

## Solutions to Homework 3

Lecturer: Inderjit Dhillon

Date Due: October 27, 2009

**Keywords:** *Classification, Logistic Regression, Perceptron, Support Vector Machines*

1. (4 points) Note that  $\frac{d\sigma(\mathbf{w}^T \mathbf{x}_n)}{d\mathbf{w}} = \sigma(\mathbf{w}^T \mathbf{x}_n)(1 - \sigma(\mathbf{w}^T \mathbf{x}_n))\mathbf{x}_n$ . Thus,

$$\frac{d \log \sigma(\mathbf{w}^T \mathbf{x}_n)}{d\mathbf{w}} = \frac{1}{\sigma(\mathbf{w}^T \mathbf{x}_n)} \frac{d\sigma(\mathbf{w}^T \mathbf{x}_n)}{d\mathbf{w}} = (1 - \sigma(\mathbf{w}^T \mathbf{x}_n))\mathbf{x}_n, \quad (1)$$

$$\frac{d \log(1 - \sigma(\mathbf{w}^T \mathbf{x}_n))}{d\mathbf{w}} = \frac{-1}{1 - \sigma(\mathbf{w}^T \mathbf{x}_n)} \frac{d\sigma(\mathbf{w}^T \mathbf{x}_n)}{d\mathbf{w}} = -\sigma(\mathbf{w}^T \mathbf{x}_n)\mathbf{x}_n. \quad (2)$$

Using (1) and (2),

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = \sum_n (z_n - y_n) \mathbf{x}_n,$$

where  $z_n = \sigma(\mathbf{w}^T \mathbf{x}_n)$ . Now, Hessian of  $E(\mathbf{w})$  w.r.t  $\mathbf{w}$  is given by:

$$H = \nabla_{\mathbf{w}}^2 E(\mathbf{w}) = \sum_n z_n(1 - z_n) \mathbf{x}_n \mathbf{x}_n^T$$

Note that  $0 < z_n < 1$ . Thus, for any  $\mathbf{u}$

$$\mathbf{u}^T H \mathbf{u} = \sum_n z_n(1 - z_n) (\mathbf{u}^T \mathbf{x}_n)^2 \geq 0. \quad (3)$$

Hence  $H \succeq 0$ , which implies  $E(\mathbf{w})$  is a convex function. Note that the equality in (3) holds if and only if  $\mathbf{u} = \mathbf{0}$  or  $\mathbf{u} \in \text{Null-Space}(X)$ , where  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$ . Thus, if  $\dim(\text{Null-Space}(X)) = 0$ ,  $H$  is strictly positive definite matrix. Hence  $E(\mathbf{w})$  is a strictly convex function and has a unique minimum.

2. (6 points)

- (a) (2 points) Using the update for  $\mathbf{w}_t$ ,

$$\mathbf{w}^{*T} \mathbf{w}_t = \mathbf{w}^{*T} \mathbf{w}_{t-1} + y_{t-1} (\mathbf{w}^{*T} \mathbf{x}_{t-1}).$$

Since  $\mathbf{w}^*$  is a separating hyperplane with margin  $\gamma$ ,  $y_{t-1} (\mathbf{w}^{*T} \mathbf{x}_{t-1}) \geq \gamma, \forall t$ . Thus,

$$\mathbf{w}^{*T} \mathbf{w}_t \geq \mathbf{w}^{*T} \mathbf{w}_{t-1} + \gamma.$$

Hence, using induction it can be easily seen that  $\mathbf{w}^{*T} \mathbf{w}_t \geq t\gamma$ . Note that  $\mathbf{w}_0 = \mathbf{0}$ , hence base case is satisfied trivially. In fact, assuming  $\mathbf{w}_0 = \mathbf{0}$  is not necessary in completing the proof. Without loss of generality, we can always find a  $\mathbf{w}_0$  such that  $\mathbf{w}^{*T} \mathbf{w}_0 \geq 0$ : first randomly initialize  $\mathbf{w}_0$ ; then change the sign of  $\mathbf{w}_0$  if the number of misclassified instances are more than half of the total number of training instances by the initial separating hyperplane.

- (b) (2 points) Using the update for  $\mathbf{w}_t$ :

$$\begin{aligned}\|\mathbf{w}_t\|^2 &= \|\mathbf{w}_{t-1} + y_{t-1}\mathbf{x}_{t-1}\|^2, \\ &= \|\mathbf{w}_{t-1}\|^2 + 2y_{t-1}\mathbf{w}_{t-1}^T\mathbf{x}_{t-1} + y_{t-1}^2\|\mathbf{x}_{t-1}\|^2, \\ &\leq \|\mathbf{w}_{t-1}\|^2 + \|\mathbf{x}_{t-1}\|^2, \\ &\leq \|\mathbf{w}_{t-1}\|^2 + R^2,\end{aligned}$$

since  $y_{t-1}\mathbf{w}_{t-1}^T\mathbf{x}_{t-1} \leq 0$  whenever an update is made,  $y_t^2 = 1$ , and  $\|\mathbf{x}_t\| \leq R \forall t$ .

Hence, using induction it can be easily shown that  $\|\mathbf{w}_t\|^2 \leq tR^2$ . Again note that base case holds trivially as  $\mathbf{w}_0 = \mathbf{0}$ .

- (c) (2 points) Using (a) and (b):

$$1 \geq \frac{\mathbf{w}^{*T}\mathbf{w}_t}{\|\mathbf{w}^*\|\|\mathbf{w}_t\|} \geq \frac{t\gamma}{\|\mathbf{w}^*\|\sqrt{t}R}$$

Thus,  $t \leq \frac{R^2\|\mathbf{w}^*\|^2}{\gamma^2}$ . This implies that the Perceptron algorithm converges to a separating hyperplane  $\mathbf{w}^*$  in at most  $\frac{R^2\|\mathbf{w}^*\|^2}{\gamma^2}$  steps.

3. (6 points)

- (a) (2 points) The sample code for least squares regression is given below.

```
function [W] = linreg(X,y)
% Solve a linear regression for classification
%
% Input:
%   X: N by d matrix
%   y: N by 1 vector
%
% Output:
%   W: d+1 by K matrix (K: number of classes)
%

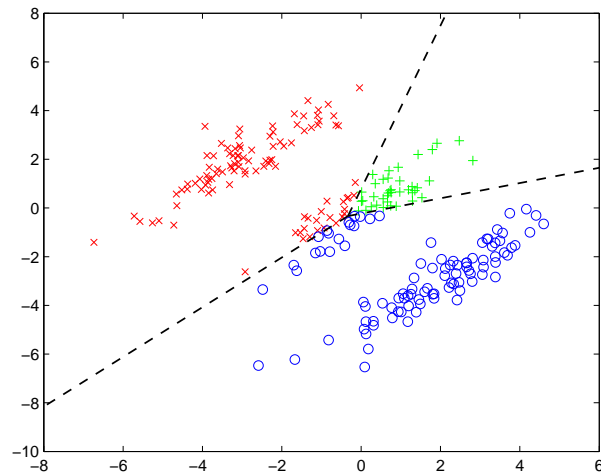
K = length(unique(y));
I = eye(K);
Y = I(y,:);
Xhat = [ones(size(y)) X];
W = Xhat'*Xhat\'(Xhat'*Y);

end % linreg function
```

Figure 1 shows the plot of data points according to the predicted labels by least squares regression, from which we observe that the region of input space assigned to the center class is very small and most of the points from that class are misclassified.

- (b) (4 points) The sample code for logistic regression is given below. Note that since the Hessian matrix  $H$  is positive semi-definite and hence rank deficient we can use the technique introduced in homework 1 to compute the inverse. In the sample code, the `pinv` Matlab function is used.

```
function [W] = logreg(X,y)
% Solve a logistic regression for classification
```



**Figure 1:** Prediction by least squares regression.

```
%
% Input:
%   X: N by d matrix
%   y: N by 1 vector
%
% Output:
%   W: d+1 by K matrix (K: number of classes)
%
```

```
K = length(unique(y));
I = eye(K);
Y = I(y,:);
Xhat = [ones(size(y)) X];
M = size(Xhat,2);

W = zeros(M,K);
maxi = 100; eta = 1e-8;
for t = 1:maxi
    Z = exp(Xhat*W);
    Z = diag(sparse(1./sum(Z,2)))*Z;
    obj = -sum(sum(Y.*log(Z)));
    if t > 1
        if abs(old_obj-obj) <= eta
            break;
        end
    end
    old_obj = obj;
    fprintf('iter %d: obj %f\n',t,obj);
    gW = Xhat'*(Z-Y);
```

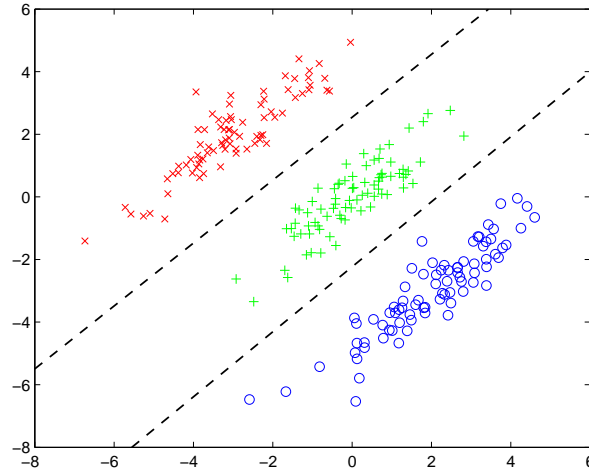
```

H = zeros(K*M);
for i = 1:K
    for j = 1:K
        Zb = Z(:,i).*(I(i,j)-Z(:,j));
        H((i-1)*M+1:i*M,(j-1)*M+1:j*M) = Xhat'*diag(Zb)*Xhat;
    end
end
W(:) = W(:)-pinv(H)*gW(:);
end

end % logreg function

```

Figure 2 shows the plot of data points according to the predicted labels by logistic regression, from which we observe that logistic regression can get a much better classification results on the training data.



**Figure 2:** Prediction by logistic regression.

4. (10 points)

(a) (2 points) As only  $\alpha_1$  and  $\alpha_2$  are updated and  $y_1^2 = 1$ ,

$$\sum_i y_i \alpha_i = 0 \Rightarrow y_1 \alpha_1 + y_2 \alpha_2 = y_1 \bar{\alpha}_1 + y_2 \bar{\alpha}_2 \Rightarrow \bar{\alpha}_1 = \alpha_1 + y_2 y_1 \alpha_2 - y_2 y_1 \bar{\alpha}_2 \geq 0. \quad (4)$$

If  $y_1 = y_2$ ,  $\alpha_1 + \alpha_2 \geq \bar{\alpha}_2$ . If  $y_1 \neq y_2$ ,  $\alpha_1 - \alpha_2 + \bar{\alpha}_2 \geq 0 \Rightarrow \bar{\alpha}_2 \geq \alpha_2 - \alpha_1$ .

(b) (2 points) From (4) it is clear that  $\alpha_1 + y_2 y_1 \alpha_2 = \gamma$ . Now,

$$W(\alpha_1, \alpha_2) = \alpha_1 + \alpha_2 - \frac{1}{2} (\alpha_1^2 K_{11} + 2s\alpha_1 \alpha_2 K_{12} + \alpha_1 y_1 v_1 + \alpha_2^2 K_{22} + \alpha_2 y_2 v_2) + \text{constant}.$$

Substituting  $\alpha_1 = \gamma - s\alpha_2$  and simplifying, we get:

$$\begin{aligned}
 W(\alpha_2) &= \gamma - s\alpha_2 + \alpha_2 - \frac{1}{2} K_{11} (\gamma - s\alpha_2)^2 - \frac{1}{2} K_{22} \alpha_2^2 \\
 &\quad - sK_{12} (\gamma - s\alpha_2) \alpha_2 - y_1 (\gamma - s\alpha_2) v_1 - y_2 \alpha_2 v_2 + \text{constant}
 \end{aligned}$$

- (c) (3 points) Setting  $\frac{dW}{d\alpha_2} = 0$ , we get:

$$\bar{\alpha}_2 = \frac{1}{d_{12}} (s(K_{11} - K_{12})\gamma + y_2(v_1 - v_2) + 1 - s)$$

Check that  $\frac{d^2W}{d\alpha_2^2} = -d_{12}$  for the above given  $\alpha_2$ . Now  $K$  is positive definite, hence  $d_{12} > 0$ .

- (d) (3 points) Note that  $\gamma = \alpha_1 + s\alpha_2$ . Substituting value of  $\gamma$  into the expression for  $\bar{\alpha}_2$  given above and simplifying, we get:

$$\bar{\alpha}_2 = \alpha_2 + \frac{y_2(E_1 - E_2)}{d_{12}},$$

where  $E_i = f(\mathbf{x}_i) - y_i = (\sum_{j=1}^N \alpha_j y_j K_{ij} + w_0) - y_i$ .

Note that  $\bar{\alpha}_2$  given by above update need not satisfy the constraints. So, we clip  $\bar{\alpha}_2$  so as to obtain a feasible dual solution.

After the update  $\bar{\alpha}_2$  should be greater than 0. Thus,

$$\bar{\alpha}_2 = \max(0, \bar{\alpha}_2).$$

Let  $y_1 = y_2$ . Using part (a),  $\bar{\alpha}_2 \leq \alpha_1 + \alpha_2$ . Thus  $\bar{\alpha}_2 = \min(\alpha_1 + \alpha_2, \bar{\alpha}_2) = \max(0, \min(\alpha_1 + \alpha_2, \bar{\alpha}_2))$ .

Now, let  $y_1 \neq y_2$ . Using part (a),  $\bar{\alpha}_2 \geq \alpha_2 - \alpha_1$ . Thus  $\bar{\alpha}_2 = \max(\bar{\alpha}_2, \alpha_2 - \alpha_1) = \max(0, \bar{\alpha}_2, \alpha_2 - \alpha_1)$ .

Using part (b),  $\bar{\alpha}_1 + s\bar{\alpha}_2 = \alpha_1 + s\alpha_2$ . This implies  $\bar{\alpha}_1 = \alpha_1 + s(\alpha_2 - \bar{\alpha}_2)$ . Note that this update gives feasible  $\bar{\alpha}_1$ .