

CS260, Winter 2017  
Problem Set 4: Linear Regression and Perceptron  
Due 2/1/2017

Ray Zhang

1    **Q1**    Jan 21, 2017

(a) Problem 1a

**Solution:** Solution to problem 1a

The log-likelihood of the data is a function of  $\beta$  and  $\sigma_n$  for all  $n = 1, 2, 3, \dots, N$ .

For a single data point:

- $y_n = x_n^T \beta + \epsilon_n$ .
- $\epsilon_n \sim N(0, \sigma_n)$ .
- $P(y_n|x_n) = N(x_n^T \beta, \sigma_n)$ .

Log probability of a single data point:

$$\log P(y_n|x_n) = -\frac{1}{2} \log 2\pi - \log \sigma_n - \frac{(y_n - x_n^T \beta)^2}{2\sigma_n^2}$$

If we express this likelihood function as a function of  $x_n$ ,  $\beta$  and  $\sigma_n$ , then we have, for all data:

$$\mathcal{L}(\beta|\sigma, x) = -\sum_{n=1}^N \frac{1}{2} \log 2\pi + \log \sigma_n + \frac{(y_n - x_n^T \beta)^2}{2\sigma_n^2}$$

(b) Problem 1b

**Solution:** Solution to problem 1b

Maximizing  $\mathcal{L}(\beta|\sigma, x)$  is the same as minimizing the negative log probability,  $-\mathcal{L}(\beta|\sigma, x)$ :

$$\operatorname{argmin}_{\beta} -\mathcal{L}(\beta|\sigma, x)$$

To find minima or maxima, we take the gradient of  $\beta$ :

$$\begin{aligned}
& \nabla_{\beta} \left( \sum_{n=1}^N \frac{1}{2} \log 2\pi + \log \sigma_n + \frac{(y_n - x_n^T \beta)^2}{2\sigma_n^2} \right) \\
&= \sum_{n=1}^N \frac{2(y_n - x_n^T \beta)}{2\sigma_n^2} x_n \\
&= \sum_{n=1}^N \frac{y_n - x_n^T \beta}{\sigma_n^2} x_n \\
&= \sum_{n=1}^N a_n x_n \\
&= X^T a = 0
\end{aligned}$$

where  $a = \Sigma^{-1}(y - X\beta)$ , where  $\Sigma_{ij} = 0$  if  $i \neq j$ , and  $\sigma_n^2$  if  $i = j = n$ .

Substituting:

$$\begin{aligned}
& X^T \Sigma^{-1}(y - X\beta) = 0 \\
&= X^T \Sigma^{-1}y - X^T \Sigma^{-1}X\beta \\
& X^T \Sigma^{-1}X\beta = X^T \Sigma^{-1}y \\
& \beta = (X^T \Sigma^{-1}X)^{-1} X^T \Sigma^{-1}y, \text{ which is what we wanted.}
\end{aligned}$$

## 2 Q2

(a) Problem 2a

**Solution:** Solution to problem 2a

We want to regularize with the following:

$(\beta_i - \beta_{i+1})^2$ , which is an L2 regularization on the size of adjacent values in  $\beta$ .

Combining this with the original L2 regularization in our cost function:

$$\mathcal{L}(\beta, \lambda_1, \lambda_2) = \|y - X\beta\|^2 + \lambda_1 \beta^T \beta + \lambda_2 (\beta - \beta_{shifted})^T (\beta - \beta_{shifted})$$

To construct  $\beta_{shifted}$ , we will use a shift matrix:  $U$ , which stands for upper shift, where  $U_{ij} = 1$  if  $i + 1 = j$ , and 0 otherwise.

To express it in differentiable form:

$$\mathcal{L}(\beta, \lambda_1, \lambda_2) = (y - X\beta)^T (y - X\beta) + \lambda_1 \beta^T \beta + \lambda_2 (\beta - U\beta)^T (\beta - U\beta)$$

(b) Problem 2b

**Solution:** Solution to problem 2b

To find the closed form, we solve for  $\beta$ :

$$\begin{aligned}
& y^T y - y^T X \beta - \beta^T X^T y + \beta^T X^T X \beta + \lambda_1 \beta^T \beta + \lambda_2 (\beta^T \beta - \beta^T U \beta - \beta^T U^T \beta - \beta^T U^T U \beta) \\
&= y^T y - 2y^T X \beta + \beta^T X^T X \beta + \lambda_1 \beta^T \beta + \lambda_2 (\beta^T \beta - 2\beta^T U \beta + \beta^T U^T U \beta) \\
& \nabla_{\beta} \mathcal{L} = -2y^T X + 2X^T X \beta + 2\lambda_1 \beta + \lambda_2 (2\beta - 2U \beta + 2U^T U \beta) = 0 \\
& y^T X = X^T X \beta + \lambda_1 \beta + \lambda_2 \beta - \lambda_2 U \beta + \lambda_2 U^T U \beta \\
& y^T X = (X^T X + \lambda_1 I + \lambda_2 (I - U + U^T U)) \beta \\
& \beta = (X^T X + \lambda_1 I + \lambda_2 (I - U + U^T U))^{-1} y^T X.
\end{aligned}$$

### 3 Q3

We will use a lagrangian to make sure that the optimization will be inside of the subspace spanned by  $A\beta = b$ :

$$\mathcal{L}(\beta, \lambda) = \|X\beta - y\|^2 + \lambda(A\beta - b)$$

Expanding this gives us:

$$\mathcal{L}(\beta, \lambda) = y^T y - 2X^T y \beta + \beta^T X^T X \beta + \lambda A \beta - \lambda b$$

Differentiating:

$$\nabla_{\beta} \mathcal{L} = -2X^T y + 2(X^T X) \beta + \lambda A = 0$$

$$X^T y = X^T X \beta + \frac{1}{2} \lambda A$$

$$\beta = (X^T X)^{-1} (X^T y - \frac{\lambda}{2} A)$$

Plugging  $\beta$  back into the constraint:

$$A\beta = b$$

$$A(X^T X)^{-1} (X^T y - \frac{\lambda}{2} A) = b$$

$$A(X^T X)^{-1} X^T y - A(X^T X)^{-1} \frac{\lambda}{2} A = b$$

Solving for  $\lambda$  gives us:

$$\frac{\lambda}{2} = (A(X^T X)^{-1} A)^{-1} (A(X^T X)^{-1} X^T y - b)$$

$$\lambda = 2(A(X^T X)^{-1} A)^{-1} (A(X^T X)^{-1} X^T y - b)$$

Plugging back to solve for  $\beta$ , we get:

$$\beta = (X^T X)^{-1} (X^T y - \frac{\lambda}{2} A)$$

$$\beta = (X^T X)^{-1} (X^T y - (A(X^T X)^{-1} A)^{-1} (A(X^T X)^{-1} X^T y - b) A)$$

## 4 Q4

Originally, our update rule was too aggressive:

$\text{sign}(w^T x_n) \neq y_n$ , then update using  $w^{t+1} = w^t + y_n x_n$ .

Our point was to move the new prediction on the other side of the hyperplane, such that:

$$y_n w^T x_n < 0, \text{ update to } y_n (w + y_n x_n)^T x_n = y_n w^T x_n + y_n^2 x_n^T x_n.$$

However, because it's too aggressive, we often get an update value where  $y_n (w^{t+1})^T x_n > 0$ . We will change our update rule like following:

$$y_n w^T x_n = \epsilon, \text{ where } \epsilon < 0.$$

We want to update such that:  $w^{t+1} = w^t + v$  for some  $v$ .

$$y_n (w + v)^T x_n = 0 = \epsilon - \epsilon$$

Thus,  $y_n v^T x_n = -\epsilon$ . That will be our lagrangian constraint.

We want to minimize  $\|v\|^2$  so that the update is minimized in terms of harshness, but we still move the hyperplane to the correct boundaries.

Thus, our lagrangian is:

$$\mathcal{L}(w, \lambda) = \|v\|^2 + \lambda (y_n v^T x_n + \epsilon)$$

$$\nabla_v \mathcal{L} = 2v + \lambda (y_n x_n) = 0$$

Thus:

$$v = \frac{\lambda (y_n x_n)}{2}$$

Plugging into the constraint:

$$y_n v^T x_n + \epsilon = 0$$

$$y_n \frac{\lambda (y_n x_n)}{2} + \epsilon = 0$$

$$y_n \frac{\lambda (y_n x_n)}{2} = -\epsilon = -(y_n w^T x_n)$$

$$\lambda = \frac{-2 y_n w^T x_n}{y_n^2 x_n^T x_n}$$

Plugging back into  $v$ :

$$v = \frac{\lambda(y_n x_n)}{2} = \frac{\frac{-2y_n w^T x_n}{y_n^2 x_n^T x_n} (y_n x_n)}{2}$$

$$v = \frac{-w^T x_n}{x_n^T x_n} x_n$$

$$v = -proj_{x_n} w$$

Thus, the conclusive result is to update by some  $v$ , where  $v$  is the projection of  $w$  onto  $x_n$ . Plugging back in to the update rule:

$$w^{t+1} = w^t - proj_{x_n} w$$

$$y_n (w^{t+1})^T x_n = y_n w^T x_n - y_n (proj_{x_n} w)^T x_n$$

$$= y_n w^T x_n + y_n \frac{-w^T x_n}{x_n^T x_n} x_n^T x_n$$

$$= y_n w^T x_n - y_n w^T x_n = 0$$

This makes intuitive sense, since we added by a positive term to push the projection onto the positive side of the hyperplane originally, here we are also pushing it by a positive term.

So the final update rule is:

$$w^{t+1} = w^t - proj_{x_n} w \text{ if } sign(w^T x_n) \neq y_n$$