Diagnosing and Correcting Balance Assessments in High Dimensions

Mark M. Fredrickson (mfredric@umich.edu) Joint work with Ben Hansen

2020-01-07

13th International Conference on Health Policy Statistics

In both randomized controlled trials and observational studies, unconfounded analysis depends on the treatment and control groups having the equivalent distributions on background variables x. In both randomized controlled trials and observational studies, unconfounded analysis depends on the treatment and control groups having the equivalent distributions on background variables x.

In randomized trials, this condition is only **guaranteed in expectation**. In observational studies, researchers attempt to **enforce this condition with matching or weighting**.

In both randomized controlled trials and observational studies, unconfounded analysis depends on the treatment and control groups having the equivalent distributions on background variables x.

In randomized trials, this condition is only **guaranteed in expectation**. In observational studies, researchers attempt to **enforce this condition with matching or weighting**.

In either case, assessing the similarity of treatment and control groups ("balance") is an important component of study design.

Notation

- There are *n* units (people, clinics, schools)
- The *p* variables are collected in x
- Each unit is assigned to treatment $Z_i = 1$ or control $Z_i = 0$.
- It will be useful to define

$$J_i = rac{Z_i - P(Z_i = 1)}{P(Z_i = 1) P(Z_i = 0)}$$

Mahalanobis Distance

Hansen and Bowers (2008) and Morgan and Rubin (2012) assess balance with a Mahalanobis distance, a normalized difference of group means **D**:

$$M = \mathbf{D}^{\prime} \operatorname{Cov} \left(\mathbf{D} \right)^{-} \mathbf{D} = \mathbf{J}^{\prime} \mathbf{x} \operatorname{Cov} \left(\mathbf{x}^{\prime} \mathbf{J} \right)^{-} \mathbf{x}^{\prime} \mathbf{J}$$

Mahalanobis Distance

Hansen and Bowers (2008) and Morgan and Rubin (2012) assess balance with a Mahalanobis distance, a normalized difference of group means **D**:

$$M = \mathbf{D}^{\prime} \operatorname{Cov} \left(\mathbf{D} \right)^{-} \mathbf{D} = \mathbf{J}^{\prime} \mathbf{x} \operatorname{Cov} \left(\mathbf{x}^{\prime} \mathbf{J} \right)^{-} \mathbf{x}^{\prime} \mathbf{J}$$

A common special case is **complete random assignment** with n_1 assigned to treatment and $n_0 = n - n_1$ to control:

$$M = \frac{n_1 n_0}{n} (\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_0)' S^2(\mathbf{x})^{-1} (\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_0)$$

where X_1 and X_0 are group means and $S^2(x)$ is the sample covariance matrix.

Since we are scaling by the inverse covariance matrix there is a convenient representation of the Mahalanobis distance using the principal components of x:

$$M = \mathbf{J}' \mathbf{x} (\mathbf{x}' \mathsf{E} (\mathbf{J}\mathbf{J}') \mathbf{x})^{-} \mathbf{x}' \mathbf{J} = \mathbf{J}' \mathbf{u} \mathbf{u}' \mathsf{E} (\mathbf{J}\mathbf{J}')^{-} \mathbf{u} \mathbf{u}' \mathbf{J},$$

where **u** is from the SVD of $\mathbf{x} = \mathbf{u}\mathbf{d}\mathbf{v}'$.

Since we are scaling by the inverse covariance matrix there is a convenient representation of the Mahalanobis distance using the principal components of x:

$$M = \mathbf{J}' \mathbf{x} (\mathbf{x}' \mathsf{E} (\mathbf{J}\mathbf{J}') \mathbf{x})^{-} \mathbf{x}' \mathbf{J} = \mathbf{J}' \mathbf{u} \mathbf{u}' \mathsf{E} (\mathbf{J}\mathbf{J}')^{-} \mathbf{u} \mathbf{u}' \mathbf{J},$$

where **u** is from the SVD of $\mathbf{x} = \mathbf{u}\mathbf{d}\mathbf{v}'$.

Going forward, we'll think of p as the rank of u' E(JJ') u.

Since we are scaling by the inverse covariance matrix there is a convenient representation of the Mahalanobis distance using the principal components of x:

$$M = \mathbf{J}' \mathbf{x} (\mathbf{x}' \mathsf{E} (\mathbf{J}\mathbf{J}') \mathbf{x})^{-} \mathbf{x}' \mathbf{J} = \mathbf{J}' \mathbf{u} \mathbf{u}' \mathsf{E} (\mathbf{J}\mathbf{J}')^{-} \mathbf{u} \mathbf{u}' \mathbf{J},$$

where **u** is from the SVD of $\mathbf{x} = \mathbf{u}\mathbf{d}\mathbf{v}'$.

Going forward, we'll think of p as the rank of u' E(JJ') u.

As an added bonus, we now have a natural way of ordering the covariates.

Inference when n > p

In an observational trial M is a useful test statistic for H_0 : $\pi = \pi_0$ (Hansen and Bowers, 2008).

In an observational trial M is a useful test statistic for H_0 : $\pi = \pi_0$ (Hansen and Bowers, 2008).

In randomized trials, *M* can be used to select particular Z that have controlled imbalance (Morgan and Rubin, 2012).

In an observational trial M is a useful test statistic for H_0 : $\pi = \pi_0$ (Hansen and Bowers, 2008).

In randomized trials, *M* can be used to select particular Z that have controlled imbalance (Morgan and Rubin, 2012).

With the number of covariates *p* being fixed and *n* tending to infinity,

$$M \stackrel{\mathsf{D}}{\rightarrow} \chi^2(p)$$

(Hansen and Bowers, 2008; Li et al., 2018)

Cerdá et al. (2012) reported an natural experiment in Medellín, Colombia:

• Intervention: Public works transportation program.

Cerdá et al. (2012) reported an natural experiment in Medellín, Colombia:

- Intervention: Public works transportation program.
- Units: 48 neighborhoods (25 with stations, 23 without) matched in pairs and two triples.

Cerdá et al. (2012) reported an natural experiment in Medellín, Colombia:

- Intervention: Public works transportation program.
- Units: 48 neighborhoods (25 with stations, 23 without) matched in pairs and two triples.
- Covariates: Survey responses, governmental data, mix of types (48 with interactions)

Cerdá et al. (2012) reported an natural experiment in Medellín, Colombia:

- Intervention: Public works transportation program.
- Units: 48 neighborhoods (25 with stations, 23 without) matched in pairs and two triples.
- Covariates: Survey responses, governmental data, mix of types (48 with interactions)
- Balance Assessment: Validating matching strategy created **comparable pairs and triples**.

Cerdá et al. (2012) Empirical Mahalanobis Distance



Figure 1: Mahalanobis distance distribution for 48 matched neighborhoods in reported in Cerdá et al. (2012)

A degenerate *M* provides no information on suitability of observational studies or equivalence at baseline in randomize trials.

A degenerate *M* provides no information on suitability of observational studies or equivalence at baseline in randomize trials.

As the previous slides demonstrated, selecting k < n principal components provides a non-degenerate distribution.

A degenerate *M* provides no information on suitability of observational studies or equivalence at baseline in randomize trials.

As the previous slides demonstrated, selecting k < n principal components provides a non-degenerate distribution.

We define reduced rank Mahalanobis distance as:

 $M_k = \mathbf{J}' \mathbf{u}_k \mathbf{u}'_k \, \mathsf{E} \, (\mathbf{J} \mathbf{J}')^- \, \mathbf{u}_k \mathbf{u}'_k \mathbf{J}$

where \mathbf{u}_k collects the first k columns of \mathbf{u} .

Picking k

Typical rules for picking k (e.g., including 80% of variance in \mathbf{x}) do not take into account the role of assignment mechanism \mathbf{Z} (Chang, 1983).

Picking k

Typical rules for picking k (e.g., including 80% of variance in \mathbf{x}) do not take into account the role of assignment mechanism \mathbf{Z} (Chang, 1983).

A natural analog in this case is maximizing the (null) variance of M_k .



Figure 2: Variance of M_k for matched neighborhoods in Medellín

Rather than approximate the distribution of M_k as $\chi^2(k)$, we could use a correction that matches both mean and variance to a scaled χ^2 :

$$M_k pprox a\chi^2(v) \Rightarrow v = rac{2k^2}{\operatorname{Var}(M_k)}, \ a = rac{k}{v}$$

Rather than approximate the distribution of M_k as $\chi^2(k)$, we could use a correction that matches both mean and variance to a scaled χ^2 :

$$M_kpprox a\chi^2(m{v}) \Rightarrow m{v} = rac{2k^2}{{
m Var}\,(M_k)}, \ a=rac{k}{v}$$

• Provides approximation to tail probabilities $P(M_k \ge m_k)$

Rather than approximate the distribution of M_k as $\chi^2(k)$, we could use a correction that matches both mean and variance to a scaled χ^2 :

$$M_kpprox a\chi^2({m v}) \Rightarrow {m v} = rac{2k^2}{{
m Var}\,(M_k)}, \ a=rac{k}{
u}$$

- Provides approximation to tail probabilities $P(M_k \ge m_k)$
- Provides approximate cutoff for selecting "rerandomization" RCT designs

Rather than approximate the distribution of M_k as $\chi^2(k)$, we could use a correction that matches both mean and variance to a scaled χ^2 :

$$M_kpprox a\chi^2({m v}) \Rightarrow {m v} = rac{2k^2}{{
m Var}\,(M_k)}, \ a=rac{k}{
u}$$

- Provides approximation to tail probabilities $P(M_k \ge m_k)$
- Provides approximate cutoff for selecting "rerandomization" RCT designs
- Alternatively, use Monte Carlo sampling from Z

Corrected Medellín Distribution



 $m_{11} = 24.5$; P ($M_{11} \ge m_{11}$): 0.01088 (1st order approx.), 0.00048 (2nd order approx.), 0.001 (empirical, 5k samples)

• **Balance assessments** compare treated and control groups for similar distributions.

- **Balance assessments** compare treated and control groups for similar distributions.
- The Mahalanobis distance statistic is useful when $n \gg p$, but becomes degenerate when $p \approx n$.

- **Balance assessments** compare treated and control groups for similar distributions.
- The Mahalanobis distance statistic is useful when $n \gg p$, but becomes degenerate when $p \approx n$.
- Reduced rank Mahalanobis distance guarantees non-degenerate distribution.

- **Balance assessments** compare treated and control groups for similar distributions.
- The Mahalanobis distance statistic is useful when $n \gg p$, but becomes degenerate when $p \approx n$.
- Reduced rank Mahalanobis distance guarantees non-degenerate distribution.
- Satterthwaite χ^2 approximation superior and relatively simple.

- **Balance assessments** compare treated and control groups for similar distributions.
- The Mahalanobis distance statistic is useful when $n \gg p$, but becomes degenerate when $p \approx n$.
- Reduced rank Mahalanobis distance guarantees non-degenerate distribution.
- Satterthwaite χ^2 approximation superior and relatively simple.
- In-progress: extend RItools package for R, formalize justification for variance maximization, implications for design (both observational and RCTs)

Thank You!