

Machine Learning: Assignment 2

Anjali Krishna Prasad collaborated with Arpita Agrawal

October 2016

1 Question 1a

$$P(y_n|x_n; b; w) = \begin{cases} \sigma(b + w^T X_n) & \text{if } y_n = 1 \\ 1 - \sigma(b + w^T X_n) & \text{if } y_n = 0 \end{cases}$$

$$\text{Where } \sigma_a = \frac{1}{1 + e^{-a}} \text{ and } a = b + w^T x_n$$

$$P(y_n|x_n; b; w) = \prod_{i=1}^n P(Y = y_i|X = x_i)$$

$$= \sum_n \sigma(b + w^T X_n)^{y_n} [1 - \sigma(b + w^T X_n)]^{1-y_n}$$

$$-\log P(y_n|x_n; b; w) = - \sum_n [y_n \log[\sigma(b + w^T X_n)] + (1 - y_n) \log[1 - \sigma(b + w^T X_n)]]$$

2 Question 1b

Append b to w $a = b + w^T X_n \rightarrow w^T X_n$

$$\text{Derivative of } \frac{\partial \sigma(a)}{\partial a} = \frac{\partial}{\partial a} \left(\frac{1}{1 + e^{-a}} \right) = \frac{-(1 + e^{-a})'}{(1 + e^{-a})^2} = \frac{e^{-a}}{(1 + e^{-a})^2} = \left(\frac{1}{1 + e^{-a}} \right) \left(1 - \frac{1}{1 + e^{-a}} \right)$$

$$\frac{\partial \sigma(a)}{\partial a} = \sigma(a)[1 - \sigma(a)]$$

$$\frac{\partial \log \sigma(a)}{\partial a} = 1 - \sigma(a)$$

$$\frac{\partial E(w)}{\partial w} = - \sum_n y_n [1 - \sigma(w^T X_n)] + (1 - y_n) [\sigma(w^T X_n)]$$

$$- \sum_n X_n [y_n - y_n \sigma(w^T X_n) - \sigma(w^T X_n) + y_n \sigma(w^T X_n)]$$

$$\sum_n X_n [\sigma(w^T X_n) - y_n]$$

$$\text{update rule for w } w^{(t+1)} = w^{(t)} - \eta \sum_n X_n [\sigma(w^T X_n) - y_n]$$

$$H = \frac{\partial^2 E}{\partial w \partial w^T} = \sum_n X_n X_n^T [1 - \sigma(w^T X_n)] [\sigma(w^T X_n)]$$

$X_n X_n^T$ has to be positive and $\sigma(w^T X_n); 1 - \sigma(w^T X_n)$ are probabilities and cannot be < 0

Product of 3 positive terms is positive therefore $v^T H v \geq 0$

Thus, positive definite. Thus, the cross-entropy error function is convex, with only one global optimum.

3 question 1c

$$P(Y = k|X = x) = \frac{\exp^{w_k^T x}}{1 + \sum_1^{k-1} \exp^{w_i^T x}} \text{ where } k = 1 \dots K-1$$

$$P(Y = k|X = x) = \frac{1}{1 + \sum_1^{k-1} \exp^{w_i^T x}} \text{ where } k = K$$

$w_K = 0$ The above formula can also be written as

$$P(Y = k|X = x) = \frac{\exp^{w_k^T x}}{\exp^{w_K^T x} + \sum_1^{k-1} \exp^{w_i^T x}}$$

$$P(Y = k|X = x) = \frac{\exp^{w_k^T x}}{\sum_1^K \exp^{w_i^T x}}$$

$$P(D) = P(y_n|x_n)$$

$$\log P(D) = \sum_{j=1}^n \log P(y_n|x_n)$$

We will change y_n to $y_n[y_{n1}y_{n2}\dots y_{nk}]^T$ a K dimensional vector using 1 of K encoding

$$y_{ji} \begin{cases} 1 & \text{if } y_j = k \\ 0 & \text{if otherwise} \end{cases}$$

$$-\log P(y_n|x_n) = -\sum_{j=1}^n \log \prod_{i=1}^k P(C_k|X_j)^{y_{ji}}$$

$$-\log P(y_n|x_n) = -\sum_{j=1}^n \sum_{i=1}^K y_{ji} \log P(C_i|X_j)$$

$$-\log P(y_n|x_n) = -\sum_{j=1}^n \sum_{i=1}^K y_{ji} \log \frac{\exp^{w_i^T x_j}}{\sum_1^K \exp^{w_i^T x_j}}$$

$$-\log P(y_n|x_n) = -\sum_{j=1}^n \sum_{i=1}^K y_{ji} [\log(\exp^{w_i^T x_j}) - \log(\sum_1^K \exp^{w_i^T x_j})]$$

$$-\log P(y_n|x_n) = -\sum_{j=1}^n \sum_{i=1}^K y_{ji} [(w_i^T x_j) - \log(\sum_1^K \exp^{w_i^T x_j})]$$

$$-\log P(y_n|x_j) = \sum_{j=1}^n \sum_{i=1}^K y_{ji} \log(\sum_1^K \exp^{w_i^T x_j}) - y_{ji}(w_i^T x_j)$$

4 Question 1d

Differentiating partially wrt w_i negative log likelihood from above answer

$$\begin{aligned} -\log P(y_n|x_n) &= \sum_j^n \sum_{i=1}^K y_{ji} \log(\sum_1^K \exp^{w_i^T x_j}) - y_{ji}(w_i^T x_j) \\ -\frac{\partial \log P(y_n|x_n)}{\partial w_i} &= \frac{\partial}{\partial w_i} (\sum_{j=1}^n \sum_{i=1}^K y_{ji} \log(\sum_1^K \exp^{w_i^T x_j}) - y_{ji}(w_i^T x_j)) \\ &= \sum_{i=1}^K \sum_{j=1}^n \left[-y_{ji} x_j + \frac{y_{ji} \exp^{w_i^T x_j} x_j}{\sum_1^K \exp^{w_i^T x_j}} \right] = \sum_{i=1}^K \sum_{j=1}^n \left[-y_{ji} x_j + y_{ji} x_j P(Y = k|X = x) \right] \\ \text{update rule for } w & w^{(t+1)} = w^{(t)} - \eta \sum_{i=1}^K \sum_{j=1}^n \left[-y_{ji} x_j + y_{ji} x_j P(Y = k|X = x) \right] \end{aligned}$$

5 Question 2a

$$f(x) = \begin{cases} p_1 \frac{1}{\sqrt{2\pi}\sigma_1} \exp \frac{-(x-\mu_1)^2}{2\sigma_1^2} & \text{if } y_n = 1 \\ p_2 \frac{1}{\sqrt{2\pi}\sigma_2} \exp \frac{-(x-\mu_2)^2}{2\sigma_2^2} & \text{if } y_n = 2 \end{cases}$$

$$\text{Likelihood function } D = \prod_{y_n=1}^2 P(y_n)P(x_n|y_n)$$

$$\log P(D) = \sum_n \log P(x_n, y_n)$$

$$\sum_{n:y_n=1} \log(p_1 \frac{1}{\sqrt{2\pi}\sigma_1} \exp \frac{-(x-\mu_1)^2}{2\sigma_1^2}) + \sum_{n:y_n=2} \log(p_2 \frac{1}{\sqrt{2\pi}\sigma_2} \exp \frac{-(x-\mu_2)^2}{2\sigma_2^2})$$

$$\sum_{n:y_n=1} (\log p_1 + \log \frac{1}{\sqrt{2\pi}\sigma_1} + \frac{-(x-\mu_1)^2}{2\sigma_1^2}) + \sum_{n:y_n=2} (\log p_2 + \log \frac{1}{\sqrt{2\pi}\sigma_2} + \frac{-(x-\mu_2)^2}{2\sigma_2^2})$$

We know that $p_2 = 1 - p_1$

$$\sum_{n:y_n=1} (\log p_1 + \log \frac{1}{\sqrt{2\pi}\sigma_1} + \frac{-(x-\mu_1)^2}{2\sigma_1^2}) + \sum_{n:y_n=2} (\log(1-p_1) + \log \frac{1}{\sqrt{2\pi}\sigma_2} + \frac{-(x-\mu_2)^2}{2\sigma_2^2})$$

Differentiating $\log P(D)$ wrt p_1 and equating to 0

$$\frac{\partial \log P(D)}{\partial p_1} = \sum_{n:y_n=1} \frac{1}{p_1} - \sum_{n:y_n=2} \frac{1}{1-p_1} = 0$$

Let n_1 = no of samples in training set with $y_n = 1$

N is the total number of samples. $N - n_1$ = no of samples in training set with $y_n = 2$

$$\frac{n_1}{p_1} - \frac{N - n_1}{(1 - p_1)} = 0$$

$$n_1 - p_1 n_1 - p_1 N + p_1 n_1 = 0$$

$$p_1 = \frac{n_1}{N}$$

$$p_2 = 1 - p_1$$

$$p_2 = 1 - \frac{n_1}{N}$$

$$p_2 = \frac{N - n_1}{N}$$

Let $N - n_1 = n_2$

$$p_2 = \frac{n_2}{N}$$

Differentiating $\log P(D)$ wrt μ_1 and equating to 0

$$\frac{\partial \log P(D)}{\partial \mu_1} = \sum_{n:y_n=1} -\frac{1}{2\sigma_1^2} [2(x_n - \mu_1)](-1) = 0$$

$$\sum_{n:y_n=1} 2x_n - 2\mu_1 = 0$$

$$\mu_1 = \frac{\sum_{n:y_n=1} x_n}{n_1}$$

Similarly Differentiating $\log P(D)$ wrt μ_2 and equating to 0

$$\mu_2 = \frac{\sum_{n:y_n=2} x_n}{n_2}$$

Differentiating $\log P(D)$ wrt σ_1 and equating to 0

$$\frac{\partial \log P(D)}{\partial \sigma_1} = \sum_{n:y_n=1} -\frac{\sqrt{2\pi}}{\sqrt{2\pi}\sigma_1} - \frac{(x_n - \mu_1)^2}{2\sigma_1^3}(-1) = 0$$

$$\sum_{n:y_n=1} \frac{\sigma_1^2 + (x_n - \mu_1)^2}{\sigma_1^3} = 0$$

$$\sigma_1^2 = \frac{\sum_{n:y_n=1} (x_n - \mu_1)^2}{n_1}$$

$$\sigma_1 = \sqrt{\frac{\sum_{n:y_n=1} (x_n - \mu_1)^2}{n_1}}$$

Similarly differentiating $\log P(D)$ wrt σ_2 and equating to 0

$$\sigma_2^2 = \frac{\sum_{n:y_n=2} (x_n - \mu_2)^2}{n_2}$$

$$\sigma_2 = \sqrt{\frac{\sum_{n:y_n=2} (x_n - \mu_2)^2}{n_2}}$$

6 Question 2b

$$P(Y = c_1|X) = \frac{P(Y = c_1)P(X|Y = c_1)}{P(Y = c_1)P(X|Y = c_1) + P(X|Y = c_2)P(Y = c_2)}$$

Dividing numerator and denominator by $P(Y = c_1)P(X|Y = c_1)$

$$P(Y = c_1|X) = \frac{1}{1 + \frac{P(X|Y=c_2)P(Y=c_2)}{P(Y=c_1)P(X|Y=c_1)}}$$

$$\frac{1}{1 + \exp(\ln \frac{P(Y=c_2)P(X|Y=c_2)}{P(Y=c_1)P(X|Y=c_1)})}$$

$$\frac{1}{1 + \exp(\ln \frac{P(Y=c_2)}{P(Y=c_1)}) * \ln \frac{P(X|Y=c_2)}{P(X|Y=c_1)}}$$

Let $P(Y = c_1) = p_1$; $P(Y = c_2) = p_2$

$$\frac{1}{1 + \exp(\ln \frac{p_2}{p_1}) * \ln \frac{P(X|Y=c_2)}{P(X|Y=c_1)}}$$

$$\frac{1}{1 + \exp[(\ln p_2 - \ln p_1) + \ln P(X|Y = c_2) - \ln P(X|Y = c_1)]}$$

$$P(X|Y = c_1) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(X - \mu_1)^T \Sigma^{-1}(X - \mu_1)\right\}$$

$$P(X|Y = c_2) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(X - \mu_2)^T \Sigma^{-1}(X - \mu_2)\right\}$$

$$\frac{1}{1 + \exp[(\ln p_2 - \ln p_1) + \ln \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp^{-\frac{1}{2}(X - \mu_2)^T \Sigma^{-1}(X - \mu_2)} - \ln \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp^{-\frac{1}{2}(X - \mu_1)^T \Sigma^{-1}(X - \mu_1)}]}$$

Lets solve the following part of the denominator separately due to space constraint

$$\ln \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp^{-\frac{1}{2}(X - \mu_2)^T \Sigma^{-1}(X - \mu_2)} - \ln \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp^{-\frac{1}{2}(X - \mu_1)^T \Sigma^{-1}(X - \mu_1)}$$

$$\ln \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} - \frac{1}{2}(X - \mu_2)^T \Sigma^{-1}(X - \mu_2) - \ln \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} + \frac{1}{2}(X - \mu_1)^T \Sigma^{-1}(X - \mu_1)$$

$$-\frac{1}{2}(X - \mu_2)^T \Sigma^{-1}(X - \mu_2) + \frac{1}{2}(X - \mu_1)^T \Sigma^{-1}(X - \mu_1)$$

Using property $(A + B)^T = A^T + B^T$ and opening the brackets, we get

$$= (\mu_2 - \mu_1)^T \Sigma^{-1} X + \frac{1}{2} [\mu_1^T \Sigma^{-1} \mu_1 - \mu_2^T \Sigma^{-1} \mu_2]$$

Getting back the entire exp term from the original equation

$$\ln p_2 - \ln p_1 + (\mu_2 - \mu_1)^T \Sigma^{-1} X + \frac{1}{2} [\mu_1^T \Sigma^{-1} \mu_1 - \mu_2^T \Sigma^{-1} \mu_2]$$

$$\text{Therefore } b = \frac{1}{2} [\mu_1^T \Sigma^{-1} \mu_1 - \mu_2^T \Sigma^{-1} \mu_2] + \ln p_2 - \ln p_1$$

$$\theta^T = [(\mu_2 - \mu_1)^T \Sigma^{-1}]$$

7 Programming assisngment

3.1

Histograms available in the my program.

```
Pearson's Correlation for column 1 is -0.387696987621
Pearson's Correlation for column 2 is 0.362987295831
Pearson's Correlation for column 3 is -0.483067421758
Pearson's Correlation for column 4 is 0.203600144696
Pearson's Correlation for column 5 is -0.424829675619
Pearson's Correlation for column 6 is 0.690923334973
Pearson's Correlation for column 7 is -0.390179110401
Pearson's Correlation for column 8 is 0.252420566225
Pearson's Correlation for column 9 is -0.385491814423
Pearson's Correlation for column 10 is -0.468849385373
Pearson's Correlation for column 11 is -0.505270756892
Pearson's Correlation for column 12 is 0.343434137151
Pearson's Correlation for column 13 is -0.73996982063
```

3.2

Linear Regression

MSE for training set 20.950144508

MSE for testing set 28.4179164975

Ridge Regression

For lambda = 0.01

MSE for training set 20.9501449001

MSE for testing set 28.4182915618

For lambda = 0.1

MSE for training set 20.9501836546

MSE for testing set 28.42168497

For lambda = 1.0

MSE for training set 20.9539918317

MSE for testing set 28.4573733573

kfold cross validation:

When we dont shuffle we get

For $\lambda = 5.4001$
 MSE for training set 21.0544652787
 MSE for testing set 28.6768029918
 Shuffling has an effect on the set.
 My values of λ is between 0 and 1
 example
 winning $\lambda = 1.0001$
 MSE for training set after using winning λ 20.9539925939
 MSE for training set after using winning λ 28.4573774989

 winning $\lambda = 0.9001$
 MSE for training set after using winning λ 20.9507674828
 MSE for training set after using winning λ 28.4332314197

 3.3
 a) Top four correlated columns with the target are 13 6 11 3
 MSE after taking top four correlated columns with the target are:
 MSE for training set 26.4066042155
 MSE for testing set 31.4962025449

 b) Select 4 features iteratively and select top 4: 13, 6, 11, 4
 MSE_train 25.1060222464
 MSE_test 34.6000723135

 brute force
 MSE_train for brute force: 25.1060222464
 MSE_test for brute force: 34.6000723135
 columns 4,6,11,13

 3.4
 MSE after polynomial feature expansion on train set: 5.05978429711
 MSE after polynomial feature expansion on test set: 14.5553049727