

Multiple Testing of Stochastic Monotonicity

Qian Wu^{*†} David M. Kaplan[‡]

May 28, 2025

Abstract

We develop multiple testing methodology to assess the evidence that an outcome variable is stochastically increasing in a covariate. Such a relationship holds globally if at each possible outcome value, the conditional CDF evaluated at that value is decreasing in the covariate. Rather than test that single global null hypothesis, we use multiple testing to separately evaluate each constituent conditional CDF inequality. From another perspective: instead of assessing only whether the mean outcome is increasing in the covariate, we assess the relationship across the entire outcome distribution. Inverting our multiple testing procedure that controls familywise error rate, we construct “inner” and “outer” confidence sets for the true set of points consistent with stochastic increasingness. Simulations show reasonable finite-sample properties. Empirically, we apply our methodology to address well-known empirical policy questions concerning socioeconomic disparities in health outcomes. Practically, we provide code implementing our methodology and replicating our results.

JEL Classification: C25, I10

Keywords: Confidence set; Familywise error rate; Health; Life satisfaction

^{*}School of Statistics and Data Science, Southwestern University of Finance and Economics, Chengdu 611130, China; wuqj@swufe.edu.cn; some initial work was completed as part of a dissertation while a doctoral student at the University of Missouri

[†]Big Data Laboratory on Financial Security and Behavior, SWUFE (Laboratory of Philosophy and Social Sciences, Ministry of Education), Chengdu 611130, China

[‡]Corresponding author; Department of Economics, University of Missouri, 615 Locust St, Columbia, MO 65211, USA; kaplandm@missouri.edu

1 Introduction

A fundamental empirical question is whether outcome variable Y is “increasing” in covariate X , which can be characterized in terms of stochastic monotonicity. For example, is life satisfaction increasing in education? We say that Y is “stochastically increasing” in X if, for any two values of X , the conditional distribution of Y at the higher value first-order stochastically dominates that at the lower value. Economically, assuming Y is scaled so higher values are better, this conditional stochastic dominance means that the conditional distribution of Y is “better” at higher X values, in the sense of higher expected utility. However, because this is a strong property, it can be statistically difficult to find strong evidence in its favor, and conversely it may nearly hold but be rejected for an economically small violation.

We contribute to the stochastic monotonicity inference literature by proposing and justifying new methodology that focuses on multiple testing, as well as allowing either continuous or ordinal/discrete outcomes. Other stochastic monotonicity tests (see below) assume a continuous outcome and only test a single null hypothesis of global stochastic monotonicity. Instead of testing the single null hypothesis that Y is stochastically increasing in X , we consider multiple testing of the conditional CDF inequalities that jointly comprise stochastic increasingness. Specifically, each inequality checks whether the conditional CDF of Y is decreasing as X increases to the next-highest value, over all values in the support of Y and any finite number of X values. Despite this being an infinite number of points when Y is continuous, our methods still control the familywise error rate.

Complementing global testing of a single null hypothesis, multiple testing addresses the aforementioned statistical concerns in the following two ways. First, if stochastic increasingness is rejected, then multiple testing shows more precisely whether many of the constituent conditional CDF inequalities are rejected, or only a few, or even just one. Second, if stochastic increasingness is not rejected, then multiple testing can be used on the reversed inequalities (conditional CDF increasing with X) to see if they can be rejected in favor of the original inequalities (decreasing with X). This gathers stronger statistical evidence in favor of stochastic increasingness than simply failing to reject the original null hypothesis, because type II error rates (false non-rejection) are not controlled. For example, due to small sample size (high statistical uncertainty), stochastic increasingness may fail to be rejected even if the estimated conditional distributions are stochastically *decreasing*.

Although the economic interpretation is the same either way, our statistical approach depends on the type of outcome variable. For ordinal or discrete outcomes, there are a finite number of inequalities, so we use a maximum t -statistic approach. This allows us to exactly

control the asymptotic familywise error rate under the least favorable null. For continuous outcomes, we achieve finite-sample control of the familywise error rate by building on the two-sample multiple testing procedure of Goldman and Kaplan (2018). They use the probability integral transform and joint distribution of order statistics to achieve finite-sample control of familywise error rate. Their method applies directly to our setting if X is binary. Otherwise, we use a Bonferroni correction to maintain the finite-sample control of familywise error rate. Although the Bonferroni correction errs on the conservative side, there are two reasons this may not have a large quantitative effect on power. First, the effect is smaller with a smaller number of X values, which applies to the education variable in our empirical illustration, for example. Second, the effect is smaller if the dependence across X is more negative. Note Bonferroni is not conservative in the extreme case of “perfect negative dependence” where a false rejection at one X implies no false rejections at any other X . In our case, if the estimated conditional CDF at $X = 2$ is higher than the true one, then it makes us more likely to reject the comparisons of $X = 1$ with $X = 2$, but less likely to reject the comparisons of $X = 2$ with $X = 3$; that is, negative dependence.

Additionally, our multiple testing procedure can be inverted into “inner” and “outer” confidence sets. The inner confidence set contains all points at which the reversed inequalities have been rejected, i.e., where there is strong evidence in favor of the original conditional CDF inequality consistent with stochastic increasingness. The outer confidence set contains all points at which the original inequalities have not been rejected. Such confidence sets are similar to those of Kaplan (2024) and to equation (1) of Armstrong and Shen (2023). The inner confidence set can be seen as a conservative estimator of the true set of inequalities consistent with stochastic increasingness, in that the inner set is contained within the true set with high probability. Conversely, the outer set contains the true set with high probability.

Our methodology is illustrated through two empirical applications. Our first application has an ordinal outcome variable, looking at the relationship between general health and education in 2022 UK data. Across almost all levels of both education and general health, we find strong evidence of health stochastically increasing with education. Our second application has a continuous outcome variable, looking at the relationship between physical health and education. Specifically, we use the physical component summary/score (PCS) computed from the SF-12, the widely used 12-item Short-Form Health Survey. While there is too much statistical uncertainty in the distributional tails to make definitive statements, there is strong evidence across most of the distribution in favor of PCS stochastically increasing with education, across all levels of education. At least for the UK context, these findings suggest that education-related health disparities are not limited to mean differences but extend to the entire distribution. In both cases, compared to a global test of stochastic

increasingness that would merely report “fail to reject,” our results are more detailed and useful.

Besides intrinsic interest in settings like the education gradient in health or intergenerational mobility (comparing child’s outcome Y with their parent’s outcome X), stochastic monotonicity also appears in certain identifying assumptions or as a testable implication thereof. For example, stochastic monotonicity is implied by the combination of (semi-)monotone treatment response and exogenous treatment selection, as discussed by Manski (1997, §3.4). As another example, Small, Tan, Ramsahai, Lorch, and Brookhart (2017) identify a weighted average treatment effect when the treatment is stochastically increasing in the instrument. Although testing is trivial in the binary–binary setting they focus on, more generally our methodology can help assess their identifying assumption of stochastic monotonicity.

Literature Our methodology contributes to the literature on stochastic monotonicity and multiple testing. Inference on stochastic monotonicity has focused on continuous Y and testing the single hypothesis of global stochastic monotonicity; for example, see Lee, Linton, and Whang (2009), Seo (2018), and Chetverikov, Wilhelm, and Kim (2021). Our reversing the direction of null hypothesis inequalities to find stronger evidence in favor of stochastic increasingness is inspired by Davidson and Duclos (2013), who make a related argument for testing the null of non-dominance against the alternative of stochastic dominance, which can be seen as a special case of stochastic monotonicity with binary X . Although we emphasize multiple testing, testing the set of stochastic monotonicity inequalities jointly would fit in the (moment) inequality literature; we essentially follow the least favorable approach with a max- t statistic as in Section 4.1.1 of Canay and Shaikh (2017), whose survey includes additional references. Our proof strategies are also similar to those of Zhao and Kaplan (2024), who instead consider multiple testing of a function’s value across a finite set of points. Kaplan and Zhao (2023) also use multiple testing with ordinal outcome variables, but they only compare two (unconditional) distributions and focus not on ordinal categories but on a latent variable’s quantiles. Surveys of the multiple testing literature can be found in Lehmann and Romano (2005b) and Romano, Shaikh, and Wolf (2010).

Paper structure Section 2 describes the setting with an ordinal or discrete outcome, and our methodology and its properties. Section 3 presents simulation results. Section 4 describes our contributions with a continuous outcome. Section 5 applies our methodology empirically. Appendix A collects proofs.

Notation and abbreviations Generally, scalars, vectors, and matrices are respectively typeset like X , \mathbf{X} , and \mathbf{X} . The indicator function is $\mathbb{1}\{\cdot\}$, with $\mathbb{1}\{A\} = 1$ if event A occurs and $\mathbb{1}\{A\} = 0$ if not, and “ \iff ” means if and only if. Acronyms used include those for confidence set (CS), continuous mapping theorem (CMT), cumulative distribution function (CDF), familywise error rate (FWER), and multiple testing procedure (MTP).

2 Results for ordinal and discrete outcomes

In this section, we describe the setting, assumptions, methodology, and asymptotic properties when the outcome variable is ordinal or discrete.

2.1 Setting

Consider random variables Y and X . The outcome Y is discrete or ordinal, with categories labeled $Y \in \{1, 2, \dots, J\}$, for finite J . For example, if the ordinal categories are “disagree,” “neutral,” and “agree,” then they are respectively coded as $Y = 1$, $Y = 2$, and $Y = 3$; or if the possible discrete values are 0, 0.5, 1, and 1.5, then these are respectively coded as $Y = 1$, $Y = 2$, $Y = 3$, and $Y = 4$. The covariate may also be discrete or ordinal, with categories labeled $X \in \{1, 2, \dots, K\}$, for finite K . Alternatively, the covariate may be a discretization of a continuous variable; for example, a continuous covariate with support $[0, 100]$ may be discretized with values in $[0, 10]$ coded as $X = 1$, $(10, 20]$ coded as $X = 2$, etc. In sum, the supports \mathcal{Y} and \mathcal{X} are

$$Y \in \mathcal{Y} \equiv \{1, \dots, J\}, \quad X \in \mathcal{X} \equiv \{1, \dots, K\}. \quad (1)$$

The population and empirical (estimated) conditional CDF values are respectively

$$F_x(y) \equiv \mathbb{P}(Y \leq y \mid X = x),$$

$$\hat{F}_x(y) \equiv \frac{\sum_{i=1}^n \mathbb{1}\{Y_i \leq y\} \mathbb{1}\{X_i = x\}}{n_x}, \quad n_x \equiv \sum_{i=1}^n \mathbb{1}\{X_i = x\}, \quad (2)$$

given iid observations over $i = 1, \dots, n$. That is, $F_x(y)$ is the proportion of the $X = x$ subpopulation with $Y \leq y$, and $\hat{F}_x(y)$ is the proportion of the $X_i = x$ subsample with $Y_i \leq y$.

We assume iid sampling and condition on the size of each subsample n_x , which is equivalent to conditioning on all the observed X_i values (like in classical linear regression results). This is also equivalent to repeated sampling of n_x iid draws of Y_i from the corresponding

conditional distribution independently for each $x = 1, \dots, K$, which is how we formalize the setting in Assumption A1.

Assumption A1. Using the notation in (1) and (2), $n_1/n_x \rightarrow \gamma_x \in (0, \infty)$ for each $x \in \mathcal{X}$, and the sample consists of n_x observations with $X_i = x$ for each $x \in \mathcal{X}$, with the corresponding Y_i values sampled iid from the corresponding population conditional distribution of $Y \mid X = x$, and all Y_i are mutually independent.

2.2 Stochastic monotonicity inequalities

Outcome Y stochastically increasing in X means that for any $x_2 > x_1$, the conditional distribution of $Y \mid X = x_2$ first-order stochastically dominates that of $Y \mid X = x_1$. In the notation of (2), this means

$$F_{x_2}(y) \leq F_{x_1}(y) \text{ for all } (x_1, x_2, y) \in \mathcal{X} \times \mathcal{X} \times \mathcal{Y} \text{ with } x_2 > x_1. \quad (3)$$

As with stochastic dominance testing, despite being economically different, strict and weak inequalities are statistically indistinguishable, so we do not emphasize the difference.

We test a subset of the inequalities in (3). We restrict attention to $y \in \{1, \dots, J-1\}$ because $F_{x_2}(J) = F_{x_1}(J) = 1$ for all x_1, x_2 . In principle, we could test all $\{(x_1, x_2) \in \mathcal{X} \times \mathcal{X} : x_1 < x_2\}$, which is $K(K-1)/2$ different values of (x_1, x_2) . However, without claiming optimality against all alternatives, we restrict attention to the $K-1$ pairs satisfying $x_2 = x_1 + 1$ because reducing the number of comparisons from $(J-1)K(K-1)/2$ to $(J-1)(K-1)$ allows a lower critical value that improves power, while still testing a subset of inequalities jointly equivalent to (3). That is, (3) holds if and only if

$$\theta_{x,y} \leq 0 \text{ for all } (x, y) \in \{1, \dots, K-1\} \times \{1, \dots, J-1\}, \quad \theta_{x,y} \equiv F_{x+1}(y) - F_x(y). \quad (4)$$

Additionally, testing (4) yields results that are easier to interpret and communicate.

For ordinal outcomes, although first-order stochastic dominance generally suggests one distribution is “better” than another, the relationship is more complex if utility depends not on the observed ordinal value but rather on a latent continuous variable. We focus here on testing the conditional ordinal distributions; for interpretation with respect to a latent variable, see Kaplan and Zhao (2023).

2.3 Multiple testing procedure

Based on (4), we consider multiple testing of the following family of hypotheses:

$$H_{0x,y}: \theta_{x,y} \leq 0 \text{ for } (x, y) \in \{1, \dots, K-1\} \times \{1, \dots, J-1\}. \quad (5)$$

The number of hypotheses is $(J-1)(K-1)$. A multiple testing procedure makes a binary decision (reject or not) for each hypothesis, hence $2^{(J-1)(K-1)}$ possible results.

Our multiple testing procedure (MTP) controls the asymptotic familywise error rate (FWER). Although there are other measures of overall false positive rate for multiple testing, like the k -FWER and false discovery proportion (Lehmann and Romano, 2005a) and the false discovery rate (Benjamini and Hochberg, 1995), we use FWER because it has a clear interpretation and allows us to invert our MTP into confidence sets. FWER is defined as (Lehmann and Romano, 2005b, §9.1)

$$\text{FWER} \equiv \text{P}(\text{reject any true } H_{0x,y}). \quad (6)$$

Conversely, $1 - \text{FWER}$ is the probability of having zero false rejections (zero type I errors). Our MTP has “strong control” because it controls asymptotic FWER regardless of which $H_{0x,y}$ are true or false (Lehmann and Romano, 2005b, §9.1).

Our MTP uses standard t -statistics but with a higher critical value that accounts for multiple testing. For any real scalar d , define

$$\hat{t}_{x,y}(d) \equiv \frac{\hat{F}_{x+1}(y) - \hat{F}_x(y) - d}{\hat{s}_{x,y}}, \quad (7)$$

where $\hat{s}_{x,y}$ is defined in (22) as an estimator of the standard deviation of $\hat{F}_{x+1}(y) - \hat{F}_x(y)$. As usual, in practice we compute the t -statistic centered at our hypothesized value $d = 0$, and asymptotic properties can be bounded by the behavior of the t -statistic centered at the true $d = \theta_{x,y} \equiv F_{x+1}(y) - F_x(y)$. As an important ingredient of our FWER derivations, Lemma 1 establishes the asymptotic normal distribution of random vector $\hat{\mathbf{t}}$ that contains all the t -statistics centered at the true $\theta_{x,y}$:

$$\hat{\mathbf{t}} \equiv (\hat{\mathbf{t}}_1, \hat{\mathbf{t}}_2, \dots, \hat{\mathbf{t}}_{J-1})', \quad \hat{\mathbf{t}}_j \equiv (\hat{t}_{K-1,j}(\theta_{K-1,j}), \hat{t}_{K-2,j}(\theta_{K-2,j}), \dots, \hat{t}_{1,j}(\theta_{1,j})). \quad (8)$$

Lemma 1. *Under Assumption A1,*

$$\hat{\mathbf{t}} \xrightarrow{d} \mathbf{t} \sim \text{N}(\mathbf{0}, \underline{\Sigma}), \quad \underline{\Sigma} = \underline{\mathbf{S}}' \underline{\mathbf{W}} \underline{\mathbf{S}},$$

where $\underline{\mathbf{W}}$ and $\underline{\mathbf{S}}$ are non-random matrices defined in (20) and (24).

We provide some brief intuition for the critical value whose validity is formally established in the proof of Theorem 2. Intuitively, the least favorable configuration here is when all $\theta_{x,y} = 0$: every null $H_{0x,y}$ is true (and thus can contribute to FWER), but just barely, so there is the highest probability of some estimated $\hat{\theta}_{x,y} > 0$ large enough to reject $H_{0x,y}$. In that case, $\hat{t}_{x,y}(0) = \hat{t}_{x,y}(\theta_{x,y})$, and we make a familywise error whenever the maximum t -statistic $\max_{x,y} \hat{t}_{x,y}(0)$ exceeds our critical value. Thus, using the asymptotic approximation, if we want $\text{FWER} = \alpha$ in this least favorable case, then the critical value should be the $(1 - \alpha)$ -quantile of the distribution of the maximum of \mathbf{t} :

$$c_\alpha \equiv (1 - \alpha)\text{-quantile of } \max_{x,y} t_{x,y}, \quad (9)$$

where the $t_{x,y}$ are the elements of \mathbf{t} and $\mathbf{t} \sim N(\mathbf{0}, \underline{\Sigma})$ by Lemma 1. In practice, the unknown $\underline{\Sigma}$ can be replaced by consistent estimator $\hat{\underline{\Sigma}}$:

$$\hat{c}_\alpha \equiv (1 - \alpha)\text{-quantile of } \max_{x,y} \tilde{t}_{x,y}, \quad \tilde{\mathbf{t}} \sim N(\mathbf{0}, \hat{\underline{\Sigma}}). \quad (10)$$

Although an analytic formula is intractable, (10) can be simulated as in our code. Alternatively, the distribution of \mathbf{t} could be approximated by bootstrap.

Method 1. *First, compute t -statistics $\hat{t}_{x,y}(0)$ as in (7). Second, given desired FWER level α , simulate the critical value \hat{c}_α in (10) following the steps in Section 2.5. Third, for each (x, y) , to test all $H_{0x,y}: \theta_{x,y} \leq 0$ as in (5), reject $H_{0x,y}: \theta_{x,y} \leq 0$ if $\hat{t}_{x,y}(0) > \hat{c}_\alpha$. Alternatively, to test the reversed family of hypotheses $H_{0x,y}^\geq: \theta_{x,y} \geq 0$, reject $H_{0x,y}$ if $\hat{t}_{x,y}(0) < -\hat{c}_\alpha$.*

Method 1 includes the reversed $H_{0x,y}^\geq$ because their rejection provides stronger evidence of $\theta_{x,y} < 0$ than does non-rejection of $H_{0x,y}$. For example, even if $\hat{\theta}_{x,y} = 1.7$ (our best guess is a positive $\theta_{x,y}$), if there is a small sample size or otherwise high uncertainty, we may still not reject $H_{0x,y}: \theta_{x,y} \leq 0$. In contrast, to reject $H_{0x,y}^\geq$, not only must we have $\hat{\theta}_{x,y} < 0$, but it must be significantly less than zero (compared to our uncertainty) for the test to control the false positive rate. In sum, non-rejection of $H_{0x,y}$ suggests the data are consistent with the hypothesis that Y is stochastically increasing in X , but rejection of $H_{0x,y}^\geq$ provides stronger evidence in favor of the inequalities that comprise stochastic increasingness of Y in X .

Theorem 2 theoretically justifies our MTP.

Theorem 2. *Under Assumption A1, Method 1 has strong control of asymptotic FWER at level α .*

2.4 Confidence sets

The MTP in Method 1 can be inverted into confidence sets for the true set

$$\mathcal{T} \equiv \{(x, y) : \theta_{x,y} \leq 0\}. \quad (11)$$

The goal is, with high asymptotic probability, for inner confidence set $\widehat{\mathcal{CS}}_{inner}$ to be contained within the true set \mathcal{T} , and for outer confidence set $\widehat{\mathcal{CS}}_{outer}$ to contain \mathcal{T} . That is, for confidence level $1 - \alpha$,

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{P}(\widehat{\mathcal{CS}}_{inner} \subseteq \mathcal{T}) &\geq 1 - \alpha, \\ \lim_{n \rightarrow \infty} \mathbb{P}(\widehat{\mathcal{CS}}_{outer} \supseteq \mathcal{T}) &\geq 1 - \alpha. \end{aligned} \quad (12)$$

Intuitively, $\widehat{\mathcal{CS}}_{inner}$ provides a conservative “estimate” of the true set \mathcal{T} , while $\widehat{\mathcal{CS}}_{outer}$ gives a larger “estimate” describing how large the true set might be. Corresponding to the earlier discussion of stronger and weaker evidence, the inner confidence set collects points for which the reversed $H_{0x,y}^{\geq}$ is rejected (strong evidence), whereas the outer confidence set collects points for which the original $H_{0x,y}$ is not rejected (weak evidence).

Theorem 3 formally establishes the property in (12) for our confidence sets described in Method 2.

Method 2. *First, run Method 1. The inner confidence set $\widehat{\mathcal{CS}}_{inner}$ collects all pairs of (x, y) for which the reversed null $H_{0x,y}^{\geq} : \theta_{x,y} \geq 0$ is rejected. The outer confidence set $\widehat{\mathcal{CS}}_{outer}$ collects all pairs of (x, y) for which the original null $H_{0x,y} : \theta_{x,y} \leq 0$ is not rejected.*

Theorem 3. *Under Assumption A1, Method 2 satisfies (12).*

If instead of \mathcal{T} in (11) we are interested in its complement $\mathcal{T}^c \equiv \{(x, y) : \theta_{x,y} > 0\}$, then the outer CS is the complement of the inner CS for \mathcal{T} , and the inner CS is the complement of the outer CS for \mathcal{T} . This follows because the event $\widehat{\mathcal{CS}}_{inner} \subseteq \mathcal{T}$ is equivalent to the inner CS complement containing \mathcal{T}^c , and the event $\widehat{\mathcal{CS}}_{outer} \supseteq \mathcal{T}$ is equivalent to the outer CS complement being contained within \mathcal{T}^c . Thus, using the notation from (12),

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{P}(\widehat{\mathcal{CS}}_{inner}^c \supseteq \mathcal{T}^c) &= \lim_{n \rightarrow \infty} \mathbb{P}(\widehat{\mathcal{CS}}_{inner} \subseteq \mathcal{T}) \geq 1 - \alpha, \\ \lim_{n \rightarrow \infty} \mathbb{P}(\widehat{\mathcal{CS}}_{outer}^c \subseteq \mathcal{T}^c) &= \lim_{n \rightarrow \infty} \mathbb{P}(\widehat{\mathcal{CS}}_{outer} \supseteq \mathcal{T}) \geq 1 - \alpha. \end{aligned}$$

2.5 Critical value simulation

We compute the critical value \hat{c}_α in (10) by simulation. Let

$$\hat{\underline{\Sigma}} = \hat{\underline{\mathbf{S}}}^{\prime} \hat{\underline{\mathbf{W}}} \hat{\underline{\mathbf{S}}}, \quad (13)$$

where $\hat{\underline{\mathbf{W}}}$ and $\hat{\underline{\mathbf{S}}}$ are the sample analogs of $\underline{\mathbf{W}}$ and $\underline{\mathbf{S}}$ from (20) and (24). The consistency result in (23) implies $\hat{\underline{\mathbf{S}}} \xrightarrow{p} \underline{\mathbf{S}}$, and by the weak law of large numbers $\hat{\underline{\mathbf{W}}} \xrightarrow{p} \underline{\mathbf{W}}$, so applying the continuous mapping theorem yields $\hat{\underline{\Sigma}} \xrightarrow{p} \underline{\Sigma}$.

Given $\hat{\underline{\Sigma}}$, the simulation proceeds as follows, as implemented in our code. First, we randomly draw a vector from $N(\mathbf{0}, \hat{\underline{\Sigma}})$. Second, we take the maximum of this vector. Third, we repeat this process N times, collecting the N maxima. Fourth, we take the $(1-\alpha)$ -quantile of these N maximum values; this is \hat{c}_α .

With large enough N , we can achieve arbitrarily small simulation error. Larger N improves accuracy but increases computation time. Given our simulation results (Section 3), we suggest $N = 1000$ as a default. In our empirical application, we use $N = 100,000$ because computation time is still only a few seconds.

3 Simulations

This section presents simulations of finite-sample FWER for our asymptotically justified multiple testing procedure when Y is ordinal or discrete. We use a variety of sample sizes and numbers of Y and X categories. We use the “least favorable” configuration that maximizes FWER by setting all null hypothesis inequalities to be binding: $\theta_{x,y} = 0$ for all (x, y) . We provide code in R (R Core Team, 2023) to replicate our results.¹

The simulation proceeds as follows. For the data generating process, Y and X respectively have J and K categories. The conditional distribution of Y is uniform across the J categories, so the conditional CDF is $F_x(y) = y/J$ for each (x, y) and all $\theta_{x,y} \equiv F_{x+1}(y) - F_x(y) = y/J - y/J = 0$. After randomly drawing a dataset having n_x observations with $X_i = x$ for each x (so total sample size $n = Kn_x$), Method 1 is run, with the critical value based on N simulations, and the results are stored. This is repeated 1000 times. Because all $H_{0x,y}$ are true, the simulated FWER is the proportion of simulated datasets for which at least one $H_{0x,y}$ is rejected. Besides FWER, we report the full range (across all simulated datasets) of simulated critical values c_α for each case.

Table 1 presents the simulation results. They are divided into three sections corresponding to the different (J, K) , but results are qualitatively similar in each section, showing the

¹https://qianjoewu.github.io/research/osm/Wu_Kaplan_OSM.zip

Table 1: Simulation results.

J	K	n_x	N	α	c_α	FWER
4	4	20	1000	0.05	[2.32, 2.59]	0.098
4	4	100	1000	0.05	[2.34, 2.58]	0.068
4	4	1000	1000	0.05	[2.36, 2.57]	0.060
4	4	10,000	1000	0.05	[2.36, 2.55]	0.045
4	4	10,000	10,000	0.05	[2.43, 2.52]	0.047
4	4	10,000	10,000	0.10	[2.18, 2.25]	0.112
4	4	10,000	10,000	0.01	[2.94, 3.08]	0.010
6	5	20	1000	0.05	[2.59, 2.84]	0.092
6	5	100	1000	0.05	[2.55, 2.81]	0.075
6	5	1000	1000	0.05	[2.57, 2.79]	0.068
6	5	10,000	1000	0.05	[2.60, 2.79]	0.053
6	5	10,000	10,000	0.05	[2.70, 2.77]	0.047
6	5	10,000	10,000	0.10	[2.45, 2.50]	0.105
6	5	10,000	10,000	0.01	[3.18, 3.32]	0.007
8	10	20	1000	0.05	[2.92, 3.18]	0.099
8	10	100	1000	0.05	[2.95, 3.18]	0.074
8	10	1000	1000	0.05	[2.94, 3.20]	0.056
8	10	10,000	1000	0.05	[2.95, 3.18]	0.046
8	10	10,000	10,000	0.05	[3.04, 3.12]	0.046
8	10	10,000	10,000	0.10	[2.82, 2.88]	0.090
8	10	10,000	10,000	0.01	[3.47, 3.61]	0.007

following patterns. First, FWER is somewhat above $\alpha = 0.05$ with the smallest sample size, closer to 0.10, suggesting that (as usual) results from smaller samples should be interpreted more cautiously. Second, as n_x increases, FWER approaches the nominal α . This reflects our procedure’s exact asymptotic FWER in this least favorable configuration; FWER would be lower in other cases. Third, compared to $N = 1000$, computing the critical value using $N = 10,000$ draws improves stability somewhat (tighter range of critical values) but improves FWER accuracy only very slightly, so in practice we suggest using $N = 1000$ as the default.

4 Results for continuous outcomes

In this section, we describe our contributions with continuous Y .

4.1 Setting and inequalities

The setting, notation, and Assumption A1 are the same as in Section 2.1, but now Y is continuous. To contrast with \mathcal{Y} from Section 2.1, here we write \mathbb{R} as the support of Y , with the understanding that the true support may be a subset of \mathbb{R} .

Analogous to (5), here we test

$$H_{0x,y}: \theta_{x,y} \leq 0 \text{ for } (x, y) \in \{1, \dots, K-1\} \times \mathbb{R}, \quad \theta_{x,y} \equiv F_{x+1}(y) - F_x(y). \quad (14)$$

This is an infinite number of hypotheses, corresponding to an infinite number of rejection decisions. However, they are highly dependent across y values, so for a given x , usually the rejected $H_{0x,y}$ correspond to only a few intervals of y values. Computationally, of course, we do not compute an infinite number of test statistics, but rather compute values corresponding to each observation and interpolate in a precise way.

4.2 Methodology

Our methodology builds on that of Goldman and Kaplan (2018). Their Method 5 includes a multiple testing procedure that directly applies to our setting when X is binary, so $K = 1$. That is, given $\theta_y \equiv F_2(y) - F_1(y)$, their Method 5 tests $H_{0y}: \theta_y \leq 0$ for each $y \in \mathbb{R}$, using order statistics to achieve strong control of finite-sample FWER (see their Theorem 9). Such multiple testing has clear economic significance. For example, let $F_2(\cdot)$ and $F_1(\cdot)$ be the wage CDFs for individuals with and without a high-school degree, respectively, in dollars per hour. Then H_{0y} means individuals without a degree have at least as high probability of wage below \$y per hour as individuals with a degree. The multiple testing procedure lets us test such comparisons across all y while controlling the probability of at least one false rejection. Reversing the direction of inequality, rejecting $H_{0y}^\geq: \theta_y \geq 0$ provides stronger evidence in favor of $F_2(y) < F_1(y)$, meaning that individuals with a degree are less likely than individuals without a degree to have wage below \$y per hour.

Extending to any finite $K \geq 2$ with FWER α , we apply their Method 5 to each of the $K - 1$ pairs of consecutive $(x, x + 1)$ values with Bonferroni-adjusted level $\alpha/(K - 1)$. In the special case $K = 2$, there is no adjustment because $K - 1 = 1$. As discussed in our introduction, the Bonferroni adjustment errs on the conservative side, but not too egregiously because the dependence across X is negative rather than positive. (It is important that we only use the Bonferroni adjustment across X and not across Y , which has an infinite number of points with strong positive dependence.) Theorem 4 formalizes this.

Theorem 4. *Consider the multiple testing procedure that applies Method 5 of Goldman*

and Kaplan (2018) with level $\alpha/(K-1)$ a total of $K-1$ times: first to $H_{0x,y}$ in (14) for $(x,y) \in \{1\} \times \mathbb{R}$, second to $H_{0x,y}$ for $(x,y) \in \{2\} \times \mathbb{R}$, etc., up to $(x,y) \in \{K-1\} \times \mathbb{R}$. Under Assumption A1 with strictly increasing conditional CDFs $F_x(\cdot)$, this procedure has strong control of finite-sample FWER at level α .

As in Section 2.4, we can invert the multiple testing procedure into confidence sets for the true set $\mathcal{T} \equiv \{(x,y) : \theta_{x,y} \leq 0\}$ from (11). Again, the outer confidence set collects the (x,y) corresponding to $H_{0x,y} : \theta_{x,y} \leq 0$ that are not rejected, whereas the inner confidence set collects the (x,y) corresponding to rejected $H_{0x,y}^\geq : \theta_{x,y} \geq 0$. The finite-sample FWER control translates to finite-sample coverage probabilities.

Theorem 5. *Consider the inner confidence set $\widehat{\mathcal{CS}}_{inner}$ that collects all pairs of (x,y) for which the reversed null $H_{0x,y}^\geq : \theta_{x,y} \geq 0$ is rejected by the multiple testing procedure in Theorem 4, and the outer confidence set $\widehat{\mathcal{CS}}_{outer}$ that collects all pairs of (x,y) for which the original null $H_{0x,y} : \theta_{x,y} \leq 0$ is not rejected. Under Assumption A1 with strictly increasing conditional CDFs $F_x(\cdot)$, these confidence sets have finite-sample coverage probability:*

$$P(\widehat{\mathcal{CS}}_{inner} \subseteq \mathcal{T}) \geq 1 - \alpha, \quad P(\widehat{\mathcal{CS}}_{outer} \supseteq \mathcal{T}) \geq 1 - \alpha.$$

Also like before, if instead of \mathcal{T} in (11) we are interested in its complement $\mathcal{T}^c \equiv \{(x,y) : \theta_{x,y} > 0\}$, then the outer CS is the complement of the inner CS for \mathcal{T} , and the inner CS is the complement of the outer CS for \mathcal{T} . That is,

$$\begin{aligned} P(\widehat{\mathcal{CS}}_{inner}^c \supseteq \mathcal{T}^c) &= P(\widehat{\mathcal{CS}}_{inner} \subseteq \mathcal{T}) \geq 1 - \alpha, \\ P(\widehat{\mathcal{CS}}_{outer}^c \subseteq \mathcal{T}^c) &= P(\widehat{\mathcal{CS}}_{outer} \supseteq \mathcal{T}) \geq 1 - \alpha. \end{aligned}$$

5 Empirical illustrations

To demonstrate our methodology, we apply it to ordinal and continuous outcome examples in Sections 5.1 and 5.2, respectively. We use publicly available data from the 2022 Understanding Society study (University of Essex, Institute for Social and Economic Research, 2024), previously known as the UK Household Longitudinal Study. We provide codes and files to replicate our results that are available online.²

In the appendix, we provide additional empirical analyses. Appendix B examines an ordinal outcome for life satisfaction, generally finding that it stochastically increases with

²Access via <https://qianjoewu.github.io/>. All code is in R (R Core Team, 2023), with help from the `MASS` package (Venables and Ripley, 2002) for multivariate normal random vector generation, and from the packages `ggplot2` (Wickham, 2016), `ggmosaic` (Jeppson and Hofmann, 2023; Jeppson, Hofmann, and Cook, 2023), and `scales` (Wickham, Pedersen, and Seidel, 2023) for plotting.

education, but with more mixed results than general health, especially at the highest satisfaction level. Appendix C examines a continuous measure of mental health, finding essentially no statistically significant associations with education.

5.1 General health and education

Numerous studies have examined the relationship between general health and educational attainment. Surveying numerous empirical papers, Grossman (2006, §4.1) concludes that education is the most important variable associated with “health,” across a variety of measures including self-reported general health. Individuals with higher education generally report better health. Potential explanatory mechanisms include education affecting preferences (like for smoking, exercise, and time preferences/patience), increasing adherence to curative treatments, and increasing health knowledge and literacy (Grossman, 2006, §4.1). Cutler and Lleras-Muney (2010) focus on the relationship between education and “health behaviors” (like smoking and exercise), finding roles for knowledge, cognitive ability, income, insurance, and social networks, all of which increase with education and in turn improve health.

This relationship between health and education potentially differs across different margins of education as well as health category. For example, the relationship between high-school dropouts and graduates may differ from the relationship between bachelor’s and graduate degree holders. In addition, the relationship for “very good” health status may differ from the relationship for “poor” health status. Our multiple testing method can detect such heterogeneous patterns because it is nonparametric and does not impose any restrictions. We assess evidence for where the level of general health is stochastically increasing in education.

We use the following variables from the Understanding Society 2022 data (University of Essex, Institute for Social and Economic Research, 2024). General health variable Y (`lmscsf1`) is an ordinal measure with five categories: “poor,” “fair,” “good,” “very good,” and “excellent,” respectively coded as 1, 2, 3, 4, and 5. Education variable X (recoded from `lmsniscd11_dv`) has four categories: “no high school degree” (highest completed level is lower secondary), “HS degree only” (highest completed is upper secondary), “bachelor’s degree” (in the codebook: “Bachelor or equivalent”), and “graduate degree” (codebook: master’s or doctorate or equivalent of either), respectively coded as 1, 2, 3, and 4. We restrict the age range to 30–65 years old for a fair comparison of people with various levels of education including graduate degrees, and to concentrate on those who are eligible for the workforce. For simplicity, observations with missing values are dropped and weights are not used. Age is non-missing for 99.96% of the sample, and among those aged 30–65, education and general

health are both observed for 92.6% of observations.

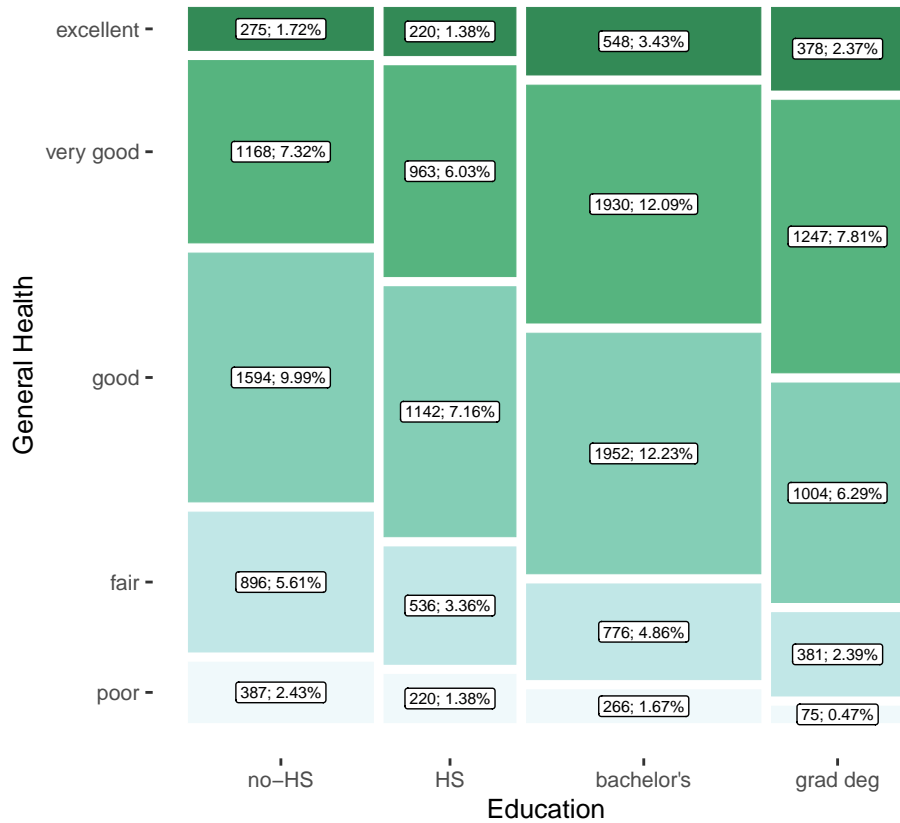


Figure 1: Mosaic plot of general health and education.

Figure 1 visualizes the data on general health and education in a mosaic plot. Each column corresponds to a category of X , and its width is proportional to the number of observations in that category. For example, the no-HS column is significantly wider than the HS column, showing that more individuals in the data do not have a high-school degree than have (only) a high-school degree. Each cell's area is proportional to the sample proportion of individuals with that particular value of (X, Y) ; the text label shows that proportion as a percentage, along with the corresponding number of observations in that cell. Because the joint probability of $(X, Y) = (x, y)$ is the product of the marginal $X = x$ probability and the conditional Y probability given $X = x$, within a column corresponding to $X = x$, each cell's height is proportional to the proportion of observations with the corresponding y value within the $X_i = x$ subsample. These conditional probabilities are scaled such that the full height of each column is probability one (or 100%). For example, in the no-HS column, the “good” health category is tall because $P(\text{good} \mid \text{no-HS})$ is relatively large. This further implies that the breaks between cells within a column show the conditional CDF values. For example, in the no-HS column, the top of the “good” cell is around $2/3$ between the bottom

and top of the column, indicating the empirical conditional CDF is around $\hat{F}_1(3) \approx 2/3$.

The mosaic plot shows stochastic monotonicity in the sample. Mathematically, for each level y , the mosaic plot shows $\hat{F}_4(y) \leq \hat{F}_3(y) \leq \hat{F}_2(y) \leq \hat{F}_1(y)$. Qualitatively, this says that the distribution of general health is “better” at higher levels of education in the sense of first-order stochastic dominance. However, we wish to learn not only about the sample but about the population relationship. Our methods help assess the strength of the evidence that these patterns also hold in the population.

Table 2: Test statistics for general health versus education ($\hat{c}_{0.05} = 2.62$).

General health y	Education category x		
	1 (vs. 2)	2 (vs. 3)	3 (vs. 4)
1: poor	−2.86	−4.16	−6.05
2: fair (or below)	−4.96	−5.85	−5.13
3: good (or below)	−4.41	−6.24	−6.58
4: very good (or below)	−1.30	−4.66	−3.12

True set: points where general health is better at next-higher education level, $\{(x, y) : F_x(y) \geq F_{x+1}(y)\}$. Gray shading: outer confidence set. Bold: inner confidence set. Confidence level 95%. Critical value computed using $N = 100,000$ random draws.

Table 2 shows our inference results, which can be interpreted as follows. The numbers are the t -statistics from (7) for testing the null hypotheses $H_{0x,y} : F_x(y) \geq F_{x+1}(y)$. Such a t -statistic is negative when $\hat{F}_x(y) > \hat{F}_{x+1}(y)$, indicating a smaller proportion of individuals with health status y or below in the $X_i = x + 1$ subsample than in the $X_i = x$ subsample, meaning better general health at the higher education level $x + 1$ than at x . If the t -statistic is negative and below the negative critical value $-\hat{c}_{0.05} = -2.62$, then there is strong evidence that this is true in the population, too. Consequently, when $\hat{t}_{x,y}(0) < -\hat{c}_{0.05}$, that (x, y) is included in the 95% inner confidence set, which collects the points with strong evidence of general health increasing with education; such t -statistics are **bold**. The 95% outer confidence set includes (x, y) as long as there is no strong evidence in the opposite direction, i.e., as long as $\hat{t}_{x,y}(0) < \hat{c}_{0.05}$; such t -statistics are **shaded gray**. Note the $\hat{t}_{x,y}(0) < -\hat{c}_{0.05}$ are thus **bold and shaded** because the outer confidence set contains the inner confidence set. For example, the table’s bottom-center number $(\hat{F}_3(4) - \hat{F}_2(4))/\hat{s}_{2,4} = -4.66$ is the t -statistic for $(x, y) = (2, 4)$, comparing individuals with only a high-school degree ($x = 2$) to those with only a bachelor’s degree ($x = 3$). Using Method 1, this means we do not reject the null that $F_2(4) \geq F_3(4)$ (hence the gray shading of the -4.66 cell), but we do reject the reversed null that $F_2(4) \leq F_3(4)$ in favor of the conclusion $F_2(4) > F_3(4)$ because $-4.66 < -\hat{c}_{0.05} = -2.62$ (hence bold **−4.66** in the table). Because there are only $J = 5$ categories of Y , we can also

interpret the conclusion $F_2(4) > F_3(4)$ as $P(Y = 5 \mid X = 2) < P(Y = 5 \mid X = 3)$, i.e., a lower probability of excellent health status in the high-school-only subpopulation than in the bachelor’s-only subpopulation.

Table 2 shows substantial and consistent evidence that higher education is associated with better general health. Almost all computed t -statistics are negative and below the negative 5% FWER critical value, collectively forming a relatively large 95% inner confidence set for the true set of points $\{(x, y) : F_x(y) \geq F_{x+1}(y)\}$ where general health is better at the higher education level. The outer confidence set includes all cells, including the t -statistic at $(x, y) = (1, 4)$ that is negative but not below $-\hat{c}_{0.05}$. This reflects the closer empirical conditional CDF values for the high-school group versus the no-high-school group for very good (or below) health status.

Overall, our results show strong evidence that general health is stochastically increasing in education for all education levels, though slightly weaker evidence that the high-school subpopulation stochastically dominates the no-high-school subpopulation. Compared to a global test that simply fails to reject the global null that health stochastically increases with education, our methodology provides much stronger and more precise results.

5.2 Physical health and education

Complementing our analysis of general health, we separately examine physical and mental health. Specifically, we consider the physical component summary/score (PCS) and mental component summary/score (MCS) that summarize responses to the SF-12, the widely used 12-item Short-Form Health Survey. As the names suggest, PCS focuses on physical well-being, while MCS captures aspects of mental health. Studying the relationship between PCS/MCS and education helps our descriptive understanding of health inequality. This can also suggest potential causal mechanisms if the relationship differs across educational levels and/or health levels. Here we focus on PCS, while Appendix C has analogous results for MCS.

We use the following variables. Our continuous Y variable is PCS (based on `lmn_sf12pcs_dv`), with higher values indicating better physical health. We rescale the original value \tilde{Y} to have range $[0, 100]$ by taking $Y = 100[\tilde{Y} - \min(\tilde{Y})]/[\max(\tilde{Y}) - \min(\tilde{Y})]$. PCS summarizes responses to questions related to physical functioning, bodily pain, general health perceptions, and physical role limitations. Our ordinal X variable is education (recoded from `lmn_nisced11_dv`). It has four categories: “no high school degree,” “HS degree only,” “bachelor’s degree,” and “graduate degree” (master’s, professional, or doctoral), respectively coded as 1, 2, 3, and 4. As in Section 5.1, we restrict the age range to 30–65 years old and

drop missing values.

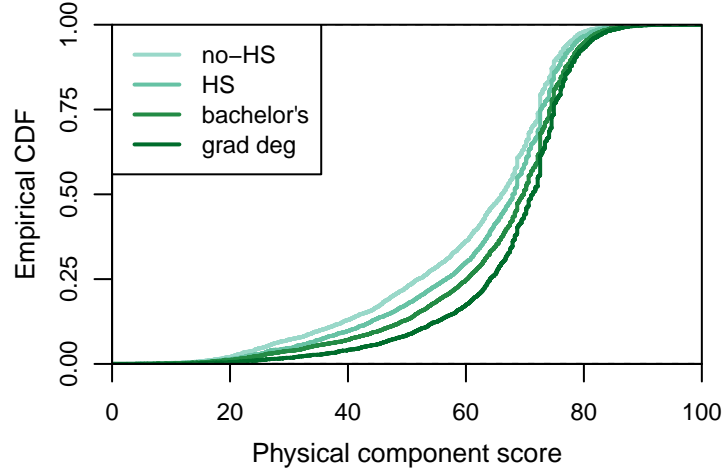


Figure 2: Empirical CDF of PCS by different education categories.

Figure 2 presents empirical CDFs of PCS for each of the four education categories. The graph shows that the empirical CDF curves are ordered such that $\hat{F}_4(\cdot) \leq \hat{F}_3(\cdot) \leq \hat{F}_2(\cdot) \leq \hat{F}_1(\cdot)$, indicating that higher education is associated with higher PCS, meaning better physical health. However, this stochastic monotonicity in the sample does not necessarily imply stochastic monotonicity in the population. Our methods help quantify the statistical uncertainty about the population relationship.

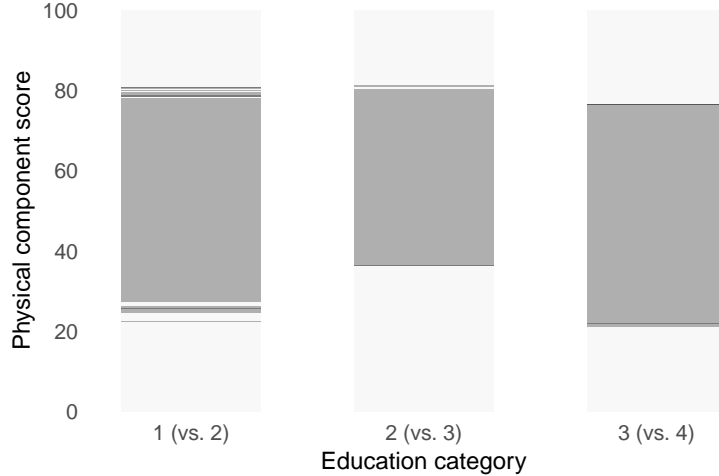


Figure 3: Inner (dark gray) and outer (light gray) 95% confidence sets.

Figure 3 shows our inner and outer confidence sets. We use a 95% confidence level for each, so given Theorem 5, our confidence sets have (at least) 95% coverage not only asymptotically but in finite samples. The true set contains all points (x, y) where the population conditional CDFs are consistent with PCS stochastically increasing in education,

$F_{x+1}(y) \leq F_x(y)$. The inner confidence set consists of the dark gray regions, which provide a conservative “estimate” of this true set, in the sense that there is a 95% ex ante (frequentist) probability of sampling a dataset for which the inner confidence set is contained within the true set. That is, the dark gray regions show where there is especially strong evidence supporting stochastic increasingness. The outer confidence set includes the dark gray regions as well as the light gray regions.

Figure 3 shows generally strong evidence of physical health stochastically increasing with education. The inner confidence set extends across a wide range of PCS values for each education comparison. This is true especially for PCS values from around 50 to 80, which covers most of the distribution: in the sample, 50 is the 8th percentile, and 80 is the 85th percentile. That is, there is strong evidence of higher education associated with better physical health across most of the physical health distribution. Although uncertainty is higher in the upper and lower tails, the outer confidence set still includes the full support $[0, 100]$ for each education category. This means that there is no statistically significant evidence *against* PCS stochastically increasing with education.

In sum, Figure 3 provides more granular and precise evidence of physical health stochastically increasing in education, compared to a global test that would merely “not reject” the null hypothesis of stochastic increasingness.

6 Conclusion

We have proposed multiple testing procedures and confidence sets to provide a richer assessment of stochastic monotonicity, by separately considering each conditional CDF inequality that together comprise stochastic monotonicity, for continuous, discrete, or ordinal outcomes. In future work, to complement our frequentist confidence sets, we plan to develop Bayesian credible sets and assess their frequentist properties. For the related *joint* testing problem, the frequentist test is much more conservative than a Bayesian test even asymptotically (e.g., Kaplan and Zhuo, 2021; Kline, 2011), but this does not translate directly to multiple testing and related confidence sets. Deriving multiple testing procedures for continuous X would also be valuable.

Acknowledgments

Thanks to the following for their helpful feedback: Alyssa Carlson, Huayan Geng, Wei Lan, Xing Ling, Zack Miller, Shawn Ni, Xiuli Sun, journal reviewers, and participants in the Asian Meeting of the Econometric Society (2022, CUHK/virtual), Missouri Valley Economic As-

sociation annual conference (2021, virtual), Economics Graduate Student Conference (2022, Washington University in St. Louis), University of Missouri Economics Research Workshop, Chinese Economists Society annual conference (2022, Guizhou/virtual), and Midwest Econometrics Group (2022, Michigan State).

References

- Armstrong, Timothy B. and Shu Shen. 2023. “Inference on optimal treatment assignments.” *The Japanese Economic Review* 74 (4):471–500. URL <https://doi.org/10.1007/s42973-023-00138-1>.
- Banks, James and Fabrizio Mazzonna. 2012. “The Effect of Education on Old Age Cognitive Abilities: Evidence from a Regression Discontinuity Design.” *The Economic Journal* 122 (560):418–448. URL <https://www.jstor.org/stable/41494443>.
- Benjamini, Yoav and Yosef Hochberg. 1995. “Controlling the false discovery rate: a practical and powerful approach to multiple testing.” *Journal of the Royal Statistical Society: Series B* 57 (1):289–300. URL <https://www.jstor.org/stable/2346101>.
- Blanchflower, David G. and Andrew J. Oswald. 2004. “Well-being over time in Britain and the USA.” *Journal of Public Economics* 88 (7–8):1359–1386. URL [https://doi.org/10.1016/S0047-2727\(02\)00168-8](https://doi.org/10.1016/S0047-2727(02)00168-8).
- Bond, Timothy N. and Kevin Lang. 2019. “The Sad Truth About Happiness Scales.” *Journal of Political Economy* 127 (4):1629–1640. URL <https://doi.org/10.1086/701679>.
- Canay, Ivan A. and Azeem M. Shaikh. 2017. “Practical and Theoretical Advances in Inference for Partially Identified Models.” In *Advances in Economics and Econometrics: Eleventh World Congress, Econometric Society Monographs*, vol. 2, edited by Ariel Pakes, Bo Honoré, Larry Samuelson, and Monika Piazzesi, chap. 9. Cambridge: Cambridge University Press, 271–306. URL <https://doi.org/10.1017/9781108227223.009>.
- Chetverikov, Denis, Daniel Wilhelm, and Dongwoo Kim. 2021. “An Adaptive Test of Stochastic Monotonicity.” *Econometric Theory* 37 (3):495–536. URL <https://doi.org/10.1017/S0266466620000225>.
- Clark, Andrew E., Paul Frijters, and Michael A. Shields. 2008. “Relative Income, Happiness, and Utility: An Explanation for the Easterlin Paradox and Other Puzzles.” *Journal of Economic Literature* 46 (1):95–144. URL <https://doi.org/10.1257/jel.46.1.95>.
- Cutler, David M. and Adriana Lleras-Muney. 2010. “Understanding differences in health behaviors by education.” *Journal of Health Economics* 29 (1):1–28. URL <https://doi.org/10.1016/j.jhealeco.2009.10.003>.
- Davidson, Russell and Jean-Yves Duclos. 2013. “Testing for Restricted Stochastic Dominance.” *Econometric Reviews* 32 (1):84–125. URL <https://doi.org/10.1080/07474938.2012.690332>.
- de New, Sonja C., Stefanie Schurer, and Dominique Sulzmaier. 2021. “Gender differences in the lifecycle benefits of compulsory schooling policies.” *European Economic Review* 140:103910. URL <https://doi.org/10.1016/j.euroecorev.2021.103910>.
- Deaton, Angus. 2008. “Income, Health, and Well-Being around the World: Evidence from the Gallup World Poll.” *Journal of Economic Perspectives* 22 (2):53–72. URL <https://doi.org/10.1215/08901431-2008-003>.

- [//doi.org/10.1257/jep.22.2.53](https://doi.org/10.1257/jep.22.2.53).
- Goldman, Matt and David M. Kaplan. 2018. “Comparing distributions by multiple testing across quantiles or CDF values.” *Journal of Econometrics* 206 (1):143–166. URL <https://doi.org/10.1016/j.jeconom.2018.04.003>.
- Grossman, Michael. 2006. “Education and nonmarket outcomes.” In *Handbook of the Economics of Education*, vol. 1, edited by E. Hanushek and F. Welch, chap. 10. Elsevier, 577–633. URL [https://doi.org/10.1016/S1574-0692\(06\)01010-5](https://doi.org/10.1016/S1574-0692(06)01010-5).
- Jeppson, Haley and Heike Hofmann. 2023. “Generalized Mosaic Plots in the ggplot2 Framework.” *The R Journal* 14 (4):50–78. URL <https://doi.org/10.32614/RJ-2023-013>.
- Jeppson, Haley, Heike Hofmann, and Dianne H. Cook. 2023. *ggmosaic: Mosaic Plots in the 'ggplot2' Framework*. URL <https://haleyjeppson.github.io/ggmosaic/>. R package version 0.3.4.
- Kaplan, David M. 2024. “Inference on Consensus Ranking of Distributions.” *Journal of Business & Economic Statistics* 42 (3):839–850. URL <https://doi.org/10.1080/07350015.2023.2252040>.
- Kaplan, David M. and Wei Zhao. 2023. “Comparing latent inequality with ordinal data.” *The Econometrics Journal* 26 (2):189–214. URL <https://doi.org/10.1093/ectj/utac030>.
- Kaplan, David M. and Longhao Zhuo. 2021. “Frequentist properties of Bayesian inequality tests.” *Journal of Econometrics* 221 (1):312–336. URL <https://doi.org/10.1016/j.jeconom.2020.05.015>.
- Kline, Brendan. 2011. “The Bayesian and frequentist approaches to testing a one-sided hypothesis about a multivariate mean.” *Journal of Statistical Planning and Inference* 141 (9):3131–3141. URL <https://doi.org/10.1016/j.jspi.2011.03.034>.
- Lee, Sokbae, Oliver Linton, and Yoon-Jae Whang. 2009. “Testing for Stochastic Monotonicity.” *Econometrica* 77 (2):585–602. URL <https://www.jstor.org/stable/40263876>.
- Lehmann, E. L. and Joseph P. Romano. 2005a. “Generalizations of the Familywise Error Rate.” *Annals of Statistics* 33 (3):1138–1154. URL <https://projecteuclid.org/euclid.aos/1120224098>.
- . 2005b. *Testing Statistical Hypotheses*. Springer Texts in Statistics. Springer, 3rd ed. URL <https://books.google.com/books?id=Y7vSVW3ebSwC>.
- Manski, Charles F. 1997. “Monotone Treatment Response.” *Econometrica* 65 (6):1311–1334. URL <https://doi.org/10.2307/2171738>.
- Oreopoulos, Philip. 2007. “Do dropouts drop out too soon? Wealth, health and happiness from compulsory schooling.” *Journal of Public Economics* 91 (11–12):2213–2229. URL <https://doi.org/10.1016/j.jpubeco.2007.02.002>.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org>.
- Romano, Joseph P., Azeem M. Shaikh, and Michael Wolf. 2010. “Multiple testing.” In *The New Palgrave Dictionary of Economics*, edited by Steven N. Durlauf and Lawrence E. Blume. London: Palgrave Macmillan, online ed., 1–5. URL <https://doi.org/10.1057/9780230226203.3826>.
- Seo, Juwon. 2018. “Tests of stochastic monotonicity with improved power.” *Journal of Econometrics* 207 (1):53–70. URL <https://doi.org/10.1016/j.jeconom.2018.04.004>.
- Small, Dylan S., Zhiqiang Tan, Roland R. Ramsahai, Scott A. Lorch, and M. Alan Brookhart. 2017. “Instrumental Variable Estimation with a Stochastic Monotonicity Assumption.”

- Statistical Science* 32 (4):561–579. URL <https://doi.org/10.1214/17-STS623>.
- University of Essex, Institute for Social and Economic Research. 2024. “Understanding Society: Calendar Year Dataset, 2022.” URL <https://doi.org/10.5255/UKDA-SN-9333-1>. Public-use data file and documentation, accessed 2025.
- van der Vaart, Aad W. 1998. *Asymptotic Statistics*. Cambridge: Cambridge University Press. URL <https://books.google.com/books?id=UEuQEM5RjWgC>.
- Venables, W. N. and B. D. Ripley. 2002. *Modern Applied Statistics with S*. New York: Springer, fourth ed. URL <https://www.stats.ox.ac.uk/pub/MASS4/>.
- Wickham, Hadley. 2016. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. URL <https://ggplot2.tidyverse.org>.
- Wickham, Hadley, Thomas Lin Pedersen, and Dana Seidel. 2023. *scales: Scale Functions for Visualization*. URL <https://scales.r-lib.org>. R package version 1.3.0, <https://github.com/r-lib/scales>.
- Zhao, Wei and David M. Kaplan. 2024. “Multiple Testing of a Function’s Monotonicity.” Working paper available at <https://kaplandm.github.io/>.

A Proofs

A.1 Proof of Lemma 1

Proof. Under Assumption A1 and the central limit theorem (e.g., van der Vaart, 1998, Prop. 2.27), as $n_x \rightarrow \infty$,

$$\sqrt{n_x}(\hat{F}_x(y) - F_x(y)) \xrightarrow{d} N(0, F_x(y)[1 - F_x(y)]). \quad (15)$$

Let $\hat{\mathbf{A}}$ be the estimator of vector \mathbf{A} that contains all conditional CDFs,

$$\mathbf{A} \equiv (F_1(1), F_1(2), \dots, F_1(J-1), \dots, F_K(1), F_K(2), \dots, F_K(J-1))'. \quad (16)$$

Under Assumption A1, by the continuous mapping theorem (CMT) (e.g., van der Vaart, 1998, Thm. 2.3),

$$\sqrt{n_1}(\hat{\mathbf{A}} - \mathbf{A}) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}), \quad (17)$$

where $\mathbf{0}$ is a vector of zeros and \mathbf{V} is the $(JK - K) \times (JK - K)$ block diagonal matrix

$$\mathbf{V} \equiv \begin{bmatrix} \mathbf{V}^1 & 0 & \dots & 0 \\ 0 & \mathbf{V}^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{V}^K \end{bmatrix} \quad (18)$$

where each $\underline{\mathbf{V}}^k$ is a $(J-1) \times (J-1)$ matrix with entry in row i , column j denoted

$$V_{ij}^k = \gamma_x F_x(\min\{i, j\})[1 - F_x(\max\{i, j\})].$$

We apply the delta method (e.g., van der Vaart, 1998, Thm. 3.1) to derive the asymptotic distribution of the full vector $\hat{\boldsymbol{\theta}}$ for testing the stochastic monotonicity inequalities. For $x \in \{1, 2, \dots, K-1\}$ and $y \in \{1, 2, \dots, J-1\}$,

$$\hat{\boldsymbol{\theta}} \equiv (\hat{\boldsymbol{\theta}}_1, \hat{\boldsymbol{\theta}}_2, \dots, \hat{\boldsymbol{\theta}}_{J-1})', \quad (19)$$

where each $\hat{\boldsymbol{\theta}}_j \equiv (\hat{\theta}_{K-1,j}, \hat{\theta}_{K-2,j}, \dots, \hat{\theta}_{1,j})$ and $\hat{\theta}_{x,y} \equiv \hat{F}_{x+1}(y) - \hat{F}_x(y)$ is the estimator of $\theta_{x,y}$ defined in (4). Then,

$$\sqrt{n_1}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \xrightarrow{d} N(\mathbf{0}, \underline{\mathbf{W}}), \quad \underline{\mathbf{W}} \equiv \underline{\mathbf{H}}' \underline{\mathbf{V}} \underline{\mathbf{H}}, \quad (20)$$

where $\underline{\mathbf{H}} = \frac{\partial}{\partial \mathbf{A}} \boldsymbol{\theta}'$ is the partial derivative of the vector $\boldsymbol{\theta}'$ with respect to \mathbf{A} in (16), and $\underline{\mathbf{V}}$ is from (18). Note each element of $\underline{\mathbf{H}}$ is either -1 , 0 , or 1 ; more specifically, within each column of $\underline{\mathbf{H}}$ (the derivative of a particular $\theta_{x,y}$ with respect to \mathbf{A}), one element is -1 , one element is 1 , and the rest equal zero.

From (20), for each $\hat{\theta}_{x,y}$, using $n_1/n_x \rightarrow \gamma_x$ from A1,

$$\begin{aligned} \sqrt{n_1}(\hat{\theta}_{x,y} - \theta_{x,y}) &= \sqrt{n_1}\{[\hat{F}_{x+1}(y) - \hat{F}_x(y)] - [F_{x+1}(y) - F_x(y)]\} \\ &= \overbrace{\sqrt{n_1/n_{x+1}}}^{\rightarrow \sqrt{\gamma_{x+1}}} \overbrace{\sqrt{n_{x+1}}[\hat{F}_{x+1}(y) - F_{x+1}(y)]}^{\text{use (15)}} - \overbrace{\sqrt{n_1/n_x}}^{\rightarrow \sqrt{\gamma_x}} \overbrace{\sqrt{n_x}[\hat{F}_x(y) - F_x(y)]}^{\text{use (15)}} \\ &\xrightarrow{d} N(0, \sigma_{x,y}^2), \\ \sigma_{x,y} &\equiv \sqrt{\gamma_x F_x(y)[1 - F_x(y)] + \gamma_{x+1} F_{x+1}(y)[1 - F_{x+1}(y)]}, \end{aligned} \quad (21)$$

where the variances are summed because the subsamples for $X_i = x$ and $X_i = x+1$ are independent given A1, and if $W \perp Z$ then $\text{Var}(W + Z) = \text{Var}(W) + \text{Var}(Z)$.

For the standard error of $\hat{\theta}_{x,y}$, which is asymptotically $\sigma_{x,y}/\sqrt{n_1}$ given (21), we use the estimator

$$\hat{s}_{x,y} \equiv \sqrt{\frac{(n_1/n_x)\hat{F}_x(y)(1 - \hat{F}_x(y)) + (n_1/n_{x+1})\hat{F}_{x+1}(y)(1 - \hat{F}_{x+1}(y))}{n_1}}. \quad (22)$$

By the consistency of the conditional CDF estimators and the CMT, as $n_1 \rightarrow \infty$,

$$\begin{aligned} \sqrt{n_1}\hat{s}_{x,y} &= \sqrt{(n_1/n_x)\hat{F}_x(y)(1 - \hat{F}_x(y)) + (n_1/n_{x+1})\hat{F}_{x+1}(y)(1 - \hat{F}_{x+1}(y))} \\ &\xrightarrow{p} \sqrt{\gamma_x F_x(y)(1 - F_x(y)) + \gamma_{x+1} F_{x+1}(y)(1 - F_{x+1}(y))} = \sigma_{x,y}, \end{aligned} \quad (23)$$

the asymptotic standard deviation in (21). Informally, the standard error estimator $\hat{s}_{x,y}$ is “consistent,” meaning formally that $\sqrt{n_1}\hat{s}_{x,y} \xrightarrow{p} \sigma_{x,y}$.

Substituting (22) into (7),

$$\hat{t}_{x,y}(\theta_{x,y}) = \frac{\hat{\theta}_{x,y} - \theta_{x,y}}{\hat{s}_{x,y}} = \frac{\sqrt{n_1}(\hat{\theta}_{x,y} - \theta_{x,y})}{\sqrt{(n_1/n_x)\hat{F}_x(y)(1 - \hat{F}_x(y)) + (n_1/n_{x+1})\hat{F}_{x+1}(y)(1 - \hat{F}_{x+1}(y))}}.$$

Thus, we can write the vector of properly centered t -statistics as $\hat{\mathbf{t}} = \sqrt{n_1}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})\hat{\mathbf{S}}$, where $\sqrt{n_1}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})$ is the left-hand side of (20), and $\hat{\mathbf{S}}$ is a diagonal matrix having elements $1/(\hat{s}_{x,y}\sqrt{n_1})$ matching the corresponding (x, y) from the $\boldsymbol{\theta}$ vector; equivalently, the row i , column j element of the matrix $\hat{\mathbf{S}}$ is $\mathbb{1}\{i = j\}\hat{W}_{ij}$. By (23), $\hat{\mathbf{S}} \xrightarrow{p} \mathbf{S}$, the diagonal matrix with corresponding (x, y) elements $1/\sigma_{x,y}$,

$$\mathbf{S} \equiv \begin{bmatrix} 1/\sigma_{K-1,1} & 0 & \cdots & 0 \\ 0 & 1/\sigma_{K-2,1} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1/\sigma_{1,J-1} \end{bmatrix}. \quad (24)$$

By (20) and CMT, $\hat{\mathbf{t}} \equiv \sqrt{n_1}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})\hat{\mathbf{S}} \xrightarrow{d} N(\mathbf{0}, \mathbf{S}'\mathbf{W}\mathbf{S})$. □

A.2 Proof of Theorem 2

Proof. When testing original hypotheses $H_{0x,y}: F_{x+1}(y) - F_x(y) \leq 0$, the set of true hypotheses is $\mathcal{T} \equiv \{(x, y) : F_{x+1}(y) - F_x(y) \leq 0\}$. For any true $H_{0x,y}$,

$$\hat{t}_{x,y}(0) = \frac{\hat{F}_{x+1}(y) - \hat{F}_x(y)}{\hat{s}_{x,y}} \leq \frac{\hat{F}_{x+1}(y) - \hat{F}_x(y) - \overbrace{[F_{x+1}(y) - F_x(y)]}^{\leq 0 \text{ given true } H_{0x,y}}}{\hat{s}_{x,y}} \equiv \hat{t}_{x,y}(\overbrace{F_{x+1}(y) - F_x(y)}^{\theta_{x,y}}). \quad (25)$$

Thus,

$$\begin{aligned} \text{FWER} &\equiv P(\text{reject } H_{0x,y} \text{ for any } (x, y) \in \mathcal{T}) \\ &= P(\hat{t}_{x,y}(0) > \hat{c}_\alpha \text{ for any } (x, y) \in \mathcal{T}) \\ &= P(\max_{(x,y) \in \mathcal{T}} \overbrace{\hat{t}_{x,y}(0)}^{\text{use (25)}} > \hat{c}_\alpha) \\ &\leq P(\max_{(x,y) \in \mathcal{T}} \hat{t}_{x,y}(\theta_{x,y}) > \hat{c}_\alpha) \\ &\leq P(\max_{\text{any } (x,y)} \hat{t}_{x,y}(\theta_{x,y}) > \hat{c}_\alpha) \end{aligned}$$

$$\begin{aligned}
& \rightarrow \text{P}\left(\max_{\text{any } (x,y)} \overbrace{t_{x,y}}^{\text{from Lemma 1}} > c_\alpha \right) \\
& = \alpha,
\end{aligned} \tag{26}$$

where the second inequality holds because $\mathcal{T} \subseteq \{1, \dots, K-1\} \times \{1, \dots, J-1\}$ (“any” (x, y)), and the convergence holds by Lemma 1 and the CMT (because max is continuous), along with $\hat{c}_\alpha \xrightarrow{p} c_\alpha$, where \hat{c}_α was defined in (10) as the $(1 - \alpha)$ -quantile of the max of a random vector following the $N(\mathbf{0}, \hat{\Sigma})$ distribution (with estimator $\hat{\Sigma}$) and c_α was defined in (9) as the $(1 - \alpha)$ -quantile of the max of $N(\mathbf{0}, \Sigma)$ (with true Σ). This last part follows because $\hat{\Sigma} \xrightarrow{p} \Sigma$ (given A1 and CMT) and applying the CMT (given that the max of a normal vector has a continuous, strictly increasing distribution function and thus continuous quantile function).

The asymptotic FWER bound for testing the reversed hypotheses $H_{0x,y}^\geq: F_{x+1}(y) - F_x(y) \geq 0$ follows essentially the same derivation with the same reason for each step below. Now defining $\mathcal{T}^\geq \equiv \{(x, y) : F_{x+1}(y) - F_x(y) \geq 0\}$,

$$\begin{aligned}
\text{FWER} & \equiv \text{P}(\text{reject } H_{0x,y}^\geq \text{ for any } (x, y) \in \mathcal{T}^\geq) \\
& = \text{P}(\hat{t}_{x,y}(0) < -\hat{c}_\alpha \text{ for any } (x, y) \in \mathcal{T}^\geq) \\
& = \text{P}\left(\min_{(x,y) \in \mathcal{T}^\geq} \hat{t}_{x,y}(0) < -\hat{c}_\alpha \right) \\
& \leq \text{P}\left(\min_{(x,y) \in \mathcal{T}^\geq} \hat{t}_{x,y}(\theta_{x,y}) < -\hat{c}_\alpha \right) \\
& \leq \text{P}\left(\min_{\text{any } (x,y)} \hat{t}_{x,y}(\theta_{x,y}) < -\hat{c}_\alpha \right) \\
& \rightarrow \text{P}\left(\min_{\text{any } (x,y)} t_{x,y} < -c_\alpha \right) = \text{P}\left(- \min_{\text{any } (x,y)} t_{x,y} > c_\alpha \right) = \text{P}\left(\max_{\text{any } (x,y)} t_{x,y} > c_\alpha \right) = \alpha,
\end{aligned}$$

also using the symmetry of the $\mathbf{t} \sim N(\mathbf{0}, \Sigma)$ distribution that implies $\max \mathbf{t} \stackrel{d}{=} -\min \mathbf{t}$. \square

A.3 Proof of Theorem 3

Proof. For the inner CS, we want to derive the probability of the event that the inner CS is contained within \mathcal{T} , which is equivalently characterized as “ $(x, y) \in \widehat{\mathcal{CS}}_{\text{inner}}$ only if $(x, y) \in \mathcal{T}$.” With reversed hypotheses $H_{0x,y}^\geq: \theta_{x,y} \geq 0$ but $\mathcal{T} \equiv \{(x, y) : \theta_{x,y} \leq 0\}$,

$$\begin{aligned}
\text{P}(\widehat{\mathcal{CS}}_{\text{inner}} \subseteq \mathcal{T}) & = \text{P}(\text{reject any } H_{0x,y}^\geq \text{ only if } (x, y) \in \mathcal{T}) \\
& = \text{P}(\text{reject any } H_{0x,y}^\geq \text{ only if } \theta_{x,y} \leq 0) \\
& = 1 - \text{P}(\text{reject any } H_{0x,y} \text{ when } \theta_{x,y} > 0) \\
& \geq 1 - \text{P}(\text{reject any } H_{0x,y} \text{ when } \theta_{x,y} \geq 0)
\end{aligned}$$

$$= 1 - \text{FWER} \geq 1 - \alpha + o(1)$$

because $\text{FWER} \leq \alpha + o(1)$ by Theorem 2.

The outer CS contains all elements in the true set \mathcal{T} if and only if we do not reject any original null hypothesis $H_{0x,y}: \theta_{x,y} \leq 0$ with $(x, y) \in \mathcal{T}$. Similar to the derivation above, the coverage probability is

$$\begin{aligned} \text{P}(\mathcal{T} \subseteq \widehat{\mathcal{CS}}_{outer}) &= \text{P}(\text{no } H_{0x,y} \text{ rejected with } (x, y) \in \mathcal{T}) \\ &= 1 - \text{P}(\text{reject any } H_{0x,y} \text{ with } \theta_{x,y} \leq 0) \\ &= 1 - \text{FWER} \\ &\geq 1 - \alpha + o(1) \end{aligned}$$

again because $\text{FWER} \leq \alpha + o(1)$ by Theorem 2. □

A.4 Proof of Theorem 4

Proof. We apply a typical Bonferroni strategy given that Method 5 of Goldman and Kaplan (2018) controls finite-sample FWER at adjusted level $\alpha/(K-1)$ for any given x . The finite-sample FWER is

$$\begin{aligned} \text{FWER} &\equiv \text{P}(\text{reject any true } H_{0x,y} \text{ over } (x, y) \in \mathcal{X} \times \mathbb{R}) \\ &= \text{P}\left(\bigcup_{x=1}^{K-1} \text{reject any true } H_{0x,y} \text{ over } y \in \mathbb{R}\right) \\ &\leq \sum_{x=1}^{K-1} \overbrace{\text{P}(\text{reject any true } H_{0x,y} \text{ over } y \in \mathbb{R})}^{\leq \alpha/(K-1) \text{ by Goldman and Kaplan (2018)}} \\ &\leq (K-1)\alpha/(K-1) = \alpha. \end{aligned} \quad \square$$

A.5 Proof of Theorem 5

Proof. The proof is identical to that of Theorem 3 except the very final inequality because here we have $\text{FWER} \leq \alpha$ instead of $\text{FWER} \leq \alpha + o(1)$. Thus, instead of $1 - \text{FWER} \geq 1 - \alpha + o(1)$, the final line becomes $1 - \text{FWER} \geq 1 - \alpha$. □

B Life satisfaction and education

Subjective well-being, including satisfaction with life, has been increasingly studied as a core dimension of human welfare (e.g., Blanchflower and Oswald, 2004; Clark, Frijters, and

Shields, 2008; Deaton, 2008). Like other measures of subjective well-being, life satisfaction is ordinal. Using a parametric latent model (like ordered probit) yields fragile results, as detailed by Bond and Lang (2019). Our method respects the ordinal nature of life satisfaction without imposing any parametric assumptions or any cardinal values.

Numerous studies have examined the relationship between educational attainment and life satisfaction, of which we discuss three examples here. Oreopoulos (2007, §5.3) uses compulsory schooling reforms in the UK as instruments and finds that one additional year of education increases the probability of later being satisfied with life by over five percentage points, suggesting that education has substantial long-run benefits for subjective well-being. In contrast, focusing on the 1947 UK compulsory schooling reform and using a regression discontinuity, Banks and Mazzonna (2012, §4) find negative although not statistically significant effects of education on quality of life (measured by the CASP-19 index), reported in their Table 7. Further highlighting heterogeneity, de New, Schurer, and Sulzmaier (2021, §7.1) find different effects within their Australian data, which includes two states’ compulsory schooling reforms whose effects are estimated by difference-in-differences. Their Figure 6 shows that the one-year increase in mandatory schooling increased life satisfaction in the state of Victoria but not South Australia, as well as smaller differences between men and women within each state.

Our study complements these. Their shared strategy to get causal identification from compulsory schooling law changes means the focus is on a relatively low margin of education, such as someone who would have stopped schooling at age 14 if not for the 1947 UK reform that required them to stay in school for one additional year, or the 1964 reform in Victoria, Australia. We do not make any causal claims but study three different and higher margins of education, including the margin between bachelor’s and graduate degrees. Additionally, we allow for heterogeneity across levels of life satisfaction, while accounting for multiple testing. In our study, the tested monotonic relationship could be potentially different when (for example) comparing high-school dropouts and graduates than when comparing bachelor’s and graduate degree holders, and the relationship may also differ across levels of life satisfaction.

We use the following variables from the Understanding Society 2022 data (University of Essex, Institute for Social and Economic Research, 2024). The outcome variable Y is satisfaction with life overall, measured by a self-reported ordinal item with seven categories: “completely dissatisfied,” “mostly dissatisfied,” “somewhat dissatisfied,” “neither sat nor dissat,” “somewhat satisfied,” “mostly satisfied,” and “completely satisfied,” coded as 1 through 7. Education variable X (recoded from `lmn_niscd11_dv`) has four categories: “no high school degree,” “HS degree only,” “bachelor’s degree,” and “graduate degree” (master’s, professional, or doctoral), respectively coded as 1, 2, 3, and 4. We restrict the

age range to 30–65 years old for a fair comparison of people with various levels of education including graduate degrees, and to concentrate on those who are eligible for the workforce. For simplicity, observations with missing values are dropped and weights are not used. Age is non-missing for 99.96% of the sample, and among those aged 30–65, education and life satisfaction are both observed for 92.1% of observations.

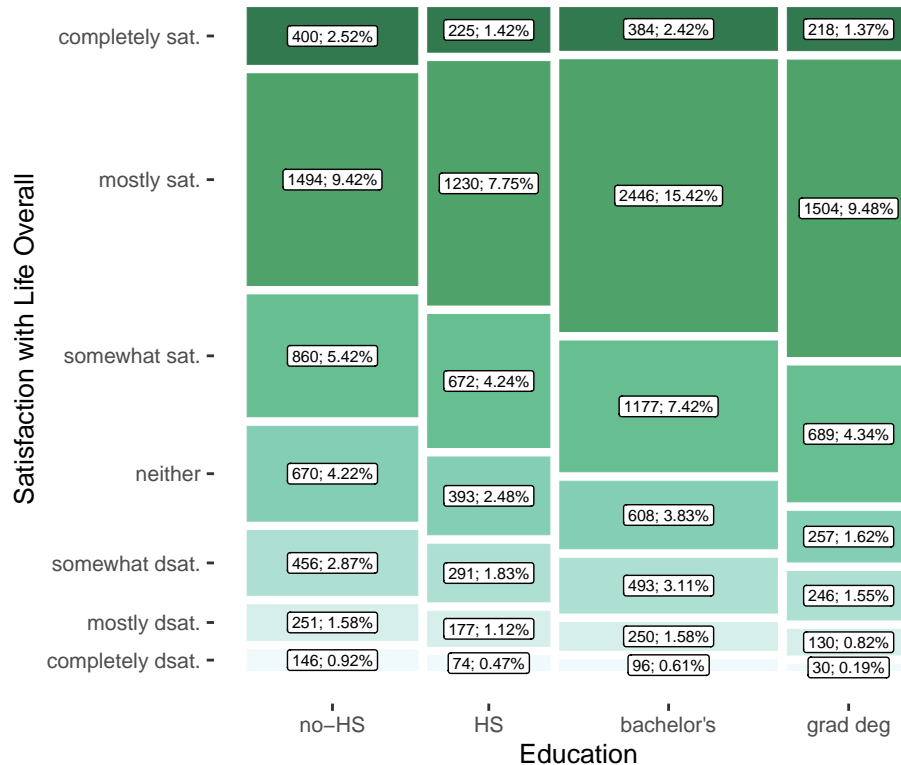


Figure 4: Mosaic plot of life satisfaction and education.

Figure 4 visualizes the data on life satisfaction and education in a mosaic plot. Each column corresponds to a category of X , and its width is proportional to the number of observations in that category. For example, the no-HS column is significantly wider than the HS column, showing that more individuals in the data do not have a high-school degree than have (only) a high-school degree. Each cell's area is proportional to the sample proportion of individuals with that particular value of (X, Y) ; the text label shows that proportion as a percentage, along with the corresponding number of observations in that cell. Because the joint probability of $(X, Y) = (x, y)$ is the product of the marginal $X = x$ probability and the conditional Y probability given $X = x$, within a column corresponding to $X = x$, each cell's height is proportional to the proportion of observations with the corresponding y value within the $X_i = x$ subsample. These conditional probabilities are scaled such that the full height of each column is probability one (or 100%). For example, in the no-HS column,

the “mostly satisfied” cell is tall because $P(\text{mostly sat'd} \mid \text{no-HS})$ is relatively large. This further implies that the breaks between cells within a column show the conditional CDF values. For example, in the no-HS column, the top of the “somewhat dsat.” cell is around $1/5$ between the bottom and top of the column, indicating the empirical conditional CDF is around $\hat{F}_1(3) \approx 1/5$.

The mosaic plot show that stochastic monotonicity is not quite satisfied even in the sample. For each level y at or below “somewhat satisfied,” the mosaic plot shows $\hat{F}_4(y) \leq \hat{F}_3(y) \leq \hat{F}_2(y) \leq \hat{F}_1(y)$. Visually, higher education generally corresponds to greater shares of individuals in higher satisfaction categories. However, in the top level of satisfaction, this relationship does not hold: the lowest education group has the highest sample proportion in the highest “completely satisfied” category, and the next-lowest education group has the next-highest such proportion. However, we wish to learn not only about the sample, but also about the population relationship. Our methods help assess the strength of the evidence that these patterns hold in the population.

Table 3: Test statistics for life satisfaction versus education ($\hat{c}_{0.05} = 2.73$).

Life satisfaction y	Education category x		
	1 (vs. 2)	2 (vs. 3)	3 (vs. 4)
1: completely dsat.	−2.54	−1.99	− 3.12
2: mostly dsat. (or below)	−1.63	− 3.11	−2.19
3: somewhat dsat. (or below)	−2.43	− 2.74	− 2.78
4: neither (or below)	− 4.58	− 3.91	− 5.21
5: somewhat sat. (or below)	− 2.74	− 3.87	− 3.68
6: mostly sat. (or below)	3.09	0.53	−0.09

True set: points where life satisfaction is better at next-higher education level, $\{(x, y) : F_x(y) \geq F_{x+1}(y)\}$. Gray shading: outer confidence set. Bold: inner confidence set. Confidence level 95%. Critical value computed using $N = 100,000$ random draws.

Table 3 shows our inference results, which can be interpreted as follows. The numbers are the t -statistics for testing the null hypotheses $H_{0x,y} : F_x(y) \geq F_{x+1}(y)$. If they exceed $\hat{c}_{0.05}$, then we reject this null, and that (x, y) is not in the outer confidence set; if instead the t -statistic is below $\hat{c}_{0.05}$, then the null is not rejected, so that (x, y) is included in the outer confidence set and shaded gray. If the t -statistic is even below $-\hat{c}_{0.05}$, then we reject the reversed null hypothesis, and (x, y) is also included in the inner confidence set; such values are **bold and shaded**. For example, the bottom-left number (3.09) in the table is for $x = 1$ and $y = 6$, corresponding to the t -statistic for the null hypothesis that $F_1(6) \geq F_2(6)$, i.e., that the population probability of being “mostly satisfied” or below ($Y \leq 6$) is higher

among individuals who have no-high school degree ($x = 1$) than among individuals who have a high school degree only ($x = 2$), or equivalently that the probability of “completely satisfied” ($Y = 7$) is higher for high-school-only ($x = 2$) than for no-high-school ($x = 1$). The corresponding t -statistic $\hat{t}_{x,y}(0)$ from (7) is $\hat{t}_{1,6}(0) = (\hat{F}_2(6) - \hat{F}_1(6))/\hat{s}_{1,6} = 3.09$, which is greater than $\hat{c}_{0.05} = 2.73$ (hence not bold, and not even shaded gray), so we reject $F_1(6) \geq F_2(6)$ in favor of $F_1(6) < F_2(6)$, or equivalently $P(Y = 7 \mid X = 1) > P(Y = 7 \mid X = 2)$. The top-right number is $\hat{t}_{3,1}(0) = -3.12$, less than $-\hat{c}_{0.05} = -2.73$ (hence in shaded gray and also bold), showing strong evidence of a lower probability of being “completely dissatisfied” ($Y \leq 1$) among individuals who have a graduate degree ($x = 4$) than among individuals who have a bachelor degree only ($x = 3$).

Table 3 generally supports life satisfaction increasing with education, but with some important caveats. In total, 10 out of the 18 points are included in the inner confidence set, where there is strong evidence of life satisfaction increasing with education. Most of these points are in the middle categories, “somewhat dissatisfied” to “somewhat satisfied.” The relatively weaker evidence in the lowest two categories is due to the smaller number of observations in those categories, as seen in the Figure 4 mosaic plot. Besides these 10 points in the inner confidence set, six of the other points have a negative t -statistic, but not below $-\hat{c}_{0.05}$; for such points, there is only suggestive evidence of life satisfaction increasing with education. One point has a positive t -statistic that is below $\hat{c}_{0.05}$: the point estimate suggests life satisfaction decreasing with education, but there is enough uncertainty that it is still consistent with life satisfaction increasing with education. The final point $(x, y) = (1, 6)$ has a t -statistic negative enough to reject life satisfaction increasing with education, so it is not even included in the outer confidence set. These last two points with positive t -statistics are both for $Y \leq 6$, which is equivalent to $P(Y = 7) = 1 - P(Y \leq 6)$. That is, the data show strong evidence that the probability of being “completely satisfied” is higher for the no-high-school group ($x = 1$) than the high-school group ($x = 2$), and weaker evidence that the probability of “completely satisfied” is higher for the high-school group ($x = 2$) than the bachelor’s group ($x = 3$).

These disaggregated results are more informative than a global test. The global test would simply reject the null that life satisfaction is stochastically increasing in education, due to the strong evidence against the global null at $(x, y) = (1, 6)$. However, the other 17 points are consistent with life satisfaction increasing with education, and at 10 of those points, there is strong evidence that indeed life satisfaction increases with education. Compared to the simple “reject” result of a global test, our method provides richer, more nuanced results.

C Mental health and education

To complement the main analysis on physical health (PCS) in Section 5.2, we present here analogous results using the mental component summary/score (MCS) from the SF-12 instrument. Our continuous Y variable is MCS (based on `lmn_sf12mcs_dv`), with higher values indicating better mental health. As with our PCS variable, we rescale the original MCS value \tilde{Y} to have range $[0, 100]$ by taking $Y = 100[\tilde{Y} - \min(\tilde{Y})]/[\max(\tilde{Y}) - \min(\tilde{Y})]$. MCS summarizes dimensions of mental health including vitality, emotional well-being, social functioning, and psychological distress. Studying its relationship with education provides insight into the socioeconomic gradient in mental well-being.

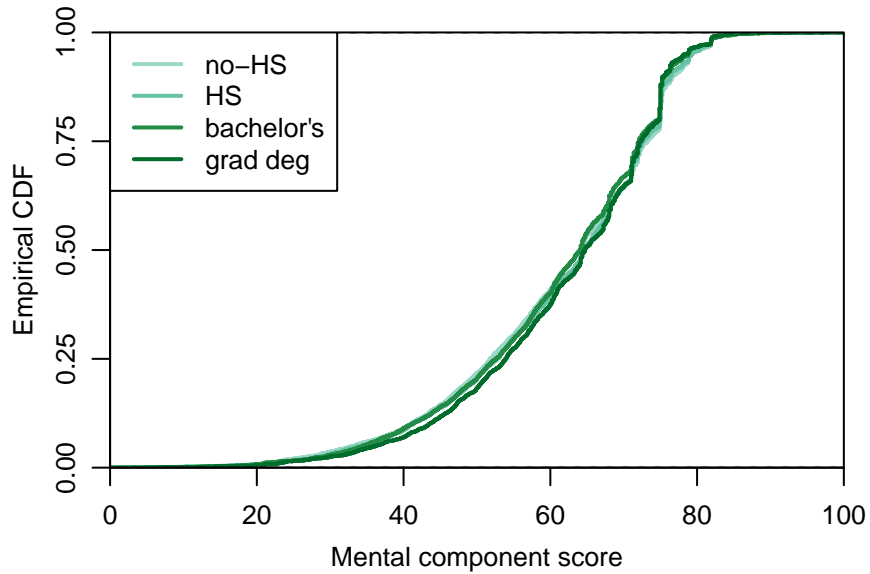


Figure 5: Empirical CDF of MCS by different education categories.

Figure 5 plots the empirical CDFs of MCS for each of the four education categories. Different from the physical component, these CDFs are very close to each other, and there are apparent crossings. Even without any formal statistical inference, the estimated magnitudes of differences across education seem to be small.

Figure 6 reports the results from our multiple testing procedure as inner and outer confidence sets. Each panel corresponds to a pairwise comparison between adjacent education levels, and the vertical axis shows MCS values. The dark gray region (a small line in the 3 vs. 4 education category) is the 95% inner confidence set, where there is strong evidence of mental health stochastically increasing with education. The light gray regions (everywhere) show the 95% outer confidence set, where the hypothesis of stochastic increasingness cannot be rejected. In contrast to the wide inner confidence set for PCS increasing with education, here the inner confidence set only includes MCS from 70.10 to 70.34 for the bachelor's vs.

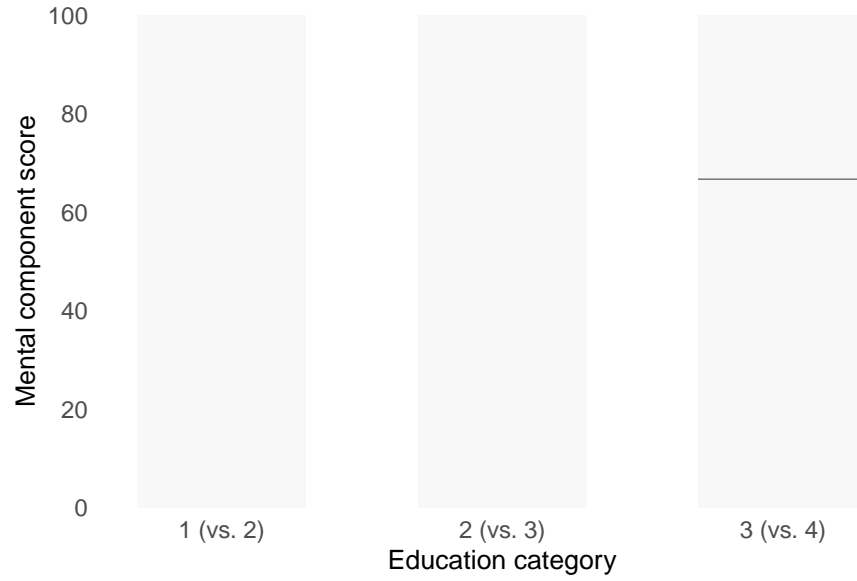


Figure 6: Inner (dark gray) and outer (light gray) 95% confidence sets.

graduate subpopulation comparison. Altogether, there is almost no strong evidence in either direction, leaving only uncertainty about the true population relationship. This is mostly due to the estimated MCS distributions being nearly the same across education categories.