

# 第一周机器学习

本文公式显示需要使用Mathjax，然后令人悲伤的是github不支持Mathjax 您可以将这篇md文件pull下来，使用您本地的markdown解析器解析 没有必要在公示显示上浪费时间，您也可以下载我本地生成的html用浏览器打开即可 或者您也可以下载我上传到github上的pdf [Mathjax开源项目地址](#)

## 绪论

### 机器学习简介

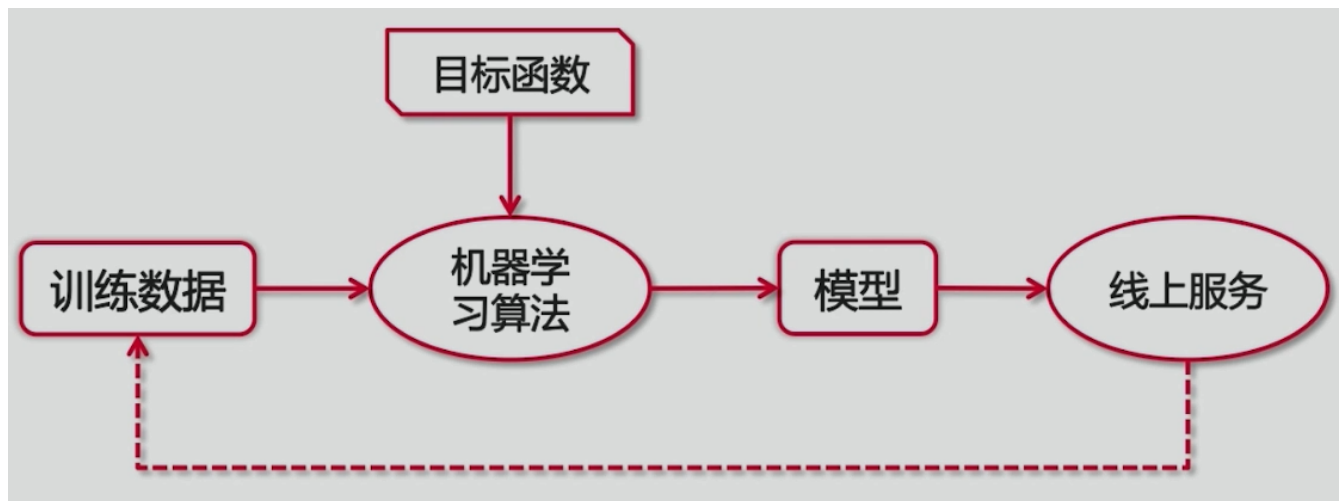
机器学习是一种将无序数据转换为价值的方法。

机​​器学习的价值-从数据中抽取规律，并用来预测未来

### 机器学习应用举例

- 分类问题-图像识别、垃圾邮件识别
- 回归问题-股价预测、房价预测
- 排序问题-点击率预估、推荐
- 生成问题-图像生成、图像风格转换、图像文字描述生成

### 机器学习的应用流程



### 机器学习岗位职责

- 数据处理(采集+去噪)
- 模型训练(特征+模型)
- 模型评估与优化(MSE、F1-score、AUC+调参)
- 模型应用(A/B测试)

## 线性回归(Linear regression)

cost function:

来源于假设误差  $\epsilon_i$  服从正态分布，然后对参数  $\theta$  进行极大似然估计，经过运算后得出  $J(\theta)$  取最小时，似然函数最大，从而推出这个式子。

此外，对这个  $J(\theta)$  求偏导，令其偏导数为0（这里涉及到矩阵偏导数计算），即可得到正规方程(normal equation)。

## 梯度下降法

### 分类

- mini-batch
- batch
- random
- SGD(动量梯度下降，有助于解决局部最值和鞍点问题)

### Code

参考代码，自己就一些细节进行优化

(<https://www.cnblogs.com/focusonepoint/p/6394339.html>)

```
1  #!/usr/bin/python
2  # -*- coding: UTF-8 -*-
3
4  import numpy as np
5  from numpy import linalg
6
7
8
9  def gradientDescent(x, y, theta, m, alpha, maxIteration):
10     '''
11     使用批处理梯度下降算法计算theta
12     '''
13     # 得到x的转置
14     # 即 x的第一行为x1 第二行为x2 第三行全部初始化为1
15     xTrains = x.transpose()
16     # theta 是一个列向量
17     for i in xrange(0,maxIteration):
18         # x矩阵(10*3)与theta(3*1)矩阵相乘
19         # hypothesis(i) = x1(i)*theta1(i) + x2(i)*theta2(2) + 1*theta0(i)
20         hypothesis = np.dot(x, theta)
21         # 作差
22         loss = hypothesis - y
23         # 当loss的范数在我们的误差允许范围内 就停止循环
24         if (linalg.norm(loss) < 1e-5):
25             break
26         # xTrains (3*10) * loss(10*1) = gradient(3*1)
27         # 计算代价函数
28         gradient = (1.0/m) * np.dot(xTrains, loss)
29         theta = theta - alpha * gradient
30     print('the number of iteration is %d' % i);
31     return theta
32
```

```

33
34 # define the prepared 训练集
35 # the meaning of column : x1,x2,y
36 dataSet = np.array([
37     [1.1,1.5,2.5],
38     [1.3,1.9,3.2],
39     [1.5,2.3,3.9],
40     [1.7,2.7,4.6],
41     [1.9,3.1,5.3],
42     [2.1,3.5,6.0],
43     [2.3,3.9,6.7],
44     [2.5,4.3,7.4],
45     [2.7,4.7,8.1],
46     [2.9,5.1,8.8],
47 ])
48
49
50 # print(dataSet)
51 m,n = np.shape(dataSet)
52 # print(m,n)
53 trainData = np.ones((m,n))
54 # 截取dataSet的前N-1列
55 trainData[:, :-1] = dataSet[:, :-1]
56 # 获取dataSet的最后一列
57 trainLabel = dataSet[:, -1]
58
59 # print(m,n)
60 theta = np.ones(n)
61 # print(theta)
62 alpha = 0.001
63
64 # the max time of iteration 这个值定义的尽量大(考虑计算机的性能)
65 maxIteration = 10000000
66 theta = gradientDescent(trainData, trainLabel, theta, m, alpha, maxIteration)
67 print('the value of theta is:')
68 print(np.round(theta,2))
69
70
71 # a test for the algorithm
72 x = np.array([
73     [3.1, 5.5],
74     [3.3, 5.9],
75     [3.5, 6.3],
76     [3.7, 6.7],
77     [3.9, 7.1]
78 ])
79
80
81 # define a predict function used to test
82 def predict(x, theta):
83     m, n = np.shape(x)
84     xTest = np.ones((m, n+1))
85     xTest[:, :-1] = x

```

```

86     yPre = np.dot(xTest, theta)
87     return yPre
88
89     print('the predicted value is')
90     yP = predict(x, theta)
91     print(np.round(yP, 2))
92
93

```

运行结果

```

1  the number of iteration is 114575
2  the value of theta is:
3  [ 0.71  1.39 -0.38]
4  the predicted value is
5  [ 9.5 10.2 10.9 11.6 12.3]
6  [Finished in 2.2s]

```

## 优化技巧

- Feature Scaling（特征缩放）
  - 归一化
    1. 线性归一化
      - $\{x\}' = \frac{x - \min(x)}{\max(x) - \min(x)}$
    2. 标准差归一化
      - $x^* = \frac{x - \overline{x}}{s}$
    3. 非线性归一化
- 多项式回归
  - $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$
  - $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$
  - $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 \sqrt{x}$
  - 上面的举例只是为了说明， $x_i$ 的取值可以不是 $x$ 的一次多项式，但是这里要注意的是特征缩放在这里显得尤为重要
- $\alpha$ 选取技巧
  - 如果 $J(\theta)$ 的值随着 $\theta$ 的取值单调递增或者出现震荡，那么 $\alpha$ 应该选的小一点

## Normal Equation（正规方程法）

### 思想

微积分思想：求导后令导数为零解方程可以求出极值点 $\theta$

对于 $\theta$ 是一个 $n$ 维向量的情况，可以利用多元函数取极值的必要条件，即偏导数为0

### 结论

### Note

1. No need to do feature scaling
2. 只适用于线性模型，不适合逻辑回归模型等其他模型
3. the pseudo inverse of matrix
  - o redundant features (x中存在线性相关的量)
  - o too many features (eg.  $m \leq n$  数据个数小于特征参数)

## Code

这里我使用上一个梯度下降法的例子作为对比,采用相同的数据对比运行结果

```
1  #!/usr/bin/python
2  # -*- coding: UTF-8 -*-
3
4  import numpy as np
5
6  def normalEqation(x, y):
7      '''
8      使用正规方程法算法计算theta
9      '''
10     # 得到x的转置
11     xTrains = x.transpose()
12     m,n = np.shape(x)
13     # theta 为 n维列向量
14     theta = np.linalg.pinv(np.dot(xTrains,x))
15     theta = np.dot(theta,xTrains)
16     theta = np.dot(theta,y)
17     return theta
18
19
20 # define the prepared 训练集
21 # the meaning of column : x1,x2,y
22 dataSet = np.array([
23     [1.1,1.5,2.5],
24     [1.3,1.9,3.2],
25     [1.5,2.3,3.9],
26     [1.7,2.7,4.6],
27     [1.9,3.1,5.3],
28     [2.1,3.5,6.0],
29     [2.3,3.9,6.7],
30     [2.5,4.3,7.4],
31     [2.7,4.7,8.1],
32     [2.9,5.1,8.8],
33 ])
34
35
36 # print(dataSet)
37 m,n = np.shape(dataSet)
38 # print(m,n)
39 trainData = np.ones((m,n))
40 # 截取dataSet的前N-1列
41 trainData[:, :-1] = dataSet[:, :-1]
```

```

42 # 获取dataSet的最后一列
43 trainLabel = dataSet[:, -1]
44
45
46
47 theta = normalEqation(trainData, trainLabel)
48 print('thec value of theta is:')
49 print(np.round(theta,2))
50
51
52 # a test for the algorithm
53 x = np.array([
54     [3.1, 5.5],
55     [3.3, 5.9],
56     [3.5, 6.3],
57     [3.7, 6.7],
58     [3.9, 7.1]
59 ])
60
61
62 # define a predict function used to test
63 def predict(x, theta):
64     m, n = np.shape(x)
65     xTest = np.ones((m, n+1))
66     xTest[:, :-1] = x
67     yPre = np.dot(xTest, theta)
68     return yPre
69
70 print('the predicted value is')
71 yP = predict(x, theta)
72 print(np.round(yP,2))

```

运行结果：

```

1 thec value of theta is:
2 [ 0.61  1.45 -0.34]
3 the predicted value is
4 [ 9.5 10.2 10.9 11.6 12.3]
5 [Finished in 0.2s]

```

由此可以知道，在特征矩阵维度不是太大情况下，对于线性回归模型，**normal equation** 是一个优先选用的方法。

## Logistic Regression

### logistic function

由于线性回归的假设函数不再适用于分类问题，因此我们需要一个函数来应用于分类问题的拟合。一般来说，回归不用在分类问题上，因为回归是连续型模型，而且受噪声影响比较大。如果非要应用进入，可以使用logistic回归。

我们可以使用logistic regression解决分类问题，Logistic回归是二分类任务的首选方法，下面讨论二分类的问题。

logistic function(sigmoid function):

这里

含义是在 $x$ 已知条件下，给定参数 $\theta$ ，事件 $y=1$ 发生的概率

logistic回归本质上是线性回归，只是在特征到结果的映射中加入了一层函数映射，即先把特征线性求和，然后使用函数 $g(z)$ 将最为假设函数来预测。 $g(z)$ 可以将连续值映射到0和1上。

对 $g(z)$ 的解释：将任意的输入映射到 $[0,1]$ 区间上，我们在线性回归中可以得到一个预测值，再将该值映射到Sigmoid函数，这样我们就实现了由值到概率的转换，也就是分类任务。

**Note:**

当 $y$ 等于1时，假设函数计算出的概率应该大于0.5，即 $\theta$ 的转置乘以 $x$ 需要大于等于0 当 $y$ 等于0时，假设函数计算出的概率应该小于0.5，即 $\theta$ 的转置乘以 $x$ 需要小于0 另外需要注意的是阈值0.5在一些情况下是可以改变的，从而获得我们所希望的特征

## cost function

这里我们将 cost function 定义为

例如， $y = 1$ 时  $h_{\theta}(x) \rightarrow 1$ ,  $\text{cost} = 0$  表示误差很小。此时，若  $h_{\theta}(x) \rightarrow 0$ ,  $\text{cost} \rightarrow \infty$  表示误差很大

## Simple Classification(简单分类算法)

Note:  $y=0$  or  $1$

这里对cost function进行优化，表示为：

这里的cost function 实际上也是由对 $\theta$ 的极大似然估计推导出来的。

by the way, remind:

ok,接着我们对  $J_{\theta}(x)$  计算偏微分

其中

将上面式子代入  $\frac{\partial}{\partial \theta_j} J(\theta)$  得

这里我们推出一个重要的结论

**Note:** 这里的  $h_{\theta}(x^{(i)})$  与 线性回归模型中的  $h_{\theta}(x^{(i)})$  定义不一样，尽管计算出来的  $\frac{\partial}{\partial \theta_j} J(\theta)$  形式相同

## Code

```
1  #!/usr/bin/python
2  # -*- coding: UTF-8 -*-
3
4  '''
5  本例程是根据学生两门课的成绩判断是否录取
6  '''
7
8
9  import numpy as np
```

```
10 import pandas as pd
11 import matplotlib.pyplot as plt
12
13 dataStr = '''
14 34.62365962451697,78.0246928153624,0
15 30.28671076822607,43.89499752400101,0
16 35.84740876993872,72.90219802708364,0
17 60.18259938620976,86.30855209546826,1
18 79.0327360507101,75.3443764369103,1
19 45.08327747668339,56.3163717815305,0
20 61.10666453684766,96.51142588489624,1
21 75.02474556738889,46.55401354116538,1
22 76.09878670226257,87.42056971926803,1
23 84.43281996120035,43.53339331072109,1
24 95.86155507093572,38.22527805795094,0
25 75.01365838958247,30.60326323428011,0
26 82.30705337399482,76.48196330235604,1
27 69.36458875970939,97.71869196188608,1
28 39.53833914367223,76.03681085115882,0
29 53.9710521485623,89.20735013750205,1
30 69.07014406283025,52.74046973016765,1
31 67.94685547711617,46.67857410673128,0
32 70.66150955499435,92.92713789364831,1
33 76.97878372747498,47.57596364975532,1
34 67.37202754570876,42.83843832029179,0
35 89.67677575072079,65.79936592745237,1
36 50.534788289883,48.85581152764205,0
37 34.21206097786789,44.20952859866288,0
38 77.9240914545704,68.9723599933059,1
39 62.27101367004632,69.95445795447587,1
40 80.1901807509566,44.82162893218353,1
41 93.114388797442,38.80067033713209,0
42 61.83020602312595,50.25610789244621,0
43 38.78580379679423,64.99568095539578,0
44 61.379289447425,72.80788731317097,1
45 85.40451939411645,57.05198397627122,1
46 52.10797973193984,63.12762376881715,0
47 52.04540476831827,69.43286012045222,1
48 40.23689373545111,71.16774802184875,0
49 54.63510555424817,52.21388588061123,0
50 33.91550010906887,98.86943574220611,0
51 64.17698887494485,80.90806058670817,1
52 74.78925295941542,41.57341522824434,0
53 34.1836400264419,75.2377203360134,0
54 83.90239366249155,56.30804621605327,1
55 51.54772026906181,46.85629026349976,0
56 94.44336776917852,65.56892160559052,1
57 82.36875375713919,40.61825515970618,0
58 51.04775177128865,45.82270145776001,0
59 62.22267576120188,52.06099194836679,0
60 77.19303492601364,70.45820000180959,1
61 97.77159928000232,86.7278223300282,1
62 62.07306379667647,96.76882412413983,1
```



63 91.56497449807442,88.69629254546599,1  
64 79.94481794066932,74.16311935043758,1  
65 99.2725269292572,60.99903099844988,1  
66 90.54671411399852,43.39060180650027,1  
67 34.52451385320009,60.39634245837173,0  
68 50.2864961189907,49.80453881323059,0  
69 49.58667721632031,59.80895099453265,0  
70 97.64563396007767,68.86157272420604,1  
71 32.57720016809309,95.59854761387875,0  
72 74.24869136721598,69.82457122657193,1  
73 71.79646205863379,78.45356224515052,1  
74 75.3956114656803,85.75993667331619,1  
75 35.28611281526193,47.02051394723416,0  
76 56.25381749711624,39.26147251058019,0  
77 30.05882244669796,49.59297386723685,0  
78 44.66826172480893,66.45008614558913,0  
79 66.56089447242954,41.09209807936973,0  
80 40.45755098375164,97.53518548909936,1  
81 49.07256321908844,51.88321182073966,0  
82 80.27957401466998,92.11606081344084,1  
83 66.74671856944039,60.99139402740988,1  
84 32.72283304060323,43.30717306430063,0  
85 64.0393204150601,78.03168802018232,1  
86 72.34649422579923,96.22759296761404,1  
87 60.45788573918959,73.09499809758037,1  
88 58.84095621726802,75.85844831279042,1  
89 99.82785779692128,72.36925193383885,1  
90 47.26426910848174,88.47586499559782,1  
91 50.45815980285988,75.80985952982456,1  
92 60.45555629271532,42.50840943572217,0  
93 82.22666157785568,42.71987853716458,0  
94 88.9138964166533,69.80378889835472,1  
95 94.83450672430196,45.69430680250754,1  
96 67.31925746917527,66.58935317747915,1  
97 57.23870631569862,59.51428198012956,1  
98 80.36675600171273,90.96014789746954,1  
99 68.46852178591112,85.59430710452014,1  
100 42.0754545384731,78.84478600148043,0  
101 75.47770200533905,90.42453899753964,1  
102 78.63542434898018,96.64742716885644,1  
103 52.34800398794107,60.76950525602592,0  
104 94.09433112516793,77.15910509073893,1  
105 90.44855097096364,87.50879176484702,1  
106 55.48216114069585,35.57070347228866,0  
107 74.49269241843041,84.84513684930135,1  
108 89.84580670720979,45.35828361091658,1  
109 83.48916274498238,48.38028579728175,1  
110 42.2617008099817,87.10385094025457,1  
111 99.31500880510394,68.77540947206617,1  
112 55.34001756003703,64.9319380069486,1  
113 74.77589300092767,89.52981289513276,1  
114 '''  
115

```

116 tmpdataList = dataStr.split()
117 dataList = []
118 for data in tmpdataList:
119     data = data.split(',')
120     dataList.append(data)
121 del tmpdataList
122
123 # define the prepared 训练集
124 # the meaning of column : x1,x2,y
125 dataSet = np.array(dataList)
126 dataSet = dataSet.astype(np.float64)
127
128 def shuffleData(dataSet):
129     # 打乱数据
130     np.random.shuffle(dataSet)
131     m,n = np.shape(dataSet)
132
133     trainData = np.ones((m,n))
134     trainData[:, :-1] = dataSet[:, :-1]
135     # 获取dataSet的最后一列 并 强制类型转换
136     trainLabel = dataSet[:, -1]
137     return trainData, trainLabel
138
139
140 # 这里我们使用matplotlib先看一下数据
141 negativeData = dataSet[dataSet[:, -1] == 0.0]
142 positiveData = dataSet[dataSet[:, -1] == 1.0]
143 trainLabel = dataSet[:, -1].astype(np.float64)
144
145
146 fig, ax = plt.subplots(figsize=(10, 5))
147 ax.scatter(positiveData[:, 0], positiveData[:, 1], s = 30, c = 'b', marker = 'o', label =
    'Admitted')
148 ax.scatter(negativeData[:, 0], negativeData[:, 1], s = 30, c = 'r', marker = 'x', label =
    'Not Admitted')
149 ax.legend()
150 ax.set_xlabel('Exam 1 Score')
151 ax.set_ylabel('Exam 2 Score')
152 plt.show()
153
154 # 下面是逻辑回归算法
155 def sigmoid(z):
156     return (1.0 / (1.0 + np.exp(-z)))
157
158
159 def model(X, theta):
160     return sigmoid(np.dot(X, theta))
161
162 # x2 x1 x0
163 # res = model(trainData, theta)
164 def cost_function(X, y, theta):
165     h_x = model(X, theta)
166     left = -y*np.log(h_x)

```

```

167     right = (1-y)*np.log(1-h_x)
168     return np.sum(left - right) / (len(X))
169
170 # x = cost_function(trainData,trainLabel,theta)
171 def gradient(X,y,theta):
172     grad = np.zeros(theta.shape)
173     error = (model(X,theta) - y).ravel()
174     for j in xrange(len(theta.ravel())):
175         term = np.multiply(error, X[:,j])
176         grad[j] = np.sum(term) / len(X)
177     return grad
178 # 3种梯度下降方法 1.批处理 2.小批处理 3.随机处理
179 # 数据量较小, 直接批处理即可
180 def batchGradientDescent(dataSet, alpha, maxIteration, thresh):
181     X,y = shuffleData(dataSet)
182     m,n = np.shape(X)
183     k = 1.0 / m
184     theta = np.zeros((n,))
185
186     trainX = X.transpose()
187     for i in xrange(0,maxIteration):
188         error = model(X, theta) - y
189         _gradient = k * np.dot(trainX, error)
190         if (np.linalg.norm(_gradient) < thresh[0]):
191             print('hit thresh1')
192             break
193         # print(gradient(X,y,theta))
194         # print(_gradient)
195         cost1 = cost_function(X,y,theta)
196         theta = theta - alpha * _gradient
197         cost2 = cost_function(X,y,theta)
198         if abs(cost2 - cost1) < thresh[1]:
199             print('hit thresh2')
200             break
201         # print(theta)
202         print('the number of iteration is %d' % (i+1))
203         # print(error)
204         return theta
205
206 # theta = batchGradientDescent(dataSet,alpha =0.001,maxIteration = 1000000,thresh =
(1e-6,1e-6))
207 # print(theta)
208
209 '''
210 hit thresh2
211 the number of iteration is 109902
212 [ 0.04771429  0.04072397 -5.13364014]
213
214 这个数据说明当迭代次数为110000次时, cost function下降就跟缓慢了
215 '''
216
217 theta = batchGradientDescent(dataSet,alpha =0.001,maxIteration = 1000000,thresh =
(0.05,1e-6))

```

```

218 print(theta)
219 # theta = batchGradientDescent(dataSet,alpha =0.001,maxIteration = 1000000,thresh =
    (1e-6,1e-6))
220 # print(theta)
221
222 '''
223 hit thresh1
224 the number of iteration is 40046
225 [ 0.02721656  0.01899417 -2.37028409]
226 [Finished in 8.2s]
227 按照梯度下降停止大概需要40000次迭代
228 '''

```

这里实际上，如果数据经过预处理以及miniBatch后获得的数据精度比较高

## Advanced optimization

Optimization algorithms:

- Gradient descent
- Conjugate gradient
- BFGS
- L-BFGS 后面三种算法不需要给出学习率  $\alpha$ ，且运算速度较快，但是算法较为复杂，选修。

## 多类别处理

遇到 $y$ 的取值不仅仅是0,1情况时，可以将一类与其余类化为两种模型，然后用划分两类的分类算法计算出 $h(x)$ ，最后每一类都对应一个 $h(x)$ ，训练出模型后，判断  $\max h_{\theta}(x)$  对应的类即为最后输出。

## 关于机器学习的一些概念补充

### 下采样与上采样

下采样，对于一个不均衡的数据，让目标值(如0和1分类)中的样本数据量相同，且以数据量少的一方的样本数量为准。

上采样就是以数据量多的一方的样本数量为标准，把样本数量较少的类的样本数量生成和样本数量多的一方相同，称为上采样。

### 交叉验证

交叉验证的基本思想是把在某种意义下将原始数据(dataset)进行分组,一部分做为训练集(train set),另一部分做为验证集(validation set or test set),首先用训练集对分类器进行训练,再利用验证集来测试训练得到的模型(model),以此来做为评价分类器的性能指标。

### 二分类模型评估方法

以正例（恐怖分子）的识别为例子

真正例（True Positive, TP）：预测值和真实值都为1 假正例（False Positive, FP）：预测值为1，真实值为0(去真)  
 真负例（True Negative, TN）:预测值与真实值都为0 假负例（False Negative, FN）：预测值为0，真实值为1(存伪)

## 召回率（也叫查全率）

正确判为恐怖分子占实际所有恐怖分子的比例。在某些情况中，我们也许需要以牺牲另一个指标为代价来最大化精度或者召回率。比如检测癌症

## 精确度(precision,也叫查准率)

在所有判为恐怖分子中，真正的恐怖分子的比例。

## 准确率（accuracy）

## 正则化(Regularization)

## 欠拟合(underfitting)和过拟合(overfitting)

### How to addressing overfitting

1. Reduce number of features

2. Regularization

- keep all the feature, but reduce magnitude/values of feature.  
it works well when we have a lot of features, each of which contributes a bit to predicting y.

3. Regularization used in linear Regression

- $J(\theta) = \frac{1}{2m} [\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2]$

$\lambda$  称为 regularization parameter Note: 加上  $\theta^2$  是一种形式，有时也可以选择加上  $|\theta|$

3.1 Gradient descent

- 其中

3.2 Normal Equation

4. Regularization used in logistic Regression

## Neural networks(神经网络)

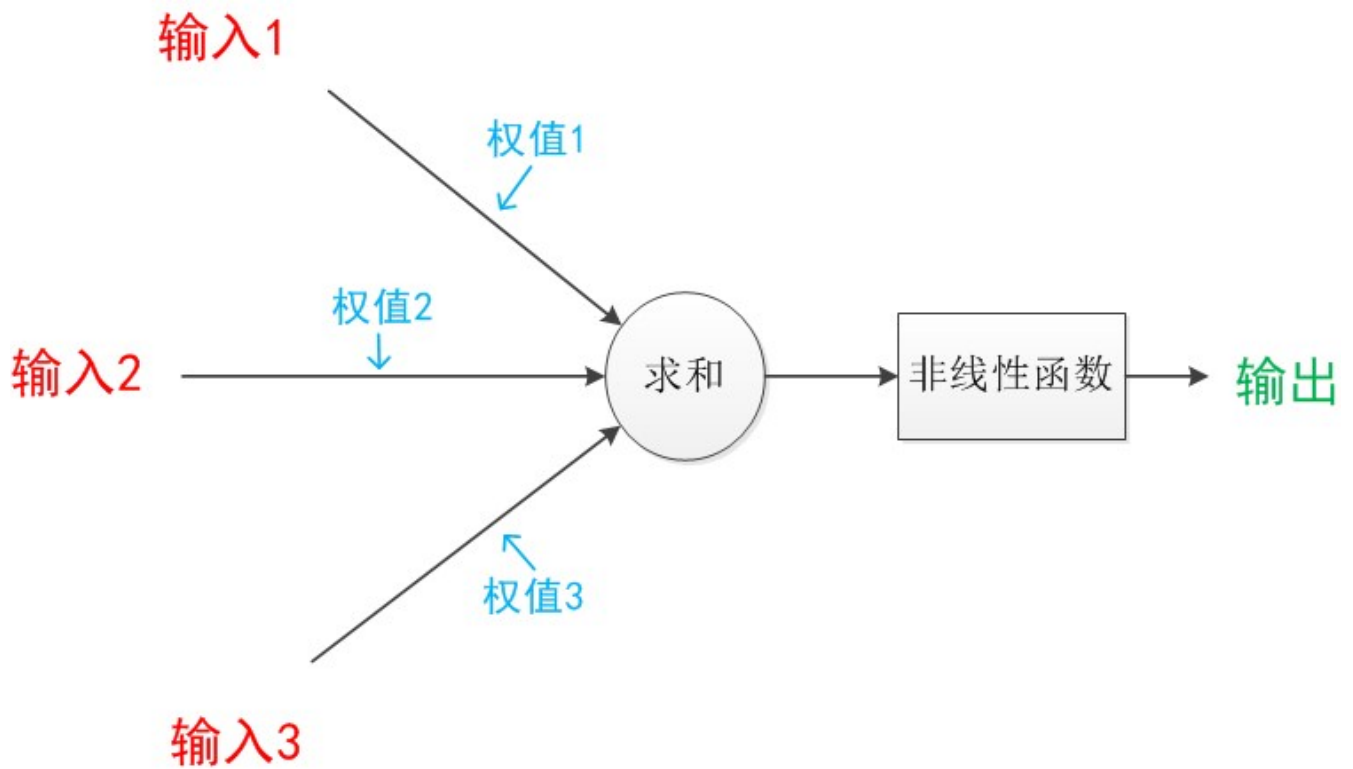
### Typeical Application（应用领域）

Example	Principle
Ad.userinfo	Online Advertising(Standard NN)
Image	Phototapping(CNN convolutional neural network)
Audio	Speech recognition(RNN recurrent neural network)
Machine translation	RNN
Autonomous driving	hybrid neural network + custom neural network

### Concepts

- Structured Data
  - Data in the database(have rows and cols)
  - 一般是离散的、有组织结构的
- Unstructured Data
  - Audio、Image、Text
  - 一般是连续的、无组织结构的

## 神经元

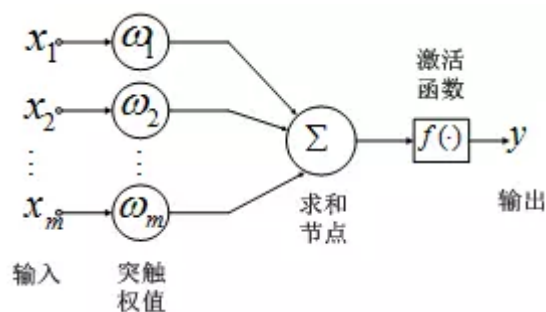
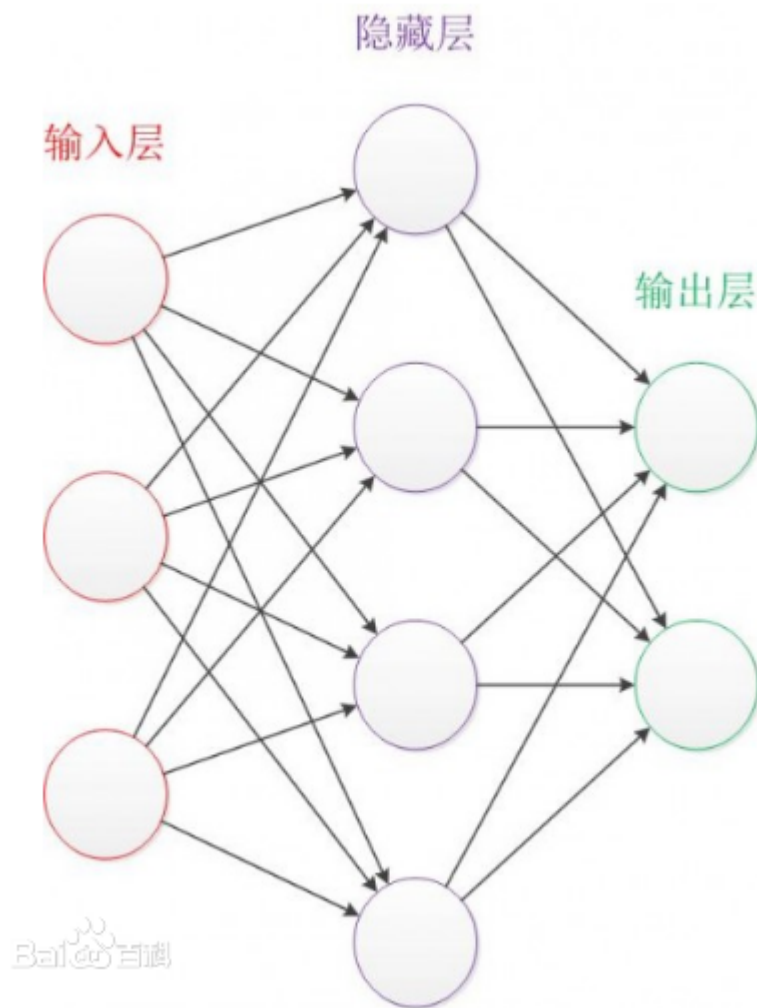


## Layer

神经网络是分层的  
一般来说

**Layer 1:** Input Layer

**Layer 2~N-1:** Hidden Layer **Layer N:** Output Layer



Note:

如果没有隐藏层，只有输入层和输出层，那么我们将这种神经网络称为“感知器”（Perceptron）。在“感知器”中，有两个层次。分别是输入层和输出层。输入层里的“输入单元”只负责传输数据，不做计算。输出层里的“输出单元”则需要对前面一层的输入进行计算。感知器只能做简单的线性分类任务，对XOR（异或）这样的简单分类任务无法解决。

Definitions:

1.  $a_i^{(j)}$  : "activation" of unit  $i$  in layer  $j$
2.  $\theta^{(j)}$  : matrix of weights controlling function mapping from layer  $j$  to layer  $j+1$

Examples:

1. 输出仅有一个的神经网络

## Forward propagation

令

则有

以上过程称为 **Forward propagation**

if network has  $s_j$  units in layer  $j$ ,  $s_{j+1}$  units in layer  $j+1$ , then  $\theta^{(j)}$  will be of dimension  $s_{j+1} \times s_j$

## Multi-class classification

若是表示多个输出，那么  $h_{\theta}(x)$  维度将大于1，变成一个向量矩阵，这个时候输出也就变成了多为

## cost function

对于

这  $m$  个样本数据训练出来的神经网络来说，我们定义：

$L$  = total number of layers in network

$s_l$  = no. of units(not counting bias unit) in layer  $l$

我们类比 **logistic regression** 的  $J(\theta)$

Neural network:

那么在神经网络中，cost function 定义为

Note:

1.  $l-1$  表示去掉输出层
2.  $i = 1 \text{ to } s_l$  表示去掉  $\theta_{j0}$  这一列
3.  $j = 1 \text{ to } s_{j+1}$  表示全部行

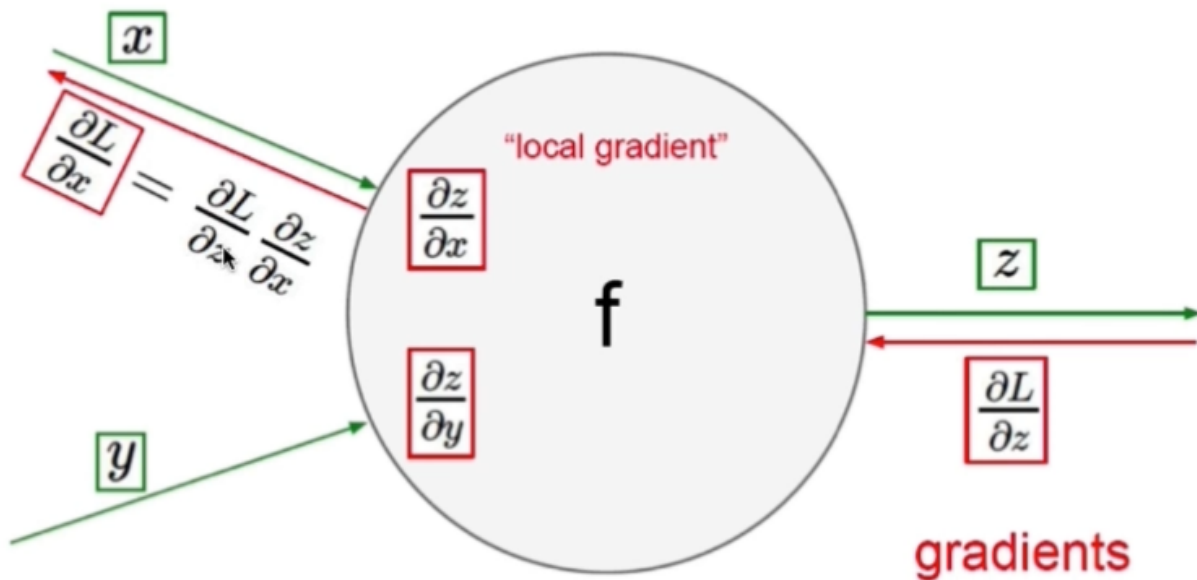
## Backpropagation algorithm(反向传播算法)

如果看不懂可以借鉴 [Blog链接](#)

1. Forward propagation
2. 为了计算导数项，引入 Back propagation algorithm  
Intuition:  $\delta^{(l)} = \text{"error" of node } j \text{ in layer } l$

直观理解





Example:

For each output unit(layer  $L = 4$ )

Step:

Training set  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

Set  $\Delta_{ij}^{(l)} = 0 \quad (\text{for all } l, i, j)$  (use to compute  $\frac{\partial \Delta}{\partial \theta_{ij}^{(l)}}(\theta)$ )

For  $i = 1$  to  $m$

set  $a^{(1)} = x^{(i)}$

Perform forward propagation to compute  $a^{(l)}$  for  $l=2, 3, \dots, L$

Using  $y^{(i)}$ , compute  $\delta^{(l)} = a^{(l)} - y^{(i)}$

Compute  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_j^{(l+1)}$

$\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$

$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)} \quad \text{if } j \neq 0 \quad D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \quad \text{if } j = 0$

$\frac{\partial \Delta}{\partial \theta_{ij}^{(l)}}(\theta) = D_{ij}^{(l)}$

## Gradient checking(梯度检测)

原理:

As for  $\vec{\theta}$

$\vec{\theta} \in \mathbb{R}^n$

$\vec{\theta} = [\theta_1, \theta_2, \dots, \theta_n]$

check that  $D_{\vec{\theta}} \approx \text{gradApprox}$

gradApprox is calculated by

$D_{\vec{\theta}}$  is calculated by Backpropation

Note:

1. 使用反向传播计算  $D_{\vec{\theta}}$

2. 使用梯度检验计算 `gradApprox`
3. 确保  $D_{\text{vect}} \approx \text{gradApprox}$
4. 不再使用 gradient checking, using backprop for learning

Important:

Be sure to disable your gradient checking code before training your classifier. If you run numerical gradient computation on every iteration of gradient descent, your code will be very slow

## Random initialization(随机初始化)

关于  $\text{vec}(\theta)$  的初始化一般具有两种方案

1.  $\text{vec}(\theta) = \text{vec}\{0\}$ 
  - After each update, parameters corresponding to inputs going into each of two hidden units are identical
2. Initialize each  $\theta_{ij}^{(l)}$  to a random value in  $[-\epsilon, \epsilon]$

显然我们选用方案2作为我们在神经网络中的  $\theta$  参数的初始化方案

## Summary

Training a neural network

1. Randomly initialize weights
2. Implement forward propagation to get  $h_{\theta}(x^{(i)})$  for any  $x^{(i)}$
3. Implement code to compute cost function  $J(\theta)$
4. Implement backprop to compute partial derivatives  $\frac{\partial}{\partial \theta_{jk}} J(\theta)$ . for  $i = 1:m$ , Perform forward propagation and back propagation using example  $(x^{(i)}, y^{(i)})$ . (Get activations  $a^{(l)}$  and delta terms  $\delta^{(l)}$  for  $l=2,3,\dots,L$ )
5. Use gradient to compute  $\frac{\partial}{\partial \theta_{jk}} J(\theta)$ . computed using backpropagation vs using numerical estimate of gradient of  $J(\theta)$
6. Use gradient descent or advanced optimization method with back propagation to try to minimize  $J(\theta)$  as a function of parameters  $\theta$

Note:

$J(\theta)$  is non-convex function in neural network, So, we can only get a local minimum.

## Code

这里我们使用 tensorflow 来逐步构建一个简单的神经网络模型。

### version 1

搜索 cifar-10 下载 python 格式的图片数据，一共有十类，这里我们使用二分类逻辑回归实现建模

```
1 #!/usr/bin/env python
2 # coding: utf-8
3
4
5 import tensorflow as tf
6 import os
```

```

7 import pickle
8 import numpy as np
9
10
11 CIFAT_DIR = '../cifar-10-batches-py'
12 print(os.listdir(CIFAT_DIR))
13
14
15 def load_data(filename):
16     """read data from data file"""
17     with open(os.path.join(filename), 'rb') as f:
18         # data = pickle.load(f, encoding='bytes')
19
20         # Python2.7代码
21         data = pickle.load(f)
22         return data['data'], data['labels']
23
24
25 class CifarData:
26     def __init__(self, filenames, need_shuffle):
27         all_data = []
28         all_labels = []
29         # 关于zip函数 具体看
30         # http://www.cnblogs.com/frydsh/archive/2012/07/10/2585370.html
31         for filename in filenames:
32             data, labels = load_data(filename)
33             for item, label in zip(data, labels):
34                 # label一共有是个类别 每个类别各 5000各
35                 # 使用该判断获取类别
36                 if label in [0, 1]:
37                     all_data.append(item)
38                     all_labels.append(label)
39
40         # 关于 vstack函数
41         # https://www.cnblogs.com/nkh222/p/8932369.html
42         self._data = np.vstack(all_data)
43         # 归一化处理
44         self._data = self._data / 127.5 - 1;
45         self._labels = np.hstack(all_labels)
46         print(self._data.shape)
47         print(self._labels.shape)
48         self._num_examples = self._data.shape[0]
49         self._need_shuffle = need_shuffle
50         self._indicator = 0
51         if self._need_shuffle:
52             self._shuffle_data()
53
54     def _shuffle_data(self):
55         # [0,1,2,3,4] => [2,1,3,4,0]
56         p = np.random.permutation(self._num_examples)
57         self._data = self._data[p]
58         self._labels = self._labels[p]
59
60     def next_batch(self, batch_size):

```

```

60     """return batch_size examples as a batch """
61     end_indicator = self._indicator + batch_size
62     if end_indicator > self._num_examples:
63         if self._need_shuffle:
64             self._shuffle_data()
65             self._indicator = 0
66             end_indicator = batch_size
67         else:
68             raise Exception("have no more examples")
69     if end_indicator > self._num_examples:
70         raise Exception('batch size is larger than all examles')
71     batch_data = self._data[self._indicator: end_indicator]
72     batch_labels = self._labels[self._indicator: end_indicator]
73     self._indicator = end_indicator
74     return batch_data, batch_labels
75
76
77 train_filenames = [os.path.join(CIFAT_DIR, 'data_batch_%d' % i) for i in range(1, 6)]
78 test_filenames = [os.path.join(CIFAT_DIR, 'test_batch')]
79
80 train_data = CifarData(train_filenames, True)
81 test_data = CifarData(test_filenames, False)
82 # batch_data, batch_labels = train_data.next_batch(10)
83 # print(batch_data, batch_labels)
84
85
86 # None 代表输入样本数是不确定的
87 x = tf.placeholder(tf.float32, [None, 3072])
88 # None
89 y = tf.placeholder(tf.int64, [None])
90 # 先构造一个 二分类器 因此输出为1
91 # (3072,1)
92 w = tf.get_variable('w', [x.get_shape()[-1], 1],
93                     initializer=tf.random_normal_initializer(0, 1))
94 # (1, )
95 b = tf.get_variable('b', [1], initializer=tf.constant_initializer(0.0))
96 # [None,3072] * [3072,1] = [None,1]
97 y_ = tf.matmul(x, w) + b
98 # [None,1]
99 p_y_1 = tf.nn.sigmoid(y_)
100 # 这里-1参数表示缺省值 保证为1列即可
101 y_resaped = tf.reshape(y, (-1, 1))
102 y_resaped_float = tf.cast(y_resaped, tf.float32)
103 # 计算loss
104 loss = tf.reduce_mean(tf.square(y_resaped_float - p_y_1))
105 predict = p_y_1 > 0.5
106 correct_prediction = tf.equal(tf.cast(predict, tf.int64), y_resaped)
107 accuracy = tf.reduce_mean(tf.cast(correct_prediction, tf.float64))
108
109 with tf.name_scope('train_op'):
110     # 这里1e-3是学习率 learning rate AdamOptimizer是梯度下降的一个变种
111     train_op = tf.train.AdamOptimizer(1e-3).minimize(loss)

```

```

112 '''
113 到此为止我们的计算图搭建完成
114 '''
115
116 init = tf.global_variables_initializer()
117 batch_size = 20
118 train_steps = 100000
119 test_steps = 100
120
121 with tf.Session() as sess:
122     sess.run(init)
123     for i in range(train_steps):
124         batch_data, batch_labels = train_data.next_batch(batch_size)
125         loss_val, accu_val, _ = sess.run(
126             [loss, accuracy, train_op],
127             feed_dict={x: batch_data, y: batch_labels})
128         if (i+1) % 500 == 0:
129             print('[Train] Step: %d, loss: %4.5f, acc: %4.5f' % (i+1, loss_val,
130 accu_val))
131         if (i+1) % 5000 == 0:
132             test_data = CifarData(test_filenames, False)
133             all_test_acc_val = []
134             for j in xrange(test_steps):
135                 test_batch_data, test_batch_labels \
136                     = test_data.next_batch(batch_size)
137                 test_acc_val = sess.run(
138                     [accuracy],
139                     feed_dict={
140                         x: test_batch_data,
141                         y: test_batch_labels
142                     })
143                 all_test_acc_val.append(test_acc_val)
144             test_acc = np.mean(all_test_acc_val)
145             print('[Test] Step: %d, acc: %4.5f ' % (i+1, test_acc))
146
147
148

```

运行结果:

```

1  [Train] Step: 98500, loss: 0.10032, acc: 0.90000
2  [Train] Step: 99000, loss: 0.10000, acc: 0.90000
3  [Train] Step: 99500, loss: 0.10080, acc: 0.90000
4  [Train] Step: 100000, loss: 0.05529, acc: 0.95000
5  (2000, 3072)
6  (2000,)
7  [Test] Step: 100000, acc: 0.81200
8
9  Process finished with exit code 0

```

**version 2**

这里我们继续使用该算法实现多分类器

```
1  #!/usr/bin/env python
2  # coding: utf-8
3
4
5  import tensorflow as tf
6  import os
7  import pickle
8  import numpy as np
9
10
11  CIFAT_DIR = '../cifar-10-batches-py'
12  print(os.listdir(CIFAT_DIR))
13
14
15  def load_data(filename):
16      """read data from data file"""
17      with open(os.path.join(filename), 'rb') as f:
18          # data = pickle.load(f, encoding='bytes')
19
20          # Python2.7代码
21          data = pickle.load(f)
22          return data['data'], data['labels']
23
24
25  class CifarData:
26      def __init__(self, filenames, need_shuffle):
27          all_data = []
28          all_labels = []
29          # 关于zip函数 具体看
30          # http://www.cnblogs.com/frydsh/archive/2012/07/10/2585370.html
31          for filename in filenames:
32              data, labels = load_data(filename)
33              for item, label in zip(data, labels):
34                  all_data.append(item)
35                  all_labels.append(label)
36
37          # 关于 vstack函数
38          # https://www.cnblogs.com/nkh222/p/8932369.html
39          self._data = np.vstack(all_data)
40          # 归一化处理
41          self._data = self._data / 127.5 - 1;
42          self._labels = np.hstack(all_labels)
43          print(self._data.shape)
44          print(self._labels.shape)
45          self._num_examples = self._data.shape[0]
46          self._need_shuffle = need_shuffle
47          self._indicator = 0
48          if self._need_shuffle:
49              self._shuffle_data()
50
51      def _shuffle_data(self):
52          # [0,1,2,3,4] => [2,1,3,4,0]
```

```

52     p = np.random.permutation(self._num_examples)
53     self._data = self._data[p]
54     self._labels = self._labels[p]
55
56     def next_batch(self, batch_size):
57         """return batch_size examples as a batch """
58         end_indicator = self._indicator + batch_size
59         if end_indicator > self._num_examples:
60             if self._need_shuffle:
61                 self._shuffle_data()
62                 self._indicator = 0
63                 end_indicator = batch_size
64             else:
65                 raise Exception("have no more examples")
66         if end_indicator > self._num_examples:
67             raise Exception('batch size is larger than all examles')
68         batch_data = self._data[self._indicator: end_indicator]
69         batch_labels = self._labels[self._indicator: end_indicator]
70         self._indicator = end_indicator
71         return batch_data, batch_labels
72
73
74 train_filenames = [os.path.join(CIFAT_DIR, 'data_batch%d' % i) for i in range(1, 6)]
75 test_filenames = [os.path.join(CIFAT_DIR, 'test_batch')]
76
77 train_data = CifarData(train_filenames, True)
78 test_data = CifarData(test_filenames, False)
79 # batch_data, batch_labels = train_data.next_batch(10)
80 # print(batch_data, batch_labels)
81
82
83 # None 代表输入样本数是不确定的
84 x = tf.placeholder(tf.float32, [None, 3072])
85 # None
86 y = tf.placeholder(tf.int64, [None])
87 # 先构造一个 二分类器 因此输出为1
88 # (3072,10)
89 w = tf.get_variable('w', [x.get_shape()[-1], 10],
90                     initializer=tf.random_normal_initializer(0, 1))
91 # (10, )
92 b = tf.get_variable('b', [10], initializer=tf.constant_initializer(0.0))
93 # [None,3072] * [3072,10] = [None,10]
94 y_ = tf.matmul(x, w) + b
95
96 # 关于softmax https://www.zhihu.com/question/23765351
97 # [[0,01,0.9,...,0.02],[ ]]
98 p_y = tf.nn.softmax(y_)
99 # 6 -->[0,0,0,0,0,1,0,0,0,0]
100 y_one_hot = tf.one_hot(y, 10, dtype=tf.float32)
101 loss = tf.reduce_mean(tf.square(y_one_hot - p_y))
102
103 '''
104 # [None,10]

```

```

104 p_y_1 = tf.nn.sigmoid(y_)
105 # 这里-1参数表示缺省值 保证为1列即可
106 y_resaped = tf.reshape(y, (-1, 1))
107 y_resaped_float = tf.cast(y_resaped, tf.float32)
108 # 计算loss
109 loss = tf.reduce_mean(tf.square(y_resaped_float - p_y_1))
110 '''
111
112 # indices
113 predict = tf.argmax(y_, 1)
114 correct_prediction = tf.equal(predict, y)
115 accuracy = tf.reduce_mean(tf.cast(correct_prediction, tf.float64))
116
117 with tf.name_scope('train_op'):
118     # 这里1e-3是学习率 learning rate AdamOptimizer是梯度下降的一个变种
119     train_op = tf.train.AdamOptimizer(1e-3).minimize(loss)
120
121 '''
122 到此为止我们的计算图搭建完成
123 '''
124
125 init = tf.global_variables_initializer()
126 batch_size = 20
127 train_steps = 10000
128 test_steps = 100
129
130 with tf.Session() as sess:
131     sess.run(init)
132     for i in range(train_steps):
133         batch_data, batch_labels = train_data.next_batch(batch_size)
134         loss_val, accu_val, _ = sess.run(
135             [loss, accuracy, train_op],
136             feed_dict={x: batch_data, y: batch_labels})
137         if (i+1) % 500 == 0:
138             print('[Train] Step: %d, loss: %4.5f, acc: %4.5f' % (i+1, loss_val,
accu_val))
139         if(i+1) % 5000 == 0:
140             test_data = CifarData(test_filenames, False)
141             all_test_acc_val = []
142             for j in xrange(test_steps):
143                 test_batch_data, test_batch_labels \
144                     = test_data.next_batch(batch_size)
145                 test_acc_val = sess.run(
146                     [accuracy],
147                     feed_dict={
148                         x: test_batch_data,
149                         y: test_batch_labels
150                     }
151                 )
152                 all_test_acc_val.append(test_acc_val)
153             test_acc = np.mean(all_test_acc_val)
154             print('[Test] Step: %d, acc: %4.5f ' % (i+1, test_acc))

```



## Note

这两部分代码都没有用到hidden layer.

实际上，**code 1** 展示的是一个神经元，这里也可以认为是逻辑回归。也就是logistic regression 看做是仅仅含有一个神经元的单 层神经网络

code 2 实际上也就是多维的logistic regreesion,其实softmax regression可以看做是含有k个神经元的一层神经网络。