

# Machine Learning, Spring 2018

## Homework 3

Due on 23:59 Apr 17, 2018  
Send to `cs282_01@163.com`  
with subject "Chinese name+student number+HW3"

### 1 Perceptron

- (1) *Proof.*  $\because \mathbf{w}(t+1) \leftarrow \mathbf{w}(t) + y(t)\mathbf{x}(t)$ , we have

$$\begin{aligned} y(t)\mathbf{w}^T(t+1)\mathbf{x}(t) &= y(t)[\mathbf{w}^T(t) + y(t)\mathbf{x}^T(t)]\mathbf{x}(t) \\ &= y(t)\mathbf{w}^T(t)\mathbf{x}(t) + |y(t)| \|\mathbf{x}(t)\|_2^2 \\ &\geq y(t)\mathbf{w}^T(t)\mathbf{x}(t) \end{aligned} \quad (1)$$

Since,  $\mathbf{x}(t)$  is misclassified by  $\mathbf{w}(t)$

$$\therefore y(t)\mathbf{w}^T(t)\mathbf{x}(t) < 0$$

So, we have  $\|\mathbf{x}(t)\|_2 \neq 0$

Therefore,  $y(t)\mathbf{w}^T(t+1)\mathbf{x}(t) > y(t)\mathbf{w}^T(t)\mathbf{x}(t)$ . □

- (2) Suppose  $\mathbf{x}(t)$  was misclassified.

$\because \mathbf{w}(t+1) \leftarrow \mathbf{w}(t) + y(t)\mathbf{x}(t)$ , we have

$$\mathbf{w}^T(t+1)\mathbf{x}(t) = \mathbf{w}^T(t)\mathbf{x}(t) + y(t)(\mathbf{x}(t) \cdot \mathbf{x}(t))$$

If  $\mathbf{x}(t)$  was incorrectly classified as negative, then  $y(t) = +1$ . It follows that the new dot product increased by  $\mathbf{x}(t) \cdot \mathbf{x}(t)$  (which is positive). The boundary moved in the right direction as far as  $\mathbf{x}(t)$  is concerned, therefore. Conversely, if  $\mathbf{x}(t)$  was incorrectly classified as positive, then  $y(t) = -1$ . It follows that the new dot product decreased by  $\mathbf{x}(t) \cdot \mathbf{x}(t)$  (which is positive). The boundary moved in the right direction as far as  $\mathbf{x}(t)$  is concerned, therefore.

### 2 Understanding logistic regression

Answer:

(1) Suppose  $p = \mathbb{P}(\mathbf{y} = 1|\mathbf{x})$

$$\begin{aligned}\text{logit}(p) &= \log \frac{p}{1-p} = \mathbf{w}^T \mathbf{x}^{(i)} \\ p &= \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x}^{(i)}}} \\ &= \sigma(\mathbf{w}^T \mathbf{x}^{(i)})\end{aligned}\tag{2}$$

(2) As long as  $\mathbf{x}$  is not zero, the squared error loss with respect to  $\mathbf{w}$  will be non-convex. It's hard to optimize. Whereas the log loss is convex.

(3)

$$\mathbb{P}(y^{(i)} = k | \mathbf{x} = \mathbf{x}^{(i)}) = \mathbb{1}\{y^{(i)} = k\} \frac{e^{f_{y^{(i)}}}}{\sum_{j=1}^3 e^{f_j}}$$

in there  $f(\mathbf{x}^{(i)}, \mathbf{W}) = \mathbf{W}^T \mathbf{x}^{(i)}, k = 1, 2, 3$

$$L(y^{(1)}, \dots, y^{(N)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}; \mathbf{W}) = \prod_{i=1}^N \mathbb{1}\{y^{(i)} = k\} \frac{e^{f_{y^{(i)}}}}{\sum_{j=1}^3 e^{f_j}}$$

The Negative Log-Likelyhood of the  $N$  samples as follows:

$$-\ln L = -\sum_{i=1}^N \mathbb{1}\{y^{(i)} = k\} (f_{y^{(i)}} - \ln \sum_{j=1}^3 e^{f_j})$$

### 3 Regularization

(1) *Proof.* For sigle sample:

$$p(t|x, \mathbf{w}, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{[t-y(x, \mathbf{w})]^2}{2\sigma^2}}$$

For all samples, likelihood is :

$$L = \prod_{i=1}^N p(t_i | x_i, \mathbf{w}, \sigma) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t_i - y_i)^2}{2\sigma^2}}$$

log-likelihood:

$$\ln L = -N \ln \sqrt{2\pi}\sigma - \frac{1}{2\sigma^2} \sum_{i=1}^N (t_i - y_i)^2$$

$\therefore$  Maximizing the log likelihood is equal to minimizing the sum-of- squares error function.  $\square$

(2) According to Bayes' theorem, for single sample:

$$\begin{aligned}
p(t, \mathbf{w}|x, \alpha, \sigma) &= p(t|\mathbf{w}, x, \alpha, \sigma) \cdot p(\mathbf{w}|x, \alpha, \sigma) \\
&= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{[t-y(x, \mathbf{w})]^2}{2\sigma^2}} \cdot \frac{1}{\sqrt{(2\pi)^D |\alpha I|}} e^{-\frac{1}{2} \mathbf{w}^T (\alpha I)^{-1} \mathbf{w}} \\
&= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{[t-y(x, \mathbf{w})]^2}{2\sigma^2}} \cdot \frac{1}{\sqrt{(2\pi\alpha)^D}} e^{-\frac{1}{2\alpha} \|\mathbf{w}\|_2^2}
\end{aligned} \tag{3}$$

For all samples, likelihood is :

$$\begin{aligned}
L &= \prod_{i=1}^N p(t_i, \mathbf{w}|x_i, \alpha, \sigma) \\
&= \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^N \left(\frac{1}{\sqrt{(2\pi\alpha)^D}}\right)^N \exp\left(\sum_{i=1}^N \left(-\frac{1}{2\sigma^2} (t_i - y_i)^2 - \frac{1}{2\alpha} \|\mathbf{w}\|_2^2\right)\right)
\end{aligned} \tag{4}$$

$\frac{1}{N}$  log-likelihood:

$$\frac{1}{N} \ln L = C - \frac{1}{N} \sum_{i=1}^N \left(\frac{1}{2\sigma^2} (t_i - y_i)^2 + \frac{1}{2\alpha} \|\mathbf{w}\|_2^2\right), C \text{ is constant}$$

The formulation of the prediction problem is:

$$\text{minimize } \frac{1}{N} \sum_{i=1}^N \left(\frac{1}{2\sigma^2} (t_i - \mathbf{w}^T x_i)^2 + \frac{1}{2\alpha} \|\mathbf{w}\|_2^2\right)$$

## 4 Program Logistic regression in matlab

(1) Fig.1. shows the norm of gradient corresponding to three algorithms. All of three algorithms were terminated when the norm of gradient small than 0.00001 (In last question, I change it to 0.001 for BFGS and add normalization, so plot maybe different as shown.) and keep in one Matlab function file named *lr\_lindq.m*. Each algorithm was called by a extra arguments *type*, 0, 1 and 2 corresponding to Negative gradient, Newton's direction and BFGS, respectively.

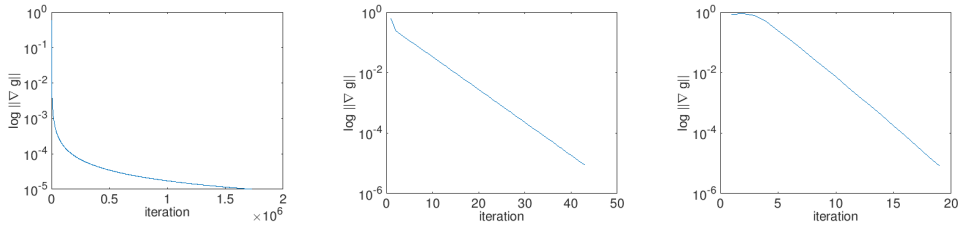


Figure 1: Negative gradient(left), Newton's direction(middle) and BFGS(right)

- (2) Terminate the BFGS when the norm of gradient small than 0.001. Accuracy of prediction in training set is equal to 0.9699.
- (3) Consider very small amount of champion data, I amplified some entries in Hessian matrix relate to champion data. Since  $H = X^T D X$ , I zoomed in the entries in  $D$  with 20000x, where is champion data, say  $y_i = 1$ . Terminate the BFGS when the norm of gradient small than 0.001. Accuracy of prediction in training set is equal to 0.9725. See code in *lr\_lindq.m* when *type* = 3.