

STA 561: Midterm I (Due March 8 at midnight)

Reminder: unlike your homeworks and projects, this one is on your own. You are allowed to use any materials you manage to dig up from class, the internet, tea leaves, or the library. However, you cannot ask living humans, sentient robots, undead humans, non-sentient humans, or your Uncle Mo. If you're unsure if a source is acceptable, the rule you should use is that a source is acceptable if the information in the source have been available when you acquired it if all humans but you vanished on Feb 28, i.e., you cannot post a question on stack exchange, create an anonymous website where people can post solutions (so that technically you found the info on the internet), or tattoo the answer on your dog and then upload photos of your dog to instagram¹.

Your exam should be submitted as pdf file generated using latex. Late submissions will not be accepted (special pre-arranged accommodations notwithstanding). Good luck and science.

¹also, don't tattoo dogs, what's the matter with you?

1. (Zabka has problems.) Suppose that we observe n i.i.d. input-output pairs $\{\mathbf{X}_i, Y_i\}_{i=1}^n$ drawn from fixed but unknown distribution P . Let (\mathbf{X}, Y) denote a generic input-output pair and assume $\mathbf{X} \in \mathbb{R}^p$ and $Y \in \mathbb{R}$. Suppose that your friend² Zabka believes that the data follow a heteroskedastic model of the form

$$Y = \mathbf{X}^\top \boldsymbol{\beta}^* + \sigma(\mathbf{X})\epsilon,$$

where ϵ is an independent error with mean zero and variance one, and $\sigma(\mathbf{x}) > 0$ is a function of \mathbf{x} . Suppose that Zabka's model is correct. He makes the following claims. State whether the claim is true or false and justify your answer (e.g., by providing a short proof or derivation).

Zabka claim 1: Because of the heteroskedastic variance, the least squares estimator $\hat{\boldsymbol{\beta}}_n = \arg \min_{\boldsymbol{\beta}} \mathbb{P}_n(Y - \mathbf{X}^\top \boldsymbol{\beta})^2$ is biased.

Zabka claim 2: It is not possible to construct a valid Wald-type confidence interval for $\boldsymbol{\beta}^*$ if we do not know or cannot estimate the function $\sigma(\mathbf{x})$.

Zabka claim 3: If we used weighted least squares, say

$$\hat{\boldsymbol{\beta}}_n^\omega = \arg \min_{\boldsymbol{\beta}} \mathbb{P}_n(Y - \mathbf{X}^\top \boldsymbol{\beta})^2 \omega(\mathbf{X}),$$

for some positive real-valued function $\omega(\mathbf{x}) > 0$ then $\hat{\boldsymbol{\beta}}_n^\omega$ will be biased unless $\omega(\mathbf{x}) \equiv 1/\sigma^2(\mathbf{x})$.

Extra credit: What is the optimal choice of $\omega(\mathbf{x})$ in this model in terms of minimizing the limiting covariance of $\sqrt{n}(\hat{\boldsymbol{\beta}}_n^\omega - \boldsymbol{\beta}^*)$? Recall that we say one matrix $\Sigma \leq \Omega$ for matrices Σ and Ω if $\Omega - \Sigma$ is nonnegative definite.

²To use the term loosely.

2. (McStreamy.) Assume the same regression setup as in the previous problem. The elastic-net estimator is given by

$$\hat{\boldsymbol{\beta}}_n^{\tau, \lambda} = \arg \min_{\boldsymbol{\beta}} \mathbb{P}_n (Y - \mathbf{X}^\top \boldsymbol{\beta})^2 + \tau \|\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1.$$

Show that this depends on the data only through the sufficient statistics $\mathbb{P}_n \mathbf{X}Y$ and $\mathbb{P}_n \mathbf{X} \mathbf{X}^\top$ and suggest an estimator that can be used with datasets that are too large to read into memory (you can assume that n is large and p is modest). How would you modify your algorithm to make it adaptive? (I.e., you have and adaptive ridge and adaptive lasso penalties.)

3. (It all goes wrong.) Construct a regression example in which $\mathbb{E}(Y|\mathbf{X} = \mathbf{x}) = \mathbf{x}^\top \boldsymbol{\beta}^*$ and the following hold: (i) $p = 100$, (ii) $\beta_j^* = 0$ for $j \geq 3$, (iii) the theoretical R^2 is 0.9, (iv) $n = 1000$, and (v) running sure independent screening (SIS) to select the top 10 variables and then fitting a model to these variables has R^2 less than 0.10 more than 50% of the time in repeated samples. That is, if one were to repeat the following many times:

- randomly draw a data set of size $n = 1000$ from your model
- run SIS to select 10 variables
- fit a linear model with the selected variables
- compute the R^2 of the fitted model,

they would find that the R^2 is less than 0.10 more than 50% of the time.

Suggest a new variant of SIS that would not have this problem (at least in your data set).

4. (Short and noisy.) Write adaptive ridge as a noise addition method. I.e., show that

$$\hat{\beta}_n^\lambda = \arg \min_{\beta} \mathbb{P}_n(Y - \mathbf{X}^\top \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2 / \left(\hat{\beta}_{n,j}^{\text{OLS}} \right)^2,$$

where β_n^{OLS} is the OLS estimator, can be written as

$$\arg \min_{\beta} \mathbb{P}_n \mathbb{E}_{\mathbf{Z}} \{Y - (\mathbf{X} + \mathbf{Z})^\top \beta\}^2,$$

for some distribution \mathbf{Z} .

Extra credit: Consider the following iterative procedure:

- (a) $\hat{\beta}_n^{(0)} = \hat{\beta}_n^{\text{OLS}}$
- (b) $\hat{\beta}_n^{(k)} = \arg \min_{\beta} \mathbb{P}_n(Y - \mathbf{X}^\top \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2 / |\hat{\beta}_j^{(k-1)}|$ for $k = 1, \dots$

show that $\hat{\beta}_n^{(k)}$ converges to the lasso solution as $k \rightarrow \infty$.

Extra credit: Implement the preceding algorithm and compare with the output of a standard lasso package (e.g., sklearn, lars, glmnet, etc.).

5. (Thresh metal.) Suppose we observe $X_1, X_2, \dots, X_n \sim_{i.i.d.} \text{Normal}(\mu, \sigma^2)$ and consider the adaptive hard-thresholding estimator $\hat{\mu}_n^{\lambda_n} = \text{sgn}(\bar{X}_n) \{|\bar{X}_n| - \lambda_n/|\bar{X}_n|\}_+$ of μ . For any estimator $\tilde{\mu}_n$, where λ_n is a tuning parameter. Define $R(\tilde{\mu}_n, \mu) = \mathbb{E}(\mu - \tilde{\mu}_n)^2$. Note* you do not need to compute any integrals.

- (a) Show that $R(\hat{\mu}_n^{\lambda_n}, 0) < R(\bar{X}_n, 0)$ provided $\lambda_n > 0$, i.e., $\hat{\mu}_n^{\lambda_n}$ has lower expected loss at $\mu = 0$.
- (b) Suppose $\mu \neq 0$. Show that if $\lambda_n \sqrt{n} \rightarrow c$ as $n \rightarrow \infty$, where $c > 0$ is a finite constant then $\limsup_{n \rightarrow \infty} R(\hat{\mu}_n^{\lambda_n}, \mu) / R(\bar{X}_n, \mu) \leq 1$, i.e., $\hat{\mu}_n^{\lambda_n}$ has expected loss that is no worse than \bar{X}_n .

Hint: under the condition on λ_n we have

$$\frac{\lambda_n}{|\bar{X}_n|} = \frac{\lambda_n \sqrt{n}}{\sqrt{n} |\bar{X}_n|}$$

where the numerator is going to the constant c and the denominator is going to ∞ almost surely. Also recall that for any v we can write $v = \text{sgn}(v)|v|$.

- (c) From the first two parts of this problem it seems as though we have found an estimator that is superior to the mean!³ However now suppose that for each n we draw $X_{n,1}, X_{n,2}, \dots, X_{n,n} \sim_{i.i.d.} \text{Normal}(\mu_n, \sigma^2)$ where $\mu_n = \kappa/\sqrt{n}$. Note that $\bar{X}_n = n^{-1} \sum_{i=1}^n X_{n,i} \sim \text{Normal}(\mu_n, \sigma^2/n)$ for all n . The reason for considering such sequences is that they allow use to study what happens in cases where it is difficult to distinguish the true μ from zero. Show what happens to $R(\hat{\mu}_n^{\lambda_n}, \mu_n) / R(\bar{X}_n, \mu_n)$ as $n \rightarrow \infty$. Your answer will depend on c and κ . What does this tell you about our new estimator?

³In fact, as the mean attains the Cramer-Rao efficiency bound we say that $\hat{\mu}_n^{\lambda_n}$ is super efficient. So, that's uh, super, right?