

Section 20

Lecture 7

Plan

- Estimators (Horvitz Thompson).
 - Generalization
 - IPW
- Using such estimators in trials.

Logistic regression

Suppose $Y \in \{0, 1\}$. Define $\beta = [\beta_1, \beta_2, \dots, \beta_k]^T$ as a vector of k parameter and consider a k dimensional covariate \mathbf{X} . Then the logistic model is defined as

$$\text{logit}(\mathbb{E}[Y_i | \mathbf{X}_i]) = \text{logit}(p_i) = \log\left(\frac{p_i}{1 - p_i}\right) = \beta^T \mathbf{X}_i,$$

or, equivalently, we can write that that Y follows a Bernoulli distribution,

$$\begin{aligned} P(Y_i = y | \mathbf{X}_i) &= p_i^y (1 - p_i)^{1-y} = \left(\frac{e^{\beta^T \mathbf{X}_i}}{1 + e^{\beta^T \mathbf{X}_i}}\right)^y \left(1 - \frac{e^{\beta^T \mathbf{X}_i}}{1 + e^{\beta^T \mathbf{X}_i}}\right)^{1-y} \\ &= \frac{e^{\beta^T \mathbf{X}_i \cdot y}}{1 + e^{\beta^T \mathbf{X}_i}}. \end{aligned}$$

Thus the likelihood is $\mathcal{L}(\beta) = \prod_{i=1}^n p_i^{Y_i} (1 - p_i)^{1 - Y_i}$, which can be solved numerically, e.g. solving the score equations (you can derive this from the log-likelihood, take derivatives wrt. β).

$$\sum_{i=1}^n \left(\frac{1}{X_i}\right) \left(Y_i - \frac{\exp(\beta^T X_i)}{1 + \exp(\beta^T X_i)}\right) = 0.$$

M-estimation, preliminaries

Consider a generic *statistical* model, and suppose we have i.i.d. random vectors Z_1, \dots, Z_n where $Z \sim \mathbb{P}_Z(z)$ from this model. Let θ be a k dimensional parameter. If θ fully characterizes $\mathbb{P}_Z(z)$, then we write $\mathbb{P}_Z(z; \theta)$. Let θ_0 denote the true value of θ . It follows that if θ fully characterizes $\mathbb{P}_Z(z)$, then the true density is $\mathbb{P}_Z(z; \theta_0)$. We are considering the (classical) statistical problem of deriving an estimator for θ .

Definition (M-estimator)

An M-estimator for θ is the solution $\hat{\theta}$ (assuming that it exists and is well defined) to the $(k \times 1)$ system of estimating equations

$$\sum_{i=1}^n M(Z_i; \hat{\theta}) = 0,$$

We say that $M(z; \theta) = \{M_1(z; \theta), \dots, M_k(z; \theta)\}^T$ is an *unbiased estimating function* for $\mathbb{E}_\theta(M(Z_i; \theta)) = 0$. The expectation is taken wrt. the distribution of Z at the law indexed by θ . From now on, we will suppress the subscript when we evaluate the expectation in the true value θ_0 , i.e. $\mathbb{E}(M(Z_i; \theta)) \equiv \mathbb{E}_{\theta_0}(M(Z_i; \theta))$.

MLE is an M-estimator

Consider a fully parametric model $\mathbb{P}_Z(z; \theta)$. Define,

$$M(z; \theta) = \frac{\delta \log(\mathbb{P}_Z(z; \theta))}{\delta \theta},$$

where the right hand side is a k dimensional vector of partial derivatives. Solving an estimating equation with this $M(z; \theta)$ yields a maximum likelihood estimator (MLE), and thus the MLE is an M-estimator.

Methods of moment estimators are M-estimators

Consider a fully parametric model $\mathbb{P}_Z(z; \theta)$. Define,

$$M_m(Z_i; \theta) = Z_i^m - \mathbb{E}_\theta(Z_i^m),$$

where $m = 1, \dots, k$, i.e. k is the dimension of θ .

Overview of properties of M-estimators

Theorem (M-estimator)

Under suitable regularity conditions, $\hat{\theta}$ is a consistent and asymptotically normal estimator,

$$\hat{\theta} \xrightarrow{P} \theta_0$$

and

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{D} \mathcal{N}(0, \Sigma),$$

where Σ is a covariance matrix.

Sufficient regularity conditions for M-estimators

Suppose that the following regularity conditions hold:

- ① $\sup_{\theta \in \Theta} |M_n(\theta) - M_0(\theta)| \xrightarrow{P} 0$.
- ② For all $\epsilon > 0$, $\inf\{|M_0(\theta)| : d(\theta, \theta_0) \geq \epsilon\} > 0 = |M_0(\theta_0)|$.
For this condition it is sufficient that there exists a unique solution, Θ is compact and M is continuous.
- ③ $M_n(\hat{\theta}_n) = o_P(1)$.

where $M_n(\theta) = \mathbb{E}_n(M(Z; \theta))$ is the expectation over the empirical distribution and $M_0(\theta) = \mathbb{E}(M(Z; \theta))$ over the true data generating law.

Proof that the conditions above are sufficient for the consistency of M-estimators

Proof.

From the 2nd condition, for all $\epsilon > 0$ there is a $\delta > 0$ such that

$$\begin{aligned} & P(d(\hat{\theta}_n, \theta_0) \geq \epsilon) \\ & \leq P(|M_0(\hat{\theta}_n)| - |M_0(\theta_0)| \geq \delta) \\ & = P(|M_0(\hat{\theta}_n)| - |M_n(\hat{\theta}_n)| + |M_n(\hat{\theta}_n)| - |M_n(\theta_0)| + |M_n(\theta_0)| - |M_0(\theta_0)| \geq \delta) \\ & \leq P(|M_0(\hat{\theta}_n)| - |M_n(\hat{\theta}_n)| \geq \frac{\delta}{3}) + P(|M_n(\hat{\theta}_n)| - |M_n(\theta_0)| \geq \frac{\delta}{3}) + \\ & \quad P(|M_n(\theta_0)| - |M_0(\theta_0)| \geq \frac{\delta}{3}). \end{aligned}$$

Condition 1 implies that the first and third probabilities go to zero.
Condition 3 implies that the second goes to zero.



Example: Smoking Cessation A on weight gain Y .

1566 cigarette smokers aged 25-74 years. The outcome weight gain measured after 10 years.

| Mean baseline characteristics | A | |
|----------------------------------|------|------|
| | 1 | 0 |
| Age, years | 46.2 | 42.8 |
| Men, % | 54.6 | 46.6 |
| White, % | 91.1 | 85.4 |
| University, % | 15.4 | 9.9 |
| Weight, kg | 72.4 | 70.3 |
| Cigarettes/day | 18.6 | 21.2 |
| Years smoking | 26.0 | 24.1 |
| Little exercise, % | 40.7 | 37.9 |
| Inactive life, % | 11.2 | 8.9 |

Miguel A Hernan and James M Robins. *Causal inference: What if?* CRC Boca Raton, FL:, 2018.

On estimation of causal effects

From slide 57, remember that in an experiment where A is randomised conditional on L – or more generally when consistency, positivity and exchangeability ($Y^a \perp\!\!\!\perp A \mid L$) hold – we have that

$$\begin{aligned}\mathbb{E}(Y^a) &= \sum_l \mathbb{E}(Y \mid L = l, A = a) \Pr(L = l) \\ &= \mathbb{E} \left[\frac{I(A = a)}{\pi(A \mid L)} Y \right],\end{aligned}$$

where $\pi(a \mid l) = P(A = a \mid L = l)$.

This equality motivates different estimators.

Regression estimator

We can also write

$$\begin{aligned}\mathbb{E}(Y^a) &= \sum_I \mathbb{E}(Y \mid L = I, A = a) \Pr(L = I) \\ &= \mathbb{E}(\mathbb{E}(Y \mid L, A = a)),\end{aligned}$$

where the outer expectation in the second line is with respect to the marginal of L . Denote

$$\mathbb{E}(Y \mid L = I, A = a) = Q(I, a).$$

$Q(I, a)$ is usually unknown, even in an experiment.

Regression estimator

Consider a parametric regression model $Q(I, a; \beta)$ of $Q(I, a)$; that is a linear or nonlinear function of (I, a) and the finite-dimensional parameter β .

We estimate β from the observed data. For example, we could in our conditional randomised trial pose a simple linear model

$$Q(I, a; \beta) = \beta_1 + \beta_2 a + \beta_3^T I,$$

which can be fitted with least squares methods.

If the outcome is binary ($Y \in \{0, 1\}$), we could fit a logistic regression model such as

$$\text{logit}\{Q(I, a; \beta)\} = \beta_1 + \beta_2 a + \beta_3^T I.$$

We can fit the logistic regression models with maximum likelihood estimators.

Definition (Correctly specified model)

A model is correctly specified if there exists a value β_0 such that $Q(I, a; \beta)$ evaluated at β_0 yields the true function $Q(I, a)$.

PS: As in any regression setting, the models we have posited may or may not be correctly specified.

Example continues

- We can estimate the conditional sample mean $\hat{\mathbb{E}}(Y | A = 1) = 4.5$ in quitters and $\hat{\mathbb{E}}(Y | A = 0) = 2.0$ in non-quitters. More specifically, the difference is

$$\hat{\mathbb{E}}(Y | A = 1) - \hat{\mathbb{E}}(Y | A = 0) = 2.5 \text{ (95% CI : 1.7, 3.4),}$$

but we will not assign a causal interpretation to the estimate.

- Let L include the baseline variables sex (0: male, 1: female), age (in years), race (0: white, 1: other), education (5 categories), intensity and duration of smoking (number of cigarettes per day and years of smoking), physical activity in daily life (3 categories), recreational exercise (3 categories), and weight (in kg).
- Suppose $A \perp\!\!\!\perp Y^a | L$.

Standardization: A natural way of estimating counterfactual outcomes

If we knew $Q(l, a)$, a natural way of estimating $\mathbb{E}(Y^a)$ is by the empirical average

$$\frac{1}{n} \sum_{i=1}^n Q(L_i, a),$$

motivated by the identification formula expression $\mathbb{E}(\mathbb{E}(Y | L, A = a))$. When we do not know $Q(l, a)$, but we assume that our model $Q(L_i, a; \beta)$ is correctly specified, we can use the outcome regression estimator to get the estimator

$$\hat{\mu}_{REG}(a) = \frac{1}{n} \sum_{i=1}^n Q(L_i, a; \hat{\beta}).$$

For example, using the linear estimator from the previous slide, we can estimate $\mathbb{E}(Y^{a=1}) - \mathbb{E}(Y^{a=0})$ by

$$\frac{1}{n} \sum_{i=1}^n Q(L_i, 1; \hat{\beta}) - \frac{1}{n} \sum_{i=1}^n Q(L_i, 0; \hat{\beta}) = \hat{\beta}_2,$$

that is, the regression parameter is the causal effect.

More broadly, our causal effects are not equal to regression coefficients

- Whereas the causal effect turned out to be equal to a regression coefficient in the previous slide, regression coefficients are not necessarily equal to our causal effect of interest.
- For example, the coefficients in the logistic regression model

$$\text{logit}\{Q(I, a; \beta)\} = \beta_1 + \beta_2 a + \beta_3^T I$$

do not necessarily translate to a causal effect of interest.

Standardization

We say that standardization is a plug-in g-formula estimator because it simply replaces the conditional means and probabilities in the g-formula by their estimates.

Section 21

Propensity score methods

Matching on the propensity score (intuitive motivation)

- In your homework you showed that, for all a ,

$$Y^a \perp\!\!\!\perp A \mid L \implies Y^a \perp\!\!\!\perp A \mid \pi(a \mid L).$$

- We could, for each treated individual (i.e. individual with $A = 1$), match this individual with an untreated individual with *similar* propensity score.
- Then crudely compare the mean in the two groups.
- This crude comparison should be fine, but...
- Potential problems:
 - What does "similar" propensity score mean?
 - How many matches should we choose?
 - Do we really get the average treatment effect?

Conditioning on the propensity score

Because, for all a ,

$$Y^a \perp\!\!\!\perp A \mid L \implies Y^a \perp\!\!\!\perp A \mid \pi(a \mid L),$$

it follows that $\mathbb{E}(Y^a \mid \pi(a \mid L) = s) = \mathbb{E}(Y \mid A = a, \pi(a \mid L) = s)$.

Thus, we could imagine estimating $\mathbb{E}(Y \mid A = a, \pi(a \mid L) = s)$. Because $\pi(a \mid L) \in (0, 1)$ is usually a continuous variable, we will usually not see individuals with the same value s among the treated and untreated.

However, we could

- fit an outcome regression to $\mathbb{E}(Y^a \mid \pi(a \mid L) = s)$, and
- match units i, j such that $d(s_i, s_j) < \delta$ for $\delta > 0$. For example percentiles. This is an *ad hoc* strategy, which works well in some practical settings, but we will not pursue it further.

Motivation for inverse probability weighting (IPW)

- We would like to adjust for confounding: intuitively, imbalance between L 's among those who are treated and untreated.
- Suppose that we find a treated subject i , who due to confounders was *unlikely* to be treated. That is, $\pi(1, L_i)$ is small.
- We *upweigh* her, so that she represents herself but also the others like herself (in terms of L) who were unexposed.
- Similarly, we upweigh untreated individuals with a small value of $\pi(0, L_i)$.
- Heuristically, we can think about the weighted sample as a pseudopopulation where we observe each individual for each exposure level. In particular, $\pi^*(0, L_i) = \pi^*(1, L_i)$ for all i in the weighted population (which we indicate by the $*$).
- In this pseudopopulation, confounders are balanced between treatment groups, and a crude comparison estimates a causal effect (Intuitively, we get a new DAG for this pseudopopulation, where the arrow from L to A is omitted).

Motivating example

Suppose the counterfactual data are:

| Group: | A | | B | | C | | |
|------------------|---|---|---|---|---|---|---|
| Response Y^1 : | 1 | 1 | 1 | 2 | 2 | 3 | 3 |
| Response Y^0 : | 0 | 0 | 0 | 1 | 1 | 2 | 2 |

and the average treatment effect $\mathbb{E}(Y^{a=1}) - \mathbb{E}(Y^{a=0}) = 1$.
but we observe:

| Group: | A | | B | | C | | | | |
|------------------|---|---|---|---|---|---|---|---|---|
| Response Y^1 : | 1 | 1 | ? | ? | 2 | ? | 3 | ? | ? |
| Response Y^0 : | ? | ? | 0 | 1 | ? | 1 | ? | 2 | 2 |

The naive contrast $\mathbb{E}(Y | A = 1) - \mathbb{E}(Y | A = 0) = \frac{7}{4} - \frac{6}{5} = 0.55$.

Example from Oliver Dukes.

Example continues

- However, from the table we see that,

$$\hat{\pi}(1, \text{group A}) = \frac{2}{3},$$

$$\hat{\pi}(1, \text{group B}) = \frac{1}{3},$$

$$\hat{\pi}(1, \text{group C}) = \frac{1}{3}.$$

- Let us estimate $\mathbb{E}(Y^{a=1})$ by a weighted average, where each observation is weighted by $\frac{1}{\hat{\pi}(1, \text{group X})}$, Group X $\in \{\text{Group A, Group B, Group C}\}$,

$$\frac{(1+1)\frac{3}{2} + 2\frac{3}{1} + 3\frac{3}{1}}{\frac{3}{2} + \frac{3}{2} + \frac{3}{1} + \frac{3}{1}} = 2.$$

and estimate $\mathbb{E}(Y^{a=0})$ by weighting each observation by $\frac{1}{\hat{\pi}(0, \text{Group X})}$, Group X $\in \{\text{Group A, Group B, Group C}\}$,

$$\frac{0\frac{3}{1} + (1+1)\frac{3}{2} + (2+2)\frac{3}{2}}{\frac{3}{1} + \frac{3}{2} + \frac{3}{2} + \frac{3}{2} + \frac{3}{2}} = 1.$$