

Multidimensional Interactive Fixed-Effects *

Hugo Freeman

October 4, 2022

[Latest version here](#)

Abstract

This paper studies a linear and additively separable model for multidimensional panel data of three or more dimensions with unobserved interactive fixed effects. Two approaches are considered to account for these unobserved interactive fixed-effects when estimating coefficients on the observed covariates. First, the model is embedded within the standard two-dimensional panel framework and restrictions are derived under which the factor structure methods in Bai (2009) lead to consistent estimation of model parameters. The second approach utilises popular machine learning techniques to develop group fixed-effects and kernel weighted fixed-effects that are more robust to the multidimensional nature of the problem. Theoretical results and simulations show the benefit of standard two-dimensional panel methods when the structure of the interactive fixed-effect term is known, but also highlight how the group fixed-effects and kernel methods perform well without knowledge of this structure. The methods are implemented to estimate the demand elasticity for beer under a handful of models for demand.

1 Introduction

Models of multidimensional data – panel data with more than two dimensions – are fast becoming popular in econometric analysis as large data sets with a multidimensional structure become available. For example, in gravity models of trade that are repeated over time one may be interested in studying trade patterns between an importer, i , an exporter, j , that is repeated

*First and foremost I would like to greatly thank my advisor, Martin Weidner, who has been exceptionally generous with his time and patience with me during my PhD and especially on this paper. I would also like to thank my secondary advisor Lars Nesheim for his guidance through my various PhD milestones and projects. For this paper specifically, additional to the above I would like to thank Ivan Fernández-Val and Andrei Zeleney for their helpful suggestions along with the audience at The International Panel Data Conference 2022, The IAAE Conference 2022, the Bristol Econometrics Study Group 2022, and the Warwick quantitative solutions & networking series 2022. This research was supported by the European Research Council grant ERC-2018-CoG-819086-PANEDA.

every quarter or year, t . One may also be interested in studying demand elasticities through consumption data that may vary by product, i , store, j , with repeated observation over week or month, t .¹ In these examples it is clear that there may exist unobserved characteristics in each dimension that can determine variation across both the dependent and independent variables that needs to be controlled for to avoid issues with endogeneity. For example, this could be shifts in taste preferences, that are unobserved by the econometrician, that may effect sales of particular products in certain stores differently over time. Thus far most analysis has addressed unobserved heterogeneity in the higher-dimensional setting by using a combination of additive scalar fixed-effects. These additive scalar fixed-effects approaches, however, can only accommodate variation in unobserved heterogeneity over a subset of dimensions with any one of the scalar fixed-effects terms. For example, in the three-dimensional model this type of fixed-effect approach can only control for variation over ij , it and jt , but not over all ijt . In the face of more complicated relationships that admit multiplicative variation across dimensions, these additive effects are unsatisfactory to control for unobserved heterogeneity. This paper develops tools to control for unobserved heterogeneity in the form of interactive fixed-effects in models of multidimensional panel data. The main body of the paper focuses on the linear and additively separable model. More generic applications of these tools are discussed in the introduction but are not formally studied.

To fix ideas consider linear parameter estimation in the following interactive fixed-effects model with three dimensions,

$$Y_{ijt} = X'_{ijt}\beta + \sum_{\ell=1}^L \varphi_{i\ell}^{(1)} \varphi_{j\ell}^{(2)} \varphi_{t\ell}^{(3)} + \varepsilon_{ijt}, \quad (1)$$

where all terms in $\sum_{\ell=1}^L \varphi_{i\ell}^{(1)} \varphi_{j\ell}^{(2)} \varphi_{t\ell}^{(3)}$ are unobserved and L is unknown. Reducing the problem to three dimensions is without loss of generality for the methods considered herein. Additive fixed-effects are omitted for brevity but are subsumed by the interactive fixed-effect term or can be removed with a simple within transformation. Let X_{ijt} be arbitrarily correlated with the unobserved interactive fixed-effects term, $\sum_{\ell=1}^L \varphi_{i\ell}^{(1)} \varphi_{j\ell}^{(2)} \varphi_{t\ell}^{(3)}$, but uncorrelated with the noise term, ε_{ijt} . The challenge to estimating β is isolating variation in X_{ijt} that is not correlated with the interactive fixed-effects term. This paper develops the multidimensional group fixed-effects and kernel weighted transformations to project out this unobserved heterogeneity and also shows settings where standard factor methods work well. The group fixed-effects method used below is similar to Bonhomme, Lamadon and Manresa (2021) and the within-cluster transformation in Freeman and Weidner (2022).

¹A non-exhaustive list of related examples can be found in the introduction of Matyas (2017) in trade, housing and prices, migration, country productivity and consumer price setting.

This paper makes two main contributions to the literature. The first is to show that the three or higher dimensional model can be couched in a standard two-dimensional panel data model and to derive sufficient conditions for consistency using these methods. The second contribution is to extend group fixed-effects methods to the multidimensional setting and introduce kernel weighted fixed-effects methods to this setting. The asymptotic results show that under certain conditions the group and kernel weighted fixed-effects can enjoy the parametric rate of convergence but that there can be some settings where the panel methods are still preferable. The panel methods do not in general achieve the parametric rate. The simulation results corroborate these theoretical findings and an empirical application that estimates the demand elasticity of beer demonstrates how these methods work in practice.

The within-cluster transformation can be motivated by considering a very simple extension to the usual within transformation to project additive fixed-effects of the form $a_{ij} + b_{it} + c_{jt}$. This is usually projected using

$$\dot{Y}_{ijt} = Y_{ijt} - \bar{Y}_{.jt} - \bar{Y}_{i.t} - \bar{Y}_{ij.} + \bar{Y}_{..t} + \bar{Y}_{.j.} + \bar{Y}_{i..} - \bar{Y}_{...}, \quad (2)$$

applied equivalently to X_{ijt} , where the variables with bars simply denote the average taken over the “dotted” index for the entire sample. That is, $\bar{Y}_{.jt} := \frac{1}{N_1} \sum_{i=1}^{N_1} Y_{ijt}$, $\bar{Y}_{..t} := \frac{1}{N_1 N_2} \sum_{i=1}^{N_1} \sum_{j=1}^{N_2} Y_{ijt}$, etc. The within-cluster transformation simply constrains the sample these averages are taken over to just within each unit’s cluster. With a slight abuse of notation, this is done using,

$$\tilde{Y}_{ijt} = Y_{ijt} - \bar{Y}_{i^*jt} - \bar{Y}_{ij^*t} - \bar{Y}_{ijt^*} + \bar{Y}_{i^*jt^*} + \bar{Y}_{i^*jt^*} + \bar{Y}_{ij^*t^*} - \bar{Y}_{i^*jt^*} \quad (3)$$

where the bar variables combined with the star indices denote means taken within that indice’s cluster. For example, \bar{Y}_{i^*jt} is the mean value of all i^* ’s assigned to i ’s cluster, \bar{Y}_{ij^*t} is the mean across both i^* in i ’s cluster and t^* in t ’s cluster, and so on. This is equivalent to including fixed-effects of the form $a_{ijg_3(t)} + b_{ig_2(j)t} + c_{g_1(i)jt}$, where $g_n(\cdot)$ maps to the cluster assignment for units in dimension n . The kernel weighted method simply uses weights rather than cluster assignments to take these averages. Hence, it is apparent that with a relatively small change in how the within transformation is performed, much more general fixed-effects can be controlled for, including the interactive fixed-effects considered in this paper.

The model for interactive fixed-effects has precedent in the standard two-dimensional panel data setting. For instance take the model considered in Bai (2009) and similar to Pesaran (2006),

$$Y_{it} = X'_{it}\beta + \sum_{\ell=1}^L \lambda_{i\ell} f_{t\ell} + e_{it}. \quad (4)$$

In that setting, Bai (2009) show that the interactive term $\sum_{\ell=1}^L \lambda_{i\ell} f_{t\ell}$ also sufficiently captures variation in additive individual and time effects without the need to specify these separately,

so these are again naturally omitted. For multidimensional applications it may be preferable to simply transform the problem in (1) to a two dimensional problem and estimate (4) directly using the transformed data. However, and as will be explained in further detail in Section 3.1, problems persist when L is large and only a subset of the unobserved heterogeneity parameters are low-dimensional. For consistent estimation of β , transforming the multidimensional array to a matrix then estimating (4) requires either: (a) all fixed-effects are low-dimensional, or; (b) that a subset of the fixed-effects are low-dimensional and the analyst knows which ones are. The requirement that the analyst has this knowledge can be highly restrictive. Furthermore, only a very slow rate of convergence can be shown for this approach. Alternatively, the within-cluster and kernel weighted transformations analysed in this paper requires only that a subset of the fixed-effect parameters are low-dimensional, though the analyst does not need to know which of the fixed-effect parameters make up this subset.

The demand elasticity for beer application uses Dominick’s supermarket data from the Chicago area from 1991-1995, where price and quantity vary over product, store, and month. The log-log and logit models are estimated. The log-log model is implemented with and without cross-elasticities to demonstrate the limited ability of the fixed-effects estimators to project relevant control variables that vary across all dimensions. The elasticities for the fixed-effects estimators, including the simple additive fixed-effects, are all similar within each model, indicating that whilst fixed-effects probably exist in this setting they are most likely not overly complicated. The estimates with and without cross-elasticities are also substantially different, which indicates that if cross-elasticities are important in this setting then even the more sophisticated fixed-effects estimators cannot project them out. This indicates that relevant control variables with high variation across all dimensions should still be included in the regression line. The estimates from the log-log with cross-elasticities closely reflect the own-price elasticities in Table 1 from Hausman, Leonard and Zona (1994).

Under sufficient regularity conditions, the methods considered in this paper may also control for variation from arbitrary functions of the fixed-effects. Similar to that considered in Zeleneev (2020) and Freeman and Weidner (2022), the functional representation of model (1) could be,

$$Y_{ijt} = X'_{ijt}\beta + h(\varphi_i^{(1)}, \varphi_j^{(2)}, \varphi_t^{(3)}) + \varepsilon_{ijt},$$

for vector-valued $\varphi_i^{(1)}$, $\varphi_j^{(2)}$ and $\varphi_t^{(3)}$. The set of fixed-effects to be transformed by the function $h(\cdot, \cdot, \cdot)$ could also extend to fixed-effects over multiple indices, e.g. α_{ij} from above. It should be noted that the setting considered in Zeleneev (2020) requires that the transformation is non-smooth, and it is not trivial to see that a “within-type” transformation will sufficiently project this type of heterogeneity. With sufficient smoothness conditions on the function transforming the fixed-effects, existing literature could be generalised to show consistency using the proposed

within-cluster transformation in the multidimensional case.

Models with discrete explanatory variables (Chernozhukov, Fernández-Val, Hahn and Newey, 2013; Hoderlein and White, 2012; Evdokimov, 2010; Fernández-Val, Freeman and Weidner, 2021), provide another interesting application of these group fixed-effects estimators. Take the following regression line for discrete valued X_{ijt} ,

$$Y_{ijt} = h\left(X_{ijt}, \varphi_i^{(1)}, \varphi_j^{(2)}, \varphi_t^{(3)}, \varepsilon_{ijt}\right).$$

Then, under sufficient smoothness conditions on the function h , the unobserved heterogeneity may also be projected out with a group fixed-effect estimator. Tensor completion techniques also have useful generalisations in this setting, for example, Tomioka, Hayashi and Kashima (2010); Li, Wang, Lu and Tang (2019); Xu (2020), for some examples of methods that consider sparse multidimensional arrays. The sparse multidimensional array problem has similar complexities to the low-rank tensor approximation problem in that they do not extend from the matrix problems in a straightforward way, hence require non-trivial extensions.

It is also important to consider unobserved heterogeneity in applications that admit discrete dependent variables. For example, for binary response variable with known $F(\cdot)$,

$$Y_{ijt} = F\left(X'_{ijt}\beta + \sum_{\ell=1}^L \varphi_{i\ell}^{(1)} \varphi_{j\ell}^{(2)} \varphi_{t\ell}^{(3)}\right). \quad (5)$$

Estimation of the unobserved heterogeneity term may then be performed with a similar iterative scheme as that proposed in Chen, Fernández-Val and Weidner (2021) or the sufficient statistic approach in Chapter 6 of Matyas (2017). The incidental parameter problem in this setting can be alleviated using methods in this paper by allowing cluster sizes to grow with data size coupled with taking grouped fixed-effects along fewer dimensions, for example in Bonhomme, Lamadon and Manresa (2021) and also Appendix C.

Menzel (2021) consider a special case of multidimensional data for bootstrapping methods where the data is D -adic. That is, each dimension of the data refers to the same set of observations, like a network graph where each index refers to an individual in the network. An example of the multidimensional version of this could be a binary indicator of a three step path, $Y_{ijk} = G_{ij}G_{jk}$, detailing if there exists a path from i to k . In any case, the type of multidimensional data considered in that work is a distinct special case of the type of data structures considered in this paper.

The technical component of this paper is highly related to the numerical analysis literature on low-rank approximations of multidimensional arrays. As pointed out in De Silva and Lim (2008), the optimisation problem of finding low-rank approximations in the tensor setting is not well-posed, hence most results in this literature rely on numerical evidence. See Kolda

and Bader (2009) for a summary of the multidimensional array decomposition problem and Vannieuwenhoven, Vandebril and Meerbergen (2012); Rabanser, Shchur and Günnemann (2017) for examples of numerical results. As such, it is necessary to innovate on this tensor low-rank problem to find appropriate analytical results. To this end, this paper utilises well-posed components of the numerical analysis literature for use in nuisance parameter applications. These applications have the advantage that they do not require the multidimensional array of fixed-effects to be reconstructed, hence do not attempt to directly solve the low-rank tensor problem. It is worth a note that Elden and Savas (2011), along with related papers, suggest a reformulation of the low multilinear rank problem that may have promising applications in econometrics, but this is left for future research.

The paper is organised as follows. Section 2 introduces the model, and notation and preliminaries; Section 3 details the estimators and associated assumptions with convergence results; Section 4 discusses the convergence results along with some alternative assumptions, and further motivates the estimation approach; Section 5 displays the simulation results; Section 6 shows the beer demand estimation empirical application; and Section 7 concludes.

2 Model

Let β^0 denote the true parameter value for the slope coefficients. The model in full dimensional generality is,²

$$\mathbf{Y} = \sum_{k=1}^K \mathbf{X}_k \beta_k^0 + \mathcal{A} + \varepsilon, \quad (6)$$

where $\mathbf{Y}, \mathbf{X}_k, \varepsilon \in \mathbb{R}^{N_1 \times N_2 \times \dots \times N_d}$. $\mathcal{A} = \sum_{\ell=1}^L \varphi_{\ell}^{(1)} \circ \dots \circ \varphi_{\ell}^{(d)}$ where $\varphi_{\ell}^{(n)} \in \mathbb{R}^{N_n}$ for each $n = 1, \dots, d$ and “ \circ ” is the outer product. L is naturally restricted to have upper bound $\min_n \{\prod_{n' \neq n} N_{n'}\}$, see Kruskal (1989). ε is a noise term uncorrelated with all \mathbf{X}_k and all unobserved fixed-effects terms. Take $i_n \in \{1, \dots, N_n\}$ for all $n \in \{1, \dots, d\}$ as the dimension specific index, where N_n is the sample size of dimension n . The regressors \mathbf{X}_k may be arbitrarily correlated with \mathcal{A} . Throughout this paper all dimensions are considered to grow asymptotically, that is $N_n \rightarrow \infty$ for all n .

Model (6) can be seen as a natural extension of the Bai (2009) model to three (or more) dimensions with the \mathcal{A} term interpreted as a “higher-dimensional” factor structure. Similar to

²For example, in index notation this model can be written as,

$$Y_{i_1, i_2, \dots, i_d} = \sum_{k=1}^K X_{i_1, i_2, \dots, i_d; k} \beta_k^0 + \mathcal{A}_{i_1, i_2, \dots, i_d} + \varepsilon_{i_1, i_2, \dots, i_d}$$

with $\mathcal{A}_{i_1, i_2, \dots, i_d} = \sum_{\ell=1}^L \varphi_{i_1 \ell}^{(1)} \dots \varphi_{i_d \ell}^{(d)}$.

this strain of the literature, all terms in \mathcal{A} are considered fixed nuisance parameters. There are potentially many extensions to the factor model setting in Bai (2009) to the higher dimension case. This paper starts with what seems the most natural extension.

The term \mathcal{A} may also incorporate additive fixed effects that vary in any strict subset of the dimensions. For example, in the three dimensional setting one may want to control for the additive terms, $a_{ij} + b_{it} + c_{jt}$. These can be controlled for using $L = \min\{N_1, N_2\} + \min\{N_1, N_3\} + \min\{N_2, N_3\}$, with the first $\min\{N_1, N_2\}$ terms $\sum_{\ell=1}^{\min\{N_1, N_2\}} \varphi_{i\ell}^{(1)} \varphi_{j\ell}^{(2)} = a_{ij}$ by setting $\varphi_{t\ell}^{(3)} = 1$ for $\ell = 1, \dots, \min\{N_1, N_2\}$, and so on for the b_{it} and c_{jt} . These could also be controlled for directly using the standard within-transformation before considering the model in (6).

This paper comprises of two main modelling approaches. The first is to embed the multidimensional model into a standard panel data model by simply flattening all arrays into matrices. The second approach uses weighted differences across each dimension to reduce each $\varphi^{(n)}$ component of \mathcal{A} separately for each n . For this reason the model assumptions are split out and stated in Section 3 alongside each estimation approach.

2.1 Notation and preliminaries

For a d -order tensor, \mathbf{A} , a factor- n flattening, denoted as $\mathbf{A}_{(n)}$, is the rearrangement of the tensor into a matrix with dimension n varying along the rows and the remaining dimensions simultaneously varying over the columns. That is, $\mathbf{A}_{(n)} \in \mathbb{R}^{N_n \times N_{n+1} N_{n+2} \dots N_1 \dots N_{n-1}}$. The Frobenius norm, $\|\cdot\|_F$, of a matrix or tensor is the entry-wise norm, $\|\mathbf{A}\|_F^2 = \sum_{i_1=1}^{N_1} \dots \sum_{i_d=1}^{N_d} A_{i_1 \dots i_d}^2$. The spectral norm, denoted $\|\cdot\|$, is the largest singular value of a matrix. For a d -order tensor, \mathbf{A} , the multilinear rank, denoted \mathbf{r} , is a vector of matrix ranks after factor- n flattening in each dimension, with each component of this vector $r_n = \text{rank}(\mathbf{A}_{(n)})$. Tensor rank, different to multilinear rank, is defined as the least number of outer products of vectors to replicate the tensor. That is, for tensor \mathbf{A} and vectors $u_\ell^{(n)} \in \mathbb{R}^{N_n}$, tensor rank is the smallest L such that $\mathbf{A} = \sum_{\ell=1}^L u_\ell^{(1)} \circ \dots \circ u_\ell^{(d)}$, where \circ is the outer product of a vector.

The n -mode product between a tensor \mathbf{A} and matrix B is denoted $\mathbf{A} \times_n B$ and has elements

$$(\mathbf{A} \times_n B)_{i_1, \dots, j, \dots, i_d} = \sum_{i_n=1}^{N_n} A_{i_1, \dots, j, \dots, i_n} B_{j, i_n},$$

which is equivalent to saying the flattening $(\mathbf{A} \times_n B)_{(n)} = B \mathbf{A}_{(n)}$. This can be referred to as “hitting” the tensor \mathbf{A} with matrix B in the n^{th} dimension, though this terminology is only stated to help with understanding of the definition.

The singular value decomposition is used in both estimation approaches in this paper so it is important to understand some of its properties. The singular value decomposition of a matrix,

$A \in \mathbb{R}^{N_1 \times N_2}$ is

$$A = U\Sigma V' = \sum_{r=1}^{\min\{N_1, N_2\}} \sigma_r u_r v_r' \quad (7)$$

where U is the matrix of left singular vectors, u_r , V is the matrix of right singular vectors, v_r , and Σ is a diagonal matrix of singular values, σ_r , with values running in descending order down the diagonal. For a rank- r matrix, the first r entries on the diagonal of Σ are strictly positive and the remaining entries are zero.

Take the approximation problem,

$$\min_{A'} \|A - A'\|_F \text{ such that } \text{rank}(A') = k. \quad (8)$$

It is well known from the Eckart-Young-Mirsky theorem that the solution to this approximation problem is the first k terms of the singular value decomposition, i.e. $\sum_{r=1}^k \sigma_r u_r v_r'$. The Eckart-Young-Mirsky theorem effectively picks out the row and column subspaces that best explain variation in the matrix A as the leading columns of the matrix U , respectively of V . The sum of squared error at the minimiser is thus $\sum_{r=k+1}^{\min\{N_1, N_2\}} \sigma_r^2$. This is commonly called a low-rank approximation and forms the cornerstone for estimation of unobserved heterogeneity in the factor model and interactive fixed-effects models in Bai and Ng (2002); Bai (2009); Moon and Weidner (2015) amongst others.

The Eckart-Young-Mirsky theorem, however, does not extend to the three or higher dimensional setting, see De Silva and Lim (2008) for details. This is why the multidimensional problem either needs to be translated to the two-dimensional setting to utilise the Eckart-Young-Mirsky theorem, or the fixed-effects parameters need to be shrunk separately, as is done with the group fixed-effects and kernel methods.

3 Estimation

This section details the three estimation approaches used. The first subsection details how to apply standard two dimensional estimators to the problem and the assumptions required for consistent estimation. The second and third subsections detail the group fixed-effect and kernel weighted approaches and the required assumptions for consistency in that setting.

3.1 Matrix low-rank approximation estimator

This section provides a description of some matrix methods that can be applied directly to the multidimensional model and stipulates the assumptions required for consistency. Note that

Kapetanios, Serlenga and Shin (2021) employ a similar approach for three-dimensional arrays in conjunction with the Pesaran (2006) common correlated effects estimator.

Consider recasting the multidimensional array problem into a two dimensional panel problem by flattening Y and X in the n -th dimension,

$$Y_{(n)} = X'_{(n)}\beta^0 + \varphi^{(n)}\Gamma'_n + \varepsilon_{(n)}$$

where $Y_{(n)}, X_{(n)}, \varepsilon_{(n)} \in \mathbb{R}^{N_n \times \prod_{n' \neq n} N_{n'}}$, $\varphi^{(n)}$ is an $N_n \times r_n$ matrix and Γ_n is an $\prod_{n' \neq n} N_{n'} \times r_n$ matrix that accounts for variation in the remaining $\varphi^{(n')}$ for all $n' \neq n$. The term r_n is indexed by the dimension n because it may vary non-trivially according to the flattened dimension. It should then be apparent that this is exactly the model described in (4), that is, the standard linear model with factor structure unobserved heterogeneity as studied in Bai (2009).

The two-dimensional estimator for a given flattening, n , optimises the following objective function,

$$R(\beta, \hat{r}_n, n) = \min_{\substack{\varphi^{(n)} \in \mathbb{R}^{N_n \times \hat{r}_n}, \\ \Gamma_n \in \mathbb{R}^{\prod_{n' \neq n} N_{n'} \times \hat{r}_n}}} \left\| Y_{(n)} - X'_{(n)}\beta - \varphi^{(n)}\Gamma'_n \right\|_F^2. \quad (9)$$

Then $\hat{\beta}_{(n^*)}^{2D} = \operatorname{argmin}_{\beta} R(\beta, \hat{r}_n, n)$ is the slope estimate for the two-dimensional setup. The analyst must choose both the dimension to flatten in, n , and the rank of the estimated interactive fixed-effects term, \hat{r}_n . It is well known that the minimum in (9) is achieved using the leading \hat{r}_n terms from the singular value decomposition of the error term, $Y_{(n)} - X'_{(n)}\beta$. This gives $\hat{\varphi}^{(n)}$ as the first \hat{r}_n columns of $\hat{U}\hat{\Sigma}$ and $\hat{\Gamma}_n$ as the first \hat{r}_n columns of \hat{V} where \hat{U} , $\hat{\Sigma}$ and \hat{V} are the terms from (7) of the singular value decomposition of $Y_{(n)} - X'_{(n)}\beta$. Because this error term is a function of β , an iteration is naturally required between estimating β and finding the singular value decomposition of the error term. This is a well studied iteration procedure, for convergence details see Bai (2009); Moon and Weidner (2015).

In the following assumptions let \hat{r}_n be the estimated number of factors for the (n) -flattening of the regression line when applying the least square methods in (9). Also, let $\mathcal{L} \subset \{1, \dots, d\}$ be a non-empty subset of the dimensions. The tensor rank parameter, L , may without loss be restricted to the upper bounded by $L \leq \min_n \{\prod_{n' \neq n} N_{n'}\}$. This is a result of elementary bounds on the tensor rank of an arbitrary tensor. In the following, the multilinear rank of \mathcal{A} is restricted such that it is low-rank along at least one of the flattenings.

Assumption 1 (Bounded norms of covariates and exogenous error).

$$(i). \quad \|X_k\|_F = O_p \left(\prod_{n=1}^d \sqrt{N_n} \right) \text{ for each } k$$

$$(ii). \quad \|\varepsilon_{(n^*)}\| = O_p \left(\max\{\sqrt{N_{n^*}}, \prod_{m \neq n^*} \sqrt{N_m}\} \right) \text{ for each } n^* \in \mathcal{L}$$

Assumption 2 (Weak exogeneity). $\text{vec}(X_k)' \text{vec}(\varepsilon) = O_p\left(\prod_{n=1}^d \sqrt{N_n}\right)$ for each k

Assumption 3 (Low multilinear rank). For some positive integer, c , $r_{n^*} < c$ for all $n^* \in \mathcal{L}$, where r_n is the n^{th} component of the multilinear rank of \mathcal{A} .

Assumption 4 (Non-singularity). Let $\sigma_s(A)$ be the s^{th} singular value for a matrix A . For each dimension $n^* \in \mathcal{L}$ that satisfies Assumption 3, there exists a $K \times 1$ unit vector δ_{n^*} such that

$$\sum_{s=2\bar{L}(n^*)+1}^{\min\{N_{n^*}, \prod_{m \neq n^*} N_m\}} \sigma_s \left(\frac{(\delta_{n^*} \cdot X_{(n^*)})(\delta_{n^*} \cdot X_{(n^*)})'}{\prod_n \sqrt{N_n}} \right) > b > 0 \quad \text{wpa1.}$$

Assumptions 1, 2 and 4 are standard regularity assumptions already well established in the literature, e.g. see Moon and Weidner (2015). Assumption 1.(i) ensures that the covariates have bounded norms, for example having bounded second moments. Assumption 1.(ii) allows for some weak correlation across dimensions, see Moon and Weidner (2015), or is otherwise implied if the noise terms are independently distributed with bounded fourth moments, see Latała (2005). Assumption 2 is implied if $X_{i_1, i_2, \dots, i_d; k \in i_1, i_2, \dots, i_d}$ are zero mean, bounded second moment and only admits weak correlation across dimensions for each $k = 1, \dots, K$. Assumption 4 simply states that, after factor projection, the set of covariates still collectively admit full-rank variation.

Assumption 3 is new and asserts that there exists at least one flattening of the interactive term, \mathcal{A} , that is low-dimensional or simply low-rank. Given that the true value for L is left mostly unrestricted at this stage, this requires that at least one of the unobserved terms $\varphi^{(n)}$ is low dimensional. Note that not all dimensions must satisfy Assumption 3 for the below result. If the correct dimension is chosen then variation from the interactive term can be sufficiently projected out using the factor model approach. This makes up the statement of the following Proposition.

Proposition 1. Let $\hat{\beta}_{(n^*)}^{2D}$ be the estimator from Bai (2009) after first flattening along dimension $n^* \in \mathcal{L}$. If Assumptions 1-4 hold, the subset \mathcal{L} is non-empty, and the estimated number of factors $\hat{r}_{n^*} \geq r_{n^*}$, then, for each $n^* \in \mathcal{L}$ satisfying Assumption 3,

$$\left\| \hat{\beta}_{(n^*)}^{2D} - \beta^0 \right\| = O_p \left(\frac{1}{\sqrt{\min\{N_{n^*}, \prod_{n \neq n^*} N_n\}}} \right). \quad (10)$$

Proposition 1 follows directly from Moon and Weidner (2015) since the flattening procedure reduces the problem to the standard linear factor model. Notice that this result only applies to estimates in the dimension(s) that satisfy the low-rank assumption in Assumption 3. That is, implicit in Proposition 1 is that the analyst has chosen the correct dimension to flatten over

when reformulating the problem as a two-dimensional panel. Assumption 3 can be relaxed to $r_{n^*} = o\left(\min\{N_{n^*}, \prod_{n \neq n^*} N_n\}\right)$ as long as the estimated number of factors is allowed to increase with data size at a faster rate than this. The constraint $\hat{r}_{n^*} \geq r_{n^*}$ can also be changed to $\hat{r}_{n^*} \geq c$, however, this is more conservative than required for the statement of the result.

The estimation procedure from Proposition 1 can also be augmented to flatten over multiple indices. For instance, the analyst may flatten such that both the rows and columns in the matrix contain multiple indices from the original array. Of course, this augmentation makes Assumption 3 harder to satisfy as it requires multiple parameters to vary in low-dimensional space. To see this take the tensor \mathcal{A} flattened over the first two indices as $\mathcal{A}_{(1,2)} \in \mathbb{R}^{N_1 N_2 \times \prod_{n \notin \{1,2\}} N_n}$. If the parameters $\varphi^{(n)}$ for $n = 3, \dots, d$ are high-dimensional, Assumption 3 is only satisfied when both $\varphi^{(1)}$ and $\varphi^{(2)}$ and their product space is low-dimensional. Clearly this is more restrictive than requiring only one of the parameter spaces to be low-dimensional. However, flattening along multiple dimensions can improve the convergence rate in Proposition 1 to $O_p\left(\frac{1}{\sqrt{\min\{N_1 N_2, \prod_{n \notin \{1,2\}} N_n\}}}\right)$, so there are benefits if this more restrictive assumption can be made. Further discussion of the matrix method results are relegated to Section 4.1, in particular some avenues to choosing the dimension to flatten over.

3.2 Group fixed-effects

This section describes the group fixed-effects estimator. Take again the model in array notation,

$$\mathbf{Y} = \sum_{k=1}^K \mathbf{X}_k \beta_k^0 + \mathcal{A} + \boldsymbol{\epsilon}.$$

A cluster assignment, \mathcal{C} , is a length d list of partition matrices, $C_n \in \{0, 1\}^{N_n \times G_n}$, where G_n is the number of clusters in dimension n and each entry of C_n is a binary indicator of a unit's membership to a given cluster. Clusters are assigned separately along each dimension. Let $\Theta_{\mathcal{C}}$ be the space of group fixed-effects parameters associated to cluster assignment \mathcal{C} . Each $\boldsymbol{\theta} \in \Theta_{\mathcal{C}}$ is an ordered set of size d of $\times_{n=1}^d N_n$ tensors. For each n in $\{1, \dots, d\}$, the tensor θ_n varies freely over dimensions $\{1, \dots, n-1\}$ and $\{n+1, \dots, d\}$ but is constant within each cluster along dimension n .³ The objective function for the group fixed-effect estimation of β under cluster assignment \mathcal{C} is

$$Q(\beta, \mathcal{C}) = \min_{\boldsymbol{\theta} \in \Theta_{\mathcal{C}}} \left\| \mathbf{Y} - \sum_{k=1}^K \mathbf{X}_k \beta_k - \sum_{n=1}^d \theta_n \right\|_F^2 \quad (11)$$

and $\hat{\beta}_{GFE, \mathcal{C}} := \operatorname{argmin}_{\beta \in \mathbb{R}^K} Q(\beta, \mathcal{C})$.

³This parameter space is exemplified in Remark 1 for the three dimensional setting for clarity.

The minimum within (11) is obtained from the within-cluster transformation in (3) from the Introduction. Remark 1 below details this objective function for the three-dimensional setting for clarity. It should be clear that the parameter space $\Theta_{\mathcal{C}}$ is indexed by cluster assignment \mathcal{C} because this assignment defines how the parameters may vary. That is, this is the estimated parameter space under a specific group fixed-effects estimator, which may only be an approximation of the true parameter space. Cluster assignments may be known or estimated, with suggestions of how to estimate these discussed in Section 4.2. In the case these are estimated from the error term, $\mathbf{Y} - \sum_{k=1}^K \mathbf{X}_k \beta_k$, there is an iteration between estimating cluster assignments and estimating slope coefficients, like in Section 3.1. This iteration is discussed in Section 4.4.

Assumption 5 (Clustering).

Let $j_n(i_n)$ be any unit in the same cluster as i_n from using cluster assignment \mathcal{C} . Then,

- (i). For all n as $N_n \rightarrow \infty$, $\frac{1}{N_n} \sum_{i_n=1}^{N_n} \left\| \varphi_{i_n}^{(n)} \right\|^2 \lesssim O_p(1)$ and,
- (ii). For a non-empty subset $\mathcal{M} \subset \{1, \dots, d\}$ take for any $n^* \in \mathcal{M}$ a sequence $\xi_{N_{n^*}} \rightarrow 0$ as $N_{n^*} \rightarrow \infty$. Then,

$$\frac{1}{N_{n^*}} \sum_{i_n=1}^{N_{n^*}} \left\| \varphi_{i_{n^*}}^{(n^*)} - \varphi_{j_{n^*}(i_{n^*})}^{(n^*)} \right\|^2 = O_p(\xi_{N_{n^*}})$$

Assumption 5.(i) restricts fixed-effects parameter space to have finite second moments. This implies that as $\{N_1, \dots, N_d\} \rightarrow \infty$, cluster allocations cannot become increasingly disparate in the underlying parameter space. Assumption 5.(ii) states that for at least one dimension the clustering procedure finds matches with asymptotically negligible difference in the underlying parameter space. Since cluster assignments are not always estimated, and actually sometimes group assignments may be given extraneously, it is useful to state Assumption 5 in generic terms that ignore these clustering mechanics. These assumptions restrict the cluster assignment to uncover closeness in the true parameter space, which implies a restriction on the underlying parameter space and on how clusters are assigned.

Below is a refinement to the regularity conditions contained within the Assumptions listed in Section 3.1 that account for the within-cluster transformation.

Assumption 6 (Regularity conditions). Let $\tilde{T}_{i_1, \dots, i_d}$ be the entries of tensor \mathbf{T} after the group fixed-effects from the minimiser of (11) are differenced out. Then,

- (i). $\left(\frac{1}{\prod_n N_n} \sum_{i_1} \cdots \sum_{i_d} \tilde{X}_{i_1, \dots, i_d} \tilde{X}'_{i_1, \dots, i_d} \right) = O_p(1)$ converges to a nonrandom positive definite matrix as $N_1, \dots, N_d \rightarrow \infty$.
- (ii). $\frac{1}{\prod_n N_n} \sum_{i_1} \cdots \sum_{i_d} \tilde{X}_{i_1, \dots, i_d} \varepsilon_{i_1, \dots, i_d} = O_p\left(\frac{1}{\sqrt{\prod_n N_n}}\right)$.

Assumption 6.(i) is very similar to Assumption 4 except that here full rank is required after the within-cluster projection rather than the factor projection. Assumption 6.(ii) is an exogeneity condition that requires weak exogeneity in the covariates after the within-cluster transformation, which can be viewed as similar to Assumption 2. This is stricter than Assumption 2 because the noise term ε can foreseeably impact cluster allocation if clusters are estimated as functionals of a residual term. This limitation is alleviated by, for instance, making sure cluster assignments are based on variables extraneous to the regression line, hence independent of ε , or perhaps through some sample splitting methods such as that proposed in Freeman and Weidner (2022).

Proposition 2 (Upper bound on group fixed-effects estimator). *Let Assumptions 5 and 6 hold for cluster allocation \mathcal{C} . Let \mathcal{M} be the set defined in Assumption 5.(ii). Then, for tensor rank L_N that may depend on sample size,*

$$\|\widehat{\beta}_{GFE,\mathcal{C}} - \beta^0\| = \sqrt{L_N} O_p \left(\prod_{n^* \in \mathcal{M}} \sqrt{\xi_{N_{n^*}}} \right) + O_p \left(\prod_{n=1}^d \frac{1}{\sqrt{N_n}} \right).$$

Discussed in Section 4.2 are methods and restrictions that restrict $\xi_{N_{n^*}}$ from Assumption 5 and Proposition 2 to $1/N_{n^*}$. This suggests that as long as $\mathcal{M} = \{1, \dots, d\}$ and L_N is bounded the parametric rate of convergence is achievable. Related to \mathcal{M} , an implicit requirement on the latent parameters from this subset of dimensions is some form of low-dimensionality in the vectors $\varphi_{i_{n^*}}^{(n^*)}$ for $n^* \in \mathcal{M}$. For a discussion on the curse of dimensionality using clustering methods, see Bonhomme, Lamadon and Manresa (2021) that suggests the dimension of these parameters should be ≤ 2 to be well-clustered. A sufficient condition for parameters in these dimension to be low-dimensional is low multilinear rank for each $n^* \in \mathcal{M}$. Hence, it is expected that the set \mathcal{M} should be a subset of \mathcal{L} , from Assumption 3. The advantage with the group fixed-effects methods is that the analyst does not need to choose which n does admit low multilinear rank, hence it is more flexible. However, since the matrix factor methods do not suffer such a large curse of dimensionality, if the low multilinear rank dimension is known then there is still some advantage in using this method, for example if the a certain multilinear rank parameter is bounded but of order 5-10.

Remark 1. *In the three dimensional setting, (3) achieves the minimum in the group fixed-effect objective function,*

$$Q(\beta, \mathbf{g}) = \min_{\alpha, \gamma, \delta} \sum_{i,j,t} (Y_{ijt} - X'_{ijt}\beta - \theta_{1,g_1(i)jt} - \theta_{2;ig_2(j)t} - \theta_{3;ijg_3(t)})^2 \quad (12)$$

where $\theta_1 \in \mathbb{R}^{\text{ncol}(g_1) \times N_2 \times N_3}$, $\theta_2 \in \mathbb{R}^{N_1 \times \text{ncol}(g_2) \times N_3}$, and $\theta_3 \in \mathbb{R}^{N_1 \times N_2 \times \text{ncol}(g_3)}$; and $\text{ncol}(\cdot)$ returns the number of columns of a matrix. \mathbf{g} is a list of group assignment for each dimension. For

example, $g_1(i)$ maps to the group identity of individual i . This is why α is restricted to vary across only $\text{ncol}(g_1)$ different values in the first dimension, which is less than N_1 .

Notice that the optimisers for θ_1 , θ_2 and θ_3 from (12) can be described as combinations of the within-cluster projection as follows,

$$\begin{aligned}\hat{\theta}_{1;g_1(i)jt} &\approx \bar{\mathcal{A}}_{i^*jt} - \bar{\mathcal{A}}_{i^*j^*t} + \bar{\mathcal{A}}_{i^*j^*t^*} \\ \hat{\theta}_{2;ig_2(j)t} &\approx \bar{\mathcal{A}}_{ij^*t} - \bar{\mathcal{A}}_{i^*j^*t} \\ \hat{\theta}_{2;ijg_3(t)} &\approx \bar{\mathcal{A}}_{ijt^*} - \bar{\mathcal{A}}_{i^*jt^*},\end{aligned}$$

though this representation is not unique.

Additional to controlling for any additive terms, this projection leaves the following interactive fixed-effects residual,

$$\tilde{\mathcal{A}}_{ijt} = \sum_{\ell=1}^L (\varphi_{i\ell}^{(1)} - \bar{\varphi}_{i^*\ell}^{(1)}) (\varphi_{j\ell}^{(2)} - \bar{\varphi}_{j^*\ell}^{(2)}) (\varphi_{t\ell}^{(3)} - \bar{\varphi}_{t^*\ell}^{(3)}). \quad (13)$$

where $\bar{\varphi}_{i^*\ell}^{(1)}$ is the group mean of $\varphi_{i\ell}^{(1)}$ for the i^* 's in i 's group, and so on for the other terms. Hence, sufficient projection of the interactive fixed-effects terms relies on the weaker condition that parameters converge to their group means, namely, $\varphi_{i\ell}^{(1)} \rightarrow \bar{\varphi}_{i^*\ell}^{(1)}$, $\varphi_{j\ell}^{(2)} \rightarrow \bar{\varphi}_{j^*\ell}^{(2)}$ or $\varphi_{t\ell}^{(3)} \rightarrow \bar{\varphi}_{t^*\ell}^{(3)}$ for each ℓ . Indeed, the group mean differencing could be seen as a weighted mean difference across the population, with equal weight given to observations within the cluster and zero weight to observations outside of each cluster. This fact is utilised for the more generic kernel weighted difference estimator in Section 3.3, which is synonymous to a Nadaraya-Watson type estimator for each fixed-effects term.

So far it has been shown that with the relatively innocuous shift from the within transformation to the within-cluster transformation, any additive terms are automatically controlled for and there are conditions to also control for the interactive term. Choice of clusters for this transformation is key to suffice this less restrictive condition. Given a set of proxies to cluster on, clustering or matching methods can be used to find these groups, for example Bonhomme, Lamadon and Manresa (2021). Developing a set of proxies to cluster on is important and is discussed in Section 4.2.

3.3 Kernel weighted fixed-effects

Let $\hat{\varphi}_{i_n}^{(n)}$ generically denote a proxy measure for unit i_n in dimension n that may be known or estimated. The use of this notation will become clear in the statement of Proposition 3 and in discussion of how to estimate these proxy measures in Section 4.2. Let \mathcal{W} be a list of weight

matrices, where the n^{th} item $W_{(n)} \in \mathbb{R}^{N_n \times N_n}$ has elements,

$$w_{i_n, j_n}^{(n)} := \frac{k\left(\frac{1}{h_n} \left\| \hat{\varphi}_{i_n}^{(n)} - \hat{\varphi}_{j_n}^{(n)} \right\| \right)}{\sum_{i'_n=1}^{N_n} k\left(\frac{1}{h_n} \left\| \hat{\varphi}_{i_n}^{(n)} - \hat{\varphi}_{i'_n}^{(n)} \right\| \right)}, \quad (14)$$

where k is a kernel function, and h_n is a bandwidth parameter. Let Δ be a d -list of $\times_{n=1}^d N_n$ tensors $\delta_n \in \mathbb{R}^{N_1 \times \dots \times N_d}$. For a given set of proxy measures and kernel function, the kernel weighted fixed-effects estimator optimises

$$S(\beta, \mathcal{W}) = \min_{\delta \in \Delta} \left\| \mathbf{Y} - \sum_{k=1}^K \mathbf{X}_k \beta_k - \sum_{n=1}^d \delta_n \times_n W_{(n)} \right\|_F^2. \quad (15)$$

Then, $\hat{\beta}_{KER, \mathcal{W}} := \operatorname{argmin}_{\beta \in \mathbb{R}^K} S(\beta, \mathcal{W})$. The notation \times_n is the n -mode product defined at the beginning of this section. Again, if proxy measures are estimated from the error term $\mathbf{Y} - \sum_{k=1}^K \mathbf{X}_k \beta_k$, then there is an iteration between slope estimation and estimation of kernel weights, much like in the other estimators presented already.

Assumption 7 (Kernels). *Denote the kernel function used as $k(\cdot)$ and let this function be finite. Then for $a \geq 0$ and $h > 0$ there exists an $\alpha > 0$ such that $k(a/h)a \lesssim O(h^\alpha)$.*

Assumption 7 refers to a bandwidth parameter, h , and restricts the kernels to penalise distance at a rate equal to or faster than $O(h^\alpha/a)$. For consistency using the kernel methods, the sequence $h \rightarrow 0$ is considered, such that an upper bound on α is the critical object of interest.

As an example of a class of kernel functions that satisfies Assumption 7, the exponential class of the form considered in Remark 2 may be utilised.

Remark 2. *For $c_1, c_2 > 0$, let $k'(a) \propto c_1 \exp(-c_2 a^2)$ for all $a \geq 0$ and $k' \in \mathcal{K}'$. Then $\operatorname{argmax}_a k'(a/h)a = h/\sqrt{2c_2}$, and,*

$$\max_a k'(a/h)a \propto \frac{c_1}{\sqrt{2c_2}} e^{-1/2} h = O(h)$$

Thus, Assumption 7 is satisfied for the exponential class of kernel functions \mathcal{K}' with $\alpha = 1$. Further, for $h \rightarrow 0$, it suffices that $\alpha \in (0, 1]$.

Assumption 7 is stated more generically than Remark 2 as there is a larger class of bounded kernel functions that satisfy the sufficient restriction for the result below. The point here is to show that Assumption 7 is satisfied for some very standard kernel functions, hence is not too restrictive.

Assumption 8 (Regularity of proxy measures). *Let $\hat{\varphi}_{i_n}^{(n)} \in \hat{\Phi}_n$ be the proxy space for the fixed-effects and let $k(\cdot)$ be a finite kernel function. Let $K_{i_n}(h_n) := \max_{j_n} k\left(\frac{1}{h_n} \left\| \hat{\varphi}_{i_n}^{(n)} - \hat{\varphi}_{j_n}^{(n)} \right\| \right)$. For $0 < e_{i_n} < K_{i_n}(h_n)$ define*

$$M_n\left(\hat{\varphi}_{i_n}^{(n)}, e_{i_n}\right) := \sum_{j=1}^{N_n} \mathbb{1}\left(k\left(\frac{1}{h_n} \left\| \hat{\varphi}_{i_n}^{(n)} - \hat{\varphi}_{j_n}^{(n)} \right\| \right) > e_{i_n}\right).$$

Then for any $e_{i_n} \in (0, K_{i_n}(h_n))$,

$$\text{plim}_{N_n \rightarrow \infty} \frac{M_n\left(\hat{\varphi}_{i_n}^{(n)}, e_{i_n}\right)}{N_n} \geq c_{i_n} \in (0, 1]. \quad (16)$$

for all $i_n \in 1, \dots, N_n$.

The upper bound on e_{i_n} , K_{i_n} , is expected to be $k(0)$ for most classes of kernels. That is, the kernel function evaluated at $\hat{\varphi}_{i_n}^{(n)} = \hat{\varphi}_{j_n}^{(n)}$ should maximise the value of the kernel function. For example, the Gaussian kernel function is maximised at $k(0) = (1/\sqrt{2\pi})$. An example of a low-level condition for Assumption 8 is presented in Remark 3.

Assumption 8 is a restriction on the data generating process of the fixed-effect proxy parameter space. This is similar to Assumption 5.5 in Altonji and Matzkin (2005), except in this case related to the fixed-effect parameter space. Note that for these to be satisfied, the probability over the support of the fixed-effect space must be strictly positive. Whilst this paper focuses on fixed-effects, that is, effects that are taken as given and not modelled as random variables, it is still useful to understand that these parameters are sampled from some space. This is the space that the restriction in Assumption 8 pertains to.

Assumption 8 places a restriction on the bounds of the kernel function and on the space of proxy measures used. The restriction on the proxy measures may be satisfied if they are generated such that the neighbourhood around each realisation grows proportionally with the sample size, such as in Remark 3,

Remark 3 (Regularity of proxy measures). *Let $\hat{\varphi}_{i_n}^{(n)} \in \hat{\Phi}_n$ and redefine $M_n^\varepsilon\left(\hat{\varphi}_{i_n}^{(n)}\right)$ as*

$$M_n^\varepsilon\left(\hat{\varphi}_{i_n}^{(n)}\right) := \sum_{j_n=1}^{N_n} \mathbb{1}\left(\hat{\varphi}_{j_n}^{(n)} \in B_\varepsilon\left(\hat{\varphi}_{i_n}^{(n)}\right)\right),$$

where $B_\varepsilon(x)$ is the ε -neighbourhood around x . Then Assumption 8 is satisfied for the dimension n for any $\varepsilon > 0$ if,

$$\text{plim}_{N_n \rightarrow \infty} \frac{M_n^\varepsilon\left(\hat{\varphi}_{i_n}^{(n^*)}\right)}{N_n} \geq c_{i_n} \in (0, 1].$$

for all $\hat{\varphi}_{i_n}^{(n)} \in \hat{\Phi}_n$.

Defined in Remark 3 is essentially the building blocks of the probability space for the fixed-effects proxy parameter space. This is a lower level restriction than Assumption 8 in that it relates to the space of the proxy measures without reference to the kernel functions used.

Assumption 9 (Regularity conditions). *Let $\tilde{T}_{i_1, \dots, i_d}$ be the entries of tensor \mathbf{T} after the kernel weighted fixed-effects from the minimiser of (15) are differenced out. Then,*

(i). $\left(\frac{1}{\prod_n N_n} \sum_{i_1} \cdots \sum_{i_d} \tilde{X}_{i_1, \dots, i_d} \tilde{X}'_{i_1, \dots, i_d} \right) = O_p(1)$ converges to a nonrandom positive definite matrix as $N_1, \dots, N_d \rightarrow \infty$.

(ii). $\frac{1}{\prod_n N_n} \sum_{i_1} \cdots \sum_{i_d} \tilde{X}_{i_1, \dots, i_d} \varepsilon_{i_1, \dots, i_d} = O_p \left(\frac{1}{\sqrt{\prod_n N_n}} \right)$.

Assumption 9 is exactly Assumption 6 but with the kernel weighted fixed-effects in place of the group fixed-effects.

Proposition 3 (Upper bound on kernel estimator). *Let the class of kernel functions used to formulate the weights and the proxy measure used in these kernel functions for the kernel weighted fixed-effects estimator satisfy Assumption 7 and 8. Also, let Assumption 9 hold for the set of regressors. Let $\left\| \varphi_{i_{n^*}}^{(n^*)} - \hat{\varphi}_{i_{n^*}}^{(n^*)} \right\|^2 = O_p(C_{n^*}^{-2})$ for $n^* \in \mathcal{M}'$ and $\left\| \varphi_{i_{n^*}}^{(n^*)} - \hat{\varphi}_{i_{n^*}}^{(n^*)} \right\|^2 = O_p(1)$ for $n' \notin \mathcal{M}'$, where \mathcal{M}' is a non-empty subset of dimensions. Let h_n be the bandwidth parameter from Assumption 8. Then, for L_N that may depend on sample size,*

$$\left\| \hat{\beta}_{KER, \mathcal{W}} - \beta^0 \right\| = \sqrt{L_N} O_p \left(\prod_{n^* \in \mathcal{M}'} \sqrt{O_p(C_{n^*}^{-2}) + O_p(C_{n^*}^{-1} h_{n^*}^\alpha) + O_p(h_{n^*}^{2\alpha})} \right) + O_p \left(\prod_{n=1}^d \frac{1}{\sqrt{N_n}} \right).$$

For $h_n^\alpha \lesssim O(C_n^{-1})$ this reduces to

$$\left\| \hat{\beta}_{KER, \mathcal{W}} - \beta^0 \right\| = \sqrt{L_N} O_p \left(\prod_{n^* \in \mathcal{M}'} O_p(C_{n^*}^{-1}) \right) + O_p \left(\prod_{n=1}^d \frac{1}{\sqrt{N_n}} \right).$$

Proposition 3 shows that the convergence rate for the kernel estimator is bounded by the convergence rate of the proxy estimates. That is, as long as the bandwidth parameter approaches zero sufficiently fast, the kernel estimator converges at a rate no worse than the convergence of the proxies when proxies are estimations of the true parameter values at or slower than $\sqrt{N_n}$ -convergence. This is expected and also a good result that the kernel method does not hinder the convergence rate from these proxies. These kernel methods do, however, suffer a curse of dimensionality since Assumption 8 becomes increasingly difficult to justify as the dimension of the fixed-effects increases. Discussions in Section 4.2 suggest $O_p(C_{n^*}^{-1}) = O_p(1/\sqrt{N_{n^*}})$. This shows, as in the group fixed-effects method, the parametric rate is attainable if $\mathcal{M}' = \{1, \dots, d\}$ and L_N is bounded.

Remark 4. *The motivation for the kernel weighted fixed-effect estimator is very similar to the group fixed-effect estimator. Take (3), stated again here in the three-dimensional setting for the kernel weighted transformation,*

$$\check{Y}_{ijt} = Y_{ijt} - \bar{Y}_{i^*jt} - \bar{Y}_{ij^*t} - \bar{Y}_{ijt^*} + \bar{Y}_{i^*j^*t} + \bar{Y}_{i^*jt^*} + \bar{Y}_{ij^*t^*} - \bar{Y}_{i^*j^*t^*}.$$

*In the Introduction, the within-cluster transformation took the average within each starred indices' cluster. For the kernel weighted difference this average is instead taken as a weighted average over the whole sample. For example, the term $\mathcal{A}_{i^*j^*t}$ from this is,*

$$Y_{i^*j^*t} = \sum_{i'=1}^{N_1} \sum_{j'=1}^{N_2} w_{i,i'}^{(1)} w_{j,j'}^{(2)} Y_{i'j't},$$

where $w_{i,i'}^{(1)}$ and $w_{j,j'}^{(2)}$ are defined in (14). The arguments for the group fixed-effects estimator then translate directly to the kernel weighted fixed-effect estimator, where smooth weights are applied instead of the binary weights implied by the within-cluster differencing.

4 Discussion of estimators

This section serves to discuss the results in Section 3, motivate further some of the chosen methods, and provide some methods to estimate cluster assignments or proxies for kernel weights. A few iteration procedures are also discussed at the end of this section.

4.1 Matrix method results

As stated already, Proposition 1 takes for granted the correct choice in dimension to flatten across before using the least square method in Bai (2009). Under Assumption 1.(ii) the singular values of the flattened normalised noise term dissipates as follows;

$$\frac{1}{\sqrt{\prod_n N_n}} \|\varepsilon_{(n)}\| = O_p \left(\frac{1}{\sqrt{\min\{N_n, \prod_{m \neq n} N_m\}}} \right).$$

Since \mathcal{A} is a collection of fixed-effects, the normalised singular values of its flattenings are $O_p(1)$, that is, the singular values are not asymptotically negligible like those of the noise term.⁴ This

⁴To see this consider the standard two dimension model and take the Frobenius norm any arbitrary component of the interactive fixed-effects term, $\lambda_r f'_r$, normalised by $1/\sqrt{NT}$

$$\frac{1}{\sqrt{NT}} \|\lambda_r f'_r\|_F = \sqrt{\frac{1}{NT} \sum_i \sum_t (\lambda_{ir} f_{tr})^2} = \sqrt{\frac{1}{N} \sum_i \lambda_{ir}^2} \sqrt{\frac{1}{T} \sum_t f_{tr}^2} = O(1).$$

The last equality comes from λ_{ir} and f_{tr} being bounded fixed-effects.

ensures that, after flattening \mathcal{A} , each of the singular values eventually dominate those of the noise term. These conditions make up similar restrictions imposed in Ahn and Horenstein (2013) that allow for the use of the eigenvalue ratio test (ER) to diagnose the number of factors. Hence, in large samples, the analyst may be able to use this test or similar to not only decide how many factors to use but also decide which dimension is likely to be low-dimensional.

Consider the factor model applied to a flattening that may not be low-rank. For a concrete example of when this can occur see the data generating process in the simulations in Section 5, where $\varphi^{(n)}$ are designed to be low-dimensional for some n , and high-dimensional otherwise. Along the dimensions of \mathcal{A} that do not conform to the low-rank assumption in Assumption 3, the tail singular values may become difficult to discern from the singular values of the noise term in small samples. This means variation from those tail factors are less likely to be projected out from the factor model unless many factors are used in this projection. If r_n for $n \notin \mathcal{L}$ is allowed to increase adversely, for example at exactly the upper bound, then factor projection may never sufficiently project all relevant factors. Also, as the number of estimated factors increases, Assumption 4 becomes harder to satisfy since variation in the set of covariates is also projected out. This demonstrates the importance of choosing the correct dimension to flatten over, which is supported by the simulation results in Section 5.

Hence, a standard factor model that estimates at least r_n factors should result in consistent estimation of the slope coefficients, see Moon and Weidner (2017). However, this relies on an important structural feature of the unobserved heterogeneity term. When flattened in the chosen dimension – the first dimension in the above example – the rank of the matrix after flattening must be low relative to data size. This implies that to successfully project out the variation in the fixed-effect term either the matrix of fixed-effects from any flattening is low-rank, or, at least one flattening leads to a low-rank matrix of fixed-effects and the analyst knows which flattening this is. To use the above example again, this means the analyst knows that $\varphi^{(1)}\Gamma'_1$ is low-rank, hence flattening in the first dimension is the correct way to recast the model to a panel data model, and so forth for the other flattenings. Whilst requiring low-rankness in at least one dimension may be an acceptable restriction, having knowledge of which dimension this low-rankness resides in is potentially more restrictive.

To understand the problem, consider the following two examples, one where the flattening is not low-rank and one where it is. First, assume $\varphi^{(1)}$ varies in a high-dimensional parameter space, e.g. with $N_1 < N_2N_3$, $\varphi^{(1)} \in \mathbb{R}^{N_1 \times N_1}$ and $\Gamma \in \mathbb{R}^{N_2N_3 \times N_1}$ with each column mutually orthogonal for both these matrices. Then the product of these matrices is full-rank and any factor projection approach will not fully control for this term. On the contrary, consider $\varphi^{(1)} \in \mathbb{R}^{N_1 \times N_1}$ where all columns are linearly dependent. Then the matrix $\varphi^{(1)}\Gamma'$ is rank-1 regardless of L and of how $\varphi^{(2)}$ and $\varphi^{(3)}$ vary, thus can be projected with a factor model estimated with 1 factor.

Hence it is important which dimension the analyst chooses to flatten over.

Well established diagnostics in Bai and Ng (2002), Ahn and Horenstein (2013) and Hallin and Liška (2007) can be used to determine the number of factors. These diagnostics can be repeated across different flattenings, which may be informative of the dimension to use for flattening. Note these procedures require an initial guess of β and relies on this guess not eradicating the factor structure in the residual; see the beginning of Section 4.4 for a concrete example of this. It should also be noted that these diagnostics are not without restrictions and can lead to spurious conclusions on the optimal number of factors. For example, the eigenvalue ratio test in Ahn and Horenstein (2013) can undershoot the number of factors when singular values decay quickly for the leading few factors. This does not interfere with the asymptotic result in that paper but can have implications in small sample estimation. Indeed, however, these diagnostics can be helpful in both the matrix recasting of the problem and the group fixed-effects estimation in the sequel.

4.2 Estimating cluster and kernel proxies

Discussed here are some important functionals of multidimensional arrays that are useful for estimating proxies to cluster on or use for kernel weights. This includes a discussion on how to uncover proxies from multidimensional array data, and a discussion on why standard matrix methods do not extend well to the multidimensional setting.

First, consider how to cluster in the within-cluster transformation. In most clustering algorithms, for example K -means or K -nearest neighbour, there is some notion of a distance metric between units considered for each cluster. To arrive at a distance there must be some space to measure that distance over. For example, using some vector u_i and the Euclidean norm of differences, $\|u_i - u_j\|$, to measure the distance between units. Algorithms to arrive at these groupings are well established when the distance metric and variable to take distance over are given. However, in this setting there is no clear variable through which to take distance over. Motivated here are methods to extract proxies that serve to measure distance across units in a way that isolates variation in each dimension of the unobserved heterogeneity term.

It is important to find proxies that isolate variation in each dimension since clustering is to be performed one index at a time. Discussed here are decompositions of multidimensional arrays that can perform this, Kolda and Bader (2009) contains a nice summary of some candidate decompositions. The method discussed here uses the higher order singular value decomposition (HOSVD), and focuses on components of this decomposition that have well formulated theoretical properties. The HOSVD is traditionally used in pursuit of a low-rank tensor decomposition by either direct truncation of left singular vectors or by some iteration approach similar to this, see for example the higher order orthogonal iteration scheme. The problem of direct truncation, however, is not well-posed because the solution to the low-rank tensor problem may not be

unique and reformulating the original tensor after the aforementioned truncation is not guaranteed to be lower tensor rank. See De Silva and Lim (2008) for an extensive explanation of the ill-posedness issues. Hence, this method cannot be used in the pursuit of analytic consistency results. Problems also arise in this setting where the reformulated tensors can be arbitrarily well approximated by a tensor of lower tensor rank, which is a result of the border rank issue of the tensor rank decomposition.

Reconsider the three dimensional model with heterogeneity of the following form

$$\mathcal{A}_{ijt} = \sum_{\ell=1}^L \varphi_{i\ell}^{(1)} \varphi_{j\ell}^{(2)} \varphi_{t\ell}^{(3)}. \quad (17)$$

As shown in De Silva and Lim (2008), the Eckart–Young–Mirsky theorem cannot be relied upon to guarantee an optimal low-rank approximation for the multidimensional array \mathcal{A} . This motivates the use of the group fixed-effects and kernel weighted fixed-effects as alternative solutions.

Also reconsider the singular value decomposition for matrices, applied to each of the n -flattenings of $\mathcal{A} \in \mathbb{R}^{N_1 \times \dots \times N_d}$ as

$$\mathcal{A}_{(n)} = U^{(n)} \Sigma_n V^{(n)'} \quad (18)$$

By the same logic as in the matrix case and formalised with the Eckart-Young-Mirsky theorem, variation over the rows of each $U^{(n)}$ explains variation over the n^{th} dimension of the multidimensional array \mathcal{A} . Thus, if $\mathcal{A}_{(n)}$ is low rank, the leading few columns of $U^{(n)}$ provide good proxies for closeness in n^{th} dimension of \mathcal{A} . If $\mathcal{A}_{(n)}$ is not low rank then these leading columns still provide the best proxies for bias reduction using the group or kernel fixed-effects. Hence, by reconsidering the tensor problem as a sequence of matrix problems, the usual singular value decomposition properties can be utilised for this reduction. This shows that a simple rearrangement of the data provides readily available techniques to measure closeness in each dimension separately.

Consider for any of the dimensions n the corresponding matrix of left singular vectors from above, $U^{(n)}$, estimated with noise ε_{ijt} . That is, each U_n are calculated from the object $\mathcal{V} = \mathcal{A} + \varepsilon$. Under reasonable regularity conditions on the noise term ε , the left singular vectors from this decomposition comprise of a signal of the underlying fixed-effect parameter and noise from ε . For example, in the three dimensional case, define the L_1 -vector $\widehat{U}_i^{(1)}$ as the i -th row of the left singular matrix of $\mathcal{A} + \varepsilon$ flattened in the first dimension. Then the vector $\widehat{U}_i^{(1)}$ may comprise of,

$$\widehat{U}_i^{(1)} = \varphi_i^{(1)} + O_p \left(\frac{1}{\sqrt{\min\{N_1, N_2 N_3\}}} \right).$$

Likewise, for any dimension n define $\widehat{U}^{(n)}$ as the matrix of singular vectors from $\mathcal{A} + \varepsilon$ flattened in the n -th dimension, where \mathcal{A} is the unobserved fixed-effects component of interest and ε is the usual idiosyncratic noise term. Then, Bai and Ng (2002) detail conditions required for the following “up-to-rotation” consistency result, which has been amended to this paper’s setting;

Lemma 1 (Theorem 1 from Bai and Ng (2002)). *For any fixed integer $k \geq 1$, there exists an $(r_n \times k)$ matrix H_n^k with $\text{rank}(H_n^k) = \min\{k, r_n\}$ and $C_n = \min\{\sqrt{N_n}, \prod_{n' \neq n} \sqrt{N_{n'}}\}$ such that for each n under some regularity conditions*

$$C_n^2 \left\| \widehat{U}_{i_n}^{(n)} - H_n^{k'} \varphi_{i_n}^{(n)} \right\|^2 = O_p(1).$$

This establishes a consistency result for estimating cluster proxies and suggests these left singular vectors are viable options to cluster in each dimension. It also makes concrete the limitation implied by the value of C_n for each index - that short indices have poorly estimated proxies. However, given that the error term displayed in (13) is multiplicative across dimension, the error from this poor approximation should become negligible as long as enough other dimension proxies are well estimated. Also, the presence of the rotation matrices, H_n^k , in Lemma 1 can be ignored since these do not change relative distances of each unit under standard distance metrics used to cluster.

4.3 Group fixed-effect convergence result

Before discussing the result from Proposition 2, an alternative restriction on the cluster assignments is proposed. If clustering is performed on a proxy measure of the fixed-effect then Assumption 5 can be stated in terms of the proxies, which forms the statement of Remark 5. This requires that the proxies form an injective mapping to the true fixed-effect parameters. An example of this are the conditions imposed in Freeman and Weidner (2022), stated in similar terms here:

Remark 5 (Clustering). *The statement of Assumption 5 can be reformulated in terms of the cluster proxies as follows. Let $\widehat{\varphi}_{i_n}^{(n)} := \widehat{\varphi}^{(n)}(\varphi_{i_n}^{(n)}) \in \mathbb{R}^{\widehat{r}_n}$ be the proxy for individual i_n used to cluster along dimension n . Then,*

(i). *For all n as $N_n \rightarrow \infty$,*

$$\frac{1}{N_n} \sum_{i_n}^{N_n} \left\| \widehat{\varphi}_{i_n}^{(n)} - \widehat{\varphi}_{j_n(i_n)}^{(n)} \right\|^2 \lesssim O_p(1)$$

(ii). *For a non-empty subset $\mathcal{M} \subset \{1, \dots, d\}$ take for any $n^* \in \mathcal{M}$ a sequence $\xi_{N_{n^*}} \rightarrow 0$ as*

$N_{n^*} \rightarrow \infty$. Then,

$$\frac{1}{N_{n^*}} \sum_{i_{n^*}}^{N_{n^*}} \left\| \hat{\varphi}_{i_{n^*}}^{(n^*)} - \hat{\varphi}_{j_{n^*}(i_{n^*})}^{(n^*)} \right\|^2 = O_p(\xi_{N_{n^*}})$$

(iii). Let $\varphi_{i_n}^{(n)} \in \Phi_n$ be the r_n -column vector of fixed effects, where Φ_n are convex sets for each n .

For each $a, b \in \Phi_n$ there exists a scalar $c_n > 0$ such that $\|a - b\| \leq c_n \cdot \|\hat{\varphi}^{(n)}(a) - \hat{\varphi}^{(n)}(b)\|$

If these alternate restrictions hold along with Assumption 6, then the bound in Proposition 2 holds for the GFE estimator.

Restrictions (i) and (ii) in Remark 5 are exactly Assumption 5.(i) and (ii) but with cluster proxies in place of the true parameter values. These are high level restrictions on the clustering mechanism that requires the mechanism to find closeness in the proxy space. Restriction (iii) in Remark 5 is an injectivity assumption on the proxy functions that demands closeness in the underlying parameter space given closeness in the proxy space. This requires that the proxies do actually provide a mapping to the true parameter space, that is, that they are reasonable proxies. An example of proxies that do this are the singular vectors from Section 4.2 that fit the requirements of Lemma 1. To see this expand the term $\|a - b\|$ and use the triangle inequality to see, $\|a - b\| \leq \|a - \hat{\varphi}^{(n)}(a)\| + \|\hat{\varphi}^{(n)}(b) - b\| + \|\hat{\varphi}^{(n)}(a) - \hat{\varphi}^{(n)}(b)\|$, where the first two terms are bound at the rate $O_p(C_n^{-1})$. Note the rotation matrices are ignored for brevity and C_n is the convergence rate from Lemma 1. Hence, asymptotically, Remark 5.(iii) can be achieved with $c_n = 1$.

This display also makes clear the bottle-neck when clustering in high-dimensional objects. The distance of the proxies, $\|\hat{\varphi}^{(n)}(a) - \hat{\varphi}^{(n)}(b)\|$, is difficult to bound using clustering methods when the dimension of the proxies are larger than two, see Graf and Luschgy (2002) and further discussion in Bonhomme, Lamadon and Manresa (2021). This implies that a low-dimensional set of proxies must bound the true parameter values for clustering methods to work well in this setting. Hence, whilst the relationship in restriction (iii) of Remark 5 may be satisfied for an arbitrarily high-dimensional set of proxies, for a reasonable family of cluster mechanisms to bound these proxies as per restrictions (i) and (ii) of this remark, restriction (iii) must also hold for a low-dimensional set of proxies. This can be highly restrictive with, for example, fixed-effects $\varphi_{i_n}^{(n)}$ that are high-dimensional.

How the sequences $\xi_{N_{n^*}}$ converges to zero and how L_N is bounded are important for the convergence result in Proposition 2. First note that for fixed L_N the first term in the result simplifies to $O_p(\prod_{n^* \in \mathcal{M}} \sqrt{\xi_{N_{n^*}}})$. Also note, if the conditions for Lemma 1 hold and clustering is based on singular vector estimates that adhere to Remark 5, then it is possible to achieve $\xi_{N_{n^*}} = O_p\left(\frac{1}{\min\{N_{n^*}, \prod_{n \neq n^*} N_n\}}\right)$. If each N_n grow at the same rate then the consistency result

is,

$$\|\widehat{\beta}_{GFE,C} - \beta^0\| = \sqrt{L_N} O_p \left(N_n^{-|\mathcal{M}|/2} \right).$$

In the worst case scenario $\sqrt{L_N}$ is upper bound by $\sqrt{L_N} \lesssim N_n^{(d-1)/2}$, which is taken from $L_N \leq \min_n \prod_{n' \neq n} N_{n'}$. The convergence result is then $O_p \left(N_n^{(d-1-|\mathcal{M}|)/2} \right)$, which is of course conservative but shows that if $|\mathcal{M}| = d$, then consistency is guaranteed albeit at the slow rate of $N_n^{1/2}$. This means that all dimensions must have good cluster assignments, which is obviously not an ideal worst case but shows the limitations of this method when L_N is unrestricted.

For the special case of $d = 3$ it can be shown that $L_N \leq \min_n \prod_{n' \neq n} r_{n'}$. From the discussion above, it is expected that $n \in \mathcal{M}$ is sufficient for $n \in \mathcal{L}$, that is, r_n is small for the set of dimensions $n \in \mathcal{M}$. This tightens the bound in Proposition 2 to $O_p \left(N_n^{\max\{-|\mathcal{M}|/2, 1-|\mathcal{M}|\}} \right)$, such that only $|\mathcal{M}| \geq 1$ is required for consistency. The analogous tensor rank bound is so far not known for the case with $d \geq 4$.

4.4 Iteration procedures

Consider taking cluster proxies from the estimated error term $\mathbf{W} = \mathbf{Y} - \mathbf{X}'\tilde{\beta}$. Define $\tilde{\beta}$ as the interim estimator used to obtain \mathbf{W} , and notice that this forms the basis of an iterative procedure, between forming clusters and estimating slope coefficients. This is illustrated in the following two-step procedure. For the below let \hat{r}_n be a hyperparameter that defines the number of singular vectors to use in the clustering stage.

1. For given $\tilde{\beta}$, take the left singular matrices from each n -flattening of $\mathbf{W} = \mathbf{Y} - \mathbf{X}'\tilde{\beta}$ to obtain $\{\hat{U}_1, \hat{U}_2, \hat{U}_3\}$.
2. Cluster on the leading \hat{r}_n columns of \hat{U}_n to generate cluster assignments in the n^{th} dimension. Use these cluster assignments in the within-cluster transformation on \mathbf{Y} and \mathbf{X} then perform pooled OLS to obtain $\hat{\beta}$.
3. Iterate steps 1 and 2 until convergence in the slope coefficients

This procedure may also be used as a debias estimator for a given initial estimate of $\tilde{\beta}$ by ignoring step 3. Iteration here may not be stable given that step 1 and 2 do not optimise the same objective function, hence for theoretical purposes it may be convenient to only consider this as a debias procedure. In practice, iterating between step 1 and 2 after some initial grid search to initialise β may be optimal.

Of course, other clustering or transformations may be used in place of the residual clustering and within-cluster transformation. In the below, two alternatives are provided. The first maintains the within-cluster transformation but considers a different set of proxies. The second

approach considers a kernel weighted transformation procedure that uses a generic set of proxies. At this stage and in the below estimator refinements the analyst may be concerned with the number of parameters required to conduct these transformation. Appendix C discusses a number of ways to reduce the size of the parameter space, including only projecting fixed-effects over a subset of dimensions and letting group sizes increase to reduce the number of groups.

Whether used as an iterative scheme or an update, the above method has some identification issues. As an illustration take the data generating process for model (1) with just one covariate,

$$X_{ijt} = -\mathcal{A}_{ijt} + \mu_{ijt},$$

where μ_{ijt} is a white noise term. Consider an initial guess of $\tilde{\beta} = 0$ when the true value is $\beta^0 = 1$. This leaves the residual term from Step 1 to base cluster assignment on as, $\mathbf{W} = \mathbf{Y} - \mathbf{X}\tilde{\beta} = \mathbf{Y}$, which reduces to $\mathbf{W} = \boldsymbol{\mu} + \boldsymbol{\epsilon}$. Thus, clustering is based solely on noise and can be reasonably described as random. The associated within-cluster transformation will not project variation in the \mathcal{A} terms that appear in both \mathbf{Y} and \mathbf{X} such that the OLS step in stage 2 produces

$$\hat{\beta} \approx \frac{Var(\mu_{ijt})}{Var(\mu_{ijt}) + Var(\mathcal{A}_{ijt})} + o_p(1).$$

For $\frac{Var(\mu_{ijt})}{Var(\mathcal{A}_{ijt})} \rightarrow 0$, $\hat{\beta} \rightarrow 0$ and the algorithm does not update the initial guess of $\tilde{\beta} = 0$. This problem also arises in the matrix methods in Section 3.1 and is a more fundamental issue with this algorithmic approach.

This example clearly displays some identification issues with the above method. Worth noting is that this may be alleviated with a grid search approach, though this can be computationally infeasible even for a moderate number of covariates since the grid grows exponentially in the number of covariates. To avoid this, proposed below is a method to extract cluster allocations from only variation in the set of covariates. As discussed below, this clustering may also be conducted on control variables extraneous to the regression line. The two-step procedure works as follows.

1. Take the left singular matrices from each n -flattening of X to obtain $\{\hat{U}_1, \hat{U}_2, \hat{U}_3\}$.
2. Cluster on the leading \hat{r}_n columns of \hat{U}_n to generate cluster assignments in the n^{th} dimension. Use these cluster assignments in the within-cluster transformation on \mathbf{Y} and \mathbf{X} then perform pooled OLS to obtain $\hat{\beta}$.

An advantage of using covariate clustering is that it can make use of control variables that are a good signal of cluster but are not included in the regression line. For example, a control variable Z_i that is constant across j and t may be a good candidate to cluster along the i dimensions but will be projected out with the within-cluster transformation, so cannot be used

directly in the pooled OLS estimation of β stage. This refinement also makes optimisation over β a convex problem, and no iteration is required because clustering is not a function of β estimates like in the first iteration procedure.

5 Simulation

Table 1 shows simulation results for the following DGP,

$$\begin{aligned} Y_{ijt} &= X_{ijt}\beta + \mathcal{A}_{ijt} + \mathcal{B}_{ijt} + \varepsilon_{ijt} \\ X_{ijt} &= \mathcal{A}_{ijt} + \mathcal{B}_{ijt} + \nu_{ijt} \end{aligned}$$

with $\mathcal{A}_{ijt} = \sum_{\ell=1}^{N_1} \varphi_{i\ell}^{(1)} \varphi_{j\ell}^{(2)} \varphi_{t\ell}^{(3)}$, $\mathcal{B}_{ijt} = \alpha_{ij} + \gamma_{it} + \delta_{jt}$. Also,

$$\begin{aligned} \varepsilon_{ijt}, \nu_{ijt}, \alpha_{ij}, \gamma_{it} \text{ and } \delta_{jt} &\overset{i.i.d.}{\sim} N(0, 1) \text{ and for each } \ell, \varphi_{i\ell}^{(1)}, \varphi_{t\ell}^{(3)} \overset{i.i.d.}{\sim} N(0, 1). \\ \varphi_{j1}^{(2)} &\overset{i.i.d.}{\sim} N(0, 1) \text{ with } \varphi_{j1}^{(2)} = \varphi_{j2}^{(2)} = \dots = \varphi_{jN_1}^{(2)} \end{aligned}$$

\mathcal{A}_{ijt} and \mathcal{B}_{ijt} are normalised to have unit variance. \mathcal{A} is specified such that it is rank 1 when flattened in the second dimension and rank N_1 when flattened in either dimension one or three. That is, the multilinear rank is $\mathbf{r} = (N_1, 1, N_1)$. This comes directly from the data generating process for each $\varphi^{(n)}$, where the matrix $\varphi^{(2)}$ is designed to be rank-1 and the matrices $\varphi^{(1)}$ and $\varphi^{(3)}$ are designed to be rank- N_1 .

In Table 1, the estimators OLS and Fixed-effects are simply the pooled OLS estimator and the pooled OLS estimator after additive fixed-effects are projected out, respectively. As expected both of these two have poor bias. The four GFE estimators perform well with reasonably low bias and standard deviation. GFE (K-means) is GFE estimator with clustering based on the K-means algorithm, with proxies taken from the residual. GFE (K-means on X) is the same estimator with proxies taken from the scalar covariate of interest. GFE (1-NN) and GFE (1-NN on X) are likewise the same estimators but using the one nearest neighbours clustering. The factor model is used after first flattening along each dimension as Factor(dim = n), where n is the dimension used for flattening. In each case, 2 factors are projected. The results show the theoretical result succinctly, where the bias is close to zero when the correct dimension is flattened over (the second dimension in this case) and very poor bias when the incorrect dimension is used (the first and third dimensions). Lastly, the kernel differencing estimator is estimated with Gaussian kernel function with bandwidths 0.5, 1 and 1.5; which are standardised to be equivalent to standard deviations of the proxy measures. All kernel estimators have comparable bias but substantially better standard deviation for bandwidth equal 1 and 1.5.

This analysis is repeated for the four dimensional case in Table 2, where the second and third dimensions admit low-dimensional unobserved interactive fixed-effects parameters. For

computational reasons, the GFE nearest neighbour estimators are omitted. The simulations suggest similar results as the three dimensional case, where the factor models perform well when flattened in the low-dimensional dimensions (second and third) and poorly in the high-dimensional dimensions (first and fourth).

3-D	Mean bias	St. dev.	MSE
OLS	0.6668	0.0033	4.45e-01
Fixed-effects	0.4997	0.0114	2.50e-01
GFE (K-means)	0.0118	0.0096	2.32e-04
GFE (K-means on X)	0.0129	0.0112	2.91e-04
GFE (1-NN)	0.0112	0.0153	3.61e-04
GFE (1-NN on X)	0.0111	0.0154	3.61e-04
Kernel (h = 0.5)	0.0030	0.0090	8.94e-05
Kernel (h = 1.0)	0.0031	0.0068	5.58e-05
Kernel (h = 1.5)	0.0037	0.0062	5.24e-05
Factor (dim = 1)	0.4319	0.0135	1.87e-01
Factor (dim = 2)	0.0030	0.0050	3.40e-05
Factor (dim = 3)	0.4319	0.0135	1.87e-01

Table 1: 3D model ($N_1 = N_2 = N_3 = 36$), with 10,000 Monte Carlo rounds. All results are in relation to β estimation.

4-D	Mean bias	St. dev.	MSE
OLS	0.6670	0.0018	4.45e-01
Fixed-effects	0.4981	0.0282	2.49e-01
GFE (K-means)	0.0012	0.0049	2.58e-05
GFE (K-means on X)	0.0013	0.0051	2.82e-05
Kernel (h = 0.5)	5.23e-05	0.0114	1.00e-04
Kernel (h = 1.0)	4.42e-05	0.0057	3.26e-05
Kernel (h = 1.5)	-1.32e-05	0.0045	2.03e-05
Factor (dim = 1)	0.3733	0.0311	1.40e-01
Factor (dim = 2)	0.0030	0.0030	1.82e-05
Factor (dim = 3)	0.0030	0.0030	1.82e-05
Factor (dim = 4)	0.3734	0.0311	1.40e-01

Table 2: 4D model ($N_1 = N_2 = N_3 = N_4 = 20$), with 10,000 Monte Carlo rounds.

A two-dimensional simulation exercise is also performed to compare the grouped fixed-effects approach to the factor model approach in a setting where theoretical results for the factor model are well known. Table 3 shows the results of this two-way setting where the data generating process is a factor model with two factors. The GFE estimators have less bias than the factor model even when the factor model overestimates the number of factors. To see this, compare the factor estimates with 2, 4 and 6 factors projected out with the GFE estimator. For increase in variance of order ≈ 4 , the GFE estimator reduces bias by an order ≈ 10 . This is a surprising improvement in estimates for a setting that is purpose designed for the factor model. Where this comparison falls down is for models with a larger number of factors because generally clustering does not perform well when the latent parameter space has dimension greater than 2.

	Mean bias	St. dev.	MSE
OLS	0.6672	0.0033	4.45e-01
Fixed-effects	0.5000	0.0043	2.50-e01
GFE	0.0002	0.0090	8.05e-05
Kernel ($h = 0.5$)	0.0003	0.0055	3.06e-05
Kernel ($h = 1.0$)	0.0003	0.0054	2.87e-05
Kernel ($h = 1.5$)	0.0004	0.0053	2.81e-05
Factor ($R = 2$)	0.0024	0.0047	2.76e-05
Factor ($R = 4$)	0.0031	0.0048	3.25e-05
Factor ($R = 6$)	0.0024	0.0049	3.03e-05

Table 3: 2D model ($N_1 = N_2 = 216$), with 10,000 Monte Carlo rounds.

6 Empirical application - demand estimation for beer

The methods proposed in this paper are applied to estimated the demand elasticity for beer. Price and quantity for beer sales is taken from the Dominick’s supermarket dataset for the years 1991-1995 and is related to supermarkets across the Chicago area. Price and quantity vary over three dimensions in this example – product (i), store (j) and month (t). Fixed-effects that interact across all three dimensions can control for taste shocks to beer consumption that differ over both product and store. Take for instance a large sporting event (temporary t shock) that changes preferences differently across locations (j) and across certain subsets of sponsored beer (i). For example, in the stadiums for the many NBA finals playoffs the Chicago Bulls played in the early 1990’s, Miller Lite beer advertisements could be seen alongside advertisements for the substitute product Canadian Club whisky. This suggests these events attracted large marketing campaign spends for these and other beer substitute brands that most likely also included price

offers at local supermarkets. Whilst the impact of these advertisements and price offers on the demand for or price of beer is not clear and, further, that it is reasonably safe to assume the econometrician does not observe the plethora of marketing campaigns around these events, the analyst would most likely still want to control for aggregate shocks like these. For this reason it is important to use methods that robustly control for unobserved fixed-effects, such as unobserved marketing campaigns, that may impact both quantity demanded and prices in unforeseen ways.

Models for demand estimation ideally account for endogenous variation in prices and quantity. The classic instrumental variable approach is to find a variable that varies exogenously to the production process but can reasonably describe price fluctuations. A popular instrument in the estimation of beer demand is the commodity price for barley, one of the product’s main ingredients, see e.g. Saleh (2014); Tremblay and Tremblay (1995); Richards and Rickard (2021). Since the price of barley is arguably not driven by the demand for it by any one supplier of beer, it can be a useful variable to instrument for price shifts. In the following, it is taken as given that the price of barley is exogenous with respect to the noise term, ε .

For validity the instrument is also required to be strong, in the sense that it is strongly correlated with price. In this dataset correlation between the price of barley, which varies over only t , and price of beer depends on how beer price is aggregated. If beer price is first integrated over i and j , such that it only varies over t , then it is highly correlated with the price of barley, at 0.61. However, if beer price is not aggregated at all it is only correlated at 0.05. This suggests there are important product and store level price drivers for beer that are not accounted for by fluctuations in the price of barley. This implies that price fluctuations in barley alone may not be viable to fully capture beer prices when considering variation over all three dimensions. For exogeneity, the price of barley must be independent of common unobserved shocks to both price and demand, which translates to being independent of $\varphi_{t,\ell}^{(3)}$ and any scalar fixed-effects that vary over t in the interactive fixed-effects model. More details are deferred to Appendix B.

An alternative method is to follow the external variable approach of Altonji and Matzkin (2005), that requires exogeneity of beer price conditional on an external variable. If the price of barley is treated as the external variable, the beer demand example requires conditional exogeneity of beer prices across time periods of similar barley prices. For example, that the price of beer is exogenous when compared to other beer prices during time periods of high barley price, and so on. Note, the external variable approach provides much more flexibility than the traditional IV approach in the multidimensional setting because it may viably be used in conjunction with the group and kernel fixed-effects approaches even when the external variable varies over a subset of dimensions. See Appendix B for an explanation of how external variables are used in this setting and how the assumption works in conjunction with the fixed-effects estimators.

The second column from Table 4 refers to the estimates for demand elasticities for the following regression model,

$$\log(quantity_{ijt}) = \log(price_{ijt})\beta + \mathcal{A}_{ijt} + \varepsilon_{ijt} \quad (19)$$

where \mathcal{A}_{ijt} is the usual interactive fixed-effects term from the prequel. This amounts to estimating the standard log-log model for demand with fixed-effects. That is,

$$quantity_{ijt} = price_{ijt}^{\beta} \exp(\mathcal{A}_{ijt} + \varepsilon_{ijt}).$$

Again, no controls are included here since they are low-dimensional and subsumed by the fixed-effects term. This model specification estimates reasonably similar elasticities as the logit case across each of the different fixed-effects estimators but relatively large differences in estimates for pooled OLS, IV and external variable estimators. The similar elasticities for the different fixed-effects estimators within Table 4 again suggests that whilst some form of fixed-effects should be included, they may not need be as complex as implied by the GFE and kernel methods.

The third column from Table 4 reports estimates of the same log-log model controlling for the average log price of other products,

$$\log(quantity_{ijt}) = \log(price_{ijt})\beta + \delta \sum_{i' \neq i} \log(price_{i'jt}) + \mathcal{A}_{ijt} + \varepsilon_{ijt}. \quad (20)$$

This model assumes homogeneous cross-elasticity over all other beer products. That is, it refers to the demand model,

$$quantity_{ijt} = price_{ijt}^{\beta} \prod_{i' \neq i} price_{i'jt}^{\delta_{ii'}} \exp(\mathcal{A}_{ijt} + \varepsilon_{ijt}),$$

where $\delta_{ii'} = \delta$ for all i and i' . Whilst this may oversimplify the system of cross-elasticities in the market for beer, it does significantly change the estimates for β in the log-log model. This suggests that cross-elasticities should probably be controlled for since β estimates do seem sensitive to their inclusion. It also shows that for a control variable with full rank variation over all dimensions, not even the more complex fixed-effects estimators can control for these. Note that most estimators returned a negative value for δ , which opposes the theory that other brands of beer, on aggregate, are substitutes. However, since prices are aggregated in such a crude way, the cross-elasticity estimates should not be taken too seriously. If interested in the cross-elasticities, then some care should be taken to segment or group products in such a way that actual substitution is being identified here, not just aggregate market forces. For this model, all fixed-effects estimates are within statistical noise of each other, this time with the external variable approach being closely aligned. These are also similar to the own-price

Estimator	$\hat{\beta}$ (St. dev.) no cross elas.	$\hat{\beta}$ (St. dev.) with cross elas.
Pooled OLS	0.06 (0.31)	0.09 (0.30)
Pooled IV	-4.76 (2.76)	-3.47 (2.30)
Pooled EV	-3.05 (0.48)	-2.74 (0.27)
Additive FE	-1.78 (0.34)	-2.80 (0.28)
Factor (dim = 1)	-1.61 (0.30)	-2.65 (0.25)
Factor (dim = 2)	-1.78 (0.33)	-2.78 (0.27)
Factor (dim = 3)	-2.09 (0.32)	-2.89 (0.27)
GFE	-1.85 (0.33)	-2.86 (0.30)
GFE (EV)	-1.84 (0.36)	-2.70 (0.30)
Kernel (Gaussian)	-1.72 (0.33)	-2.58 (0.29)
Kernel (EV, Gaussian)	-1.97 (0.47)	-2.81 (0.43)

Table 4: Log-log demand elasticities (73 products, 41 stores, 57 months).

Standard deviations were bootstrapped by resampling along each dimension separately. In the first dimension, product 1 is fixed across bootstrap samples. Column 2 displays estimates for the model (19) with no cross elasticities. Column 3 displays estimates for the model (20), which controls for cross elasticities.

elasticity estimates from Table 1 in Hausman, Leonard and Zona (1994). IV is estimated with very high variation in both log-log models, which may be due to barley being a weak instrument.

Table 5 refers to estimates from the standard logit demand model,

$$\log(\text{quantity}_{ijt}) - \log(\text{quantity}_{1jt}) = \text{price}_{ijt}\beta + \mathcal{A}_{ijt} + \varepsilon_{ijt}$$

where \mathcal{A}_{ijt} is the usual interactive fixed-effects and no control variables are included since the set of available controls are rank-deficient and automatically projected out with standard scalar fixed-effects and from differencing out the outside option. The outside option is encoded as product number 1 and is the aggregate consumption of products with small quantities consumed. This serves the purpose of creating an outside option to do the necessary logit demand transformation as well as to avoid issues related to an unbalanced panel for the many niche products with sparse consumption amounts. Own price elasticity is calculated as $\eta_{ijt} = \text{price}_{ijt}\beta(1 - \text{quantity}_{ijt}/\sum_{ijt} \text{quantity}_{ijt})$ and the mean elasticity is taken as the mean of this measure for each estimator. The pooled instrumental variable and external variable estimates estimate relatively large elasticities. However, all of the fixed-effects approaches estimate statistically similar slope coefficients and elasticities at the mean. This implies that whilst some fixed-effects may exist in the true model for demand, they are unlikely complex enough to require the high-dimensional projections from the GFE or kernel methods. To robustly test for

the existence of fixed-effects in an IV model there must be an instrument with variation over all dimensions such that fixed-effects can be projected out alongside the IV model. This of course also takes for granted that the IIA logit model is the true model for demand.

Estimator	Coefficient (bootstrap st. dev.)	Elasticity at mean
Pooled OLS	-0.60 (0.04)	-3.26 (0.22)
Pooled IV	-0.72 (0.04)	-3.91 (0.22)
Pooled EV	-0.71 (0.05)	-3.86 (0.27)
Additive FE	-0.32 (0.05)	-1.74 (0.27)
Factor (dim = 1)	-0.29 (0.04)	-1.58 (0.22)
Factor (dim = 2)	-0.32 (0.05)	-1.74 (0.27)
Factor (dim = 3)	-0.37 (0.05)	-2.01 (0.27)
GFE	-0.32 (0.05)	-1.74 (0.27)
GFE (EV)	-0.31 (0.05)	-1.68 (0.27)
Kernel (Gaussian)	-0.30 (0.05)	-1.63 (0.27)
Kernel (EV, Gaussian)	-0.34 (0.08)	-1.85 (0.44)

Table 5: Logit demand estimates (73 products, 41 stores, 57 months).

Standard deviations were bootstrapped by resampling along each dimension separately. In the first dimension, product 1 is fixed across bootstrap samples as the outside option and the remaining products are resampled with replacement.

7 Conclusion

This paper shows methods to generalise the interactive fixed-effect to the multidimensional case with more than two dimensions. Theoretical results show that standard matrix methods can be applied to this setting but require additional knowledge of the data generating process. The multiplicative interactive error from the group fixed-effects and kernel methods show a potential improvement on the asymptotic rate of convergence and suggest a more robust approach to projecting fixed-effects. Simulations corroborate these theoretical results and show the relative advantage of using a standard factor model when the structure of the interactive term is known. They also show the robustness of the group fixed-effects estimator to not having this same knowledge. Inference in these models is still an open question for further research.

The model is applied to a simple demand model for beer consumption. The application demonstrated how the GFE and kernel methods integrate well with the external variable approach in Altonji and Matzkin (2005) as apposed to instrumental variable approaches that do

not allow for fixed-effect estimation when the instrument is rank-deficient. The application shows that whilst some fixed-effects should likely be included in the model for beer demand, they are unlikely to be overly complicated to justify the GFE or kernel methods. This is a useful analysis, as it provides a robustness check for the specification of fixed-effects in model specifications.

References

- Ahn, S. C. and A. R. Horenstein (2013). Eigenvalue ratio test for the number of factors. *Econometrica* 81(3), 1203–1227.
- Altonji, J. G. and R. L. Matzkin (2005). Cross section and panel data estimators for nonseparable models with endogenous regressors. *Econometrica* 73(4), 1053–1102.
- Bai, J. (2009). Panel data models with interactive fixed effects. *Econometrica* 77(4), 1229–1279.
- Bai, J. and S. Ng (2002, January). Determining the number of factors in approximate factor models. *Econometrica* 70(1), 191–221.
- Bonhomme, S., T. Lamadon, and E. Manresa (2021). Discretizing unobserved heterogeneity. *Econometrica* (Forthcoming).
- Chen, M., I. Fernández-Val, and M. Weidner (2021). Nonlinear factor models for network and panel data. *Journal of Econometrics* 220(2), 296–324.
- Chernozhukov, V., I. Fernández-Val, J. Hahn, and W. Newey (2013). Average and quantile effects in nonseparable panel models. *Econometrica* 81(2), 535–580.
- De Silva, V. and L.-H. Lim (2008). Tensor rank and the ill-posedness of the best low-rank approximation problem. *SIAM Journal on Matrix Analysis and Applications* 30(3), 1084–1127.
- Elden, L. and B. Savas (2011). Perturbation theory and optimality conditions for the best multilinear rank approximation of a tensor. *SIAM journal on matrix analysis and applications* 32(4), 1422–1450.
- Evdokimov, K. (2010). Identification and estimation of a nonparametric panel data model with unobserved heterogeneity. *Department of Economics, Princeton University*.
- Fernández-Val, I., H. Freeman, and M. Weidner (2021). Low-rank approximations of nonseparable panel models. *The Econometrics Journal* 24(2), C40–C77.
- Freeman, H. and M. Weidner (2022). Linear panel regressions with two-way unobserved heterogeneity. *arXiv preprint arXiv:2109.11911*.
- Graf, S. and H. Luschgy (2002). Rates of convergence for the empirical quantization error. *The Annals of Probability* 30(2), 874–897.

- Hallin, M. and R. Liška (2007). Determining the number of factors in the general dynamic factor model. *Journal of the American Statistical Association* 102(478), 603–617.
- Hausman, J., G. Leonard, and J. D. Zona (1994). Competitive analysis with differentiated products. *Annales d'Economie et de Statistique*, 159–180.
- Hoderlein, S. and H. White (2012). Nonparametric identification in nonseparable panel data models with generalized fixed effects. *Journal of Econometrics* 168(2), 300–314.
- Kapetanios, G., L. Serlenga, and Y. Shin (2021). Estimation and inference for multi-dimensional heterogeneous panel datasets with hierarchical multi-factor error structure. *Journal of Econometrics* 220(2), 504–531.
- Kolda, T. G. and B. W. Bader (2009). Tensor decompositions and applications. *SIAM review* 51(3), 455–500.
- Kruskal, J. B. (1989). Rank, decomposition, and uniqueness for 3-way and n-way arrays. In *Multiway data analysis*, pp. 7–18.
- Latała, R. (2005). Some estimates of norms of random matrices. *Proceedings of the American Mathematical Society* 133(5), 1273–1282.
- Li, X., A. Wang, J. Lu, and Z. Tang (2019). Statistical performance of convex low-rank and sparse tensor recovery. *Pattern Recognition* 93, 193–203.
- Matyas, L. (2017). The econometrics of multi-dimensional panels. *Advanced studies in theoretical and applied econometrics*. Berlin: Springer.
- Menzel, K. (2021). Bootstrap with cluster-dependence in two or more dimensions. *Econometrica* 89(5), 2143–2188.
- Moon, H. R. and M. Weidner (2015). Linear regression for panel with unknown number of factors as interactive fixed effects. *Econometrica* 83(4), 1543–1579.
- Moon, H. R. and M. Weidner (2017). Dynamic linear panel regression models with interactive fixed effects. *Econometric Theory* 33(1), 158–195.
- Pesaran, M. H. (2006, 07). Estimation and inference in large heterogeneous panels with a multifactor error structure. *Econometrica* 74(4), 967–1012.
- Rabanser, S., O. Shchur, and S. Günnemann (2017). Introduction to tensor decompositions and their applications in machine learning. *arXiv preprint arXiv:1711.10781*.
- Richards, T. J. and B. J. Rickard (2021). Dynamic model of beer pricing and buyouts. *Agribusiness* 37(4), 685–712.
- Saleh, J. C. (2014). Simple estimators for cross price elasticity parameters with product differentiation: panel data methods and testing. *Revista de la Competencia y la Propiedad Intelectual* 10(18), 15–40.

- Tomioka, R., K. Hayashi, and H. Kashima (2010). Estimation of low-rank tensors via convex optimization. *arXiv preprint arXiv:1010.0789*.
- Tremblay, C. H. and V. J. Tremblay (1995). Advertising, price, and welfare: evidence from the us brewing industry. *Southern Economic Journal*, 367–381.
- Vannieuwenhoven, N., R. Vandebril, and K. Meerbergen (2012). A new truncation strategy for the higher-order singular value decomposition. *SIAM Journal on Scientific Computing* 34(2), A1027–A1052.
- Xu, A.-B. (2020). Tensor completion via a low-rank approximation pursuit. *arXiv preprint arXiv:2004.08872*.
- Zeleneev, A. (2020). Identification and estimation of network models with nonparametric unobserved heterogeneity. *Department of Economics, Princeton University*.

A Proofs

Proof of Proposition 2. In the following, let $\text{vec}(\tilde{\mathbf{X}})$ be the $\prod_n N_n \times K$ matrix of vectorised covariates after the within-cluster transformations where each column is a vectorised transformed covariate. The vec operator on other variables is the standard vectorisation operator. Also let $N = \prod_n N_n$ and the subscript $i = 1, \dots, N$ be the index for the vectorised data when i has no subscript. Then,

$$\begin{aligned}\beta_{GFE,C} &= \left(\text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\mathbf{X}}) \right)^{(-1)} \text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\mathbf{Y}}) \\ &= \left(\text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\mathbf{X}}) \right)^{(-1)} \text{vec}(\tilde{\mathbf{X}})' \left(\text{vec}(\tilde{\mathbf{X}}) \beta^0 + \text{vec}(\tilde{\mathcal{A}}) + \text{vec}(\tilde{\varepsilon}) \right) \\ &= \beta^0 + \left(\text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\mathbf{X}}) \right)^{(-1)} \text{vec}(\tilde{\mathbf{X}})' \left(\text{vec}(\tilde{\mathcal{A}}) + \text{vec}(\tilde{\varepsilon}) \right),\end{aligned}$$

such that,

$$\begin{aligned}\|\beta_{GFE,C} - \beta^0\| &= \left\| \left(\text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\mathbf{X}}) \right)^{(-1)} \text{vec}(\tilde{\mathbf{X}})' \left(\text{vec}(\tilde{\mathcal{A}}) + \text{vec}(\tilde{\varepsilon}) \right) \right\| \\ &\leq \|\kappa_N\| + \|\omega_N\|\end{aligned}$$

where

$$\begin{aligned}\|\kappa_N\| &:= \left\| \left(\text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\mathbf{X}}) \right)^{(-1)} \text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\mathcal{A}}) \right\|; \\ \|\omega_N\| &:= \left\| \left(\text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\mathbf{X}}) \right)^{(-1)} \text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\varepsilon}) \right\|.\end{aligned}$$

The terms κ_N and ω_N are dealt with separately.

First to bound κ_N . Notice,

$$\|\kappa_N\| \leq \left\| \left(\text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\mathbf{X}}) \right)^{(-1)} \right\|_F \left\| \text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\mathcal{A}}) \right\|$$

Focus on the right hand part, and let $\langle \cdot, \cdot \rangle_F$ be the Frobenius inner product,

$$\left\| \text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\mathcal{A}}) \right\| = \left\| \begin{bmatrix} \langle \tilde{X}_1, \tilde{\mathcal{A}} \rangle_F \\ \vdots \\ \langle \tilde{X}_K, \tilde{\mathcal{A}} \rangle_F \end{bmatrix} \right\| \leq \left\| \begin{bmatrix} \sum_{i=1}^N |\tilde{X}_{i,1} \tilde{\mathcal{A}}_i| \\ \vdots \\ \sum_{i=1}^N |\tilde{X}_{i,K} \tilde{\mathcal{A}}_i| \end{bmatrix} \right\| \quad (\text{A.1})$$

where the triangle inequality is used entry-wise. By Hölder's inequality

$$\sum_{i=1}^N |\tilde{X}_{i,k} \tilde{\mathcal{A}}_i| \leq \left\| \text{vec}(\tilde{\mathbf{X}}_k) \right\| \left\| \text{vec}(\tilde{\mathcal{A}}) \right\| \quad \text{for each } k = 1, \dots, K$$

This bounds the norm in (A.1) as,

$$\left\| \text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\mathcal{A}}) \right\| \leq \sqrt{\sum_{k=1}^K \left\| \text{vec}(\tilde{\mathbf{X}}_k) \right\|^2} \left\| \text{vec}(\tilde{\mathcal{A}}) \right\|.$$

From Assumption 1.(i) there is $\sqrt{\sum_{k=1}^K \left\| \text{vec}(\tilde{\mathbf{X}}_k) \right\|^2} = O_p(\sqrt{\prod_n N_n})$. This leaves $\left\| \text{vec}(\tilde{\mathcal{A}}) \right\|$.

Take $g_n(i_n)$ as the indices in i_n 's cluster such that $|g_n(i_n)|$ is the cluster size. Also, let $\bar{\varphi}_{i_n^*}^{(n)}$ be the cluster average for i_n 's cluster. Then,

$$\begin{aligned} \left\| \text{vec}(\tilde{\mathcal{A}}) \right\|^2 &= \left\| \tilde{\mathcal{A}} \right\|_F^2 = \sum_{i_1, \dots, i_d} \left(\sum_{\ell=1}^L \prod_{n=1}^d \left(\varphi_{i_n, \ell}^{(n)} - \bar{\varphi}_{i_n^*, \ell}^{(n)} \right) \right)^2 \\ (\text{Jensen's inequality}) \quad &\leq L^2 \sum_{i_1, \dots, i_d} \sum_{\ell=1}^L \frac{1}{L} \prod_{n=1}^d \left(\varphi_{i_n, \ell}^{(n)} - \bar{\varphi}_{i_n^*, \ell}^{(n)} \right)^2 \\ &\leq L \left(\prod_{n=1}^d N_n \right) \prod_{n=1}^d \frac{1}{N_n} \sum_{i_n} \left\| \varphi_{i_n}^{(n)} - \bar{\varphi}_{i_n^*}^{(n)} \right\|^2. \end{aligned}$$

Expand the term, $\left\| \varphi_{i_n}^{(n)} - \bar{\varphi}_{i_n^*}^{(n)} \right\|^2$,

$$\begin{aligned} \left\| \varphi_{i_n}^{(n)} - \bar{\varphi}_{i_n^*}^{(n)} \right\|^2 &= \left\| \varphi_{i_n}^{(n)} - \frac{1}{|g_n(i_n)|} \sum_{j_n \in g_n(i_n)} \varphi_{j_n}^{(n)} \right\|^2 \\ &\leq \frac{1}{|g_n(i_n)|^2} \left(\sum_{j_n \in g_n(i_n)} \left\| \varphi_{i_n}^{(n)} - \varphi_{j_n}^{(n)} \right\| \right)^2 \\ &\leq \max_{j_n \in g_n(i_n)} \left\| \varphi_{i_n}^{(n)} - \varphi_{j_n}^{(n)} \right\|^2 \quad (\text{A.2}) \end{aligned}$$

Then by Assumption 5,

$$\left\| \text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\mathcal{A}}) \right\| \leq \sqrt{L} \left(\prod_{n=1}^d N_n \right) O_p \left(\prod_{n \in \mathcal{M}} \sqrt{\xi_{N_n}} \right)$$

Lastly, Assumption 6.(i) implies the left hand term of $\|\kappa_N\|$ is $O_p(1/\prod_n N_n)$. This leaves

$$\|\kappa_N\| = \sqrt{L} O_p \left(\prod_{n^* \in \mathcal{M}} \sqrt{\xi_{N_{n^*}}} \right)$$

Finally, to bound $\|\omega_N\|$. Note that

$$\|\omega_N\| \leq \left\| \left(\text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\mathbf{X}}) \right)^{(-1)} \right\|_F \left\| \text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\boldsymbol{\varepsilon}}) \right\|.$$

Use Assumption 2 to bound the right hand term, $\left\| \text{vec}(\tilde{\mathbf{X}})' \text{vec}(\tilde{\boldsymbol{\varepsilon}}) \right\| = O_p(\sqrt{\prod_n N_n})$. Then, as above, the left hand term is $O_p(1/\prod_n N_n)$ such that

$$\|\omega_N\| = O_p \left(\frac{1}{\sqrt{\prod_n N_n}} \right).$$

■

Proof of Remark 5. Begin from A.2 in the proof of Proposition 2. The right hand terms, $\left\| \varphi_{i_n}^{(n)} - \varphi_{j_n}^{(n)} \right\|^2$, are bound as,

$$\left\| \varphi_{i_n}^{(n)} - \varphi_{j_n}^{(n)} \right\|^2 \leq c_n^2 \left\| \hat{\varphi}_{i_n}^{(n)} - \hat{\varphi}_{j_n}^{(n)} \right\|^2,$$

by Remark 5.(iii). The result then follows immediately by applying the conditions Remark 5.(i) and (ii) after this inequality. ■

Proof of Proposition 3. The interactive fixed-effect approximation error can be summarised as,

$$\left\| \text{vec}(\tilde{\mathcal{A}}) \right\|^2 = \sum_{i_1, \dots, i_d} \left(\sum_{\ell=1}^L \prod_{n=1}^d \left(\varphi_{i_n, \ell}^{(n)} - \bar{\varphi}_{i_n^*, \ell}^{(n)} \right) \right)^2.$$

By similar steps as the proof of Proposition 2 this can be bound by

$$\left\| \text{vec}(\tilde{\mathcal{A}}) \right\|^2 \leq L \left(\prod_{n=1}^d N_n \right) \prod_{n=1}^d \frac{1}{N_n} \sum_{i_n} \left\| \varphi_{i_n}^{(n)} - \bar{\varphi}_{i_n^*}^{(n)} \right\|^2.$$

Concentrate on the last term, $\frac{1}{N_n} \sum_{i_n} \left\| \varphi_{i_n}^{(n)} - \bar{\varphi}_{i_n^*}^{(n)} \right\|^2$,

$$\frac{1}{N_n} \sum_{i_n} \left\| \varphi_{i_n}^{(n)} - \bar{\varphi}_{i_n^*}^{(n)} \right\|^2 \leq \frac{1}{N_n} \sum_{i_n} \frac{\left(\sum_{j_n} k \left(\frac{1}{h_n} \left\| \hat{\varphi}_{i_n}^{(n)} - \hat{\varphi}_{j_n}^{(n)} \right\| \right) \left\| \varphi_{i_n}^{(n)} - \varphi_{j_n}^{(n)} \right\| \right)^2}{\left(\sum_{j_n} k \left(\frac{1}{h_n} \left\| \hat{\varphi}_{i_n}^{(n)} - \hat{\varphi}_{j_n}^{(n)} \right\| \right) \right)^2}, \quad (\text{A.3})$$

where elementary norm bounds are used to bound the left hand side.

Use as shorthand $\hat{a}_{ij}^{(n)} := \|\hat{\varphi}_{i_n}^{(n)} - \hat{\varphi}_{j_n}^{(n)}\|$. Expand $\|\varphi_{i_n}^{(n)} - \varphi_{j_n}^{(n)}\|$ around the proxies for each fixed-effect term and bound using the triangle inequality as,

$$\begin{aligned} \|\varphi_{i_n}^{(n)} - \varphi_{j_n}^{(n)}\| &\leq \|\varphi_{i_n}^{(n)} - \hat{\varphi}_{i_n}^{(n)}\| + \|\varphi_{j_n}^{(n)} - \hat{\varphi}_{j_n}^{(n)}\| + \|\hat{\varphi}_{i_n}^{(n)} - \hat{\varphi}_{j_n}^{(n)}\| \\ &= O_p(C_n^{-1}) + \hat{a}_{ij}^{(n)}. \end{aligned}$$

Then by the Cauchy-Schwarz inequality the term (A.3) can be bound as,

$$\begin{aligned} \frac{1}{N_n} \sum_{i_n} \|\varphi_{i_n}^{(n)} - \hat{\varphi}_{i_n}^{(n)}\|^2 &\leq \frac{1}{N_n} \sum_{i_n} \frac{\left(\sum_{j_n} k(\hat{a}_{ij}^{(n)}/h_n) \left(O_p(C_n^{-1}) + \hat{a}_{ij}^{(n)} \right) \right)^2}{\left(\sum_{j_n} k(\hat{a}_{ij}^{(n)}/h_n) \right)^2} \\ \text{(Assumption 8)} \quad &= \frac{1}{N_n} \sum_{i_n} \left(\sum_{j_n} k(\hat{a}_{ij}^{(n)}/h_n) \left(O_p(C_n^{-1}) + \hat{a}_{ij}^{(n)} \right) \right)^2 O_p(N_n^{-2}) \\ &= \frac{1}{N_n} \sum_{i_n} \left(O_p(C_n^{-1}) \sum_{j_n} k(\hat{a}_{ij}^{(n)}/h_n) + \sum_{j_n} k(\hat{a}_{ij}^{(n)}/h_n) \hat{a}_{ij}^{(n)} \right)^2 O_p(N_n^{-2}) \end{aligned} \quad (\text{A.4})$$

Where Assumption 8 implies that $\left(\sum_{j_n} k(\hat{a}_{ij}^{(n)}/h_n) \right)^2$ converges in probability at the rate N_n^2 to a strictly positive constant such that the inverse is also convergent by the continuous mapping theorem. Note that the discontinuity of the inverse function is at $\left(\sum_{j_n} k(\hat{a}_{ij}^{(n)}/h_n) \right)^2 / N_n^2 = 0$, which has zero probability by Assumption 8. This shows, $\left(\sum_{j_n} k(\hat{a}_{ij}^{(n)}/h_n) \right)^{-2} = O_p(N_n^{-2})$.

Since the class of kernels are bounded, $\sum_{j_n} k(\hat{a}_{ij}^{(n)}/h_n) = O(N_n)$ and from Assumption 7 there is $\sum_{j_n} k(\hat{a}_{ij}^{(n)}/h_n) \hat{a}_{ij}^{(n)} = O(h_n^\alpha)$. Thus, (A.4) continues,

$$\begin{aligned} \frac{1}{N_n} \sum_{i_n} \|\varphi_{i_n}^{(n)} - \hat{\varphi}_{i_n}^{(n)}\|^2 &\leq \frac{1}{N_n} \sum_{i_n} (O_p(C_n^{-1} N_n) + O(h_n^\alpha N_n))^2 O_p(N_n^{-2}) \\ &= O_p(C_n^{-2}) + O_p(C_n^{-1} h_n^\alpha) + O_p(h_n^{2\alpha}) \end{aligned} \quad (\text{A.5})$$

Summarising leaves

$$\frac{1}{N_n} \sum_{i_n} \|\varphi_{i_n}^{(n)} - \hat{\varphi}_{i_n}^{(n)}\|^2 \leq O_p(C_n^{-2}) + O_p(C_n^{-1} h_n^\alpha) + O_p(h_n^{2\alpha}).$$

By similar arguments in the proof of Proposition 2 this makes up the first component of Proposition 3. The second component comes directly from extending Assumption 6 to the kernel estimator, as per the proof of Proposition 2. \blacksquare

B Demand application: fixed-effects with external variables

This section provides an explanation of how the external variable approach can work in conjunction with the GFE and kernel fixed-effects approach and briefly how the instrumental variable approach is limited in the presence of fixed-effects.

Take the interactive fixed-effects model and consider the usual IV approach. For simplicity, take the two equation model for IV estimation as,

$$\begin{aligned} Y &= X\beta + \mathcal{A} + \mathcal{B} + \varepsilon \\ X &= Z\delta + \mathcal{A} + \mathcal{B} + \nu \end{aligned} \tag{A.6}$$

where, as before, $\mathcal{A} = \sum_{\ell=1}^L \varphi_{\ell}^{(1)} \circ \varphi_{\ell}^{(2)} \circ \varphi_{\ell}^{(3)}$ and \mathcal{B} is a collection of scalar fixed-effects. If Z varies over only one dimensions, like the price of barley, then it is not possible to project these fixed-effects out. This is because once X is projected onto Z in the first stage, it then too only varies over t and reduces the second stage problem to one dimension. Hence, Z must be mean independent of $\mathcal{A} + \mathcal{B}$ to maintain the exogeneity condition. In the beer demand application this insists that Z_t must be independent of any time related shocks that impact both demand and supply for beer; that is, the usual IV exogeneity condition with respect to the time fixed-effect in the interactive term. For mean independence with \mathcal{B} it is required that Z_t is independent of any scalar fixed-effect that varies over t . For independence with \mathcal{A} , it suffices that Z_t is independent of $\varphi_{t\ell}^{(3)}$ for all ℓ .

The external variable model is implemented using the external variable as a kernel proxy to difference out a weighted mean in the t dimension. The advantage here is that this method does not project the covariates onto the external variable space like with instrumental variables, hence does not reduce the covariates dimensionality. This means it can be used in conjunction with standard fixed-effects methods or GFE and kernel methods. The exogeneity condition here, though, is different to the standard instrumental variable condition.

Required in the external variable setting is some form of continuous mapping from the fixed-effect parameter in the t dimension and the external variable. For example, there could be a Lipschitz continuous mapping as,

$$\|\varphi^{(3)}(z) - \varphi^{(3)}(z_0)\| \leq C\|z - z_0\|.$$

Then, as long as Z_t is sampled from a distribution with strictly positive probability over its support, closeness in the external variable can be found, which in turns bounds variation in the fixed-effect parameter. Hence, with some form of contraction from Z_t to $\varphi_t^{(3)}$, \mathcal{A} can be projected out. Because other fixed-effects can be included, \mathcal{B} can be projected out with simple additive fixed-effects and the other terms in \mathcal{A} can also be differenced out using the GFE or

kernel methods. This leaves the remainder error, from GFE or kernel,

$$\sum_{\ell=1}^L \left(\varphi_{i,\ell}^{(1)} - \bar{\varphi}_{i^*,\ell}^{(1)} \right) \left(\varphi_{j,\ell}^{(2)} - \bar{\varphi}_{j^*,\ell}^{(2)} \right) \left(\varphi_{\ell}^{(3)}(Z_t) - \bar{\varphi}_{\ell}^{(3)}(Z_{t^*}) \right).$$

Thus, the error reduction from the first two dimensions, over i and j , is maintained and any endogeneity related to the third dimension can be dealt with by finding closeness in the external variable, much like is done with the usual proxies considered throughout.

C Reducing the number of estimated parameters

Analysts may be concerned with the number of parameters implied by the least squares problem (12). In practice, this equation implies a total of $N_1 N_2 g(N) + N_1 g(N) N_3 + g(N) N_2 N_3$ parameters, where $g(N)$ is the number of groups in each dimension that may depend on total data size $N = \prod_n N_n$. This implies the number of fixed-effects parameters with respect to total data size is

$$\frac{g(N) \sum_{n=1}^d \prod_{n' \neq n} N_{n'}}{\prod_n N_n} = g(N) \cdot O\left(\frac{1}{\min_{n \in \{1, \dots, d\}} N_n}\right) \quad (\text{A.7})$$

Hence, in the linear setting, the loss of degrees of freedom is negligible as long as the group size $g(N)$ does not grow too fast with respect to total data size. However, this makes estimation in non-linear settings like (5) problematic because of the incidental parameter bias, see in Chen, Fernández-Val and Weidner (2021). For this reason it is useful to consider versions of the within-cluster transformation that do not require so many parameters. The following is a non-exhaustive list of methods to reduce the number of estimated parameters.

The first approach to consider is to simply ensure the group sizes are small with respect to data size. To do this consider $g_n := g(N_n)$ as the number of groups in dimension n . Take a similar calculation to (A.7) to obtain the total number of parameters $\sum_n g_n \prod_{n' \neq n} N_{n'}$. It should then be clear that as long as $g_n = o(N_n)$, the number of estimated parameters is small with respect to total data size and the incidental parameter problem is asymptotically negligible. However, the condition that $g_n = o(N_n)$ may be highly restrictive. For example, if a sample of the unobserved parameter space is very disparate then this condition restricts the analyst to make poor approximations of the fixed-effects terms as each $\varphi_{i_n^*,\ell}^{(n^*)} - \varphi_{j_n^*(i_n^*),\ell}^{(n^*)}$ will be very large. This is why it is important to consider the alternatives provided below. As can be seen in (13), the approximation error is multiplicative across dimensions, which means the analyst needs only to approximate a subset of these well. This fact is utilised in the below displays.

Consider clusters along just one dimension. The within-cluster transformation associated

with this is simply,

$$\tilde{A}_{ijt} = A_{ijt} - A_{i^*jt} = \sum_{\ell=1}^L (\varphi_{i\ell}^{(1)} - \bar{\varphi}_{i^*\ell}^{(1)}) \varphi_{j\ell}^{(2)} \varphi_{t\ell}^{(3)}.$$

Under some high-level assumptions on the unobserved fixed-effects, $\bar{\varphi}_{i^*\ell}^{(1)} = \varphi_{i\ell}^{(1)} + O\left(\frac{1}{N_1}\right)$. Also, the term $\bar{\varphi}_{i^*\ell}^{(1)}$ may have to be estimated - call the estimate $\hat{\varphi}_{i^*\ell}^{(1)}$. Again, under some high-level assumptions, this could be estimated as $\hat{\varphi}_{i^*\ell}^{(1)} = \bar{\varphi}_{i^*\ell}^{(1)} + O_p\left(\frac{1}{\sqrt{N_2 N_3}}\right)$. Combining this leaves the estimated $\tilde{A}_{ijt} = O_p\left(\frac{1}{\min\{N_1, \sqrt{N_2 N_3}\}}\right)$. So selection of which dimension, d^* , to cluster and difference over solves the optimisation $d^* = \operatorname{argmax}_{d \in \{1,2,3\}} \min\{N_d, \prod_{n,m \neq d; n \neq m} \sqrt{N_n N_m}\}$. This procedure requires $N_{n \neq d^*} N_{m \notin \{d^*, n\}} \times g(d^*)$ parameters to estimate, where $g(d^*)$ is the number of groups for dimension d^* . Of course, choice of d^* may also incorporate the number of parameters required for estimation. Note that this method does not automatically project the additive terms from \mathcal{B} , so this should be performed after an initial within projection.

This logic can be extended to a difference across two-dimensions as,

$$\tilde{\tilde{A}}_{ijt} = A_{ijt} - A_{i^*j^*t} = \sum_{\ell=1}^L \left(\varphi_{i\ell}^{(1)} \varphi_{j\ell}^{(2)} - \bar{\varphi}_{i^*\ell}^{(1)} \bar{\varphi}_{j^*\ell}^{(2)} \right) \varphi_{t\ell}^{(3)}.$$

By the same reasoning as above this leads to

$$\hat{A}_{ijt} = O_p \left(\left(\min_{d \in \{1,2\}} \min \left\{ N_d, \sqrt{N_{\{1,2\} \setminus d} N_3} \right\} \right)^{-1} \right)$$

The optimal dimensions to cluster and difference on is

$$\{d_1^*, d_2^*\} = \operatorname{argmax}_{d_1, d_2} \min_{d \in \{d_1, d_2\}} \min \left\{ N_d, \sqrt{N_{\{d_1, d_2\} \setminus d} N_{n \notin \{d_1, d_2\}}} \right\}.$$

This requires $g(d_1^*)g(d_2^*) \times N_{n \notin \{d_1^*, d_2^*\}}$ parameters.

Take a further difference to obtain

$$\tilde{\tilde{\tilde{A}}}_{ijt} = (A_{ijt} - A_{i^*j^*t}) - (A_{ijt^*} - A_{i^*j^*t^*}) = \sum_{\ell=1}^L (\varphi_{i\ell}^{(1)} \varphi_{j\ell}^{(2)} - \bar{\varphi}_{i^*\ell}^{(1)} \bar{\varphi}_{j^*\ell}^{(2)}) (\varphi_{t\ell}^{(3)} - \bar{\varphi}_{t^*\ell}^{(3)}).$$

This reduces to

$$\tilde{\tilde{\tilde{A}}}_{ijt} = O_p \left(\left(\min_{d \in \{1,2\}} \min \left\{ N_d, \sqrt{N_{\{1,2\} \setminus d} N_3} \right\} \min \left\{ N_3, \sqrt{N_1 N_2} \right\} \right)^{-1} \right),$$

which is smaller than the two cluster difference. d^* can be found similarly. This requires $g(d_1^*)g(d_2^*) \times N_{n \notin \{d_1^*, d_2^*\}} + N_{n \notin \{d_1^*, d_2^*\}} \min_{m \in \{d_1^*, d_2^*\}} N_m g(d_{n \neq m}^*)$.

The above parameter reduction exercises and specifically the choice of which dimension(s) to cluster on are also subject to the proxies used for clustering. For example, along some dimensions

there may exist observable characteristics that provide a good signal of individual unobserved fixed-effect cluster. Diagnostics discussed in Section 3.1 also uncover which dimension exhibits low-rank variation, making it a good candidate for single dimension clustering. The d^* 's above are given as guides in applications where there is no obvious dimension to concentrate on when parameter reduction is required. It should also be clear that more estimated parameters can lead to tighter asymptotic rates of decay in the unobserved remainder term, which becomes obvious in the asymptotic results discussed later. One last consideration when choosing from these reduced parameter options is the implication on the additive fixed-effects terms, where not all additive terms are automatically projected with each of these reduction methods.