

EKT-812 Problem Set 1: Suggested Solutions

There are 160 points available on this problem set, and it will be graded out of 120. That is, there are 40 bonus points available.

Due Date: Wednesday, February 20

Probability Review

1. Consider the following joint distribution:

$\downarrow X \ Y \rightarrow$	1	2	3	4
0	0.1	0.05	0.025	0.025
1	0.07	0.13	0.04	0.06
2	0.1	0.1	0.25	0.05

- Find the first two moments - the mean and variance - of the marginal distributions of X and Y . Check that the identity $V[X] = E[X^2] - E[X]^2$ holds.
- Confirm that $E[XY] = 3.19$. Use this in combination with your calculations in (a) above to find $\text{cov}(X, Y)$.
- Find the conditional probability mass function $P(Y = y|X = x)$ and use it to calculate the conditional mean function $E[Y|X = x]$ for each possible realization of X .
- Calculate $E[E[Y|X = x]]$. With respect to which distribution is the outer expectation taken? The inner expectation? What do you notice about your answer as compared to the means you calculated in (a) above?

[3 + 2 + 3 + 2 = 10 pts]

2. Consider a population of workers, each with two possible skills: call them s_0 and s_1 . Think of s_0 as, say, programming skill and s_1 as language skill, which is useful for becoming e.g. a lawyer. Assume that earnings are determined as follows: there is a prevailing wage $w_0 > 0$ in the tech sector and $w_1 > 0$ in the legal sector. Workers are paid by skill, so more skilled workers earn more, and in particular,

$$Y_0 = w_0 e^{s_0} \tag{1}$$

$$Y_1 = w_1 e^{s_1} \tag{2}$$

where Y_0 and Y_1 are *potential* earnings in each sector.

Assume that the joint distribution of skills in the population is bivariate normal, with $E[s_0] = E[s_1] = 0$. (Here skills are on a logarithmic scale, so the fact that they can be negative is not a problem.) Let the standard deviation of s_0 be σ_0 , the standard deviation of s_1 be σ_1 , and let $\text{cov}(s_0, s_1) = \sigma_{01}$.

- Each worker has a choice about whether to become a programmer or a lawyer, given his or her endowment of skills (s_0, s_1) , but no one can work as both (and no one decides not to work at all). What is a reasonable decision rule for a worker facing this problem?

Answer: We can probably assume that people prefer more money to less. Then, they would choose to become a lawyer (let's use $D = 1$ as an indicator for their occupational choice) when $Y_1 \geq Y_0$. This happens if and only if $\log Y_1 \geq \log Y_0$; so

$$D = 1 \iff s_1 - s_0 \geq \log(w_0/w_1)$$

- Suppose workers behave according to the decision rule you wrote down in (a) above. Find an expression for the fraction of workers who become lawyers.

Hint: The sum of two normal random variables is also normally distributed.

Answer: We want to find $P(D = 1)$. Given the above decision rule, this is $P(s_1 - s_0 \geq \log(w_0/w_1))$. As a shorthand, let $\Delta = s_0 - s_1$ and let $\sigma_\Delta^2 = \sigma_0^2 + \sigma_1^2 - 2\sigma_{01}$. The distribution of Δ is normal with mean zero and variance σ_Δ^2 . Then the supply of workers to the legal sector is

$$P(D = 1) = P(s_1 - s_0 \geq \log(w_0/w_1)) = P\left(\frac{\Delta}{\sigma_\Delta} \geq \frac{\log(w_0/w_1)}{\sigma_\Delta}\right) = 1 - \Phi\left(\frac{\log(w_0/w_1)}{\sigma_\Delta}\right)$$

where Φ is the standard normal CDF.

[5 + 10 = 15 points]

3. Consider the random variable X with density

$$f_X(x) = \begin{cases} \frac{x}{9} & \text{if } x \in [0, 3] \\ \frac{6}{9} - \frac{x}{9} & \text{if } x \in (3, 6] \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

- (a) Find the CDF of X , $F_X(x)$. Draw the density f_X and the CDF F_X on the same set of axes.
- (b) Find the inverse CDF. (The inverse CDF is the function $x = G(p)$ such that $F_X(G(p)) = p$, i.e. given a probability p , G returns the value x such that $P(X \leq x) = p$.)

[5 + 5 = 10 points]

4. Table 1.2 of Chapter 1 of Deaton (1997), which is included in the assignment repository, contains summary statistics about the joint distribution of income and consumption for Ivorian households over the two years 1985-1986. Use the notation y_{it} for the income of household i in year t , and c_{it} for the consumption of household i in year t .

- (a) What is the correlation between consumption in the two years?

Hint: use the fact that $V[c_{i,t+1} - c_{i,t}] = V[c_{i,t+1}] + V[c_{i,t}] - 2\text{cov}(c_{i,t+1}, c_{i,t})$.

Answer: We have $V[\Delta c] = 987^2$; $V[c_{t+1}] = 1513^2$; $V[c_t] = 1236^2$. So,

$$\text{cov}(c_{t+1}, c_t) = \frac{-1}{2} \{987^2 - 1513^2 - 1236^2\}. \quad (4)$$

And, dividing by $\sqrt{V[c_t]V[c_{t+1}]} = 1513 \cdot 1236$, we get

$$\begin{aligned} \rho &= \frac{-1}{2} \left\{ \frac{987^2}{1513 \cdot 1236} - \frac{1513}{1236} - \frac{1236}{1513} \right\} \\ &\approx 0.76. \end{aligned} \quad (5)$$

- (b) Suppose these data come from a simple random sample of households. Further, suppose there is some measurement error in the income data, so that

$$y_{it} = y_{it}^* + \varepsilon_{it} \quad (6)$$

where y_{it}^* true income, and ε_{it} is “noise” that is independent of y_{it}^* each period, with $E[\varepsilon_{it}] = 0$ and $V[\varepsilon_{it}] = \sigma^2$.

However, measurement error may be persistent over time (perhaps because it is driven by negligence on the part of the reporting household), so let ρ be the correlation between $\varepsilon_{i,t+1}$ and ε_{it} .

Using the notation $\Delta y_{i,t+1} = y_{i,t+1} - y_{it}$ for the observed change in income for household i , and similarly for true income y_{it}^* , how would you calculate the variance in true income growth from panel data on y_{it} , if you knew ρ and σ^2 ?

Answer: The variance in observed income is $V[\Delta y] = V[\Delta y^* + \Delta \varepsilon] = V[\Delta y^*] + V[\Delta \varepsilon]$, by the independence of ε and y^* . Then, notice that $V[\Delta \varepsilon] = V[\varepsilon_{t+1}] + V[\varepsilon_t] - 2\text{cov}(\varepsilon_{t+1}, \varepsilon_t) = 2\sigma^2 - 2\sigma^2\rho = 2\sigma^2(1 - \rho)$. So if we have panel data on y and we know ρ and σ^2 , we can compute $V[\Delta y^*]$ as $V[\Delta y] - 2\sigma^2(1 - \rho)$.

- (c) Consider two methods of estimating the mean change in income, $E[\Delta y_{i,t+1}^*]$, over this period: (i) taking two independent cross-sectional surveys, each of size n , and (ii) collecting a two-period panel, also of size n . For method (i), say the households surveyed at time t are labelled $i = 1, \dots, n$ and those surveyed at $t + 1$ are labelled $i = n + 1, \dots, 2n$. The first method leads to the estimator

$$\bar{\Delta} = \frac{1}{n} \sum_{i=n+1}^{2n} y_{i,t+1} - \frac{1}{n} \sum_{i=1}^n y_{i,t} \quad (7)$$

while the second leads to the estimator

$$\hat{\Delta} = \frac{1}{n} \sum_{i=1}^n \Delta y_{i,t+1} \quad (8)$$

Under what conditions - on σ^2 , ρ , and the joint distribution of $(y_{i,t+1}^*, y_{i,t}^*)$ - will the first method be more precise than the second? (Here, “precise” means having a lower variance.) Does your answer depend on the sample size? If so, why?

Answer: We have

$$\begin{aligned} V[\bar{\Delta}] &= n^{-1}V[y_{i,t+1}] + n^{-1}V[y_{it}] \\ &= n^{-1}(V[y_{t+1}^*] + V[y_t^*] + 2\sigma^2) \end{aligned} \quad (9)$$

where the first equality follows from the independence of the two cross-sections. The variance of the panel estimator is

$$\begin{aligned} V[\hat{\Delta}] &= n^{-1}V[\Delta y] \\ &= n^{-1}(V[\Delta y^*] + 2\sigma^2(1 - \rho)) \end{aligned} \quad (10)$$

Thus, $V[\bar{\Delta}] \leq V[\hat{\Delta}]$ whenever

$$V[y_{t+1}^*] + V[y_t^*] + 2\sigma^2 \leq V[\Delta y^*] + 2\sigma^2(1 - \rho) \quad (11)$$

which we can rearrange to give

$$\rho\sigma^2 \leq -\text{cov}(y_{t+1}^*, y_t^*) \iff \rho\sigma^2 + \text{cov}(y_{t+1}^*, y_t^*) \leq 0 \quad (12)$$

i.e when the sum of the autocovariances (of the measurement error ε_t and true income y_t^*) is less than zero, we should prefer two independent samples to a panel.¹

To explore the intuition for this result, imagine measurement errors are independent over time, so $\rho = 0$. Then we should prefer a panel to two independent cross-sections whenever $\text{cov}(y_{t+1}^*, y_t^*) > 0$, i.e true incomes are persistent (which is typically the case). This is because the cross-sectional variation in y^* represents mainly permanent differences across households (and measurement error). Thus, the changes in income are less variable than the levels.

On the other hand, if $\rho < 0$ so the measurement error is mean-reverting, repeated observations on the same households are less useful for measuring changes. In this case, the observed changes in y will tend to be dominated by fluctuations in measurement error rather than real changes in y^* .

[5 + 5 + 10 = 20 points]

¹You didn’t need to prove this, but you should confirm that both estimators are *unbiased*, meaning that their expected values are $E[\Delta y^*]$.

Data Manipulation in R

Filtering, Sorting, and Generating New Variables

5. Do exercise 1 from section 5.2.4 of (???)

[5 points]

6. Do exercises 3-4 of section 5.3.1 of (???)

[$2 \times 5 = 10$ points]

7. Do exercise 2 of section 5.5.2 of (???)

[5 points]

Grouped Summaries and Filters

8. Do exercise 6 of section 5.7.1 of (???)

[5 points]

9. Do exercise 4 of section 5.6.7 of (???)

[5 points]

Reshaping Data

10. Do exercise 2-4 of section 12.3.3 of (???)

[$3 \times 5 = 15$ points]

Counterfactuals

Included in this repository are six short texts. Read each one carefully, and, identify the claims being made. Classify them as either causal claims, or non-causal claims. If they are not causal, what are they - value judgements? A statement of fact? Something else?

For those claims you think are causal, comment on whether the authors present any evidence for their claim. If so, what might be some concerns or alternative interpretations of that evidence? If they do not present any evidence for their claims, what sort of information would you need to decide whether it was true or false?

Please note: You are welcome to have any views you like about these topics, but I will deduct points for extraneous opinion. The point of this exercise is to read the texts carefully and critically, and to be able to recognize causal statements (including implicit ones). It is **not** to discuss the particular issues mentioned in the texts.

[$6 \times 10 = 60$ points]

References