# On randomization-based and regression-based inferences for $2^K$ factorial designs

Jiannan Lu

*Microsoft Corporation, One Microsoft Way, Redmond, WA 98052, USA*

### ARTICLE INFO

### ABSTRACT

We extend the randomization-based causal inference framework in Dasgupta et al. (2015) for general $2^K$ factorial designs, and demonstrate the equivalence between regression-based and randomization-based inferences. Consequently, we justify the use of regression-based methods in $2^K$ factorial designs from a finite-population perspective.

© 2016 Published by Elsevier B.V.

## 1. Introduction

Factorial designs, originally introduced for agricultural experiments (Fisher, 1935; Yates, 1937), have gained more popularity in recent times because of their abilities to investigate multiple treatment factors simultaneously. As pointed out by Ding (2014), although rooted in randomization theory (e.g., Kempthrone, 1952), factorial designs have been dominantly analyzed by regression methods in practice. Unfortunately, however, regression-based inference might not be suitable under certain circumstances. For example, several researchers (e.g., Miller, 2006; Lu et al., 2015) have pointed out that in many randomized experiments we cannot treat the experimental units as a random sample drawn from a hypothetical super-population, and should instead restrict the scope of inference to the finite-population of the experimental units themselves. Realizing the inherent deficiencies of regression-based inference, Dasgupta et al. (2015) advocated conducting randomization-based inference for factorial designs by utilizing the concept of potential outcomes (Neyman, 1990; Rubin, 1974). The proposed framework for balanced $2^K$ factorial designs is flexible, interpretable and applicable to both finite-population and super-population settings.

Given the advantages of randomization-based inference, it is necessary to generalize the framework in Dasgupta et al. (2015) for more general, i.e., unbalanced, $2^K$ factorial designs. Moreover, it is of great importance to reconcile randomization-based and regression-based inferences, i.e., the point estimators of the factorial effects and their corresponding confidence regions. However, although the equivalence between randomization-based and regression-based inferences for randomized treatment-control studies (i.e., $2^1$ factorial designs) has been well established in the existing literature (Schochet, 2010; Samii and Aronow, 2012; Lin, 2013), similar discussions for $2^K$ factorial designs appear to be absent. In this paper, we fulfill the aforementioned two-fold task.

The paper proceeds as follows. Section 2 extends the randomization-based inference framework in Dasgupta et al. (2015) to general $2^K$ factorial designs. Section 3 demonstrates the equivalence between randomization-based and regression-based inferences for $2^K$ factorial designs. Section 4 considers extensions, and Section 5 concludes and discusses possible future directions.

*E-mail address:* jiannl@microsoft.com.

## 2. Randomization-based inference for general $2^K$ factorial designs

### 2.1. $2^K$ factorial designs

Consider $K$ distinct factors, each with two levels -1 and 1. We construct the model matrix (Wu and Hamada, 2009) $H = (h_0, \ldots, h_{2^K-1})$ as follows (Espinosa et al., 2015):

- let $h_0 = 1_{2^K}$;
- for $k = 1, \ldots, K$, construct $h_k$ by letting its first $2^{K-k}$ entries be $-1$, the next $2^{K-k}$ entries be 1, and repeating $2^{k-1}$ times;
- for $k = K + 1, \ldots, K + \binom{K}{2}$, let $h_k = h_{k_1} \cdot h_{k_2}$, where $k_1, k_2 \in \{1, \ldots, K\}$;
  - $\ldots$
- let $h_{J-1} = h_1 \cdot \cdots \cdot h_K$.

For $j = 1, \ldots, 2^K$, let $\tilde{h}_j$ denote the $j$th row of the model matrix $H$. A well-known fact is that the model matrix $H$ is orthogonal, i.e.,

$$HH' = (\tilde{h}_j \tilde{h}_{j'}')_{2^k \times 2^K} = 2^K I_{2^K}, \qquad H'H = \sum_{j=1}^{2^K} \tilde{h}_j' \tilde{h}_j = 2^K I_{2^K}. \tag{1}$$

The $j$th row of $\tilde{H} = (h_1, \ldots, h_K)$ is the $j$th treatment combination $z_j$, and the columns of $H$ define the factorial effects. To be specific, the first column $h_0$ corresponds to the null effect, the next $K$ columns $h_1, \ldots, h_K$ correspond to the main effects of the $K$ factors, the next $\binom{K}{2}$ columns $h_{K+1}, \ldots, h_{K+\binom{K}{2}}$ correspond to the two-way interactions, etc., and eventually the last column $h_{J-1}$ corresponds to the $K$-factor interaction.

**Example 1.** For $2^2$ factorial designs, the model matrix is:

$$H = \begin{matrix} & \overset{h_0}{} & \overset{h_1}{} & \overset{h_2}{} & \overset{h_3}{} \\ \tilde{h}_0 & \\ \tilde{h}_1 & \\ \tilde{h}_2 & \\ \tilde{h}_3 & \end{matrix} \begin{pmatrix} 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & 1 & 1 & 1 \end{pmatrix}.$$

The four treatment combinations are $z_1 = (-1, -1)$, $z_2 = (-1, 1)$, $z_3 = (1, -1)$ and $z_4 = (1, 1)$. We represent the main effects of factors 1 and 2 by $h_1 = (-1, -1, 1, 1)'$ and $h_2 = (-1, 1, -1, 1)'$ respectively, and the two-way interaction by $h_3 = (1, -1, 1, -1)'$.

### 2.2. Randomization-based inference

For consistency, we adopt the notations in Dasgupta et al. (2015). Let $N \geq 2^{K+1}$ be the number of experimental units. Under the Stable Unit Treatment Value Assumption (Rubin, 1980), for unit $i$, we denote its potential outcome under treatment combination $z_j$ as $Y_i(z_j)$, for all $j = 1, \ldots, 2^K$. Let $Y_i = \{Y_i(z_1), \ldots, Y_i(z_{2^K})\}'$, and we define the factorial effect vector of unit $i$ as

$$\tau_i = \frac{1}{2^{(K-1)}} H' Y_i. \tag{2}$$

Having defined the potential outcomes and factorial effects on the individual-level, we shift focus to the population-level. For all $j$, we let

$$\bar{Y}(z_j) = \frac{1}{N} \sum_{i=1}^{N} Y_i(z_j)$$

be the average potential outcome under treatment combination $z_j$, across all experimental units. Let $\bar{Y} = \{\bar{Y}(z_1), \ldots, \bar{Y}(z_{2^K})\}'$, and we define the population-level factorial effect vector as

$$\tau = \frac{1}{N} \sum_{i=1}^{N} \tau_i = \frac{1}{2^{(K-1)}} H' \bar{Y}. \tag{3}$$

We consider general $2^K$ factorial designs. For $j = 1, \ldots, 2^K$, we randomly assign $n_j \geq 2$ units to treatment $z_j$. Note that $\sum_{j=1}^{2^K} n_j = N$. For unit $i$, we let

$$W_i(z_j) = \begin{cases} 1, & \text{if unit } i \text{ is assigned treatment } z_j, \\ 0, & \text{otherwise.} \end{cases}$$