# Homework #2

Yongjin Shin

20090488

Industrial Management Engineering, POSTECH

dreimer@postech.ac.kr

March 21, 2018

**Problem 1.** Provide detailed description of ridge regression, including algorithm derivation. Implement your ridge regression to obtain the regularization path that was discussed in class.

**Solution** If the weights $w_i$s are not constrained, they can explode and hence are susceptible to very high variance. To control variance, we might regularize the coefficients to impose the ridge constraint:

$$minimize \quad \Sigma_{i=1}^{n}(y_i - w^T z_i)^2 \quad s.t. \quad \Sigma_{j=1}^{D} w_j^2 \leq B$$
$$\Leftrightarrow \quad minimize \quad (y - Zw)^T(y - Zw) \quad s.t. \quad \Sigma_{j=1}^{D} w_j^2 \leq B$$

where Z is assumed to be standardized (mean 0, unit variance) and y is assumed to be centered.

Then, the bounded constrained form is equivalent to loss function:

$$\begin{aligned}
L(w)_{l_2} = minmize \quad & \frac{1}{2}||y - Zw||_2^2 + \frac{\lambda}{2}||w||_2^2 \\
= \quad & \frac{1}{2}(y - Zw)^T(y - Zw) + \frac{\lambda}{2}w^T w \\
= \quad & \frac{1}{2}(y^T y - y^T Zw - w^T Z^T y + w^T Z^T Zw) + \frac{\lambda}{2}w^T w
\end{aligned} \quad (1)$$

Since the equation 1 is convex, and hence has a unique solution. Take derivatives, we can obtain:

$$\begin{aligned}
\frac{\partial L(w)_{l_2}}{\partial w} &= \frac{\partial}{\partial w}\left[\frac{1}{2}||y - Zw||_2^2 + \frac{\lambda}{2}||w||_2^2\right] \\
&= \frac{\partial}{\partial w}\left[\frac{1}{2}(y^T y - y^T Zw - w^T Z^T y + w^T Z^T Zw) + \frac{\lambda}{2}w^T w\right] \\
&= -Z^T y + Z^T Zw + \lambda w \\
&= -Z^T y + (Z^T Z + \lambda I)w
\end{aligned} \quad (2)$$

Therefore, the solution is now seen to be:

$$\hat{w}_\lambda^{ridge} = (Z^T Z + \lambda I)^{-1} Z^T y \tag{3}$$

$\lambda$ is the shrinkage parameter which controls the size of the coefficients. As $\lambda$ goes to 0, we obtain the least squares solutions. On the other hand, as $\lambda$ goes to $infinity$, we have $\hat{w}_{\lambda=\infty}^{ridge} = 0$ Notice that even if $Z^T Z$ is not invertible, inclusion of $\lambda$ makes problem non-singular. Ridge regression can be implemented as following algorithm and the full script **Appendix A**:

---

**Algorithm 1:** Ridge Regression

---
**input** : Standardized $Z$ (mean 0, unit variance), Centered $y$, $\lambda$
**output:** $\hat{w}_{lambda}^{ridge} = [\hat{w}_1, ..., \hat{w}_D]^T$, $df_\lambda$

**1** $Compute \quad tmp = (Z^T Z + \lambda I)^{-1} Z^T$;

**2** $\hat{w}_\lambda^{ridge} = tmp * y$;
**3** $df_\lambda = trace(Z * tmp)$;

---

A smoother matrix S is a linear operator satisfying: $\hat{y} = Sy$. The effective degrees of freedom (or effective number of parameters) is defined as:
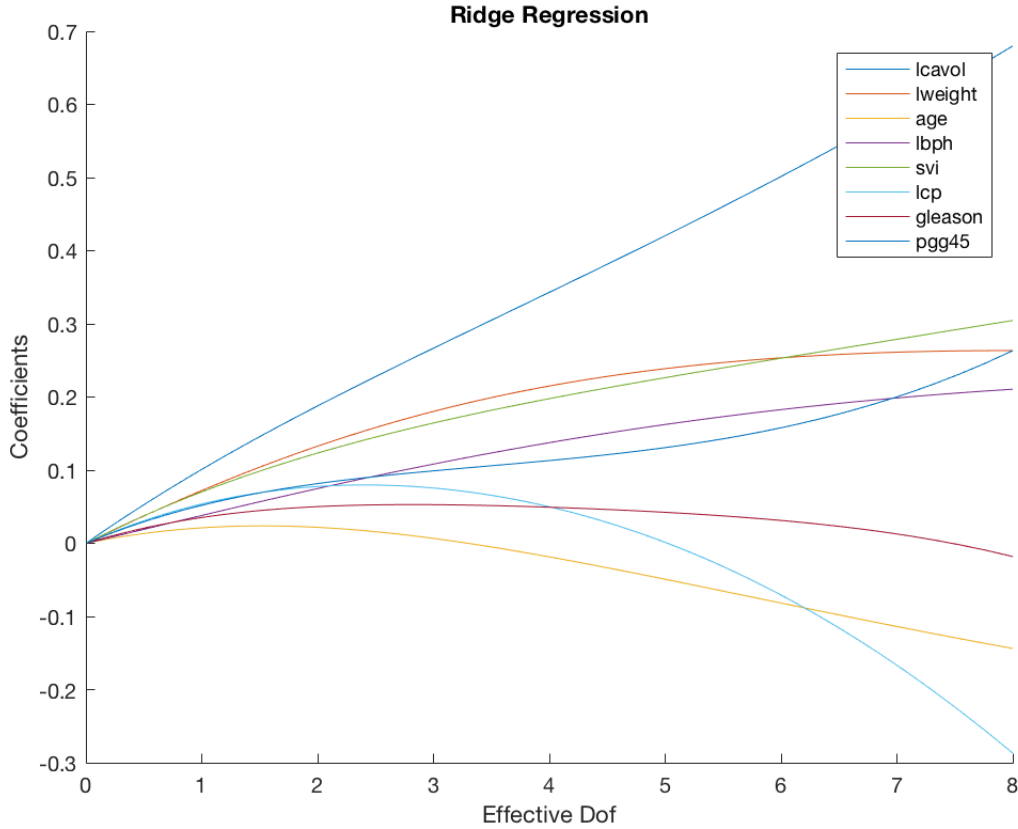
$$df(S) = tr(S) \tag{4}$$

In ridge regression, the fits are given by:

$$\hat{y} = Z(Z^T Z + \lambda I)^{-1} Z^T y \tag{5}$$

So the effective degrees of freedom (effective number of parameters) in ridge regression are given by:

$$df(\lambda) = tr(S_\lambda) = tr\left[Z(Z^T Z + \lambda I)^{-1} Z^T\right] = \Sigma_{i=1}^{D} \frac{d_i^2}{d_i^2 + \lambda} \tag{6}$$

The result of ridge regression path can be given as the below.

**Ridge Regression**

**Problem 2.** Provide detailed description of LASSO regression, including the derivation of the coordinate descent algorithm. Implement your LASSO regression to obtain the regularization path that was discussed in class.

**Solution** LASSO coefficients are the solutions to the $l_1$ optimization problem:

$$minimize \quad (y - Zw)^T(y - Zw) \quad s.t. \quad \Sigma_{i=1}^{D}|w_i| \leq B \tag{7}$$

This is equivalent to loss function:

$$
\begin{aligned}
L(w)_{l_1} &= \Sigma_{i=1}^{n}(y_i - w^T z_i)^2 + \lambda\Sigma_{i=1}^{D}|w_i| \\
&= (y - Zw)^T(y - Zw) + \lambda||w||_1
\end{aligned}
\tag{8}
$$

The coordinate descent will be used, which is to optimize the parameter one bye one tactic. Hence, $Zw$ can be separated with two terms; one with $d$th column and the other without $d$th column.

$$
\begin{aligned}
L(w)_{l_1} &= \frac{1}{2}\Sigma_{i=1}^{n}(y_i - w^T z_i)^2 + \lambda\Sigma_{i=1}^{D}|w_i| \\
&= \frac{1}{2}\Sigma_{i=1}^{n}(y_i - w_d^T z_{i,d} - w_{-d}^T z_{i,-d})^2 + \lambda\Sigma_{j=1}^{D}|w_j|
\end{aligned}
\tag{9}
$$

3

Now we only need to calculate the derivative of loss function with respect to weight $w$.

$$\frac{\partial L(w)_{l_1}}{\partial w_d} = \Sigma_{i=1}^n (y_i - w_d^T z_{i,d} - w_{-d}^T z_{i,-d})(-z_{i,d}) + \lambda \frac{\partial |w_d|}{\partial w_d}$$

$$= w_d^T \left( \Sigma_{i=1}^n z_{i,d}^2 \right) - \Sigma_{i=1}^n (y_i - w_{-d}^T z_{i,-d}) z_{i,d} + \lambda \frac{\partial |w_d|}{\partial w_d} \qquad (10)$$

$$= w_d^T \alpha_d + \beta_d + \lambda \frac{\partial |w_d|}{\partial w_d}$$

Computing $\frac{\partial ||w_d||_1}{\partial w_d}$ can be solved by sub-differentials.

$$\frac{\partial f}{\partial x} = \begin{cases} -1 & if \quad x < 0 \\ [-1, 1] & if \quad x = 0 \\ 1 & if \quad x > 0 \end{cases} \qquad (11)$$

Thus, the estimate of $w_d$ given the other parameters is calculated as:

$$\hat{w}_d^{lasso} = \begin{cases} \frac{\beta_d + \lambda}{\alpha_d} & if \quad \beta_d < -\lambda \\ 0 & if \quad \beta_d \in [-\lambda, \lambda] \\ \frac{\beta_d - \lambda}{\alpha_d} & if \quad \beta_d > \lambda \end{cases} \qquad (12)$$

LASSO can be implemented as following algorithm and the full script **Appendix B**:

---

**Algorithm 2:** The Cooridnate descent algorithm for LASSO

    **input** : Standardized $Z$ (mean 0, unit variance), Centered $y$, $\lambda$
    **output:** $\hat{w}^{lasso} = [\hat{w}_1, ..., \hat{w}_D]^T$

**1**   $\hat{w}_\lambda^{ridge} = ridge(Z, y, \lambda)$ ;                `/* initialize parameters w */`

**2**   **while** *Not Converged* **do**

**3**     $w_{old} = w$ ;                `/* Save the previous parameters into` $w_{old}$ `*/`

**4**     *Compute*    $\alpha_d = \Sigma_{i=1}^n z_{i,d}^2$;

**5**     *Compute*    $\beta_d = \Sigma_{i=1}^n (y_i - w_{-d}^T z_{i,-d}) z_{i,d}$;

**6**     **if** $w_d < -\lambda$ **then**

**7**        $w_d = \frac{\beta_d + \lambda}{\alpha_d}$;

**8**     **else if** $w_d > \lambda$ **then**

**9**        $w_d = \frac{\beta_d - \lambda}{\alpha_d}$;

**10**    **else**

**11**       $w_d = 0$;

**12**    **end**

      `/* Check the convergence condition`                         `*/`

**13**    **if** $max(abs(w_{old} - w)) < threshold$ **then**

**14**       return $\hat{w}^{lasso}$

**15**    **else**

**16**       *Continue*    *looping*

**17**    **end**

**18** **end**

---

If $t_0 = \Sigma_{i=1}^D |\hat{w}_i^{OLS}|$ (equivalently, $\lambda = 0$), we obtain no shrinkage. The path of solution is indexed by a fraction of shrinkage factor of $t_0$. The result of LASSO path can be given as the below.

Lasso Regression