# A Paradoxical Result in Estimating Regression Coefficients

研究方法

112354010 陳品華

# Contents

# 1 Introduction

➢ This article presents a **counterintuitive result** regarding the **estimation of a regression slope coefficient** ($\beta_1$).

➢ The **precision of the slope estimator** can **deteriorate** when **additional information is used** to estimate its value.

   A. pooled estimate of the variance $Var(\tilde{\beta}_1)$ > not pooled estimate of the variance $Var(\hat{\beta}_1)$

   B. actual variance is known $Var(\tilde{\beta}_1^*)$ > actual variance is unknown $Var(\hat{\beta}_1)$

## 2 Simple Linear Regression Model

$$y = \beta_0 + \beta_1 x + \varepsilon \,,\, E(\varepsilon) = 0 \,,\, \text{Var}(\varepsilon) = \sigma^2 \text{ with } x \text{ independent of } \varepsilon$$

$$\text{where} \quad E(y|x) = \beta_0 + \beta_1 x \,,\, \text{Var}(y|x) = \sigma^2$$

$$E \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix} \,,\, \text{cov} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{pmatrix}$$

$$\beta_1 = \frac{\sigma_{xy}}{\sigma_x^2} \,,\, \beta_0 = \mu_y - \beta_1 \mu_x$$

# 2 Estimating the $Var(\hat{\beta}_1)$ (Least Squares Estimator)

$$\hat{\beta}_1 = \frac{s_{xy}}{s_x^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$Var(\hat{\beta}_1) = \frac{1}{n-3}\left(\frac{\sigma_y^2}{\sigma_x^2} - \beta_1^2\right)$$

$<proof>$

$$Var(\hat{\beta}_1) = Var\left(E(\hat{\beta}_1|x)\right) + E(Var(\hat{\beta}_1|x))$$

$$= Var(\beta_1) + E\left(\frac{\sigma^2}{S_x^2(n-1)}\right)$$

$$= 0 + \frac{\sigma^2}{\sigma_x^2} E\left(\frac{(n-1)S_x^2}{\sigma_x^2}\right)^{-1}$$

$$= \frac{\sigma_y^2 - \beta_1^2 \sigma_x^2}{(n-3)\sigma_x^2}$$

$$= \frac{1}{n-3}\left(\frac{\sigma_y^2}{\sigma_x^2} - \beta_1^2\right)$$

# 2 Estimating the $Var(\tilde{\beta}_1^*)$ ($\sigma_x^2$ is known)

$$\tilde{\beta}_1^* = \frac{S_{xy}}{\sigma_x^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{(n-1)\sigma_x^2}$$

$$Var(\tilde{\beta}_1^*) = \frac{1}{n-1}\left(\frac{\sigma_y^2}{\sigma_x^2} + \beta_1^2\right)$$

$$<proof> \quad Var(\tilde{\beta}_1^*) = E\left(Var(\tilde{\beta}_1^*|x)\right) + Var(E(\tilde{\beta}_1^*|x))$$

$$= E\left(\frac{\sigma^2 S_x^2}{S_x^4(n-1)}\right) + Var\left(\frac{S_x^2}{\sigma_x^2}\beta_1\right)$$

$$= \frac{\sigma^2}{\sigma_x^2(n-1)} + \frac{2\beta_1^2}{n-1}$$

$$= \frac{1}{n-1}\left(\frac{\sigma_y^2}{\sigma_x^2} + \beta_1^2\right)$$

# 2  Estimating the $Var(\tilde{\beta}_1)$ (2-Samples pooled estimator)

| Treatment 1 | Treatment 2 |
|:---:|:---:|
| $x_{11}$ <br> $y_{11}$ | $x_{12}$ <br> $y_{12}$ |
| $x_{21}$ <br> $y_{21}$ | $x_{22}$ <br> $y_{22}$ |
| . <br> . <br> . | . <br> . <br> . |
| $x_{n_1 1}$ <br> $y_{n_1 1}$ | $x_{n_2 2}$ <br> $y_{n_2 2}$ |

$x_{c1}$  $\sigma_x^2$          $x_{c2}$  $\sigma_x^2$

$y_1$  $\sigma_{x1}^2$          $y_2$  $\sigma_{x2}^2$

$n_1$  $\sigma_{y1}^2$          $n_2$  $\sigma_{y2}^2$

$$\tilde{\beta}_1 = \frac{n_1 + n_2 - 2}{n_1 - 1} \left(\sum_{j=1}^{2} x'_{cj} \, x_{cj}\right)^{-1} x'_{c1} y_1$$

$$Var(\tilde{\beta}_1) \geq \frac{\sigma_1^2}{\sigma_x^2(n_1 - 1)} \left(\frac{n_1 + n_2 - 2}{n_1 + n_2}\right)^2 + \frac{2\beta_1^2}{n_1 - 1}\left(\frac{n_2 - 1}{n_1 + n_2}\right)$$

# 3 Comparison

| Least Squares Estimator | $\sigma_x^2$ is known | 2-samples pooled estimator |
|---|---|---|
| $Var(\hat{\beta}_1) = \dfrac{1}{n-3}\left(\dfrac{\sigma_y^2}{\sigma_x^2} - \beta_1^2\right)$ | $Var(\tilde{\beta}_1^*) = \dfrac{1}{n-1}\left(\dfrac{\sigma_y^2}{\sigma_x^2} + \beta_1^2\right)$ | $Var(\tilde{\beta}_1) \gtrsim \dfrac{1}{n_1-1}\left(\dfrac{\sigma_{y1}^2}{\sigma_x^2} + \dfrac{n_2-n_1}{n_1+n_2}\beta_1^2\right)$ |

# 3 Comparison between $\hat{\beta}_1$(LSE) and $\tilde{\beta}_1^*$ ($\sigma_x^2$ is known)

In an extreme special case : $y = \beta_0 + \beta_1 x$

| $\hat{\beta}_1 = \beta_1$ | $\tilde{\beta}_1^* = \beta_1 \dfrac{S_x^2}{\sigma_x^2}$ |
|---|---|
| perfect estimation | not a perfect estimator whenever $S_x^2 \neq \sigma_x^2$ |

# 3 Comparison between $Var(\hat{\beta}_1)$ and $Var(\tilde{\beta}_1)$

$$\boxed{Var(\hat{\beta}_1) < Var(\tilde{\beta}_1^*)}$$

Another insight : Cauchy - Schwartz inequality : $\sigma_{xy}^2 \leq \sigma_x^2 \sigma_y^2$

$$\beta_1^2 = \frac{\sigma_{xy}^2}{\sigma_x^4} \leq \frac{\sigma_y^2}{\sigma_x^2}$$

| | |
|---|---|
| $Var(\hat{\beta}_1)$ approaches to 0 | $Var(\tilde{\beta}_1^*)$ approaches to $\frac{2\sigma_y^2}{\sigma_x^2(n-1)}$ |

$$\boxed{Var(\hat{\beta}_1) < Var(\tilde{\beta}_1^*)}$$

$$\hat{\beta}_1 = \frac{s_{xy}}{s_x^2} = \frac{s_{xy}}{\sigma_x^2}\left(\frac{\sigma_x^2}{S_x^2}\right) = \tilde{\beta}_1^* w \quad, where \; w = \frac{\sigma_x^2}{S_x^2}$$

when $\tilde{\beta}_1^*$ overestimates $\beta_1$, the ratio $w$ will tend to pull $\hat{\beta}_1$ down toward the true slope $\beta_1$

# 3 Comparison between $Var(\tilde{\beta}_1^*)$ and $Var(\tilde{\beta}_1)$

| | |
|---|---|
| $Var(\tilde{\beta}_1^*) = \dfrac{1}{n-1}\left(\dfrac{\sigma_y^2}{\sigma_x^2} + \beta_1^2\right)$ | $Var(\tilde{\beta}_1) \gtrsim \dfrac{1}{n_1-1}\left(\dfrac{\sigma_{y1}^2}{\sigma_x^2} + \dfrac{n_2-n_1}{n_1+n_2}\beta_1^2\right)$ |

| pooled estimator $\tilde{\beta}_1$ | Relationship between $Var(\tilde{\beta}_1^*)$ and $Var(\tilde{\beta}_1)$ |
|---|---|
| $n_1$ is fixed<br>$n_2 \longrightarrow \infty$ | $Var(\tilde{\beta}_1)$ would converge to $Var(\tilde{\beta}_1^*)$ |
| finite $n_1$ and $n_2$ | $Var(\tilde{\beta}_1) < Var(\tilde{\beta}_1^*)$ |

# 3 Comparison between $Var(\hat{\beta}_1)$ and $Var(\tilde{\beta}_1)$

| | |
|---|---|
| $Var(\hat{\beta}_1) = \dfrac{1}{n-3}(\dfrac{\sigma_y^2}{\sigma_x^2} - \beta_1^2)$ | $Var(\tilde{\beta}_1) \gtrsim \dfrac{1}{n_1-1}(\dfrac{\sigma_{y1}^2}{\sigma_x^2} + \dfrac{n_2-n_1}{n_1+n_2}\beta_1^2)$ |

| pooled estimator $\tilde{\beta}_1$ | Relationship between $Var(\hat{\beta}_1)$ and $Var(\tilde{\beta}_1)$ |
|---|---|
| $n_1 \gg n_2$ | $Var(\hat{\beta}_1) \approx Var(\tilde{\beta}_1)$ |

# Simulation

**population** :

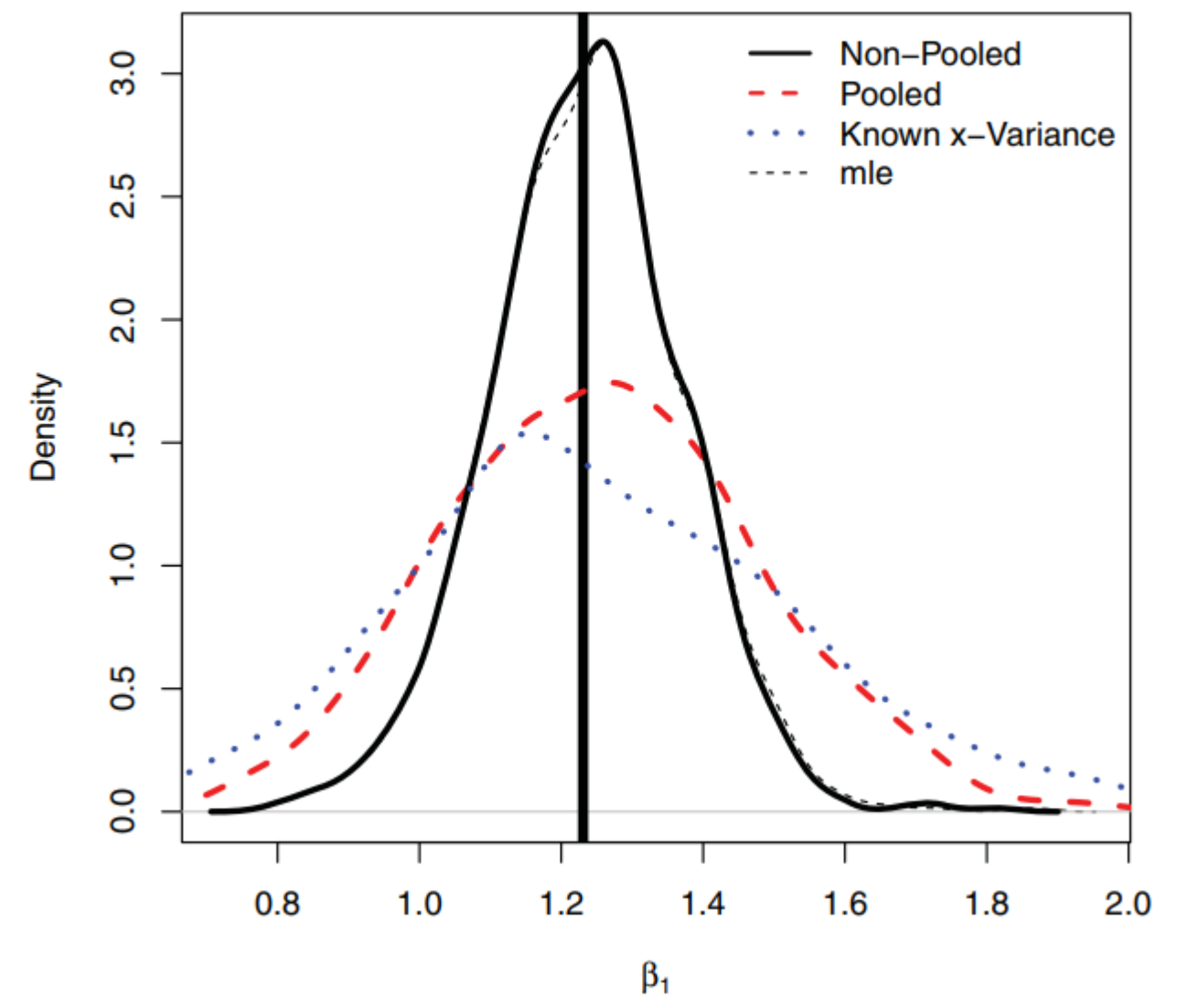group $1 : (x, y) \sim BN(1, 2, 1.3^2, 2^2, 0.8)$

group $2 : (x, y) \sim BN(1, 2.5, 1.3^2, 2.5^2, 0.7)$

**sample** :

$n_1 = 50$

$n_2 = 50$
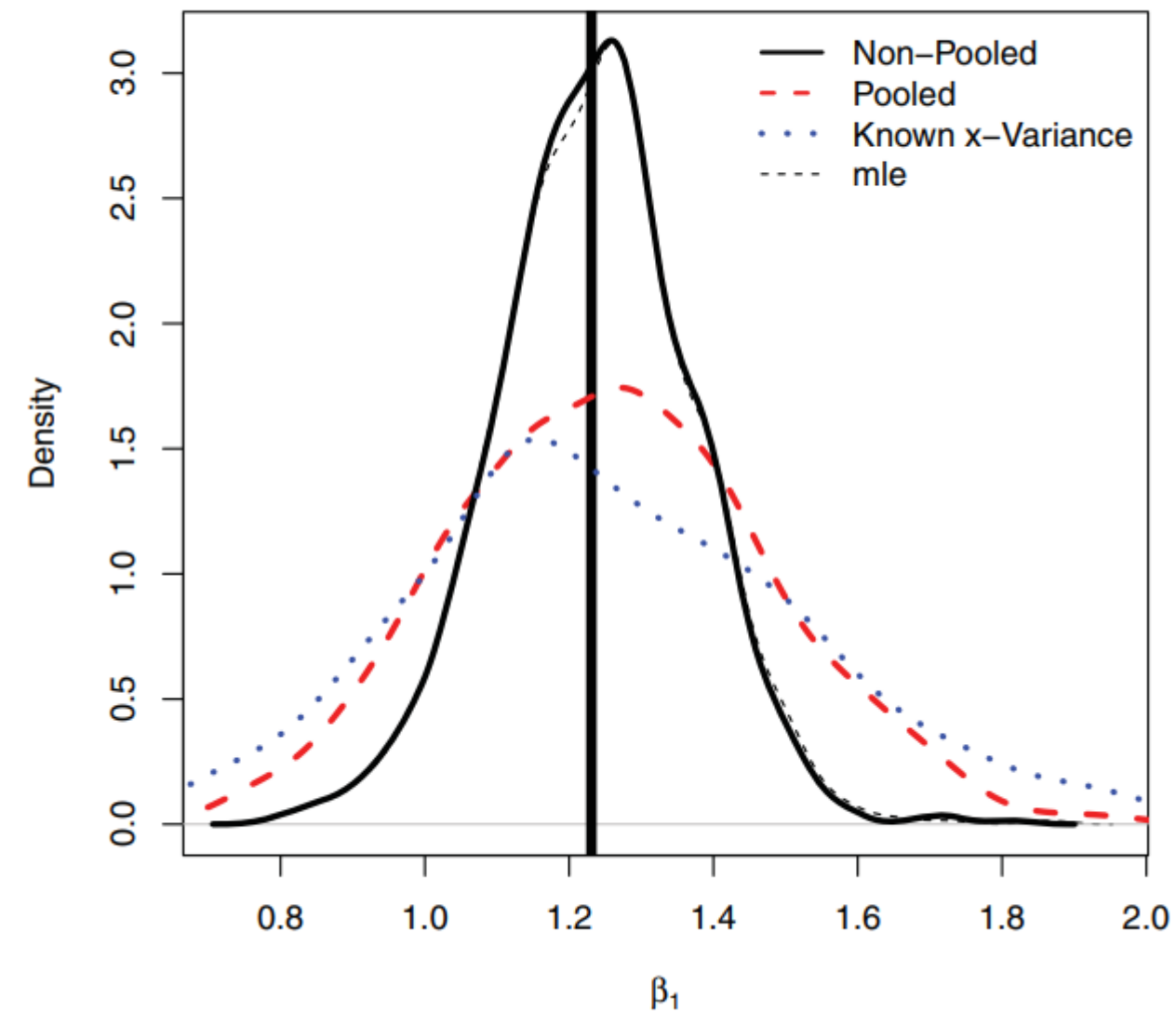
**Slope Estimates: Two Independent Samples**

# 4 Simulation

Group 1 ($\beta_1 = 1.23$)

1. $\sigma_{LSE} = 0.1372$

2. $\sigma_{MLE} = 0.1395$

3. $\sigma_{pooled} = 0.2197$

4. $\sigma_{\sigma_x^2 \text{ is known}} = 0.2851$

**Slope Estimates: Two Independent Samples**

# 5  Conclusion

➢ Using a **known variance** to estimate $\beta_1$ may lead to **an increase in the variance** of the estimator, which is **contrary to conventional statistical wisdom.**

➢ The authors point out that this paradoxical result has important implications for interpreting the results of randomized experiments and propose a new method for better estimating the relationship between predictor variables and outcomes.

# Thanks!