

Working Paper No. 13-09

## Orthogonalization of Categorical Data:

# Orthogonalization of Categorical Data: How to Fix a Measurement Problem in Statistical Distance Metrics

Ross Shalpinberg<sup>1</sup>

November 10, 2013

Policy makers depend on economists, statisticians, and other social scientists to make accurate observations and draw solid conclusions from quantitative analysis. Econometrics, for example, has come a long way in the 20th century and guides many decisions made today. On the other hand, some statistical procedures have not had significant development, but instead lied and their original assumptions are forgotten. The appropriateness of many of these measurements has come into question, and while criticism is often leveled, little is done to correct them. In reality, there is a proliferation of measurement problems being committed everyday. This problem involves the use of statistical distance metrics to measure social phenomena. For example, measurements which would routinely be used to answer questions like: by how much have the imports of the United States changed in the past year? By how much has racial diversity changed in the past decade? Does greater ethno-linguistic diversity lead to civil conflict? These and similar questions rely on accurate multi-variate distance metrics. However all dis-

*“Many multivariate statistical methods can be regarded as techniques for investigating a sample space in which each sample member is represented by a point.”* John C. Gower (1977), pg 13.

*“Measurement is a big part of mobilizing for impact. You set goals, and then you use data to make sure you’re making progress toward it. This is crucial in business—and it’s just as important in the fight against poverty and disease”* Bill H. Gates (2013), pg 52.

## 1 Introduction: The Problem

Before introducing any formal mathematics, consider the five following measurement puzzles:

**Puzzle 1:** Consider a two-country world where Country  $C$  exports half corn and half corn meal. Country  $D$  exports half corn and half computers. Which one has the most diverse exports? Measures of export diversification indicate that both countries are exactly equally diverse.

**Puzzle 2:** In City  $A$  exactly 5 percent of the labor force are Economics Professors. In City  $B$ , exactly 5 percent of the labor force are Research Economists. The Location Quotient doesn’t recognize cross-discipline similarities, so between the two cities, City  $A$  is classified as being relatively sparse in Research Economists and City  $B$  is classified as being relatively dense in Research Economists.

**Puzzle 3:**





still maintaining the original structure of the data. That procedure is the subject of this paper and is detailed in the Methodology section.

How big is this problem and the corresponding bias? That depends on the data, but a rough estimate is given by finding the average value of the data's similarity matrix, where  $s_{ij}$  is the  $(i,j)$  of a similarity matrix:

$$bias = \frac{\sum_{i=1}^n \sum_{j=1}^n s_{ij}}{n^2} \quad (1)$$

Using 4-digit SITC international trade data for the year 2000, this number is  $\frac{481801.6}{772^2} = 0.808$ . In other words, the average export product  $x_i$  is, on average 80.8 percent like product  $x_j$ . However, all current distance metrics, and hence all standard trade metrics, implicitly assume that similarity is zero between all categories. This is clearly not true, and without zero similarity between categories, the standard multi-dimensional metrics are not valid.

The question naturally arises: how widespread is the problem? Well it exists in every branch of

procedure and so may be unable to use this procedure until a similarity-calculating procedure is found. This could be an ideal subject for future research.

To preview the proposed orthogonalization procedure, one can see it as a change of coordinate systems. I take as given a set of data vectors and the measure of similarity between every pair of its dimensions. The basic idea is that the similarity between dimensions can increase which reduces distance between individual dimensions in a vector. This is best seen in the spherical coordinate system (see Spiegel 1959, Munkres 1991). The measure of the angle from the vector to an axis is given by  $\cos^{-1}$ . The orthogonalization procedure then uses the change of coordinates to find the length of this vector along each axis. The rectangular coordinate system is what most empirical measures are based upon, at least those with concepts like angle and distance. So in order for a quantitative measure to be valid, we must change the coordinate system to what the measure is assuming. This is the basic idea of the paper.

## 2 Literature

The aforementioned problem of heterogeneity between dimensions is, as far this author can tell, completely unacknowledged when working with shares data. That said, the problem is recognized when working with quantitative variables which are not in the form of shares, and has been the focus of substantial research. I can identify 9 distinct orthogonalization procedures each of which are based on two basic methods, of which there are undoubtedly more. The first method, found overwhelming in Statistics and Econometrics, involves the use of a correlation or covariance matrix to find orthogonal dimensions. The second method, found in Mathematics and applied in Computer Science and Physics, involves knowing exactly how the system behaves in a non-stochastic fashion and having perfect measurements.

Two ideas distinguish my problem and solution from the rest of the literature. The first is that the data which I am examining always exists on a unit simplex. Thus the range of possible values that variables may take is relatively limited, and no negative values are allowed. This eliminates the use of correlation and covariance matrices since these procedures commonly produce negative

values. Secondly, my over-arching argument rests on the idea that the true coordinates of these observations are not known, but detailed information exists in the form of similarity values which can be used to find the true location.



idea behind principal coordinates is that it takes a distance measure between all pairs of *observations* and gives them coordinates; according to Gower (1977, pg.19), "We can ask how the coordinates of points with the given distances be found." On the other hand, the orthogonalization method that I discuss in this paper adjusts for a similarity measure between *variables*. The idea of this adjustment between variables is that, after adjustment, similarity between observations (or other measures) can be measured, based on the adjusted variables. Consider that principal coordinates analysis takes a matrix of similarities between observations as given, and then adjusts the variables to fit those similarities. In contrast, the orthogonalization procedure described herein is quite the opposite in that it seeks to create a similarity matrix between observations based on the given similarity between the variables.

Third, and also similar to principal components analysis, is factor analysis. In factor analysis a researcher attempts to identify unobserved, latent categorical variables. In this case the covariance between dimensions leads to recognition of a previously unidentified latent variable. So this statis-

exports.

Fifth, regression analysis and analysis of variance, are chiefly concerned with accounting for covariance. Each variable is treated as a dimension, and the covariance between dimensions can greatly affect the estimate of the mean value of a regressor on the regressand. In an abstract sense this is similar because the practitioner realizes that variables are not completely independent of one another, and so the covariance, or angle between dimensions, is included in the process by design.

Not including important covariates leads to omitted variable bias: the magnitude of a parameter is inaccurate. This is geometrically equivalent to a parameter value being projected onto an  $n$ -plane but not parallel to its coordinate axis, with the angle between its proposed axis and the actual projection proportional to the correlation with the omitted variable. This is exactly the argument that I am making for distance measures.

## 2.2 Exact Methods for Orthogonalization

The second class of orthogonalization procedures is based on mathematical procedures for rotation, have no stochastic assumption, and the underlying data generating process has no latent variables. The ever-present implicit assumption that I am trying to upend here is that  $n$ -space coordinates are all a priori known. For this reason the Gram-Schmidt process, the Householder Transformation, and the Givens Rotation can all be ruled out as potential orthogonalization techniques because they all make this assumption. I am not going to detail each method, because none of them can work due to this assumption. Again, each assumes that the coordinates of a vector are known, whereas I only assume partial information about the coordinates is known.

## 2. Other Literature

Sixth, this paper has ties to Measure Theory. The main point of measure theory concerns distinguishing a measurement of an attribute from the attribute itself. Consider common commodities like heat, corn, and computers. How different are these things? As economists, we don't particularly care about heat, corn, or computers in themselves, but rather about the implied underlying pro-

able to measure is not necessarily the same as what we need to measure to form general theoretical statements.

Seventh, this paper closely relates to Index Number Theory, however, this paper has nothing direct to say about prices. In Index Number Theory, one can typically identify two distinct approaches: the Axiomatic Approach versus the Economic Approach. What is the point to having these two different approaches? The point is that the data does not line up exactly with theory be-

### 3 Methodology

The problem, stated in yet another way, is that the categories in which much data is classified is ad hoc, with some categories more alike than others. To fix this problem, one first needs a measure of similarity between all dimensions, to which I will defer to other papers. For example in International Trade see Hidalgo, et al (2007) or for a more general treatment see Dauxois and Kuentz (2002). Second, according to Gentle (2007), this similarity data is best viewed as representing the angle between dimensions. With this in mind, the orthogonalization procedure is then to take each data share  $x_{c,i}$  and project it onto an orthogonal coordinate system, Euclidean  $n$ -space. Then one can apply a number of distance metrics. This projection is best viewed as a change from spherical<sup>3</sup> to rectangular coordinates for each individual dimension.

#### 3.1 Similarity Matrices and Angle Between Dimensions

of vectors: “The cosine of the angle between two vectors is related to the correlation between the vectors, so a matrix of the cosine of the angle between the columns of a given matrix could also be

between the z-axis and the x-y plane, in radians. Let  $\theta$  be the angle between the x-axis and the z-y plane. Then given the values for spherical coordinates  $(r; \phi; \theta)$ , the corresponding rectangular coordinates  $(x_1; y_1; z_1)$  can be found by :

$$x = r \sin \phi \cos \theta \tag{3}$$

$$y = r \sin \phi \sin \theta \tag{4}$$

$$z = r \cos \phi \tag{5}$$

The above equations are a projection of a vector in spherical coordinates into the rectangular coordinate system. These should be familiar to the reader and are typically first encountered in multivariate Calculus.

### • The Orthogonalization Procedure: Change of Coordinates

The most promising method to obtain an orthogonal coordinate system is to use a change from hyperspherical to rectangular coordinates. I use the algorithm described in Lin (1995)<sup>4</sup>. The basic idea here is to treat each dimension of a vector as its own vector. Then, because the angle of each dimension is known in regards to every other dimension, and using a trigonometric-based algorithm, one can project the length of the vector onto each dimension, repeat for each entry in the vector, and sum them up at the end.

Define a vector of shares data by  $x$  which has  $n$  rows indexed by  $i$ . The associated  $n$  by  $n$  similarity matrix is  $S$ , with elements  $s_{ij}$  here rows are indexed by  $i$  and columns indexed by  $j$ . Redefine  $s_{ij}$  in terms of degrees and convert it from a similarity matrix to a distance matrix:

$$d_{ij} = (1 - s_{ij})90 \tag{6}$$

For every  $i$ , define each entry in the vector  $x_i$  as a radius. Define each column entry  $j$  in row  $i$  of matrix  $S$  as the angle formed by the vector  $i$  to dimension  $j$ . To convert to rectangular coordinates,

---

<sup>4</sup>I thank Professor Jeanne Huffot for providing me with an equivalent algorithm.

align the numeraire good as the first good. This represents a simple rotation of the coordinate system<sup>5</sup>.

$$x_{1;j} = x_1 \cos(\theta_{1;1}) \quad (7)$$

So, similarly, find the projection of the second good onto each axis. Do this for each of the  $n$  goods using the following algorithm.

$$\begin{aligned} x_{2;j} &= x_1 \cos \theta_{1;j} \\ x_{3;j} &= x_1 \sin \theta_{1;j} \cos \theta_{2;j} \\ x_{4;j} &= x_1 \sin \theta_{1;j} \sin \theta_{2;j} \cos \theta_{3;j} \end{aligned} \quad (8)$$

$$\begin{aligned} x_{n-2;j} &= r \sin \theta_{1;j} \sin \theta_{2;j} \sin \theta_{3;j} \dots \sin \theta_{n-3;j} \cos \theta_{n-2;j} \\ x_{n-1;j} &= r \sin \theta_{1;j} \sin \theta_{2;j} \sin \theta_{3;j} \dots \sin \theta_{n-2;j} \cos \theta_{n-1;j} \\ x_{n;j} &= r \sin \theta_{1;j} \sin \theta_{2;j} \sin \theta_{3;j} \dots \sin \theta_{n-1;j} \sin \theta_{n;j} \end{aligned}$$

The previous algorithm is adapted from Lin (1995).<sup>6</sup> Note that the pattern of the hyperspherical algorithm is such that the vast majority of terms are sine and each line ends with cosine except for the very last line which ends in sine.

Repeat the above procedure for all  $i$  and then define for all  $j$ :

$$x_j =$$

And finally, because this is shares data and exists on a unit simplex, the sum of the entries must add to 1. Define total unadjusted shares (TUS) as:

$$TUS = \sum_{j=1}^n x_j \quad (11)$$

And normalize each entry using  $TUS$ :

$$x_j^a = \frac{x_j}{TUS} \quad (12)$$

The above equations outline the orthogonalization procedure for a single data vector. Likely a researcher would be comparing many different data vectors and would need to complete this procedure for each vector. This is the end of the orthogonalization procedure.

#### 4 Distance Metrics

The following is an introduction to a subset of common distance metrics used in many different statistical and social science fields. Many more distance metrics exist, and as with the literature review, this list is by no means exhaustive. In various fields these distance metrics go by specific



Where only the positive root is used. When  $p = 1$ , the distance metric is known as either Manhattan, or City-Block Distance:

$$D_{Manhattan} = \sum_{i=1}^n |x_{c;i} - x_{d;i}| \quad (14)$$

City-block distance gets its name from the fact that to get from one point to another in a city grid one must follow the streets. Particularly in Manhattan, streets intersect at right angles, so the absolute value in the difference Manhattan (Manhattan street) - 384 (dimension 0 (Manhattan) - Manhattan) - 384 (total) - 384 (area) - 3

$$D_{Euclidean} = \sqrt{\sum_{i=1}^n (x_{c;i} - x_{d;i})^2} \quad (15)$$

When  $p$  approaches infinity, one Chebyshev's Distance:

$$D_{Chebyshev} = \max_{i=1}^n |x_{c;i} - x_{d;i}| \quad (16)$$

The second class of distance estimators scales the coordinate values. Canberra Distance typifies this class:

$$D_{Canberra} = \sum_{i=1}^n \frac{|x_{c;i} - x_{d;i}|}{|x_{c;i}| + |x_{d;i}|} \quad (17)$$

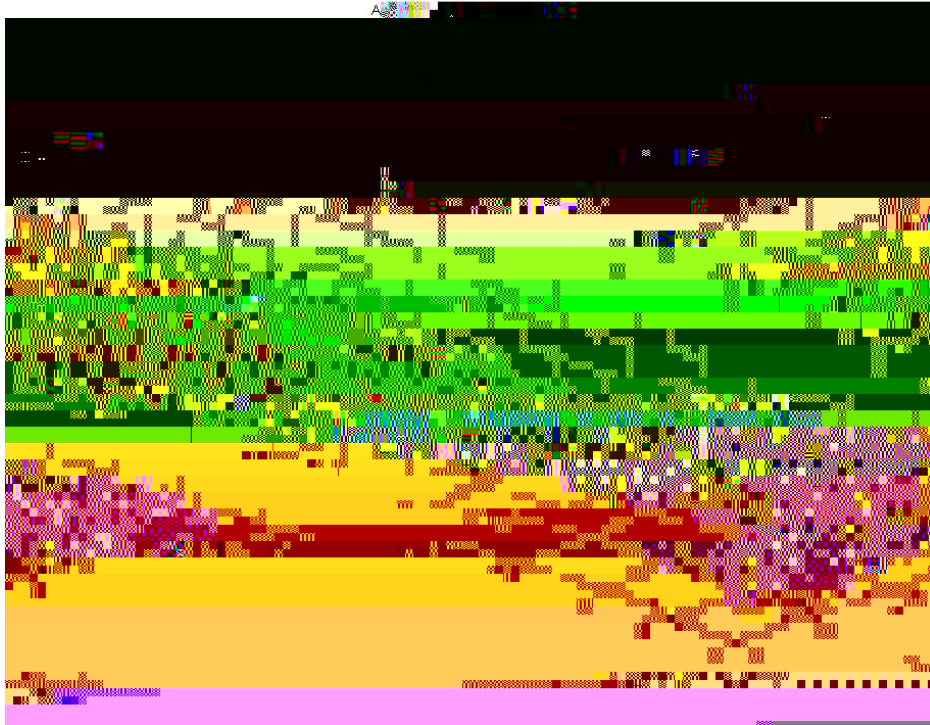
The third class of distance includes the Czekanowski Coefficient, which goes beyond a myriad

## 4 Simulation

My above qualitative argument for the need for an orthogonalization procedure is hopefully persuasive. However, I find it useful to present a very general example using a series of simple simulations. I will consider a three-dimensional world where a single observation  $x_i$  is composed of  $k$  shared attributes, where the sum of  $k$  attributes is one. Using a random number generator, I will assign values to the  $k$  attributes as well as to the  $\frac{k^2}{2}$   $k$  similarities between attributes. This is equivalent to finding a random point in a random  $k$ -space. I will then calculate the Euclidean distance to the origin first ignoring the similarities, and then compare this to the Euclidean distance using the orthogonalization procedure. I repeat this for varying values of  $n$  and  $k$ , with the results displayed in Figures 2 and 1. Here the number of observations are  $n = 1; 2; \dots; 120^7$  which are plotted along the x-axis, and the number of dimensions is  $k = 2; 3; \dots; 160,^8$  plotted along the y-axis. The z-axis (vertical) represents the measured distance on  $k$  dimensions between a point and the origin, averaged for  $n$  observations. Figure 2 ignores the similarities between dimensions and computes Euclidean distance in the normal way. Compare these average values to those in Figure 1 which do take into account the similarities between dimensions and thus compute the true average distances. Figure 3 plots the simple difference between the two surfaces.

By definition, the distances using the Law of Cosines are correct, it is the Euclidean distances,

Figure 1: Actual Distance using the Law of Cosines



pletely invariant to the number of observations, and are completely dependent on the number of dimensions used. So these measures depend more on the number of dimensions used rather than the actual values in the shares data.

## Applications

The following section outlines a few examples in the literature where the orthogonalization method can potentially yield great benefit. I plan to academically pursue these topics in the near future. I have drafted or am working on proposals on all of the following topics.

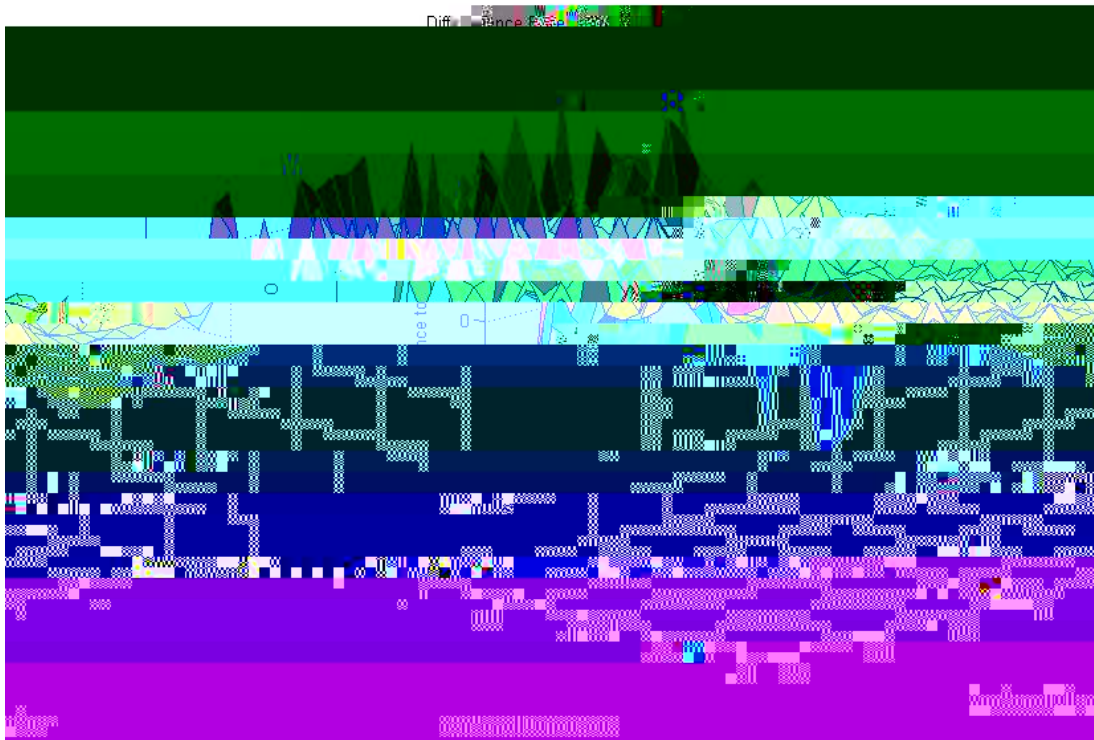
### 5.1 Application: Price Indices

The computation of index numbers suffers from three primary challenges. The first is that the data is in the form of categories, which naturally do not obey the laws of arithmetic. The second is that the weights of categories change over time. These first two challenges are commonly referred to as the “Index Number Problem.” The third is that classification and categorical ambiguity

Figure 2: Estimate using Euclidean Distance



Figure 3: Difference Between Actual and Estimated



egorizing people as red, yellow, brown, black, and white (Funderburg 2013, 83). The simplicity of this categorization has, quite understandably, been contentiously opposed. The offense is likely not so much in the names used, as it is in the broadness of each category. When being given a label, most individuals would likely want to be recognized as closely as possible to the category in which they self-identify. To this end, a researcher may feel compelled to divide the categories into smaller subcategories. The only problem is that the index monotonically increases with the number of categories. While it's possible that this measure is correct at an aggregation level, the point is that it's not clear which level of aggregation is appropriate. In particular, think about multi-racial people. In computing the IQV, most researchers treat a bi- or multi-racial person as being in a completely different category. However, a multi-racial person is really, by definition, a combination

$$IQV = \frac{1}{n-1} \sum_{k=1}^P (p_k)^2 \quad (22)$$

Where  $p_k$  is the share of group  $k$  in the total population. This is a normalized version of Euclidean distance from the origin. So with cross-category observations, the bi- or multi-racial observations can be partially grouped into categories, changing the computation in a way that is not immediately clear.

## 5. Application: Development and Political Institutions: The Index of Ethno-linguistic Fractionalization

Incredibly similar to the IQV is a measure known as the Index of Ethno-linguistic Fractionalization (ELF), which is applied extensively in the literature in the fields of political science and economic development. The equation is given by :

$$ELF = 1 - \sum_{k=1}^K p_k^2 \quad (23)$$

Where  $K \geq 2$  and  $p_k^2$  is the share of ethnic group  $k$  in the total population. This is a version of non-normalized Euclidean distance from the origin.

The basic idea behind the Index is, just like the IQV, to measure diversity. The problem, as clearly defined by Laitin and Posner (2001) is two-fold. First, a researcher needs to be careful about the level of aggregation used in defining ethnic groups. Second, not all ethnic groups are equally unlike. To this end, Bossert, D'Ambrosio and Ferrara (2005) define a Herfindahl Index that coincides with the Generalized Index of Ethno-linguistic Fractionalization Index which accounts for similarities between categories. Indeed they almost define the Law of Cosines distance metric in Knippenberg (2013), but stop short of taking a geometric interpretation of distance metrics. So I am pleased that this problem of non-zero similarity between categories has been recognized before in this literature, but has not had an orthogonalization procedure applied.



Bill Gates (2013, pg 52).



the same notation as above, denote total exports of country  $c$  as  $X_c$ , where  $X_c = \sum_{i=1}^n X_{c;i}$ . Define  $x_{c;i}$  to be the share of good  $i$  in total exports of country  $c$ , where  $x_{c;i} = \frac{X_{c;i}}{X_c}$ , and, consequently,  $\sum_{i=1}^n x_{c;i} = 1$ . Equivalently for a second, but with the subscript  $d$ , and for the world, with the subscript  $w$ . All measurements except the last two are assumed to be taken in the same time period, so time subscripts are otherwise suppressed.

The Hirschman-Herfindahl Index:

$$HHI_{c;d} = \sum_{i=1}^n \frac{1}{X_c^2} x_{c;i}^2 \quad (24)$$

The export similarity Index, Finger and Kreinin (1979):

$$FK_{c;d} = \sum_{i=1}^n \min(x_{c;i}, x_{d;i}) \quad (25)$$

The Grubel-Lloyd Index, Grubel and Lloyd (1971):

$$GL_{c;d} = 1 - \frac{\sum_{i=1}^n |x_{c;i} - x_{d;i}|}{\sum_{i=1}^n (x_{c;i} + x_{d;i})} \quad (26)$$

Two definitions are common for the Export Diversification Index. The first follows directly from Finger and Kreinin (1979):

$$DX1_c = \sum_{i=1}^n \min(x_{c;i}, x_{w;i}) \quad (27)$$

Where the subscript  $w$  stands for “world”. The more common definition is exactly the same as the Hirschman Index:

$$DX2_c = \sum_{i=1}^n \frac{1}{X_c^2} x_{c;i}^2 \quad (28)$$

The Trade Compatibility Index, Michael (1993)<sup>11</sup>:

The Export Specialization Index:

$$ES_c = \frac{X_{c;i}}{m_{d;i}} \quad (30)$$

Changes in Global Demand for Major Exports:

$$CGD_c = \sum_{i=1}^n S_{i,0} (X_{i;t} - X_{i,0}) \quad (31)$$

Changes in Global Market Share for Major Exports:

$$CGMS_c = (S_{i;t} - S_{i,0}) M_{g,t} \quad (32)$$

And lastly, the Thiel Index of export concentration:

$$T_c = \frac{1}{n} \sum_{i=1}^n \frac{X_i}{\sum_{i=1}^n X_i} \ln \frac{X_i}{\sum_{i=1}^n X_i} \quad (33)$$

As the reader can see, each trade statistic treats each export (or import) product as a separate dimension, and there is no system of weights or compensation for dimensions being more or less alike.

These trade statistics can be classified in several ways. The Hirschman Index is of the absolute type: the others describe a country's export shares as some distance from the origin. All of the others are of the relative type. The export diversification (Finger and Kreinin) tells the Manhattan distance between a country's export shares and the world export shares. The rest give the distance between two country's export shares: the Grubel-Lloyd gives the distance exactly in terms of Canberra distance, and the rest of the trade statistics are of the relative type: they tell the distance between two non-origin points. The export similarity and export diversification measures (both based on the work of Finger and Kreinin) are nearly identical to the Czekanowski Coefficient, except that they are already in terms of shares, whereas the Czekanowski Coefficient converts to shares after summing the values.

### 7. Simple Example: x International Trade

Consider a two-country world with two goods: guns and butter. The first country, denoted by  $c$ , produces 20 percent butter and 80 percent guns, while the second country, denoted by  $d$ , produces

70 percent butter and 30 percent guns:

$$y_c = \begin{matrix} y_{c;b} \\ y_{c:g} \end{matrix} = \begin{matrix} 0.2 \\ 0.8 \end{matrix} ; y_d = \begin{matrix} y_{d;b} \\ y_{d:g} \end{matrix} = \begin{matrix} 0.7 \\ 0.3 \end{matrix} ;$$

Then suppose the production of guns and butter share some common attributes. For example, both need land: butter producers more so to raise dair cows and but guns producers also need land for placing factories. Both also need metal: butter producers need metal for producing churns and vats, and but gun producers need metal relative more to produce stocks and barrels. By some external measurement process we know the the similarity between guns and butter to be 0.8, or 80 percent of the inputs are alike. Then the similarity matrix, denoted by  $S$  with individual elements

$s_{b,b}$ ,  $s_{b,g}$ ,  $s_{g,b}$ , and  $s_{g,g}$

Figure 4: nadjusted Shares

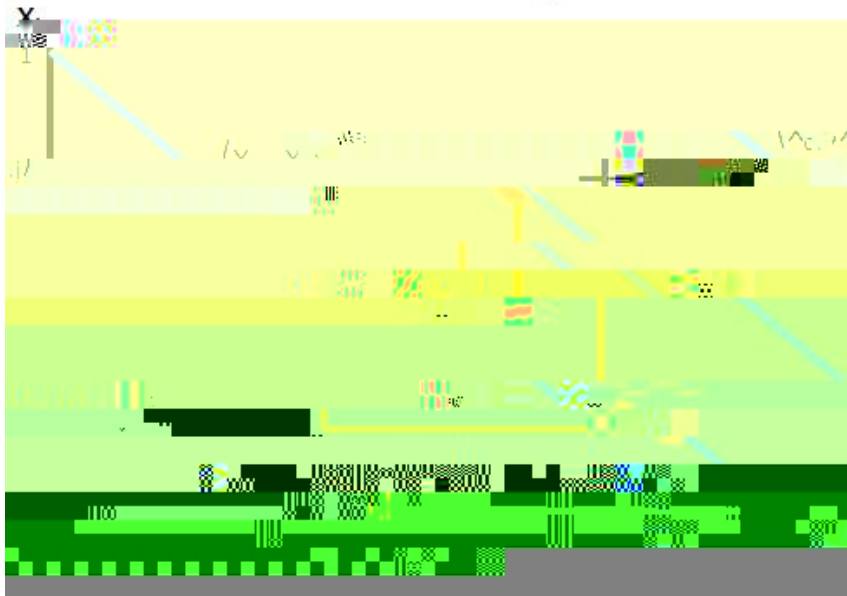
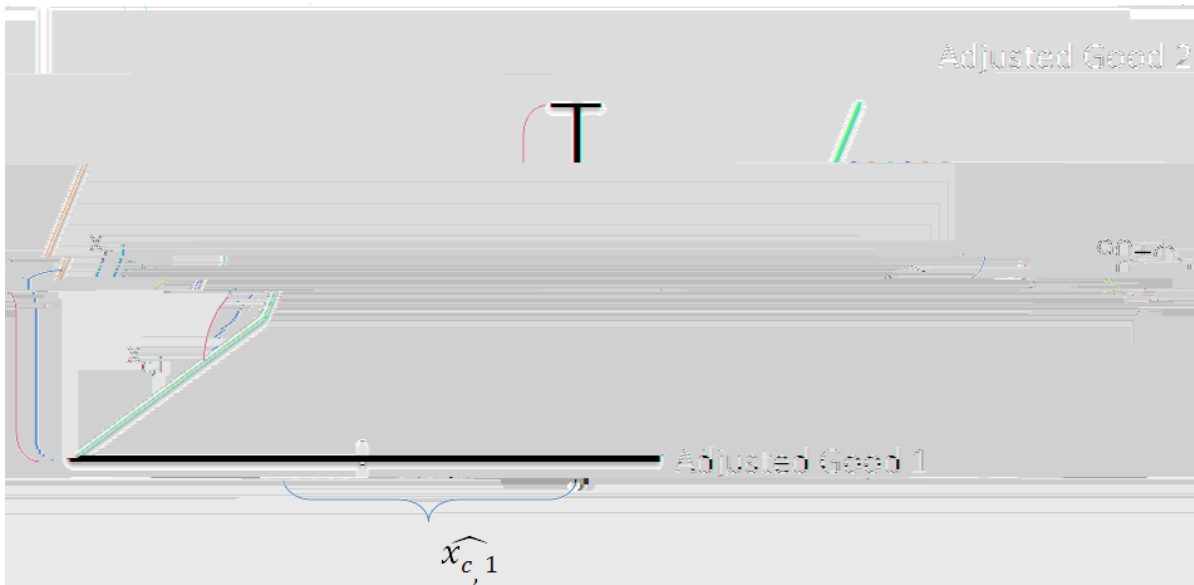


Figure 5: Projection onto Principal Axes



$$y_{c,2} = y_{c,2,1} + y_{c,2,2} = 0 + 0.247 = 0.247$$

Lastl , because this is shares data, the the sum of the shares must equal 1:

$$y_{c,1} + y_{c,2} = 0.961 + 0.247 = 1.208$$

And then:

$$y_{c,1}^b = \frac{0.961}{1.208} = 0.796$$

$$y_{c,2}^b = \frac{0.247}{1.208} = 0.204$$

Equivalentl for countr  $d$ : Project the  $y_{d,b}$  vector onto the  $y_{d,1}$  and  $y_{d,2}$  axes:

$$y_{d,1,1} = y_{d,b} \cos(0^\circ) = 0.7(1) = 0.7$$

$$y_{d,2,1} = y_{d,b} \sin(0^\circ) = 0.7(0) = 0$$

Similarl , projecting the  $y_{d,g}$  vector onto the  $y_{d,1}$  and  $y_{d,2}$  vector space fields:

$$y_{d,1,2} = y_{d,g} \cos(18^\circ) = 0.3(0.951) = 0.285$$

$$y_{d,2,2} = y_{d,g} \sin(18^\circ) = 0.3(0.309) = 0.093$$

And the last step for countr  $c$  is to add together the results of the two projections:

$$y_{d,1}^b = y_{d,1,1} + y_{d,1,2} = 0.7 + 0.285 = 0.985$$

$$y_{d,2}^b = y_{d,2,1} + y_{d,2,2} = 0 + 0.093 = 0.093$$

Lastl , because this is shares data, the the sum of the shares must equal 1:

$$y_{d,1}^b + y_{d,2}^b = 0.985 + 0.093 = 1.078$$

And then:

$$y_{d,1}^c = \frac{0.985}{1.078} = 0.914$$

And notice that for the unadjusted vectors, the normalized Euclidean distance could have been:

$$dist_{c;d}^a = \frac{1}{\sqrt{2}} \sqrt{(0.2 - 0.7)^2 + (0.8 - 0.3)^2} = \frac{0.707}{\sqrt{2}} = 0.5; \quad (35)$$

As an aside, the Law of Cosines distance metric in Knippenberg (2012) finds the Euclidean distance between the unadjusted vectors which is equivalent to the distance between the adjusted but non-normalized vectors, see Appendix B for the proof.

The reader can hopefully see that when similarity is zero, then  $\cos(0) = 1$ , allowing the orthogonalization process to return the original vectors of guns and butter. So, to take the analysis a step further, assume that the export share vector of each country is in exactly the same proportion as their production vectors:

$$x_c = y_c = \begin{pmatrix} 0.2 \\ 0.8 \end{pmatrix}; \quad x_d = y_d = \begin{pmatrix} 0.7 \\ 0.3 \end{pmatrix};$$

To abstract from any confounding effects, assume that each country has equal economic output, that these are the only two countries in the world, and that each exports goods equal to 1 normalized unit of value. Abstracting away from any theory on how the countries are trading or on their quantities of that trade, the empirical international trade literature suggests a number of measures.

Using the original, unadjusted trade vectors, the composition of bilateral trade is given by  $GL_{c;d}^{un} = 0.5$ . <sup>0</sup> Contrastinal trade ls4d[(:)(ectorsis)-s1nal trade

1945, 1944):

$$H_c = \sum_{i=1}^n \frac{x_{c,i}^2}{X_c} \quad (3)$$

Where  $x_{c,i} = X_c$  is the share of good  $i$  in the export bundle of country  $c$ . Using the original data, this comes out to be  $H_c^{un} = 0.68$  and  $H_d^{un} = 0.58$ . And using the adjusted data vectors this comes out as:  $H_c^a = 0.584$  and  $H_d^a = 0.777$ . Again, the economic significance of the differences between these two measures is subjective, but what is interesting is that the ordering has reversed. Where in the unadjusted index, country  $c$  is the more concentrated country, in the adjusted index, the more concentrated country is now  $d$ .

## 7. The Product Space

Here I demonstrate the change-of-coordinates orthogonalization procedure in a high-dimensional example: that of the product space of international trade. The product space is an idea conceived and visualized by Hidalgo, et al. (2007), who use export shares to find a measure of similarity between export product categories and then map them using a network analysis approach. I take their analysis a step further by using the similarity measures to adjust the original country export vectors, and I show that the measurements, while clearly correlated, are very different. Because of the computational intensity of the orthogonalization procedure<sup>12</sup>, I have only produced estimates on the export similarity measure. A detailed treatment of the consequences of changing the export similarity measure can be found in Knippenberg (2012), here I insert the new export similarity measure into a gravity equation of international trade and find very different results from previous studies.

Export similarity was first conceived by Finger and Kreinin (1979) as a simple measure for comparing export content across either countries or time. I denote this measure as  $FK_{c,d}$  and it is defined in equation (22). I use a version of  $FK_{c,d}$ , which is derived in Sun and Ng (2000), and is given in equation (19). The measure has been used in hundreds of academic papers on international trade.

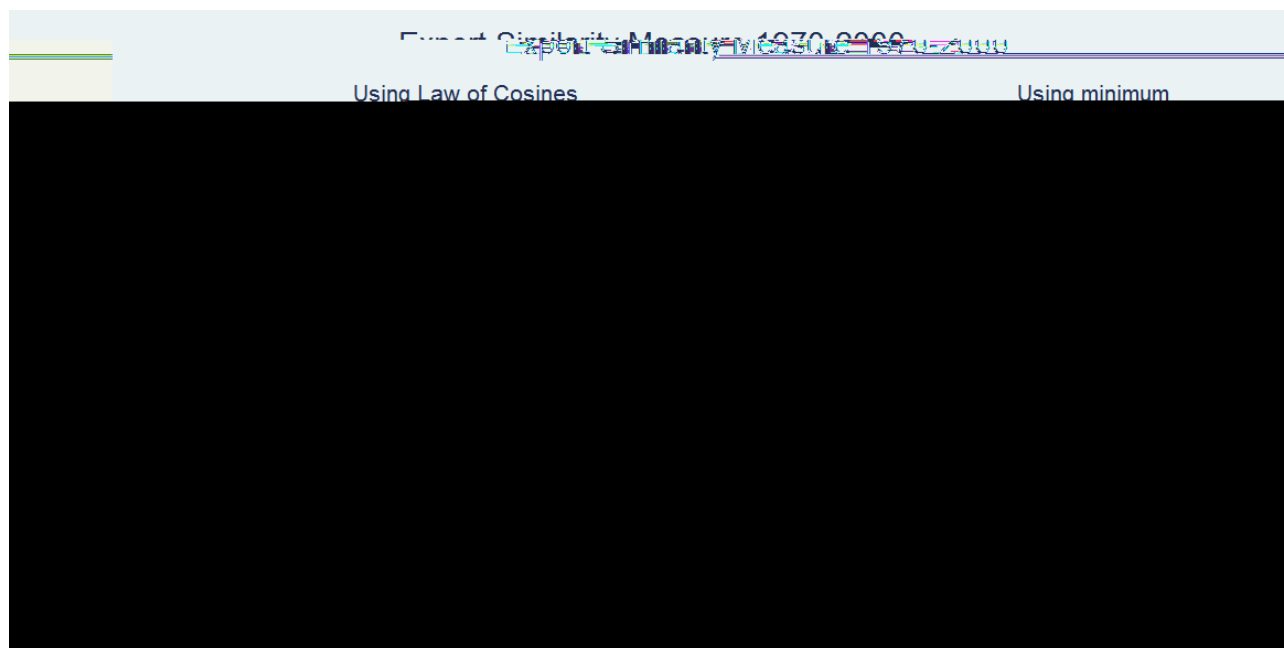
---

<sup>12</sup>Computing this variable for 47,653 observations took approximately four weeks on a desktop computer with a quad-core 3.3Ghz processor.

The steps taken to arrive at these export similarity indices are as follows. First I downloaded the export data from Feenstra's website. The data is 4-digit SITC trade data with 799 categories. I am using 5-year intervals from 1970 to 2000 for 133 countries. Second, I transform export values into export shares. Third, I follow Hidalgo, et al. (2007) to calculate similarity between export categories. Fourth, using this similarity matrix and export shares, I apply the orthogonalization procedure to obtain the adjusted export vectors. Lastly, I apply a Euclidean Distance algorithm



Figure 6: Histograms of Export Similarity Measures



to situations in which similarity between variates or correlation is already accounted for, such as regression analysis or principal components analysis. Given the nature of international trade shares data, this orthogonalization procedure is clearly applicable. Furthermore, the Trade Theory models assume an equal marginal rate of substitution between varieties of a good. However, if two varieties are more similar than either are to a third, then equal marginal rates of substitution cannot mathematically hold. After applying this orthogonalization procedure, the marginal rates of substitution between the adjusted goods should be equal because the variables are orthogonal to one another. This would make the data consistent with the theory, and is a promising area for future research.

This procedure works only when a bivariate notion of “similarity” or “distance” is computable, as these similarity measures directly feed into the equation. This procedure is not applicable where similarity is not defined or calculable. Finding a way to calculate this similarity in many different contexts is an area for future research where notions of covariance, correlation, may be very important. Furthermore, a simple lack of a way to calculate similarity doesn’t make the previous distance metrics any more valid - they are still computed using the incorrect coordinate system.



## Conclusion

I like the following quote from a linear algebra textbook: “Physical Laws must be independent of any particular coordinate system used in describing them mathematically, if they are to be valid” Spiegel (1959 pg 11). It reminds me that just because you can measure something doesn’t mean that what you have measured must necessarily obey the laws of our theory: sometimes a researcher has to manipulate data to make sense of it. In the case of shares data, often it is in the wrong coordinate system and must be converted to the proper system before familiar measures can be applied, like measures of distance in the rectangular coordinate system. I have argued throughout this paper that arbitrary classifications are not automatically defined by the rectangular coordinate system. However the rectangular coordinate system is the only requirement for applying familiar statistical distance metrics. In other words, the principle axes of the coordinate system are rarely the same as the axes of the data, so distance metrics cannot be immediately applied.

Besides the justification of the orthogonalization procedure, the previous paragraphs have also laid out areas for future research. The more mundane of these include re-estimating the effects of unbiased indices on outcomes. For example in trade, this could include the effect of export similarity or diversification on bilateral trade (Knippenberg 2012), or likewise the effect of the Grubel-Lloyd or Herfindahl Indices on various response variables. Theoretical research, on the other hand, holds even more promising avenues. As touched upon earlier, the continuum of goods assumption in International Trade can be re-visited: after normalizing the goods vectors, each adjusted good should have equal marginal rates of substitution, as each represents an orthogonal underlying good. I have written this paper in an attempt to state as general as possible about its applications: the extensive examples in international trade are merely a consequence of my own experience. The concepts described herein have wide applicability in all areas of empirical research and I look forward to conducting these applications in the near future.

## References

- [1] Axel, Sheldon Ja . (1997).

- [18] Lin, Chii-Dong. (1995) "Hyperspherical Coordinate Approach to Atomic and Other Coulombic Three-Body Systems." *Physics Reports* 257: 1-83.
- [19] Malaney, Pia. (1995) "The Index Number Problem: A Differential Geometric Approach." *Ph.D. Thesis*, Harvard University, Department of Economics, December 1995.
- [20] Marsaglia, George. (1972). "Picking a Point from the Surface of a Sphere". *The Annals of Mathematical Statistics* 43(2): 45-48.
- [21] Michael, Michael. (1995). "Trade Preferential Agreements in Latin America: An Experimental Assessment" World Bank Policy Research Paper 1583, March 2000.
- [22] Mikic, Mia. (2005). "Commonly Used Trade Indicators: A Note" ERT Capacity Building Workshop on Trade Research.
- [23] Munkres, James. (1991). *An Introduction to Manifolds* Westview Press, 1991, Chapter six.
- [24] Ng, Francis. (2002). "Appendix B: Trade Indicators and Indices" in *Development, Trade and the WTO: A Handbook* Worldbank Bank Publications, Washington, D.C, pg 585-588.
- [25] Panda, Avneet, Edward Y. Chang and Arun Qamra. (2000). "Hypersphere Indexer". *Lecture Notes in Computer Science*, Volume 4080, 2000.
- [26] Spiegel, Murray. (1959). *Vector Analysis With an Introduction to Tensor Analysis* McGraw-Hill, New York, New York. ISBN 07-0 0228-X.
- [27] Sun, Guang-hen and Ye Kang-g. (2000). "The measurement of structural differences between economies: An axiomatic characterization." *Economic Theory* 1 : 313-321.

## A Appendix : Proof to Equivalence of Finger-Kreinin and Sun-Ng Distance Measures

This section provides a proof that the export similarity measures from Finger and Kreinin (1979) (FK) and Sun and Ng (2000) (SN) are perfectly negatively correlated. Because of the minimum function in FK and the absolute function in SN, this proof is not conducive to deduction, but an inductive argument is easier to show. Define FK and SN according to their authors:

$$FK = \sum_{i=1}^n \min\left(\frac{X_{c,i}}{X_c}; \frac{X_{d,i}}{X_d}\right); \quad (37)$$

and:

$$SN = \sum_{i=1}^n \frac{|X_{c,i} - X_{d,i}|}{2} \quad (38)$$

### Proposition:

Let  $n$  denote the number of export products. Let  $c$  and  $d$  be any two countries. Denote export share of good  $i$  in country  $c$  as  $X_{c,i}$ , where  $i = 1; \dots; n$ . Because  $X_{c,i}$  is an export share,

$$\sum_{i=1}^n X_{c,i} = 1; \quad (39)$$

is satisfied by the definition of a share. The same equation also holds for any other country  $d$ . Let the sums  $FK$  and  $SN$  be defined as above, then the following equality always holds:

$$SN = 1 - FK \quad (40)$$

### Proof:

#### A.1 Case 1.1

Let  $n = 2$  and Let  $X_{c,1} = X_{d,1}$ , then because  $X_{c,1} + X_{c,2} = 1$  and  $X_{d,1} + X_{d,2} = 1$ , it must also be true that  $X_{c,2} = X_{d,2}$ . In this case,

$$\begin{aligned} FK &= \min(X_{c,1}; X_{d,1}) + \min(X_{c,2}; X_{d,2}) \\ &= X_{c,1} + X_{c,2} \\ &= 1 \end{aligned} \quad (41)$$

Similarly,

$$SN = \frac{X_{c,1} - X_{d,1}}{2} + \frac{X_{c,2} - X_{d,2}}{2} \quad (42)$$

By assumption,  $X_{c,1} - X_{d,1} = 0$  and since  $X_{c,2} = X_{d,2}$ , then  $X_{c,2} - X_{d,2} = 0$ . Therefore,  $SN = 0$  and  $FK = 1$ .  $\square$



#### A.4 Case 2.2

Let  $n \geq 2$  and  $x_{c;i} > x_{d;i}$  for  $i = 1, \dots, j$ . Let  $x_{c;i} > x_{d;i}$  for  $i = 1, \dots, k$ . Let  $x_{c;i} > x_{d;i}$  for  $i = 1, \dots, l$ . Where  $j + k + l = n$ , and  $j, k, l \geq 0$ . Then by definition, Equation (37) implies:

$$FK = \sum_{i=1}^j x_{c;i} + \sum_{i=1}^k x_{d;i} + \sum_{i=1}^l x_{c;i} \quad (56)$$

Or equivalently where the last summation is replaced by  $x_{d;i}; i = 1, \dots, l$ . By the shares definition, Equation (39) implies for country  $c$ :

$$\sum_{i=1}^j x_{c;i} + \sum_{i=1}^k x_{c;i} + \sum_{i=1}^l x_{c;i} = 1; \quad (57)$$

as well as for country  $d$ :

$$\sum_{i=1}^j x_{d;i} + \sum_{i=1}^k x_{d;i} + \sum_{i=1}^l x_{d;i} = 1; \quad (58)$$

And by the definition SN (38):

$$SN = \frac{1}{2} \left( \sum_{i=1}^j (x_{c;i} - x_{d;i}) + \sum_{i=1}^k (x_{d;i} - x_{c;i}) + \sum_{i=1}^l (x_{c;i} - x_{d;i}) \right); \quad (59)$$

Distributing through the summations and rearranging yields:

$$SN = \frac{1}{2} \left( \sum_{i=1}^j x_{c;i} - \sum_{i=1}^k x_{c;i} + \sum_{i=1}^l x_{c;i} - \sum_{i=1}^j x_{d;i} + \sum_{i=1}^k x_{d;i} - \sum_{i=1}^l x_{d;i} \right); \quad (60)$$

Rearranging (57) implies:

$$\sum_{i=1}^j x_{c;i} + \sum_{i=1}^l x_{c;i} = \sum_{i=1}^k x_{d;i} + \sum_{i=1}^l x_{d;i}$$



Simplifying:

$$SN = 1 \prod_{i=1}^n X_{c;i} \prod_{i=1}^n X_{d;i} \prod_{i=1}^n X_{d;i} \quad (65)$$

Then substituting in the definition of  $F$ , Equation (56), yields the desired result:

$$SN = 1 \quad FK: \quad (66)$$

Thus the relationship holds for both  $n = 2$  and  $n > 2$ , proving the proposition by induction.

## B Appendix : Proof of Equivalence Between Orthogonalization and the n-Dimensional Law of Cosines

**Pr o p o s i t i o n:**

In an  $n$ -Hilbert space, the norm distance  $\|x_1 - x_2\|$  with similarity matrix  $S$  equals  $\|Sx_1 - Sx_2\|$  with similarity matrix  $I$ , the identity matrix.

I will only prove this equivalence for the two-good case. The notation needed to prove



However, this property does not apply to a heterogeneous space such as the product space for two reasons. First of all, as evidenced in Hidalgo et al (2007), the Product Space is extremely heterogeneous: some areas of the product space are dense and others are disparate. In terms of the previous equation, this difference can be seen in the fact that some areas are more densely populated than others.

from: nor(anoms.)6oy7 orstri(o)-1s(in)-