

Machine Learning (2016 Fall) HW2 Report

Spam Classification by Logistic Regression

Student : 沈怡廷 r05943001

Kaggle_ID : r05943001_Panda

A. Logistic Regression Function

```
81 """ logistic regression """
82 def logReg(X, y, X_test=None, y_test=None, lr=1e-7, batch=1, lamb=0,
83           epoch=10000, print_every=100):
84     # Initialize weight
85     W = np.random.randn(X.shape[1]) / X.shape[1] / X.shape[0]
86     L_train = []
87     A_train = []
88     L_test = []
89     A_test = []
90
91     # AdaGrad
92     G = np.zeros(W.shape)
93
94     for i in range(epoch):
95         # batch
96         b = 0
97         idx = []
98
99         for j in np.random.permutation(X.shape[0]):
100             idx.append(j)
101             b += 1
102
103         if b >= batch:
104             # prediction error
105             pred_err = sigmoid(X[idx].dot(W)) - y[idx]
106
107             # calculate gradient
108             grad_X = X[idx].T.dot(pred_err)
109             grad_regular = lamb * W * batch / X.shape[0]
110             grad = grad_X + grad_regular
111
112             # calculate weight
113             G += grad**2
114             W -= grad*lr/np.sqrt(G+1e-8)
115
116             # reset parameter
117             b = 0
118             idx = []
119
120         # prevent inf and nan
121         cross_Entropy = croEntropy(sigmoid(X.dot(W)), y)
122         cross_Entropy[np.isinf(cross_Entropy)] = 1000
123         cross_Entropy[np.isnan(cross_Entropy)] = 0
124
125         # calculate loss and accuracy
126         Loss_train = cross_Entropy.sum() + lamb*(W**2).sum()
127         Accuracy_train = calcAccuracy(X, W, y)
128
129         L_train.append(Loss_train)
130         A_train.append(Accuracy_train)
131
132     # testing
```

```
107         # calculate gradient
108         grad_X = X[idx].T.dot(pred_err)
109         grad_regular = lamb * W * batch / X.shape[0]
110         grad = grad_X + grad_regular
111
112         # calculate weight
113         G += grad**2
114         W -= grad*lr/np.sqrt(G+1e-8)
115
116         # reset parameter
117         b = 0
118         idx = []
119
120     # prevent inf and nan
121     cross_Entropy = croEntropy(sigmoid(X.dot(W)), y)
122     cross_Entropy[np.isinf(cross_Entropy)] = 1000
123     cross_Entropy[np.isnan(cross_Entropy)] = 0
124
125     # calculate loss and accuracy
126     Loss_train = cross_Entropy.sum() + lamb*(W**2).sum()
127     Accuracy_train = calcAccuracy(X, W, y)
128
129     L_train.append(Loss_train)
130     A_train.append(Accuracy_train)
131
132     # testing
133     if X_test is not None and y_test is not None:
134         cross_Entropy = croEntropy(sigmoid(X_test.dot(W)), y_test)
135         cross_Entropy[np.isinf(cross_Entropy)] = 1000
136         cross_Entropy[np.isnan(cross_Entropy)] = 0
137
138         Loss_test = cross_Entropy.sum() + lamb*(W**2).sum()
139         Accuracy_test = calcAccuracy(X_test, W, y_test)
140
141         L_test.append(Loss_test)
142         A_test.append(Accuracy_test)
143
144     # print out the progress
145     if i % print_every == 0:
146         if X_test is not None and y_test is not None:
147             print('\tePOCH: %d; loss: %.4f; Acc_train: %.4f; Acc_test: %.4f' %
148                   (i, Loss_train, Accuracy_train, Accuracy_test))
149         else:
150             print('\tePOCH: %d; loss: %.4f; Acc_train: %.4f' %
151                   (i, Loss_train, Accuracy_train))
152
153     print('\nfinal loss: %.4f' % L_train[-1])
154     print('final training accuracy: %.4f' % A_train[-1])
155     print('final testing accuracy: %.4f' % A_test[-1])
156
157     return W, L_train, A_train, L_test, A_test
158
```

B. Describe your another method, and which one is best

我實作的另外一個方法為「**probabilistic generative model**」，即上課所提到，假設資料分布呈現某種機率模型，利用 training data 的資料找出 average 及 variance，即可用來判斷 testing data。以下為其程式碼：

```

159 """ generative model """
160 def probGen(X, y):
161     data_class1 = [x for idx, x in enumerate(X) if y[idx]==1]
162     data_class2 = [x for idx, x in enumerate(X) if y[idx]==0]
163
164     # transform to numpy array
165     data_class1_array = np.matrix(data_class1)
166     data_class2_array = np.matrix(data_class2)
167
168     # delete answer column
169     data_class1_array = np.delete(data_class1_array, 57, 1)
170     data_class2_array = np.delete(data_class2_array, 57, 1)
171
172     # calculate number
173     N1 = data_class1_array.shape[0]
174     N2 = data_class2_array.shape[0]
175     N = N1 + N2
176
177     # calculate mean
178     u1 = np.mean(data_class1_array, axis=0)
179     u2 = np.mean(data_class2_array, axis=0)
180
181     # minus mean
182     data_class1_array_zero_mean = data_class1_array - u1
183     data_class2_array_zero_mean = data_class2_array - u2
184
185     # square
186     data_class1_square = data_class1_array_zero_mean.T*(data_class1_array_zero_mean)
187     data_class2_square = data_class2_array_zero_mean.T*(data_class2_array_zero_mean)
188
189     # variance
190     var1 = data_class1_square/N1
191     var2 = data_class2_square/N2
192     var = var1*N1/N + var2*N2/N
193
194     # avoid singular
195     var_mod = var + np.eye(var.shape[0])*1e-6
196     var_inv = np.linalg.inv(var_mod)
197
198     # weight
199     W = (u1 - u2)*(var_inv)
200     b = 0.5*u2*var_inv*(u2.T) - 0.5*u1*var_inv*(u1.T) + np.log(N1/N2)
201     W = np.concatenate((W, b), axis=1)
202
203     # accuracy
204     Accuracy = calcAccuracy(X, W.T, y)
205
206     return W.T, Accuracy
207

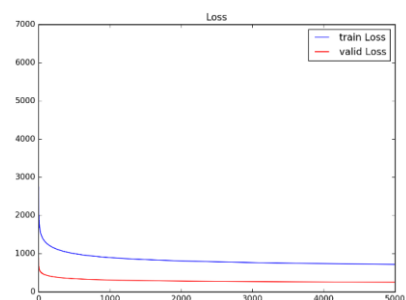
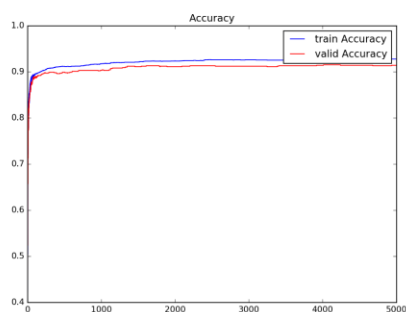
```

與 logistic regression 的比較如下圖所示，可以發現到，logistic regression 無論在 training data 或 testing data 上，表現均比較優異，主要原因或許是資料數量不夠多，機率模型有可能會有過多的錯誤判斷，導致效果較差。

	Training	Validation	Public Testing
Logistic	0.9280	0.9251	0.9333
Generative	0.8923	N/A	0.87667

C. TA depend on your other discussion and detail

a. Loss and accuracy during training



b. More about my logistic regression implementation

Initial weight and bias : random

Training set : random pick 3/4 of training data

Validation set : the rest 1/4 of training data

Optimization : AdaGrad

Batch size : whole training set (without SGD)

Learning rate : 0.01

Epoch : 5000

Regularization : None ($\lambda = 0$)