

Machine Learning(IDC410) Report: Programming Exercise 2

Rimjhim Goel
Roll No: MS18133

Contents

1	Question 1: Generating the data set	2
2	Question 2: Logistic regression	2
3	Question 3:	2
3.1	Impact of variation of n and θ on logistic regression	2
3.1.1	Variation of θ	2
3.1.2	Variation of n	3
3.2	Derivation of partial derivative of cost function	4
4	Question 4: L1 and L2 regularization	4
4.1	L1/Lasso regression	5
4.2	L2/Ridge regression	5
5	List of files submitted	6

1 Question 1: Generating the data set

```
1 import numpy as np
2 from scipy.stats import bernoulli
3 import matplotlib.pyplot as plt
4
5 def gen_data(theta,n,m):
6     X = np.random.randn(n,m+1)
7     X[:,0]=1 # first column set to 1
8     beta = np.random.randn(m+1,1) # random coefficients
9     bern = np.transpose(bernoulli.rvs(size=n,p=theta)) #Bernoulli distribution
10    Y = 1/(1+ np.exp(-(np.matmul(X,beta))))
11
12    for i in range(len(Y)): #classification of Y
13        if Y[i]>0.5:
14            Y[i] = 1
15        else:
16            Y[i]= 0
17    for i in range(len(Y)): # flipping according to Bernoulli distribution
18        if bern[i]==1:
19            Y[i] = 1-Y[i]
20
21    return X,beta,Y
```

2 Question 2: Logistic regression

Log likelihood/ cost function -

$$\log C = \sum_{n=1}^N y_n \log(\sigma(\beta^T X_n)) + (1 - y_n) \log(1 - \sigma(\beta^T X_n))$$

Gradient-

$$\frac{\partial \log C}{\partial \beta} = X^T (Y - Y_{\text{predicted}})$$

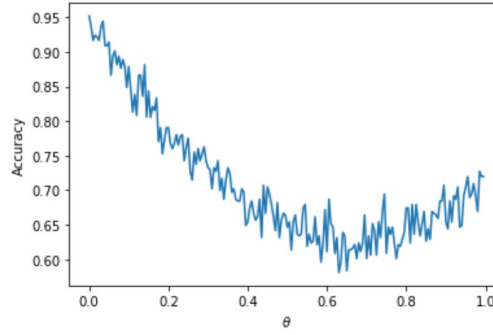
```
1 def logistic_regression(X,Y,k,tau,lam):
2     z = float(len(X))
3     b = np.random.randn(int(len(beta)),1) #random initialisation of the parameters
4     previous_cost = 0 #initial cost set to 0
5     cost = np.array([])
6
7
8     for i in range(k):# k is no. of iterations
9         y_pred = 1.0/(1+ np.exp(-(np.matmul(X,b))))# predicted y
10        temp = np.subtract(y_pred,Y) # Y-Xb
11        temp2 = np.subtract(np.ones_like(Y),Y)
12        # Cost function/ Cross Entropy
13        current_cost = -(1/z)*np.matmul(np.transpose(Y),np.log(y_pred,out=np.zeros_like(y_pred), where=(y_pred!=0))) - np.matmul(np.transpose(temp2),np.log(temp2,out=np.zeros_like(temp2), where=(temp2!=0)))
14        cost = np.append(cost,current_cost)
15        if abs(current_cost-previous_cost)<=tau: # tau - threshold condition on change in cost function
16            break
17
18        previous_cost = current_cost
19        grad = (1/z)*(np.matmul(np.transpose(X),temp))# gradient calculation
20        b = b - lam*grad #improved b, delta is step size/learning parameter
21    return(y_pred,cost[-1])
```

3 Question 3:

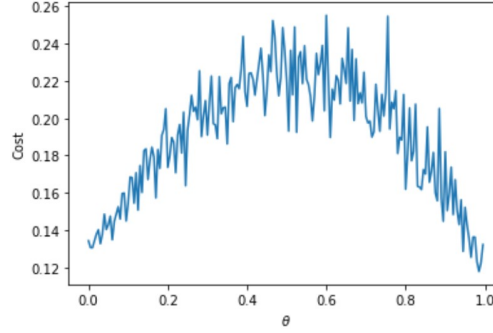
3.1 Impact of variation of n and θ on logistic regression

3.1.1 Variation of θ

For investigating how θ affects the accuracy of the model, other variables(n,m,X, β) were kept fixed for each theta value. $y_{\text{predicted}}$ is classified as 1 if probability is greater than 0.5 and 0 otherwise. Accuracy of the



(a) Accuracy vs θ



(b) Cost vs θ

model is calculated as-

$$Accuracy = 1 - \frac{\sum_{i=1}^n (y_i - y_{predicted_i})}{n}$$

where n is the total no. of elements in the y vector.

Since Bernoulli distribution varies in each iteration, accuracy is averaged over a number of iterations for each θ .

For $\theta \rightarrow 0$ and $\theta \rightarrow 1$, the cost function is minimum. Accuracy is maximum for $\theta \rightarrow 0$ and then starts reducing and is minimum for $\theta \rightarrow 0.5$ and then starts increasing again. This is because for $\theta \rightarrow 0.5$, the randomness in the data is maximum as probability of flipping and not flipping is equal. For $\theta = 0$ and $\theta = 1$, either none of the values are flipped or all are flipped. This means that the randomness in the data is less as compared to the θ values in between. Hence, the accuracy is maximum for $\theta \rightarrow 0$ and then again increases for $\theta \rightarrow 1$.

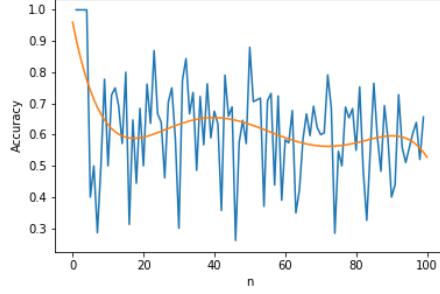
3.1.2 Variation of n

For observing the effect of n on the model, θ was kept zero i.e. no random flipping due to noise. There are lot of fluctuations in the accuracy because of inherent randomness due to random initialisation of β in the logistic regression model. However, the accuracy is averaged over a number of iterations for each n.

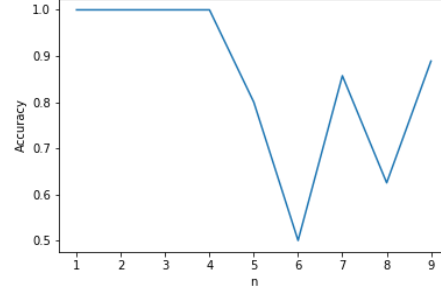
The plot shows accuracy = 1 for very small values of n. This is because it is easier to fit lesser data points. A perfect fit does not imply efficient machine learning. Small testing data sets underestimate the parameters and lead to overfitting. Such algorithms produce poor results when applied on training sets.

Thus an optimum sample size is necessary to represent all the parameters.

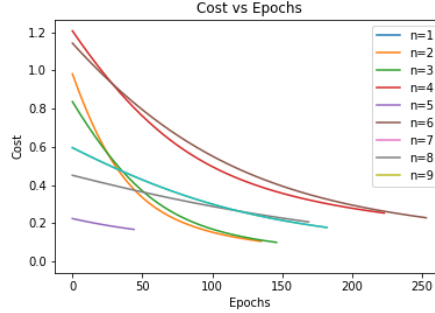
The cost converges to a minimum value for each n but there is no relation between n and the number of epochs it takes to converge.



(a) Accuracy vs n



(b) Accuracy for small n



(c) Cost vs Epochs

3.2 Derivation of partial derivative of cost function

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

Let p_n be the probability of being in a certain class.

$$p_n = \sigma(\beta^T X_n)$$

Maximum likelihood principle:

$$p(y_n) = p_n^{y_n} (1 - p_n)^{(1-y_n)}$$

$$p(y_n) = \sigma(\beta^T X_n)^{y_n} (1 - \sigma(\beta^T X_n))^{(1-y_n)}$$

Likelihood function:

$$C = \prod_{n=1}^N \sigma(\beta^T X_n)^{y_n} (1 - \sigma(\beta^T X_n))^{(1-y_n)}$$

Log likelihood:

$$\log C = \sum_{n=1}^N y_n \log(\sigma(\beta^T X_n)) + (1 - y_n) \log(1 - \sigma(\beta^T X_n))$$

Maximising the log likelihood principle:

$$\frac{\partial \log C}{\partial \beta} = \sum_{n=1}^N y_n \frac{1}{\sigma(\beta^T X_n)} \frac{\partial \sigma(\beta^T X_n)}{\partial \beta} - (1 - y_n) \frac{1}{1 - \sigma(\beta^T X_n)} \frac{\partial \sigma(\beta^T X_n)}{\partial \beta}$$

$$= \sum_{n=1}^N y_n (1 - \sigma(\beta^T X_n)) x_n - (1 - y_n) \sigma(\beta^T X_n) x_n$$

$$\frac{\partial \log C}{\partial \beta} = X^T (Y - Y_{\text{predicted}})$$

4 Question 4: L1 and L2 regularization

Regularization is used to avoid overfitting of the data by adding additional terms to the cost function. These additional terms shrink the estimated parameters towards zero. Regularisation is essentially adding bias to the model weights so that it does not overfit.

4.1 L1/Lasso regression

L1 norm is added

$$\|\beta\| = |\beta_1| + |\beta_2| + \dots + |\beta_{n+1}|$$

$$Cost = Cost + \lambda \sum_{d=1}^{m+1} |\beta_d|$$

$$\frac{\partial Cost}{\partial \beta} = -X^T(Y - X\beta) + \lambda I' \text{sign}(\beta)$$

where I' is an identity matrix of size $m + 1$ and the first element is zero (since β_0 is not penalised) and λ is the tuning parameter which represents the amount of regularisation. Larger the value of λ , greater the penalty on the parameter estimates which shrinks the estimates towards zero.

L1 regularisation reduces the number of features (β) by eliminating not-so-important features. This is used with data sets having large number of features to reduce complexity.

```

1  #lasso regression L1
2  def logistic_regression_L1(X,Y,k,tau,lam,alpha):
3      z = float(len(X))
4      b = np.random.randn(int(len(beta)),1) #random initialisation of the parameters
5      previous_cost = 0 #initial cost set to 0
6      cost = np.array([])
7
8      b_L2 = 0
9
10     for i in range(1,len(b)):
11         b_L2 += abs(b[i]) # L1 norm
12
13     I_L2 = np.ones((len(beta),1))
14     I_L2[:,0]=0 # I' matrix
15     I_L2 = I_L2*np.sign(b)
16
17     for i in range(k):# k is no. of iterations
18
19         y_pred = 1.0/(1+ np.exp(-(np.matmul(X,b))))# predicted y
20
21         temp = np.subtract(y_pred,Y) # Y-Xb
22         temp2 = np.subtract(np.ones_like(Y),Y)
23         # cost/cross entropy calculation
24         current_cost = -(1/z)*np.matmul(np.transpose(Y),np.log(y_pred,out=np.
zeros_like(y_pred), where=(y_pred!=0))) - np.matmul(np.transpose(temp2),np.log(
temp2,out=np.zeros_like(temp2), where=(temp2!=0)))
25         # for L1 regression
26         current_cost = current_cost + (1/z)*(alpha*b_L2)
27         cost = np.append(cost,current_cost)
28
29         if abs(current_cost-previous_cost)<=tau: # tau - threshold condition on
change in cost function
30             break
31         previous_cost = current_cost
32
33         grad = (1/z)*(np.matmul(np.transpose(X),temp))# gradient calculation
34         # gradient calculation for L1 regression, alpha is tuning parameter
35         grad = grad - alpha*I_L2
36         b = b - lam*grad #improved b, delta is step size/learning parameter
37     return(y_pred,cost[-1])

```

4.2 L2/Ridge regression

square of L2 norm is added

$$\|\beta\|^2 = |\beta_1|^2 + |\beta_2|^2 + \dots + |\beta_{n+1}|^2$$

$$Cost = Cost + \frac{\lambda}{2} \sum_{d=1}^{m+1} |\beta_d|^2$$

$$\frac{\partial Cost}{\partial \beta} = -X^T(Y - X\beta) + \lambda I' \beta$$

L2 regularisation does not eliminate β , it just makes their contribution smaller. This is used when the parameters are highly correlated.

```

1  #ridge regression L2
2
3  def logistic_regression_L2(X,Y,k,tau,lam,alpha):
4      z = float(len(X))
5      b = np.random.randn(int(len(beta)),1) #random initialisation of the parameters
6      previous_cost = 0 #initial cost set to 0
7      cost = np.array([])
8      iter_ =[]
9
10     b_L2 = 0
11     for i in range(1,len(b)):
12         b_L2 += b[i]*b[i] #square of L2 norm
13
14     I_L2 = np.ones((len(beta),1))
15     I_L2[:,0]=0 #I'matrix
16     I_L2 = I_L2*b
17
18     for i in range(k):# k is no. of iterations
19         iter_.append(i)
20         y_pred = 1.0/(1+ np.exp(-(np.matmul(X,b))))# predicted y
21         temp = np.subtract(y_pred,Y) # Y-Xb
22         temp2 = np.subtract(np.ones_like(Y),Y)
23         # cost/cross entropy calculation
24         current_cost = -(1/z)*np.matmul(np.transpose(Y),np.log(y_pred,out=np.
25         zeros_like(y_pred), where=(y_pred!=0))) - np.matmul(np.transpose(temp2),np.log(
26         temp2,out=np.zeros_like(temp2), where=(temp2!=0)))
27         current_cost = current_cost + (1/z)*(alpha*b_L2) #cost with L2 term
28         cost = np.append(cost,current_cost)
29         if abs(current_cost-previous_cost)<=tau: # tau - threshold condition on
30         change in cost function
31             break
32         previous_cost = current_cost
33         #gradient calculation
34         grad = (1/z)*(np.matmul(np.transpose(X),temp))# gradient calculation
35         # gradient for L2 regression, alpha is tuning parameter
36         grad = grad - alpha*I_L2
37         b = b - lam*grad #improved b, delta is step size/learning parameter
38     return(y_pred,cost[-1])

```

5 List of files submitted

1. .ipynb notebook containing code
2. pdf report