

STATS305A - Lecture 13

John Duchi

Scribed by Michael Howes

11/04/21

Contents

1	Announcements	1
2	Risk and model selection	1
2.1	Bias variance tradeoff	2
3	Ridge regression	2
3.1	Idea	2
3.2	Ridge estimator	3
3.3	Regularization path	6
4	Principal component analysis (PCA) and regression	6

1 Announcements

- Exercise 3 to be released today.
- An email announcement will be sent when exercise 3 is available.

2 Risk and model selection

Recall that we are now looking at model selection and prediction. We are starting with the more general model $y = f(x) + \varepsilon$ where $f(x) = \mathbb{E}[y|x]$ and we want to fit a predictor $h : \mathcal{X} \rightarrow \mathbb{R}$.

We can think of having a sequence of model spaces $\mathcal{H}_1 \subset \mathcal{H}_2 \subset \dots \subset \mathcal{H}_k \subset \dots$. We want to find h such that the quantity

$$\frac{1}{n} \sum_{i=1}^n (y_i - h(x_i))^2,$$

is small but h is not too “complex”. Typically we think of \mathcal{H}_k as growing in complexity as k increases. Thus we wish to choose $h \in \mathcal{H}_k$ minimizing

$$\sum_{i=1}^n (y_i - h(x_i))^2 + \text{penalty}(h),$$

or

$$\sum_{i=1}^n (y_i - h(x_i))^2 + \text{penalty}(\mathcal{H}_k).$$

Talking about this last time led to Mallows's C_p statistic which was

$$C_p = \|y - X\beta\|_2^2 + 2\sigma^2 p,$$

where $\beta \in \mathbb{R}^p$ and $2\sigma^2 p$ is a complexity penalty. Mallows's method is to choose the model that minimizes C_p .

2.1 Bias variance tradeoff

For any $x \in \mathcal{X}$ and any model procedure \hat{h} we have

$$\begin{aligned} \mathbb{E}[(\hat{h}(x) - f(x))^2] &= \mathbb{E}\left[(\hat{h}(x) - \mathbb{E}[\hat{h}(x)] + \mathbb{E}[\hat{h}(x)] - f(x))^2\right] \\ &= \left(\mathbb{E}[\hat{h}(x)] - f(x)\right)^2 + \mathbb{E}\left[(\hat{h}(x) - \mathbb{E}[\hat{h}(x)])^2\right] \\ &= \text{Bias}(\hat{h}(x))^2 + \text{Var}(\hat{h}(x)). \end{aligned}$$

Recall that we defined the in-sample risk to be

$$R_{in}(\hat{h}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(\hat{h}(x_i) - f(x_i))^2],$$

where the expectation is over $Y_i = f(x_i) + \varepsilon_i$.

Example 1. (actually everything we will look at) Let $\mu = [f(x_i)]_{i=1}^n = \mathbb{E}[Y]$. Suppose

$$\hat{y} = Hy = [\hat{h}(x_i)]_{i=1}^n,$$

for some matrix H that is not necessarily a projection matrix or symmetric. If $\varepsilon \sim (0, \sigma^2 I_n)$, then

$$\begin{aligned} nR_{in}(\hat{h}) &= \mathbb{E}[\|\hat{y} - \mu\|_2^2] \\ &= \mathbb{E}[\|Hy - \mu\|_2^2] \\ &= \mathbb{E}[\|H\mu - H\varepsilon - \mu\|_2^2] \\ &= \mathbb{E}[\|H\varepsilon - (I - H)\mu\|_2^2] \\ &= \|(I - H)\mu\|_2^2 + \mathbb{E}[\|H\varepsilon\|_2^2] \\ &= \text{Bias}(H)^2 + \sigma^2 \text{tr}(H^T H). \end{aligned}$$

Likewise for the residual sum of squares (RSS) we have

$$\begin{aligned} \mathbb{E}[RSS] &= \mathbb{E}[\|\hat{y} - y\|_2^2] \\ &= \mathbb{E}[\|y - Hy\|_2^2] \\ &= \mathbb{E}[\|(I - H)\mu + (I - H)\varepsilon\|_2^2] \\ &= \|(I - H)\mu\|_2^2 + n\sigma^2 - 2\sigma^2 \text{tr}(H) + \sigma^2 \text{tr}(H^T H) \\ &= nR_{in}(\hat{h}) + n\sigma^2 - 2\sigma^2 \text{tr}(H). \end{aligned}$$

3 Ridge regression

3.1 Idea

If X is “ill-conditioned”, then some errors in y can induce changes in $\hat{\beta}$ that are too large if we are using OLS.

We can formalize this with the SVD of X . Suppose $X = U\Gamma V^T$ is the SVD of X and $\Gamma = \text{diag}(\gamma_1, \dots, \gamma_d)$. Then

$$\begin{aligned}\hat{\beta} &= (X^T X)^{-1} X^T y \\ &= V\Gamma^{-1} U^T y.\end{aligned}$$

Suppose $U = [u_1, \dots, u_d] \in \mathbb{R}^{n \times d}$. Consider two alternative responses

$$\begin{aligned}y' &= y + tu_1, \quad t \in \mathbb{R}, \\ y'' &= y + tu_d, \quad t \in \mathbb{R}.\end{aligned}$$

We thus have

$$\begin{aligned}\hat{\beta}(y') &= \hat{\beta}(y) + \frac{tv_1}{\gamma_1} \\ \hat{\beta}(y'') &= \hat{\beta}(y) + \frac{tv_d}{\gamma_d}.\end{aligned}$$

If γ_d is very small (say $\gamma_d = \frac{1}{1000}$), then $\hat{\beta}(y'')$ will be very far from $\hat{\beta}(y)$ even if t is small. This issue arises when $\frac{\gamma_1}{\gamma_d}$ is very large. The quantity $\text{cond}(X) = \frac{\gamma_1}{\gamma_d}$ is called the *condition number* of X . In numerical analysis, $\text{cond}(X) \geq 10^6$ leads to problems. In statistics things are more sensitive and we run into trouble if $\text{cond}(X) \gtrsim 100$.

3.2 Ridge estimator

One solution to the sensitivity of $\hat{\beta}$ is to regularize $\hat{\beta}$ by forcing $\|\hat{\beta}\|_2^2$ to be small.

Definition 1. For $\lambda \geq 0$, define the *ridge estimator* to be

$$\hat{\beta}_\lambda := \underset{b}{\text{argmin}} \left\{ \|y - Xb\|_2^2 + \lambda \|b\|_2^2 \right\}.$$

Since $\partial_b \left(\|y - Xb\|_2^2 + \lambda \|b\|_2^2 \right) = X^T Xb - X^T y + \lambda b$, we have that

$$\hat{\beta}_\lambda = (X^T X + \lambda I)^{-1} X^T y.$$

This method of estimating β is called ridge regression or Tikhonov regularization.

Fact 1. If $X = U\Gamma V^T$, then

$$\hat{y}_\lambda = X\hat{\beta}_\lambda = \sum_{j=1}^d \frac{\gamma_j^2}{\gamma_j^2 + \lambda} u_j u_j^T y = H_\lambda y,$$

where

$$H_\lambda = U \text{diag} \left(\frac{\gamma_j^2}{\gamma_j^2 + \lambda} \right) U^T.$$

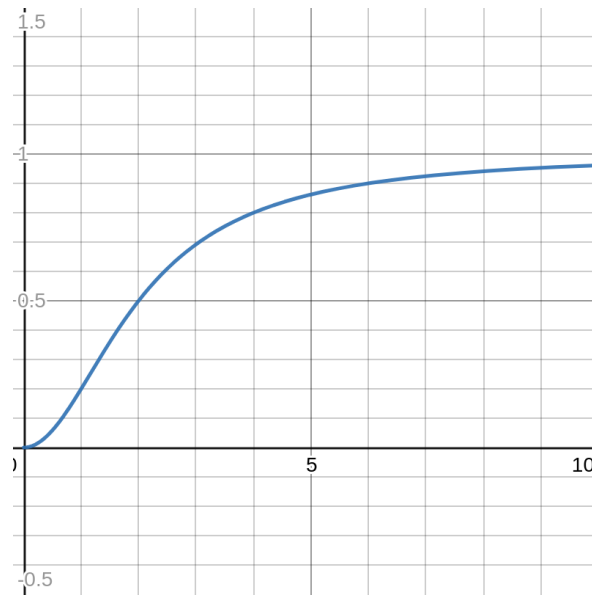
Also

$$\frac{\partial}{\partial \lambda} R_{in}(\hat{\beta}_\lambda)|_{\lambda=0} = -\frac{2\sigma^2 \text{tr}(\Gamma^{-2})}{n} < 0.$$

Some thoughts on the above facts before we prove them. The normal OLS estimator satisfies

$$\begin{aligned}\hat{y} &= Hy \\ &= UU^T y \\ &= \sum_{j=1}^d u_j u_j^T y \\ &= \Pi_{\text{range}(X)} y.\end{aligned}$$

On the other hand \hat{y}_λ is the projection of y onto $\text{range}(X) = \text{span}\{u_i\}$ but then shrunk by $\frac{\gamma_j^2}{\gamma_j^2 + \lambda}$. The function $\gamma_j \mapsto \frac{\gamma_j^2}{\gamma_j^2 + \lambda}$ looks like this ($\lambda = 4$ was used in this plot):



When $\gamma_j = 0$, $\frac{\gamma_j^2}{\gamma_j^2 + \lambda} = 0$ and as $\gamma_j \rightarrow \infty$, $\frac{\gamma_j^2}{\gamma_j^2 + \lambda} \rightarrow 1$. When $\gamma_j = \sqrt{\lambda}$, $\frac{\gamma_j^2}{\gamma_j^2 + \lambda} = \frac{1}{2}$. The idea is that for small values of γ_j , there is too much variance and we shrink away. When γ_j is large we have lot of “juice” in that direction and we do not shrink much.

Remark 1. In the above we are shrinking towards 0 but typically we either replace y with $y - \bar{y}\mathbf{1}$ or use

$$\hat{\beta}_\lambda = \underset{(b_0, b)}{\operatorname{argmin}} \left\{ \|Xb + b_0\mathbf{1} - y\|_2^2 + \lambda \|b\|_2^2 \right\}.$$

Either of these changes mean that we shrink toward \bar{y} instead of 0 and we do not regularize the intercept term. Note that if we replace y with $y - \bar{y}\mathbf{1}$, our estimator becomes

$$\hat{y}_\lambda = \bar{y}\mathbf{1} + U \operatorname{diag} \left(\frac{\gamma_j^2}{\gamma_j^2 + \lambda} \right) U^T (y - \bar{y}\mathbf{1}),$$

thus we can see that we shrink towards \bar{y} if there is little info in X .

The statement $\frac{\partial}{\partial \lambda} R_{in}(\hat{\beta}_\lambda)|_{\lambda=0}$ says that there is always a ridge estimator that beats OLS estimation. Thus introducing a bit of bias always improves estimation (even if $y = X\beta + \varepsilon$ is true).

Proof of Fact. We know $\hat{y}_\lambda = X\hat{\beta}_\lambda$ and that $\hat{\beta}_\lambda = (X^T X + \lambda I)^{-1} X^T y$, then $\hat{y}_\lambda = H_\lambda y$ where

$$\begin{aligned} H_\lambda &= X(X^T X + \lambda I)^{-1} X^T \\ &= U\Gamma V^T (V\Gamma^2 V^T + \lambda I)^{-1} V\Gamma U^T \\ &= U\Gamma V^T (V(\Gamma^2 + \lambda I)V^T)^{-1} V\Gamma U^T \\ &= U\Gamma(\Gamma^2 + \lambda I)^{-1} \Gamma U^T \\ &= U \operatorname{diag} \left(\frac{\gamma_j^2}{\gamma_j^2 + \lambda} \right) U^T. \end{aligned}$$

Thus we have proved the first part of the fact. For the second part define $\dot{H}_\lambda = \frac{\partial}{\partial \lambda} H_\lambda$, then

$$\dot{H}_\lambda = U \operatorname{diag} \left(\frac{d}{d\lambda} \left(\frac{\gamma_j^2}{\gamma_j^2 + \lambda} \right) \right) U^T = -U \operatorname{diag} \left(\frac{\gamma_j^2}{(\gamma_j^2 + \lambda)^2} \right) U^T.$$

Recall that

$$\begin{aligned} nR_{in}(\hat{\beta}_\lambda) &= \|(I - H_\lambda)\mu\|_2^2 + \sigma^2 \operatorname{tr}(H_\lambda^2) \\ &= \mu^T (I - 2H_\lambda + H_\lambda^2) \mu + \sigma^2 \operatorname{tr}(H_\lambda^2) \end{aligned}$$

Thus

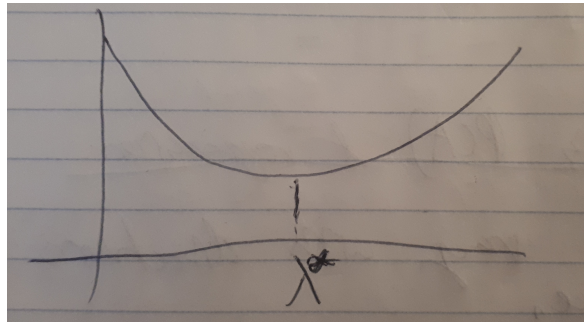
$$\frac{\partial}{\partial \lambda} nR_{in}(\hat{\beta}_\lambda) = -2\mu^T \dot{H}_\lambda \mu + 2\mu^T H_\lambda \dot{H}_\lambda \mu + 2\sigma^2 \operatorname{tr}(H_\lambda \dot{H}_\lambda).$$

When $\lambda = 0$, $H_\lambda = UU^T$ and $H_\lambda \dot{H}_\lambda = \dot{H}_\lambda = \dot{H}_0$. Thus

$$\begin{aligned} \frac{\partial}{\partial \lambda} \left(\frac{n}{2} R_{in}(\hat{\beta}_\lambda) \right) \Big|_{\lambda=0} &= -\mu^T \dot{H}_0 \mu + \mu^T \dot{H}_0 \mu + \sigma^2 \operatorname{tr}(\dot{H}_0) \\ &= \sigma^2 \operatorname{tr}(\dot{H}_0) \\ &= -\sigma^2 \operatorname{tr}(U\Gamma^{-2}U^T) \\ &= -\sigma^2 \operatorname{tr}(\Gamma^{-2}). \end{aligned}$$

□

John said that the matrix derivatives are not important for this course but it's good for us to see the calculation and understand the implication that the OLS estimator is worse than a ridge estimator. Too much regression can be a bad thing. If $\lambda \rightarrow \infty$, then $\hat{y}_\lambda = \bar{y}\mathbf{1}$ which is not a very good predictor. The risk as a function of λ tends to look something like this:



A natural question is can we find λ^* ? The number λ^* is the value of λ with the smallest risk. We know that

$$\mathbb{E}[RSS_\lambda] = nR_{in}(\hat{\beta}_\lambda) + n\sigma^2 - 2\sigma^2 \operatorname{tr}(H_\lambda).$$

Thus if we can estimate σ^2 with $\hat{\sigma}^2$, then we can choose λ to minimize

$$\mathbb{E}[RSS_\lambda] + 2\hat{\sigma}^2 \operatorname{tr}(H_\lambda).$$

On the homework we will look at different estimators of $\hat{\sigma}^2$.

3.3 Regularization path

We often want to compute $\hat{\beta}_\lambda$ for all $\lambda \geq 0$. This gives us a path of solutions parametrized by λ which is called a *regularization path*. This can be done very efficiently if we have the SVD of X . Suppose $X = U\Gamma V^T$, then

$$\begin{aligned}\hat{\beta}_\lambda &= (X^T X + \lambda I)^{-1} X^T y \\ &= V (\Gamma^2 + \lambda I)^{-1} \Gamma U^T y \\ &= \sum_{j=1}^d v_j \frac{\gamma_j}{\gamma_j^2 + \lambda} u_j^T y.\end{aligned}$$

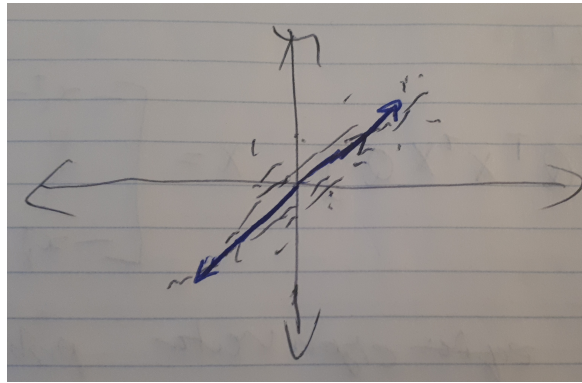
If we store $\tilde{y} = U^T y$, then we can compute

$$\hat{\beta}_\lambda = V \operatorname{diag} \left(\frac{\gamma_j}{\gamma_j^2 + \lambda} \right) \tilde{y},$$

which can be done in $O(d^2)$ time.

4 Principal component analysis (PCA) and regression

Idea: Often y should vary with directions in x that have high variance. This is the idea of *principal component analysis (PCA)*: find directions in X -space that vary the most (equivalently find direction in X -space that most accurately reconstruct our data $\{x_i\}$). For example if our data looks like the below picture, we wish to find the blue line.



Suppose we have a subspace $\mathcal{S} = \operatorname{span}\{q_1, \dots, q_k\}$ where $q_i \in \mathbb{R}^d$ are orthonormal, that is

$$q_i^T q_j = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{if } i \neq j. \end{cases}$$

Then the projection onto \mathcal{S} is $\Pi_{\mathcal{S}} = QQ^T$ where $Q = [q_1, \dots, q_k] \in \mathbb{R}^{d \times k}$. Then for all $x \in \mathbb{R}^d$,

$$\begin{aligned}\operatorname{dist}^2(x, \mathcal{S}) &= \|x - \Pi_{\mathcal{S}} x\|_2^2 \\ &= \|(I - QQ^T)x\|_2^2 \\ &= \|x\|_2^2 - \|Q^T x\|_2^2.\end{aligned}$$

This gives us a procedure for PCA. We wish to find a k -dimensional subspace which is closest to all our data x_1, \dots, x_n . Thus we wish to solve

$$\begin{aligned} \min_{Q^T Q = I_k} \sum_{i=1}^n \text{dist}^2(x_i, \text{span}(Q)) &= \min_{Q^T Q = I_k} \sum_{i=1}^n \|(I - QQ^T)x_i\|_2^2 \\ &= \min_{Q^T Q = I_k} \sum_{i=1}^n \|x_i\|_2^2 - \|Q^T x_i\|_2^2 \\ &= \max_{Q^T Q = I_k} \sum_{i=1}^n \|Q^T x_i\|_2^2 \\ &= \max_{Q^T Q = I_k} \text{tr}(Q^T X^T X Q), \end{aligned}$$

where

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \in \mathbb{R}^{n \times d}.$$

Remark 2. This is an eigenvalue/eigenvector problem. If $X = U\Gamma V^T$, then our problem becomes

$$\max_{Q^T Q = I_k} \text{tr}(Q^T V \Gamma^2 V^T Q).$$

Our idea: put all “juice” for the first k directions in V into Q since

$$\Gamma^2 = \text{diag}(\gamma_1^2, \dots, \gamma_d^2),$$

where $\gamma_1 \geq \dots \geq \gamma_d$. So $Q = [v_1, \dots, v_k]$. For reference this concept is called:

- Variational characterization of eigenvectors.
- Rayleigh-Ritz eigenvector characterization.