

Berry–Esseen bounds for design-based causal inference
with possibly diverging treatment levels and varying group sizes

Lei Shi

UC Berkeley, Biostatistics

January 16, 2025

joint work with Peng Ding (Berkeley Statistics)

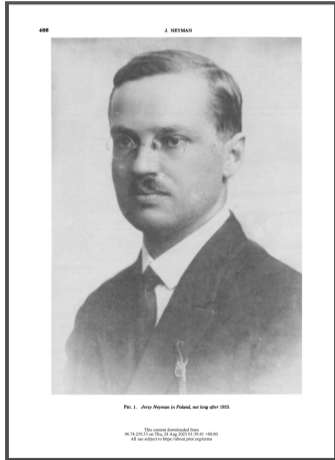
<https://arxiv.org/abs/2209.12345>

What we shall cover in the presentation

- ▶ **Technical aspects** for analyzing **completely randomized experiments** in **general settings**
- ▶ **“Technical aspects”**: bridge asymptotics and finite sample inference
 - ▶ Berry–Esseen bound (BEB): a finite sample characterization of central limit theorems
- ▶ **“Completely randomized experiments”** (CRE): the most basic design in statistics
 - ▶ Design-based inference: handle uncertainty from random sampling instead of distributional modelling
- ▶ **“General settings”**: go beyond classical multi-armed completely randomized experiments
 - ▶ Diverging treatment levels and varying group sizes: regimes not fully covered by classical literature and requiring new technical tools

Neyman (1923 Polish/1990 English)

- ▶ On the application of probability theory to agricultural experiments (100 years old!)



466

J. NEYMAN

On the Application of Probability Theory to Agricultural Experiments. Essay on Principles. Section 9.

Jerzy Szpiro-Neyman

Translated and edited by D. M. Dabrowska and T. P. Speed from the Polish original, which appeared in *Roczniki Nauk Rolniczych Tom X (1923) 1-51 (Annals of Agricultural Sciences)*

Abstract. In the portion of the paper translated here, Neyman introduces a model for the analysis of field experiments conducted for the purpose of comparing a number of crop varieties, which makes use of a double-indexed array of unknown potential yields, one index corresponding to varieties and the other to plots. The yield corresponding to only one variety will be observed on any given plot, but through an urn model embodying sampling without replacement from this doubly indexed array, Neyman obtains a formula for the variance of the difference between the averages of the observed yields of two varieties. This variance involves the variance over all plots of the potential yields and the correlation coefficient r between the potential yields of the two varieties on the same plot. Since it is impossible to estimate r directly, Neyman advises taking $r = 1$, showing that in practice this may lead to using too large an estimated standard deviation, when comparing two variety means.

Key words and phrases: Field experiment, varieties, unknown potential yields, urn model, sampling without replacement, correlation.

[Numbers in brackets correspond to page numbers in the original text.]

I will now discuss the design of a field experiment involving plots. I should mention that this is a task for an agricultural person however, because mathematicians operate only with general designs. In designing this experiment, let us consider a field divided into n equal plots and let

$$U_1, U_2, \dots, U_n$$

be the true yields of a particular variety on each of these plots. If all the members U_i are equal, each of them may be called the average yield of the field. Otherwise the average yield may be thought of as the arithmetic mean

$$\bar{u} = \frac{\sum_{i=1}^n U_i}{n}.$$

D. M. Dabrowska is Assistant Professor, Division of Biostatistics, School of Public Health, University of California, Los Angeles, California 90024-1722. T. P. Speed is Professor and Chair, Department of Statistics, University of California, Berkeley, California 94720.

466

This content downloaded from 96.74.254.21 on Thu, 24 Aug 2016 19:26:41 UTC All use subject to https://about.jstor.org/terms

Neyman (1923 Polish/1990 English)

- ▶ Well cited in causal inference literature for “potential outcomes”
 - ▶ often with Rubin (1974), sometimes called the Neyman–Rubin model
- ▶ Neyman also introduced the “design-based inference” for experiments
 - ▶ N units and Q treatment levels: $N \times Q$ fixed potential outcomes
 - ▶ treatment assignment: random permutation (the urn model)
 - ▶ inference based solely on the randomness of the treatment
 - ▶ “unbiased” estimation and “conservative” confidence interval
- ▶ “this paper represents the first attempt to evaluate, formally or informally, the repeated-sampling properties of statistics over their nonnull randomization distributions”

Slightly more general setup than Neyman (1923/1990)

- ▶ Experiment with N units and Q treatment levels
- ▶ $N \times Q$ potential outcomes: $\{Y_i(q) : i = 1, \dots, N; q = 1, \dots, Q\}$

i	$Y_i(1)$	$Y_i(2)$	\dots	$Y_i(Q)$
1	$Y_1(1)$	$Y_1(2)$	\dots	$Y_1(Q)$
\vdots	\vdots	\vdots	\ddots	\vdots
N	$Y_N(1)$	$Y_N(2)$	\dots	$Y_N(Q)$

- ▶ mean $\bar{Y}(q) = N^{-1} \sum_{i=1}^N Y_i(q)$
- ▶ vectorized mean $\bar{\mathbf{Y}} = (\bar{Y}(1), \dots, \bar{Y}(Q))^T$
- ▶ covariance $S(q, q') = (N - 1)^{-1} \sum_{i=1}^N \{Y_i(q) - \bar{Y}(q)\} \{Y_i(q') - \bar{Y}(q')\}$
- ▶ covariance matrix $\mathbf{S} = (S(q, q'))_{q, q'=1, \dots, Q}$

Slightly more general setup than Neyman (1923/1990)

- ▶ Parameter of interest $\gamma = F^T \bar{Y}$
 - ▶ F is $Q \times H$ contrast matrix, with columns orthogonal to $\mathbf{1}_Q$
 - ▶ Examples:
 - (i) ATE: $F = (1, -1)^T$;
 - (ii) Factorial effects (Dasgupta et al. 2015; Zhao & Ding 2022)
- ▶ Complete randomization of $\mathbf{Z} = (Z_1, \dots, Z_N)$: N balls into Q urns
 - ▶ fixed sample sizes N_1, \dots, N_Q with $\sum_{q=1}^Q N_q = N$
 - ▶ random permutation of N_1 1s, \dots , N_q Q s
 - ▶ $\mathbb{P}(\mathbf{Z} = \mathbf{z}) = N_1! \cdots N_Q! / N!$ for all possible values of $\mathbf{z} = (z_1, \dots, z_N)$.
- ▶ Observed outcome $Y_i = Y_i(Z_i) = \sum_{q=1}^Q Y_i(q) \mathbf{1}\{Z_i = q\}$
- ▶ *Randomization model*: fixed potential outcomes, random \mathbf{Z}

Basic statistics under the randomization model

- ▶ Sample mean $\hat{Y}_q = N_q^{-1} \sum_{Z_i=q} Y_i$
- ▶ Vectorized sample mean $\hat{Y} = (\hat{Y}_1, \dots, \hat{Y}_Q)^\top$ has covariance matrix

$$\text{cov}\{\hat{Y}\} = V_{\hat{Y}} = \text{diag}\{N_q^{-1} S(q, q)\}_{q \in [Q]} - N^{-1} S$$

- ▶ Covariance estimator $\hat{V}_{\hat{Y}} = \text{diag}\{N_q^{-1} \hat{S}(q, q)\}_{q \in [Q]}$
 - ▶ sample variance $\hat{S}(q, q)$, no sample covariance $\hat{S}(q, q')$
 - ▶ **conservative** due to the term $-N^{-1} S$
- ▶ Point estimation for $\gamma = F^\top \bar{Y}$: $\hat{\gamma} = F^\top \hat{Y}$ is unbiased
- ▶ Conservative sandwich covariance estimation: $\hat{V}_{\hat{\gamma}} = F^\top \hat{V}_{\hat{Y}} F$

Inference in CRE: established results and subtleties

- ▶ Inference on γ relies on more results
 - ▶ CLT of $\hat{\gamma}$ and consistency (or conservativeness) of $\hat{V}_{\hat{\gamma}}$
- ▶ Most existing literature focuses on the “**small Q and large N_q 's**” regime
 - ▶ treatment-control setting has a rich literature: Freedman (2008), Lin (2013), Imbens and Rubin (2015)
 - ▶ multi-armed experiment with a few treatment levels: Li and Ding (2017), Zhao and Ding (2023), Dasgupta et al (2015), Pashley (2023)
- ▶ With many treatment levels (Q) and small group sizes (N_q), inference is non-trivial
 - ▶ CLT has different regimes
 - ▶ Might need new construction of variance estimator
 - ▶ Consistency of variance estimation requires new proof

A canonical example: 2^K factorial design

- ▶ K binary factors generate $Q = 2^K$ treatment levels
 - ▶ treatment levels $q = 1, \dots, Q \iff$ factor levels: $z_1, \dots, z_K = 0, 1$
- ▶ Potential outcomes $Y_i(q) \iff Y_i(z_1, \dots, z_K)$
- ▶ $\gamma = F^\top \bar{Y}$ may contain a subset of the factorial effects
 - ▶ Wu and Hamada (2021 book) and Dasgupta et al (2015)
 - ▶ recall \bar{Y} is the vector of mean potential outcomes
 - ▶ F has orthogonal columns; each column has half Q^{-1} and half $-Q^{-1}$
- ▶ Even moderate K generates large Q
- ▶ Factorial experiments may or may not have replications

Real world examples from literature

- ▶ Example 1: agricultural screening trials.
 - ▶ Brownie and Boos (1994): discussed one study that compares different plant varieties in resisting aphid infestation. The study involves $Q = 35$ plant varieties and $N_q = 4$ replications within each treatment arm. [Large Q small N_q 's]
 - ▶ Casler (2015): “Numerous special situations exist for which there is a strong temptation or need to **devote all resources toward multiple treatments and none to replication** or independent observations of those treatments”. [Unreplicated designs]
- ▶ Example 2: partially nested experiments and provider effect in behavioral study
 - ▶ Bauer et al. (2008): participants suffering from depression might be assigned to one of two study arms: cognitive-behavioral group therapy (CBT) or control. Individuals assigned to CBT are administered treatment within small groups. Control participants, in contrast, are not placed into groups and have no particular relationship to one another. [Mixture regimes]

Summary of the general regimes

Regime	Q	N_q	CLT, variance estimation, and BEB
(R1)	Small	Large	CLT and variance estimation; no BEB
(R2)	Large	Large	Seems similar to (R1) but not studied
(R3)	Large	Small but $N_q \geq 2$	Not studied
(R4)	Large	$N_q = 1$	Not studied; variance estimation is nontrivial
(R5)	Mixture of the above		Not studied

- ▶ (R1)–(R4): nearly uniform design with roughly the same sample sizes across treatment groups: $N_q = c_q N_0$ with bounded c_q for some N_0
- ▶ (R5): general design with varying group sizes
- ▶ Question: can we establish an inference scheme that unifies the above regimes?
⇒ general BEBs and variance estimation?

A BEB based on BEB for linear permutational statistic

- ▶ Standardize $\hat{\gamma}$ as $\tilde{\gamma} = V_{\hat{\gamma}}^{-1/2}(\hat{\gamma} - \gamma)$
- ▶ Unifying the regimes: bound $|\mathbb{P}\{b^\top \tilde{\gamma} \leq t\} - \Phi(t)|$ with population quantities?
- ▶ Write $\tilde{\gamma}$ as a linear permutational statistic $\Gamma = (\Gamma_1, \dots, \Gamma_H)^\top$ with

$$\Gamma_h = \sum_{i=1}^N M_h(i, \pi(i)) : \text{where } \pi \text{ is random permutation}$$

- ▶ M_h : a set of matrices satisfying certain standardization conditions (details in paper)
- ▶ Main Theorem of Bolthausen (1984): There exists an absolute constant $C > 0$, such that

$$\sup_{t \in \mathbb{R}} |\mathbb{P}\{\Gamma_1 \leq t\} - \Phi(t)| \leq \frac{C}{N} \sum_{i,j \in [N]} |M_1(i,j)|^3.$$

BEB #1: BEB for linear contrasts

- ▶ Apply Bolthausen (1984) to obtain a general BEB: There exists $C > 0$ such that for any vector b with $\|b\|_2 = 1$, we have

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P}\{b^\top \tilde{\gamma} \leq t\} - \Phi(t) \right| \leq C \left\| b^\top V_{\hat{\gamma}}^{-1/2} F^\top \right\|_\infty \cdot \max_{q \in [Q]} N_q^{-1} M_N(q)$$

where $M_N(q) = \max_{i \in [M]} |Y_i(q) - \bar{Y}(q)|$ is the maximum absolute deviation from the mean for $Y_i(q)$'s (Hajek 1960; Li and Ding 2017)

- ▶ The above BEB is general but (i) not uniform over b ; (ii) not intuitive for interpretation

BEB #1: BEB for linear contrasts

- ▶ Condition on trade-off between outcomes and contrast: recall

$$V_{\hat{\gamma}} = F^T V_{\hat{Y}} F = F^T \text{Diag} \{ N_q^{-1} S(q, q) \} F - N^{-1} F^T S F.$$

Assume that there exists σ_F such that $V_{\hat{\gamma}} = F^T V_{\hat{Y}} F \succeq \sigma_F^{-2} F^T \text{Diag} \{ N_q^{-1} S(q, q) \} F$

- ▶ There exists $C > 0$ such that

$$\sup_{\|b\|_2=1} \sup_{t \in \mathbb{R}} \left| \mathbb{P} \{ b^T \tilde{\gamma} \leq t \} - \Phi(t) \right| \leq C \max_{i \in [N], q \in [Q]} \min \{ \text{I}(i, q), \text{II}(i, q) \}$$

where

$$\text{I}(i, q) = \sigma_F \left| \frac{Y_i(q) - \bar{Y}(q)}{\sqrt{N_q S(q, q)}} \right|, \quad \text{II}(i, q) = \frac{\|F(q, \cdot)\|_2 \cdot N_q^{-1} |Y_i(q) - \bar{Y}(q)|}{\sqrt{\varrho_{\min} \{ F^T V_{\hat{Y}} F \}}}$$

Comment on the previous BEB: Additional condition

- ▶ We imposed an additional condition: $V_{\hat{\gamma}} = F^\top V_{\hat{Y}} F \succeq \sigma_F^{-2} F^\top \text{Diag} \{ N_q^{-1} S(q, q) \} F$
- ▶ Means that $F^\top \text{Diag} \{ N_q^{-1} S(q, q) \} F$ controls $V_{\hat{\gamma}}$ (from both up and below)
- ▶ Rules out those cases that involve extreme choices of F and S and lead to ill-conditioned covariance structure.
- ▶ holds in most “interesting” cases
 - ▶ Two-arm randomized experiments: rule out the scenario where the potential outcomes are perfectly negatively correlated (i.e., there exists a constant $c > 0$ such that $Y_i(0) = -cY_i(1)$ for all $i \in [N]$)
 - ▶ More examples in the paper: uncorrelated potential outcomes, testing sharp null, ...

Comment on the previous BEB: Two regimes in the BEB

- ▶ Term I is more useful with large N_q :

$$I(i, q) = \sigma_F \left| \frac{Y_i(q) - \bar{Y}(q)}{\sqrt{N_q S(q, q)}} \right|$$

- ▶ Term II is more useful with small N_q (but dense F):

$$II(i, q) = \frac{\|F(q, \cdot)\|_2 \cdot N_q^{-1} |Y_i(q) - \bar{Y}(q)|}{\sqrt{\varrho_{\min}\{F^\top V_{\hat{Y}} F\}}}$$

⇒ Think about one contrast case, i.e., $H = 1$;

⇒ $\|F(q, \cdot)\|_2 \leq \|F\|_\infty$, $\varrho_{\min}\{F^\top V_{\hat{Y}} F\}$ is around $\|F\|_2^2$.

A BEB for nearly uniform design

- ▶ Condition on the contrast: $\|F\|_\infty \leq cQ^{-1}$ and $\varrho_{\min}\{F^\top F\} \geq c'Q^{-1}$
 - ▶ this condition is motivated by factorial effects in **factorial designs**
 - ▶ BEB should not depend on scaling of F
 - ▶ F cannot be sparse if many N_q 's are small
 - ▶ F cannot be degenerate: if degenerate, then consider subset
- ▶ There exists $C > 0$ such that

$$\sup_{\|b\|_2=1} \sup_{t \in \mathbb{R}} \left| \mathbb{P}\{b^\top \tilde{\gamma} \leq t\} - \Phi(t) \right| \leq C\sigma_F \frac{\max_{q \in [Q]} M_N(q)}{\{\min_{q \in [Q]} S(q, q)\}^{1/2}} \sqrt{\frac{H}{N}}$$

- ▶ recall $M_N(q)$ is the maximum deviation from the mean; $\min_{q \in [Q]} S(q, q)$ is a scaling factor
- ▶ H is the number of contrast in $F =$ dimension of γ
- ▶ **Depends on $N = O(Q \cdot N_0)$, not N_0 or Q**

A BEB for general designs

- ▶ Partition treatment arms into “L(arge)” and “S(mall)” based on N_q
- ▶ Partition “S” into “R(eplicated)” and “U(nreplicated)”
- ▶ Partition $F^\top = (F_S^\top, F_L^\top)$; partition $F_S^\top = (F_U^\top, F_R^\top)$
- ▶ Condition on F_S : $\|F_S\|_\infty \leq c|Q_S|^{-1}$ and $\varrho_{\min}\{F_S^\top F_S\} \geq c'|Q_S|^{-1}$
- ▶ There exists $C > 0$ such that

$$\begin{aligned} & \sup_{\|b\|_2=1} \sup_{t \in \mathbb{R}} \left| \mathbb{P}\{b^\top \tilde{\gamma} \leq t\} - \Phi(t) \right| \\ & \leq C \sigma_F \max \left\{ \max_{q \in Q_L} \frac{M_N(q)}{\sqrt{N_q S(q, q)}}, \frac{\max_{q \in Q_S} M_N(q)}{\{\min_{q \in Q_S} S(q, q)\}^{1/2}} \cdot \sqrt{\frac{H}{N_S}} \right\} \end{aligned}$$

Design-based causal inference: the big picture

- ▶ Wald-type inference based on $\hat{T} = (\hat{\gamma} - \gamma)^\top \hat{V}_{\hat{\gamma}}^{-1} (\hat{\gamma} - \gamma)$ and χ_H^2
- ▶ Two standard steps in statistics
 - ▶ Standardized statistic $T = (\hat{\gamma} - \gamma)^\top V_{\hat{\gamma}}^{-1} (\hat{\gamma} - \gamma) \approx \chi_H^2$
 - ▶ Covariance estimation $\hat{V}_{\hat{\gamma}}$: conservative for the true covariance
- ▶ Two regimes depending on H : the dimension of F
 - ▶ small, fixed H : $T \approx \chi_H^2$
 - ▶ large, diverging H : $T \approx \chi_H^2 \approx H + \sqrt{2H} \cdot \mathcal{N}(0, 1)$
- ▶ With many N_q 's being 1, covariance estimation is non-trivial
- ▶ All the above requires new technical results

Design-based causal inference: notation and conditions

- ▶ Define $T_0 = \xi_H^\top \xi_H \sim \chi_H^2$ where $\xi_H \sim \mathcal{N}(0, I_H)$
- ▶ Moment conditions on the potential outcomes: for all q
 - ▶ there exists $\Delta > 0$ such that $N^{-1} \sum_{i=1}^N \{Y_i(q) - \bar{Y}(q)\}^4 \leq \Delta^4$
 - ▶ there exists $\nu > 0$ such that $M_N(q) \leq \nu$
 - ▶ there exists $\underline{S} > 0$ such that $S(q, q) \geq \underline{S}$
 - ▶ for simplicity, assume **bounded** $\Delta, \nu, \underline{S}$; can allow them to diverge slowly
- ▶ Important regimes
 - ▶ with replications
 - ▶ without replications
 - ▶ mixture of the above

Design-based inference: BEB over convex sets

- ▶ Need to bound the distributional distance $\sup_{t \in \mathbb{R}} |\mathbb{P}(T \leq t) - \mathbb{P}(T_0 \leq t)|$.
- ▶ Inherently a BEB for quadratic forms and not implied by BEB #1 (for linear)
- ▶ (BEB over convex sets) Assume $|M_h(i, j)| \leq B_N$. There exists a universal constant $C > 0$, such that

$$\begin{aligned} & \sup_{A \in \mathcal{A}} |\mathbb{P}\{\Gamma \in A\} - \mathbb{P}\{\xi_H \in A\}| \\ & \leq CH^{13/4}NB_N(B_N^2 + N^{-1}) + CH^{3/4}B_N + CH^{13/8}N^{1/4}B_N^{3/2} + CH^{11/8}N^{1/2}B_N^2. \end{aligned} \tag{1}$$

When $B_N = O(N^{-1/2})$, $\sup_{A \in \mathcal{A}} |\mathbb{P}\{\Gamma \in A\} - \mathbb{P}\{\xi_H \in A\}| \leq \frac{CH^{13/4}}{N^{1/2}}$.

- ▶ Established $O(N^{-1/2})$ rates using Fang and Röllin (2015), based on Stein coupling

Design-based inference: nearly uniform design with $N_q \geq 2$

- ▶ BEB #2: There exists $C > 0$ such that

$$\sup_{t \in \mathbb{R}} |\mathbb{P}(T \leq t) - \mathbb{P}(T_0 \leq t)| \leq \frac{C \max_{q \in [Q]} M_N(q)^3}{\{\min_{q \in [Q]} S(q, q)\}^{3/2}} \cdot \frac{H^{19/4}}{N^{1/2}}$$

- ▶ Conservative variance estimation: recall $\widehat{V}_{\widehat{Y}} = \text{diag}\{N_q^{-1} \widehat{S}(q, q)\}_{q \in [Q]}$
- ▶ Valid Wald-type inference if $H^{19/2}/N \rightarrow 0$
 - ▶ work with “small Q large N_q 's” and “large Q and small N_q 's”
 - ▶ not too many contrasts; particularly useful for 2^K factorial design: $K = \log N$ and $H = O(K^2)$ for main effects and two-way interactions

Design-based inference: uniform design with $N_q = 1$

- ▶ BEB for T the same; covariance estimation challenging
- ▶ Strategy one: mimicking the variance estimation for sample mean
 - ▶ write $\hat{\gamma} = F^\top \hat{Y} = Q^{-1} \sum_q QF(q, \cdot)^\top Y_q$, with observed outcome Y_q
 - ▶ covariance estimation: $\hat{V}_{\hat{\gamma}} = \mu_Q^{-1} \sum_q (QF(q, \cdot)^\top Y_q - \hat{\gamma})^{\otimes 2}$
 - ▶ correction factor $\mu_Q = Q(Q - 2)$
- ▶ Strategy two: grouping outcomes to estimate the variances
 - ▶ partition the levels into disjoint groups. $\langle g \rangle$ group for treatment q , with group mean $\hat{Y}_{\langle g \rangle}$
 - ▶ diagonal covariance estimator with $\hat{V}_{\hat{\gamma}}(q, q) = \mu_{\langle g \rangle} (Y_q - \hat{Y}_{\langle g \rangle})^2$
 - ▶ correction factor $\mu_{\langle g \rangle} = (1 - 2N^{-1})^{-1} (1 - |\langle g \rangle|^{-1})^{-2}$
- ▶ Both **conservative** but in different ways (detailed results in the paper)
- ▶ More principled covariance estimation is still an **open question**

Design-based inference: design with varying group sizes

- ▶ BEB holds, depending on the partition based on group sizes
- ▶ Covariance estimation, depending on the partition “U” and “R+L”
- ▶ Wald-type inference is conservative
- ▶ Combination of the results for previous regimes
- ▶ Details omitted

Design-based inference: some open questions

- ▶ BEB for **many contrasts**, e.g. $H \approx N$ in analysis of variance
- ▶ BEB for **studentized statistics**: $\hat{V}_{\hat{\gamma}}^{-1/2} \hat{\gamma}$
 - ▶ a non-sharp bound used in Shi, Ding and Wang (2023)
 - ▶ it may be possible to obtain a better bound using Stein's method
- ▶ **Concentration inequalities** for design-based inference: more statistical applications?
 - ▶ Bloniarz et al (2016) and Lei and Ding (2021) used some
 - ▶ S. Chatterjee used Stein's method to derive results for permutations
- ▶ Statistical issues
 - ▶ **fractional factorial design**: not all treatment levels are present in the experiment, but can assume away higher-order interactions
 - ▶ **more user-friendly statistical procedures**: regression-based analysis, covariate adjustment, more complicated designs

Related papers

- ▶ Li and Ding (2017) General forms of finite population central limit theorems with applications to causal inference. *JASA*
- ▶ Zhao and Ding (2023) Covariate adjustment in multi-armed, possibly factorial experiments. *JRSSB*
- ▶ Shi and Ding (2022) Berry–Esseen bounds for design-based causal inference with possibly diverging treatment levels and varying group sizes. *ArXiv*
- ▶ Shi Ding and Wang (2023) Forward screening and post-screening inference in factorial designs. *ArXiv*