

Machine Learning 2016 Fall

hw2 – Spam Classification

b03902035 黃兆緯

1. Logistic Regression Function

```
18 """
19 Fit training data
20 params:
21 X:          Numpy array with shape (n_samples, n_features)
22 y:          Numpy array with shape (n_samples, 1)
23 """
24 def fit(self, X, y):
25     if len(y.shape) == 1:
26         y = np.array([[num] for num in y])
27
28     if len(X) == 0 or len(X) != len(y):
29         print('length of X and y not equal')
30         return
31
32     init_w = np.array([-1.43] + [0.01 for i in range(len(X[0])-1)])
33
34     self.w = init_w
35     w_acc = np.array([1e-5 for i in range(len(X[0]))])
36
37     # run several passes
38     for i in range(self.iteration+1):
39         # compute gradients of w and sum over all training data
40         wgrad = self.compute_grad(X, y)
41
42         # compute summation of past gradients, for adagrad
43         w_acc = w_acc + wgrad**2
44
45         # update parameters, using adagrad
46         if self.adagrad:
47             self.w = self.w - self.rate*(1.0/np.sqrt(w_acc))*wgrad
48         else:
49             self.w = self.w - self.rate*wgrad
50
51         # compute and print training error/validation error every 1000 pass
52         if i % 1000 == 0:
53             # E_in(training set error)
54             train_ans = (self.predict(X)>0.5).astype(float)
55             train_error = np.mean(np.absolute(y[:,0] - train_ans))
56             print('iteration %d, \t\ttrain error: %f' % (i, train_error))
```

```
58 """
59 Predict output values for given data
60 params:
61 x:          features to predict
62
63 returns:
64 y:          output value
65 """
66 def predict(self, x):
67     return expit(np.dot(x, self.w))
68
69 """
70 Compute gradients given data
71 params:
72 X:          training features
73 y:          ground truth
74
75 returns:
76 wgrad:      gradients for weights
77 """
78 def compute_grad(self, X, y):
79     delta = y[:,0] - self.predict(X)
80     wgrad = -delta.dot(X)
81     return wgrad
```

(1) predict

$$y = \text{sigmoid}(\sum (w_i \cdot x_i))$$

(2) compute grad

計算 w 的 gradient, $grad_i = \sum (-x_{ji} \cdot (y_j - \bar{y}_j))$

(3) update

使用 gradient descent(沒有使用 SGD, 一次計算全部 sample), 並用 Adagrad 加強 gradient descent 的功能

2. method 2: Multi-layer Perceptrons

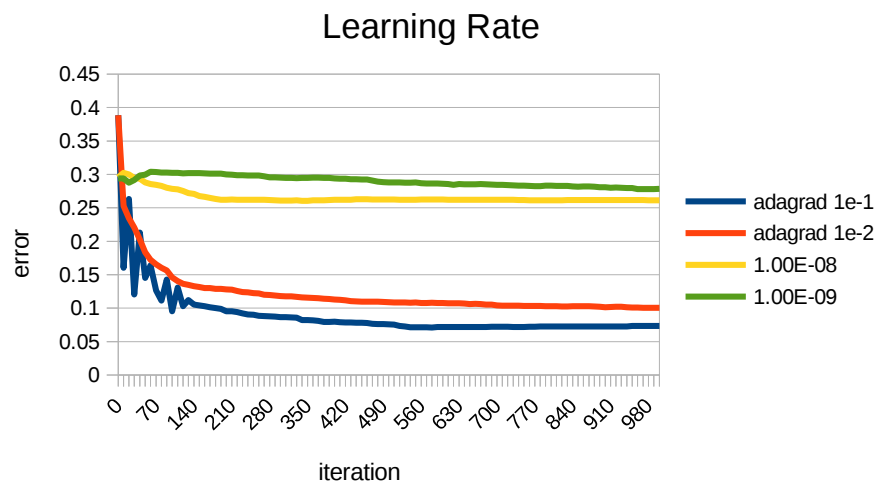
```
49 def fit(self, X, y):
50     self.initialize_coef(X, y)
51
52     num_examples = len(X)
53
54     if len(y.shape) == 1:
55         y = np.array([[num] for num in y])
56
57     # when should we stop?
58     last_error = 1000000.0
59     consecutive_increase = 0
60
61     for num_iter in range(self.max_passes):
62         a = self.forward(X)
63
64         error = loss(a[-1], y)
65         values = np.sum(np.array([np.dot(s.ravel(), s.ravel()) for s in self.W]))
66         error_reg = error + (0.5 * self.reg) * values / num_examples
67
68         if num_iter % 50 == 0:
69             train_error = np.mean(np.absolute((a[-1]>0.5).astype(float)-y))
70             #print('\t', num_iter, '\t', error_reg, '\t', train_error)
71
72             if error_reg > last_error:
73                 consecutive_increase += 1
74             elif error_reg < last_error:
75                 consecutive_increase = 0
76
77             last_error = error_reg
78
79             if consecutive_increase == 2:
80                 #print('consecutive increasing, stop training')
81                 break
82
83     deltas = [np.array([0] for i in range(len(self.dims)))]
84     dw, db = [0 for i in range(len(self.W))], [0 for i in range(len(self.b))]
85
86     last = len(self.dims)-2
87
88     deltas = [np.array([0] for i in range(len(self.dims)))]
89     dw, db = [0 for i in range(len(self.W))], [0 for i in range(len(self.b))]
90
91     last = len(self.dims)-2
92     deltas[last] = a[-1] - y
93
94     for i in range(len(self.dims)-2, 0, -1):
95         deltas[i-1] = np.dot(deltas[i], self.W[i].T)
96         derivatives(a[i], deltas[i-1], self.activation_type[i-1])
97         dW[i-1] = (np.dot(a[i-1].T, deltas[i-1]) + self.reg * self.W[i-1]) / num_examples
98         db[i-1] = np.mean(deltas[i-1], 0)
99
100     grads = dW + db
101     self.optimizer.update_params(grads)
102
103 def forward(self, X):
104     z, a = [0 for i in range(len(self.dims))], [X]+[0 for i in range(len(self.dims)-1)]
105
106     for i in range(len(self.hidden_layer_size)+1):
107         z[i+1] = a[i].dot(self.W[i]) + self.b[i]
108         if (i + 1) != len(self.hidden_layer_size)+1:
109             a[i+1] = activation(z[i+1], self.activation_type[i])
110         else:
111             a[i+1] = z[i+1]
112
113     return a
```

第二個方法我選擇實作 DNN，主要的技巧如下：

- (1) Forward pass
從 input layer 開始，計算下一層 neuron 的 input $input = activation(b + output \cdot w)$ ，直到 output layer
- (2) Back propagation
從 output layer 開始，先計算 $delta_{last} = y - \bar{y}$ ，接著利用每一層間的 w 計算前一層的結果在 $delta$ 上佔的權重，再利用這個權重更新 w
- (3) update
使用 ADAM optimizer，參考 sklearn 的實作

3. Discussion

(1) Logistic Regression Learning Rate



沒有使用 adagrad 時，很明顯 w 更新地很慢，甚至會卡住無法下降，使用 adagrad 可以很好地解決這個問題。learning rate 愈大時 error 下降愈快，但 rate 太大時可能會有波動(如 adagrad 1e-1)。

(2) Neural Net Hidden-layer Size

size	(30,1)	(30,30,1)	(30,30,2)	(30,30,5)	(30,30,10)	(30,20,5)	(50,30,5)
accuracy	0.928	0.938	0.943	0.945	0.945	0.946	0.945

針對 DNN 的 hidden layer size 做了一些實驗，這裡的 accuracy 是拿 training data 做 5-fold cross validation 的平均 accuracy。可以看到三層 hidden layer 表現明顯較兩層好，而每層 hidden layer 的維度則沒有太大影響。考慮 training 速度之後，選用(30, 30, 2)來做其他實驗。

(3) Feature Selection

由於在 hw1 中，feature selection 是很重要的一個環節，因此在這次作業也嘗試選擇 feature。

- (a) Chi-Square
對於每個 feature 和 class 做 chi-square test，移除值最小的兩維(剛好是最後兩維，平均大寫序列長度和總大寫序列長度)，重新 train 後發現 accuracy 變低。
 - (b) Pearson correlation
對於每個 feature 和 class 算 pearson correlation，若 feature 和 class 具有線性關係，correlation 應該較大，因此移除 correlation 最小的數維，但實驗後也發現 accuracy 降低。
- 經過以上兩個實驗後，找不到比較好取 feature 的方法，因此都使用所有維度進行實驗。

(4) Comparison between different models

model	Logistic Regression	MLP	Random Forest n_trees=9	Random Forest n_trees=99	Gradient Boosting n_trees=9	Gradient Boosting n_trees=99
accuracy	0.928	0.945	0.943	0.952	0.928	0.950

除了 DNN，這次還實作了 Random Forest，另外用 sklearn 的 GradientBoostingClassifier 一起做了實驗。除了 Logistic Regression 之外的幾個 model 都有差不多的 accuracy，最後最好的 model 就是用 DNN model + Random Forest model 進行投票得來的。