

3.7.0 → Overview

→ overview (recall that OLS relies on ordinary least squares) to determine the coefficient estimates.

In the next subsection, we present alternatives that vary slightly from a OLS fit, namely:

→ shrinkage (regularization) methods → ridge + lasso

→ dimension reduction methods → PCA + partial least squares

→ these methods may assume that variables used will be centered or scaled

3.7.1 → Standardizing variables

→ centered variable is the result of subtracting the sample mean from a variable

→ a scaled variable is the result of dividing a variable by its sample standard deviation

→ a standardized variable is the result of first centering a variable, then scaling it.

→ Across various resources, it is common to see that scaling uses the bias-corrected (degree of n rather than n-1) estimate for the unbiased version. In many cases, this choice does not have a large impact on the final results.

→ If an OLS problem is not explicit, we normalize scaling w/ the biased sample standard deviation.

→ Example → original model: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon \Rightarrow \hat{\beta}_0 = 59.13$

$$\begin{aligned} \hat{\beta}_1 &= 2.29 \\ \hat{\beta}_2 &= -0.08 \end{aligned}$$

Different

→ Centered predictors: $y = \beta_0 + \beta_1(x_1 - \bar{x}_1) + \beta_2(x_2 - \bar{x}_2) + \epsilon \Rightarrow \hat{\beta}_0 = 40$

$$\begin{aligned} \hat{\beta}_1 &= -0.293 \\ \hat{\beta}_2 &= -0.007 \end{aligned}$$

Same

→ only the β_0 element, but it leads to the same fitted equation

$$\begin{aligned} y &= 40 - 0.293(x_1 - \bar{x}_1) - 0.007(x_2 - \bar{x}_2) \\ &= 40 - 0.293x_1 + 0.293\bar{x}_1 - 0.007x_2 + 0.007\bar{x}_2 \\ &= 59.13 - 0.293x_1 - 0.007x_2 \quad \checkmark \end{aligned}$$

Same from example

2 reasons

→ Scaled model: $y = \beta_0 + \beta_1 \left(\frac{x_1}{\bar{x}_1}\right) + \beta_2 \left(\frac{x_2}{\bar{x}_2}\right) + \epsilon \Rightarrow \hat{\beta}_0 = 59.13$

$$\begin{aligned} \hat{\beta}_1 &= -4.57 \\ \hat{\beta}_2 &= -0.03 \end{aligned}$$

Different -- --

→ Again, same fitted equation \Rightarrow model is scale invariant

$$\begin{aligned} y &= 59.13 - 4.57 \left(\frac{x_1}{\bar{x}_1}\right) - 0.03 \frac{x_2}{\bar{x}_2} \\ &= 59.13 - 2.29 x_1 - 0.007 x_2 \quad \checkmark \end{aligned}$$

→ when using OLS, note that:

→ Centering after scaling, the predictors doesn't automatically change the fitted equation

→ when centered predictors are used, the intercept estimate will equal the sample mean of the response, while the rest of the estimated regression coefficients remain unchanged

→ when scaled predictors are used, the intercept estimate does not change. However, the other coefficient estimates will change by a factor equal to the sample standard of the associated explanatory variable

→ standardized model: $y = \beta_0 + \beta_1 \left(\frac{x_1 - \bar{x}_1}{\bar{x}_1}\right) + \beta_2 \left(\frac{x_2 - \bar{x}_2}{\bar{x}_2}\right) + \epsilon \Rightarrow \hat{\beta}_0 = 40$

$$\begin{aligned} \hat{\beta}_1 &= -4.57 \\ \hat{\beta}_2 &= -0.03 \end{aligned}$$

... same fitted equation again \checkmark ...

→ Training / testing

→ Suppose we're given training set w/ $S = (x_1, \dots, x_n)$ & $S = (y_1, \dots, y_n)$,

when we predict to test set, we should use those SAME $\hat{\beta}$'s

to avoid making different standardizations

3.7.2 → Ridge regression

→ overview Recall that ordinary least squares (OLS) aims to find regression coefficients that minimize SSE

$$\boxed{\text{SSE} = \sum (y_i - \hat{y}_i)^2}$$

$\hat{y}_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$

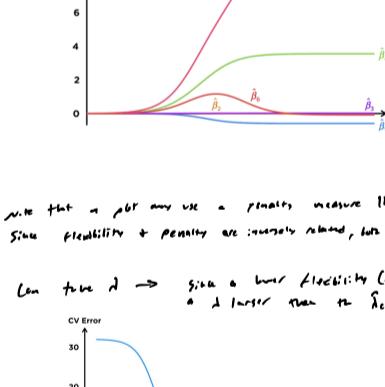
→ As p (i.e. flexibility) grows, SSE gets closer to having no bias. However, by the bias-variance tradeoff, this is at the expense of increasing the variance . One way to reduce variance (while allowing a bit more bias) is to restrict the possible values of the regression estimates. Specifically, we restrict by forcing the coefficient estimates to be closer to or shrink towards 0.

→ The first shrinkage method is called ridge regression. It minimizes the same SSE expression w/ an added restriction of

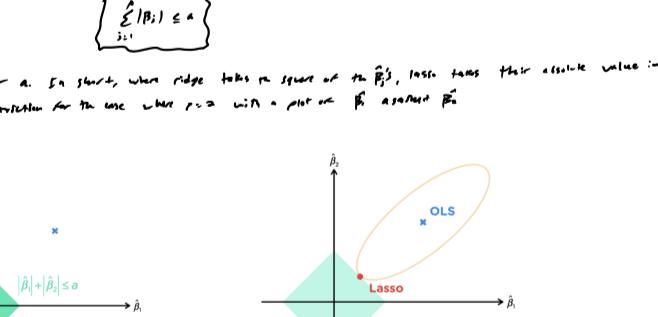
$$\boxed{\|\beta\|^2 \leq a}$$

for some constant a called the budget parameter. Notice that this restricts all the coefficient estimates except the intercept.

→ we can depict this restriction for the case where $p=2$. First consider a plot of possible $\hat{\beta}_1$ & $\hat{\beta}_2$ values.



→ In the plot above, the cross marks the values $\hat{\beta}_1$ & $\hat{\beta}_2$ under OLS. Since OLS minimizes the SSE, any other point in this plot will result in a larger SSE relative to the OLS. We illustrate SSE as a third dimension using contours. For each contour, every point along the contour has the same SSE value at those $\hat{\beta}_1$ & $\hat{\beta}_2$ coordinates. The further the center is from the origin, the larger the SSE.



→ The plots above include the region $\hat{\beta}_1^2 + \hat{\beta}_2^2 \leq a$ shaded in green. Ridge regression produces the coefficient estimates that minimize the SSE while restricted to the circle. Note the circle doesn't include the origin. The center that falls touches the edge where close satisfies the ridge regression criterion. Hence, their intersection is the closer to the origin (in terms of SSE) among all other points within the circle, & defines the ridge estimates.

→ The plots demonstrate how the values of $\hat{\beta}_1$ & $\hat{\beta}_2$ would tend to shrink towards 0 for ridge regression. But it is possible that no shrinkage occurs, as when the green circle is large enough to include the OLS contours.

→ In addition, the regression can be interpreted directly into an optimization problem. Ridge regression seeks to minimize the expression:

$$\boxed{\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \|\beta\|^2}$$

where λ is the tuning parameter that controls the strength of the shrinkage - note that:

as well also, that no shrinkage penalty $\Rightarrow \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_p$ will equal the old estimates.

→ as λ approaches 0, the entire expression is minimized when $\hat{\beta}_1, \dots, \hat{\beta}_p = 0$.

→ This means λ is inversely related to flexibility. λ allows us to select a flexibility level between the full model (large) & the flexible null model (0). It means we're finding the best value of λ in Cross-validation.

→ Furthermore, since might want to minimize $\sum_{i=1}^n \hat{\beta}_i^2$ instead, it can be simplified using L2 norm notation. First, let

$$\beta = \begin{pmatrix} \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_p \end{pmatrix}$$

→ Then the L2 norm of β is

$$\boxed{\|\beta\|_2 = \sqrt{\sum_{i=1}^p \hat{\beta}_i^2}}$$

→ Consequently, the restriction inequality & the expression to be minimized for ridge regression can be simplified as follows:

$$\boxed{\|\beta\|_2 \leq a}$$

$$\boxed{\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \|\beta\|_2^2}$$

→ as well also, that no shrinkage penalty $\Rightarrow \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_p$ will equal the old estimates.

→ as λ approaches 0, the entire expression is minimized when $\hat{\beta}_1, \dots, \hat{\beta}_p = 0$.

→ other important details regarding Ridge regression include:

→ The x_i 's should be scaled variables w/ the original. This puts $\hat{\beta}_1, \dots, \hat{\beta}_p$ on the same scale, which is crucial b/c the shrinkage parameter puts equal weight on each one. That means the ridge estimates are scale invariant.

→ As λ increases, $\|\beta\|_2^2$ must decrease to minimize SSE plus penalty. In ridge regression, it is possible for an individual $\hat{\beta}_i$ to increase in absolute value (i.e. move away from 0) as λ increases.

→ Practically, none of the ridge estimates will equal 0 for a large enough λ . This means ridge regression does not drop variables from a model.

→ just as with OLS when dealing w/ high dimensions, it's sum of the predictors will have a meaningful estimated coefficient.

→ Example

→ 6 variables → after scaling the predictors, running the ridge procedure, & reverting to scale back to the original units, we see how the ridge coefficients change as a function of flexibility in the graph below.

→ Shrinkage appears from right to left (increasing penalty from right to left)

($\lambda \rightarrow 0 \Rightarrow$ OLS estimate $\lambda \rightarrow \infty \Rightarrow$ null model)

→ reciprocal to inversely related to flexibility for plotting

→ note that in plot may use a penalty measure like λ in the horizontal axis instead of flexibility.

Since flexibility & penalty are inversely related, but this would lose descriptive details losing the same information.

→ On first look $\lambda \rightarrow$ gives a lower flexibility (i.e. higher penalty) w/ no loss in model accuracy or precision (separability).

→ In fact, when λ is too small, we'll end up with a lot of zero coefficients. This means ridge regression does not drop variables from a model.

→ just as with OLS when dealing w/ high dimensions, it's sum of the predictors will have a meaningful estimated coefficient.

→ note that in plot may use a penalty measure like λ in the horizontal axis instead of flexibility.

Since flexibility & penalty are inversely related, but this would lose descriptive details losing the same information.

→ On first look $\lambda \rightarrow$ gives a lower flexibility (i.e. higher penalty) w/ no loss in model accuracy or precision (separability).

→ In fact, when λ is too small, we'll end up with a lot of zero coefficients. This means ridge regression does not drop variables from a model.

→ just as with OLS when dealing w/ high dimensions, it's sum of the predictors will have a meaningful estimated coefficient.

→ note that in plot may use a penalty measure like λ in the horizontal axis instead of flexibility.

Since flexibility & penalty are inversely related, but this would lose descriptive details losing the same information.

→ On first look $\lambda \rightarrow$ gives a lower flexibility (i.e. higher penalty) w/ no loss in model accuracy or precision (separability).

→ In fact, when λ is too small, we'll end up with a lot of zero coefficients. This means ridge regression does not drop variables from a model.

→ just as with OLS when dealing w/ high dimensions, it's sum of the predictors will have a meaningful estimated coefficient.

→ note that in plot may use a penalty measure like λ in the horizontal axis instead of flexibility.

Since flexibility & penalty are inversely related, but this would lose descriptive details losing the same information.

→ On first look $\lambda \rightarrow$ gives a lower flexibility (i.e. higher penalty) w/ no loss in model accuracy or precision (separability).

→ In fact, when λ is too small, we'll end up with a lot of zero coefficients. This means ridge regression does not drop variables from a model.

→ just as with OLS when dealing w/ high dimensions, it's sum of the predictors will have a meaningful estimated coefficient.

→ note that in plot may use a penalty measure like λ in the horizontal axis instead of flexibility.

Since flexibility & penalty are inversely related, but this would lose descriptive details losing the same information.

→ On first look $\lambda \rightarrow$ gives a lower flexibility (i.e. higher penalty) w/ no loss in model accuracy or precision (separability).

→ In fact, when λ is too small, we'll end up with a lot of zero coefficients. This means ridge regression does not drop variables from a model.

→ just as with OLS when dealing w/ high dimensions, it's sum of the predictors will have a meaningful estimated coefficient.

→ note that in plot may use a penalty measure like λ in the horizontal axis instead of flexibility.

Since flexibility & penalty are inversely related, but this would lose descriptive details losing the same information.

→ On first look $\lambda \rightarrow$ gives a lower flexibility (i.e. higher penalty) w/ no loss in model accuracy or precision (separability).

→ In fact, when λ is too small, we'll end up with a lot of zero coefficients. This means ridge regression does not drop variables from a model.

→ just as with OLS when dealing w/ high dimensions, it's sum of the predictors will have a meaningful estimated coefficient.

→ note that in plot may use a penalty measure like λ in the horizontal axis instead of flexibility.

Since flexibility & penalty are inversely related, but this would lose descriptive details losing the same information.

→ On first look $\lambda \rightarrow$ gives a lower flexibility (i.e. higher penalty) w/ no loss in model accuracy or precision (separability).

→ In fact, when λ is too small, we'll end up with a lot of zero coefficients. This means ridge regression does not drop variables from a model.

→