM/SVM.1.png', bbox_inches='tight')
1071     plt.close()
1072
1073     plt.scatter(ss, cc, c=resAcc)
1074     plt.xscale('log')
1075     plt.yscale('log')
1076     plt.xlabel('C')
1077     plt.ylabel('Sigma')
1078     plt.clim(np.min(resAcc), np.max(resAcc))
1079     plt.colorbar().set_label('Accuracy')
1080     plt.scatter(bestC, bestSigma, marker='x', color='r')
1081     plt.savefig('../Results/SVM/SVM.2.png', bbox_inches='tight')
1082     plt.close()
1083
1084     return bestSVM, bestAcc, bestC, bestSigma
1085
1086 def threadMethodLinear(x, xVal, y, yVal, c):
1087     print('Testing c ' + str(c))
1088     s = svm.SVC(kernel='linear', C=c)
1089     s.fit(x, y)
1090
1091     print('Completed c ' + str(c))
1092     return s, accuracy_score(yVal, s.predict(xVal))
1093
1094 def evaluateSVMLinear(x, xVal, y, yVal):
1095     startTime = time.time()
1096
1097     cs = np.array([0.01, 0.03, 0.1, 0.3, 1, 3, 10, 30])
1098
1099     numCs = cs.shape[0]
1100
1101     resAcc = np.zeros(numCs)
1102     resSVM = np.empty_like(resAcc, dtype=object)
1103     threads = np.empty_like(resAcc, dtype=object)
1104
1105     for i in np.arange(numCs):
1106         threads[i] = ThreadWithReturnValue(target=threadMethodLinear, args
=(x, xVal, y, yVal, cs[i],))
1107         threads[i].start()
1108
1109     for i in np.arange(numCs):
1110         resSVM[i], resAcc[i] = threads[i].join()
1111
1112     bestAcc = np.max(resAcc)
1113     w = np.where(resAcc == bestAcc)
1114     bestC = cs[w[0]]
1115     bestSVM = resSVM[w[0]][0]

```

```

1116
1117     print()
1118     print("Best accuracy: " + str(bestAcc))
1119     print("Best C: " + str(bestC))
1120
1121     endTime = time.time()
1122     print('Seconds elapsed of test: ' + str(endTime - startTime))
1123
1124     plt.plot(cs, resAcc)
1125     plt.xscale('log')
1126     plt.xlabel('C')
1127     plt.ylabel('Accuracy')
1128     plt.scatter(bestC, bestAcc, marker='x', color='r')
1129     plt.savefig('../Results/SVM/SVM.LIN.png', bbox_inches='tight')
1130     plt.close()
1131
1132     return bestSVM, bestAcc, bestC

```

src/evaluateSVM.py

```

1000 from enum import Enum
1001
1002 import numpy as np
1003 import matplotlib.pyplot as plt
1004 import seaborn as sns
1005
1006 #region Enums
1007 class CapShape(Enum):
1008     b = 0
1009     c = 1
1010     x = 2
1011     f = 3
1012     s = 4
1013     p = 5
1014     o = 6
1015
1016 class Surface(Enum):
1017     i = 1
1018     g = 2
1019     y = 3
1020     s = 4
1021     h = 5
1022     l = 6
1023     k = 7
1024     t = 8
1025     w = 9
1026     e = 10
1027     d = 11 # TO DO: wat
1028     f = 12
1029
1030 class Color(Enum):
1031     n = 1
1032     b = 2
1033     g = 3

```

```

1034         r = 4
1035         p = 5
1036         u = 6
1037         e = 7
1038         w = 8
1039         y = 9
1040         l = 10
1041         o = 11
1042         k = 12
1043         f = 13
1044
1045     class GillAttach(Enum):
1046         empty = 0
1047         a = 1
1048         x = 2
1049         d = 3
1050         e = 4
1051         s = 5
1052         p = 6
1053         f = 7
1054
1055     class GillSpacing(Enum):
1056         empty = 0
1057         c = 1
1058         d = 2
1059         f = 3
1060
1061     class StemRoot(Enum):
1062         empty = 0
1063         b = 1
1064         s = 2
1065         c = 3
1066         u = 4
1067         e = 5
1068         z = 6
1069         r = 7
1070         f = 8
1071
1072     class VeilType(Enum):
1073         empty = 0
1074         p = 1
1075         u = 2
1076
1077     class RingType(Enum):
1078         empty = 0
1079         c = 1
1080         e = 2
1081         r = 3
1082         g = 4
1083         l = 5
1084         p = 6
1085         s = 7
1086         z = 8

```

```

1087         y = 9
1088         m = 10
1089         f = 11
1090
1091     class Habitat(Enum):
1092         g = 0
1093         l = 1
1094         m = 2
1095         p = 3
1096         h = 4
1097         u = 5
1098         w = 6
1099         d = 7
1100
1101     class Season(Enum):
1102         s = 0
1103         u = 1
1104         a = 2
1105         w = 3
1106 #endregion
1107
1108     def prepareData(data):
1109         data['class'] = [int(i == 'e') for i in data['class']]
1110
1111         for i in CapShape._member_names_:
1112             data['cap-shape-' + i] = np.where(data['cap-shape'] == i, 1, 0)
1113         data.pop('cap-shape')
1114
1115         for i in Surface._member_names_:
1116             data['cap-surface-' + i] = np.where(data['cap-surface'] == i, 1, 0)
1117         data.pop('cap-surface')
1118
1119         for i in Color._member_names_:
1120             data['cap-color-' + i] = np.where(data['cap-color'] == i, 1, 0)
1121         data.pop('cap-color')
1122
1123         data['does-bruise-or-bleed'] = [int(i == 't') for i in data['does-
bruise-or-bleed']]
1124
1125         for i in GillAttach._member_names_:
1126             data['gill-attachment-' + i] = np.where(data['gill-attachment'] ==
i, 1, 0)
1127         data.pop('gill-attachment')
1128
1129         for i in GillSpacing._member_names_:
1130             data['gill-spacing-' + i] = np.where(data['gill-spacing'] == i, 1,
0)
1131         data.pop('gill-spacing')
1132
1133         for i in Color._member_names_:
1134             data['gill-color-' + i] = np.where(data['gill-color'] == i, 1, 0)
1135         data.pop('gill-color')
1136

```

```

1137     for i in StemRoot._member_names_:
1138         data['stem-root-' + i] = np.where(data['stem-root'] == i, 1, 0)
1139     data.pop('stem-root')
1140
1141     for i in Surface._member_names_:
1142         data['stem-surface-' + i] = np.where(data['stem-surface'] == i, 1,
1143         0)
1144     data.pop('stem-surface')
1145
1146     for i in Color._member_names_:
1147         data['stem-color-' + i] = np.where(data['stem-color'] == i, 1, 0)
1148     data.pop('stem-color')
1149
1150     for i in VeilType._member_names_:
1151         data['veil-type-' + i] = np.where(data['veil-type'] == i, 1, 0)
1152     data.pop('veil-type')
1153
1154     for i in Color._member_names_:
1155         data['veil-color-' + i] = np.where(data['veil-color'] == i, 1, 0)
1156     data.pop('veil-color')
1157
1158     data['has-ring'] = [int(i == 't') for i in data['has-ring']]
1159
1160     for i in RingType._member_names_:
1161         data['ring-type-' + i] = np.where(data['ring-type'] == i, 1, 0)
1162     data.pop('ring-type')
1163
1164     for i in Color._member_names_:
1165         data['spore-print-color-' + i] = np.where(data['spore-print-color']
1166         == i, 1, 0)
1167     data.pop('spore-print-color')
1168
1169     for i in Habitat._member_names_:
1170         data['habitat-' + i] = np.where(data['habitat'] == i, 1, 0)
1171     data.pop('habitat')
1172
1173     for i in Season._member_names_:
1174         data['season-' + i] = np.where(data['season'] == i, 1, 0)
1175     data.pop('season')
1176
1177     print("\nINITIAL VARIABLES: " + str(len(list(data.columns))))
1178     print("\nCOLUMNS WITH ALL 0 VALUES:\n")
1179
1180     for i in (list(data.columns)):
1181         if (data[i] == 0).all():
1182             print(i)
1183             data.pop(i)
1184
1185     print("\nTOTAL VARIABLES: " + str(len(list(data.columns))))
1186     print()
1187
1188 def createGraphMetricalValue(data, edibles, poisonous, tag, title, xLabel,
1189 yLabel, fileName):

```

```

1187     plt.plot(np.arange(len(data[tag][edibles])), np.sort(data[tag][edibles
1188 ])[:-1], label='Edible')
1189     plt.plot(np.arange(len(data[tag][poisonous])), np.sort(data[tag][
1190 poisonous])[:-1], c='#ff7f0e', label='Poisonous')
1191     plt.title(title)
1192     plt.xlabel(xLabel)
1193     plt.ylabel(yLabel)
1194     plt.legend()
1195     plt.savefig('../Results/Analysis/' + fileName, bbox_inches='tight')
1196     plt.close()
1197
1198 def createGraphBinaryValue(data, edibles, poisonous, tag, labels, width,
1199 title, fileName):
1200     x = np.arange(len(labels)) # the label locations
1201
1202     edib = [len(np.where(data[tag][edibles])[0]), len(np.where(data[tag][
1203 edibles] == 0)[0])]
1204     pois = [len(np.where(data[tag][poisonous])[0]), len(np.where(data[tag][
1205 poisonous] == 0)[0])]
1206
1207     fig, ax = plt.subplots()
1208     ax.bar(x - width/2, edib, width, label='Edible')
1209     ax.bar(x + width/2, pois, width, label='Poisonous')
1210     plt.ylabel('Mushrooms')
1211     plt.title(title)
1212     plt.xticks(x, labels)
1213     plt.legend()
1214     plt.savefig('../Results/Analysis/' + fileName, bbox_inches='tight')
1215     plt.close()
1216
1217 def createGraphEnumValue(data, edibles, poisonous, tags, labels, width,
1218 title, fileName):
1219     x = np.arange(len(labels)) # the label locations
1220
1221     edib = []
1222     pois = []
1223
1224     for tag in tags:
1225         edib.append(len(np.where(data[tag][edibles])[0]))
1226         pois.append(len(np.where(data[tag][poisonous])[0]))
1227
1228     fig, ax = plt.subplots(figsize=(16,9))
1229     ax.bar(x - width/2, edib, width, label='Edible')
1230     ax.bar(x + width/2, pois, width, label='Poisonous')
1231     plt.ylabel('Mushrooms')
1232     plt.title(title)
1233     plt.xticks(x, labels)
1234     plt.legend()
1235     plt.savefig('../Results/Analysis/' + fileName, bbox_inches='tight')
1236     plt.close()
1237
1238 def analyzeData(data):
1239     print(data.head())

```

```

1234
1235 print(data.describe())
1236
1237 edibles = np.where(data['class'])[0]
1238 poisonous = np.where(data['class'] == 0)[0]
1239
1240 numEdibles = len(edibles)
1241 numPoisonous = len(poisonous)
1242
1243 classMush = [numEdibles, numPoisonous]
1244
1245 # 00 number of edibles and poisonous
1246 plt.bar(np.arange(2), classMush[1].set_color('#ff7f0e'))
1247 plt.xticks(np.arange(2), ['Edibles: ' + str(len(edibles)), 'Poisonous: '
1248 + str(len(poisonous))])
1249 plt.ylabel('Quantity')
1250 plt.title('Number of poisonous and edibles mushrooms')
1251 plt.savefig('../Results/Analysis/00EdiblesNumb.png', bbox_inches='tight')
1252 plt.close()
1253
1254 # 01 cap diameter
1255 createGraphMetricalValue(data, edibles, poisonous, 'cap-diameter',
1256 'Cap diameter of poisonous and edible mushrooms', 'Mushroom', 'Diameter', '01CapDiameter.png')
1257
1258 # 02 cap shape
1259 createGraphEnumValue(data, edibles, poisonous,
1260 ['cap-shape-b', 'cap-shape-c', 'cap-shape-x', 'cap-shape-f', 'cap-shape-s', 'cap-shape-p', 'cap-shape-o'],
1261 ['Bell', 'Conical', 'Convex', 'Flat', 'Sunken', 'Spherical', 'Other'],
1262 0.35,
1263 'Cap shape by class', '02CapShape.png')
1264
1265 # 03 cap surface
1266 createGraphEnumValue(data, edibles, poisonous,
1267 ['cap-surface-i', 'cap-surface-g', 'cap-surface-y', 'cap-surface-s', 'cap-surface-h', 'cap-surface-l',
1268 'cap-surface-k', 'cap-surface-t', 'cap-surface-w', 'cap-surface-e', 'cap-surface-d'],
1269 ['Fibrous', 'Grooves', 'Scaly', 'Smooth', 'Shiny', 'Leatherly',
1270 'Silky', 'Sticky', 'Wrinkled', 'Fleshy', 'd(Unknown)'], 0.35,
1271 'Cap surface by class', '03CapSurface.png')
1272
1273 # 04 cap color
1274 createGraphEnumValue(data, edibles, poisonous,
1275 ['cap-color-n', 'cap-color-b', 'cap-color-g', 'cap-color-r', 'cap-color-p', 'cap-color-u',
1276 'cap-color-e', 'cap-color-w', 'cap-color-y', 'cap-color-l', 'cap-color-o', 'cap-color-k'],
1277 ['Brown', 'Buff', 'Gray', 'Green', 'Pink', 'Purple',
1278 'Red', 'White', 'Yellow', 'Blue', 'Orange', 'Black'], 0.35,
1279 'Cap color by class', '04CapColor.png')

```



```

1278
1279 # 05 does bruise
1280 createGraphBinaryValue(data, edibles, poisonous, 'does-bruise-or-bleed',
1281
1282 ["Bruise or Bleed", "Doesn't"], 0.35,
1283 'Does bruise or bleed by class', '05BruiseOrBleed.png')
1284
1285 # 06 gill attachment
1286 createGraphEnumValue(data, edibles, poisonous,
1287 ['gill-attachment-empty', 'gill-attachment-a', 'gill-attachment-x', '
1288 gill-attachment-d',
1289 'gill-attachment-e', 'gill-attachment-s', 'gill-attachment-p', 'gill-
1290 attachment-f'],
1291 ['Empty value', 'Adnate', 'Adnexed', 'Decurrent',
1292 'Free', 'Sinuate', 'Pores', 'None'], 0.35,
1293 'Gill attachment by class', '06GillAttachment.png')
1294
1295 # 07 gill spacing
1296 createGraphEnumValue(data, edibles, poisonous,
1297 ['gill-spacing-empty', 'gill-spacing-c', 'gill-spacing-d', 'gill-
1298 spacing-f'],
1299 ['Empty value', 'Close', 'Distant', 'None'], 0.35,
1300 'Gill spacing by class', '07GillSpacing.png')
1301
1302 # 08 gill color
1303 createGraphEnumValue(data, edibles, poisonous,
1304 ['gill-color-n', 'gill-color-b', 'gill-color-g', 'gill-color-r', 'gill-
1305 color-p', 'gill-color-u',
1306 'gill-color-e', 'gill-color-w', 'gill-color-y', 'gill-color-o', 'gill-
1307 color-k', 'gill-color-f'],
1308 ['Brown', 'Buff', 'Gray', 'Green', 'Pink', 'Purple',
1309 'Red', 'White', 'Yellow', 'Orange', 'Black', 'None'], 0.35,
1310 'Gill color by class', '08GillColor.png')
1311
1312 # 09 stem height
1313 createGraphMetricalValue(data, edibles, poisonous, 'stem-height',
1314 'Stem height of poisonous and edible mushrooms', 'Mushroom', 'Height',
1315 '09StemHeight.png')
1316
1317 # 10 stem width
1318 createGraphMetricalValue(data, edibles, poisonous, 'stem-width',
1319 'Stem width of poisonous and edible mushrooms', 'Mushroom', 'Width', '
1320 10StemWidth.png')
1321
1322 # 11 stem root
1323 createGraphEnumValue(data, edibles, poisonous,
1324 ['stem-root-empty', 'stem-root-b', 'stem-root-s',
1325 'stem-root-c', 'stem-root-r', 'stem-root-f'],
1326 ['Empty value', 'Bulbous', 'Swollen',
1327 'Club', 'Rooted', 'None'], 0.35,
1328 'Stem root by class', '11StemRoot.png')
1329
1330 # 12 stem surface

```

```

1323 createGraphEnumValue(data, edibles, poisonous,
1324 [ 'stem-surface-i', 'stem-surface-g', 'stem-surface-y', 'stem-surface-s',
1325 'stem-surface-h', 'stem-surface-k', 'stem-surface-t', 'stem-surface-f'
1326 ],
1327 [ 'Fibrous', 'Grooves', 'Scaly', 'Smooth',
1328 'Shiny', 'Silky', 'Sticky', 'None'], 0.35,
1329 'Stem surface by class', '12StemSurface.png')
1330
1331 # 13 stem color
1332 createGraphEnumValue(data, edibles, poisonous,
1333 [ 'stem-color-n', 'stem-color-b', 'stem-color-g', 'stem-color-r', 'stem-
1334 color-p', 'stem-color-u',
1335 'stem-color-e', 'stem-color-w', 'stem-color-y', 'stem-color-l', 'stem-
1336 color-o', 'stem-color-k', 'stem-color-f'],
1337 [ 'Brown', 'Buff', 'Gray', 'Green', 'Pink', 'Purple',
1338 'Red', 'White', 'Yellow', 'Blue', 'Orange', 'Black', 'None'], 0.35,
1339 'Stem color by class', '13StemColor.png')
1340
1341 # 14 veil type
1342 createGraphEnumValue(data, edibles, poisonous,
1343 [ 'veil-type-empty', 'veil-type-u'],
1344 [ 'Partial', 'Universal'], 0.35,
1345 'Veil type by class', '14VeilType.png')
1346
1347 # 15 veil color
1348 createGraphEnumValue(data, edibles, poisonous,
1349 [ 'veil-color-n', 'veil-color-u', 'veil-color-e',
1350 'veil-color-w', 'veil-color-y', 'veil-color-k'],
1351 [ 'Brown', 'Purple', 'Red',
1352 'White', 'Yellow', 'Black'], 0.35,
1353 'Veil color by class', '15VeilColor.png')
1354
1355 # 16 has ring
1356 createGraphBinaryValue(data, edibles, poisonous, 'has-ring',
1357 ["Has ring", "Doesn't"], 0.35,
1358 'Does it has a ring by class', '16HasRing.png')
1359
1360 # 17 ring type
1361 createGraphEnumValue(data, edibles, poisonous,
1362 [ 'ring-type-empty', 'ring-type-e', 'ring-type-r', 'ring-type-g', 'ring-
1363 type-l',
1364 'ring-type-p', 'ring-type-z', 'ring-type-m', 'ring-type-f'],
1365 [ 'Empty value', 'Evanescent', 'Flaring', 'Grooved', 'Large',
1366 'Pendant', 'Zone', 'Movable', 'None'], 0.35,
1367 'Ring type by class', '17RingType.png')
1368
1369 # 18 spore print color
1370 createGraphEnumValue(data, edibles, poisonous,
1371 [ 'spore-print-color-n', 'spore-print-color-g', 'spore-print-color-r', '
1372 spore-print-color-p',
1373 'spore-print-color-u', 'spore-print-color-w', 'spore-print-color-k'],
1374 [ 'Brown', 'Gray', 'Green', 'Pink',

```

```

1370     'Purple', 'White', 'Black'], 0.35,
1371     'Spore print color by class', '18SporePrintColor.png')
1372
1373 # 19 habitat
1374 createGraphEnumValue(data, edibles, poisonous,
1375     ['habitat-g', 'habitat-l', 'habitat-m', 'habitat-p',
1376     'habitat-h', 'habitat-u', 'habitat-w', 'habitat-d'],
1377     ['Grasses', 'Leaves', 'Meadows', 'Paths',
1378     'Heaths', 'Urban', 'Waste', 'Woods'], 0.35,
1379     'Habitat by class', '19Habitat.png')
1380
1381 # 19 season
1382 createGraphEnumValue(data, edibles, poisonous,
1383     ['season-s', 'season-u', 'season-a', 'season-w'],
1384     ['Spring', 'Summer', 'Autumn', 'Winter'], 0.35,
1385     'Season by class', '20Season.png')
1386
1387 # heatmap best
1388 plt.figure(figsize=(6, 8))
1389 sns.heatmap(data.corr()[['class']].sort_values('class', ascending=False)
1390     [1:31], annot=True, vmin=-1, vmax=1)
1391 plt.title('Heat map of top 30 most correlated variables to being edible')
1392 plt.savefig('../Results/Analysis/heatMap30Best.png', bbox_inches='tight')
1393 plt.close()
1394
1395 # heatmap worst
1396 plt.figure(figsize=(6, 8))
1397 sns.heatmap(data.corr()[['class']].sort_values('class', ascending=False)
1398     [-30:], annot=True, vmin=-1, vmax=1)
1399 plt.title('Heat map of top 30 least correlated variables to being edible')
1400 plt.savefig('../Results/Analysis/heatMap30Worst.png', bbox_inches='tight')
1401 plt.close()
1402
1403 # collage
1404 # data.hist(figsize=(10, 5))
1405 # plt.tight_layout()
1406 # plt.savefig('../Results/Analysis/collageWithAllGraphs.png',
1407     bbox_inches='tight')
1408 # plt.close()

```

src/prepareData.py

```

1000 from threading import Thread
1001
1002 class ThreadWithReturnValue(Thread):
1003     def __init__(self, group=None, target=None, name=None,
1004         args=(), kwargs={}, Verbose=None):
1005         Thread.__init__(self, group, target, name, args, kwargs)
1006         self._return = None
1007     def run(self):

```

```

1008         if self._target is not None:
1009             self._return = self._target(*self._args,
1010                                         **self._kwargs)
1011     def join(self, *args):
1012         Thread.join(self, *args)
1013     return self._return

```

src/threadRetVal.py

```

1000 import numpy as np
1001 import matplotlib.pyplot as plt
1002
1003 import sklearn.svm as svm
1004 from sklearn.metrics import accuracy_score
1005
1006 from pandas.io.parsers import read_csv
1007
1008 from prepareData import analyzeData, prepareData
1009 from evaluateSVM import evaluateSVM, evaluateSVMLinear
1010 from evaluateNeuronal import evalNN, getNumAcertadosNN
1011 from evaluateLogistic import evalLogisticReg, getNumAcertadosLog
1012
1013 from threadRetVal import ThreadWithReturnValue
1014
1015 import time
1016
1017 def loadCSV(fileName):
1018     data = read_csv(fileName, sep=';', on_bad_lines='skip')
1019     data.fillna('empty', inplace=True)
1020     prepareData(data)
1021     return data
1022
1023 def evalLogistic(xTrain, xVal, xTest, yTrain, yVal, yTest):
1024     print("\nCOMENCING TRAINING OF LOGISTIC REGRESION\n")
1025
1026     th, pol, acc, cost, l, exp = evalLogisticReg(xTrain, xVal, yTrain, yVal
1027 )
1028
1029     xPolTest = pol.fit_transform(xTest)
1030
1031     acertados = getNumAcertadosLog(th, xPolTest, yTest)
1032     accuracy = acertados*100/np.shape(xTest)[0]
1033     print('Accuracy over Test sample: ' + str(accuracy) + "%")
1034
1035     return accuracy, acc, cost, l, exp
1036
1037 def evalNueralThread(ocultas, num_entradas, num_etiquetas, xTrain, xVal,
1038 xTest, yTrain, yVal, yTest, tagsTrain, tagsVal):
1039     print("\nTesting with " + str(ocultas))
1040     num_ocultas = np.array(ocultas)
1041     th, acc, cost, reg, laps = evalNN(
1042         num_entradas, num_ocultas, num_etiquetas,
1043         xTrain, xVal, yTrain, yVal, tagsTrain, tagsVal)
1044

```

```

1043     acertados = getNumAcertadosNN(xTest, yTest, th)
1044
1045     accuracy = acertados*100/np.shape(xTest)[0]
1046     print("Accuracy over Test sample: " + str(accuracy) + "%")
1047     return accuracy, acc, cost, reg, laps
1048
1049 def evalNeuronal(xTrain, xVal, xTest, yTrain, yVal, yTest, n):
1050     print("\nCOMENCING TRAINING OF NEURONAL NETWORK\n")
1051
1052     ocultas = np.array([[20], [40], [60], [80], [100], [120], [60, 20],
1053                        [60, 40]])
1054     numOcultas = ocultas.shape[0]
1055
1056     resAcc = np.zeros(numOcultas)
1057     resAccVal = np.zeros(numOcultas)
1058     resCost = np.zeros(numOcultas)
1059     resReg = np.zeros(numOcultas)
1060     resLaps = np.zeros(numOcultas)
1061
1062     startTime = time.time()
1063
1064     tagsTrain = np.zeros((len(yTrain), 2))
1065     for i in range(len(yTrain)):
1066         tagsTrain[i][int(yTrain[i])] = 1
1067
1068     tagsVal = np.zeros((len(yVal), 2))
1069     for i in range(len(yVal)):
1070         tagsVal[i][int(yVal[i])] = 1
1071
1072     num_etiquetas = 2
1073     num_entradas = n
1074
1075     threads = np.empty(numOcultas, dtype=object)
1076
1077     for i in np.arange(numOcultas):
1078         threads[i] = ThreadWithReturnValue(target=evalNueronalThread, args
1079         =(ocultas[i], num_entradas, num_etiquetas, xTrain, xVal, xTest, yTrain,
1080         yVal, yTest, tagsTrain, tagsVal,))
1081         threads[i].start()
1082
1083     for i in np.arange(numOcultas):
1084         resAcc[i], resAccVal[i], resCost[i], resReg[i], resLaps[i] =
1085         threads[i].join()
1086
1087     bestAcc = np.max(resAcc)
1088     bestOcIndex = np.where(bestAcc == resAcc)[0]
1089     bestCost = resCost[bestOcIndex]
1090     bestReg = resReg[bestOcIndex]
1091     bestLaps = resLaps[bestOcIndex]
1092     bestAccVal = resAccVal[bestOcIndex]
1093
1094     print()
1095     print('Best neurons in hidden layer: ' + str(ocultas[bestOcIndex]))

```

```

1092     print('Best Accuracy: ' + str(bestAcc) + '%')
1093     print("Accuracy of train: " + str(bestAccVal) + "%")
1094     print('Best cost: ' + str(bestCost))
1095     print('Best reg: ' + str(bestReg))
1096     print('Best laps: ' + str(bestLaps))
1097
1098     endTime = time.time()
1099     print('Seconds elapsed of test: ' + str(endTime - startTime))
1100
1101     ocultasStr = []
1102     for i in np.arange(numOcultas):
1103         ocultasStr.append(str(ocultas[i]))
1104
1105     plt.plot(np.arange(numOcultas), resAcc)
1106     plt.xticks(np.arange(numOcultas), ocultasStr)
1107     plt.savefig(' ../ Results/NeuronalNetwork/NN.png', bbox_inches='tight')
1108     plt.close()
1109
1110     return ocultas[bestOcIndex], bestAcc, bestAccVal, bestCost, bestReg,
bestLaps
1111
1112 def evalSVM(xTrain, xVal, xTest, yTrain, yVal, yTest):
1113     print("\nCOMENCING TRAINING OF SVM\n")
1114
1115     s, acc, c, sig = evaluateSVM(xTrain, xVal, yTrain, yVal)
1116     accuracy = accuracy_score(yTest, s.predict(xTest)) * 100
1117     print('Accuracy over Test sample: ' + str(accuracy) + '%')
1118
1119     return accuracy, acc, c, sig
1120
1121 def evalSVMLinear(xTrain, xVal, xTest, yTrain, yVal, yTest):
1122     print("\nCOMENCING TRAINING OF SVM LINEAR\n")
1123
1124     s, acc, c = evaluateSVMLinear(xTrain, xVal, yTrain, yVal)
1125     accuracy = accuracy_score(yTest, s.predict(xTest)) * 100
1126     print('Accuracy over Test sample: ' + str(accuracy) + '%')
1127     return accuracy, acc, c
1128
1129 def createDataSets(x, y, m, trainPerc, valPerc, testPerc):
1130     if(trainPerc + valPerc + testPerc > 1.0):
1131         print("ERROR: Percentages given not valid")
1132         exit(-1)
1133
1134     print('Size of Training set: ' + str(int(trainPerc * m)))
1135     print('Size of Validation set: ' + str(int(valPerc * m)))
1136     print('Size of Teseting set: ' + str(int(testPerc * m)))
1137
1138     valPerc += trainPerc
1139     testPerc += valPerc
1140
1141     train = int(trainPerc * m)
1142     val = int(valPerc * m)
1143     test = int(testPerc * m)

```

```

1144
1145     xTrain = x[:train]
1146     yTrain = y[:train]
1147
1148     xVal = x[train:val]
1149     yVal = y[train:val]
1150
1151     xTest = x[val:test]
1152     yTest = y[val:test]
1153
1154     return xTrain, yTrain, xVal, yVal, xTest, yTest
1155
1156 def main():
1157     dataR = loadCSV("../Data/MushroomDataset/secondary_data_shuffled.csv")
1158
1159     #analyzeData(dataR)
1160
1161     data = dataR.to_numpy()
1162     y = data[:, 0]
1163     x = data[:, 1:]
1164
1165     m = y.shape[0]
1166     n = x.shape[1]
1167
1168     trainPerc = 0.2
1169     valPerc = 0.2
1170     testPerc = 0.2
1171
1172     # 0.2 0.2 0.2
1173     # Seconds taken for the evaluation: 8767.091146945953
1174     # Minutes taken for the evaluation: 146.11818578243256
1175     # Hours taken for the evaluation: 2.4353030963738758
1176
1177     xTrain, yTrain, xVal, yVal, xTest, yTest = createDataSets(x, y, m,
1178     trainPerc, valPerc, testPerc)
1179
1180     startTime = time.time()
1181
1182     logAcc, logAccVal, logCost, logL, logExp = evalLogistic(xTrain, xVal,
1183     xTest, yTrain, yVal, yTest)
1184
1185     nnHiddenLayer, nnAcc, nnAccVal, nnCost, nnReg, nnLaps = evalNeuronal(
1186     xTrain, xVal, xTest, yTrain, yVal, yTest, n)
1187
1188     svmLAcc, svmLAccVal, svmLC= evalSVMLinear(xTrain, xVal, xTest, yTrain,
1189     yVal, yTest)
1190
1191     svmAcc, svmAccVal, svmC, svmSig = evalSVM(xTrain, xVal, xTest, yTrain,
1192     yVal, yTest)
1193
1194     print('\nRESULTS OF LOGISTIC REGRESSION\n')
1195     print("Best accuracy: " + str(logAcc))
1196     print("Accuracy of train: " + str(logAccVal) + "%")

```

```

1192     print('Best cost: ' + str(logCost))
1193     print('Best lambda: ' + str(logL))
1194     print('Best exponent: ' + str(logExp))
1195
1196     print('\nRESULTS OF NEURONAL NETWORK\n')
1197     print('Best Accuracy: ' + str(nnAcc) + '%')
1198     print("Accuracy of train: " + str(nnAccVal) + "%")
1199     print('Best neurons in hidden layer: ' + str(nnHiddenLayer))
1200     print('Best cost: ' + str(nnCost))
1201     print('Best reg: ' + str(nnReg))
1202     print('Best laps: ' + str(nnLaps))
1203
1204     print('\nRESULTS OF SVM LINEAR\n')
1205     print("Best accuracy: " + str(svmLAcc))
1206     print("Accuracy of train: " + str(svmLAccVal) + "%")
1207     print("Best C: " + str(svmLC))
1208
1209     print('\nRESULTS OF SVM\n')
1210     print("Best accuracy: " + str(svmAcc))
1211     print("Accuracy of train: " + str(svmAccVal) + "%")
1212     print("Best C: " + str(svmC))
1213     print("Best Sigma: " + str(svmSig))
1214
1215     endTime = time.time()
1216     deltaTime = endTime - startTime
1217
1218     print()
1219     print('Seconds taken for the evaluation: ' + str(deltaTime))
1220     print('Minutes taken for the evaluation: ' + str(deltaTime/60))
1221     print('Hours taken for the evaluation: ' + str(deltaTime/3600))
1222
1223     plt.figure(figsize=(14, 10))
1224     plt.bar(np.arange(4), [logAcc, nnAcc, svmLAcc, svmAcc])
1225     plt.title('Accuracy comparasion between learning algorithms')
1226     plt.ylabel('Accuracy')
1227     xlabels = ['Logistic Regresion\nAccuracy: ' + str(logAcc) + '%', '
Neural Network\nAccuracy : ' + str(nnAcc) + '%', 'SVM Linear\nAccuracy:
' + str(svmLAcc) + '%', 'SVM Gaussain\nAccuracy: ' + str(svmAcc) + '%']
1228     plt.xticks(np.arange(4), xlabels)
1229     plt.savefig('../Results/FinalGraph.png', bbox_inches='tight')
1230     plt.close()
1231
1232 if __name__ == "__main__":
1233     main()

```

src/main.py