

STOR 767 Spring 2019 Hw5

Due on 03/18/2019 in Class

Problem Set

Remark. This homework focuses on additive models, model assessment and selection, and support vector machines.

Instruction.

- **Theoretical Part and Computational Part** are respectively credited 60 points. At most 100 points in total will be accounted for this homework.
- Submission of handwritten solution for the **Theoretical Part** of this homework is allowed.
- Please use **RMarkdown** to create a formatted report for the **Computational Part** of this homework.
- Some of the problems are selected or modified from the textbook [Friedman et al., 2009].

Theoretical Part

1. (10 pt) (Naive Bayes and Logistic GAM, Textbook Ex. 6.9) What's the differences between the naive Bayes model and a generalized additive Logistic regression model in terms of (a) model assumptions, and (b) estimation? If all the variables are discrete, what can you say about the Logistic GAM?
2. (15 pt) (Optimism, Textbook Ex. 7.4, 7.5) Let $\mathcal{Y} = \{Y_i\}_{i=1}^n$ be a training sample, $\mathcal{Y}^{\text{new}} = \{Y_i^{\text{new}}\}_{i=1}^n \stackrel{iid}{=} \mathcal{Y}$ be an independent copy of \mathcal{Y} , $\{\hat{Y}_i\}_{i=1}^n$ be the in-sample prediction based on \mathcal{Y} . Recall the in-sample prediction error and its training estimate

$$\text{Err}_{\text{in}} := \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathcal{Y}, \mathcal{Y}^{\text{new}}} \ell(Y_i^{\text{new}}, \hat{Y}_i), \quad \overline{\text{err}} := \frac{1}{n} \sum_{i=1}^n \ell(Y_i, \hat{Y}_i),$$

and the optimism

$$\text{op} := \text{Err}_{\text{in}} - \mathbb{E}_{\mathcal{Y}} \overline{\text{err}}.$$

Consider the squared-error loss $\ell(y, \hat{y}) := (y - \hat{y})^2$.

(I) Show that

$$\text{op} = \frac{2}{n} \sum_{i=1}^n \text{Cov}_{\mathcal{Y}}(\hat{Y}_i, Y_i).$$

(II) Assume $\text{Var}(Y_i) = \sigma^2$ ($1 \leq i \leq n$). Write \mathcal{Y} in vector form $\mathbf{Y} \in \mathbb{R}^n$. Let $\mathbf{S} \in \mathbb{R}^{n \times n}$ be a (fixed) smoother matrix, $\hat{\mathbf{Y}} := \mathbf{S}\mathbf{Y}$ be the linear-smoother in-sample prediction vector. Show that

$$\text{op} = \frac{2}{n} \text{Tr}(\mathbf{S}) \sigma^2.$$

3. (15 pt) (Bootstrap Prediction Error) Suppose $\mathcal{Y} := \{Y_1 = 1, Y_2 = 2, Y_3 = 6\}$ where $n = 3$. Consider a linear model

$$Y_i = \theta + \epsilon_i \quad (i = 1, 2, 3)$$

with $\epsilon_1, \epsilon_2, \epsilon_3 \stackrel{iid}{\sim} (0, \sigma^2)$ and squared-error loss $\ell(y, \hat{y}) := (y - \hat{y})^2$.

- (a) Consider Bootstrap on \mathcal{Y} . Enumerate all possible unordered **Bootstrap bags**¹ and their Bootstrap probabilities. For example, $\{1, 1, 2\}$ is a possible Bootstrap bag with probability $3/27$. Indicate the **out-of-bag (OOB)** sample points² for each unordered Bootstrap bag.
- (b) For each Bootstrap sample \mathcal{Y}_b^* , derive the least-square prediction rule and its prediction on \mathcal{Y} as $\{\hat{Y}_{bi}^*\}_{i=1}^n$. Compare the training error $\overline{\text{err}}$, Bootstrap prediction error estimate

$$\text{err}_b^* := \sum_{i=1}^n \ell(Y_i, \hat{Y}_{bi}^*), \quad \widehat{\text{Err}}_{\text{boot}} := \lim_{B \rightarrow +\infty} \frac{1}{nB} \sum_{b=1}^B \text{err}_b^*,$$

and the OOB prediction error estimate³

$$\text{err}_{\text{oob},b}^* := \sum_{Y_i \in \mathcal{Y} \setminus \mathcal{Y}_b^*} \ell(Y_i, \hat{Y}_{bi}^*), \quad p_{\text{oob},n} := \left(1 - \frac{1}{n}\right)^n, \quad \widehat{\text{Err}}_{\text{oob}} := \lim_{B \rightarrow +\infty} \frac{1}{np_{\text{oob},n}B} \sum_{b=1}^B \text{err}_{\text{oob},b}^*.$$

4. (20 pt) (SVM) Let $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \subseteq \mathbb{R}^p \times \{\pm 1\}$ be a training sample. Consider the large-margin linear classification problem

$$\begin{aligned} \max_{(\mathbf{w}, b) \in \mathbb{R}^{p+1}} \quad & \gamma \\ \text{s.t.} \quad & y_i(b + \mathbf{w}^T \mathbf{x}_i) \geq \gamma \quad (1 \leq i \leq n) \\ & \|\mathbf{w}\|_2 = 1 \end{aligned} \tag{1}$$

- (a) Show that (1) is equivalent to

$$\begin{aligned} \min_{(\mathbf{w}, b) \in \mathbb{R}^{p+1}} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_i(b + \mathbf{w}^T \mathbf{x}_i) \geq 1 \quad (1 \leq i \leq n) \end{aligned} \tag{2}$$

- (b) Introduce Lagrangian variables $\boldsymbol{\alpha} \in \mathbb{R}_+^n$ to inequality constraints in (2) and write down the Lagrangian function $L(\mathbf{w}, b; \boldsymbol{\alpha})$ [see [Boyd and Vandenberghe, 2004](#), Chapter 5]. Use strong duality

$$\begin{aligned} \min_{(\mathbf{w}, b) \in \mathbb{R}^{p+1}} \max_{\boldsymbol{\alpha} \in \mathbb{R}_+^n} L(\mathbf{w}, b; \boldsymbol{\alpha}) &= \min_{(\mathbf{w}, b) \in \mathbb{R}^{p+1}} L_{\mathcal{P}}(\mathbf{w}, b) \quad (\text{primal problem (2)}) \\ &= \max_{\boldsymbol{\alpha} \in \mathbb{R}_+^n} \min_{(\mathbf{w}, b) \in \mathbb{R}^{p+1}} L(\mathbf{w}, b; \boldsymbol{\alpha}) = \max_{\boldsymbol{\alpha} \in \mathbb{R}_+^n} L_{\mathcal{D}}(\boldsymbol{\alpha}) \quad (\text{dual problem (3)}) \end{aligned}$$

¹We call a size- n sample with replacement from \mathcal{Y} as a Bootstrap bag.

²Let \mathcal{Y}^* be a Bootstrap bag from \mathcal{Y} , then $\mathcal{Y} \setminus \mathcal{Y}^*$ is the OOB sample.

³ $\widehat{\text{Err}}_{\text{oob}}$ is the same as the leave-one-out Bootstrap estimate $\widehat{\text{Err}}^{(1)}$ introduced in [Friedman et al. \[2009, Equation \(7.56\)\]](#), where they only differ in the order of summation (summing over Bootstrap bags and sample points).

to derive the Lagrangian dual problem

$$\begin{aligned}
\max_{\boldsymbol{\alpha} \in \mathbb{R}^n} \quad & L_{\mathcal{D}}(\boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\
\text{s.t.} \quad & \alpha_i \geq 0 \quad (1 \leq i \leq n) \\
& \sum_{i=1}^n \alpha_i y_i = 0
\end{aligned} \tag{3}$$

where the primal problem solution given dual optima $\boldsymbol{\alpha}^*$ is

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i.$$

- (c) Use KKT conditions to argue that $\text{supp}(\boldsymbol{\alpha}^*) := \{1 \leq i \leq n : \alpha_i^* \neq 0\}$ indicates the support vectors. Show how to solve for b^* . What's the support hyperplanes and margin?
- (d) (Kernel Trick) Let K be a positive semidefinite (PSD) kernel on \mathbb{R}^p generating a reproducing kernel Hilbert space (RKHS) \mathcal{H}_K , admitting eigen expansion

$$K(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^{\infty} \gamma_j \phi_j(\mathbf{x}) \phi_j(\mathbf{x}'). \quad (\mathbf{x}, \mathbf{x}' \in \mathbb{R}^p)$$

For any $f \in \mathcal{H}_K$, there exists $\{\theta_j\}_{j=1}^{\infty} \subseteq \mathbb{R}$ such that

$$f(\mathbf{x}) = \sum_{j=1}^{\infty} \theta_j \phi_j(\mathbf{x}), \quad \|f\|_{\mathcal{H}_K}^2 = \sum_{j=1}^{\infty} \frac{\theta_j^2}{\gamma_j} < +\infty.$$

Consider an RKHS analog to (2)

$$\begin{aligned}
\min_{f \in \mathcal{H}_K, b \in \mathbb{R}} \quad & \frac{1}{2} \|f\|_{\mathcal{H}_K}^2 \\
\text{s.t.} \quad & y_i [b + f(\mathbf{x}_i)] \geq 1 \quad (1 \leq i \leq n)
\end{aligned} \tag{4}$$

Show that the Lagrangian dual problem now becomes⁴

$$\begin{aligned}
\max_{\boldsymbol{\alpha} \in \mathbb{R}^n} \quad & L_{\mathcal{D}}(\boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \\
\text{s.t.} \quad & \alpha_i \geq 0 \quad (1 \leq i \leq n) \\
& \sum_{i=1}^n \alpha_i y_i = 0
\end{aligned} \tag{5}$$

where the primal problem solution given dual optima $\boldsymbol{\alpha}^*$ is⁵

$$f^*(\mathbf{x}) = \sum_{i=1}^n \alpha_i^* y_i K(\mathbf{x}_i, \mathbf{x}). \quad (\mathbf{x} \in \mathbb{R}^p)$$

⁴It greatly reduces the nonlinear problem to simply replace $[\langle \mathbf{x}_i, \mathbf{x}_j \rangle]_{n \times n}$ by the kernel matrix $[K(\mathbf{x}_i, \mathbf{x}_j)]_{n \times n}$ and motivates

⁵It hence also shows that $f^* \in \text{span}\{K(\mathbf{x}_i, \cdot)\}_{i=1}^n$, which is a generic result for loss minimization over RKHS [Wahba, 1990,

Friedman et al., 2009, Ex 5.15].

Computational Part

1. (20 pt) **Backfitting and Coordinate Descent in LASSO** [Wu and Lange, 2008, Friedman et al., 2010]

Recall that the univariate LASSO regression $\{Y_i\}_{i=1}^n$ on standardized regressor $\{X_i\}_{i=1}^n$ with $\sum_{i=1}^n X_i = 0$, $\frac{1}{n} \sum_{i=1}^n X_i^2 = 1$ ⁶ is soft-thresholding

$$\operatorname{argmin}_{\alpha, \beta \in \mathbb{R}} \frac{1}{2n} \sum_{i=1}^n (Y_i - \alpha - \beta X_i)^2 + \lambda |\beta| = (\bar{Y}, \mathcal{S}(\hat{\beta}_{\text{LS}}; \lambda))$$

where $\hat{\beta}_{\text{LS}} = \frac{1}{n} \sum_{i=1}^n X_i(Y_i - \bar{Y})$ is the ordinary least-square estimate, \mathcal{S} is a soft-thresholding operator

$$\mathcal{S}(z; \lambda) := \operatorname{sign}(z)(|z| - \lambda)_+ = \begin{cases} z - \lambda, & z > \lambda \\ 0, & -\gamma < z \leq \lambda \\ z + \lambda, & z \leq -\lambda \end{cases}$$

Derive the cyclic backfitting algorithm to solve multivariate LASSO regression given a standardized covariate matrix $\mathbf{X} = [X_{ij}]_{n \times p} \in \mathbb{R}^{n \times p}$ with $\sum_{i=1}^n X_{ij} = 0$, $\frac{1}{n} \sum_{i=1}^n X_{ij}^2 = 1$, response vector $\mathbf{Y} \in \mathbb{R}^n$ and ℓ^1 -regularization parameter $\lambda > 0$. Write an **R** function `lasso` and compare it with `glmnet` on your simulated

$$n = p = 100, \quad \{X_{ij}\}_{i,j=1}^{100}, \{Y_i\}_{i=1}^{100} \stackrel{iid}{\sim} \mathcal{N}(0, 1), \quad \lambda = 1/10.$$

Hint. The algorithm derived above is implemented in `glmnet` [Friedman et al., 2010]. In order to get exactly the same result from `glmnet`, standardize the data on your own to avoid internal scaling since `glmnet` would report coefficients in the original scale. Specify `lambda = 1/10` in `glmnet` to avoid internal generated λ sequence. Set the `thresh` option to `1e-20` to get an accurate fit.

2. (20 pt) (Textbook Ex. 7.9) **Prostate Cancer Data**

Carry out a best-subset regression analysis on the *Prostate Cancer Data* as Hw2 has done, while using AIC, BIC, 5-fold and 10-fold CVs, and Bootstrap .632 estimates of prediction error to tune the best size of subsets. Discuss the results.

3. (20 pt) **South African Heart Disease Data**

Perform Support Vector Machine analysis on the *South African Heart Disease Data* with various kernels and compare the prediction performance with the results using LDA, QDA, and Logistic regression in Hw3. Remember to tune the bandwidth parameters in nonlinear kernels using cross-validation.

References

Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge University Press, http://web.stanford.edu/~boyd/cvxbook/bv_cvxbook.pdf, 2004. 2

⁶It admits with the internal standardization of `glmnet`. Note that `scale` function scales as $\frac{1}{n-1} \sum_{i=1}^n X_i^2 = 1$.

- Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*. Springer-Verlag, <https://web.stanford.edu/~hastie/Papers/ESLII.pdf>, second edition, 2009. 1, 2, 3
- Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software*, 33(1):1, 2010. 4
- Grace Wahba. *Spline models for observational data*, volume 59. Siam, 1990. 3
- Tong Tong Wu and Kenneth Lange. Coordinate descent algorithms for lasso penalized regression. *The Annals of Applied Statistics*, 2(1):224–244, 2008. 4