

①

Statistical Learning III

- Ridge regression computation
- Comparison Ridge, Lasso, Best-subset
- Probability Theory remarks
- Bayesian Statistics
 - a) Intro
 - b) Laplace's rule of succession
 - c) MAP interpretation of Ridge, Lasso

I. Ridge regression computation

$$\hat{\beta}_{\text{ridge}} = \underset{\beta}{\operatorname{arg\min}} \text{RSS}(\beta) + \lambda \sum_{j=1}^p \beta_j^2 \quad T = (y_1, \dots, y_N)^\top$$

Can interpret as least squares with extra fake training data

$$(y_{N+1}), \dots, (y_{N+p-1})$$

$$y_{N+j} = 0$$

$$x_{N+j} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ \sqrt{\lambda} \\ 0 \\ \vdots \\ 0 \end{pmatrix} \leftarrow (j+1)\text{-th position for } j=1, \dots, p-1$$

$$\text{Then } (y_{N+j} - x_{N+j}^\top \beta)^2 = \lambda \beta_{j+1}^2 \text{ so}$$

$$\hat{\beta}_{\text{ridge}} = \underset{\beta}{\operatorname{arg\min}} \sum_{i=1}^{N+p-1} (y_i - x_i^\top \beta)^2 = \text{least squares with extra fake data.}$$

(8)

5. Comparison of Ridge, Lasso, Best-subset Selection

as Ridge, Lasso vs Best-subset

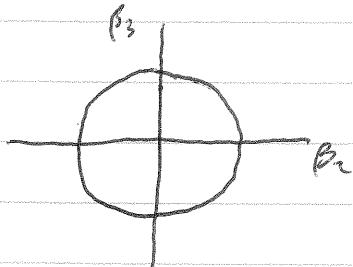
↑
small changes in data cause

small changes in estimate, so less variance than best-subset.

b) Ridge vs Lasso

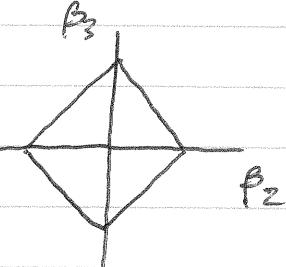
Ridge

$$\sum_{j=2}^p \beta_j^2 \leq t$$



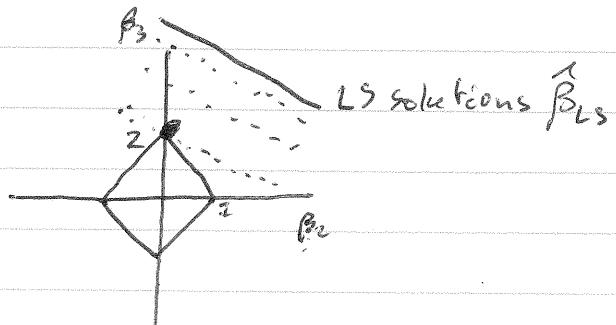
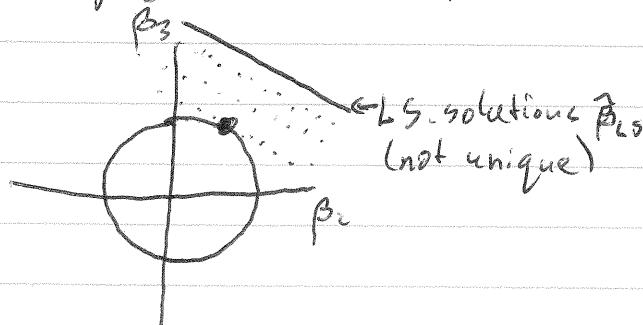
Lasso

$$\sum_{j=1}^p |\beta_j| \leq t$$



* Figure 3.11

* For highly correlated features:



β_2 and β_3 get about equal weight

minor noise determines which of corners 1 and 2 is chosen:
one variable gets all weight,
the other 0.

3. Probability Theory Remarks

Bayes' rule:

$$\Pr(A|B) = \frac{\Pr(B|A) \cdot \Pr(A)}{\Pr(B)}$$

Often used in Bayesian ~~and~~ and standard frequentist statistics

Densities are slippery:

- * Location of maximum depends on choice of parametrization
- * Uniform density in one parametrization is not uniform in another parametrization

(see slides)

(2)

2. Bayesian Statistics

a) Intro

Probability model: $\mathcal{P} = \{P_\theta(x, y) | \theta \in \Theta\}$ or $\mathcal{P} = \{P_\theta(y) | \theta \in \Theta\}$

Frequentist statistics (standard):

- optimal parameter θ^* is fixed, but unknown
- diversity of methods
- need proofs/experiments to justify methods

Bayesian statistics:

- pretend that true parameter θ^* is a random variable distributed according to prior distribution $\pi(\theta)$ that we know.

$T = Y_1, \dots, Y_N$ (no features for simplicity)

$\Pr(T, \theta) = P_\theta(T) \cdot \pi(\theta)$ is joint distribution of data and parameters

Since we know the full joint distribution, can simply use probability theory to compute any probability we are interested in. ← single method

If we estimate different probabilities this way, they are beautifully consistent.

~~But~~ Bayesian premise is too strong:

- prior π is chosen for computational or information theoretic properties in practice, so cannot ~~blindly~~ assume θ^* is random sample from π .
- need proofs/experiments to justify Bayesian methods

(3)

frequentist Modern motivation:

- If we choose π right, then often works really well (both in theory and in practice).
- Learns faster if true θ^* has high prior probability, slower if true θ^* has small prior probability,
so can use π to express prior knowledge about our data.

Posterior Distribution:

$$\pi(\theta|T) := \Pr(\theta|T) = \frac{\Pr(\theta, T)}{\Pr(T)} = \frac{P_\theta(T) \cdot \pi(\theta)}{\Pr(T)}$$

- Often puts its probability mass closer and closer to θ^* as $N \rightarrow \infty$. (vd Vaart et al.)
- Expresses uncertainty about θ

Predictive Distribution:

$$\Pr(Y|T) = \int_0^\infty P_\theta(Y) \cdot \pi(\theta|T) d\theta = \frac{\Pr(Y, T)}{\Pr(T)}$$

↑ ↑
 new sample often better predictions
 outside of than frequentist $P_{\hat{\theta}}(Y)$
 training set because $\pi(\theta|T)$ keeps track of uncertainty
better than fixed single choice $\hat{\theta}$.

b] Example: Laplace Rule of Succession

Bernoulli model: $P_\theta(Y) = \begin{cases} \theta & \text{for } Y=1 \\ 1-\theta & \text{for } Y=0 \end{cases} \quad \theta \in [0, 1]$

Maximum likelihood: $\hat{\theta} = \frac{n_1}{N} \rightarrow P_{\hat{\theta}}(Y=1) = \hat{\theta} = \frac{n_1}{N}$ ← dangerous for prediction if $n_1=0$

Suppose $\pi(\theta)=1$ is uniform prior density, ← not the same as "no prior knowledge" because depends on parametrisation

Then $\Pr(Y=1|T) = \frac{n_1 + 1}{N + 2}$

(Laplace, 1814)

$$\Pr(Y=0|T) = \frac{n_0 + 1}{N + 2}$$

↙ beta function

Proof: $\frac{\Pr(Y, T)}{\Pr(T)} = \frac{\int P_\theta(Y, T) \cdot \pi(\theta) d\theta}{\int P_\theta(T) \cdot \pi(\theta) d\theta} = \frac{\int \theta^{n_1+Y} (1-\theta)^{n_0+2-Y} d\theta}{\int \theta^{n_1} (1-\theta)^{n_0} d\theta} = \frac{n_1 + 1}{N + 2}$ □

(4)

c) MAP interpretation of Ridge and Lasso

Maximum a Posteriori (MAP) parameters:

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} \pi(\theta | T) = \underset{\theta}{\operatorname{argmax}} \frac{p_{\theta}(T) \cdot \pi(\theta)}{p_T(T)}$$

$$= \underset{\theta}{\operatorname{argmax}} p_{\theta}(T) \cdot \pi(\theta)$$

- maximizes posterior density, so for continuous parameters depends on parametrisation
- "real" Bayesians prefer prediction with predictive distribution

Ridge / Lasso:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \text{RSS}(\beta) + \lambda \text{pen}(\beta)$$

$$\text{ridge: pen}(\beta) = \sum_{j=2}^p \beta_j^2 \quad \text{lasso: pen}(\beta) = \sum_{j=2}^p |\beta_j|$$

Suppose Gaussian noise:

$$y = x^T \beta + \varepsilon \quad \varepsilon \sim N(0, \sigma^2)$$

$$p_{\beta}(y_1, \dots, y_N | x_1, \dots, x_N) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - x_i^T \beta)^2}{2\sigma^2}}$$

$$\hat{\beta}_{\text{MAP}} = \underset{\beta}{\operatorname{argmax}} p_{\beta}(y_1, \dots, y_N | x_1, \dots, x_N) \cdot \pi(\beta)$$

$$= \underset{\beta}{\operatorname{argmin}} -\log p_{\beta}(y_1, \dots, y_N | x_1, \dots, x_N) - \log \pi(\beta)$$

$$= \underset{\beta}{\operatorname{argmin}} N \cdot \left(-\log \frac{1}{\sqrt{2\pi\sigma^2}} \right) + \frac{1}{2\sigma^2} \text{RSS}(\beta) - \log \pi(\beta)$$

$$= \underset{\beta}{\operatorname{argmin}} \text{RSS}(\beta) - 2\sigma^2 \log \pi(\beta)$$

Suppose data have been pre-processed such that we can assume that the intercept $\hat{\beta}_0 = 0$.

(5)

Ridge: Choose π s.t.

$$\beta_1 = 0 \text{ with prob. 1}$$

$$(\beta_2, \dots, \beta_p) \sim N(0, \sigma^2 I)$$

$$-\log \pi(\beta) = \sum_{j=2}^p -\log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(\beta_j - 0)^2}{2\sigma^2}} \right)$$

$$= (p-1)(-\log \frac{1}{\sqrt{2\pi\sigma^2}}) + \sum_{j=2}^p \frac{\beta_j^2}{2\sigma^2}$$

$$\hat{\beta}_{MAP} = \underset{\beta}{\operatorname{argmin}} \text{RSS}(\beta) + \frac{\sigma^2}{2\sigma^2} \sum_{j=2}^p \beta_j^2$$

is ridge with $\lambda = \frac{\sigma^2}{2\sigma^2}$

Is also posterior mean $E[\beta]_{\pi(\beta|T)}$.

Lasso: Choose π s.t.

$$\beta_1 = 0 \text{ with prob. 1.}$$

$$(\beta_2, \dots, \beta_p) \sim \pi \underset{j=2}{\overset{p}{\prod}} \frac{1}{2\sigma} e^{-\frac{|\beta_j|}{\sigma}}$$

$$-\log \pi(\beta) = \sum_{j=2}^p -\log \left(\frac{1}{2\sigma} \cdot e^{-\frac{|\beta_j|}{\sigma}} \right)$$

$$= (p-1)(-\log \frac{1}{2\sigma}) + \sum_{j=2}^p \frac{|\beta_j|}{\sigma}$$

$$\hat{\beta}_{MAP} = \underset{\beta}{\operatorname{argmin}} \text{RSS}(\beta) + \frac{2\sigma^2}{\sigma} \sum_{j=2}^p |\beta_j|$$

is Lasso with $\lambda = \frac{2\sigma^2}{\sigma}$.

Remarks:

- "real" Bayesian prefers predicting with predictive distribution
- MAP + CV to determine λ is very unBayesian.