Taylor & Francis
Taylor & Francis Group

Check for updates

# Causal Interaction in Factorial Experiments: Application to Conjoint Analysis

Naoki Egami [a] and Kosuke Imai [b,c]

[a]Department of Politics, Princeton University, Princeton, NJ; [b]Department of Government, and Department of Statistics, Harvard University, Cambridge, MA; [c]Department of Politics and Center for Statistics and Machine Learning, Princeton University, Princeton, NJ

## ABSTRACT

We study causal interaction in factorial experiments, in which several factors, each with multiple levels, are randomized to form a large number of possible treatment combinations. Examples of such experiments include conjoint analysis, which is often used by social scientists to analyze multidimensional preferences in a population. To characterize the structure of causal interaction in factorial experiments, we propose a new causal interaction effect, called the *average marginal interaction effect* (AMIE). Unlike the conventional interaction effect, the relative magnitude of the AMIE does not depend on the choice of baseline conditions, making its interpretation intuitive even for higher-order interactions. We show that the AMIE can be nonparametrically estimated using ANOVA regression with weighted zero-sum constraints. Because the AMIEs are invariant to the choice of baseline conditions, we directly regularize them by collapsing levels and selecting factors within a penalized ANOVA framework. This regularized estimation procedure reduces false discovery rate and further facilitates interpretation. Finally, we apply the proposed methodology to the conjoint analysis of ethnic voting behavior in Africa and find clear patterns of causal interaction between politicians' ethnicity and their prior records. The proposed methodology is implemented in an open source software package. Supplementary materials for this article, including a standardized description of the materials available for reproducing the work, are available as an online supplement.

## 1. Introduction

Statistical interaction among treatment variables can be interpreted as causal relationships when the treatments are randomized in an experiment. Causal interaction plays an essential role in the exploration of heterogenous treatment effects. This article develops a framework for studying causal interaction in randomized experiments with a factorial design, in which there are multiple factorial treatments with each having several levels. A primary goal of causal interaction analysis is to identify the combinations of treatments that induce large additional effects beyond the sum of effects separately attributable to each treatment.

Our motivating application is conjoint analysis, which is a type of randomized survey experiment with a factorial design (Luce and Tukey 1964). Conjoint analysis has been extensively used in marketing research to investigate consumer preferences and predict product sales (e.g., Green, Krieger, and Wind 2001; Marshall and Bradlow 2002). In a typical conjoint analysis, respondents are asked to evaluate pairs of product profiles where several characteristics of a commercial product such as price and color are randomly chosen. Because these product characteristics are represented by factorial variables, conjoint analysis can be seen as an application of randomized factorial design. Thus, the causal estimands and estimation methods proposed in this article are widely applicable to any factorial experiments with many factors.

Recently, conjoint analysis has also gained its popularity among medical and social scientists who study multidimensional preferences among a population of individuals (e.g., Marshall et al. 2010; Hainmueller and Hopkins 2015). In this article, we focus on the latter use of conjoint analysis by estimating population average causal effects. Specifically, we analyze a conjoint analysis about coethnic voting in Africa to examine the conditions under which voters prefer political candidates of the same ethnicity (see Section 2 for the details of the experiment and Section 6 for our empirical analysis).

One important limitation of conjoint analysis, as currently conducted in applied research, is that causal interactions are largely ignored. This is unfortunate because studies of multidimensional choice necessarily involve the consideration of interaction effects. However, the exploration of causal interactions in conjoint analysis is often difficult for two reasons. First, the relative magnitude of the conventional causal interaction effect depends on the choice of baseline condition. This is problematic because many factors used in conjoint analysis do not have natural baseline conditions (e.g., gender, racial group, religion, occupation). Second, a typical conjoint analysis has several factors with each having multiple levels. This means that we must apply a regularization method to reduce false discovery and facilitate interpretation. Yet, the lack of invariance property means that the results of standard regularized estimation will depend on the choice of baseline conditions.

To overcome these problems, we propose an alternative definition of causal interaction effect that is invariant to the choice of baseline condition, making its interpretation intuitive even for higher-order interactions (Sections 3 and 4). We call this new causal quantity of interest, the *average marginal interaction effect* (AMIE), because it marginalizes the other treatments rather than conditioning on their baseline values as done in the conventional causal interaction effect. The proposed approach enables researchers to effectively summarize the structure of causal interaction in high-dimension by decomposing the total effect of any treatment combination into the separate effect of each treatment and their interaction effects.

Finally, we also establish the identification condition and develop estimation strategies for the AMIE (Section 5). We propose a nonparametric estimator of the AMIE and show that this estimator can be recast as an ANOVA with weighted zero-sum constraints (Scheffe 1959). Exploiting this equivalence relationship, we apply the method proposed by Post and Bondell (2013) and directly regularize the AMIEs within the ANOVA framework by collapsing levels and selecting factors. Because the AMIE is invariant to the choice of baseline condition, our regularization also has the same invariance property. This also enables a proper regularization of the conditional average effects, which can be computed using the AMIEs. Without the invariance property, the results of regularized estimation will depend on the choice of baseline conditions. All of our theoretical results and estimation strategies are shown to hold for causal interaction of any order. The proposed methodology is implemented via an open-source software package, *FindIt: Finding Heterogeneous Treatment Effects* (Egami, Ratkovic, and Imai 2017), which is available for download at the Comprehensive R Archive Network (CRAN; *https://cran.r-project.org/package=FindIt*).

Our article builds on the causal inference and experimental design literatures that are concerned about interaction effects (see, e.g., Cox 1984; Jaccard and Turrisi 2003; de González and Cox 2007; VanderWeele and Knol 2014). In addition, we draw upon the recent articles that provide the potential outcomes framework for causal inference with factorial experiments and conjoint analysis (Hainmueller, Hopkins, and Yamamoto 2014; Dasgupta, Pillai, and Rubin 2015; Lu 2016a, 2016b). Indeed, the AMIE is a direct generalization of the average marginal effect studied in this literature that can be used to characterize the causal heterogeneity of a high-dimensional treatment.

Finally, this article is also related to the literature on heterogenous treatment effects, in which the goal of analysis is to find an optimal treatment regime. Much of this literature, however, focuses on the interaction between a single treatment and pretreatment covariates (e.g., Hill 2012; Green and Kern 2012; Wager and Athey 2017; Grimmer, Messing, and Westwood 2017) or a dynamic setting where a sequence of treatment decisions is optimized (e.g., Murphy 2003; Robins 2004). We emphasize that if the goal of analysis is to find an optimal treatment regime, rather than to understand the structure of causal heterogeneity, the marginalized causal quantities such as the one proposed in this article may be of little use. In such settings, researchers typically estimate the causal effects of specific treatment combinations (e.g., Imai and Ratkovic 2013).

## 2. Conjoint Analysis of Ethnic Voting

Conjoint analysis has a long history dating back to the theoretical article by Luce and Tukey (1964). In terms of its application, it has been widely used by marketing researchers over the last 40 years to measure consumer preferences and predict product sales (Green and Rao 1971; Green, Krieger, and Wind 2001; Marshall and Bradlow 2002). It has also become a popular statistical tool in the medical and social sciences (e.g., Marshall et al. 2010; Hainmueller and Hopkins 2015) to study multidimensional preferences of a variety of populations such as patients and voters.

Conjoint analysis can be considered as an application of factorial randomized experiments. For example, in a typical conjoint analysis used for marketing research, respondents evaluate a commercial product whose several characteristics such as price and color, etc., are randomly selected. Factorial variables represent these characteristics with several levels (e.g., \$1, \$5, \$10 for price, and red, green, and blue for color). Similarly, in political science research, conjoint analysis may be used to evaluate candidates where factors may represent their party identification, race, gender, and other attributes.

In this article, we examine a recent conjoint analysis conducted to study coethnic voting in Uganda (Carlson 2015). Coethnic voting refers to the tendency of some voters to prefer political candidates whose ethnicity is the same as their own. Researchers have observed that coethnic voting occurs frequently among African voters, but the identification of causal effects is often difficult because the ethnicity of candidates is often correlated with other characteristics that may influence voting behavior. To address this problem, the original author conducted a conjoint analysis, in which respondents were asked to choose one of the two hypothetical candidates whose attributes were randomly assigned.

For the experiment, a total of 547 respondents were sampled from villages in Uganda. We analyze a subset of 544 observations after removing three observations with missing data. Each respondent was given the description of three pairs of hypothetical presidential candidates. They were then asked to cast a vote for one of the candidates within each pair. These hypothetical candidates are characterized by a total of four factors shown in Table 1: Coethnicity (2 levels), Record (7 levels), Platform (3 levels), and Degree (2 levels).

While the levels of all factors are randomly and independently selected for each hypothetical candidate, the distribution of candidate ethnicity depends on the local ethnic diversity so that enough respondents share the same ethnicity as their assigned hypothetical candidates. The original analysis was based on a mixed effects logistic regression with a respondent random effect. While previous studies showed that many voters unconditionally favor coethnic candidates, Carlson (2015) found that voters tend to favor only coethnic candidates with good prior record.

We focus on two methodological challenges of the original analysis. First, the author tests the existence of causal interaction between Coethnicity and Record, but does not explicitly estimate causal interaction effects. We propose a definition of causal interaction effects in randomized experiments with a factorial design and show how to estimate them. Second, the

**Table 1.** Levels of four factors from the conjoint analysis in Carlson (2015).

| Factors | Levels | |
|---|---|---|
| Coethnicity | Yes | a coethnic of a respondent |
| | No | not a coethnic of a respondent |
| Record | Yes/Village | politician for a village with good prior record |
| | Yes/District | politician for a district with good prior record |
| | Yes/MP | member of parliament with good prior record |
| | No/Village | politician for a village without good prior record |
| | No/District | politician for a district without good prior record |
| | No/MP | member of parliament without good prior record |
| | No/Business | businessman without good prior record |
| Platform | Job | promise to create new jobs |
| | Clinic | promise to create clinics |
| | Education | promise to improve education |
| Degree | Yes | masters degree in business, law, economics, or development |
| | No | bachelors degree in tourism, horticulture, forestry or theater |

author dichotomized two factors, `Record` and `Platform`, which have more than two levels and does not have a natural baseline condition. We show how to use a data-driven regularization method when estimating causal interaction effects in a high-dimensional setting. Our reanalysis of this experiment appears in Section 6.

## 3. Two-Way Causal Interaction

In this section, we introduce a new causal quantity, the *average marginal interaction effect* (AMIE), and show that, unlike the conventional causal interaction effect, it is invariant to the choice of baseline condition. The invariance property enables simple interpretation and effective regularization even when there are many factors. While this section focuses on two-way causal interaction for the sake of simplicity, all definitions and results will be generalized beyond two-way interaction in Section 4.

### 3.1. The Setup

Consider a simple random sample of $n$ units from the target population $\mathcal{P}$. Let $A_i$ and $B_i$ be two factorial treatment variables of interest for unit $i$ where $L_A$ and $L_B$ be the number of ordered or unordered levels for factors $A$ and $B$, respectively. We use $a_\ell$ and $b_m$ to represent levels of the two factors where $\ell = \{0, 1, \ldots, L_A - 1\}$ and $m = \{0, 1, \ldots, L_B - 1\}$. The support of treatment variables $A$ and $B$, therefore, is given by $\mathcal{A} = \{a_0, a_1, \ldots, a_{L_A-1}\}$ and $\mathcal{B} = \{b_0, b_1, \ldots, b_{L_B-1}\}$, respectively.

We call a combination of factor levels $(a_\ell, b_m)$ a *treatment combination*. Thus, in the current set-up, the total number of unique treatment combinations is $L_A \times L_B$. Let $Y_i(a_\ell, b_m)$ denote the potential outcome variable of unit $i$ if the unit receives the treatment combination $(a_\ell, b_m)$. For each unit, only one of the potential outcome variables can be observed, and the realized outcome variable is denoted by $Y_i = \sum_{a_\ell \in \mathcal{A}, b_m \in \mathcal{B}} \mathbf{1}\{A_i = a_\ell, B_i = b_m\} Y_i(a_\ell, b_m)$, where $\mathbf{1}\{A_i = a_\ell, B_i = b_m\}$ is an indicator variable taking the value 1

when $A_i = a_\ell$ and $B_i = b_m$, and taking the value 0 otherwise. In this article, we make the stability assumption, which states that there is neither interference between units nor different versions of the treatment (Cox 1958; Rubin 1990).

In addition, we assume that the treatment assignment is randomized.

$$\{Y_i(a_\ell, b_m)\}_{a_\ell \in \mathcal{A}, b_m \in \mathcal{B}} \perp\!\!\!\perp \{A_i, B_i\} \quad \text{for all } i = 1, \ldots, n. \quad (1)$$

$$\Pr(A_i = a_\ell, B_i = b_m) > 0 \quad \text{for all } a_\ell \in \mathcal{A} \text{ and } b_m \in \mathcal{B}. \quad (2)$$

This assumption rules out the use of fractional factorial designs where certain combinations of treatments have zero probability of occurrence. In some cases, however, researchers may wish to eliminate certain treatment combinations for substantive reasons. The standard recommendation is to set the probability for those treatment combinations to small nonzero values under a full factorial design so that the assumption continues to hold (see Hainmueller, Hopkins, and Yamamoto 2014, footnote 18). Another possibility is to restrict one's analysis to a subset of data and hence the corresponding subset of estimands so that the assumption is satisfied.

Under this set-up, we review two noninteractive causal effects of interest. First, we define the *average combination effect* (ACE), which represents the average causal effect of a treatment combination $(A_i, B_i) = (a_\ell, b_m)$ relative to a prespecified baseline condition $(a_0, b_0)$ (e.g., Dasgupta, Pillai, and Rubin 2015):

$$\tau_{AB}(a_\ell, b_m; a_0, b_0) \equiv \mathbb{E}\{Y_i(a_\ell, b_m) - Y_i(a_0, b_0)\}, \quad (3)$$

where $a_\ell, a_0 \in \mathcal{A}$ and $b_m, b_0 \in \mathcal{B}$.

Another causal quantity of interest is the *average marginal effect* (AME). For each unit, we define the marginal effect of treatment condition $A_i = a_\ell$ relative to a baseline condition $a_0$ by averaging over the distribution of the other treatment $B_i$. Then, the AME is the population average of this unit-level marginal effect (e.g., Hainmueller, Hopkins, and Yamamoto 2014; Dasgupta, Pillai, and Rubin 2015):

$$\psi_A(a_\ell, a_0) \equiv \mathbb{E}\left[\int \{Y_i(a_\ell, B_i) - Y_i(a_0, B_i)\} \, dF(B_i)\right], \quad (4)$$

where $a_\ell, a_0 \in \mathcal{A}$ and $B_i$ is another factor whose distribution function is $F(B_i)$. The AME of $b_m$ relative to $b_0$, that is, $\psi_B(b_m, b_0)$, can be defined similarly.

We emphasize that while these two causal quantities require the specification of baseline conditions, the relative magnitude is not sensitive to this choice. For example, if we sort the ACEs by their relative magnitude, the resulting order does not depend on the values of the treatment variables selected for the baseline conditions $(a_0, b_0)$. The same property is applicable to the AMEs where the choice of baseline condition $a_0$ does not alter their relative magnitude.

### 3.2. The Average Marginal Interaction Effect

We propose a new two-way causal interaction effect, called the *average marginal interaction effect* (AMIE), which is useful for randomized experiments with a factorial design. For each unit, the marginal interaction effect represents the causal effect induced by the treatment combination beyond the sum of the marginal effects separately attributable to each treatment. The

AMIE is the population average of this unit-level marginal interaction effect. Specifically, the two-way AMIE of treatment combination $(a_\ell, b_m)$, with baseline condition $(a_0, b_0)$, is defined as

$$
\begin{aligned}
\pi_{AB}&(a_\ell, b_m; a_0, b_0) \\
&\equiv \mathbb{E}\Bigg[ Y_i(a_\ell, b_m) - Y_i(a_0, b_0) \\
&\quad - \int \{Y_i(a_\ell, B_i) - Y_i(a_0, B_i)\} dF(B_i) \\
&\quad - \int \{Y_i(A_i, b_m) - Y_i(A_i, b_0)\} dF(A_i) \Bigg] \\
&= \tau_{AB}(a_\ell, b_m; a_0, b_0) - \psi_A(a_\ell, a_0) - \psi_B(b_m, b_0),
\end{aligned}
\tag{5}
$$

where $a_\ell, a_0 \in \mathcal{A}$ and $b_m, b_0 \in \mathcal{B}$, $\pi_{AB}(a_\ell, b_m; a_0, b_0)$ is the AMIE, and $\psi(\cdot, \cdot)$ is the AME defined in Equation (4).

The AMIE is closely connected to the conventional definition of the *average interaction effect* (AIE). In the causal inference literature (e.g., Cox 1984; VanderWeele 2015; Dasgupta, Pillai, and Rubin 2015), researchers define the AIE of treatment combination $(a_\ell, b_m)$ relative to baseline condition $(a_0, b_0)$ as,

$$
\begin{aligned}
\xi_{AB}&(a_\ell, b_m; a_0, b_0) \\
&\equiv \mathbb{E}\{Y_i(a_\ell, b_m) - Y_i(a_0, b_m) - Y_i(a_\ell, b_0) + Y_i(a_0, b_0)\},
\end{aligned}
\tag{6}
$$

where $a_\ell, a_0 \in \mathcal{A}$ and $b_m, b_0 \in \mathcal{B}$.

Similar to the AMIE, the AIE has an *interactive effect interpretation*, representing the additional average causal effect induced by the treatment combination beyond the sum of the average causal effects separately attributable to each treatment. This interpretation is based on the following algebraic equality:

$$
\begin{aligned}
\xi_{AB}(a_\ell, b_m; a_0, b_0) &= \tau_{AB}(a_\ell, b_m; a_0, b_0) \\
&\quad - \mathbb{E}\{Y_i(a_\ell, b_0) - Y_i(a_0, b_0)\} \\
&\quad - \mathbb{E}\{Y_i(a_0, b_m) - Y_i(a_0, b_0)\}.
\end{aligned}
\tag{7}
$$

The difference between the AMIE and the AIE is that the former subtracts the AMEs from the ACE while the latter subtracts the sum of two separate effects due to $A_i = a_\ell$ and $B_i = b_m$ while holding the other treatment variable at its baseline value, that is, $A_i = a_0$ or $B_i = b_0$.

In addition, the AIE has a *conditional effect interpretation*,

$$
\begin{aligned}
\xi_{AB}(a_\ell, b_m; a_0, b_0) &= \mathbb{E}\{Y_i(a_\ell, b_m) - Y_i(a_0, b_m)\} \\
&\quad - \mathbb{E}\{Y_i(a_\ell, b_0) - Y_i(a_0, b_0)\},
\end{aligned}
$$

which denotes the difference in the average causal effect of $A_i = a_\ell$ relative to $A_i = a_0$ between the two scenarios, one when $B_i = b_m$ and the other when $B_i = b_0$. When such conditional effects are of interest, the AMIE can be used to obtain them. For example, we have

$$
\mathbb{E}\{Y_i(a_\ell, b_0) - Y_i(a_0, b_0)\} = \psi_A(a_\ell; a_0) + \pi_{AB}(a_\ell, b_0; a_0, b_0).
\tag{8}
$$

Clearly, the scientific question of interest should determine the choice between the AMIE and AIE. In Section 6, we illustrate how to use the AMIEs for estimating the average conditional effects when necessary.

Finally, the AMIE and the AIE are linear functions of one another. This result is presented below as a special case of Theorem 1 presented in Section 4.

*Result 1 (Relationships Between the Two-Way AMIE and the Two-Way AIE).* The two-way average marginal interaction effect (AMIE), defined in Equation (5), equals the following linear function of the two-way average interaction effects (AIEs), defined in Equation (6):

$$
\begin{aligned}
\pi_{AB}(a_\ell, b_m; a_0, b_0) &= \xi_{AB}(a_\ell, b_m; a_0, b_0) \\
&\quad - \sum_{a \in \mathcal{A}} \Pr(A_i = a)\, \xi_{AB}(a, b_m; a_0, b_0) \\
&\quad - \sum_{b \in \mathcal{B}} \Pr(B_i = b)\, \xi_{AB}(a_\ell, b; a_0, b_0).
\end{aligned}
$$

Likewise, the AIE can be expressed as the following linear function of the AMIEs:

$$
\begin{aligned}
\xi_{AB}(a_\ell, b_m; a_0, b_0) &= \pi_{AB}(a_\ell, b_m; a_0, b_0) - \pi_{AB}(a_\ell, b_0; a_0, b_0) \\
&\quad - \pi_{AB}(a_0, b_m; a_0, b_0).
\end{aligned}
$$

Result 1 implies that all the AMIEs are zero if and only if all the AIEs are zero. Thus, testing the absence of causal interaction can be done by an $F$-test, investigating either all the AIEs or all the AMIEs are zero. All causal estimands introduced in this section are identifiable under the assumption of randomized treatment assignment (i.e., Equations (1) and (2)).

### 3.3. Invariance to the Choice of Baseline Condition

One advantage of the AMIE over the AIE is its invariance to the choice of baseline condition. That is, the relative difference of any pair of AMIEs remains unchanged even if one chooses a different baseline condition. Most causal effects, including the ACE and the AME, have this invariance property. In contrast, the relative magnitude of any two AIEs depends on the choice of baseline condition unless all AIEs are zero. The invariance property is important because without it researchers cannot systematically compare interaction effects of different treatment combinations. We state this as Result 2, which is a special case of Theorem 2 presented in Section 5.

*Result 2 (Invariance and Lack Thereof to the Choice of Baseline Condition).* The average marginal interaction effect (AMIE), defined in Equation (5), is interval invariant. That is, for any $(a_\ell, b_m) \neq (a_{\ell'}, b_{m'})$ and $(a_0, b_0) \neq (a_{\tilde{\ell}}, b_{\tilde{m}})$, the following equality holds,

$$
\begin{aligned}
\pi_{AB}&(a_\ell, b_m; a_0, b_0) - \pi_{AB}(a_{\ell'}, b_{m'}; a_0, b_0) \\
&= \pi_{AB}(a_\ell, b_m; a_{\tilde{\ell}}, b_{\tilde{m}}) - \pi_{AB}(a_{\ell'}, b_{m'}; a_{\tilde{\ell}}, b_{\tilde{m}}).
\end{aligned}
$$

Note that the above difference of the AMIEs is also equal to another AMIE, $\pi_{AB}(a_\ell, b_m; a_{\ell'}, b_{m'})$.

In contrast, the average interaction effect (AIE), defined in Equation (6) does not have the invariance property. That is, the following equality does not generally hold,

$$
\begin{aligned}
\xi_{AB}&(a_\ell, b_m; a_0, b_0) - \xi_{AB}(a_{\ell'}, b_{m'}; a_0, b_0) \\
&= \xi_{AB}(a_\ell, b_m; a_{\tilde{\ell}}, b_{\tilde{m}}) - \xi_{AB}(a_{\ell'}, b_{m'}; a_{\tilde{\ell}}, b_{\tilde{m}}).
\end{aligned}
$$

In addition, the AIE is interval invariant if and only if all the AIEs are zero.

The sensitivity of the AIEs to the choice of baseline condition can be further illustrated by the fact that the AIE of any treatment combination pertaining to one of levels in the baseline condition is equal to zero. That is, if $(a_0, b_0)$ is the baseline condition, then $\xi_{AB}(a_0, b_m; a_0, b_0) = \xi_{AB}(a_\ell, b_0; a_0, b_0) = 0$. If the researchers are only interested in the conditional effect interpretation of the AIEs, these zero AIEs are not of interest. However, this restriction is problematic for the interactive effect interpretation especially when no natural baseline condition exists. In such circumstances, zero AIEs make it impossible to explore all relevant causal interaction effects. To the contrary, researchers need not to restrict their quantities of interest when using the AMIE, which can take a nonzero value even when one treatment is set to the baseline condition. For example, the AMIE can be positive if the effect of the second treatment is large when the first treatment is set to its baseline value.

While it is invariant to the choice of baseline condition, the AMIE critically depends on the distribution of treatments, that is, $P(A, B)$. This is because the AMIE is a function of the AMEs, which are themselves obtained by marginalizing out other treatments. This dependency of causal quantities is not new. The potential outcomes framework for $2^k$ factorial experiments introduced by Dasgupta, Pillai, and Rubin (2015), for example, defines causal estimands based on the uniform distribution of treatments. Many applied researchers independently randomize multiple treatments and then estimate the AME of each treatment by simply ignoring the other treatments. This estimation procedure implicitly conditions on the empirical distribution of treatment assignments.

Although the uniform or empirical distribution would be a reasonable default choice for many experimentalists, researchers can improve the external validity of their experiment by using a treatment distribution based on the target population (Hainmueller, Hopkins, and Yamamoto 2014). This is important for the conjoint analysis, in which treatments are often characteristics of people. In our empirical application (see Section 2), for example, researchers could obtain the detailed information about the attributes of actual candidates and use it as the basis of treatment distribution.

## 4. Generalization to Higher-Order Interaction

In this section, we generalize the two-way AMIE introduced in Section 3 to higher-order causal interaction with more than two factors. We prove that a higher-order AMIE retains the same desirable properties and intuitive interpretation.

### 4.1. The Setup

Suppose that we have a total of $J$ factorial treatments denoted by an vector $\mathbf{T}_i = (T_{i1}, T_{i2}, \ldots, T_{iJ})$ where $J \geq 2$ and each factor $T_{ij}$ has a total of $L_j$ levels. Without loss of generality, let $\mathbf{T}_i^{1:K}$ be a subset of $K$ treatments of interest where $K \leq J$ whereas $\mathbf{T}_i^{(K+1):J}$ denotes the remaining $(J - K)$ factorial treatment variables, which are not of interest. As before, we assume that the treatment assignment is randomized.

*Assumption 1 (Randomized Treatment Assignment).*

$$Y_i(\mathbf{t}) \perp\!\!\!\perp \mathbf{T}_i \quad \text{and} \quad \Pr(\mathbf{T}_i = \mathbf{t}) > 0 \quad \text{for all } \mathbf{t}.$$

In addition, we assume that $J$ factorial treatments are independent of one another.

*Assumption 2 (Independent Treatment Assignment).*

$$T_{ij} \perp\!\!\!\perp \mathbf{T}_{i,-j} \quad \text{for all } j \in \{1, 2, \ldots, J\},$$

where $\mathbf{T}_{i,-j}$ denotes the $(J - 1)$ factorial treatments excluding $T_{ij}$.

Assumption 2 is not required for some of the results obtained below, but it considerably simplifies the notation.

We now generalize the definition of the two-way ACE given in Equation (3) by accommodating more than two factorial treatments of interest $\mathbf{T}_i^{1:K}$ while allowing for the existence of additional treatments $\mathbf{T}_i^{(K+1):J}$, which are marginalized out.

*Definition 1 (The K-Way Average Combination Effect).* The $K$-way average combination effect (ACE) of treatment combination $\mathbf{T}_i^{1:K} = \mathbf{t}^{1:K}$ relative to baseline condition $\mathbf{T}_i^{1:K} = \mathbf{t}_0^{1:K}$ is defined as,

$$\tau_{1:K}(\mathbf{t}^{1:K}; \mathbf{t}_0^{1:K}) \equiv \mathbb{E}\left[ \int \left\{ Y_i(\mathbf{T}_i^{1:K} = \mathbf{t}^{1:K}, \mathbf{T}_i^{(K+1):J}) \right. \right. $$
$$\left. \left. - Y_i(\mathbf{T}_i^{1:K} = \mathbf{t}_0^{1:K}, \mathbf{T}_i^{(K+1):J}) \right\} dF(\mathbf{T}_i^{(K+1):J}) \right].$$

The generalization of the AME defined in Equation (4) to this setting is straightforward. For example, the AME of $T_{i1}$ is obtained by marginalizing the remaining factors $\mathbf{T}_i^{2:J}$ out.

### 4.2. The K-Way Average Marginal Interaction Effect

We now extend the definition of the two-way AMIE, given in Equation (5), to higher-order causal interaction and discuss its relationships with the conventional higher-order causal interaction effect. We define the $K$-way AMIE as the additional effect of treatment combination beyond the sum of all lower-order AMIEs.

*Definition 2 (The K-Way Average Marginal Interaction Effect).* The $K$-way average marginal interaction effect (AMIE) of treatment combination $\mathbf{T}_i^{1:K} = \mathbf{t}^{1:K}$, relative to baseline condition, $\mathbf{T}_i^{1:K} = \mathbf{t}_0^{1:K}$, is given by,

$$\pi_{1:K}(\mathbf{t}^{1:K}; \mathbf{t}_0^{1:K}) \equiv \mathbb{E}\left\{ \tau_{1:K}^{(i)}(\mathbf{t}^{1:K}; \mathbf{t}_0^{1:K}) - \sum_{k=1}^{K-1} \sum_{\mathcal{K}_k \subseteq \mathcal{K}_K} \pi_{\mathcal{K}_k}^{(i)}(\mathbf{t}^{\mathcal{K}_k}; \mathbf{t}_0^{\mathcal{K}_k}) \right\}$$
$$= \tau_{1:K}(\mathbf{t}^{1:K}; \mathbf{t}_0^{1:K}) - \sum_{k=1}^{K-1} \sum_{\mathcal{K}_k \subseteq \mathcal{K}_K} \pi_{\mathcal{K}_k}(\mathbf{t}^{\mathcal{K}_k}; \mathbf{t}_0^{\mathcal{K}_k}),$$

where $\mathcal{K}_k \subseteq \mathcal{K}_K = \{1, \ldots, K\}$ such that $|\mathcal{K}_k| = k$ with $k = 1, \ldots, K$, $\tau_{1:K}^{(i)}(\mathbf{t}^{1:K}; \mathbf{t}_0^{1:K})$ is the unit-level combination effect, and $\pi_{1:K}^{(i)}(\mathbf{t}^{1:K}; \mathbf{t}_0^{1:K})$ is the unit-level $K$-way marginal interaction effect.

This definition reduces to Equation (5) when $K = 2$ because the one-way AMIE is equal to the AME, that is, $\pi_1(t; t_0) = \psi_1(t, t_0)$.

As in the two-way case, the $K$-way AMIE is closely related to the $K$-way AIE. To generalize the two-way AIE given in Equation (6), we first define the two-way AIE of treatment combination $\mathbf{t}^{1:2} = (t_1, t_2)$, relative to baseline condition $\mathbf{t}_0^{1:2} = (t_{01}, t_{02})$ by marginalizing the remaining treatments $\mathbf{T}^{3:J}$. The unit-level two-way interaction effect and the two-way AIE are defined as

$$\xi_{1:2}(\mathbf{t}^{1:2}; \mathbf{t}_0^{1:2}) \equiv \mathbb{E}\left[ \int \left\{ Y_i(t_1, t_2, \mathbf{T}_i^{3:J}) - Y_i(t_{01}, t_2, \mathbf{T}_i^{3:J}) \right. \right.$$
$$\left. \left. - Y_i(t_1, t_{02}, \mathbf{T}_i^{3:J}) + Y_i(t_{01}, t_{02}, \mathbf{T}_i^{3:J}) \right\} dF\left(\mathbf{T}_i^{3:J}\right) \right].$$

In addition, define the *conditional* two-way AIE by fixing the level of another treatment $T_{i3}$ at $t^*$.

$$\xi_{1:2}(\mathbf{t}^{1:2}; \mathbf{t}_0^{1:2} \mid T_{i3} = t^*)$$
$$\equiv \mathbb{E}\left[ \int \{ Y_i(t_1, t_2, t^*, \mathbf{T}_i^{4:J}) - Y_i(t_{01}, t_2, t^*, \mathbf{T}_i^{4:J}) \right.$$
$$\left. - Y_i(t_1, t_{02}, t^*, \mathbf{T}_i^{4:J}) + Y_i(t_{01}, t_{02}, t^*, \mathbf{T}_i^{4:J}) \} dF(\mathbf{T}_i^{4:J}) \right].$$

Then, the three-way AIE can be defined as the difference between the ACE of treatment combination $\mathbf{t}^{1:3} = (t_1, t_2, t_3)$ and the sum of all conditional two-way and one-way AIEs while conditioning on the baseline condition $\mathbf{t}_0^{1:3} = (t_{01}, t_{02}, t_{03})$,

$$\xi_{1:3}(\mathbf{t}^{1:3}; \mathbf{t}_0^{1:3})$$
$$= \tau_{1:3}(\mathbf{t}^{1:3}; \mathbf{t}_0^{1:3}) - \left\{ \xi_{1:2}(\mathbf{t}^{1:2}; \mathbf{t}_0^{1:2} \mid T_{i3} = t_{03}) \right.$$
$$+ \xi_{2:3}(\mathbf{t}^{2:3}; \mathbf{t}_0^{2:3} \mid T_{i1} = t_{01}) + \xi_{1,3}(\mathbf{t}^{1,3}; \mathbf{t}_0^{1,3} \mid T_{i2} = t_{02}) \}$$
$$- \left\{ \xi_1(t_1; t_{01} \mid \mathbf{T}_i^{2:3} = \mathbf{t}_0^{2:3}) + \xi_2(t_2; t_{02} \mid \mathbf{T}_1^{1,3} = \mathbf{t}_0^{1,3}) \right.$$
$$\left. + \xi_3(t_3; t_{03} \mid \mathbf{T}_i^{1:2} = \mathbf{t}_0^{1:2}) \right\}. \tag{9}$$

Note that the one-way conditional AIEs are equivalent to the average effects of single treatments while holding the other treatments at their base level. For example, $\xi_1(t_1; t_{01} \mid \mathbf{T}_i^{2:3} = \mathbf{t}_0^{2:3})$ is equal to $\tau_{1:3}(t_1, \mathbf{t}_0^{2:3}; \mathbf{t}_0)$. We also note that $\xi_1(t_1; t_{01}) = \psi_1(t_1; t_{01}) = \pi_1(t_1; t_{01})$ holds. In this way, we can generalize the AIE to higher-order causal interaction.

*Definition 3 (The K-Way Average Interaction Effect).* The $K$-way average interaction effect (AIE) of treatment combination $\mathbf{T}_i^{1:K} = \mathbf{t}^{1:K} = (t_1, \ldots, t_K)$ relative to baseline condition $\mathbf{T}_i^{1:K} = \mathbf{t}_0^{1:K} = (t_{01}, \ldots, t_{0K})$ is given by,

$$\xi_{1:K}(\mathbf{t}^{1:K}; \mathbf{t}_0^{1:K})$$
$$\equiv \mathbb{E}\left\{ \tau_{1:K}^{(i)}(\mathbf{t}^{1:K}; \mathbf{t}_0^{1:K}) - \sum_{k=1}^{K-1} \sum_{\mathcal{K}_k \subseteq \mathcal{K}_K} \xi_{\mathcal{K}_k}^{(i)}(\mathbf{t}^{\mathcal{K}_k}; \mathbf{t}_0^{\mathcal{K}_k} \mid \mathbf{T}_i^{\mathcal{K}_K \backslash \mathcal{K}_k} = \mathbf{t}_0^{\mathcal{K}_K \backslash \mathcal{K}_k}) \right\}$$
$$= \tau_{1:K}(\mathbf{t}^{1:K}; \mathbf{t}_0^{1:K}) - \sum_{k=1}^{K-1} \sum_{\mathcal{K}_k \subseteq \mathcal{K}_K} \xi_{\mathcal{K}_k}(\mathbf{t}^{\mathcal{K}_k}; \mathbf{t}_0^{\mathcal{K}_k} \mid \mathbf{T}_i^{\mathcal{K}_K \backslash \mathcal{K}_k} = \mathbf{t}_0^{\mathcal{K}_K \backslash \mathcal{K}_k}),$$

where the second summation is taken over the set of all possible $\mathcal{K}_k \subseteq \mathcal{K}_K = \{1, 2, \ldots, K\}$ such that $|\mathcal{K}_k| = k$, $\tau_{1:K}^{(i)}(\mathbf{t}^{1:K}; \mathbf{t}_0^{1:K})$ is the unit-level combination effect, and $\xi_{\mathcal{K}_k}^{(i)}(\mathbf{t}^{\mathcal{K}_k}; \mathbf{t}_0^{\mathcal{K}_k} \mid \mathbf{T}_i^{\mathcal{K}_K \backslash \mathcal{K}_k} = \mathbf{t}_0^{\mathcal{K}_K \backslash \mathcal{K}_k})$ represents the unit-level interaction effect.

While both estimands have similar interpretations, the $K$-way AMIE differs from the $K$-way AIE in important ways. First, the AMIE is expressed as a function of its lower-order effects whereas the AIE is based on the lower-order *conditional* AIEs rather than the lower-order AIEs. This implies that we can

decompose the $K$-way ACE as the sum of the $K$-way AMIE and all lower-order AMIEs.

$$\tau_{1:K}(\mathbf{t}^{1:K}; \mathbf{t}_0^{1:K}) = \sum_{k=1}^{K} \sum_{\mathcal{K}_k \subseteq \mathcal{K}_K} \pi_{\mathcal{K}_k}(\mathbf{t}^{\mathcal{K}_k}; \mathbf{t}_0^{\mathcal{K}_k}). \tag{10}$$

The decomposition is useful for understanding how interaction effects of various order relate to the overall effect of treatment combination. However, because of conditioning on the baseline value, a similar decomposition is not applicable to the AIEs.

Second, in the experimental design literature, the $K$-way AIE is often interpreted as a conditional interaction effect (see, e.g., Jaccard and Turrisi 2003; Wu and Hamada 2011). For example, the three-way AIE of treatment combination $\mathbf{T}_i^{1:3} = \mathbf{t}^{1:3} = (t_1, t_2, t_3)$ relative to baseline condition $\mathbf{T}_i^{1:3} = \mathbf{t}_0^{1:3} = (t_{01}, t_{02}, t_{03})$, given in Equation (9), can be rewritten as the difference in the conditional two-way AIEs where the third factorial treatment is either set to $t_3$ or $t_{03}$,

$$\xi_{1:3}(\mathbf{t}^{1:3}; \mathbf{t}_0^{1:3}) = \xi_{1:2}(\mathbf{t}^{1:2}; \mathbf{t}_0^{1:2} \mid T_{i3} = t_3) - \xi_{1:2}(\mathbf{t}^{1:2}; \mathbf{t}_0^{1:2} \mid T_{i3} = t_{03}).$$

Lemma 1 shows that this equivalence relationship can be generalized to the $K$-way AIE (see Appendix A.1).

Unfortunately, as recognized by others (see, e.g., Wu and Hamada 2011, p. 112), although it is useful when $K = 2$, this conditional interpretation faces difficulty when $K$ is greater than two. For example, the three-way AIE has the conditional effect interpretation, characterizing how the conditional two-way AIE varies as a function of the third factorial treatment. However, according to this interpretation, the two-way AIE, which varies according to the second treatment of interest, itself describes how the main effect of one treatment changes as a function of another treatment. This means that the three-way AIE is the conditional effect of another conditional effect, making it difficult for applied researchers to gain an intuitive understanding.

Finally, as in the two-way case, we can express the $K$-way AMIE and $K$-way AIE as linear functions of one another. The next theorem summarizes this result.

*Theorem 1 (Relationships Between the K-Way AMIE and the K-Way AIE).* Under Assumption 2, the $K$-way average marginal interaction effect (AMIE), given in Definition 2, equals the following linear function of the $K$-way average interaction effects (AIEs), given in Definition 3. That is, for any $\mathbf{t}^{1:K}$ and $\mathbf{t}_0^{1:K}$, we have

$$\pi_{1:K}(\mathbf{t}^{1:K}; \mathbf{t}_0^{1:K}) = \xi_{1:K}(\mathbf{t}^{1:K}; \mathbf{t}_0^{1:K})$$
$$+ \sum_{k=1}^{K-1} (-1)^k \sum_{\mathcal{K}_k \subseteq \mathcal{K}_K} \int \xi_{\mathcal{K}_k}(\mathbf{T}^{\mathcal{K}_k}, \mathbf{t}^{\mathcal{K}_K \backslash \mathcal{K}_k}; \mathbf{t}_0^{\mathcal{K}_K}) dF(\mathbf{T}^{\mathcal{K}_k}),$$

where $\mathcal{K}_k \subseteq \mathcal{K}_K = \{1, \ldots, K\}$ such that $|\mathcal{K}_k| = k$ with $k = 1, \ldots, K$. Likewise, but without requiring Assumption 2, the $K$-way AIE can be written as the following linear function of the $K$-way AMIEs:

$$\xi_{1:K}(\mathbf{t}^{1:K}; \mathbf{t}_0^{1:K}) = \sum_{k=1}^{K} (-1)^{K-k} \sum_{\mathcal{K}_k \subseteq \mathcal{K}_K} \pi_{\mathcal{K}_k}(\mathbf{t}^{\mathcal{K}_k}, \mathbf{t}_0^{\mathcal{K}_K \backslash \mathcal{K}_k}; \mathbf{t}_0^{\mathcal{K}_k}, \mathbf{t}_0^{\mathcal{K}_K \backslash \mathcal{K}_k}).$$

Proof is in Appendix A.2. All causal estimands introduced above are identifiable under Assumption 1. We propose nonparametric unbiased estimators in Section 5.

### 4.3. Invariance to the Choice of Baseline Condition

As is the case for the two-way AMIE, the $K$-way AMIE is invariant to the choice of baseline condition. In contrast, the $K$-way AIEs lack this invariance property. The next theorem generalizes Result 2 to the $K$-way causal interaction.

*Theorem 2 (Invariance and Lack Thereof to the Choice of Baseline Condition).* The $K$-way average marginal interaction effect (AMIE), given in Definition 2, is interval invariant. That is, for any treatment combination $\mathbf{t}^{1:K} \neq \tilde{\mathbf{t}}^{1:K}$ and control condition $\mathbf{t}_0^{1:K} \neq \tilde{\mathbf{t}}_0^{1:K}$, the following equality holds,

$$
\begin{aligned}
&\pi_{1:K}(\mathbf{t}^{1:K}; \mathbf{t}_0^{1:K}) - \pi_{1:K}(\tilde{\mathbf{t}}^{1:K}; \mathbf{t}_0^{1:K}) \\
&= \pi_{1:K}(\mathbf{t}^{1:K}; \tilde{\mathbf{t}}_0^{1:K}) - \pi_{1:K}(\tilde{\mathbf{t}}^{1:K}; \tilde{\mathbf{t}}_0^{1:K}).
\end{aligned}
$$

In contrast, the average interaction effect (AIE), given in Definition 3 does not possess the invariance property. That is, the following equality does not generally hold,

$$
\begin{aligned}
&\xi_{\mathcal{K}_K}(\mathbf{t}^{\mathcal{K}_K}; \mathbf{t}_0^{\mathcal{K}_K}) - \xi_{\mathcal{K}_K}(\tilde{\mathbf{t}}^{\mathcal{K}_K}; \mathbf{t}_0^{\mathcal{K}_K}) \\
&= \xi_{\mathcal{K}_K}(\mathbf{t}^{\mathcal{K}_K}; \tilde{\mathbf{t}}_0^{\mathcal{K}_K}) - \xi_{\mathcal{K}_K}(\tilde{\mathbf{t}}^{\mathcal{K}_K}; \tilde{\mathbf{t}}_0^{\mathcal{K}_K}).
\end{aligned} \tag{11}
$$

Proof is in Appendix A.3.

## 5. Estimation and Regularization

In this section, we show how to estimate the AMIE using the general notation introduced in Section 4. For the sake of simplicity, our discussion focuses on the two-way AMIE but we show that all the results presented here can be generalized to the $K$-way AMIE. We first introduce nonparametric estimators based on difference in sample means. We then prove that the AMIE can also be nonparametrically estimated using ANOVA with weighted zero-sum constraints (Scheffe 1959).

While ANOVA is mainly used for a balanced design, our approach is applicable to the unbalanced design as well so long as Assumptions 1 and 2 hold. Finally, we show how to directly regularize the AMIEs by collapsing levels and selecting factors (Post and Bondell 2013). Because of the invariance property of the AMIEs, this regularization method is also invariant to the choice of baseline condition. The proposed method reduces false discovery and facilitates interpretation when there are many factors and levels.

### 5.1. Difference-in-Means Estimators

In the causal inference literature, the following difference-in-means estimators have been used to nonparametrically estimate the ACE and AME (e.g., Hainmueller, Hopkins, and Yamamoto 2014; Dasgupta, Pillai, and Rubin 2015):

$$
\begin{aligned}
\hat{\tau}_{jj'}(\ell, m; 0, 0) &= \frac{\sum_{i=1}^{n} Y_i \mathbf{1}\{T_{ij} = \ell, T_{ij'} = m\}}{\sum_{i=1}^{n} \mathbf{1}\{T_{ij} = \ell, T_{ij'} = m\}} \\
&\quad - \frac{\sum_{i=1}^{n} Y_i \mathbf{1}\{T_{ij} = 0, T_{ij'} = 0\}}{\sum_{i=1}^{n} \mathbf{1}\{T_{ij} = 0, T_{ij'} = 0\}}, \\
\hat{\psi}_j(\ell; 0) &= \frac{\sum_{i=1}^{n} Y_i \mathbf{1}\{T_{ij} = \ell\}}{\sum_{i=1}^{n} \mathbf{1}\{T_{ij} = \ell\}} - \frac{\sum_{i=1}^{n} Y_i \mathbf{1}\{T_{ij} = 0\}}{\sum_{i=1}^{n} \mathbf{1}\{T_{ij} = 0\}}.
\end{aligned}
$$

These estimators are unbiased only when the treatment assignment distribution of an experimental study is used to define the AMEs and AMIEs. Then, Definition 2 naturally implies the following nonparametric estimator of the two-way AMIE:

$$
\hat{\pi}_{jj'}(\ell, m; 0, 0) = \hat{\tau}_{jj'}(\ell, m; 0, 0) - \hat{\psi}_j(\ell; 0) - \hat{\psi}_{j'}(m; 0).
$$

Similarly, the nonparametric estimator of higher-order AMIE can be constructed. It is important to emphasize that these nonparametric estimators do not assume the absence of higher-order interactions (Hainmueller, Hopkins, and Yamamoto 2014).

### 5.2. Nonparametric Estimation with ANOVA

Alternatively, the AMIEs can be estimated nonparametrically using ANOVA with weighted zero-sum constraints, which is a convex optimization problem (Scheffe 1959). For example, the two-way AMIE considered above can be estimated by the saturated ANOVA whose objective function is as follows,

$$
\begin{aligned}
\sum_{i=1}^{n} \Bigg( &Y_i - \mu - \sum_{j=1}^{J} \sum_{\ell=0}^{L_j - 1} \beta_\ell^j \mathbf{1}\{T_{ij} = \ell\} \\
&- \sum_{j=1}^{J-1} \sum_{j' > j} \sum_{\ell=0}^{L_j - 1} \sum_{m=0}^{L_{j'} - 1} \beta_{\ell m}^{jj'} \mathbf{1}\{T_{ij} = \ell, T_{ij'} = m\} \\
&- \sum_{k=3}^{J} \sum_{\mathcal{K}_k \subset \mathcal{K}_J} \sum_{\mathbf{t}^{\mathcal{K}_k}} \beta_{\mathbf{t}^{\mathcal{K}_k}}^{\mathcal{K}_k} \mathbf{1}\{\mathbf{T}_i^{\mathcal{K}_k} = \mathbf{t}^{\mathcal{K}_k}\} \Bigg)^2,
\end{aligned} \tag{12}
$$

where $\mu$ is the global mean, $\beta_\ell^j$ is the coefficient for the first-order term for the $j$th factor with $\ell$ level, $\beta_{\ell m}^{jj'}$ is the coefficient for the second-order interaction term for the $j$th and $j'$th factors with $\ell$ and $m$ levels, respectively, and more generally $\beta_{\mathbf{t}^{\mathcal{K}_k}}^{\mathcal{K}_k}$ is the coefficient for the interaction term for a set of $k$ factors $\mathcal{K}_k$ when their levels are equal to $\mathbf{t}^{\mathcal{K}_k}$. Note that as in Section 4, we have $|\mathcal{K}_k| = k$ and $\mathcal{K}_J = \{1, 2, \ldots, J\}$. We emphasize that the nonparametric estimation requires all interaction terms up to $J$-way interaction. See Section 5.3 for efficient parametric estimation.

We minimize the objective function given in Equation (12) subject to the following weighted zero-sum constraints where the weights are given by the marginal distribution of treatment assignment,

$$
\sum_{\ell=0}^{L_j - 1} \Pr(T_{ij} = \ell) \beta_\ell^j = 0 \quad \text{for all } j, \tag{13}
$$

$$
\sum_{\ell=0}^{L_j - 1} \Pr(T_{ij} = \ell) \beta_{\ell m}^{jj'} = 0 \quad \text{for all } j \neq j'
$$

$$
\text{and } m \in \{0, 1, \ldots, L_{j'} - 1\}, \tag{14}
$$

$$
\sum_{\ell=0}^{L_j - 1} \Pr(T_{ij} = \ell) \mathbf{1}\{t_j = \ell\} \beta_{\mathbf{t}^{\mathcal{K}_k}}^{\mathcal{K}_k} = 0 \quad \text{for all } j, \mathbf{t}^{\mathcal{K}_k},
$$

$$
\text{and } \mathcal{K}_k \subset \mathcal{K}_J \text{ such that } k \geq 3 \text{ and } j \in \mathcal{K}_k. \tag{15}
$$

Finally, the next theorem shows that the difference in the estimated ANOVA coefficients represents a nonparametric estimate of the AMIE.

*Theorem 3 (Nonparametric Estimation with ANOVA).* Under Assumptions 1 and 2, differences in the estimated coefficients from ANOVA based on Equations (12)–(15) represent nonparametric unbiased estimators of the AME and the AMIE:

$$\mathbb{E}(\hat{\beta}_\ell^j - \hat{\beta}_0^j) = \psi_j(\ell; 0), \quad \mathbb{E}(\hat{\beta}_{\ell m}^{jj'} - \hat{\beta}_{00}^{jj'}) = \pi_{jj'}(\ell, m; 0, 0),$$
$$\mathbb{E}(\hat{\beta}_{\mathbf{t}^{\mathcal{K}_k}}^{\mathcal{K}_k} - \hat{\beta}_{\mathbf{t}_0^{\mathcal{K}_k}}^{\mathcal{K}_k}) = \pi_{\mathcal{K}_k}(\mathbf{t}^{\mathcal{K}_k}; \mathbf{t}_0^{\mathcal{K}_k}).$$

Proof is given in Appendix A.4. These estimators are asymptotically equivalent to their corresponding difference-in-means estimators when the treatment assignment distribution of an experimental study is used as weights. The proposed ANOVA framework, however, allows researchers to use any treatment assignment distributions to define the AME and the AMIE so long as Assumptions 1 and 2 hold.

### 5.3. Regularization

A key advantage of this ANOVA-based estimator in Section 5.2 over the difference-in-means estimator in Section 5.1 is that we can directly regularize the AMIEs in a penalized regression framework. The regularization is especially useful for reducing false positives and facilitating interpretation when the number of factors is large.

We apply the regularization method (Grouping and Selection using Heredity in ANOVA or GASH-ANOVA) proposed by Post and Bondell (2013), which places penalties on difference in coefficients of the ANOVA regression. As shown above, these differences correspond to the AMEs and AMIEs. While there exist other regularization methods for categorical variables (e.g., Yuan and Lin 2006; Meier, Van De Geer, and Bühlmann 2008; Zhao, Rocha, and Yu 2009; Huang et al. 2009; Huang, Breheny, and Ma 2012; Lim and Hastie 2015), these methods regularize coefficients rather than their differences. In addition, GASH-ANOVA collapses levels and selects factors by jointly considering the AMEs and AMIEs rather than the AMEs alone. This is attractive because many social scientists believe large interaction effects can exist even when marginal effects are small. The method also collapses levels in a mutually consistent manner.

Finally, because the AMEs and AMIEs are invariant to the choice of baseline condition, this regularization method also inherits the invariance property, which is not generally the case (Lim and Hastie 2015). In particular, even if one is interested in conditional average causal effects, regularization should be based on the AMEs and AMIEs because of their invariance property. As shown in Equation (8), we can compute the conditional average effects directly from these quantities.

To illustrate the application of GASH-ANOVA, consider a situation of practical interest in which we assume the absence of causal interaction higher than the second order. That is, in Equation (12), we assume $\beta_{\mathbf{t}^{\mathcal{K}_k}}^{\mathcal{K}_k} = 0$ for all $k \geq 3$. GASH-ANOVA collapses two levels within a factor by directly and jointly regularizing the AMEs and AMIEs that involve those two levels. Define the set of all the AMEs and AMIEs that involve levels $\ell$ and $\ell'$ of the $j$th factor as follows,

$$\phi^j(\ell, \ell') = \left\{ \left| \beta_\ell^j - \beta_{\ell'}^j \right| \right\} \bigcup \left\{ \bigcup_{j' \neq j} \bigcup_{m=0}^{L_{j'}-1} \left| \beta_{\ell m}^{jj'} - \beta_{\ell' m}^{jj'} \right| \right\}.$$

Finally, the penalty is given by,

$$\sum_{j=1}^{J} \sum_{\ell, \ell'} w_{\ell \ell'}^j \max\{\phi^j(\ell, \ell')\} \leq c,$$

where $c$ is the cost parameter and $w_{\ell \ell'}^j$ is the adaptive weight of the following form,

$$w_{\ell \ell'}^j = \left[ (L_j + 1)\sqrt{L_j} \max\{\bar{\phi}^j(\ell, \ell')\} \right]^{-1},$$

where $(L_j + 1)\sqrt{L_j}$ is the standardization factor (Bondell and Reich 2009), and $\bar{\phi}^j(\ell, \ell')$ represents the corresponding set of all AMEs and AMIEs estimated without regularization. Post and Bondell (2013) showed that, when combined with Equations (12)–(15), the resulting optimization problem is a quadratic programming problem. They also prove that the method has the oracle property.

## 6. Empirical Analysis

We apply the proposed method to the conjoint analysis of coethnic voting described in Section 2. Although conjoint analysis is based on the randomization of multiple factors, it differs from factorial experiments in that respondents evaluate pairs of randomly selected profiles. Thus, we only observe which profile they prefer within a given pair but do not know how much they like each profile. As shown below, this particular feature of conjoint analysis leads to a modified formulation of ANOVA model. As explained in Section 5, we can apply the standard ANOVA (possibly with regularization) to estimate the AMEs and AMIEs in a typical factorial experiment. Our analysis finds clear patterns of causal interaction between the `Record` and `Coethnicity` variables as well as between the `Record` and `Platform` variables.

### 6.1. A Statistical Model of Preference Differentials

Our empirical application is based on the choice-based conjoint analysis, in which respondents are asked to evaluate three pairs of hypothetical presidential candidates in turn. Let $Y_i(\mathbf{t})$ be the potential preference by respondent $i$ for a hypothetical candidate characterized by a vector of attributes $\mathbf{t}$. In this experiment, $\mathbf{t}$ is a four-dimensional vector, based on the values of factorial treatments shown in Table 1 where each factor $T_{ij}$ has $L_j$ levels (i.e., {Coethnicity, Record, Platform, Degree}).

Given the limited sample size, we assume the absence of three-way or higher-order causal interaction and use the following ANOVA regression model of potential outcomes with all one-way effects and two-way interactions:

$$Y_i(\mathbf{t}) = \mu + \sum_{j=1}^{4} \sum_{\ell=0}^{L_j-1} \beta_\ell^j \mathbf{1}\{t_{ij} = \ell\}$$
$$+ \sum_{j=1}^{3} \sum_{j' > j} \sum_{\ell=0}^{L_j-1} \sum_{m=0}^{L_{j'}-1} \beta_{\ell m}^{jj'} \mathbf{1}\{t_{ij} = \ell, t_{ij'} = m\} + \epsilon_i(\mathbf{t}).$$

(16)

The results in Section 5.2 imply that the coefficients in this model represent the AIEs and AMIEs.

In this conjoint analysis, respondents evaluate a pair of hypothetical candidates with different attributes. This means that

we only observe whether respondent $i$ prefers a candidate with attributes $\mathbf{T}_i^*$ over another candidate with attributes $\mathbf{T}_i^\dagger$. Thus, based on the model of preference given in Equation (16), we construct a linear probability model of preference differential,

$$\Pr(Y_i(\mathbf{T}_i^*) > Y_i(\mathbf{T}_i^\dagger) \mid \mathbf{T}_i^*, \mathbf{T}_i^\dagger)$$

$$= \tilde{\mu} + \sum_{j=1}^{4} \sum_{\ell=0}^{L_j-1} \beta_\ell^j (\mathbf{1}\{T_{ij}^* = \ell\} - \mathbf{1}\{T_{ij}^\dagger = \ell\})$$

$$+ \sum_{j,j'} \sum_{\ell,m} \beta_{\ell m}^{jj'} (\mathbf{1}\{T_{ij}^* = \ell, T_{ij'}^* = m\} - \mathbf{1}\{T_{ij}^\dagger = \ell, T_{ij'}^\dagger = m\}),$$

where $\tilde{\mu} = 0.5$ if a position within a pair does not matter. Note that the independence of irrelevant alternatives is assumed. If we additionally assume the difference in errors follow independent Type I extreme value distributions, the model becomes the conditional logit model, which is popular in conjoint analysis (McFadden 1974).

We minimize the sum of squared residuals, subject to the constraints given in Equations (13) and (14) where $\Pr(T_{ij} = \ell)$ represents the marginal distribution of $T_{ij}^*$ and $T_{ij}^\dagger$ together. We also apply the regularization method discussed in Section 5.3. To be consistent with the original dummy coding, we treat Record and Platform as ordered categorical variables and place penalties on the differences between adjacent levels rather than the differences based on every pairwise comparison. We use the order of levels as shown in Table 1. We choose the uniform distribution for treatment assignment and select the value of the cost parameter $c$ based on the minimum mean squared error criterion in 10-fold cross-validation.

Since the inference for a regularization method that collapses levels of factorial variables is not established in the literature (Bühlmann and Dezeure 2016), we focus on the stability of selection (e.g., Breiman 1996; Meinshausen and Bühlmann 2010). In particular, we estimate the selection probability for each AME and AMIE using one minus the proportion of 5000 bootstrap replicates in which all coefficients for the corresponding factor or factor interaction are estimated to be zero (Efron 2014; Hastie, Tibshirani, and Wainwright 2015). Although we do not control the family wise error rate, we follow Meinshausen and Bühlmann (2010) and use 90% cutoff as our default.

Another possible inferential approach is sample splitting where we collapse levels and select factors using training data and then estimate and compute confidence intervals for the AMEs and AMIEs using test data (Wasserman and Roeder 2009; Athey and Imbens 2016; Chernozhukov et al. 2018). Although we do not present the results based on this approach here, it can be implemented through our open-source software package, FindIt.

### 6.2. Findings

We begin by reporting the ranges of the estimated AMEs and AMIEs and their selection probability to determine significant factors and factor interactions, respectively. As shown in Table 2, three factors—Record, Platform, and Coethnicity— are found to be significant factors whereas Degree is not. In terms of the AMIEs, the interaction Coethnicity ×

**Table 2.** Ranges of the estimated average marginal effects (AMEs) and estimated average marginal interaction effects (AMIEs). The estimated selection probability of the AME (AMIE) is one minus the proportion of 5000 bootstrap replicates in which all coefficients for the corresponding factor (factor interaction) are estimated to be zero.

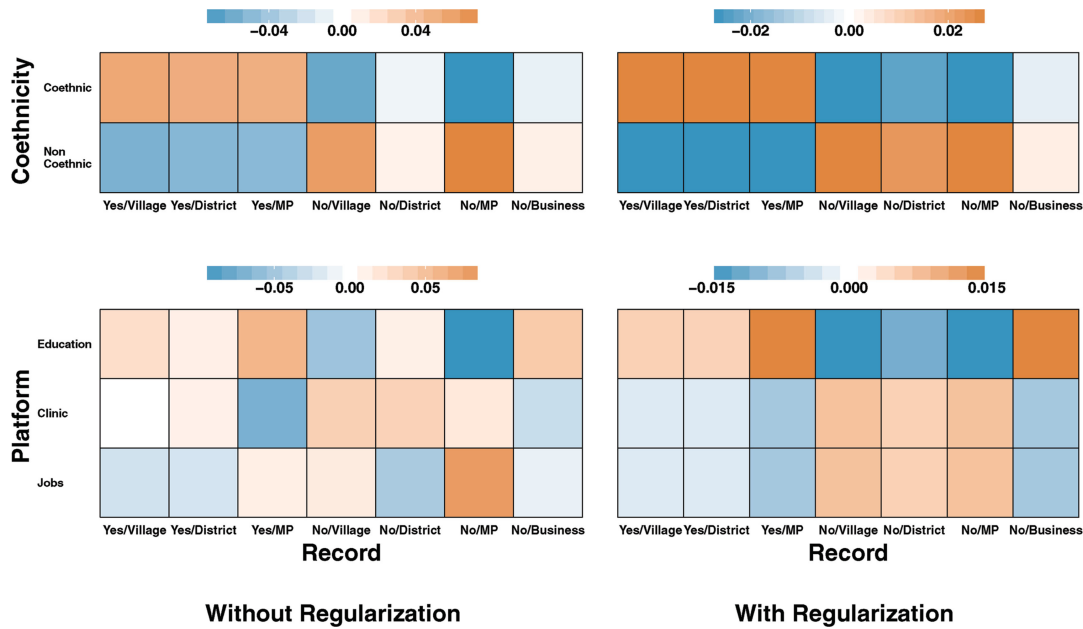| | Range | Selection prob. |
|---|---|---|
| AME | | |
| Record | 0.122 | 1.00 |
| Coethnicity | 0.053 | 1.00 |
| Platform | 0.023 | 1.00 |
| Degree | 0.000 | 0.58 |
| AMIE | | |
| Coethnicity × Record | 0.054 | 1.00 |
| Record × Platform | 0.030 | 1.00 |
| Platform × Coethnicity | 0.008 | 0.99 |
| Record × Degree | 0.000 | 0.60 |
| Coethnicity × Degree | 0.000 | 0.60 |
| Platform × Degree | 0.000 | 0.60 |

Record, which is the basis of the main finding in the original article, is estimated to have a large range of 5.4 percentage point, and is selected with probability one. The range of this AMIE is as great as that of the AME of Coethnicity and is greater than that of Platform. Additionally, the proposed method selects the causal interactions, Record × Platform and Platform × Coethnicity, with probability close to one. We focus on the two largest causal interactions, Coethnicity × Record and Record × Platform.

Next, we examine the estimated AMEs presented in Table 3. For the Record variable, under the 90% selection probability rule, we collapse a total of original seven levels into three levels—{Yes/Village, Yes/District, Yes/MP}, {No/Village, No/District , No/MP}, and {No/Businessman}. This partition suggests that politicians with good record are preferred over those without it including businessman. Similarly, we find two groups in the Platform variable—{Jobs, Clinic} and {Education}—where voters appear to favor candidates with the education platform on average.

We now investigate two significant causal interactions, Coethnicity × Record and Record × Platform. Figure 1 visualizes all estimated AMIEs within each factor interaction. The cells with warmer red (colder blue) color

**Table 3.** The estimated average marginal effects (AMEs). The estimated selection probability is the proportion of 5000 bootstrap replicates in which the difference between two adjacent levels is estimated to be different from zero.

| Factor | AME | Selection prob. |
|---|---|---|
| Record | | |
| Yes/Village | 0.122 | )0.64 |
| Yes/District | 0.122 | )0.80 |
| Yes/MP | 0.101 | )1.00 |
| No/Village | 0.047 | )0.76 |
| No/District | 0.051 | )0.84 |
| No/MP | 0.047 | )0.99 |
| No/Businessman | base | |
| Platform | | |
| Jobs | −0.023 | )0.80 |
| Clinic | −0.023 | )0.97 |
| Education | base | |
| Coethnicity | 0.053 | 1.00 |
| Degree | 0.000 | 0.57 |

**Figure 1.** The estimated AMIEs for `Coethnicity × Record` (the first row) and `Platform × Record` (the second row). The first and second columns show the estimated AMIEs without and with regularization, respectively.

represents a greater (smaller) AMIE than the average AMIE within that factor interaction. The estimates with regularization (right column) show clearer patterns for causal interaction than those without regularization (left column).

First, regarding the `Coethnicity × Record` interaction (upper panel of the figure), for example, we find that being coethnic gives an average bonus of 5.3 percentage point if a candidate is an MP with good record beyond the average effect of coethnicity (selec. prob. = 1). In contrast, being coethnic has an additional penalty of 4.6 percentage points when a candidate is a district level politician without good record (selec. prob. = 0.98). As shown in Equation (8), we can compute the average conditional effect as the sum of the AME and AMIE. As expected, while the conditional average effect of being coethnic for an MP candidate with good record is 10.7 percentage point (selec. prob. = 1), this effect is almost zero for an MP candidate without good record. These findings support the argument of Carlson (2015).

The decomposition shown in Equation (10) can be used to understand the ACE. As an illustration, we decompose the ACE of {Coethnic, No/Business} relative to {Non-coethnic, No/MP}, which is a estimated negative effect of 2.4 percentage points (selec. prob. = 0.89), as follows,

$$
\underbrace{\tau(\texttt{Coethnic, No/Business; Non-coethnic, No/MP})}_{-2.4}
$$
$$
= \underbrace{\psi(\texttt{Coethnic; Non-coethnic})}_{5.3}
$$
$$
+ \underbrace{\psi(\texttt{No/Business; No/MP})}_{-4.7}
$$
$$
+ \underbrace{\pi(\texttt{Coethnic, No/Business; Non-coethnic, No/MP})}_{-3.0}.
$$

We observe that while the average effect of being coethnic is 5.3 percentage points, being a businessman, relative to being an MP without good record, yields an average effect of negative 4.7 percentage points. In addition, being a coethnic businessman

has an additional penalty of 3 percentage points relative to non-coethnic MP without good record. All three estimates are selected with probability close to one.

Finally, we examine the `Platform × Record` interaction, which was not discussed in the original study. We find two distinct groups: (1) politicians with record, businessmen without record and (2) politicians without record. Candidates in the second group appear to receive an additional penalty by promising to improve education. Specifically, the estimated AMIE of {Education, No/MP} relative to {Job, No/MP} is −2.3 percentage point (selec. prob. = 0.99). In fact, the average conditional effect of `Education` relative to `Job` given `No/MP` is about zero (selec. prob. = 0.75). These results suggest that even though promising to improve education is effective on average (the estimated AME of `Education` relative to `Job` is 2.3 percentage point (selec. prob. = 0.98)), it has no effect for politicians without record.

## 7. Concluding Remarks

In this article, we propose a new causal interaction effect for randomized experiments with a factorial design, in which there exist many factors with each having several levels. We call this quantity, the average marginal interaction effect (AMIE). Unlike the conventional causal interaction effect, the AMIE is invariant to the choice of baseline. This enables us to provide a simpler interpretation even in a high-dimensional setting. We show how to nonparametrically estimate the AMIE within the ANOVA regression framework. The invariance property also enables us to apply a regularization method by directly penalizing the AMIEs. This reduces false discovery and facilitates interpretation.

We emphasize that the AMIE, which is a generalization of the average marginal effect studied in the literature on factorial experiments, critically depends on the distribution of treatments. For example, in a well-known audit study of labor market discrimination where researchers randomize

the information on the resume of a fictitious job applicant (e.g., Bertrand and Mullainathan 2004), the average effect of applicant's race requires the specification of other attributes such as education levels and prior job experiences. In the real world, these characteristics may be correlated with race and act as an effect modifier. Thus, ideally, researchers should obtain the target population distribution of treatments, for example, the characteristics of job applicants in a relevant labor market, and use it as the basis for treatment randomization. This will improve the external validity of experimental studies.

Finally, our method is motivated by and applied to conjoint analysis, a popular survey experiment with a factorial design. The methodological literature on conjoint analysis has largely ignored the role of causal interaction. The method proposed in this article allows researchers to effectively explore significant causal interaction among several factors. Although not investigated in this article, future research should investigate interaction between treatments and pretreatment covariates. It is also of interest to develop sequential experimental designs in the context of factorial experiments so that researchers can efficiently reduce the number of treatments.

## Supplementary Materials

In the supplementary materials, we provide proofs of all the theorems presented in the article.

## Acknowledgments

## ORCID

Naoki Egami http://orcid.org/0000-0002-5491-2174
Kosuke Imai http://orcid.org/0000-0002-2748-1022

## References

Athey, S., and Imbens, G. (2016), "Recursive Partitioning for Heterogeneous Causal Effects," *Proceedings of the National Academy of Sciences*, 113, 7353–7360. [537]

Bertrand, M., and Mullainathan, S. (2004), "Are Emily and Greg More Employable Than Lakisha and Jamal?: A Field Experiment on Labor Market Discrimination," *American Economic Review*, 94, 991–1013. [539]

Bondell, H. D., and Reich, B. J. (2009), "Simultaneous Factor Selection and Collapsing Levels in ANOVA," *Biometrics*, 65, 169–177. [536]

Breiman, L. (1996), "Heuristics of Instability and Stabilization in Model Selection," *The Annals of Statistics*, 24, 2350–2383. [537]

Bühlmann, P., and Dezeure, R. (2016), Discussion of "Regularized Regression for Categorical Data" by Tutz and Gertheiss, *Statistical Modelling*, 16, 205–211. [537]

Carlson, E. (2015), "Ethnic Voting and Accountability in Africa: A Choice Experiment in Uganda," *World Politics*, 67, 353–385. [530,538]

Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., and Robins, J. (2018), "Double/Debiased Machine Learning for Treatment and Structural Parameters," *The Econometrics Journal*, 21, C1–C68. [537]

Cox, D. R. (1958), *Planning of Experiments*, New York: Wiley. [531]

——— (1984), "Interaction," *International Statistical Review*, 52, 1–24. [530,532]

Dasgupta, T., Pillai, N. S., and Rubin, D. B. (2015), "Causal Inference From $2^k$ Factorial Designs by Using Potential Outcomes," *Journal of the Royal Statistical Society*, Series B, 77, 727–753. [530,531,532,533,535]

de González, A. B., and Cox, D. R. (2007), "Interpretation of Interaction: A Review," *The Annals of Applied Statistics*, 1, 371–385. [530]

Efron, B. (2014), "Estimation and Accuracy After Model Selection," *Journal of the American Statistical Association*, 109, 991–1007. [537]

Egami, N., Ratkovic, M., and Imai, K. (2017), "FindIt: Finding Heterogeneous Treatment Effects," available at the Comprehensive R Archive Network (CRAN), available at *https://CRAN.R-project.org/package=FindIt*. [530,539]

Green, D. P., and Kern, H. L. (2012), "Modeling Heterogeneous Treatment Effects in Survey Experiments With Bayesian Additive Regression Trees," *Public Opinion Quarterly*, 76, 491–511. [530]

Green, P. E., Krieger, A. M., and Wind, Y. (2001), "Thirty Years of Conjoint Analysis: Reflections and Prospects," *Interfaces*, 31, 56–73. [529,530]

Green, P. E., and Rao, V. R. (1971), "Conjoint Measurement for Quantifying Judgmental Data," *Journal of Marketing Research*, 8, 355–363. [530]

Grimmer, J., Messing, S., and Westwood, S. J. (2017), "Estimating Heterogeneous Treatment Effects and the Effects of Heterogeneous Treatments With Ensemble Methods," *Political Analysis*, 25, 413–434. [530]

Hainmueller, J., and Hopkins, D. J. (2015), "The Hidden American Immigration Consensus: A Conjoint Analysis of Attitudes Toward Immigrants," *American Journal of Political Science*, 59, 529–548. [529,530]

Hainmueller, J., Hopkins, D. J., and Yamamoto, T. (2014), "Causal Inference in Conjoint Analysis: Understanding Multidimensional Choices via Stated Preference Experiments," *Political Analysis*, 22, 1–30. [530,531,533,535]

Hastie, T., Tibshirani, R., and Wainwright, M. (2015), *Statistical Learning With Sparsity: The Lasso and Generalizations*, Boca Raton, FL: CRC Press. [537]

Hill, J. L. (2012), "Bayesian Nonparametric Modeling for Causal Inference," *Journal of Computational and Graphical Statistics*, 20, 217–240. [530]

Huang, J., Breheny, P., and Ma, S. (2012), "A Selective Review of Group Selection in High-Dimensional Models," *Statistical Science*, 27, 481–499. [536]

Huang, J., Ma, S., Xie, H., and Zhang, C.-H. (2009), "A Group Bridge Approach for Variable Selection," *Biometrika*, 96, 339–355. [536]

Imai, K., and Ratkovic, M. (2013), "Estimating Treatment Effect Heterogeneity in Randomized Program Evaluation," *Annals of Applied Statistics*, 7, 443–470. [530]

Jaccard, J., and Turrisi, R. (2003), *Interaction Effects in Multiple Regression*, Thousand Oaks, CA: Sage Publications. [530,534]

Lim, M., and Hastie, T. (2015), "Learning Interactions via Hierarchical Group-Lasso Regularization," *Journal of Computational and Graphical Statistics*, 24, 627–654. [536]

Lu, J. (2016a), "Covariate Adjustment in Randomization-Based Causal Inference for 2k Factorial Designs," *Statistics & Probability Letters*, 119, 11–20. [530]

——— (2016b), "On Randomization-Based and Regression-Based Inferences for 2k Factorial Designs," *Statistics & Probability Letters*, 112, 72–78. [530]

Luce, R. D., and Tukey, J. W. (1964), "Simultaneous Conjoint Measurement: A New Type of Fundamental Measurement," *Journal of Mathematical Psychology*, 1, 1–27. [529,530]

Marshall, D., Bridges, J. F., Hauber, B., Cameron, R., Donnalley, L., Fyie, K., and Johnson, F. R. (2010), "Conjoint Analysis Applications in Health: How are Studies Being Designed and Reported?" *The Patient: Patient-Centered Outcomes Research*, 3, 249–256. [529,530]

Marshall, P., and Bradlow, E. T. (2002), "A Unified Approach to Conjoint Analysis Models," *Journal of the American Statistical Association*, 97, 674–682. [529,530]

McFadden, D. (1974), "Conditional Logit Analysis of Qualitative Choice Behavior," in *Frontiers in Econometrics*, ed. P. Zarembka, New York: Academic Press, pp. 105–142. [537]

Meier, L., Van De Geer, S., and Bühlmann, P. (2008), "The Group Lasso for Logistic Regression," *Journal of the Royal Statistical Society*, Series B, 70, 53–71. [536]

Meinshausen, N., and Bühlmann, P. (2010), "Stability Selection," *Journal of the Royal Statistical Society*, Series B, 72, 417–473. [537]

Murphy, S. A. (2003), "Optimal Dynamic Treatment Regimes (with discussions)," *Journal of the Royal Statistical Society*, Series B, 65, 331–366. [530]

Post, J. B., and Bondell, H. D. (2013), "Factor Selection and Structural Identification in the Interaction ANOVA Model," *Biometrics*, 69, 70–79. [530,535,536]

Robins, J. M. (2004), "Optimal Structural Nested Models for Optimal Sequential Decisions," in *Proceedings of the Second Seattle Symposium in Biostatistics: Analysis of Correlated Data*, New York: Springer, pp. 189–326. [530]

Rubin, D. B. (1990), Comments on "On the Application of Probability Theory to Agricultural Experiments. Essay on Principles. Section 9" by J. Splawa-Neyman translated from the Polish and edited by D. M. Dabrowska and T. P. Speed, *Statistical Science*, 5, 472–480. [531]

Scheffe, H. (1959), *The Analysis of Variance*, New York: Wiley. [530,535]

VanderWeele, T. (2015), *Explanation in Causal Inference: Methods for Mediation and Interaction*, Oxford: Oxford University Press. [532]

VanderWeele, T. J., and Knol, M. J. (2014), "A Tutorial on Interaction," *Epidemiologic Methods*, 3, 33–72. [530]

Wager, S., and Athey, S. (2017), "Estimation and Inference of Heterogeneous Treatment Effects Using Random Forests," *Journal of the American Statistical Association*. [530]

Wasserman, L., and Roeder, K. (2009), "High Dimensional Variable Selection," *Annals of Statistics*, 37, 2178–2201. [537]

Wu, C. J., and Hamada, M. S. (2011), *Experiments: Planning, Analysis, and Optimization* (Vol. 552), New York: Wiley. [534]

Yuan, M., and Lin, Y. (2006), "Model Selection and Estimation in Regression With Grouped Variables," *Journal of the Royal Statistical Society*, Series B, 68, 49–67. [536]

Zhao, P., Rocha, G., and Yu, B. (2009), "The Composite Absolute Penalties Family for Grouped and Hierarchical Variable Selection," *The Annals of Statistics*, 37, 3468–3497. [536]