

Statistical Learning II

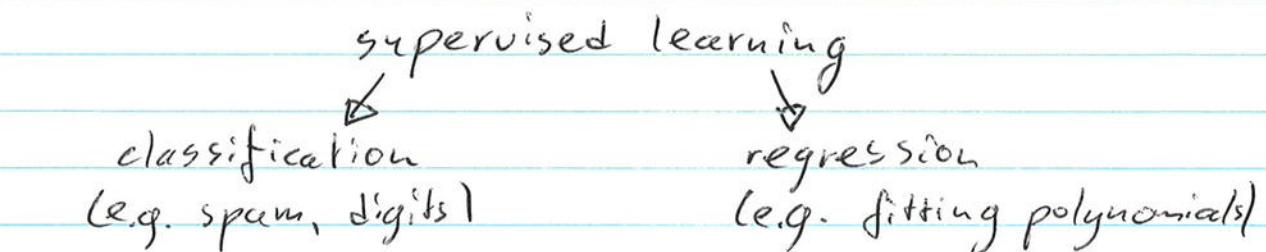
8-11-2019

1. Organization
2. Warm up a) Statistical Decision Theory
b) ERM \rightarrow least squares
3. K-fold Cross-validation
4. Model Selection
 - a) Intro
 - b) Best-subset
 - c) Ridge Regression
 - d) Lasso
5. Homework 1

1. Organization

- Anyone more time on exam?
- Usis/Blackboard/website

2. Warm up



$$T = \begin{pmatrix} y_1 \\ x_1 \end{pmatrix}, \dots, \begin{pmatrix} y_N \\ x_N \end{pmatrix}$$

$\hookrightarrow f \in \mathcal{F}$ and predict $\hat{y} = \hat{f}(x)$ for new x

(2)

2a Statistical Decision Theory

$EPE(f) = \mathbb{E}_{x,y} [L(Y, f(x))]$ measures quality of f

regression: $L(Y, f(x)) = (Y - f(x))^2$ "squared error"

classification:

$$L(Y, f(x)) = \begin{cases} 0 & \text{if } f(x) = Y \\ 1 & \text{if } f(x) \neq Y \end{cases}$$

"0/1-loss"

Bayes optimal predictor:

$$f_B = \underset{f}{\operatorname{argmin}} EPE(f)$$

regression: $f_B(x) = \mathbb{E}[Y|x]$

classification: $f_B(x) = \underset{g}{\operatorname{argmax}} \Pr(g=Y|x)$

Estimators:

\hat{f} depends on $T \Rightarrow EPE(\hat{f})$ depends on T
 \Rightarrow evaluate $\mathbb{E}_T [EPE(\hat{f})]$

2b) Empirical Risk Minimization

$$\hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N L(Y_i, f(x_i))$$

classification: - minimize #mistakes

- usually cannot compute efficiently
 for 0/1-loss.

(unless #mistakes = 0 is achievable on data T)

regression: - minimize sum of squared errors
 - can compute if \mathcal{F} is linear model

8

4. k-Fold Cross-Validation

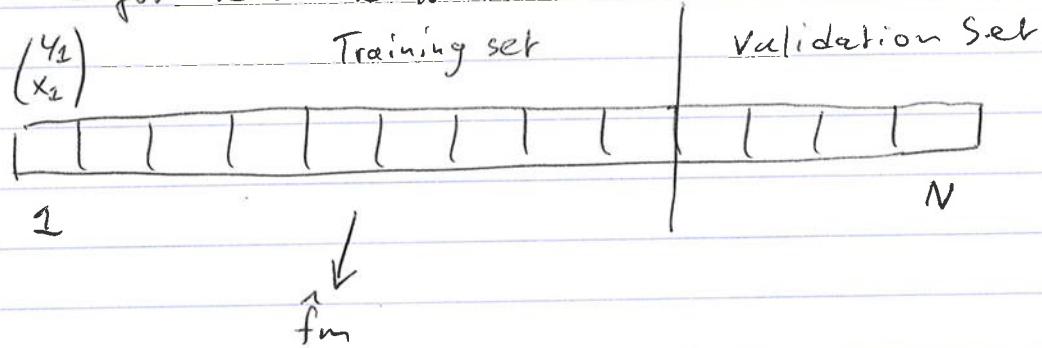
Hyperparameter $m \in \{1, \dots, M\}$

E.g. * m is degree of polynomial in linear regression.
* m is k in k -nearest neighbour

Training set $\rightarrow \hat{f}_m$ for each m .

What happens if I look at the fit on the training data to choose m ? (i.e., I use ERM)?

- for degree of polynomial.
- for k in k -NN.



Hold-out estimate:

- Randomly set aside validation set V
- Construct \hat{f}_m on rest of training data for all m
- Use ERM on validation set:

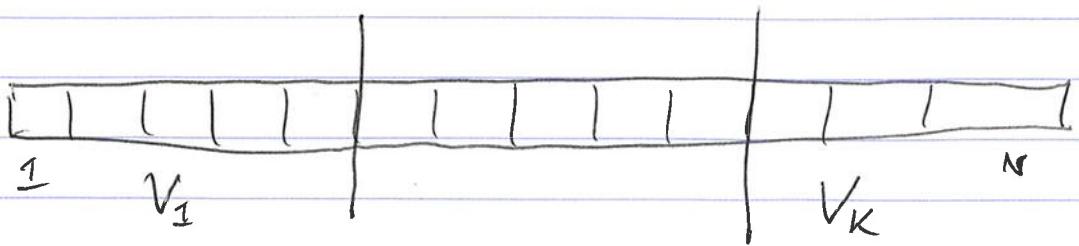
$$\hat{m} = \underset{m}{\operatorname{argmin}} \frac{1}{|V|} \sum_{i \in V} L(y_i, \hat{f}_m(x_i))$$

- Works well if ~~validation~~ V is large enough,
but that reduces training set a lot

(9)

K-Fold Cross validation:

Try to get away with smaller validation set by averaging multiple choices.



For $k = 1, \dots, K$:

- train \hat{f}_m on all data except V_k (for each m)

- evaluate on V_k :

$$\hat{L}_m^k = \frac{1}{|V_k|} \sum_{i \in V_k} L(Y_i, \hat{f}_m(X_i))$$

$$\hat{m} = \arg \min_m \frac{1}{K} \sum_{k=1}^K \hat{L}_m^k \quad (\text{average over all choices of validation set})$$

Often train final \hat{f}_m on all data, or take the ~~average~~ average of the ~~estimates~~ K estimators found for \hat{m} during cross-validation.

$K=N$: "leave-one-out cross-validation"

(story about Merlin Stone and LOO.)

Best K depends on:

Complicated,
so optimal
choice
only known in
special cases

- * Bias: estimates \hat{f}_m trained on $N - N/K$ examples might perform better if trained on all data, so small K might be pessimistic: overestimate EPE
- * Variance: - larger validation set size N/K reduces variance
 - but fewer folds K increases variance
 - stability of estimators \hat{f}_m also has effect

In practice: K=5 or K=10 is good in many cases (see book for references)

5. Model Selection for Regression Intro

Suppose many features p . Maybe even $p \gg N$?

often not satisfied with least squares, because we want:

1. better prediction accuracy:

least squares often has low bias, but high variance

2. better interpretability:

of many features, which ones are most important?

E.g. If we have 20 000 genes, want to know
which ones are (most) relevant for prediction.

(2)

2. Best-subset Selection

Want to use only m most important features.

Select $m \leq p$ features as follows:

1. Run least squares for each of the $\binom{p}{m} \approx \binom{p}{m}^m$ subsets of m features.
2. Choose the subset such that $RSS(\hat{\beta})$ is smallest.

Different ways to choose m . E.g. use cross-validation.

Problems:

- * Computation: exponentially many subsets in m
- * Variance: discrete choice of subset; if multiple subsets approximately equally good, then which one of those is chosen varies a lot if we would sample a training set multiple times.

Computable approximations:

- * Forward Stepwise: start with no variables,
greedily choose one variable to add
until we have m variables
- * Backward Stepwise: start with all variables,
greedily choose one variable to remove
until we have m left.

(3)

3. Ridge regression

a) Definition

We set the first feature to 1: $x = \begin{pmatrix} 1 \\ x_2 \\ x_3 \\ \vdots \\ x_p \end{pmatrix}$

$$\text{ridge regression: } \hat{\beta}_{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \text{ RSS}(\beta) + \lambda \sum_{j=2}^p \beta_j^2$$

$$= \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N \left(Y_i - \beta_1 - \sum_{j=2}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=2}^p \beta_j^2$$

alternative formulation: ~~if β is unique~~ If β is unique:

$$\hat{\beta}_{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \text{ RSS}(\beta)$$

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

~~OR β is unique then~~

One-to-one relation between λ and t ~~effortless~~

b) Motivation

* More stable/less variance than least squares
for highly correlated features

Example: two features the same: $x_2 = x_3$

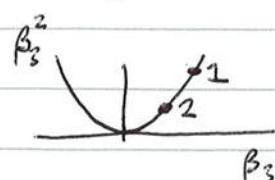
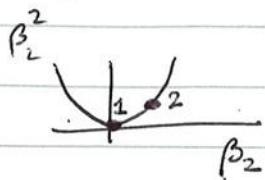
Then $\text{RSS}(\beta) = \sum_{i=1}^N \left(Y_i - \beta_1 - \sum_{j=2}^p x_{ij} \beta_j \right)^2$ the same for

1.] $\beta = (\beta_1, \beta_2, \beta_3, \dots, \beta_p)$ and

2.] $\beta = (\beta_1, \beta_2 + \alpha, \beta_3 - \alpha, \dots, \beta_p)$ for any α .

(4)

so least squares does not distinguish



ridge regression prefers equal weights (situation 2)

For highly correlated features it makes the weights approximately equal

* $\hat{\beta}_{\text{ridge}}$ is always uniquely defined for $d > 0$

Because $\text{RSS}(\beta) + \lambda \sum_{j=2}^p \beta_j^2$ is strictly convex.

Warning: need to normalize features

- If we make feature x_j 10 times as small, then least squares makes $\hat{\beta}_j$ 10 times as big and its predictions don't change.
- This does not hold for ridge regression, so need to center and rescale features to standard range

c) Computation

- Can interpret as least squares with fake training data.

$$\begin{pmatrix} Y_{N+1} \\ X_{N+1} \end{pmatrix}, \dots, \begin{pmatrix} Y_{N+p-1} \\ X_{N+p-1} \end{pmatrix}$$

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

$$(Y_{N+j} - X_{N+j}^\top \beta)^2 = \lambda \beta_{j+1}^2$$

(5)

$$U = \begin{pmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{N1} & \dots & x_{Np} \end{pmatrix} \quad V = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}$$

with fake training data:

$$\bar{U} = \begin{pmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_N & x_{Np} \\ 0 & \sqrt{\lambda} & 0 & 0 & \dots \\ 0 & 0 & \sqrt{\lambda} & 0 & \dots \\ 0 & 0 & 0 & \sqrt{\lambda} & \dots \\ 0 & 0 & 0 & 0 & \sqrt{\lambda} \end{pmatrix} \quad \bar{V} = \begin{pmatrix} y_1 \\ \vdots \\ y_N \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\hat{\beta}_{LS} = (U^T U)^{-1} U^T V$$

$$\hat{\beta}_{\text{ridge}} = (\bar{U}^T \bar{U})^{-1} \bar{U}^T \bar{V}$$

$$= (U^T U + \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \lambda & 0 & 0 & 0 & 0 \\ 0 & 0 & \lambda & 0 & 0 & 0 \\ 0 & 0 & 0 & \lambda & 0 & 0 \\ 0 & 0 & 0 & 0 & \lambda & 0 \\ 0 & 0 & 0 & 0 & 0 & \lambda \end{pmatrix})^{-1} U V$$

$$= (U^T U + \lambda I)^{-1} U V$$

~~$$A = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$~~

$$A = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

Different from book, because deals with intercept by pre-processing the data, but is equivalent.

(7)

4. Lasso

a) Definition

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

If $\hat{\beta}$ is unique, then

$$\hat{\beta}_{\text{Lasso}} = \underset{\beta}{\operatorname{argmin}} \text{RSS}(\beta) \\ \text{subject to } \sum_{j=2}^p |\beta_j| \leq t$$

One-to-one relation: $t \leftrightarrow \lambda$

b) Motivation

- * Sets many weights β_j to 0 to find only most important features: good for prediction and interpretation

- * $p \gg N$ allowed

NB Need to normalize features

c) Computation

- Exist multiple efficient algorithms
- LARS computes solutions for all values of t in one pass

(8)

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

as Ridge, Lasso vs Best-subset

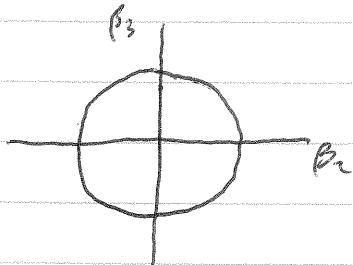
↑
small changes in data cause

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

b) Ridge vs Lasso

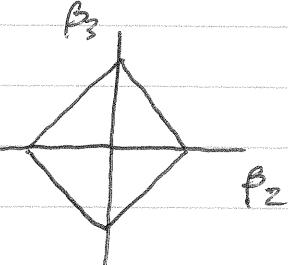
Ridge

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



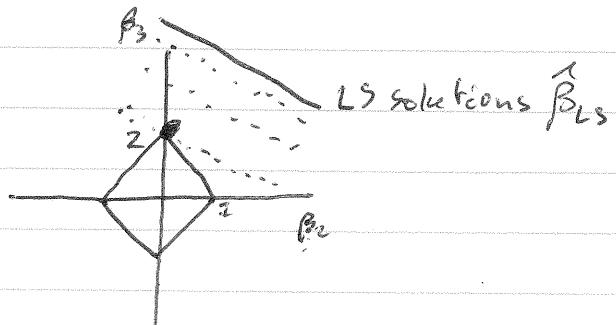
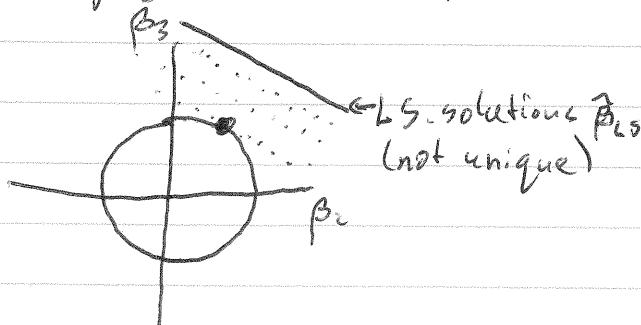
Lasso

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



* 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.

6. Homework 1

- * Goal: try out linear regression methods on real data: predicting housing prices
- * One of the features is racist
Bad? Depends on what you do with predictions.
- * Statistician's responsibility goes beyond technical aspects. Management does not know details of the methods used!