

2017 Qual Section 2

②

2. (25 points) Suppose that the pair (X, Y) is distributed such that $X \sim \text{normal}(0, 1)$ and $Y \sim \text{Bernoulli}(\theta)$, $0 < \theta < 1$. For example, X could be the log of the level of a certain biomarker and Y an indicator of some disease. Assume that the value of θ is known and given (e.g. we know that the disease prevalence is 0.001).

Here we try to answer the following question: Given the above specifications, what is the largest possible correlation between X and Y ?

Notation: Define $\rho = \text{corr}(X, Y)$. Use $\phi(\cdot)$ and $\Phi(\cdot)$ to denote the standard normal pdf and cdf, respectively. Define any new notation you use.

- (a) (1 point) Is $\rho = 1$ possible? Explain.
- (b) (6 points) Find the joint distribution for (X, Y) that maximizes $\text{corr}(X, Y)$ subject to the model stated above. Show that it (that distribution) has the property that $E[X|Y = 1] = \theta^{-1}\phi(\Phi^{-1}(1 - \theta))$.
- (c) (6 points) Obtain an explicit expression for $\rho^* = \text{corr}(X, Y)$ within the joint distribution found in the previous part. Compute the numerical value of ρ^* for the case $\theta = 0.001$.
- (d) (6 points) Now suppose we are interested in various diseases with different prevalences (θ) ranging in $(0, 1)$. Find the value of θ that leads to the largest possible value of ρ^* , and compute that largest value, to be denoted ρ^{**} (compute its numerical value).
- (e) (3 points) Note: This part is totally independent of the previous parts, even though the models have some similarity. You can reuse results obtained above if needed.
Suppose that the iid pairs (X_i, Y_i) , $i = 1, \dots, n$, are distributed such that $X_i \sim \text{normal}(0, 1)$, $Y_i \sim \text{Bernoulli}(\theta)$, $0 < \theta < 1$, and $\text{corr}(X_i, Y_i) = \rho$, where both θ and ρ are unknown parameters. Develop an estimating equation for (ρ, θ) based on the vectors $Z_i := (T_i, Y_i)^\top$, $i = 1, \dots, n$, where $T_i := X_i Y_i$. Obtain the estimates $(\hat{\rho}, \hat{\theta})$ in explicit form.
- (f) (3 points) Based on $\hat{\rho}$ from the previous part, develop a large-sample (as $n \rightarrow \infty$) 95% confidence interval for ρ . The interval should not depend on any unknown parameters. Describe and justify your procedure clearly.

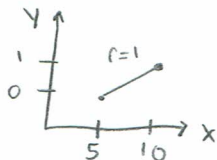
Note: $\Phi(-3.09) \approx 0.001$, $\Phi(-2.33) \approx 0.01$, $\Phi(-1.96) \approx 0.025$, $\Phi(-1.64) \approx 0.05$, $\Phi(-1.28) \approx 0.1$

2a) $I > \rho = 1$ possible? Explain

Given $X \sim N(0,1)$ and $Y \sim \text{Bern}(\theta)$, $0 < \theta < 1$.

Take the following example:

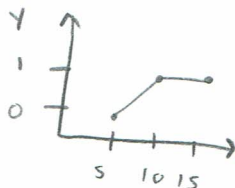
$$X = (5, 10), Y = (0, 1)$$



For any two points with positive slope, you can connect these two points w/ a line of perfect fit and hence have $r=1$.

However, for three or more points this is not possible b/c $Y \sim \text{Bern}(\theta)$.

Take, for example, $X = (5, 10, 15)$; $Y = (0, 1, 1)$.



As can be seen by this example,

Since $Y \sim \text{Bern}(\theta)$, there will never be a case

where $r=1$ b/c we can never create a linear fit between X and Y that perfectly fits all the data.]

2b) Find the joint distn. for (X, Y) that maximizes $\text{corr}(X, Y)$ subject to the model stated above. Show that the chosen distribution has the property that $E[X|Y=1] = \theta^{-1} \phi(\Phi^{-1}(1-\theta))$.

$$\text{Corr}(X, Y) = \rho = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} = \frac{E[XY] - \overbrace{E[X]}^0 \cdot \overbrace{E[Y]}^\theta}{\sqrt{1 \cdot \theta(1-\theta)}} = \frac{\overbrace{E[XY]}^{\text{want to maximize this}}}{\underbrace{\sqrt{\theta(1-\theta)}}_{\text{fixed}}}$$

$$\begin{aligned} \text{Maximize } E[XY] &= \overbrace{E[X|Y=0]}^{y=0} \cdot P(Y=0) + \overbrace{E[X|Y=1]}^{y=1} \cdot P(Y=1) \\ &= 0 \cdot \underbrace{E[X|Y=0]}_{(1-\theta)} \cdot P(Y=0) + 1 \cdot \underbrace{E[X|Y=1]}_{\theta} \cdot P(Y=1) \\ &= \theta E[X|Y=1] \\ &\quad \text{Need to maximize this} \end{aligned}$$

Notice that the joint pmf of (X, Y) is concentrated on the lines $Y=0$ & $Y=1$, like in this picture. Notice that, in order to maximize $E[X|Y=1]$ (the expected value of X given that $Y=1$), there is some value of $x > c$ \rightarrow the positive density of $f_X(x)$ is maximized (so that the mean is greatest).

Thus, need a value of c such that we

$$\text{have } \{(x, y) : x < c, y=0 \text{ \& } x > c, y=1\}$$

$$\text{Thus, } c \text{ must satisfy } \underbrace{P(Y=1)}_{\theta} = P(X > c) \quad \sim N(0,1)$$

$$\Rightarrow \theta = 1 - P(X \leq c) \Rightarrow \theta = 1 - \Phi(c) \Rightarrow \theta - 1 = -\Phi(c) \Rightarrow 1 - \theta = \Phi(c)$$

$$\Rightarrow c = \Phi^{-1}(1-\theta).$$

$$\begin{aligned} \text{Since } E[X|Y=1] &= E[X|X > c] = \int_x x \cdot P(X|X > c) dx = \int_x x \cdot \frac{P(X, X > c)}{P(X > c)} dx \\ &= \int_x x \cdot \frac{P(x) \cdot \mathbb{I}(x > c)}{\theta} dx = \int_c^\infty x \cdot \frac{\phi(x)}{\theta} dx = \int_c^\infty \frac{x \phi(x)}{1 - \Phi^{-1}(c)} dx \end{aligned}$$

$$= \frac{1}{1 - \Phi^{-1}(c)} \int_c^\infty x \phi(x) dx \quad (*)$$

$$\text{Note that } \frac{d}{dx}(\phi(x)) = \frac{d}{dx} \left[\frac{1}{\sqrt{2\pi}} e^{-x^2/2} \right] = -\frac{x}{\sqrt{2\pi}} e^{-x^2/2} = -x \phi(x)$$

$$\Rightarrow \frac{d\phi(x)}{dx} = -x \phi(x) \Rightarrow d\phi(x) = -x \phi(x) dx, \text{ Sub into } (*) \text{ to get:}$$

$$E[X|Y=1] = \frac{1}{1 - \Phi^{-1}(c)} \int_{-\infty}^c d\phi(x) = \frac{\phi(c)}{1 - \Phi^{-1}(c)} = \frac{\phi(\Phi^{-1}(1-\theta))}{\theta} \quad \checkmark$$

$$\text{where } \text{Corr}(X, Y) \text{ is maximized for } f(x, y) = \phi(x) \mathbb{I}(x > \Phi^{-1}(1-\theta)) + \phi(x) \mathbb{I}(x \leq \Phi^{-1}(1-\theta))$$

2c) obtain an explicit expression for $\rho^* = \text{Corr}(X, Y)$ within the joint distribution found in the previous part. Compute the numerical value of ρ^* for $\theta = 0.001$.

From last part, know that the chosen distribution has the property that

$$E[X|Y=1] = \frac{1}{\theta} \phi(\Phi^{-1}(1-\theta)).$$

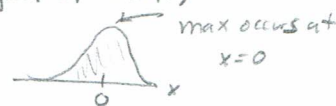
Then, since $\text{Corr}(X, Y) = \rho^* = \frac{\theta E[X|Y=1]}{\sqrt{\theta(1-\theta)}} = \frac{\theta \cdot \frac{1}{\theta} \phi(\Phi^{-1}(1-\theta))}{\sqrt{\theta(1-\theta)}}$

$$\Rightarrow \text{Corr}(X, Y) = \rho^* = \frac{\phi(\Phi^{-1}(1-\theta))}{\sqrt{\theta(1-\theta)}}$$

For $\theta = 0.001$, have $\text{Corr}(X, Y) = \rho^* = \frac{\phi(\Phi^{-1}(1-0.001))}{\sqrt{0.001(1-0.001)}} = \frac{\phi(3.09)}{\sqrt{0.001(1-0.001)}}$? Can't do that shit in my head. I'm not a calculator!

2d) Now, suppose we are interested in various diseases w/ different prevalences ranging from (0, 1). Find the value of θ that leads to the largest possible value of ρ^* and label it ρ^{**} (compute numerical value).

From part c), $\rho^* = \frac{\phi(\Phi^{-1}(1-\theta))}{\sqrt{\theta(1-\theta)}} = \frac{\phi(c)}{\sqrt{\theta(1-\theta)}}$. Note that $\phi(x)$ is the pdf of $N(0, 1)$



$$\Rightarrow \phi(x) \text{ is maximized at } x=0. \Rightarrow c=0 \Rightarrow \Phi^{-1}(1-\theta)=0 \Rightarrow 1-\theta=\Phi(0) \Rightarrow \theta=1-\Phi(0)$$

$$\Rightarrow \theta=0.5$$

$$\text{Thus, } \rho^{**} = \frac{\phi(0)}{\sqrt{0.5(1-0.5)}} = \frac{1}{2} \phi(0) = \frac{1}{2} \cdot \frac{1}{\sqrt{2\pi}} e^0 = \boxed{\frac{1}{\sqrt{2\pi}}}$$

2017 Quid, Section 2

2e) Suppose iid pairs (X_i, Y_i) , $i = 1, \dots, n \rightarrow X_i \sim N(0, 1)$ and $Y_i \sim \text{Bern}(\theta)$, $0 < \theta < 1$.

Given $\text{Corr}(X_i, Y_i) = \rho$, where both θ and ρ are unknown parameters.

Want to develop an estimating eqn. for (ρ, θ) based on $Z_i = (T_i, Y_i)^T = (X_i, Y_i, Y_i)^T$.

Obtain estimates $(\hat{\rho}, \hat{\theta})$ in explicit form.

$$\rho = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}} = \frac{E[XY] - E[X]E[Y]}{\sqrt{\text{Var}(X) \text{Var}(Y)}} \approx \frac{\frac{1}{n} \sum_i X_i Y_i}{\sqrt{\frac{1}{n} \sum_i Y_i (1 - \frac{1}{n} \sum_i Y_i)}} = \frac{\bar{Y}}{\sqrt{\bar{Y}(1-\bar{Y})}}$$

$$\text{Thus, } \hat{\rho} = \frac{\bar{Y}}{\sqrt{\bar{Y}(1-\bar{Y})}} \text{ where } \hat{\theta} = \bar{Y}$$

2f) Based on \hat{p} from the previous part, develop a large sample ($n \rightarrow \infty$)

95% CI for p .

NOTE: The interval should not depend on any unknown parameters.

$$95\% \text{ CI}(p) = (\hat{p} - 1.96 \sqrt{\text{Var}(\hat{p})}, \hat{p} + 1.96 \sqrt{\text{Var}(\hat{p})})$$

From the previous part, know that $\hat{p} = \frac{\bar{T}}{\sqrt{\bar{T}(1-\bar{T})}}$, Will use delta method on $(\frac{\bar{T}}{\bar{Y}})$ to derive asymptotic variance of \hat{p} .

Find dist. of \bar{T} : Know by CLT that $\sqrt{n}(\bar{T}_n - E[T_i]) \xrightarrow{d} N(0, \text{Var}[T_i])$

Find dist. of \bar{Y} : Know by CLT that $\sqrt{n}(\bar{Y}_n - \frac{E[Y_i]}{\theta}) \xrightarrow{d} N(0, \frac{\text{Var}[Y_i]}{\theta(1-\theta)})$

Then, by multivariate CLT, have

$$\sqrt{n} \left(\begin{pmatrix} \bar{T}_n \\ \bar{Y}_n \end{pmatrix} - \begin{pmatrix} E[T_i] \\ \theta \end{pmatrix} \right) \xrightarrow{d} N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \text{Var}(T_i) & \rho \sqrt{\text{Var}(T_i) \theta(1-\theta)} \\ " & \theta(1-\theta) \end{pmatrix} \right)$$

since $\rho = \frac{\text{Cov}(T_i, Y_i)}{\sqrt{\text{Var}(T_i) \text{Var}(Y_i)}} \Rightarrow \text{Cov}(T_i, Y_i) = \dots$

Then, by delta method, $\sqrt{n} (g(\bar{T}_n, \bar{Y}_n) - g(\frac{E[T_i]}{a}, \frac{\theta}{b})) \xrightarrow{d} N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \nabla g^T \Sigma \nabla g \right)$

$$\text{where } g(a, b) = a(b(1-b))^{1/2} \Rightarrow \nabla g(a, b) = \begin{pmatrix} \frac{d}{da} [a(b(1-b))^{1/2}] \\ \frac{d}{db} [a(b(1-b))^{1/2}] \end{pmatrix} = \begin{pmatrix} \frac{1}{\sqrt{b(1-b)}} \\ \frac{-a(1-2b)}{2[b(1-b)]^{3/2}} \end{pmatrix}$$

$$= \begin{pmatrix} \frac{1}{\sqrt{\theta(1-\theta)}} \\ \frac{-E[T_i](1-2\theta)}{2[\theta(1-\theta)]^{3/2}} \end{pmatrix}$$

$$\text{Then, } \sqrt{n} \left(\frac{\bar{T}}{\sqrt{\bar{T}(1-\bar{T})}} - \frac{E[T_i]}{\sqrt{\theta(1-\theta)}} \right) \xrightarrow{d} N \left(0, \begin{pmatrix} \frac{1}{\sqrt{\theta(1-\theta)}} & \frac{-E[T_i](1-2\theta)}{2[\theta(1-\theta)]^{3/2}} \end{pmatrix} \begin{pmatrix} \text{Var}(T_i) & \rho \sqrt{\text{Var}(T_i) \theta(1-\theta)} \\ \rho \sqrt{\text{Var}(T_i) \theta(1-\theta)} & \theta(1-\theta) \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{\theta(1-\theta)}} \\ \frac{-E[T_i](1-2\theta)}{2[\theta(1-\theta)]^{3/2}} \end{pmatrix} \right)$$

≈ 2
cont'd.

2 f cont'd.

where $\tilde{\sigma}^2 =$

$$\begin{aligned}
 & \left(\frac{\text{Var}(T_i)}{\theta(1-\theta)} - \frac{E[T_i](1-2\theta)\rho\sqrt{\text{Var}(T_i)}}{2[\theta(1-\theta)]^{3/2}} \right) \left(\frac{\rho\sqrt{\text{Var}(T_i)}}{\theta(1-\theta)} - \frac{E[T_i](1-2\theta)(\theta(1-\theta))}{2[\theta(1-\theta)]^{3/2}} \right) \left(\frac{1}{\theta(1-\theta)} - \frac{E[T_i](1-2\theta)}{2[\theta(1-\theta)]^{3/2}} \right) \\
 &= \frac{\text{Var}(T_i)}{\theta(1-\theta)} - \frac{E[T_i](1-2\theta)\rho\sqrt{\text{Var}(T_i)}}{2[\theta(1-\theta)]^{3/2}} - \frac{E[T_i](1-2\theta)\rho\sqrt{\text{Var}(T_i)}}{2[\theta(1-\theta)]^{3/2}} + \frac{E[T_i]^2(1-2\theta)^2}{4[\theta(1-\theta)]^2} \\
 &= \frac{4\text{Var}(T_i)[\theta(1-\theta)] + E[T_i]^2(1-2\theta)^2}{4(\theta(1-\theta))^2} - \frac{2E[T_i](1-2\theta)\rho\sqrt{\text{Var}(T_i)}}{2[\theta(1-\theta)]^{3/2}}
 \end{aligned}$$

Replace θ w/ \bar{y} , $E[T_i]$ w/ \bar{T} , $\text{Var}(T_i) = \underbrace{E[T_i^2] - E[T_i]^2}_{\approx} \text{w/ } \frac{1}{n} \sum T_i^2 - \left(\frac{1}{n} \sum T_i \right)^2 = \frac{1}{n} \sum T_i^2 - \bar{T}^2$
 and ρ w/ $\frac{\bar{T}}{\sqrt{\bar{y}(1-\bar{y})}}$, we get

$$\underbrace{\tilde{\sigma}^2}_{\text{asymptotic variance of } \hat{\rho}} = \frac{4 \left[\frac{1}{n} \sum T_i^2 - \bar{T}^2 \right] [\bar{y}(1-\bar{y})] + \bar{T}^2 (1-2\bar{y})^2}{4(\bar{y}(1-\bar{y}))^2} - \frac{2\bar{T}(1-2\bar{y}) \frac{\bar{T}}{\sqrt{\bar{y}(1-\bar{y})}} \sqrt{\frac{1}{n} \sum T_i^2 - \bar{T}^2}}{2[\bar{y}(1-\bar{y})]^{3/2}}$$

$$\text{Then, } 95\% \text{ CI}(\rho) = \left(\hat{\rho} - 1.96 \sqrt{\tilde{\sigma}^2}, \hat{\rho} + 1.96 \sqrt{\tilde{\sigma}^2} \right) = \left(\hat{\rho} - 1.96 \tilde{\sigma}, \hat{\rho} + 1.96 \tilde{\sigma} \right)$$