

## Problem Set 2

Due: November 9th, 2015

Problem sets are due at 5PM in the GSIs mailbox. You may work in groups, but each student should turn in their own write-up (including individually commented and executed code if applicable). A copy of any computer code should also be e-mailed to the GSI prior to the due date.

## 1 Pencil and paper problems

1. Let  $\mathbf{Z} \sim \mathcal{N}\left(\theta, \frac{\sigma^2}{N} I_K\right)$ . Show that  $\mathbb{E}[g(\mathbf{Z})(\mathbf{Z} - \theta)] = \frac{\sigma^2}{N} \mathbb{E}[\nabla_{\mathbf{Z}} g(\mathbf{Z})]$  (HINT: Use integration by parts).

2. Consider the following estimate of the Risk  $\mathbb{E}\left[\|\hat{\theta} - \theta\|^2\right]$  associated with the weakly differentiable estimate  $\hat{\theta}$ :

$$\hat{R}(\mathbf{Z}) = K\sigma^2 + 2\sigma^2 \sum_{k=1}^K \frac{\partial g(\mathbf{Z})}{\partial Z_k} + \sum_{k=1}^K (\hat{\theta}_k - Z_k)^2.$$

Show that this risk estimate is unbiased:  $\mathbb{E}_{\theta}[\hat{R}(\mathbf{Z})] = R(\hat{\theta}, \theta)$ .

3. (Wasserman, 2006). Let  $\mathbf{Z} \sim \mathcal{N}\left(\theta, \frac{\sigma^2}{N} I_K\right)$  and consider the following “soft threshold” estimator of  $\theta$ :

$$\hat{\theta}_k = \text{sgn}(Z_k) (|Z_k| - \lambda)_+, \quad k = 1, \dots, K.$$

In words this estimator shrinks the MLE of  $\theta_k$  toward zero when it is large (in absolute value) and shrinks it exactly to zero when it is small (in absolute value). Use Stein’s Unbiased Risk Estimate (SURE) to show that

$$\hat{R}_{\text{SURE}}(\mathbf{Z}, \lambda) = \frac{K}{N} \sigma^2 - \frac{2\sigma^2}{N} \sum_{k=1}^K \mathbf{1}(|Z_k| \leq \lambda) + \sum_{k=1}^K \min(Z_k^2, \lambda^2).$$

Provide a concrete prediction problem where you would expect the risk properties of the soft threshold estimator to be attractive.

4. (Wasserman, 2006) Let  $\mathbf{Z} \sim \mathcal{N}\left(\theta, \frac{\sigma^2}{N} I_K\right)$  and  $\mathcal{M}$  be the class of ordered subsets

$$\{\emptyset, \{1\}, \{1, 2\}, \{1, 2, 3\}, \dots, \{1, \dots, K\}\}$$

and consider the estimator  $\hat{\theta}_k(M) = Z_k \mathbf{1}(k \in M)$  for  $k = 1, \dots, K$ . Use SURE to show that

$$\hat{R}_{\text{SURE}}(\mathbf{Z}, M) = \frac{\sigma^2}{N} |M| + \sum_{k \in M^c} \left(Z_k^2 - \frac{\sigma^2}{N}\right)$$

with  $|M|$  denoting the cardinality of  $M$  and  $M^c$  the absolute complement of  $M$  in the universe  $\mathcal{M}$ .

## 2 Calorie demand

Download the **RPS\_calorie\_data.out** dataset from the course webpage. For this problem set you will require only two columns of the dataset, those with the headings **Y0tc** and **X0te**. The first equals the log of total calories consumed in a household and the latter equals the log of total expenditure. To learn more about the dataset see Section 4 of Graham and Powell (2012, *Econometrica*). Subramanian and Deaton (1996, *Journal of Political Economy*) provide more background on calorie demand analysis.

1. Let  $Y$  denote log calories and  $X$  denote log expenditure. Assume that

$$m(x) = \mathbb{E}[Y | X = x] = \sum_{k=1}^K \alpha_k g_k(x)$$

where  $g_k(x) = x^{k-1}$ .

2. Using the power series basis described above and the Gram-Schmidt algorithm construct a new basis that is orthogonal to the design points (set  $K = 12$ ). Let  $W_i$  denote the  $K \times 1$  vector of orthonormal basis functions for the  $i^{\text{th}}$  household. Compute the least squares fit

$$m(X_i) = W_i' \hat{\theta}$$

with

$$\hat{\theta} = \left[ \sum_{i=1}^N W_i W_i' \right]^{-1} \times \left[ \sum_{i=1}^N W_i Y_i \right].$$

Plot this function onto a scatter of the unsmoothed data.

3. Now use the shrinkage estimator of Efromovich (1999) as described in lecture to estimate  $m(X_i)$ . Plot this function onto a scatter of the unsmoothed data. Comment on your findings.
4. Now compute the soft threshold estimate of  $m(X_i)$  (as defined above). Plot this function onto a scatter of the unsmoothed data.

## 3 Income and Geography

Download the dataset used in Hall and Jones (1997, *Quarterly Journal of Economics*) from

<http://web.stanford.edu/~chadj/HallJones400.asc>

From the second table in this file extract the logYL and Latitude columns, retaining only complete cases. Repeat the analysis in Problem 2 above with  $Y$  equal to logYL and  $X$  equal to Latitude.