

Yao Thibaut Kpegli

Licence 3 in Economics, ENS Paris-Saclay
2023–2024

April 2024

Outline

- 1 MLE
- 2 Binary outcome
- 3 PSM
- 4 Application

Intuition

Context

- 1 You have a sample $\mathbf{y} = \{y_1, \dots, y_n\}$
- 2 You know that the sample comes from a random variable Y with a vector of parameters $\boldsymbol{\theta} \in R^K$ whose true value is $\boldsymbol{\theta}_0$
- 3 You don't know the true value $\boldsymbol{\theta}_0$

Objective

Provide an estimate of $\boldsymbol{\theta}_0$

Intuition

$\hat{\boldsymbol{\theta}}_{MLE}$ = the value of $\boldsymbol{\theta}$ that is such that the probability of having observed \mathbf{y} is the highest possible.

Intuition

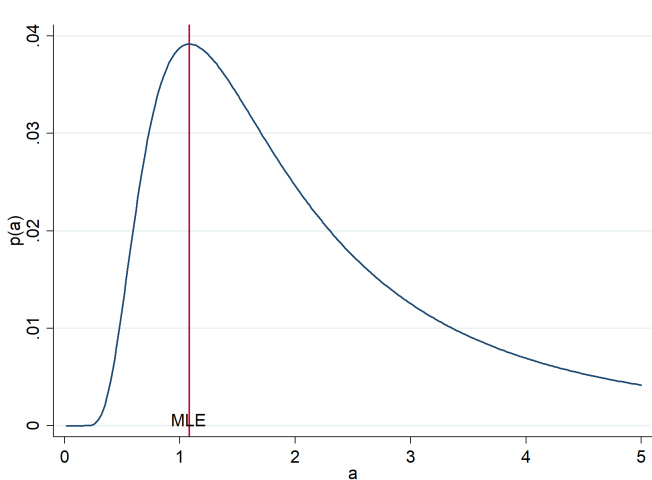
Example (a continuous case)

- ➊ Assume that lifetime of an electronic equipment is a *r.v* following an exponential distribution with parameter $a > 0$.
- ➋ The density of exponential distribution with parameter $a > 0$ is $f(y; a) = \frac{1}{a} \exp\left(-\frac{y}{a}\right)$
- ➌ We have observed randomly the lifetime 3 times, thereby constituting a sample $y_1 = 1, y_2 = 0.5$ and $y_3 = 1.75$.
- ➍ You want to estimate a , that is the vector of parameters is simply $\theta = a$
- ➎ The joint density of having observed $\{y_1, y_2, y_3\}$ is

$$p(a) = \frac{1}{a} \exp\left(-\frac{1}{a}\right) \times \frac{1}{a} \exp\left(-\frac{0.5}{a}\right) \times \frac{1}{a} \exp\left(-\frac{1.75}{a}\right) = \frac{1}{a^3} \exp\left(-\frac{3.25}{a}\right)$$

Intuition

Figure: probability density of observing the sample as function of a



Intuition

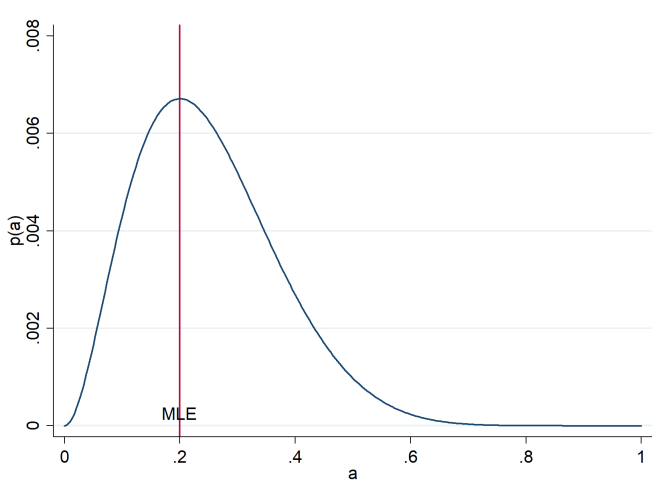
Example (a discrete case)

- ➊ Assume that the fact that a resident of a town has a specific disease is a *r.v* following a Bernoulli distribution with parameter $a \in (0, 1)$, i.e. the probability of having the disease is a
- ➋ We have randomly 10 people, thereby constituting a sample $\{y_1, y_2, \dots, y_{10}\} = \{0, 1, 0, 0, 1, 0, 0, 0, 0, 0\}$.
- ➌ You want to estimate a , that is the vector of parameters is simply $\theta = a$
- ➍ The joint probability of observing $\{y_1, \dots, y_{10}\}$ is

$$p(a) = (1-a)a(1-a)(1-a)a(1-a)(1-a)(1-a)(1-a)(1-a) = a^2(1-a)^8$$

Intuition

Figure: probability of observing the sample as function of a



Notations and Definitions

Definition (Likelihood function)

The likelihood function \mathcal{L} is:

$$\begin{array}{ccc} R^K & \longrightarrow & [0,1] \\ \theta & \mapsto & \mathcal{L}(\boldsymbol{\theta}; \mathbf{y}) = g(\mathbf{y}; \boldsymbol{\theta}) \end{array}$$

Definition (Log-likelihood function)

The log-likelihood function is:

$$\log \mathcal{L}(\boldsymbol{\theta}; \mathbf{y}) = \sum_{i=1}^n \log f(y_i; \boldsymbol{\theta})$$

Notations and Definitions

Example

- ① if $Y_i \rightsquigarrow \mathcal{B}(p)$, then $\boldsymbol{\theta} = p$ and

$$\log \mathcal{L}(\boldsymbol{\theta}; \mathbf{y}) = \left(\sum_{i=1}^n y_i \right) \log(p) + \left(n - \sum_{i=1}^n y_i \right) \log(1 - p)$$

- ② if $Y_i \rightsquigarrow \mathcal{N}(m, \sigma^2)$, then $\boldsymbol{\theta} = [m, \sigma^2]'$ and

$$\log \mathcal{L}(\boldsymbol{\theta}; \mathbf{y}) = -\frac{1}{2} \sum_{i=1}^n \left(\log \sigma^2 + \log(2\pi) + \frac{(y_i - m)^2}{\sigma^2} \right)$$

Notations and Definitions

Definition (Score function)

The score function is:

$$S(y; \boldsymbol{\theta}) = \frac{\partial \log f(y, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

Example

- ① For $Y \rightsquigarrow \mathcal{B}(p)$

$$S(y, \boldsymbol{\theta}) = \frac{y}{p} - \frac{1-y}{1-p}$$

- ② For $Y \rightsquigarrow \mathcal{N}(m, \sigma^2)$

$$S(y, \boldsymbol{\theta}) = \begin{bmatrix} \frac{y-m}{\sigma^2} \\ \frac{1}{2\sigma^2} \left(\left(\frac{y-m}{\sigma} \right)^2 - 1 \right) \end{bmatrix}$$

Notations and Definitions

Proposition 1 (Score expectation)

The expectation of the score is zero

Proof:

$$E[S(Y, \theta)] = \int \frac{\partial \log f(y; \theta)}{\partial \theta} f(y; \theta) dy = \int \frac{\partial f(y; \theta)}{\partial \theta} \frac{1}{f(y; \theta)} f(y; \theta) dy =$$

$$\frac{\partial \int f(y; \theta) dy}{\partial \theta} = \frac{\partial(1)}{\partial \theta} = 0$$

Definition (Information matrix)

$$\mathcal{I}_Y(\theta) = -E\left(\frac{\partial^2 \log f(Y, \theta)}{\partial \theta \partial \theta'}\right)$$

Notations and Definitions

Remark (additivity of the information matrix)

The information matrix of two independent experiments is:

$$\mathcal{I}_{X,Y}(\boldsymbol{\theta}) = \mathcal{I}_X(\boldsymbol{\theta}) + \mathcal{I}_Y(\boldsymbol{\theta})$$

Proposition 2 (Variance of the score)

The variance of the score is equal to the information matrix

$$V\left(S(Y; \boldsymbol{\theta})\right) \equiv E\left(\left(\frac{\partial \log f(Y, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right) \left(\frac{\partial \log f(Y, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right)'\right) = \mathcal{I}_Y(\boldsymbol{\theta})$$

Proof:

Note that $\frac{\partial^2 \log f(y, \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} = \frac{\partial^2 f(y; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \frac{1}{f(y; \boldsymbol{\theta})} - \frac{\partial \log f(y, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \frac{\partial \log f(y; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}'}$

Also $E\left[\frac{\partial^2 f(y; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \frac{1}{f(y; \boldsymbol{\theta})}\right] = \frac{\partial^2(1)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} = \mathbf{0}$. This leads to the result.

Notations and Definitions

Example

- ① For $Y \rightsquigarrow \mathcal{B}(p)$

$$\mathcal{I}_Y(\boldsymbol{\theta}) = E \left[\frac{Y}{p^2} + \frac{1-Y}{(1-p)^2} \right] = \frac{1}{p(1-p)}$$

- ② For $Y \rightsquigarrow \mathcal{N}(m, \sigma^2)$

$$\mathcal{I}_Y(\boldsymbol{\theta}) = E \begin{bmatrix} \frac{1}{\sigma^2} & \frac{y-m}{\sigma^4} \\ \frac{y-m}{\sigma^4} & \frac{(y-m)^2}{\sigma^6} - \frac{1}{2\sigma^4} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{2\sigma^4} \end{bmatrix}$$

Notations and Definitions

Theorem 1 (Cauchy-Schwarz inequality)

Let X and Y be two random variables. Then:

$$|Cov(X, Y)| \leq \sqrt{V(X)V(Y)}$$

Proof: Let $Z = Y - \frac{Cov(X, Y)}{V(X)}X$. Then,

$Cov(X, Z) = Cov(X, Y) - Cov(X, X) \frac{Cov(X, Y)}{V(X)} = 0$. Then,

$$V(Y) = V\left(Z + \frac{Cov(X, Y)}{V(X)}X\right) = V(Z) + \left(\frac{Cov(X, Y)}{V(X)}\right)^2 V(X) \geq \frac{(Cov(X, Y))^2}{V(X)}.$$

Hence, $V(Y)V(X) \geq (Cov(X, Y))^2$ which leads to the result.

Notations and Definitions

Theorem 2 (Fréchet-Darmois-Cramér-Rao bound)

Consider an unbiased estimator $\hat{\theta}(\mathbf{Y})$ of θ . The variance of the estimator $\hat{\theta}(\mathbf{Y})$ has a lower bound:

$$V\left(\hat{\theta}(\mathbf{Y})\right) \geq \mathcal{I}_Y(\theta)^{-1} \equiv B_F(\theta)$$

Proof: First, $\text{Cov}\left(\hat{\theta}(\mathbf{Y}), S(Y; \theta)\right) = E\left[(\hat{\theta}(\mathbf{Y}) - \theta)S(Y; \theta)\right] =$
 $E\left[\hat{\theta}(\mathbf{Y})S(Y; \theta)\right] = \int \hat{\theta}(\mathbf{y})S(Y; \theta)f(y; \theta)dy = \int \hat{\theta}(\mathbf{y})\frac{\partial f(y; \theta)}{\partial \theta}\frac{1}{f(y; \theta)}f(y; \theta)dy =$
 $\frac{\partial \int \hat{\theta}(\mathbf{y})f(y; \theta)dy}{\partial \theta} = \frac{\partial E\left(\hat{\theta}(\mathbf{Y})\right)}{\partial \theta} = \frac{\partial \theta}{\partial \theta} = \mathbf{1}$. Hence, by Cauchy-Schwarz inequality
 we have $V\left(\hat{\theta}(\mathbf{Y})\right)V\left(S(Y; \theta)\right) \geq \mathbf{1}$. The result follows as $V\left(S(Y; \theta)\right) = \mathcal{I}_Y(\theta)$.

Maximum Likelihood Estimation (MLE)

Definition (Identification)

The vector of parameters θ is identifiable if, for any other vector θ^* :

$$\theta^* \neq \theta \implies \log \mathcal{L}(\theta^*; \mathbf{y}) \neq \log \mathcal{L}(\theta; \mathbf{y})$$

Definition (Likelihood equation)

Necessary condition for maximizing the likelihood function:

$$\frac{\partial \log \mathcal{L}(\theta; \mathbf{y})}{\partial \theta} = \sum_{i=1}^n S(y_i; \theta) = \mathbf{0}$$

Definition (Maximum Likelihood Estimator (MLE))

The maximum likelihood estimator θ_{MLE} is the vector θ that maximizes the likelihood function. Formally:

$$\theta_{MLE} = \arg \max_{\theta} \log \mathcal{L}(\theta; \mathbf{y}) = \arg \max_{\theta} \mathcal{L}(\theta; \mathbf{y})$$

Maximum Likelihood Estimator (MLE)

Example

- ① For $Y \rightsquigarrow \mathcal{B}(p)$

$$\theta_{MLE} = \frac{1}{n} \sum_{i=1}^n y_i$$

- ② For $Y \rightsquigarrow \mathcal{N}(m, \sigma^2)$

$$\theta_{MLE} = \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n y_i \\ \frac{1}{n} \sum_{i=1}^n (y_i - m_{MLE})^2 \end{bmatrix}$$

Regularity conditions

Regularity conditions

- ① The support of Y does not depends on θ
- ② θ_0 is identified
- ③ The log-likelihood function is continuous in θ
- ④ $E(\log f(Y; \theta_0))$ exists
- ⑤ The log-likelihood function is twice continuously differentiable
- ⑥ The information matrix at θ_0 $\mathcal{I}_Y(\theta_0) = -E\left(\frac{\partial^2 \log f(Y, \theta_0)}{\partial \theta \partial \theta'}\right)$ exists and is nonsingular

Proposition 3 (Properties of MLE)

Under regularity conditions, the MLE is (i) **consitent**, (ii) **asymptotically normally distributed** and (iii) **asymptotically efficient**.

Consistency

Kullback-Liebler divergence

If $f_{\theta_0}(y)$ and $f_{\theta_1}(y)$ are two densities, the Kullback-Leibler divergence of f_{θ_1} w.r.t f_{θ_0} is

$$KL(f_{\theta_1} \| f_{\theta_0}) = E_{\theta_0} \left[\log \frac{f(Y, \theta_0)}{f(Y, \theta_1)} \right] = \int f(y, \theta_0) \log \frac{f(y, \theta_1)}{f(y, \theta_0)} dy$$

Proposition 3-0

- ① $KL(f_{\theta_1} \| f_{\theta_0}) \geq 0$
- ② $KL(f_{\theta_1} \| f_{\theta_0}) = 0$ iff $f_{\theta_0} = f_{\theta_1}$

Proof: First, $-\log(x)$ is a convex function. By Jensen's inequality,

$$KL(f_{\theta_1} \| f_{\theta_0}) = E_{\theta_0} \left[-\log \frac{f(Y, \theta_1)}{f(Y, \theta_0)} \right] \geq -\log E_{\theta_0} \left[\frac{f(Y, \theta_1)}{f(Y, \theta_0)} \right] =$$

$$-\log \int f(y, \theta_1) dy = 0.$$

Second, because $-\log(x)$ is strictly convex, equality holds if and only if

$f(y, \theta_1)/f(y, \theta_0)$ is constant.

Consistency

Proposition 3-1 (Consistency of MLE)

Under regularity conditions, θ_{MLE} converge in probability to the true value θ_0 :

$$\text{plim } \theta_{MLE} = \theta_0$$

Informal argument:

$$\begin{aligned} \theta_{MLE} &= \arg \max_{\theta} \frac{1}{n} \sum_{i=1}^n \log f(Y_i; \theta) \\ &= \arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n \log f(Y_i; \theta) \\ &= \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \log f(Y_i; \theta_0) - \frac{1}{n} \sum_{i=1}^n \log f(Y_i; \theta) \\ &\underset{n \rightarrow +\infty}{\simeq} \arg \min_{\theta} E_{\theta_0} \log f(Y; \theta_0) - E_{\theta_0} \log f(Y; \theta) \\ &\underset{n \rightarrow +\infty}{\simeq} \arg \min_{\theta} KL(f_{\theta_0} \| f_{\theta}) = \theta_0 \\ &\underset{n \rightarrow +\infty}{\simeq} \theta_0 \end{aligned}$$

with $\underset{n \rightarrow +\infty}{\simeq}$ stands for Law of Large Numbers

Consistency

Example

- ① For $Y \rightsquigarrow \mathcal{B}(p)$

$$\theta_{MLE} = \frac{1}{n} \sum_{i=1}^n Y_i \xrightarrow{p} E[Y] = p$$

- ② For $Y \rightsquigarrow \mathcal{N}(m, \sigma^2)$

$$\theta_{MLE} = \begin{bmatrix} m_{MLE} \\ \sigma_{MLE}^2 \end{bmatrix}$$

$$m_{MLE} = \frac{1}{n} \sum_{i=1}^n Y_i \xrightarrow{p} E[Y] = m$$

$$\sigma_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n Y_i^2 - \left(\frac{1}{n} \sum_{i=1}^n Y_i \right)^2 \xrightarrow{p} E[Y^2] - E[Y]^2 = \sigma^2$$

Consistency

Remark

- ① Regularity conditions are sufficient but not necessary conditions to have consistency
- It is therefore possible for the MLE to be consistent even in situations that do not meet regularity conditions
- e.g: if $Y \rightsquigarrow \mathcal{U}[0, \theta]$, we have $\boldsymbol{\theta}_{MLE} = \max\{Y_1, \dots, Y_n\}$

$$\lim_{n \rightarrow +\infty} P(|\boldsymbol{\theta}_{MLE} - \boldsymbol{\theta}| > \epsilon) = \lim_{n \rightarrow +\infty} \left(1 - \frac{\epsilon}{\theta}\right)^n = 0$$

Asymptotic normality

Proposition 2 (Asymptotic normality)

The MLE estimator θ_{MLE} is normally distributed asymptotically:

$$\sqrt{n}(\theta_{MLE} - \theta_0) \xrightarrow{dist.} \mathcal{N}(0, \mathcal{I}_Y(\theta_0)^{-1})$$

Proof:

- First order Taylor expansion of the first derivative of

$$\log \mathcal{L}(\theta; \mathbf{y})$$

$$\frac{\partial \log \mathcal{L}(\theta; \mathbf{y})}{\partial \theta} \simeq \frac{\partial \log \mathcal{L}(\theta_0; \mathbf{y})}{\partial \theta} + \frac{\partial^2 \log \mathcal{L}(\theta_0; \mathbf{y})}{\partial \theta \partial \theta'} (\theta - \theta_0)$$

- Evaluate at $\theta = \theta_{MLE}$

$$0 \simeq \frac{\partial \log \mathcal{L}(\theta_0; \mathbf{y})}{\partial \theta} + \frac{\partial^2 \log \mathcal{L}(\theta_0; \mathbf{y})}{\partial \theta \partial \theta'} (\theta_{MLE} - \theta_0)$$

Asymptotic normality

- Rearrange

$$\sqrt{n}(\boldsymbol{\theta}_{MLE} - \boldsymbol{\theta}_0) \simeq \sqrt{n} \frac{\frac{\partial \log \mathcal{L}(\boldsymbol{\theta}_0; \mathbf{y})}{\partial \boldsymbol{\theta}}}{-\frac{\partial^2 \log \mathcal{L}(\boldsymbol{\theta}_0; \mathbf{y})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}} = \sqrt{n} \frac{\frac{1}{n} \sum_{i=1}^n \frac{\partial \log f(y_i; \boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}}}{-\frac{1}{n} \sum_{i=1}^n \frac{\partial^2 \log f(y_i; \boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}}$$

- Central Limit Theorem

$$A_n = \frac{1}{n} \sum_{i=1}^n \frac{\partial \log f(y_i; \boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \xrightarrow{\text{dist.}} \mathcal{N}\left(E[S(Y, \boldsymbol{\theta}_0)], \frac{V(S(Y, \boldsymbol{\theta}_0))}{n}\right)$$

- Law of Large Number

$$B_n = -\frac{1}{n} \sum_{i=1}^n \frac{\partial^2 \log f(y_i; \boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \xrightarrow{p} -E\left[\frac{\partial^2 \log f(Y; \boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}\right] = \mathcal{I}_Y(\boldsymbol{\theta}_0)$$

- Remembering that $E[S(Y, \boldsymbol{\theta}_0)] = 0$ and $V(S(Y, \boldsymbol{\theta}_0)) = \mathcal{I}_Y(\boldsymbol{\theta}_0)$, Slutsky Theorem leads to

$$\sqrt{n}(\boldsymbol{\theta}_{MLE} - \boldsymbol{\theta}_0) = \sqrt{n} \frac{A_n}{B_n} \xrightarrow{\mathcal{L}} \mathcal{N}\left(0, \mathcal{I}_Y(\boldsymbol{\theta}_0)^{-1}\right)$$

Asymptotically efficient

Proposition (Asymptotically efficient)

θ_{MLE} is asymptotically efficient, *i.e.* achieves the FDCR lower bound for consistent estimators.

Proof: We have

$$\sqrt{n}(\theta_{MLE} - \theta_0) \xrightarrow{dist.} \mathcal{N}(0, \mathcal{I}_Y(\theta_0)^{-1})$$

This means that

$$V(\theta_{MLE}) = I(\theta_0)^{-1}$$

with $I(\theta_0) = n\mathcal{I}_Y(\theta_0)$ the information matrix associated to $\{Y_1, \dots, Y_n\}$.

Remark:

- The asymptotic variance-covariance matrix $I(\theta_0)$ of the MLE depends on the unknown value of θ_0
- In practice, the matrix is evaluated at θ_{MLE}

Sum up

Sum up

- ① You have a sample $\mathbf{y} = (y_1, \dots, y_n)$
- ② You know that the sample comes from a random variable Y with a vector of parameters $\boldsymbol{\theta} \in R^K$ whose true value $\boldsymbol{\theta}_0$ is unknown
- ③ The log-likelihood of one observation y_i is computed analytically (as a function of $\boldsymbol{\theta}$): $l_i(\boldsymbol{\theta}; y_i) = \log f(y_i; \boldsymbol{\theta})$
- ④ The log-likelihood of the sample is $\log \mathcal{L}(\boldsymbol{\theta}; \mathbf{y}) = \sum_i l_i(\boldsymbol{\theta}; y_i)$
- ⑤ The The MLE estimator results from the optimization problem

$$\boldsymbol{\theta}_{MLE} = \arg \max_{\boldsymbol{\theta}} \log \mathcal{L}(\boldsymbol{\theta}; \mathbf{y})$$

- ⑥ We have: $\boldsymbol{\theta}_{MLE} \rightsquigarrow \mathcal{N}(0, I(\boldsymbol{\theta}_0)^{-1})$, where $I(\boldsymbol{\theta}_0)$ is estimated as $I(\boldsymbol{\theta}_{MLE})$.

Outline

- 1 MLE
- 2 Binary outcome
- 3 PSM
- 4 Application

Context and Objectives

Context

- 1 You have observed a sample $\mathbf{y} = \{y_1, \dots, y_n\}$ and y_i has only two possible values (say 0 and 1)

$$y_i = \begin{cases} 1 & \text{with probability } p_i \\ 0 & \text{with probability } 1 - p_i \end{cases} \quad (1)$$

- 2 You have also observed a vector of characteristics \mathbf{x}_i ($K, 1$) associated to the individual i
- 3 You suspect that \mathbf{x}_i determines p_i :
 - and you assume this is through a function g depending on vectors of parameters $\boldsymbol{\theta}$ ($K, 1$)

$$p_i = g(\boldsymbol{\theta}' \mathbf{x}_i)$$

- $E[y_i | \mathbf{x}_i] = p_i = g(\boldsymbol{\theta}' \mathbf{x}_i)$
- $V(y_i | \mathbf{x}_i) = p_i(1 - p_i) = g(\boldsymbol{\theta}' \mathbf{x}_i)(1 - g(\boldsymbol{\theta}' \mathbf{x}_i))$

Context and Objectives

Objectives

Provide an estimate of θ

Example

Examples of binary context:

- treated or untreated

Context and Objectives

Objectives

Provide an estimate of θ

Example

Examples of binary context:

- treated or untreated
- buy or not a transportation ticket
- declares to tax administration the right level of income or not
- living in the city or in the countryside
- wear a mask or not / has covid or not
- trust or not in trust interaction
- bet or not
- success/failure (exams)
- ...

Linear Probability Model (LPM)

Linear Probability Model (LPM)

- 1 OLS regression of y on x

$$y_i = \theta' x_i + \epsilon_i$$

- In LPM, we then have $g(\theta' x_i) = \theta' x_i$
- 2 Under the the assumptions of **conditional-mean-zero** and **non-correlated errors**, such a regression could be consistent
- 3 But, at least three problems
 - **heteroskedasticity**: $V(\epsilon_i | x_i) = V(y_i | x_i) = \theta' x_i (1 - \theta' x_i)$
 - **discrete error**: $\epsilon_i | x_i = (1 - \theta' x_i, -\theta' x_i; \theta' x_i, 1 - \theta' x_i)$, so error cannot be normal
 - **unrestricted probability**: estimated probability $\hat{p}_i = \hat{\theta}' x_i$ may be outside the range $[0, 1]$

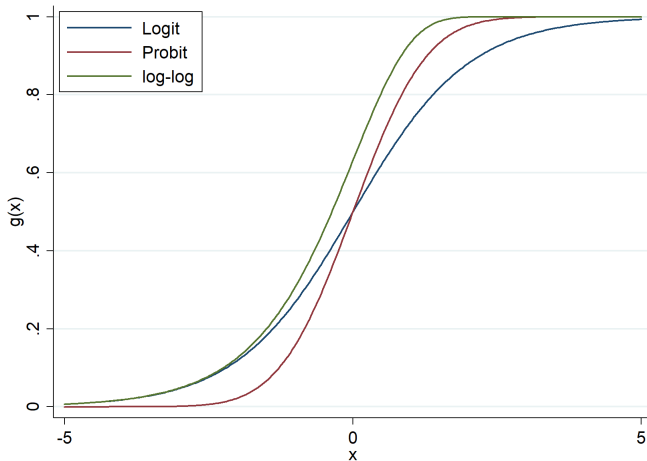
Alternative models

Table: Four common specifications of $g(\cdot)$

Model	Function $g(z)$	Derivative
LPM	z	1
Logit	$\frac{\exp(z)}{1 + \exp(z)}$	$\frac{\exp(z)}{(1 + \exp(z))^2}$
Probit	$\int_{-\infty}^z \phi(t) dt$	$\phi(z)$
log-log	$1 - \exp(-\exp(z))$	$\exp(-\exp(z)) \exp(z)$
$\phi(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{z^2}{2})$ is the density function of $\mathcal{N}(0, 1)$		

Alternative models

Figure: Probit, Logit and log-log functions



Interpretation in terms of latent variable

Interpretation in terms of latent variable

- ① A continuous but unobservable variable y_i^* :

$$y_i^* = \boldsymbol{\theta}'\mathbf{x}_i + \epsilon_i$$

with ϵ_i following normal or logistic distribution

$$y_i = \begin{cases} 1 & \text{if } y_i^* > 0 \\ 0 & \text{if } y_i^* \leq 0 \end{cases} \quad (2)$$

- ② $p_i = P(y_i^* > 0) = P(\epsilon_i > -\boldsymbol{\theta}'\mathbf{x}_i) = 1 - g(-\boldsymbol{\theta}'\mathbf{x}_i) = g(\boldsymbol{\theta}'\mathbf{x}_i)$

Relationship with Random Utility Models

Relationship with Random Utility Models

- ① What precedes relates to Random Utility Models (RUM)
- ② Agent (i) chooses $y_i = 1$ if the utility associated with this choice ($U_{i,1}$) is greater than the one of $y_i = 0$ ($U_{i,0}$)
- ③ The random utility:

$$U_{i,j} = V_{i,j} + \epsilon_{i,j}$$

- where $V_{i,j}$ is the deterministic component of the utility associated with choice $j \in \{0, 1\}$ and $\epsilon_{i,j}$ is a random (agent-specific) component.
- ④ considering that $g(\cdot)$ is the c.d.f of $\epsilon_{i,0} - \epsilon_{i,1}$, then:

$$p_i = P(V_{i,1} + \epsilon_{i,1} > V_{i,0} + \epsilon_{i,0}) = g(V_{i,1} - V_{i,0})$$

- ⑤ In the simple case, $V_{i,j} = \theta'_j x_i$, we have

$$p_i = g(\theta' x_i) \quad \text{with} \quad \theta' = \theta'_1 - \theta'_0$$

Estimation

Estimation

- ➊ These models can be estimated by Maximum Likelihood approaches (see previous seance).
- ➋ (y_i, \mathbf{x}_i) are assumed to be independent across entities i
- ➌ We have y_i that follows a Bernoulli distribution:

$$P(Y_i = y_i) = p_i^{y_i} (1 - p_i)^{1-y_i} \quad \text{with} \quad p_i = g(\boldsymbol{\theta}' \mathbf{x}_i) \quad (3)$$

- ➍ Log-likelihood of entity i :

$$l_i(\boldsymbol{\theta}; \mathbf{y}, \mathbf{x}) = y_i \log(g(\boldsymbol{\theta}' \mathbf{x}_i)) + (1 - y_i) \log(1 - g(\boldsymbol{\theta}' \mathbf{x}_i)) \quad (4)$$

- ➎ Log-likelihood of the sample:

$$\mathcal{L}(\boldsymbol{\theta}; \mathbf{y}, \mathbf{x}) = \sum_{i=1}^N l_i(\boldsymbol{\theta}; \mathbf{y}, \mathbf{x}) \quad (5)$$

Estimation

Estimation

- 1 F.O.C of the optimization program

$$\frac{\partial \log \mathcal{L}(\boldsymbol{\theta}; \mathbf{y})}{\partial \boldsymbol{\theta}} = \mathbf{0}, \text{ that is:}$$

$$\frac{\partial \log \mathcal{L}(\boldsymbol{\theta}; \mathbf{y})}{\partial \boldsymbol{\theta}} = \sum_{i=1}^n \frac{g'(\boldsymbol{\theta}' \mathbf{x}_i) (y_i - g(\boldsymbol{\theta}' \mathbf{x}_i))}{g(\boldsymbol{\theta}' \mathbf{x}_i) (1 - g(\boldsymbol{\theta}' \mathbf{x}_i))} \mathbf{x}_i = \mathbf{0}$$

- Nonlinear equation that generally has to be numerically solved

- 2 We have

$$\boldsymbol{\theta}_{MLE} \xrightarrow{dist.} \mathcal{N}(\mathbf{0}, I(\boldsymbol{\theta}_0)^{-1})$$

$$\text{where } I(\boldsymbol{\theta}_0) \approx - \frac{\partial^2 \log \mathcal{L}(\boldsymbol{\theta}_{MLE}; \mathbf{y}, \mathbf{x})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}$$

Marginal Effects

Marginal Effects

- 1 The marginal effect is the effect on the probability that $y_i = 1$ of a marginal increase of $x_{i,k}$:

$$m_{i,k} \equiv \frac{\partial p_i}{\partial x_{i,k}} = \underbrace{g'(\boldsymbol{\theta}' \mathbf{x}_i)}_{>0} \theta_k$$

- 2 The sign of the marginal effect $m_{i,k}$ is the one of θ_k
- 3 It can be estimated by $\hat{m}_{i,k} = g'(\boldsymbol{\theta}'_{MLE} \mathbf{x}_i) \theta_{MLE,k}$
 - the marginal effect $m_{i,k}$ depends on \mathbf{x}_i , then it is specific to each entity i
- 4 Two solutions to have "aggregate" marginal effects
 - Marginal Effect at the Mean (MEM) : $g'(\boldsymbol{\theta}'_{MLE} \bar{\mathbf{x}}) \theta_{MLE,k}$
 - Average Marginal Effect (AME) : $\frac{1}{n} \sum_{i=1}^n \hat{m}_{i,k}$

Goodness of fit

McFadden's **pseudo** - R^2

- The **pseudo** - R^2 :

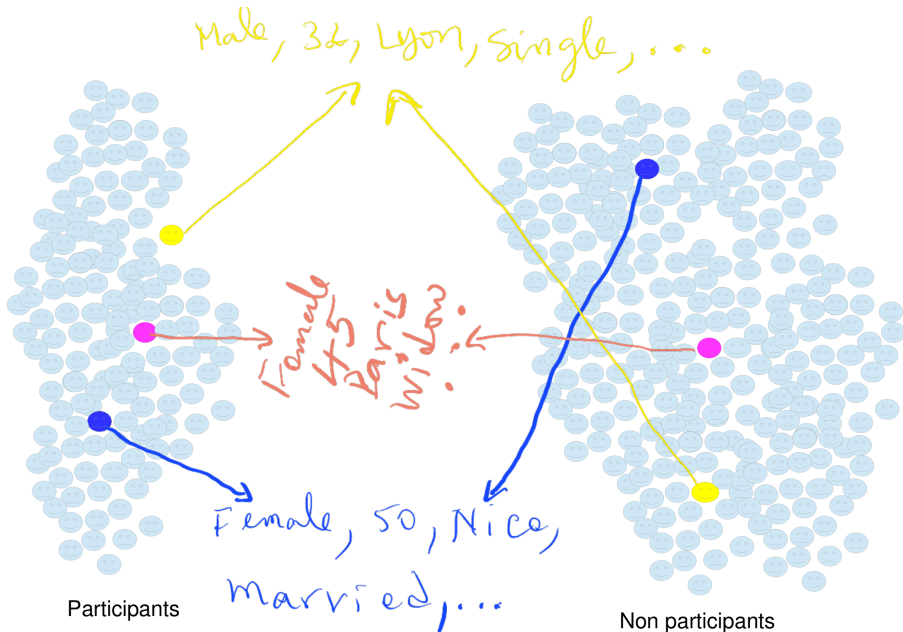
$$\text{pseudo} - R^2 = 1 - \frac{\mathcal{L}(\theta; \mathbf{y}, \mathbf{x})}{\mathcal{L}_0(\mathbf{y})}$$

- with $\mathcal{L}_0(\mathbf{y})$ the (maximum) log-likelihood that would be obtained for a model containing only a constant term (i.e. with $x_i = 1$ for all i).
- Intuitively, **pseudo** - R^2 will be 0 if the explanatory variables do not allow to predict the outcome variable (y).

Outline

- 1 MLE
- 2 Binary outcome
- 3 PSM
- 4 Application

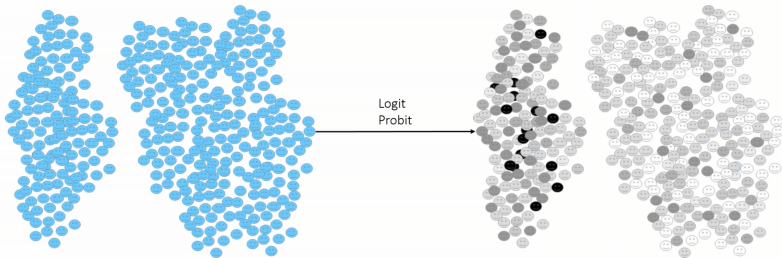
- Applicable when:
 - Treated and untreated have not been randomly assigned
 - Rich set of available information on both treated and untreated
- Aims to approximate the results of random assignment by searching within **the sample of untreated individuals** for those who are **similar** to the **treated individuals**



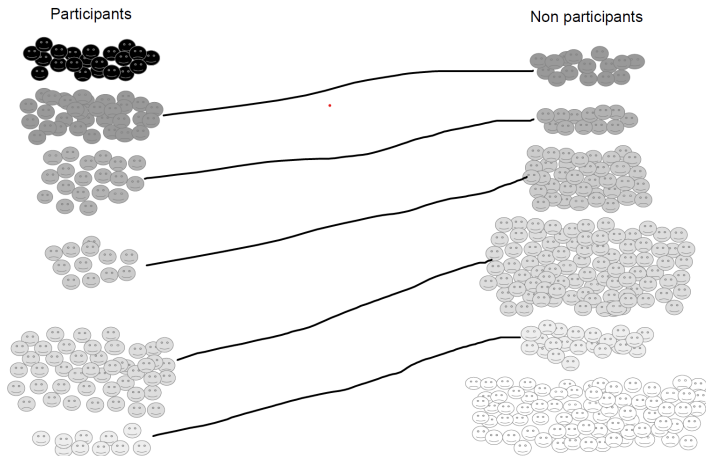
- **X**: rich set of information (age, education, sex, place of residence of observation, etc...)
- **Score**: $p_i \equiv P(\mathbf{T}_i = 1) = g(\boldsymbol{\theta}' \mathbf{X}_i)$

Participants + Non participants

score computed for participants and non participants

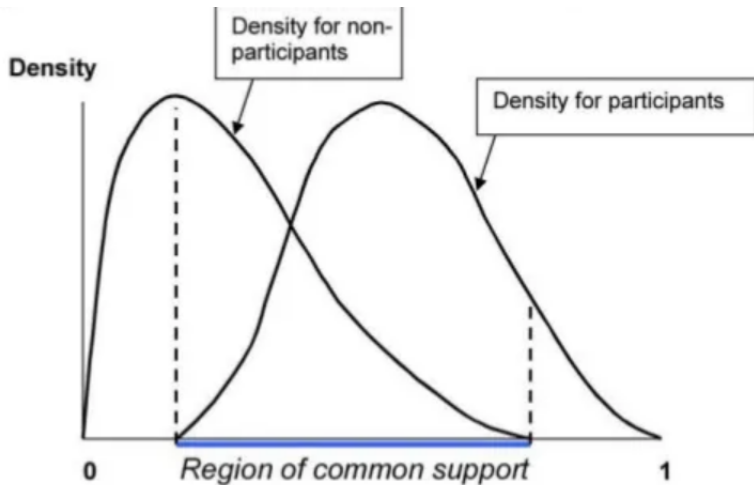


The probability of participation increases with the darkness of the dots.



Treatment effect is computed as the difference in means between matched treated and untreated

Condition of common support



Condition of sample balancing

The average values of observable characteristics \mathbf{X} calculated for treated and untreated should be close for values close to the propensity score

Condition of CIA

Conditional Independence Assumption (CIA): treatment assignment is independent of potential outcomes after conditioning on the set of observed characteristics:

$$E[Y_1|T = 1, X] = \underbrace{E[Y_1|T = 0, X]}_{\text{Unobserved}}$$

$$\underbrace{E[Y_0|T = 1, X]}_{\text{Unobserved}} = E[Y_0|T = 0, X]$$

- Difference in means between matched treated and untreated individuals:

$$\underbrace{E[Y_1|T=1, X] - E[Y_0|T=0, X]}_{\text{Observed}} = \underbrace{E[Y_1|T=1, X] - E[Y_0|T=1, X]}_{\text{ATT}} + \underbrace{E[Y_0|T=1, X] - E[Y_0|T=0, X]}_{\text{Bias}}$$

- Under CIA, the bias is null and the simple difference corresponds to ATT

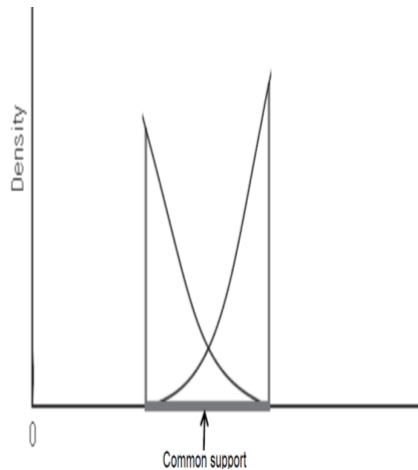
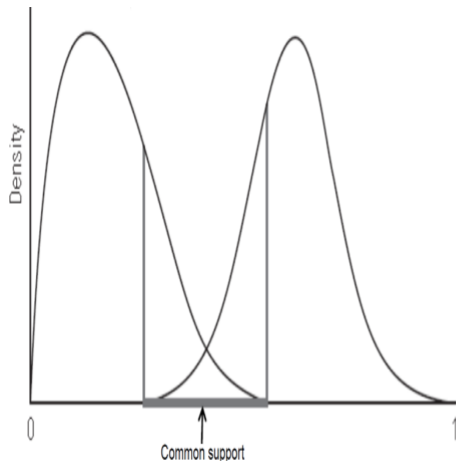
Implementation

- Select a set of conditioning variables and compute the propensity score (logit, probit, LPM, etc.)
- Check the common support
- Match individuals and check the sample balancing
- Estimate the average treatment effects of interest

Limits: unobservable characteristics

- Selection issues are only due to observable characteristics
 - After controlling for observable characteristics, no selection issues.
 - What about unobservable characteristics (namely self-selection) ?

Limits: common support and local effect



Outline

- 1 MLE
- 2 Binary outcome
- 3 PSM
- 4 Application

Application

Blattman , C., Annan, J. (2010). The consequences of child soldiering. The review of economics and statistics, 92(4), 882-898.

Thank you for your attention !