

1 Background

1.1 Predictive Shrinkage

Suppose we are given an $n \times p$ centered data matrix X and an $n \times 1$ vector of responses Y —the pair (X, Y) constitutes the training data. Assume the linear model $Y_i = \alpha + (\beta^*)^\top X_i + \varepsilon_i$, where $\varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$. Define the ordinary least squares (OLS) estimate

$$(\hat{\alpha}, \hat{\beta}) := \arg \min_{\alpha, \beta} \left\{ \|Y - \mathbf{1}\alpha - X\beta\|_2^2 \right\}, \quad (1)$$

$$\text{given by } \hat{\alpha} = \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i \text{ and } \hat{\beta} = (X^\top X)^{-1} X^\top (Y - \bar{Y}\mathbf{1}). \quad (2)$$

If we have test data $x \sim (0, \Sigma)$ and $y = \alpha + x^\top \beta^* + \varepsilon$, the OLS prediction is $\hat{y} = \hat{\alpha} + \hat{\beta}^\top x$. The bias-variance decomposition of prediction mean square error (PMSE) is

$$\mathbb{E}[(y - \hat{y})^2] = \left(1 + \frac{1}{n}\right) \sigma^2 + \mathbb{E}[(\beta^* - \hat{\beta})^\top \Sigma (\beta^* - \hat{\beta})] \quad (3)$$

$$= \left(1 + \frac{p+1}{n} + \frac{1}{n} \text{tr}(\mathbb{E}[(\Sigma - S) S^{-1}])\right) \sigma^2, \quad (4)$$

where $S = n^{-1} X^\top X$. If the design of the training set is fixed (so S is constant) and we assume the test set follows the distribution of the training set in the sense that $\Sigma = S$, then the last term vanishes and so the PMSE is $(1 + \frac{1}{n} + \frac{p}{n}) \sigma^2$. Alternatively, if the rows X_i of X are all i.i.d. $\mathcal{N}(0, \Sigma)$, then nS is a Wishart matrix, and so $\mathbb{E}[(nS)^{-1}] = \frac{\Sigma^{-1}}{\nu}$ where $\nu = n - p - 1$, yielding a larger overall PMSE of $(1 + \frac{1}{n} + \frac{p}{\nu}) \sigma^2$.

Under the first set of assumptions, where $\Sigma = S$ exactly, we can write the last term in equation (3) as $\mathbb{E}[(\beta^* - \hat{\beta})^\top \Sigma (\beta^* - \hat{\beta})] = \mathbb{E}[\|\hat{\xi} - \xi\|_2^2]$, where $\hat{\xi} = \Sigma^{1/2} \hat{\beta} \sim \mathcal{N}(\xi^*, (\sigma^2/n) I_p)$ and $\xi^* = \Sigma^{1/2} \beta^*$. This is a normal-means estimation problem, so when $p > 2$ we can achieve lower MSE $\mathbb{E}[\|\hat{\xi} - \xi^*\|_2^2] < \mathbb{E}[\|\hat{\xi} - \xi^*\|_2^2]$ with the James-Stein estimate

$$\tilde{\xi} = \left(1 - \frac{(p-2)(\hat{\sigma}^2/n)\nu}{(\nu+2)\|\hat{\xi}\|_2^2}\right) \hat{\xi}, \quad (5)$$

yielding the shrunk regression coefficients $\tilde{\beta} = \hat{K} \hat{\beta}$, where $\hat{K} = \left(1 - \frac{(p-2)(\hat{\sigma}^2/n)\nu}{(\nu+2)\hat{\beta}^\top S \hat{\beta}}\right)$. It follows that $\tilde{y} = \hat{\alpha} + \hat{K} \hat{\beta}^\top x$ has strictly better PMSE than OLS \hat{y} . Since $\hat{K} < 1$ we are left to conclude that the OLS predictions on held-out data were too large in magnitude. Pre-shrunk predictors of this form were first studied by Copas (1983), who also provided the Stein-shrinkage interpretation.

Another form of shrinkage is *ridge regression*

$$(\hat{\alpha}, \hat{\beta}(\lambda)) := \arg \min_{\alpha, \beta} \left\{ \|Y - \mathbf{1}\alpha - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \right\}, \quad (6)$$

$$\text{given by } \hat{\alpha} = \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i \text{ and } \hat{\beta}(\lambda) = (X^\top X + \lambda I_p)^{-1} X^\top (Y - \bar{Y}\mathbf{1}). \quad (7)$$

In particular, $\hat{\beta}(0) = \hat{\beta}$ is OLS and $\lim_{\lambda \uparrow \infty} \hat{\beta}(\lambda) = 0$. Note $\hat{\beta}(\lambda) = (X^\top X + \lambda I_p)^{-1} (X^\top X) \hat{\beta}$ shrinks $\hat{\beta}$ towards zero in a manner that accounts for the variability in the covariates X . For $\lambda > 0$ it is not the case that $\hat{\beta}(\lambda) = K \hat{\beta}$ for some K , which is to say the family of estimators $\tilde{\beta}(K, \lambda) = K \hat{\beta}(\lambda)$ is not overparametrized. We allow any $K > 0$, since for λ large, $\hat{\beta}(\lambda)$ tends to overshrink. The estimator minimizes the following loss functions:

$$\tilde{\beta}(K, \lambda) := \arg \min_{\beta} \left\{ \|K(Y - \mathbf{1}\bar{Y}) - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \right\} \quad (8)$$

$$= \arg \min_{\beta} \left\{ \|Y - \mathbf{1}\bar{Y} - X\beta\|_2^2 + \frac{\lambda}{K} \|\beta\|_2^2 + (K^{-1} - 1) \|X\beta\|_2^2 \right\}. \quad (9)$$

The first minimization performs ridge regression on the centered design X and shrunk (or inflated) centered responses $K(Y - \mathbf{1}\bar{Y})$. The second is like ridge with an additional penalty. If $K < 1$ then $(K^{-1} - 1)\|X\beta\|_2^2$ penalizes large $X\beta$, and if $K > 1$ it penalizes small $X\beta$.

When $\lambda > 0$, the ridge regression estimate $\hat{\beta}(\lambda)$ is biased, but has the advantage of being well-defined even in the high dimensional case $n < p$. Whether we can benefit from the additional pre-factor K depends on how reliably we can detect whether we need to inflate or shrink. In the case where β^* is known to be sparse, we can ask the same set of questions about the magnitude of our predictions after Lasso, where the regularization term in (6) is replaced with $\lambda\|\beta\|_1$. We start with simulation studies to get a sense of in what cases we can significantly improve predictions by including K .