

Second Midterm Review Sheet (Part 1 of 2)

*Ec240a – Second Half, Fall 2021*

In preparing for the exam review your lecture notes, assigned course readings, problem sets and prepare answers to the questions on the review sheet(s). You may bring to the exam a single  $8.5 \times 11$  inch sheet of paper with notes on it. No calculation aides are allowed in the exam (e.g., calculators, phones, laptops etc.). You may work in pencil or pen, however I advise the use of pencil. Please be sure to bring sufficient blue books with you to the exam if possible. Scratch paper will be provided.

[1] You observe a simple random sample of size  $N$  from the population

$$Y_0 \sim N(\mu, \sigma^2)$$

as well as a second, independent, simple random sample, also of size  $N$ , from the population

$$Y_1 \sim N(\mu, 4\sigma^2).$$

The value of  $\sigma^2$  is known. Consider the family of estimates of  $\mu$

$$\hat{\mu}(c_0, c_1) = c_0 \bar{Y}_0 + c_1 \bar{Y}_1,$$

where  $\bar{Y}_0 = \frac{1}{N} \sum_{i=1}^N Y_{0i}$  and  $\bar{Y}_1 = \frac{1}{N} \sum_{i=1}^N Y_{1i}$ .

[a] Show that mean squared error equals

$$\mathbb{E} \left[ (\hat{\mu}(c_0, c_1) - \mu)^2 \right] = \frac{c_0^2 \sigma^2}{N} + \frac{c_1^2 4\sigma^2}{N} + (1 - c_0 - c_1)^2 \mu^2. \quad (1)$$

[b] Derive the oracle estimator (within the family) which minimizes (1).

[c] Show that

$$\hat{R}(c_0, c_1) = \frac{c_0^2 \sigma^2}{N} + \frac{c_1^2 4\sigma^2}{N} + (1 - c_0 - c_1)^2 \frac{1}{2} \left\{ \bar{Y}_0^2 + \bar{Y}_1^2 - \frac{5\sigma^2}{N} \right\} \quad (2)$$

is an unbiased estimate of (2). Can you propose another unbiased risk estimate? Why would you prefer one unbiased risk estimate over another?

[d] Describe in *words* how one might use (2) to construct an implementable estimator of  $\mu$ .

[2] Let  $m(Z) = \mathbb{E}[X|Z]$  and consider the linear regression

$$\mathbb{E}^*[Y|X, m(Z), A] = \alpha_0 + \beta_0 X + \gamma_0 m(Z) + A.$$

[a] Show that

$$\mathbb{E}^*[m(Z)|X] = \delta_0 + \xi_0 X$$

with

$$\begin{aligned} \delta_0 &= (1 - \xi_0) \mathbb{E}[X] \\ \xi_0 &= \frac{\mathbb{V}(\mathbb{E}[X|Z])}{\mathbb{E}[\mathbb{V}(X|Z)] + \mathbb{V}(\mathbb{E}[X|Z])}. \end{aligned}$$

[b] Assume the population under consideration is working age adults who grew up in the San Francisco Bay Area. Let  $Y$  denote a adult log income, let  $X$  denote the log income of one's parents as a child and let  $Z$  be a vector of dummy variables denoting an individual's neighborhood of residence as a child. Provide an interpretation of  $\xi_0$  as a measure of residential stratification by income.

[c] Establish the notation  $\rho = \text{corr}(A, X)$ ,  $\mu_A = \mathbb{E}[A]$ ,  $\mu_X = \mathbb{E}[X]$ ,  $\sigma_A^2 = \mathbb{V}(A)$  and  $\sigma_X^2 = \mathbb{V}(X)$ . Show that

$$\mathbb{E}^*[Y|X] = \alpha_0 + \gamma_0(1 - \xi_0)\mu_X + \left(\mu_A - \rho \frac{\sigma_A}{\sigma_X} \mu_X\right) + \left\{ \beta_0 + \gamma_0 \xi_0 + \rho \frac{\sigma_A}{\sigma_X} \right\} X.$$

[d] Your research assistant computes an estimate of  $\mathbb{E}^*[Y|X]$  using random sample from San Francisco. She computes a separate estimate using a random sample from New York City. Assume that there is more residential stratification by income in New York than in San Francisco. How would you expect the intercept and slope coefficients to differ across the two regression fits?

[3] Let  $Y$  equal tons of banana's harvested in a given season for a randomly sampled Honduran banana planation. Output is produced using labor and land according to  $Y = AL^{\alpha_0} D^{1-\alpha_0}$ , where  $L$  is the number of employed workers and  $D$  is the size of the plantation in acres and we assume that  $0 < \alpha_0 < 1$ . The price of a unit of output is  $P$ , while that of a unit of labor is  $W$ . These prices may vary across plantations (e.g., due to transportation costs, labor market segmentation etc.). We will treat  $D$  as a fixed factor;  $A$  captures sources of plantation-level differences in farm productivity due to unobserved differences in, for example, soil quality and managerial capacity. Plantation owners choose the level of employed labor to maximize profits. The observed values of  $L$  are therefore solutions to the optimization problem:

$$L = \arg \max_l P \cdot A l^{\alpha_0} D^{1-\alpha_0} - W \cdot l.$$

[a] Show that the amount of employed labor is given by

$$L = \left\{ \alpha_0 \frac{P}{W} A \right\}^{\frac{1}{1-\alpha_0}} D. \quad (3)$$

[b] Let  $a_0 = \frac{1}{1-\alpha_0} \ln \alpha_0 + \frac{1}{1-\alpha_0} \mathbb{E}[\ln A]$ ,  $b_0 = \frac{1}{1-\alpha_0}$ , and  $V = \frac{1}{1-\alpha_0} \{\ln A - \mathbb{E}[\ln A]\}$ . Show that the log of the labor-land ratio is given by

$$\ln \left( \frac{L}{D} \right) = a_0 + b_0 \ln \left( \frac{P}{W} \right) + V \quad (4)$$

and that, letting  $c_0 = \mathbb{E}[\ln A]$  and  $U = \ln A - \mathbb{E}[\ln A]$ , the log of planation yield (output per unit of land) is given by

$$\ln \left( \frac{Y}{D} \right) = c_0 + \alpha_0 \ln \left( \frac{L}{D} \right) + U. \quad (5)$$

[c] Briefly discuss the content and plausibility of the restriction

$$\mathbb{E}[\ln A | \ln(P/W)] = \mathbb{E}[\ln A]. \quad (6)$$

[d] Using (4), (5) and (6) show that the coefficient on  $\ln(L/D)$  in  $\mathbb{E}^*[\ln(Y/D) | \ln(L/D)]$  equals

$$\alpha_0 + (1 - \alpha_0) \frac{\mathbb{V}(\ln A)}{\mathbb{V}(\ln A) + \mathbb{V}(\ln(P/W))}.$$

Provide some economic intuition for this result.

[e] Using (4), (5) and (6) show that the coefficient on  $\ln(L/D)$  in  $\mathbb{E}^*[\ln(Y/D)|\ln(L/D), V]$  equals  $\alpha_0$ . Provide some economic intuition for this result.

[f] Assume that all plantations face the same output price ( $P$ ) and labor cost ( $W$ ). What value does the coefficient on  $\ln(L/D)$  in  $\mathbb{E}^*[\ln(Y/D)|\ln(L/D)]$  equal now? Why?

[4] Let  $Y$  be a scalar random variable,  $X$  a  $K$  vector of covariates (which includes a constant), and  $W$  a vector of additional covariates (which excludes a constant). Consider the long (linear) regression

$$\mathbb{E}^*[Y|W, X] = X'\beta_0 + W'\gamma_0. \quad (7)$$

Next define the short and auxiliary regressions

$$\mathbb{E}^*[Y|X] = X'b_0 \quad (8)$$

$$\mathbb{E}^*[W|X] = \Pi_0 X. \quad (9)$$

[a] Let  $V = W - \mathbb{E}^*[W|X]$  be the projection error associated with the auxiliary regression. Show that

$$\begin{aligned} \mathbb{E}^*[Y|V, X] &= \mathbb{E}^*[Y|X] + \mathbb{E}^*[Y|1, V] - \mathbb{E}[Y] \\ &= \mathbb{E}^*[Y|X] + \mathbb{E}^*[Y|V] \end{aligned}$$

where  $\mathbb{E}^*[Y|1, V]$  denotes the linear regression of  $Y$  onto a constant and  $V$ , while  $\mathbb{E}^*[Y|V]$  denotes the corresponding regression without a constant (HINT: Observe that  $\mathbb{C}(X, V) = 0$ ).

[b] Next show that  $\mathbb{E}^*[Y|V, X] = \mathbb{E}^*[Y|W, X]$  and hence that the coefficient on  $V$  in  $\mathbb{E}^*[Y|V, X]$  coincides with that on  $W$  in  $\mathbb{E}^*[Y|W, X]$ .

[c] Let  $U = Y - \mathbb{E}^*[Y|X]$  be the projection error associated with the short regression. Derive the coefficient on  $V$  in the linear regression of  $U$  onto  $V$  (excluding a constant).

[d] Discuss the possible practical value of the results shown in [b] and [c] above.

[5] This question is about the Rambly Shambly Hex Bolt Corporation. Let  $Y_t$  be the number of hex bolts produced by Rambly Shambly Hex in year  $t$ ,  $M_t$  tons of steel used in production,  $K_t$  total factory capital stock, and  $L_t$  total person-hours worked. We assume that

$$Y_t = A_t M_t^\alpha K_t^\beta L_t^\gamma.$$

[a] Rambly Shambly Hex is owned by an eccentric billionaire who chooses  $M_t$ ,  $K_t$  and  $L_t$  each year randomly using an eternally unchanging roulette-wheel-like-device (i.e., inputs are chosen independently of each other and independently of  $A_t$ ). Further assume that the distribution of  $A_t$  is i.i.d. over time. Show that under this input choice mechanism that

$$\mathbb{E}^*[\ln Y_t | \ln M_t, \ln K_t, \ln L_t] = \lambda + \alpha \ln M_t + \beta \ln K_t + \gamma \ln L_t$$

with  $\lambda = \mathbb{E}[\ln A_t]$ . Is this same result likely to hold if Rambly Shambly instead chose input levels to maximize profits? Why or why not?

[b] Further show that under the completely random input choice scheme described above that:

$$\mathbb{E}^* [\ln Y_t | \ln M_t, \ln K_t, \ln L_t] = \mathbb{E}^* [\ln Y_t | \ln M_t] + \mathbb{E}^* [\ln Y_t | \ln K_t] + \mathbb{E}^* [\ln Y_t | \ln L_t] - 2\mathbb{E} [\ln Y_t] .$$

[c] Let  $X_t = (\ln M_t, \ln K_t, \ln L_t)'$  and  $\sigma^2 = \mathbb{V}(\ln A_t)$ . Argue that under the completely random input choice scheme:

$$\sqrt{T} (\hat{\theta} - \theta_0) \xrightarrow{D} N(0, \Lambda) ,$$

for  $\theta = (\alpha, \beta, \gamma)'$ ,  $\Lambda = \sigma^2 \mathbb{V}(X_t)^{-1}$  and  $\hat{\theta}$  estimated by the OLS fit of  $\ln Y_t$  onto a constant and  $X_t$  for  $t = 1, \dots, T$ . What are the values of the off-diagonal elements of  $\mathbb{V}(X_t)$ ?

[d] For  $T = 3,859$ , an OLS fit gives

$$\hat{\theta} = \begin{pmatrix} 0.32 \\ 0.36 \\ 0.42 \end{pmatrix}, \quad \hat{\Lambda} = \begin{pmatrix} 5 & 1/200 & 2/1000 \\ 1/200 & 3 & 3/1000 \\ 2/1000 & 3/1000 & \frac{198}{100} \end{pmatrix} .$$

Construct a Wald Statistic (carefully explaining each step in the construction) for the null hypothesis of constant returns to scale (i.e.,  $H_0 : \alpha + \beta + \gamma = 1$ ). What is the appropriate reference distribution and critical value for a two-sided test with size  $\alpha = 0.05$ ? Do you reject the null?

[6] Let  $A \in \{a_l, a_h\}$  denote an individual's unobserved 'entrepreneurial acumen', and  $X$  be a binary indicator taking a value of one if an individual completed an undergraduate degree and zero otherwise. Let  $Y$  equal annual earnings. The following table gives the conditional mean of  $Y$  for each of the four possible 'entrepreneurial acumen' and schooling combinations

	$A = a_h$	$A = a_l$
$X = 1$	\$45,000	\$35,000
$X = 0$	\$50,000	\$15,000

Assume that  $m(x, a) = \mathbb{E}[Y | X = x, A = a]$  is a structural function in the following sense: in subpopulations homogenous in 'entrepreneurial acumen',  $m(x, a)$ , traces out how average earnings would change with external manipulations in college completion behavior. The population frequency of each of the four schooling and 'entrepreneurial acumen' combinations is

	$A = a_h$	$A = a_l$
$X = 1$	0.20	0.10
$X = 0$	0.05	0.65

[a] While on the elevator in Evans Hall you heard a grumpy individual (possibly a professor) claim "the best students should just start a tech firm in their parents' garages, we can train the rest to become corporate lawyers". Comment with reference to the population described above.

[b] Calculate the average annual earnings level in this economy,  $\mathbb{E}[Y]$ , and the averages conditional on college completion,  $\mathbb{E}[Y | X = 1]$ , and not,  $\mathbb{E}[Y | X = 0]$ .

[c] Calculate average earnings in a counterfactual world where  $\Pr(X = 1 | A = a_l) = 1$  and  $\Pr(X = 1 | A = a_h) =$

0.

[d] What is the expected earnings gain associated with college completion for a random draw from the population?

[e] Let  $W = 1$  if an individual operated a lemonade stand at some point during childhood and zero otherwise. Assume that (i)  $0 < \Pr(X = 1 | W = w) < 1$  for  $w \in \{0, 1\}$  and (ii) that  $X$  is conditionally independent of  $A$  given  $W$ . Show that for  $q(x, w) = \mathbb{E}[Y | X = x, W = w]$  we have  $\mathbb{E}[q(x, W)] = \mathbb{E}[m(x, A)]$ . How would your answer change if  $\Pr(X = 1 | W = 1)$  were equal to one?

[f] Maintaining the assumptions of part [e] above show that

$$\mathbb{E} \left[ \frac{\mathbf{1}(X = x) Y}{\Pr(X = x | W)} \right] = \mathbb{E}[m(x, A)].$$

Provide an intuitive discussion of this result.

[g] Available is a random sample of size  $N$  from the population of high school graduates. For each unit we observe  $Z = (W, X, Y)'$ . Let

$$R_1 = (\mathbf{1}(X = 0) \mathbf{1}(W = 0), \mathbf{1}(X = 0) \mathbf{1}(W = 1), \mathbf{1}(X = 1) \mathbf{1}(W = 0), \mathbf{1}(X = 1) \mathbf{1}(W = 1))'$$

where  $\mathbf{1}(\bullet)$  denotes the indicator function and

$$\mathbf{S} = \begin{pmatrix} Y \\ W \\ X \\ WX \end{pmatrix}, \mathbf{R} = \begin{pmatrix} R_1' & \mathbf{0}_3' \\ \mathbf{0}_3 \mathbf{0}_4' & I_3 \end{pmatrix}$$

with  $\mathbf{0}_k$  a  $k \times 1$  vector of zeros and  $I_k$  a  $k \times k$  identity matrix. Establish the following notation:  $\mu_{xw} = q(x, w)$ ,  $\sigma_{xw}^2 = \mathbb{V}(Y | X = x, W = w)$ ,  $p_x = \Pr(X = x)$ ,  $q_w = \Pr(W = w)$ , and  $r_{xw} = \Pr(X = x, W = w)$ . Assume (i)  $\sigma_{xw}^2$  is finite and (ii) that  $r_{xw} > 0$  for all four  $x$  and  $w$  combinations.

Consider the estimate

$$\hat{\beta} = \left[ \frac{1}{N} \sum_{i=1}^N \mathbf{R}_i' \mathbf{R}_i \right]^{-1} \times \left[ \frac{1}{N} \sum_{i=1}^N \mathbf{R}_i' \mathbf{S}_i \right].$$

[i] Show that  $\hat{\beta} \rightarrow \beta_0$  and provide an expression for each of the seven elements of  $\beta_0$  in terms of the notation established above (i.e. in terms of  $\mu_{xw}$ ,  $\sigma_{xw}^2$  etc.). How would your analysis change if it were the case that  $r_{11} = 0$ ?

[ii] Show that  $\sqrt{N}(\hat{\beta} - \beta_0)$  converges in distribution to a normal random variable. Provide an explicit expression for the covariance matrix of this normal distribution in terms of the notation established above (i.e. in terms of  $\mu_{xw}$ ,  $\sigma_{xw}^2$  etc.).

[iii] Using the elements of  $\hat{\beta}$  construct estimates of  $\mathbb{E}(q(1, W))$  and  $\mathbb{E}(q(0, W))$ . Establish the consistency of these estimates.

[iv] You are interested in the joint hypothesis that  $\mu_{10} = \mu_{00}$  and  $\mu_{11} = \mu_{01}$ . Discuss the substance of this hypothesis in light of the empirical set-up developed in parts [a] to [f] above. Show that this hypothesis may be represented as a restriction of the form  $C\beta_0 = c$  for some matrix  $C$  and column vector  $c$ .

[v] Your professor provides you will the following (consistent) estimates:  $\hat{\mu}_{00} = 10,000$ ,  $\hat{\mu}_{01} = 20,000$ ,

$\hat{\mu}_{10} = 50,000$ ,  $\hat{\mu}_{11} = 10,000$ ,  $\hat{\sigma}_{00}^2 = 1,000$ ,  $\hat{\sigma}_{01}^2 = 2,000$ ,  $\hat{\sigma}_{10}^2 = 500$ ,  $\hat{\sigma}_{11}^2 = 1,000$ ,  $\hat{r}_{00} = 0.2$ ,  $\hat{r}_{01} = 0.3$ , and  $\hat{r}_{10} = 0.25$ . Construct a test statistic for the hypothesis described in part [iii] above and compare it to the appropriate critical value. For your reference the 0.95 quantiles of  $\chi^2$  random variables with parameters 1, 2 and 3 are, respectively, 3.84, 5.99 and 7.81.