

## STA 561: Homeworks 7 & 8 (Due April 17 at midnight)

Reminder: work together! Share ideas, brainstorm, explain/verify your answers but write up your own work. Your homework should be submitted as pdf file generated using either latex or an python notebook.

1. (A simple decision model.) Consider the following generative model for  $(X, A, Y)$

- $X^{1/5} \sim \text{Normal}(0, 1)$
- $A = 1_{|X| \leq \tau}$
- $Y = \beta_0 + \beta_1 A + \beta_2 X + \beta_3 AX + \epsilon$ , where  $\epsilon \sim \text{Normal}(0, 1)$ .

Suppose that  $\beta_0 = \beta_1 = \beta_2 = 1$  and  $\beta_3 = 0.5$ . What is the optimal decision rule? I.e., mapping  $\pi : \mathbb{R} \rightarrow \{0, 1\}$  such that if decisions are assigned according to  $\pi$  the value  $V(\pi)$  is maximized. Generate 1000 data sets of size  $n = 500$  from this model for  $\tau = 0.01$  and  $\tau = 0.025$  and use OLS to estimate  $\beta_0, \dots, \beta_3$ . In what proportion of your data sets was the p-value for  $\beta_3$  significant? What's happening here? Are the standard causal assumptions verified?

2. (Run on sentence, run on.) Suppose that we have a black-box regression model that inputs data of the form  $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$  and outputs an estimator  $\hat{f}_n(\mathbf{x})$  of  $\mathbb{E}(Y|\mathbf{X} = \mathbf{x})$ . Our goal in this problem to explore approximating  $\hat{f}_n$  with a kernel. For this problem you'll be exploring two approaches for constructing kernels: (i) born-again random forests in which you will generate many inputs  $\mathbf{Z}_1, \dots, \mathbf{Z}_B$  (where  $B$  is large) from the convex hull of the support of  $\mathbf{X}_1, \dots, \mathbf{X}_n$ , create outputs  $\hat{f}_n(\mathbf{Z}_1), \dots, \hat{f}_n(\mathbf{Z}_B)$ , then fit  $\left\{(\mathbf{Z}_b, \hat{f}_n(\mathbf{Z}_b))\right\}_{b=1}^B$  using a random forest from which you can extract the random forest kernel; (ii) ad-hoc kernels in which you will repeat the following steps for  $b = 1, \dots, B$

- bootstrap the data
- randomly select a subset of the predictors (columns of  $\mathbf{X}$ ) by drawing  $M$  entries without replacement from  $\{1, \dots, p\}$
- apply the black box to the bootstrapped and column-subset data to obtain  $\hat{f}_n^{(b)}$
- draw random seeds  $y_1, \dots, y_L$  uniformly from  $Y_1, \dots, Y_n$  and corresponding Voronoi partition of  $\mathbb{R}$

then for any  $\mathbf{x} \in \mathbb{R}^p$  define  $A_n^{(b)}(\mathbf{x})$  to the partition to which  $\widehat{f}_n^{(b)}(\mathbf{x})$  belongs and define the kernel distance between two points  $\mathbf{x}, \mathbf{z}$  to be

$$K(\mathbf{x}, \mathbf{z}) \triangleq \frac{1}{B} \sum_{b=1}^B 1_{\mathbf{x} \in A_n^{(b)}(\mathbf{z})}.$$

Note that the ad hoc kernel depends on  $L$ ,  $K$ , and  $B$  which you will need to tune/adjust.

Implement the two kernel methods and conduct a simulation study comparing local linear models fit using these kernel functions when the black box model is random forests and boosting (implemented in `xgboost`).