

Influencia de datos discordantes en el modelo Tucker3

*José Luis Vásquez Pérez*¹ y *Emma Cambillo Moyano*²

Resumen: El presente artículo aborda la evaluación del Modelo Tucker3 mediante el análisis de conjuntos de datos macroeconómicos de los departamentos del Perú durante el periodo 2007-2021. El enfoque se centra en la cuantificación del impacto de observaciones discordantes, particularmente en relación con Lima y Callao, sobre la estructura de la solución, la determinación de las dimensiones y la representación de las vías en sus respectivos componentes. Los resultados destacan el relevante papel de Lima y Callao en la variabilidad total, aunque este impacto no conlleva modificaciones estructurales en el modelo. No obstante, se observa una influencia en la importancia relativa de los componentes en las diferentes vías del Modelo Tucker3.

Palabras clave: Tucker3, análisis de tres vías, datos discordantes, macroeconomía

Influence of discordant data on the Tucker3 model

Abstract: The present article addresses the assessment of the Tucker3 Model through the analysis of macroeconomic datasets from the departments of Peru spanning the period 2007-2021. The focus is on quantifying the impact of discordant observations, particularly in relation to Lima and Callao, on the structure of the solution, the determination of dimensions, and the representation of pathways in their respective components. The results underscore the significant role of Lima and Callao in the total variability, although this impact does not entail structural modifications in the model. Nevertheless, an influence is observed in the relative importance of components across different pathways in the Tucker3 Model.

Keywords: Tucker3, multiway, macroeconomics

Recibido: 11/12/2023 *Aceptado:* 06/07/2024 *Publicado online:* 30/12/2024

¹UNMSM, Facultad de Ciencias Matemáticas. e-mail: joluvasquez@gmail.com

²UNMSM, Facultad de Ciencias Matemáticas. e-mail: ecambillom@unmsm.edu.pe

1. Introducción

En el ámbito del análisis multivariante, la integración y representación precisa de los datos en modelos estadísticos son fundamentales para garantizar interpretaciones confiables, especialmente cuando los arreglos de los datos se extienden a través de múltiples dimensiones. Este artículo se enfoca en el análisis tridimensional de datos, considerando individuos, variables y ocasiones, los cuales son sintetizados a través del Modelo Tucker3. Se examina específicamente la presencia de datos discordantes en conjuntos de datos macroeconómicos departamentales del Perú durante el periodo 2007-2021.

El Modelo Tucker3 ha demostrado su eficacia al abordar conjuntos de datos con estructuras multidimensionales, ofreciendo una descomposición que revela las interacciones subyacentes. No obstante, la presencia de observaciones atípicas puede introducir distorsiones significativas en los resultados, comprometiendo la validez de las conclusiones obtenidas.

Al seleccionar datos macroeconómicos departamentales del Perú a lo largo de un periodo extenso, nuestro objetivo va más allá de comprender las fluctuaciones económicas a lo largo del tiempo. También nos proponemos abordar desafíos específicos asociados con la heterogeneidad geográfica. Este análisis técnico se enfoca en cuantificar el impacto de datos discordantes en el Modelo Tucker3, con especial atención a la región de Lima y Callao. En este contexto, nuestro artículo contribuirá a la literatura al proporcionar una evaluación cuantitativa de la influencia de datos discordantes en el Modelo Tucker3, aplicada de manera específica a datos macroeconómicos departamentales en el contexto peruano.

2. Metodología

El modelo Tucker3 constituye una generalización del análisis de componentes principales y la descomposición de valores singulares, diseñado para tratar arreglos de tres vías (Andersson y Bro, 1998; Kroonenberg, 2008; Kiers, 2000). Este enfoque se aplica a conjuntos de datos que contemplan las puntuaciones de individuos en diversas variables y bajo diversas condiciones, siendo estas últimas, específicamente, las distintas ocasiones (Amaya y Pacheco, 2002; Paredesa y cols., 2018; Gallo, 2015; Dell'Anno y Amendola, 2015; Bautista Mendoza, 2009). En este sentido, el modelo Tucker3 facilita la síntesis de información clave de los datos de tres vías, permitiendo condensar la esencia de la información en un número reducido de componentes en cada dimensión (Kroonenberg, 1983, 2008; Timmerman y Kiers, 2000).

El modelo Tucker3 describe la descomposición de una matriz de tres dimensiones $I \times J \times K$ de la siguiente manera:

$$x_{ijk} = \sum_{p=1}^P \sum_{q=1}^Q \sum_{r=1}^R a_{ip} b_{jq} c_{kr} g_{pqr} + e_{ijk} \quad (1)$$

donde $i = 1, \dots, I$, $j = 1, \dots, J$, $k = 1, \dots, K$, y a_{ip} , b_{jq} , y c_{kr} representan los elementos de las matrices $\mathbf{A}(I \times P)$, $\mathbf{B}(J \times Q)$ y $\mathbf{C}(K \times R)$, respectivamente. La matriz central $\mathbf{G}(P \times Q \times R)$ contiene los elementos g_{pqr} . El término e_{ijk} denota los errores asociados con x_{ijk} .

El modelo se ajusta a un conjunto de datos minimizando la suma de los cuadrados de los términos de error, expresada como:

$$f(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{G}) = \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K \left(x_{ijk} - \sum_{p=1}^P \sum_{q=1}^Q \sum_{r=1}^R a_{ip} b_{jq} c_{kr} g_{pqr} \right)^2 \quad (2)$$

El ajuste se calcula como la suma de los cuadrados de las diferencias entre los datos observados y las aproximaciones del modelo.

La estimación óptima de mínimos cuadrados, según Kroonenberg (1983), se logra cuando el

ajuste es igual a la suma de los cuadrados de las aproximaciones de los datos, representado por:

$$\hat{x}_{ijk} = \sum_{p=1}^P \sum_{q=1}^Q \sum_{r=1}^R a_{ip} b_{jq} c_{kr} g_{pqr}$$

Esta estimación se utiliza frecuentemente en una proporción de ajuste, definida como:

$$\frac{\sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K \hat{x}_{ijk}^2}{\sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K x_{ijk}^2}$$

Esta proporción proporciona una medida de la calidad del ajuste en relación con la variabilidad total de los datos (Kiers y Der Kinderen, 2003). Por otro lado, las matrices **A**, **B** y **C** se asumen ortonormales en columnas, lo que puede no afecta al ajuste óptimo; entonces el ajuste viene dado por la suma de los elementos centrales al cuadrado (Kroonenberg y De Leeuw, 1980).

Ahora, un paso importante es seleccionar la cantidad de componentes que se utilizarán. En el modelo Tucker3 el análisis se realiza por combinaciones de componentes, de manera que una estrategia es considerar un conjunto de valores y considerar la solución más útil, esto en el sentido de parsimonia, interpretabilidad o estabilidad. Timmerman y Kiers (2000) proponen un procedimiento denominado *DIFFIT* basado en la comparación de los ajustes relacionado con el número total de los componentes utilizados en el análisis, esto es simplemente una analogía a la prueba de Cattell (Cattell, 1966):

1. Se recopilan los valores de los ajustes en diversas combinaciones de soluciones del modelo Tucker3.
2. Se comparan soluciones con el mismo número total de componentes ($S = P + Q + R$) y se retienen únicamente aquellas soluciones que, para un valor dado de S , ofrecen el mejor ajuste.
3. Se calculan las diferencias como la disparidad entre la mejor solución con S componentes y aquella que tiene $S-1$ componentes. A partir de estos resultados, se obtiene un subconjunto de soluciones para las cuales:

$$dif_s > dif_{s+j}, \text{ para todo } j > 0$$

Estas soluciones se denotan como $m = 1, \dots, M$, y los diferentes valores de S están dados por $t(m)$, con sus respectivos valores dif asociados dados por $dif_{t(m)}$.

4. Se calcula

$$b_{t(m)} = \frac{dif_{t(m)}}{dif_{t(m+1)}}$$

para determinar la relación entre el aumento del ajuste resultante del componente $t(m)$ -ésimo y el del siguiente componente de interés, el $t(m+1)$ -ésimo.

5. Se seleccionan solo aquellos números de componentes para los cuales:

$$dif_{t(m)} > \frac{\|\mathbf{X}\|^2}{S_{max}}$$

Donde S_{max} es el número total máximo de componentes sensibles. De las soluciones restantes, se elige aquella que tenga el valor más bajo de $b_{t(m)}$ y se denota el número total asociado de componentes como S_c .

6. Se selecciona el número de componentes asociados con el mejor ajuste entre todos los modelos utilizando un total de componentes $S_c = P + Q + R$.

3. Resultados

La aplicación del modelo Tucker3 a datos macroeconómicos ha sido poco explorada; no obstante, vale la pena destacar los esfuerzos de investigadores como Bautista Mendoza (2009), Amaya y Pacheco (2002), así como Paredesa y cols. (2018). Estos estudios se han enfocado en analizar estructuras de tres vías utilizando datos colombianos. Además, Gallo (2015) ha investigado el consumo de energía en varios países europeos, mientras que Dell'Anno y Amendola (2015) han utilizado datos del metro de Barcelona en sus análisis.

En la presente aplicación, los datos utilizados provienen del Instituto Nacional de Estadística e Informática del Perú y constituyen una muestra de 24 departamentos del país ¹. Esta muestra abarca 12 variables macroeconómicas que representan el Valor Bruto Agregado en diversas actividades económicas ², durante un período que comprende los años 2007 a 2021. En consecuencia, se dispone de un conjunto de datos con dimensiones $24 \times 12 \times 15$.

No obstante, la muestra presenta una característica notable en relación con los individuos. Lima y Callao tienen un impacto significativo en todas las variables consideradas, lo que podría introducir distorsiones en la solución del modelo debido a la variabilidad que estos dos departamentos podrían introducir. De hecho, esta situación podría afectar el rendimiento del algoritmo de mínimos cuadrados alternantes, según se sugiere en Pravdova y cols. (2001). Por este motivo, se llevarán a cabo dos subconjuntos de datos: uno que incluirá a Lima y Callao ($24 \times 12 \times 15$) y otro que los excluirá ($23 \times 12 \times 15$). Se realizará una comparación de los resultados obtenidos de ambas muestras para determinar y evaluar la influencia de la presencia de Lima y Callao.

Ambos modelos fueron estimados mediante el algoritmo de mínimos cuadrados alternantes, utilizando un criterio de convergencia establecido en $1e - 10$ con un máximo de 100 iteraciones. Siguiendo los seis pasos previamente mencionados, se procede al cálculo de la suma de los componentes S y a la determinación de las diferencias entre la suma de cuadrados de los residuos mínimos para s y $s + j$. Este proceso permite obtener el coeficiente $b_{t(m)}$. Los resultados obtenidos se detallan en la tabla (1), que presenta la selección de los mejores ajustes bajo el criterio *Diff-Diff*. Se analizan las combinaciones de los componentes P , Q , y R . Con base en estos resultados, se selecciona un valor correspondiente al mínimo $b_{t(m)}$, el cual se alcanza con $S = 5$ y una combinación de $P = 2, Q = 2, R = 1$.

En base a la configuración inicial ($24 \times 12 \times 15$), se logró una elección de $S = 5$ con $P = 2, Q = 2$, y $R = 1$, alcanzando un nivel de ajuste del 99,2%. En contraste, la muestra alternativa ($23 \times 12 \times 15$), que excluye a Lima y Callao, exhibió una disposición de los componentes en todas las dimensiones con valores de $P = 2, Q = 2$, y $R = 1$, capturando el 86,6% de la variabilidad total en los datos.

En un primer análisis, se evidencia que la estructura de la solución permaneció inalterada, ya que en ambas situaciones se obtuvo el mismo número de componentes en las tres dimensiones ($S = 5$), como se ilustra en la figura (1). Sin embargo, cabe resaltar una discrepancia del 12,6% en puntos porcentuales entre ambas soluciones, siendo esta diferencia atribuible de manera específica a la influencia de Lima y Callao.

En cuanto a las contribuciones, se destaca que el modelo que incorpora a Lima y Callao presenta dos componentes ($P = 2$) en la dimensión de individuos, representando el 97,4% y el 1,8%, respectivamente. Asimismo, en la dimensión de variables, esta configuración se desglosa en dos componentes ($Q = 2$) con contribuciones del 97,4% y el 1,8%. Finalmente, la dimensión

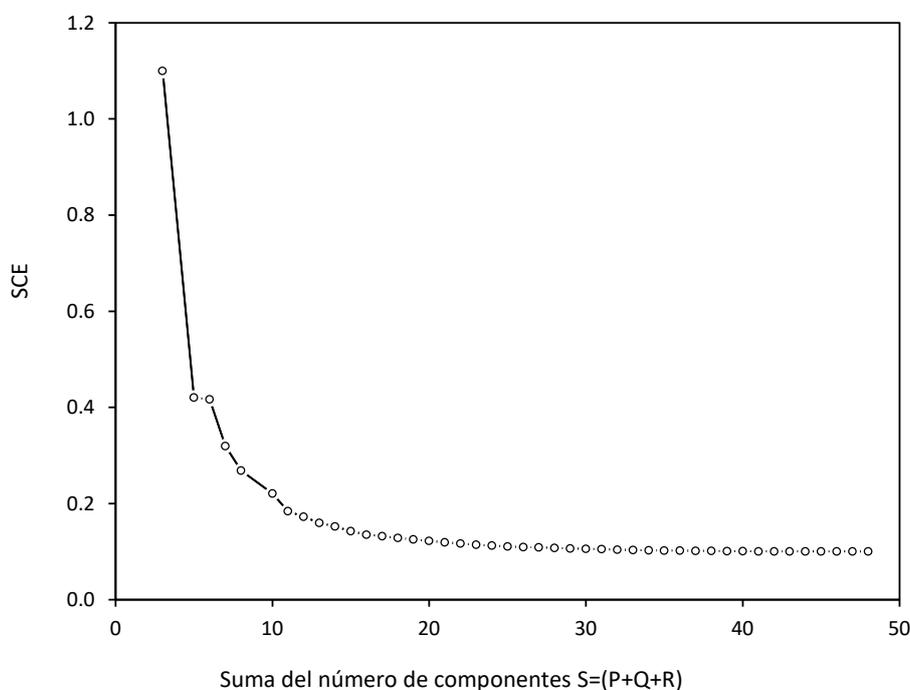
¹Se consideran los siguientes departamentos: Amazonas, Áncash, Apurímac, Arequipa, Ayacucho, Cajamarca, Cusco, Huancavelica, Huánuco, Ica, Junín, La Libertad, Lambayeque, Loreto, Madre de Dios, Moquegua, Pasco, Piura, Puno, San Martín, Tacna, Tumbes, Ucayali y Lima y Callao

²Estas variables describen el Valor Bruto Agregado (VBA) de 12 actividades económicas: agricultura; pesca y acuicultura; extracción de petróleo, gas, minerales y servicios conexos; manufactura; Electricidad, Gas y Agua; Construcción; comercio, mantenimiento y reparación de vehículos automotores y motocicletas; transporte, almacenamiento, correo y mensajería; Alojamiento y Restaurantes; Telecomunicaciones y otros Servicios de Información; administración pública y defensa; otros Servicios

Cuadro 1: Comparación de los resultados del criterio *Diff-Diff*

S=P+Q+R	P	Q	R	Incluye Lima y Callao		Excluye Lima y Callao	
				SCEmin	b	SCEmin	b
3	1	1	1	1.1	7.0	1.1	4.1
5	2	2	1	0.4	0.0	0.5	0.0
6	2	2	2	0.4	1.9	0.5	1.4
7	3	3	1	0.3	1.1	0.3	0.1
8	3	3	2	0.3	1.3	0.3	3.2
10	4	4	2	0.2	2.8	0.2	1.2
11	5	3	3	0.2	0.8	0.2	0.3
12	5	4	3	0.2	1.3	0.2	3.8
13	6	4	3	0.2	0.7	0.2	0.8
14	6	5	3	0.2	1.4	0.2	1.1
15	7	4	4	0.1	2.0	0.1	1.2

Figura 1: Scree Plot del modelo tucker3



que aborda las ocasiones cuenta con una única componente ($R = 1$) y una contribución del 99,2%.

En relación con el modelo que excluye a Lima y Callao, como se mencionó previamente, presenta la misma combinación de componentes; sin embargo, se observa una disminución en la contribución de cada componente. En la dimensión de individuos, con sus dos componentes ($P = 2$), estas representan ahora un 66,5% y un 20,1%, respectivamente. En cuanto a las variables, las dos componentes también abarcan la misma variabilidad que en el caso de los individuos. Por último, en la dimensión de ocasiones, la componente única ($R = 1$) captura un 86,6% de la variabilidad.

Evidentemente, Lima y Callao juegan un papel preponderante tanto en el nivel total de contribución como en los componentes de las tres dimensiones, como se ilustra en la tabla (2). La exclusión de estos elementos resulta en una mayor relevancia de la segunda componente en las

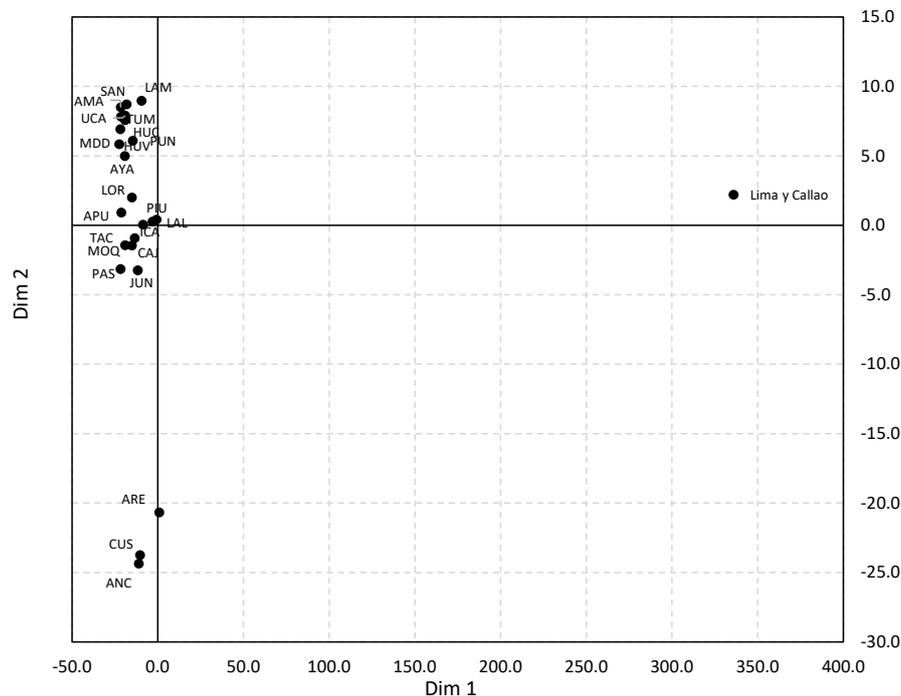
Cuadro 2: Contribución de componentes a la inercia total de los dos modelos

	Total		Sin Lima y Callao		Diferencia	
	Comp. 1	Comp. 2	Comp. 1	Comp. 2	Comp. 1	Comp. 2
Total	99.2		86.6		12.6	
A	97.4	1.8	66.5	20.1	30.9	-18.4
B	97.4	1.8	66.5	20.1	30.9	-18.4
C	99.2	—	86.6	—	12.6	—

dimensiones de individuos y variables. Esto sugiere que, aunque la omisión de Lima y Callao no incide en la solución ni en el número de componentes en las dimensiones, su impacto se refleja en una diferencia aproximada de 18,4 puntos porcentuales. En contraste con la descripción de Pravdova y cols. (2001), la muestra no exhibe una alta proporción de datos atípicos, siendo inferior al 40 %. De hecho, Lima y Callao emergen como un factor significativo en la variabilidad total, aunque no afectan la solución per se.

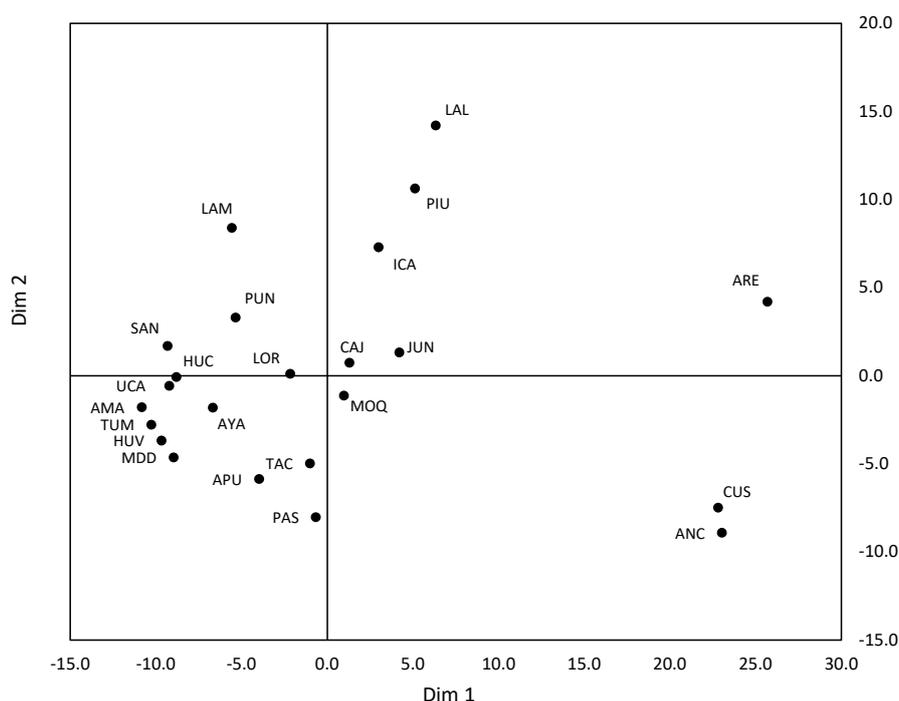
En relación con las triadas, que son componentes de la matriz \mathbf{G} , se observa una consistencia en magnitud con respecto a los componentes. Esto se debe principalmente a que las ocasiones se representan mediante una única componente ($R = 1$). En el escenario que incluye a Lima y Callao, esta dimensión exhibe elementos contribuyentes del 97,4 % y 1,8 % en su matriz diagonal. Por otro lado, en la situación sin Lima y Callao, el primer elemento disminuye a 66,5 %, mientras que el segundo elemento aumenta a 20,1 %.

Figura 2: Departamentos del Perú, incluyendo Lima y Callao, en el plano (1,2)



La situación mencionada se aprecia de manera más evidente en la representación de los individuos (departamentos) en los planes factoriales de los componentes seleccionados. En este sentido, la figura (2) ilustra la representación de los individuos considerando Lima y Callao. En dicha representación, se destaca de manera exagerada en ambas componentes, lo cual genera una distorsión en las relaciones entre los demás departamentos, disminuyendo su visibilidad. Por otro lado, la figura (3) presenta la representación de los individuos excluyendo a Lima y

Figura 3: Departamentos del Perú, excluyendo Lima y Callao, en el plano (1,2)



Callao de la solución del modelo. En este caso, se logra una visualización más clara de las relaciones entre los departamentos, ya que se elimina la exageración anteriormente mencionada. Es importante señalar que la exclusión de Lima y Callao no altera de manera significativa las relaciones entre los demás departamentos en el plano. No obstante, en términos comparativos, se sugiere aislar a Lima y Callao para facilitar la comparación entre los diferentes departamentos. Excluyendo a Lima y Callao, se pueden destacar claramente grupos definidos en las relaciones entre los departamentos. Por un lado, Áncash y Cusco revelan una afinidad significativa, la cual podría vincularse con su relación respecto a la componente 1. Por otro lado, Lambayeque, Piura e Ica forman un grupo destacado debido a sus similitudes en actividades económicas. De manera notable, se observa una relación interesante entre Cajamarca, Junín y Cusco.

Al analizar la representación de las variables en los planos (1,2), las figuras (4) y (5) ilustran claramente el impacto que tuvo la inclusión o exclusión de Lima y Callao en dicha representación. En términos de patrones, se observa que ambas representaciones conservan las tendencias, las agrupaciones y las distancias. La distinción entre ambas radica en la magnitud de las dimensiones y en el hecho de que una figura representa la rotación de la otra.

En este contexto, la dispersión de las actividades económicas revela una marcada influencia de sectores como el petróleo, la minería y el gas. Posteriormente, los servicios destacan como una actividad de mayor importancia en relación con la dimensión 1. Se aprecian otras relaciones significativas con respecto a la electricidad, la pesca, los hoteles, las telecomunicaciones y el transporte.

4. Conclusión

El análisis de los arreglos en el modelo Tucker3 revela una influencia significativa de Lima y Callao en todas las variables, destacando su papel preponderante al contribuir con un aumento de 12.6 puntos porcentuales en la variabilidad total. La exclusión de estos elementos no altera la estructura de la solución; ambos arreglos mantienen la misma dimensión en los componentes

Figura 4: Variables, incluyendo Lima y Callao, en el plano (1,2)

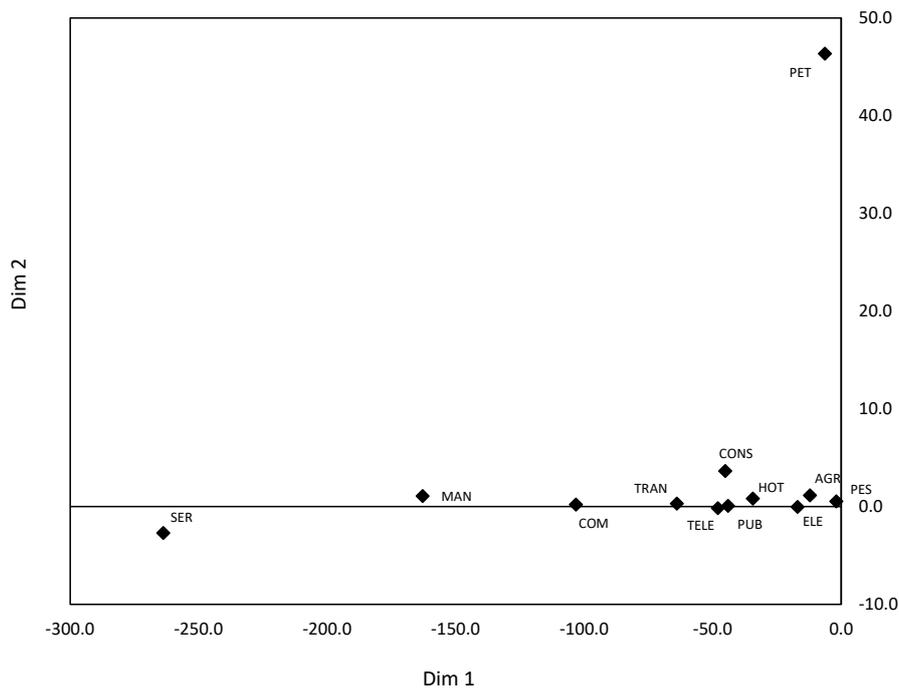
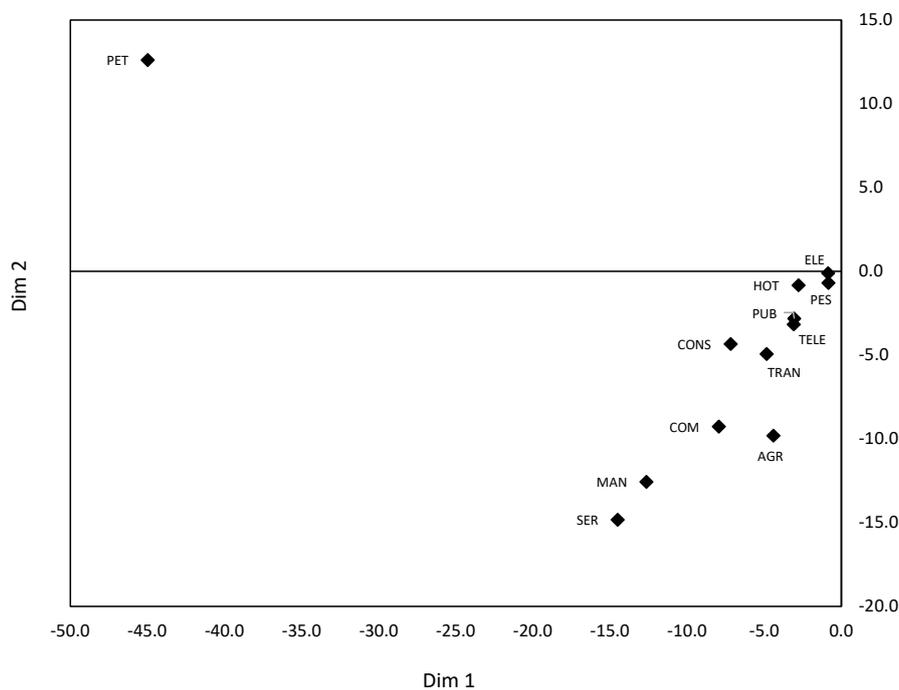


Figura 5: Variables, incluyendo Lima y Callao, en el plano (1,2)



($S = 5, P = 2, Q = 2, R = 1$). Sin embargo, esta exclusión impacta la contribución de los componentes en la vía de los individuos ($-30,9$ puntos porcentuales) y las variables ($+18,4$ puntos porcentuales), generando una diferencia sustancial en el nivel total de contribución.

Aunque la presencia de Lima y Callao no afecta la solución, su exclusión ofrece una mayor aplicabilidad para la segunda componente en la vía de los individuos y las variables. Este hallazgo resalta la importancia de realizar un análisis diferenciado que considere las particularidades

regionales, como la influencia de áreas geográficas, y la relevancia de evaluar la estabilidad y validez bajo diversas configuraciones de arreglos del modelo.

En cuanto al análisis de la representación de los individuos en los planes factoriales, se evidencia que la presencia de Lima y Callao genera una distorsión significativa, obstaculizando la observación de las relaciones entre los demás departamentos. La exclusión de Lima y Callao resulta en una mejora sustancial en la visualización, proporcionando una representación más clara de las interacciones entre los departamentos.

Es crucial señalar que la exclusión de Lima y Callao no provoca cambios significativos en las relaciones entre los demás departamentos, indicando estabilidad en dichas interacciones. No obstante, con el fin de facilitar la comparación entre departamentos, se recomienda aislar Lima y Callao. Esta exclusión permite identificar grupos definidos en las relaciones entre departamentos, como la afinidad entre Áncash y Cusco, posiblemente vinculada a su influencia en la componente 1, así como el grupo conformado por Lambayeque, Piura e Ica, caracterizado por similitudes en actividades económicas. También se destaca la relación interesante entre Cajamarca, Junín y Cusco. En conjunto, estos resultados ofrecen valiosas perspectivas para comprender las dinámicas interdepartamentales y subrayan la importancia de considerar distintas configuraciones en el análisis.

Finalmente, los resultados presentados indican que el papel preponderante de Lima y Callao como un individuo no ha generado un impacto significativo en la estructura general de la solución del modelo. Se observa que los mismos componentes persisten en cada una de las vías, lo que se refleja tanto en la dispersión de individuos en el plano de los componentes de la primera vía como en la distribución de las variables en su plano correspondiente. No obstante, es importante señalar que esta presencia también contribuye a un nivel notablemente elevado en la variabilidad total del arreglo de las tres vías.

Referencias bibliográficas

- Amaya, J. L., y Pacheco, P. N. (2002). Análisis factorial dinámico mediante el método tucker3. *Revista Colombiana de Estadística*, 25(1), 43–57.
- Andersson, C. A., y Bro, R. (1998). Improving the speed of multi-way algorithms:: Part i. tucker3. *Chemometrics and intelligent laboratory systems*, 42(1-2), 93–103.
- Bautista Mendoza, G. R. (2009). Comparación de los métodos tucker 3 y análisis factorial múltiple para el análisis de datos tres vías. *Departamento de Estadística*.
- Cattell, R. B. (1966). The scree test for the number of factors. *Multivariate behavioral research*, 1(2), 245–276.
- Dell’Anno, R., y Amendola, A. (2015). Social exclusion and economic growth: An empirical investigation in european economies. *Review of Income and Wealth*, 61(2), 274–301.
- Gallo, M. (2015). Tucker3 model for compositional data. *Communications in Statistics-Theory and Methods*, 44(21), 4441–4453.
- Kiers, H. A. (2000). Towards a standardized notation and terminology in multiway analysis. *Journal of Chemometrics: A Journal of the Chemometrics Society*, 14(3), 105–122.
- Kiers, H. A., y Der Kinderen, A. (2003). A fast method for choosing the numbers of components in tucker3 analysis. *British Journal of Mathematical and Statistical Psychology*, 56(1), 119–125.
- Kroonenberg, P. M. (1983). *Three-mode principal component analysis: Theory and applications* (Vol. 2). DSWO press.
- Kroonenberg, P. M. (2008). *Applied multiway data analysis*. John Wiley & Sons.
- Kroonenberg, P. M., y De Leeuw, J. (1980). Principal component analysis of three-mode data by means of alternating least squares algorithms. *Psychometrika*, 45(1), 69–97.
- Paredesa, J. E. R., Lozanoa, J. M. B., Restrepoa, J. M., y Londonob, D. A. (2018). Análisis de datos en tres vías para la evaluación de la dinámica económica de los departamentos en colombia durante el periodo 2000-2015. three way data analysis for the evaluation of the economic dynamic of departments of colombia during.
- Pravdova, V., Estienne, F., Walczak, B., y Massart, D. (2001). A robust version of the tucker3 model. *Chemometrics and Intelligent Laboratory Systems*, 59(1-2), 75–88.
- Timmerman, M. E., y Kiers, H. A. (2000). Three-mode principal components analysis: Choosing the numbers of components and sensitivity to local optima. *British journal of mathematical and statistical psychology*, 53(1), 1–16.