

Elements of Statistical Learning

Solutions to the Exercises

Yu Zhang, sjtuzy@gmail.com

November 25, 2009

Exercise 2.6 Consider a regression problem with inputs x_i and outputs y_i , and a parameterized model $f_\theta(x)$ to be fit by least squares. Show that if there are observations with tied or identical values of x , then the fit can be obtained from a reduced weighted least squares problem.

Proof For known heteroskedasticity (e.g., grouped data with known group sizes), use weighted least squares (WLS) to obtain efficient unbiased estimates. Fig.1 explains “Observations with tied or identical values of x ”. In

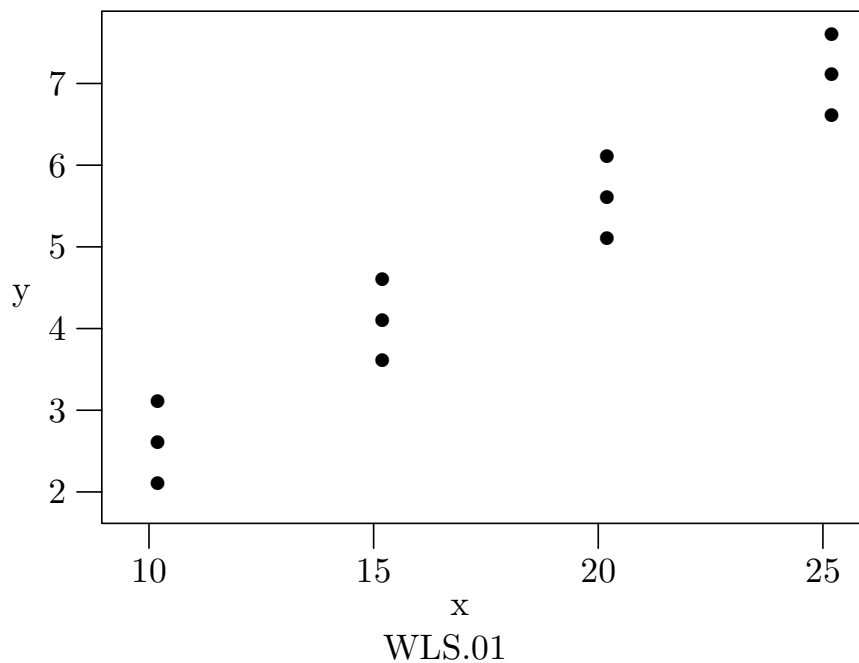


Figure 1: Observations with tied

the textbook, section 2.7.1 also explain this problem.

If there are multiple observation pairs $x_i, y_{il}, l = 1, 2, \dots, N_i$ at each value of x_i , the risk is limited as follows:

$$\begin{aligned}
\operatorname{argmin}_{\theta} \sum_i \sum_{l=1}^{N_i} (f_{\theta}(x_i) - y_{il})^2 &= \operatorname{argmin}_{\theta} \sum_i (N_i f_{\theta}(x_i)^2 - 2 \sum_{l=1}^{N_i} f_{\theta}(x_i) y_{il} \\
&\quad + \sum_{l=1}^{N_i} y_{il}^2) \\
&= \operatorname{argmin}_{\theta} \sum_i N_i \{(f_{\theta}(x_i) - \bar{y}_i)^2 + \text{Constant}\} \\
&= \operatorname{argmin}_{\theta} \sum_i N_i \{(f_{\theta}(x_i) - \bar{y}_i)^2\} \tag{1}
\end{aligned}$$

which is a weighted least squares problem.

Exercise 3.19 Show that $\|\hat{\beta}^{\text{ridge}}\|$ increases as its tuning parameter $\lambda \rightarrow 0$. Does the same property hold for the lasso and partial least squares estimates? For the latter, consider the “tuning parameter” to be the successive steps in the algorithm.

Proof Let $\lambda_2 < \lambda_1$ and β_1, β_2 be the optimal solution. We denote the loss function as follows

$$f_{\lambda}(\beta) = \|\mathbf{Y} - \mathbf{X}\beta\| + \lambda \|\beta\|_2^2$$

Then we have

$$\begin{aligned}
f_{\lambda_1}(\beta_2) + f_{\lambda_2}(\beta_1) &\geq f_{\lambda_2}(\beta_2) + f_{\lambda_1}(\beta_1) \\
\lambda_1 \|\beta_2\|_2^2 + \lambda_2 \|\beta_1\|_2^2 &\geq \lambda_2 \|\beta_2\|_2^2 + \lambda_1 \|\beta_1\|_2^2 \\
(\lambda_1 - \lambda_2) \|\beta_2\|_2^2 &\geq (\lambda_1 - \lambda_2) \|\beta_1\|_2^2 \\
\|\beta_2\|_2^2 &\geq \|\beta_1\|_2^2
\end{aligned}$$

So $\|\hat{\beta}^{\text{ridge}}\|$ increases as its tuning parameter $\lambda \rightarrow 0$.

Similarly, in lasso case, $\|\hat{\beta}\|_1$ increase as $\lambda \rightarrow 0$. But this can't guarantee the l_2 -norm increase. Fig.2 is a direct view of this property.

In partial least square case, it can be shown that the PLS algorithm is equivalent to the conjugate gradient method. This is a procedure that iteratively computes approximate solutions of $|\beta A = b|$ by minimizing the quadratic function

$$\frac{1}{2} \beta^{\top} A \beta - \mathbf{b}^{\top} \beta$$

along directions that are $|\mathbf{A}|$ -orthogonal (Ex.3.18). The approximate solution obtained after m steps is equal to the PLS estimator obtained after p iterations. The canonical algorithm can be written as follows:

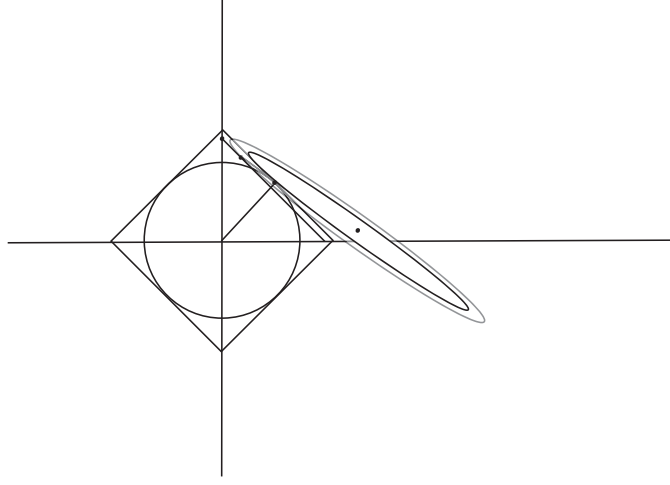


Figure 2: Lasso

1. Initialization $\beta_0 = 0, \mathbf{d}_0 = \mathbf{r}_0 = \mathbf{b} - \mathbf{A}\beta_0 = \mathbf{b}$
2. $a_i = \frac{\mathbf{d}_i^H \mathbf{r}_i}{\mathbf{d}_i^H \mathbf{A} \mathbf{d}_i}$
3. $\beta_{i+1} = \beta_i + a_i \mathbf{d}_i$
4. $\mathbf{r}_{i+1} = \mathbf{b} - \mathbf{A}\beta_{i+1} (= \mathbf{r}_i - a_i \mathbf{A} \mathbf{d}_i)$
5. $b_i = -\frac{\mathbf{r}_{i+1}^H \mathbf{A} \mathbf{d}_i}{\mathbf{d}_i^H \mathbf{A} \mathbf{d}_i}$
6. $\mathbf{d}_{i+1} = \mathbf{r}_{i+1} + b_i \mathbf{d}_i$

The squared norm of β_{j+1} can be written as

$$\|\beta_{j+1}\|^2 = \|\mathbf{r}_j\|^2 + 2a_j^2 \|\mathbf{d}_j\|^2 + a_j \mathbf{d}_j^H \beta_j \quad (2)$$

We need only shown that $a_j \mathbf{d}_j^H \beta_j > 0$.

$$\because \mathbf{r}_{j+1}^H \mathbf{d}_{j+1} = \mathbf{r}_{j+1}^H \mathbf{r}_{j+1} + b_j \mathbf{r}_{j+1}^H \mathbf{d}_j \text{ and } \langle \mathbf{r}_{j+1}, \mathbf{d}_j \rangle = 0 \quad (3)$$

$$\therefore a_j = \frac{\mathbf{d}_j^H \mathbf{r}_j}{\mathbf{d}_j^H \mathbf{A} \mathbf{d}_j} = \frac{\|\mathbf{r}_j\|^2}{\|\mathbf{d}_j\|_{\mathbf{A}}^2} > 0 (\mathbf{A} \text{ is positive definite.}) \quad (4)$$

$$\because \beta_j = \sum_{i=0}^j a_i \mathbf{d}_i \quad (5)$$

$$\therefore \mathbf{d}_j^H \beta_j = \sum_{i=0}^j a_i \mathbf{d}_j^H \mathbf{d}_i \quad (6)$$

Now we need to show that $\mathbf{d}_j^H \mathbf{d}_i > 0, i \neq j$. By Step 6, we have $\mathbf{d}_j = \mathbf{r}_j + \sum_{i=0}^{j-1} (\prod_{k=i}^{j-1} b_k) \mathbf{r}_i$.

$$\because \langle \mathbf{r}_i, \mathbf{r}_j \rangle = 0, i \neq j \quad (7)$$

$$\therefore b_k > 0 \rightarrow \mathbf{d}_j^H \mathbf{d}_i > 0 \quad (8)$$

$$\because \mathbf{A} \mathbf{d}_i = a_i^{-1} (\mathbf{r}_i - \mathbf{r}_{i+1}) \quad (9)$$

$$\therefore b_j = -\frac{\mathbf{r}_{j+1}^H (\mathbf{r}_j - \mathbf{r}_{j+1})}{a_j \|\mathbf{d}_j\|_A^2} = \frac{\|\mathbf{r}_{j+1}\|^2}{a_j \|\mathbf{d}_j\|_A^2} > 0 \quad (10)$$

So in PLS case, the $\|\beta\|_2^2$ increase with m.

Exercise 4.1 Show how to solve the generalized eigenvalue problem $\max a^\top \mathbf{B} a$ subject to $a^\top \mathbf{W} a = 1$ by transforming to a standard eigenvalue problem.

Proof Hence \mathbf{W} is the the common covariance matrix, we have \mathbf{W} is semi-positive definite. WLOG \mathbf{W} is positive definite, we have

$$\mathbf{W} = \mathbf{P}^2 \quad (11)$$

Let $b = \mathbf{P} a$, then the problem is

$$a^\top \mathbf{B} a = b^\top \mathbf{P}^{-1} \mathbf{B} \mathbf{P}^{-1} b = b^\top \mathbf{B}^* b \quad (12)$$

subject to $a^\top \mathbf{W} a = 1 = b^\top \mathbf{B}^* b = 1$. Now the problem transform to a standard eigenvalue problem.

Exercise 4.2 Suppose we have features $x \in \mathbb{R}^p$, a two-class response, with class sizes N_1, N_2 , and the target coded as $-N/N_1, N/N_2$.

(d) Show that this result holds for any (distinct) coding of the two classes.

W.L.O.G any distinct coding y'_1, y'_2 . Let $\hat{\beta} = (\beta, \beta_0)^\top$. Compute the partial deviation of the $\text{RSS}(\hat{\beta})$ we have

$$\frac{\partial \text{RSS}(\hat{\beta})}{\partial \beta} = -2 \sum_{i=1}^N (y_i - \beta_0 - \beta^\top x_i) x_i^\top = 0 \quad (13)$$

$$\frac{\partial \text{RSS}(\hat{\beta})}{\partial \beta_0} = -2 \sum_{i=1}^N (y_i - \beta_0 - \beta^\top x_i) = 0 \quad (14)$$

It yields that

$$\beta^\top \sum_i (x_i - \bar{x}_i) x_i^\top = \sum_i y_i x_i \quad (15)$$

$$\beta_0 = \frac{1}{N} \sum_i (y_i - \beta^\top x_i) \quad (16)$$

Then we get

$$\beta^\top \sum_i (x_i - \bar{x}_i) x_i^\top = \sum_i (y_i - \frac{1}{N} \sum_i y_i) x_i^\top \quad (17)$$

$$= N_1 y'_1 \hat{\mu}_1 + N_2 y'_2 \hat{\mu}_2 - \frac{N_1 y'_1 + N_2 y'_2}{N_1 + N_2} (N_1 \hat{\mu}_1 + N_2 \hat{\mu}_2) \quad (18)$$

$$= \frac{1}{N} (N_1 N_2 y'_2 (\hat{\mu}_2 - \hat{\mu}_1) - N_1 N_2 y'_1 (\hat{\mu}_2 - \hat{\mu}_1)) \quad (19)$$

$$= \frac{N_1 N_2}{N} (y'_2 - y'_1) (\hat{\mu}_2 - \hat{\mu}_1) \quad (20)$$

Hence the proof in (c) still holds.

Exercise 4.3 Suppose we transform the original predictors X to \hat{Y} via linear regression. In detail, let $\hat{Y} = X(X^\top X)^{-1} X^\top Y = X\hat{B}$, where Y is the indicator response matrix. Similarly for any input $x \in \mathbb{R}^p$, we get a transformed vector $\hat{y} = \hat{B}^\top x \in \mathbb{R}^K$. Show that LDA using \hat{Y} is identical to LDA in the original space.

Proof Prof.Zhang had given the solution about this problem, but only tow student notice this problem is tricky. Here I give the main part of the solution.

- First y_1, \dots, y_K must independent, then we have $r(\Sigma_y) = K$.
- Assume that $\hat{B}_{p \times K}, r(\hat{B}) = K < p, r(\Sigma_x) = p$. Note that

$$r(\hat{B}(\hat{B}^\top \Sigma_x \hat{B})^{-1} \hat{B}^\top) < r(\hat{B}) < r(\Sigma_x)$$

. Since $\hat{B}(\hat{B}^\top \Sigma_x \hat{B})^{-1} \hat{B}^\top \neq \Sigma_x^{-1}$ when $\dim(Y) < \dim(X)$.

- This problem are equal to prove that $\hat{B}(\hat{B}^\top \Sigma_x \hat{B})^{-1} \hat{B}^\top (\mu_k^x - \mu_l^x) = \Sigma_x^{-1} (\mu_k^x - \mu_l^x)$. Let $P = \Sigma_x \hat{B}(\hat{B}^\top \Sigma_x \hat{B})^{-1} \hat{B}^\top$, we have

$$P^2 = P \quad (21)$$

Hence P is projection matrix. Note that $Y = \{y_1, y_2, \dots, y_K\}$ is indi-

cator response matrix, we have

$$\begin{aligned}
\mu_k^x &= \frac{1}{N_k} X^\top y_k \\
\Sigma_x &= \sum_{k=1}^K \sum_{g_i=k} (x_i - \mu_i^k)(x_i - \mu_i^k)^\top / (N - K) \\
&= \frac{1}{N - K} \left(\sum_{k=1}^K \sum_{g_i=k} x_i x_i^\top - \sum_{k=1}^K N_k \mu_k^x (\mu_k^x)^\top \right) \\
&= \frac{1}{N - K} \left(X^\top X - \sum_{k=1}^K X^\top y_k y_k^\top \right) \\
&= \frac{1}{N - K} (X^\top X - X^\top Y Y^\top X)
\end{aligned}$$

Note that

$$P(N_1 \mu_1^x, N_2 \mu_2^x, \dots, N_K \mu_K^x) = P X^\top Y \quad (22)$$

So we need only to shown that $P X^\top Y = X^\top Y$.

$$P X^\top Y = \Sigma_x B (B^\top \Sigma_x B) B^\top X^\top Y$$

By the definition of B , we have

$$\begin{aligned}
B^\top X^\top Y &= Y^\top X (X^\top X)^{-1} X^\top Y \\
\Sigma_x B &= \frac{1}{N - K} ((X^\top X)(X^\top X)^{-1} X^\top Y - X^\top Y Y^\top X (X^\top X)^{-1} X^\top Y) \\
&= \frac{1}{N - K} X^\top Y (I - Y^\top X (X^\top X)^{-1} X^\top Y) \\
(B^\top \Sigma_x B)^{-1} &= (N - K) (Y^\top X (X^\top X)^{-1} (X^\top X - X^\top Y Y^\top X) (X^\top X)^{-1} X^\top Y)^{-1} \\
&= (N - K) (Y^\top X (X^\top X)^{-1} X^\top Y [I - Y^\top X (X^\top X)^{-1} X^\top Y])^{-1}
\end{aligned}$$

Let $Q = Y^\top X (X^\top X)^{-1} X^\top Y$ we have

$$B^\top X^\top Y = Q \quad (23)$$

$$\Sigma_x B = \frac{1}{N - K} X^\top Y (I - Q) \quad (24)$$

$$(B^\top \Sigma_x B)^{-1} = (N - K) (Q(I - Q))^{-1} \quad (25)$$

Hence $Q(I - Q)$ is invertible matrix, we have

$$K = r(Q(Q - I)) \leq r(Q) \Rightarrow r(Q) = K \quad (26)$$

So Q is invertible matrix, then it yields

$$(B^\top \Sigma_x B)^{-1} = (N - K) (I - Q)^{-1} Q^{-1} \quad (27)$$

Combine with (16),(17),(20), we have

$$PX^\top Y = \frac{1}{N-K} X^\top Y (I-Q)(N-K)(I-Q)^{-1} Q^{-1} Q \quad (28)$$

$$= X^\top Y \quad (29)$$

Since we have $\hat{B}(\hat{B}^\top \Sigma_x \hat{B})^{-1} \hat{B}^\top (\mu_k^x - \mu_l^x) = \Sigma_x^{-1} (\mu_k^x - \mu_l^x)$.