



## Análisis multivariante aplicado a encuesta nacional de ocupación y empleo (ENOE-INEGI)

Alvarez Sabouret Martín  
Facultad de Ciencias Económicas de la Universidad Nacional de Tucumán, Argentina  
*martinras92@gmail.com*

### Contenido

INTRODUCCIÓN.....	2
DESARROLLO DEL PROBLEMA.....	2
PREGUNTAS DE INVESTIGACIÓN .....	2
OBJETIVO GENERAL .....	2
OBJETIVOS ESPECÍFICOS.....	2
MARCO METODOLÓGICO .....	3
MARCO TEÓRICO.....	3
<b>DESARROLLO</b> .....	6
ANÁLISIS DESCRIPTIVO.....	8
Gráfico de dispersión matricial empleando variables sexo, ingresos mensuales, tipo de empleo .....	8
Gráfico de dispersión empleando ingresos mensuales, años de escolaridad, tipo de empleo.	9
Gráfico de dispersión con línea de tendencia empleando las variables Ingresos mensuales, Nivel educativo, años de educación.....	10
En la siguiente gráfica se representa el nivel de estudio discriminado en función del sexo ...	10
BOXPLOTS POR ESTADO .....	11
APLICACIÓN DE ANÁLISIS DE CORRESPONDENCIAS MÚLTIPLE (MCA).....	12
MCA empleando las variables sexo, posición ocupada y categoría de ingresos .....	12
Figura: MCA1 .....	13
MCA empleando variables ingresos y nivel educativo .....	15
MCA con variables SEXO INGRESOS y NIVEL EDUCATIVO .....	18
CLUSTER .....	20
REGRESIÓN LINEAL MÚLTIPLE .....	24
La ecuación de regresión lineal múltiple quedaría como sigue: .....	25
APÉNDICE .....	26
BIBLIOGRAFÍA .....	31



## INTRODUCCIÓN

La ciencia de datos es hoy en día la herramienta fundamental para la explotación de datos y la generación de conocimiento. Entre los objetivos que persigue se encuentra la búsqueda de modelos que describan patrones y comportamientos a partir de los datos con el fin de tomar decisiones o hacer predicciones. Es un área que ha experimentado un enorme crecimiento al extenderse el acceso a grandes volúmenes de datos e incluso su tratamiento en tiempo real, requiriendo de técnicas sofisticadas que puedan tratar con los problemas prácticos como escalabilidad, robustez ante errores, adaptabilidad con modelos dinámicos.

Ante el acceso a la encuesta nacional de ocupación y empleo del instituto INEGI, me surgió el interés de poder analizar la misma y el hecho de poder emplear herramientas de análisis multivariante y el uso de R-Studio aprendidos a lo largo del cursado de Análisis Cuantitativo II.

## DESARROLLO DEL PROBLEMA

Ante la obtención de una base de datos de acceso público, se presenta el desafío poder trabajar sobre la misma empleando técnicas como ser análisis de correspondencias y correlación múltiples, como también estadísticas descriptivas mediante el empleo de software estadístico RStudio.

## PREGUNTAS DE INVESTIGACIÓN

¿Existen correspondencias entre las variables categóricas analizadas?

¿Cuáles son aquellas variables más explicativas en el fenómeno de ingresos salariales?

¿Existen diferencias en función del sexo de los individuos?

¿Acaso se cumple la premisa de que, “a mayores años de estudio mejores ingresos”?

## OBJETIVO GENERAL

El objetivo es obtener información estadística sobre las características de ocupación y empleo de la muestra bajo estudio, así como información sociodemográfica y económica, con el fin de analizar la estructura laboral y ocupacional de la misma

## OBJETIVOS ESPECÍFICOS

Aplicar estadística descriptiva, correspondencias múltiples, clústers, regresión lineal múltiple en busca de realizar un análisis más exhaustivo.

Lograr captar información sobre el tamaño, composición y distribución de la fuerza de trabajo y de los niveles de participación en la actividad económica, sus características y de los niveles de ingresos, estudios y cargo y nivel de salarios de las personas.



## MARCO METODOLÓGICO

El trabajo se lleva a cabo empleando una investigación del tipo cuantitativo, no experimental de corte transversal dado que se emplea una muestra sobre un trimestre específico.

Se toma la base de datos ENOE 2021 (encuesta nacional de ocupación y empleo) perteneciente al INEGI (Instituto nacional de estadística y geografía, México), la cual cuenta con más de 10mil registros con 11 variables que se desarrollarán más adelante.

Se pretende llevar a cabo un tipo de investigación exploratorio y explicativo con la intención de abordar la situación laboral de los sujetos encuestados en el censo en busca de patrones y relaciones entre las variables del DataFrame.

Con respecto a las técnicas de análisis, se emplearán estadística descriptiva, análisis multivariante como ser el análisis de correspondencia múltiple, clústers, regresión lineal múltiple con la implementación del software estadístico R.

## MARCO TEÓRICO

### **Análisis multivariante** (Hair, Anderson, Tatham, & Black, 1999)

Gran parte de esta creciente comprensión y pericia en el análisis de datos ha venido a través del estudio de la estadística y de la inferencia estadística.

Igualmente, importante, sin embargo, ha sido el dilatado conocimiento y aplicación de un grupo de técnicas estadísticas conocidas como análisis multivariante.

Las técnicas del análisis multivariante están siendo ampliamente aplicadas a la industria,

administración y centros de investigación de ámbito universitario. Por otra parte, pocos campos de investigación o estudio han fracasado en integrar las técnicas multivariantes en su «caja de herramientas» analítica.

los directivos de empresa o los funcionarios de la administración pública, sea cual sea su entorno, que tienen que desarrollar sus conocimientos del análisis multivariante para conseguir una mejor comprensión de los complejos fenómenos de sus ámbitos de trabajo.

Cualquier investigador que examina sólo relaciones de dos variables y que evita el análisis multivariante está ignorando poderosas herramientas que podrían suministrar información potencialmente útil.

Esos métodos hacen posible plantear preguntas específicas y precisas de considerable complejidad en marcos idóneos, lo que posibilita llevar a cabo investigaciones teóricamente significativas y evaluar los efectos de las variaciones paramétricas ocurridas de forma natural en el contexto en que normalmente ocurren. De esta forma, se pueden preservar las correlaciones naturales entre las múltiples influencias sobre el comportamiento y se pueden estudiar estadísticamente los efectos aislados de esas influencias sin provocar el típico aislamiento de esos individuos o variables.

En un sentido amplio, se refiere a todos los métodos estadísticos que analizan simultáneamente medidas múltiples de cada individuo u objeto sometido a



investigación. Cualquier análisis simultáneo de más de dos variables puede ser considerado aproximadamente como un análisis multivariante. En sentido estricto, muchas técnicas multivariantes son extensiones del análisis univariante (análisis de distribuciones de una sola variable) y del análisis bivariante (clasificaciones cruzadas, correlación, análisis de la varianza y regresiones simples utilizadas para analizar dos variables).

### **Análisis de correspondencias** (Hair, Anderson, Tatham, & Black, 1999)

El análisis de correspondencias es una técnica de interdependencia recientemente desarrollada que facilita tanto la reducción dimensional de una clasificación de objetos (por ejemplo, productos, personas, etc.) sobre un conjunto de atributos y el mapa perceptual de objetos relativos a estos atributos. Los investigadores se enfrentan constantemente a la necesidad de «cuantificar datos cualitativos» que encuentran en variables nominales. El análisis de correspondencias difiere de otras técnicas de interdependencia discutidas antes en su capacidad para acomodar tanto datos no métricos como relaciones no lineales.

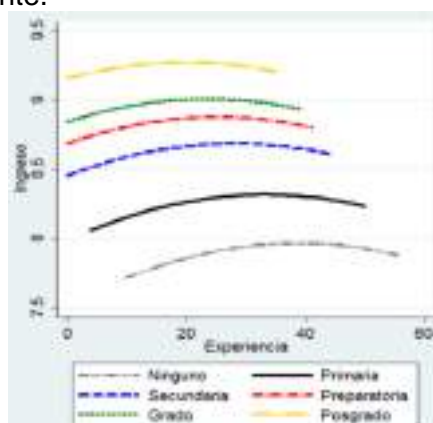
En su forma más básica, el análisis de correspondencias emplea una tabla de contingencia.

### **Ecuación de Mincer** (Mincer, Jacob, 1958)

La ecuación de Mincer es, en economía del trabajo y en economía de la educación, una relación matemática que relaciona el número de años de estudio y el número de años de experiencia con el salario de un individuo en el mercado laboral. Trabajo, nombrado en memoria de Jacob Mincer. Se utiliza para medir la compensación que hacen los agentes entre la búsqueda de una educación superior y la entrada directa al mercado laboral. La ecuación se puede notar con ingresos, ingresos para una persona sin educación o experiencia, el número de años de escolaridad y el número de años de experiencia.

$$\ln W = \ln W_0 + \rho S + \beta_1 X + \beta_2 X^2$$

Dónde las variables tienen los significados siguientes  $W$  son los ingresos  $W_0$  son los ingresos de alguien sin ninguna educación y ni experiencia);  $S$  son los años de escolarización;  $X$  son los años de experiencia laboral potencial en el mercado. Los parámetros  $\rho$ ,  $\beta_1$ ,  $\beta_2$  pueden ser interpretados como los retornos a escolarización y experiencia, respectivamente.



Fuente: ["Función de ingreso de Mincer \(wikimedia.org\)"](https://es.wikipedia.org/wiki/Funci%C3%B3n_de_ingreso_de_Mincer)



### **Análisis clúster** (Hair, Anderson, Tatham, & Black, 1999)

El análisis clúster es una técnica analítica para desarrollar subgrupos significativos de individuos u objetos. De forma específica, el objetivo es clasificar una muestra de entidades (personas u objetos) en un número pequeño de grupos mutuamente excluyentes basados en similitudes entre las entidades. En el análisis clúster, a diferencia del análisis discriminante, los grupos no están predefinidos.

Por consiguiente, se usa la técnica para identificar los grupos.

Habitualmente, el análisis clúster implica al menos dos etapas. La primera es la medida de alguna forma de similitud o asociación entre las entidades para determinar cuántos grupos existen en realidad en la muestra. La segunda etapa es describir las personas o variables para determinar su composición. Este paso puede llevarse a cabo aplicando el análisis discriminante a los grupos identificados por la técnica clúster.

que es la tabulación cruzada de dos variables categóricas. A continuación, transforma los datos no métricos en un nivel métrico y realiza una reducción dimensional (similar al análisis factorial) y un mapa perceptual (similar al análisis multidimensional).

### **Clustering** (García, Molina, & Bustamante, 2018)

También llamadas *técnicas de agrupamiento* o *técnicas de segmentación*, permiten la identificación de tipologías o grupos donde los elementos guardan gran similitud entre sí y muchas diferencias con los de otros grupos. Así se puede segmentar el colectivo de clientes, el conjunto de valores e índices financieros, el espectro de observaciones astronómicas, el conjunto de zonas forestales, el conjunto de empleados y de sucursales u oficinas, etc. La segmentación está teniendo mucho interés desde hace ya tiempo dadas las importantes ventajas que aporta al permitir el tratamiento de grandes colectivos de forma pseudo particularizada, en el más idóneo punto de equilibrio entre el tratamiento individualizado y aquel totalmente masificado.

Las herramientas de *clustering* pueden estar basadas en múltiples tipos de técnicas, como estadística, empleo de algoritmos matemáticos o generación de reglas o de redes neuronales para la trata miento de registros. Para otro tipo de elementos a agrupar o segmentar, como texto y documentos, se usan técnicas de reconocimiento de conceptos.

Esta técnica suele servir de punto de partida para después hacer un análisis de clasificación sobre los clústeres.

La principal característica de esta familia de técnicas es la utilización de una medida de similaridad que, en general, está basada en los atributos que describen a los objetos, y se define usualmente por proximidad en un espacio multidimensional. Para datos numéricos, suele ser preciso preparar los datos antes de realizar *datamining* sobre ellos, de manera que en primer lugar se someten a un proceso de estandarización.

### **Regresión múltiple** (Hair, Anderson, Tatham, & Black, 1999)

La regresión múltiple es el método de análisis apropiado cuando el problema del investigador incluye una única variable métrica dependiente que se supone está relacionada con una o más variables métricas independientes. El objetivo del análisis de la regresión múltiple es predecir los cambios en la variable dependiente en respuesta a cambios en varias de las variables independientes.



Este objetivo se consigue muy a menudo a través de la regla estadística de los mínimos cuadrados.

La regresión múltiple es útil siempre que el investigador esté interesado en predecir la cantidad o la magnitud de la variable dependiente. Por ejemplo, se puede hacer la predicción de los gastos mensuales de cenar fuera de casa (variables dependientes) con información referente a la renta familiar, su tamaño y la edad del cabeza de familia (variables independientes). De la misma forma, el investigador puede intentar predecir las ventas de una compañía a partir de información sobre sus gastos en publicidad, el número de vendedores y el número de tiendas que distribuyen sus productos.

## **DESARROLLO**

El trabajo se realizó sobre la base de datos proporcionada de manera pública por (INEGI) Instituto Nacional de Estadística y Geografía, es uno de los órganos constitucionales autónomos de México con gestión, personalidad jurídica y patrimonio propios, responsable de normar y coordinar el Sistema Nacional de Información Estadística y Geografía.

Se encarga de realizar los censos nacionales; integrar el sistema de cuentas nacionales y estatales (es decir, el flujo de producción, consumo y distribución de la actividad económica); y, desde 2011, elabora los índices nacionales de Precios al Consumidor, e Índice Nacional de Precios al Productor.

Específicamente se trabajó sobre la sección de Encuesta Nacional de Ocupación y Empleo (ENOE), la cual busca medir el salario mensual (ingreso mensual) de los mexicanos y es. Esta encuesta se realiza cada trimestre. Se tomó una muestra de la encuesta realizada en el cuarto trimestre de 2021.

### **Aclaraciones sobre las variables**

A la hora de definir la variable sexo, el instituto informa lo siguiente: La cuestión de si en las estadísticas de hombres y mujeres se debe referir a “género” o a “sexo” ha quedado esclarecida en los foros internacionales que abordan este tema. Según la oficina de Estadísticas de Suecia, la palabra “sexo” hace referencia a las diferencias biológicas entre hombres y mujeres; mientras que “género” es una construcción social (normas, costumbres y prácticas) y una condición de las diferencias entre los sexos y de las relaciones sociales entre hombres y mujeres, ya que la identidad social de género depende de factores ideológicos, históricos, culturales, religiosos, étnicos y económicos. En México, para las estadísticas sobre establecimientos y empresas, particularmente las de los Censos Económicos, cuando presentan el personal ocupado que trabaja en las unidades económicas, se asume que la desagregación de esa variable en hombres y mujeres, es una distinción de sexo. El eje de este producto son los hombres y las mujeres, desde el punto de vista de su inserción en los puestos de trabajo que ocupan en las unidades económicas. Por lo anterior, es importante señalar que la desagregación de hombres y mujeres que se hace en este documento está en



**XII Muestra Académica de Trabajos de Investigación de la Licenciatura en Administración**

función de las diversas categorías de personal ocupado total que se distinguen en los Censos Económicos.

El data frame puntualmente posee 10.280 registros con 11 variables las cuales son:

1. **Estado**, el cual indica el área geográfica de residencia del sujeto entrevistado.
2. **Sex**, discriminando por sexo (mujer, hombre).
3. **Edad**.
4. **Pos\_ocu**, la misma clasifica al individuo según sea:
  - a. “Rabajador subordinado y remunerado”;
  - b. “Empleador”; o
  - c. “Trabajador por cuenta propia”.
5. **Ing\_salarios**: indica el rango en el que se encuentra el salario percibido en comparación con el salario mínimo. Por lo cual se tienen las siguientes categorías:
  - a. Hasta un salario mínimo.
  - b. Más de 1 hasta 2 salarios mínimos.
  - c. Mas de 2 hasta 3 salarios mínimos.
  - d. Más de 3 hasta 5 salarios mínimos.
  - e. Más de 5 salarios mínimos.
6. **Niv\_edu**: indica el nivel educativo:
  - a. Primaria incompleta
  - b. Primaria completa
  - c. Secundaria completa
  - d. Medio superior y superior
7. **Anios\_esc**: cantidad de años de estudios, en la muestra analizada el intervalo de los sujetos tiene un mínimo de 0años y máximo de 22años.
8. **Hesocup**: la cual indica las horas semanales destinadas al trabajo.
9. **Ingreso\_mensual**: el salario percibido en moneda local de forma mensualizada.
10. **Num\_trabajos**: indica la cantidad de trabajos que posee la persona siendo el máximo de trabajos de la muestra de dos empleos.
11. **Tipo\_empleo**: categorizando por Formal o Informal.

A continuación, se muestra una captura de la base de datos para tener una mejor noción de la misma.

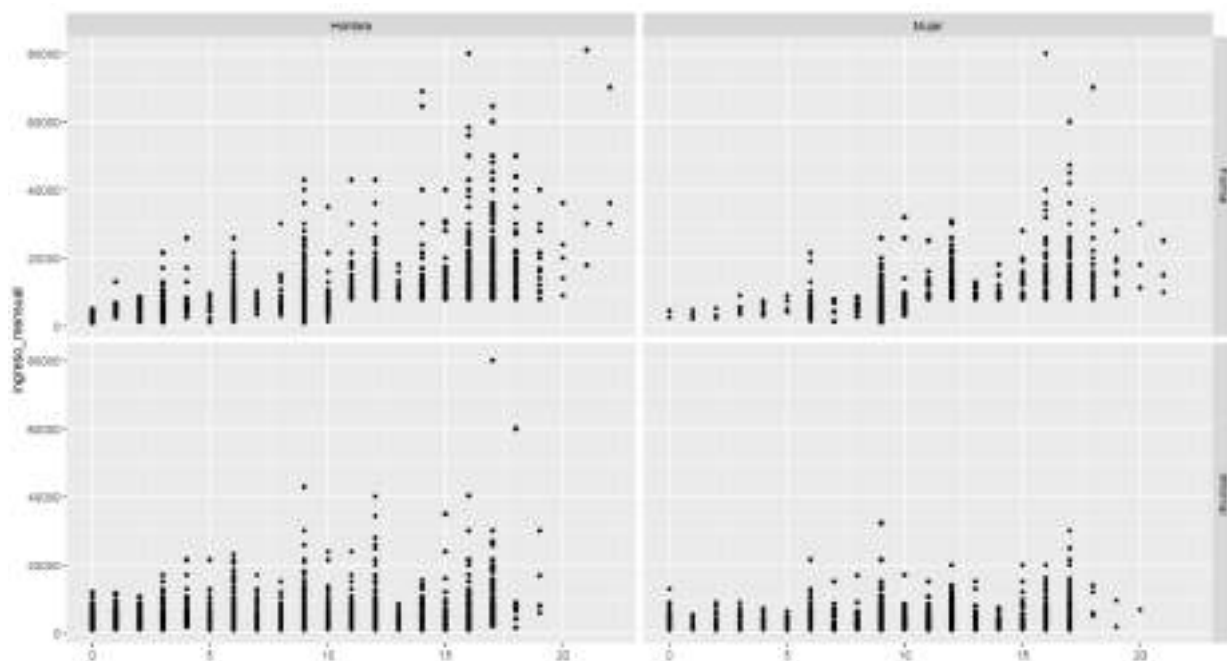
estado	sex	edad	pos_ocu	ing_salarios	niv_edu	anio_esc	hesuca	ingreso_mensual	num_trabajos	tipo_empleo
2	Trabajo	Mujer	36	Trabajadores subordinados y remunerados	Más de 1 hasta 2 salarios mínimos	Secundaria completa	21	45	2070 Uta	Informal
3	Eranga	Hombre	51	Trabajadores subordinados y remunerados	Más de 3 hasta 5 salarios mínimos	Medio superior y superior	17	66	12800 Uta	Formal
4	Salaco	Hombre	25	Trabajadores subordinados y remunerados	Más de 3 hasta 5 salarios mínimos	Medio superior y superior	25	48	12800 Uta	Formal
5	Tobacco	Mujer	50	Trabajadores subordinados y remunerados	Más de 1 hasta 2 salarios mínimos	Secundaria completa	9	46	3870 Uta	Informal
6	Tobacco	Mujer	41	Trabajadores por cuenta propia	Hasta un salario mínimo	Medio superior y superior	17	5	1800 Uta	Informal
7	Nuevo León	Mujer	36	Trabajadores subordinados y remunerados	Más de 1 hasta 2 salarios mínimos	Secundaria completa	9	50	4885 Uta	Formal



## ANÁLISIS DESCRIPTIVO

En el siguiente gráfico de dispersión se trabajó con las variables **sexo**, **ingresos**, y **tipo de empleo**. Se puede apreciar las diferencias salariales entre los hombres y las mujeres, donde en términos generales los ingresos de los hombres son más altos sobre todo cuando el tipo de empleo es **Informal**

Gráfico de dispersión matricial empleando variables sexo, ingresos mensuales, tipo de empleo



Fuente: elaboración propia

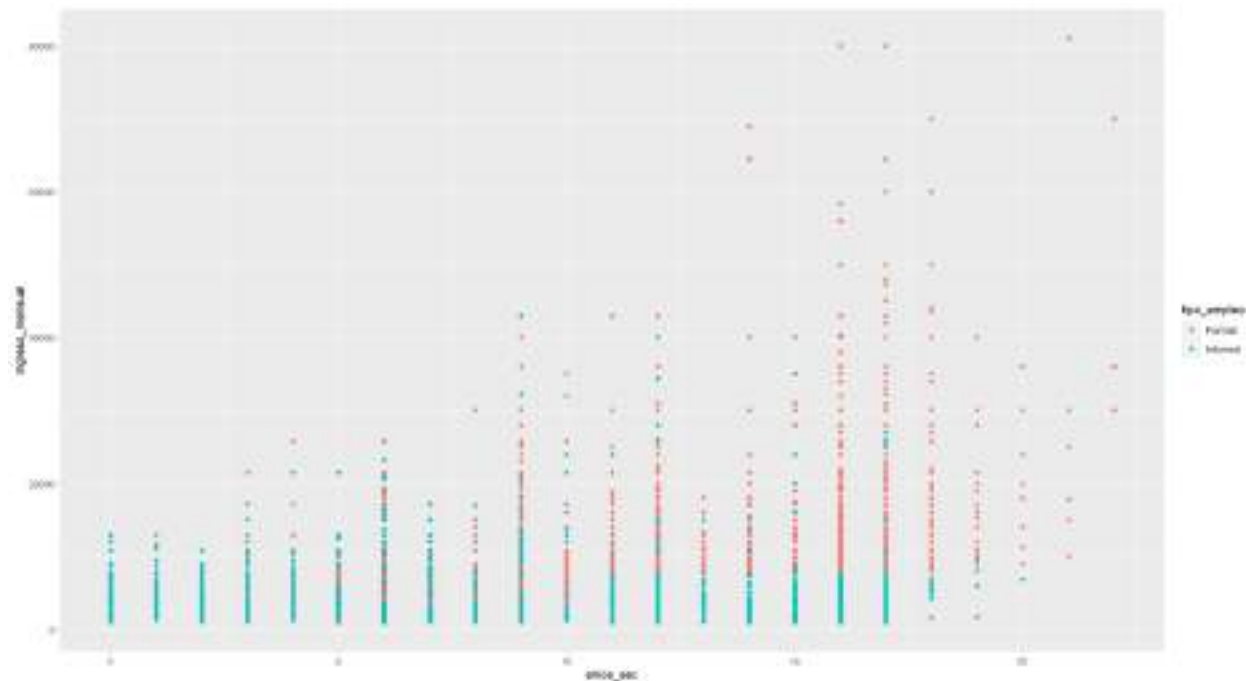
Al realizar un gráfico de dispersión empleando las variables años de estudios, tipo de empleo e ingresos mensuales, se aprecia una correlación positiva entre los años de estudio y los ingresos y una relación negativa entre tipo de empleo “informal” y años de estudio, es decir, a mayores años de estudio se aprecia mayores ingresos, a su vez a mayores años de estudio el tipo de empleo informal va decreciendo tomando más preponderancia el empleo formal. El fenómeno se corresponde tal como lo predicen las ecuaciones de Mincer.





**XII Muestra Académica de Trabajos de Investigación de la Licenciatura en Administración**

Gráfico de dispersión empleando ingresos mensuales, años de escolaridad, tipo de empleo.

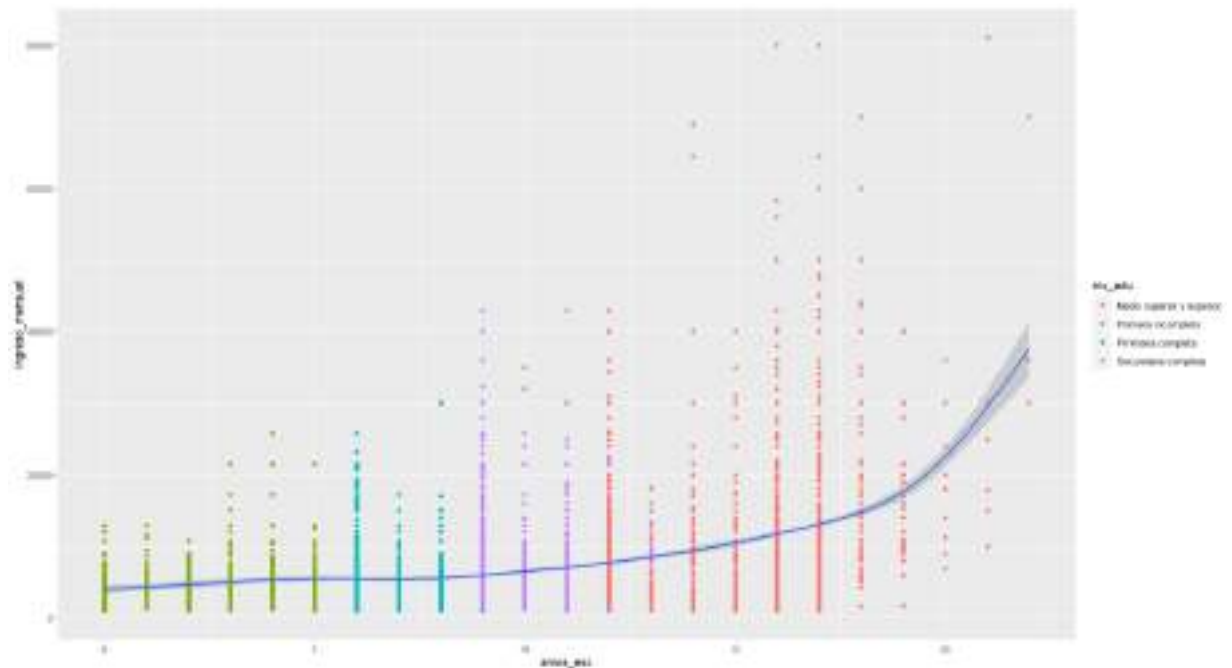


Fuente: elaboración propia

Si incorporamos la categorización en base al nivel educativo y una línea de tendencia, se evidencia la correlación entre los años de estudio y el nivel educativo y la tendencia positiva a poseer mayores ingresos mensuales donde la curva comienza a manifestarse exponencialmente a partir de los 12 años de estudio.

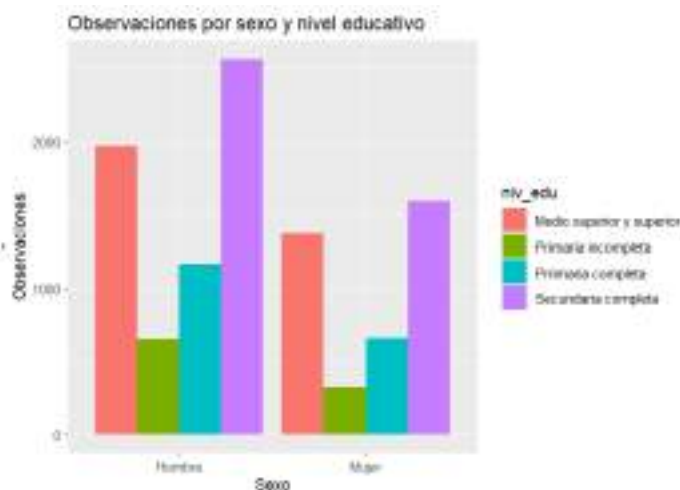


Gráfico de dispersión con línea de tendencia empleando las variables Ingresos mensuales, Nivel educativo, años de educación.



Fuente: elaboración propia

En la siguiente gráfica se representa el nivel de estudio discriminado en función del sexo



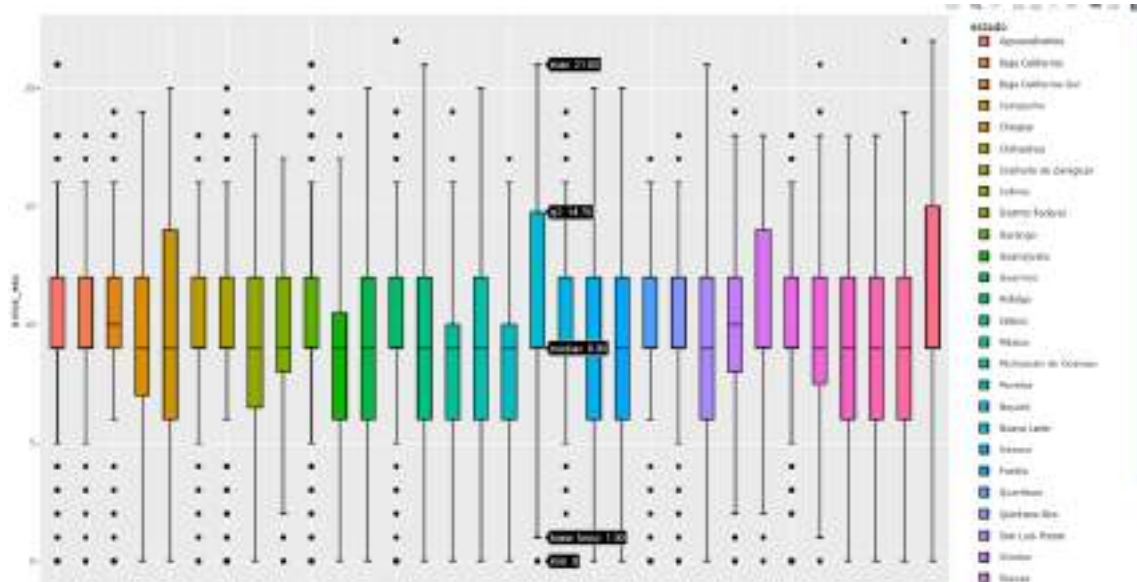
Fuente: elaboración propia

Con la intención de visualizar cuáles son aquellos estados que poseen mayor cantidad de años de estudio, se empleó la función “ggplotly” que permite visualizar

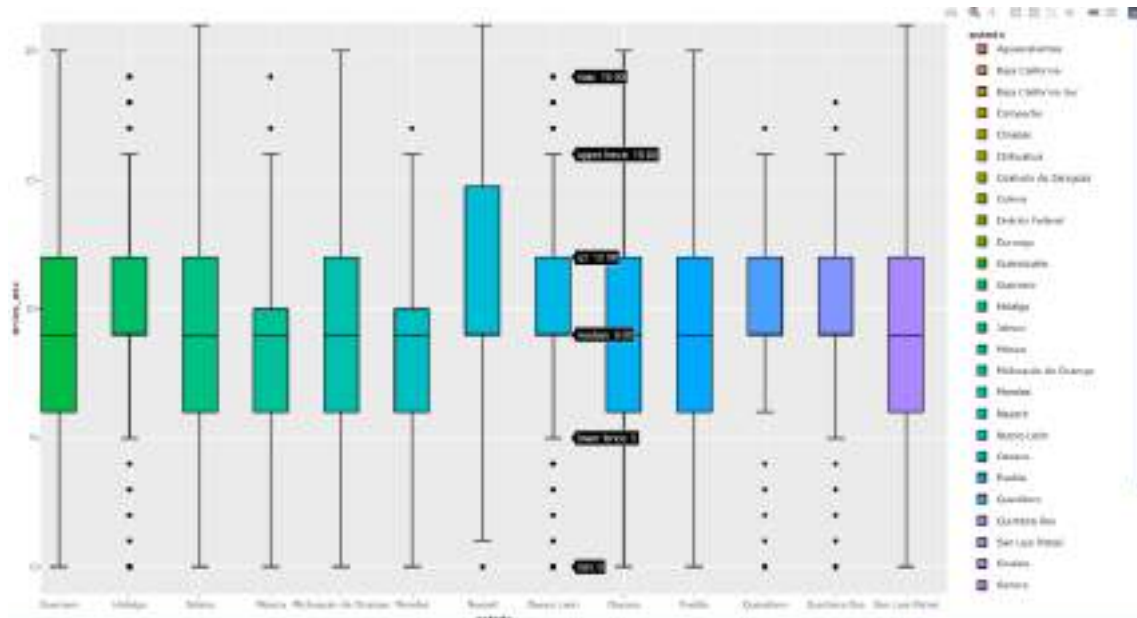


claramente dentro de los diagramas de caja representados la información pertinente a la mediana, los cuartiles, los máximos y mínimos y si existiesen outliers.

### BOXPLOTS POR ESTADO



Fuente: elaboración propia



Fuente: elaboración propia

Analizando los estados, se aprecia que dentro de los que poseen mayores años de estudio lo encabezan de mayor a menor cantidad de años: la Ciudad de México, Nuevo León y Sonora.



Mientras que los estados que presentan la menor cantidad de años de estudios de sus habitantes, ordenados de menor a mayor son los siguientes:

Chiapas, seguido de Oaxaca, Guerrero, Michoacán, Veracruz, Guanajuato, Zacatecas, Hidalgo, Yucatán y San Luis Potosí.

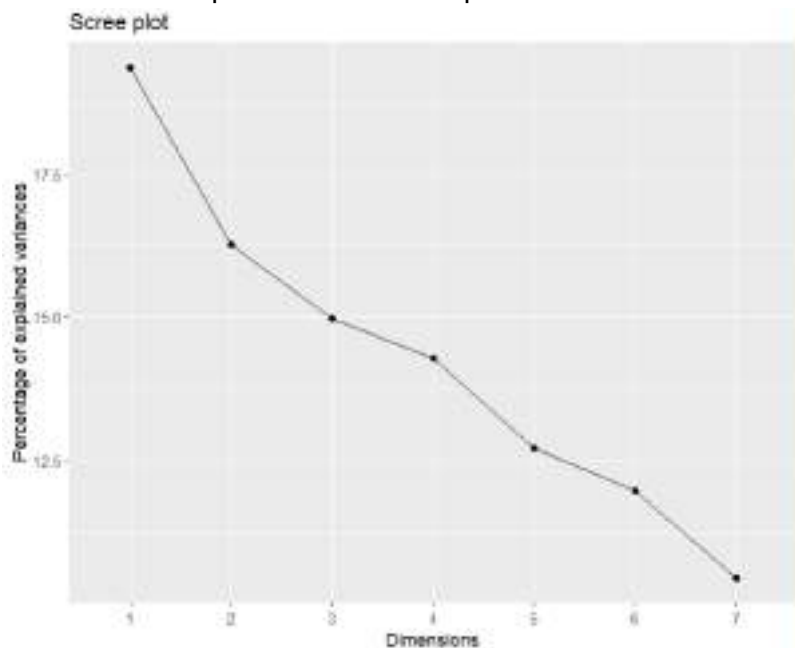
## APLICACIÓN DE ANÁLISIS DE CORRESPONDENCIAS MULTIPLE (MCA)

MCA empleando las variables sexo, posición ocupada y categoría de ingresos

	Dim 1	Dim 2	Dim 3	Dim 4	Dim 5
sex	0.3685904	0.003683828	0.1475355	2.458417e-31	0.1845988
pos_ocu	0.4103937	0.577227422	0.1228955	1.561746e-25	0.2998977
ing_salarios	0.5750452	0.557411953	0.7770748	1.000000e+00	0.4057946

Fuente: elaboración propia

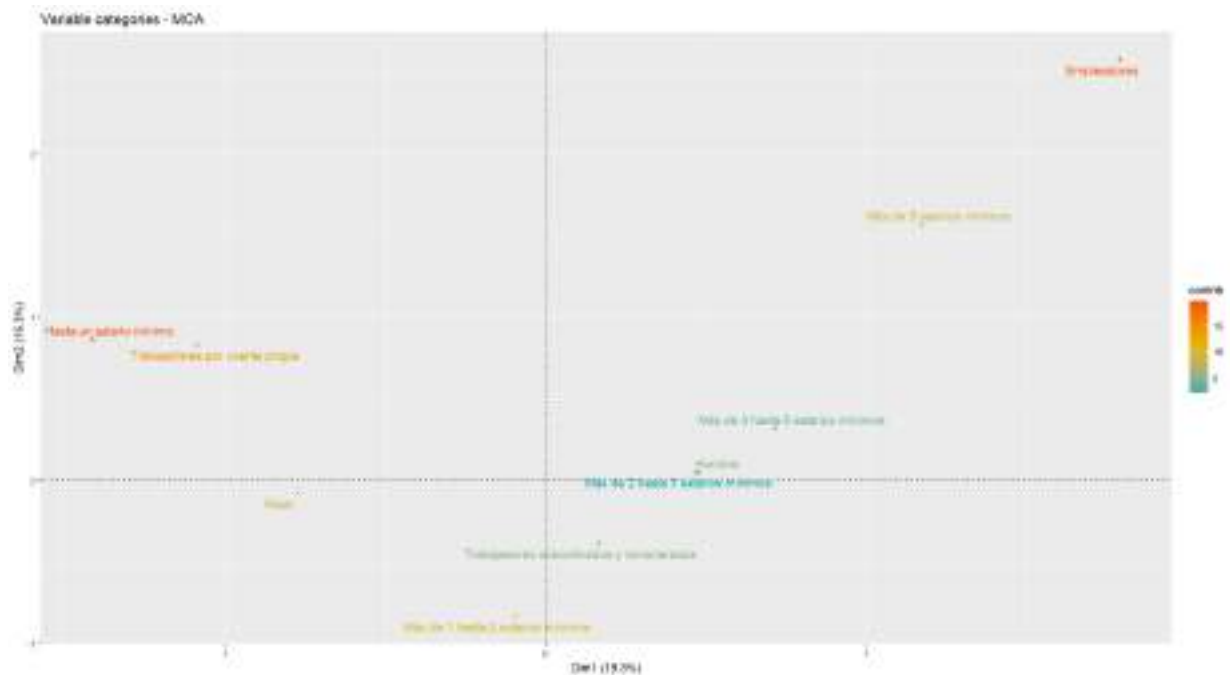
Con el empleo de las variables mencionadas, se procede a el cálculo de la significancia que aportan las dimensiones generadas. A continuación, se aprecia el porcentaje de explicación de las variaciones, Con las dimensiones 1 y 2 se acumula una explicación del 40% aproximadamente.



Fuente: elaboración propia



Figura: MCA1



Observaciones de la figura MCA1, recordemos que las variables en juego son el sexo, la ocupación en el empleo y el rango de ingreso salarial mensual. Las dimensiones bajo estudio explican alrededor del 40% del fenómeno.

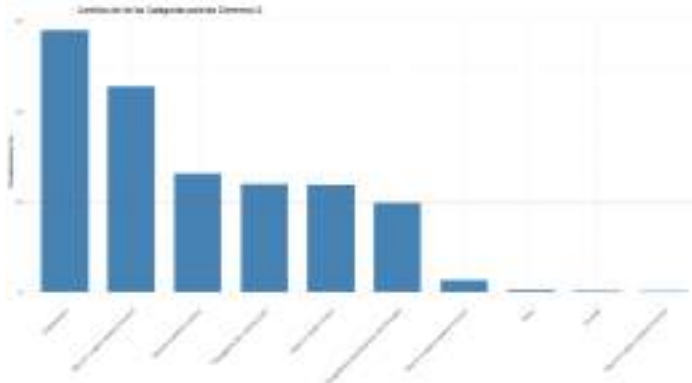
Podemos apreciar la correspondencia que existe entre los individuos que son empleadores quienes poseen mayor tendencia a tener un rango de ingresos superior al de 5 salarios mínimos.

Al otro extremo se encuentran aquellos trabajadores por cuenta propia que se corresponden con los ingresos más bajos, es decir, "hasta un salario mínimo".

Si bien cuesta realizar otras correspondencias dado que las demás variables intermedias en base al mapa de calor, no tienen grandes tasas de contribución a las dimensiones. Sin embargo, observamos que los hombres se encuentran más cercanos a percibir más de dos salarios mínimos con mayor proximidad a la posición ocupada en el empleo "subordinados y remunerados". Mientras que las mujeres se encuentran en el punto medio entre la posición ocupada como trabajadores por cuenta propia y subordinados remunerados y el rango de ingresos entre 1 hasta 3 salario mínimos.

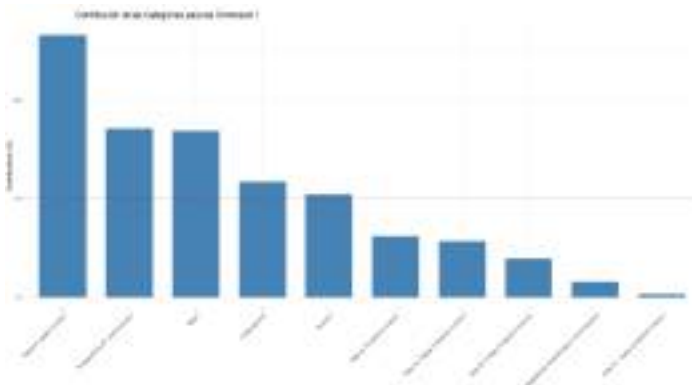


### Contribuciones de las variables a la dimensión 2



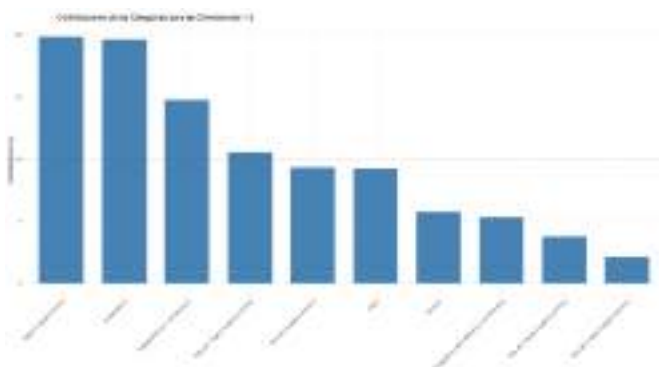
Fuente: elaboración propia

### Contribución de las variables a la Dimensión 1



Fuente: elaboración propia

### Contribución de las variables a ambas dimensiones



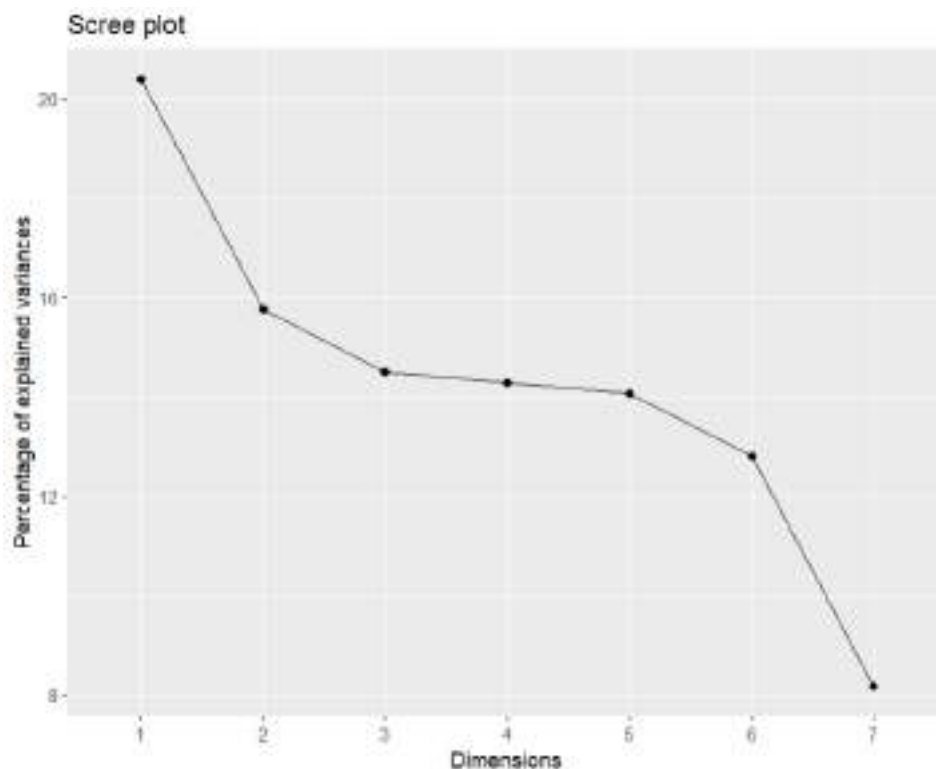
Fuente: elaboración propia



## MCA empleando variables ingresos y nivel educativo

	Dim 1	Dim 2	Dim 3	Dim 4	Dim 5
ing_salarios	0.7137726	0.5520722	0.5074478	1.000000e+00	0.4925522
niv_edu	0.7137726	0.5520722	0.5074478	2.906889e-24	0.4925522

Fuente: elaboración propia

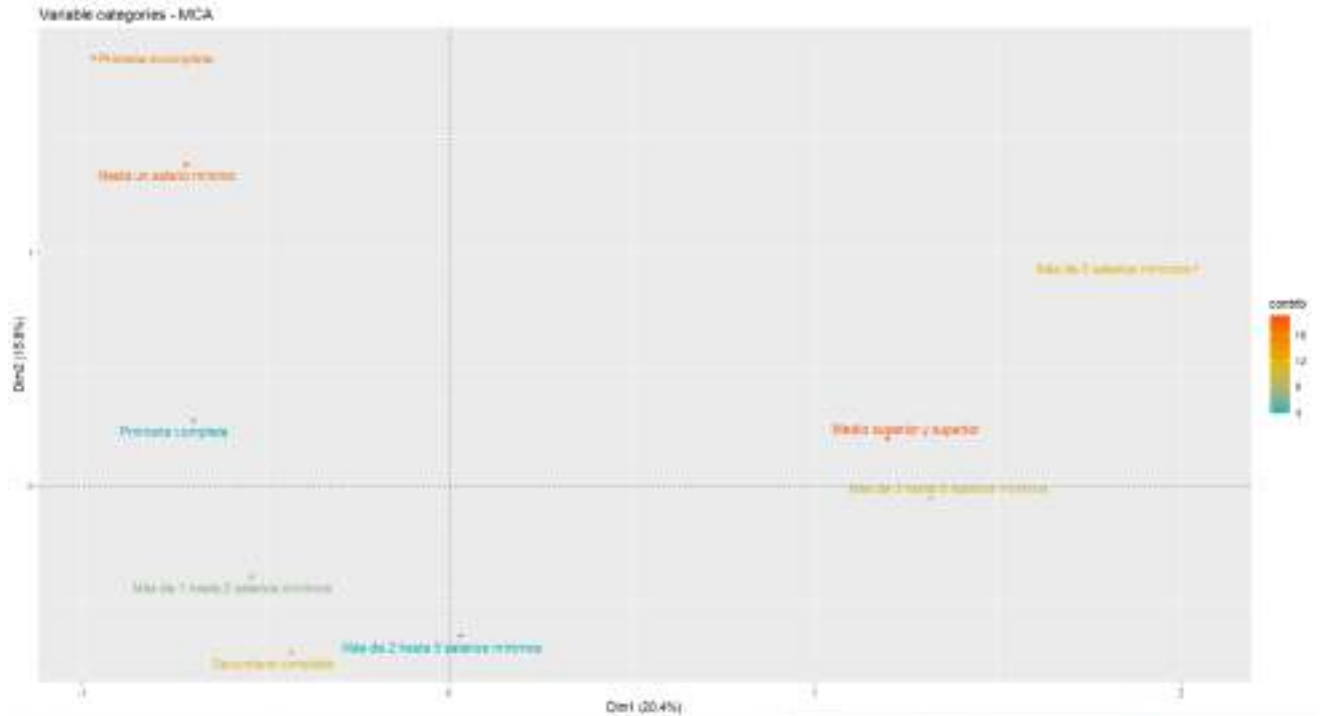


Fuente: elaboración propia

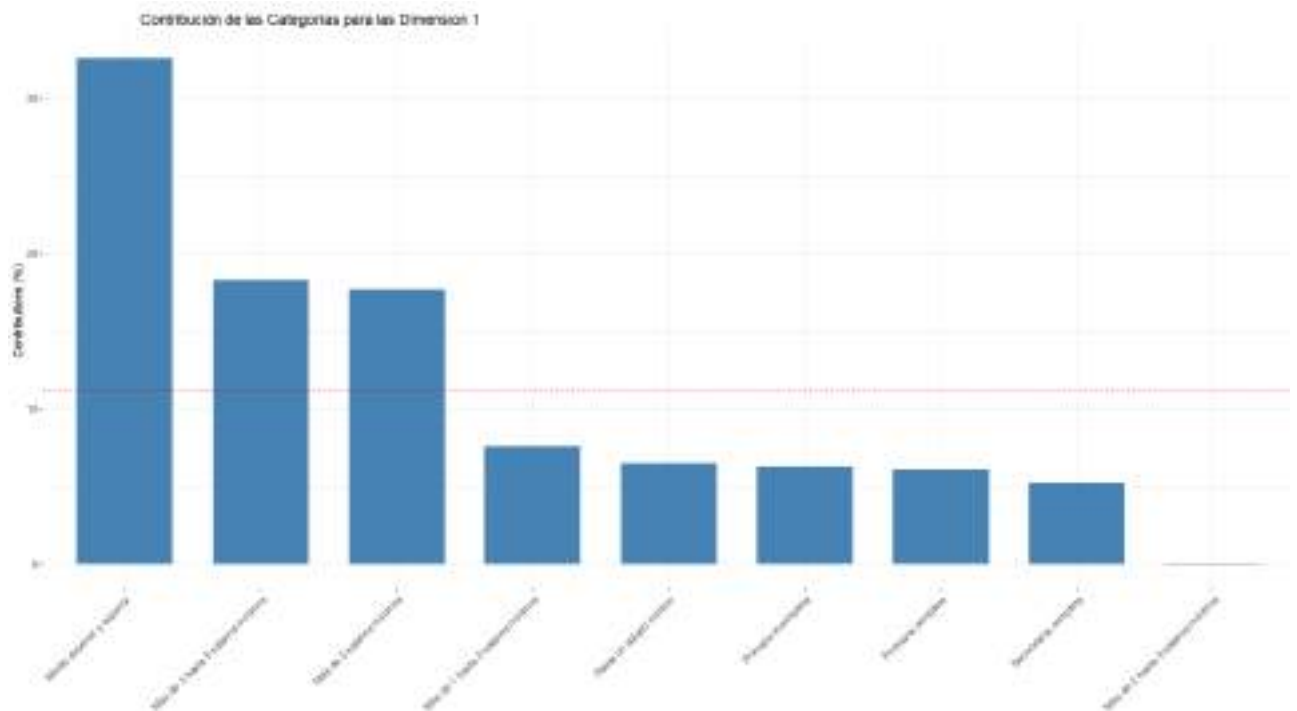
En esta ocasión se analizan las variables ingresos con el nivel educativo, vemos como se corresponde con los estudios de Mincer al respecto del nivel de estudios con el salario esperado, apreciándose una correspondencia clara entre el nivel de estudio medio-superior y superior y los ingresos de 3 a 5 salarios mínimos y superiores a 5 salarios mínimos.



**XII Muestra Académica de Trabajos de Investigación de la Licenciatura en Administración**



Fuente: elaboración propia

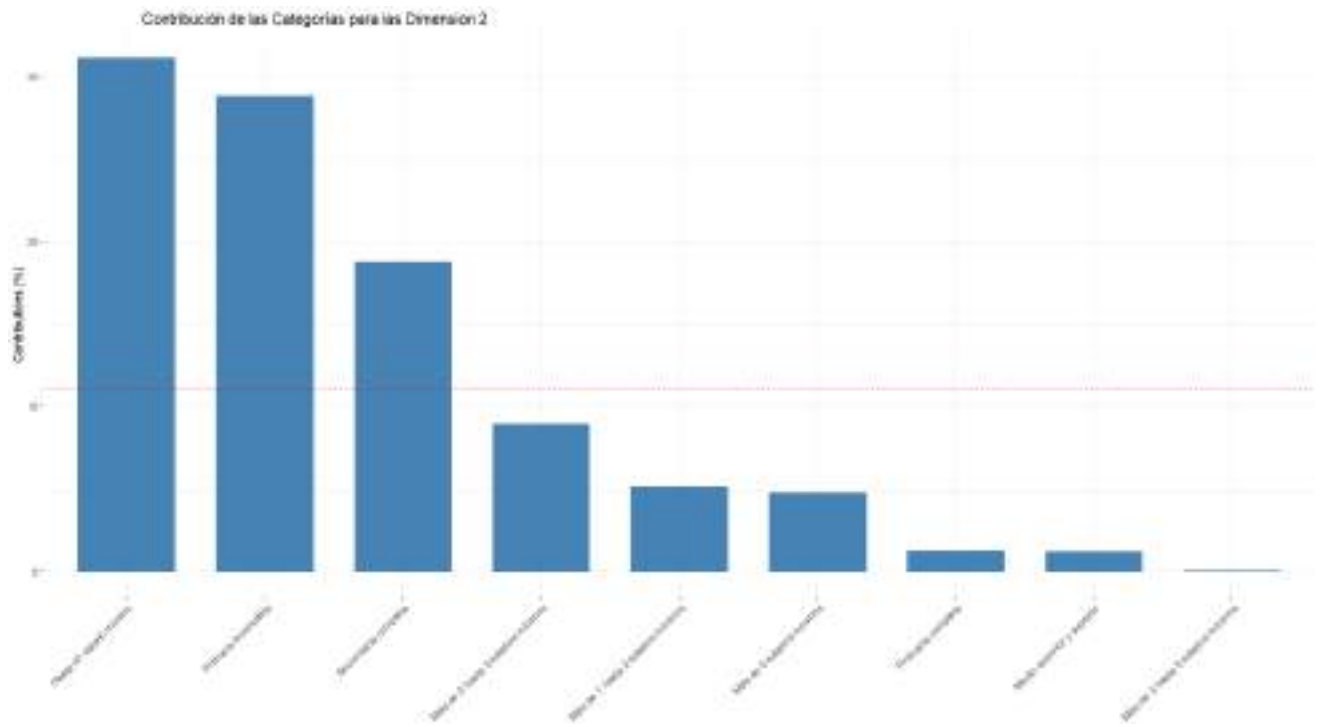


Fuente: elaboración propia

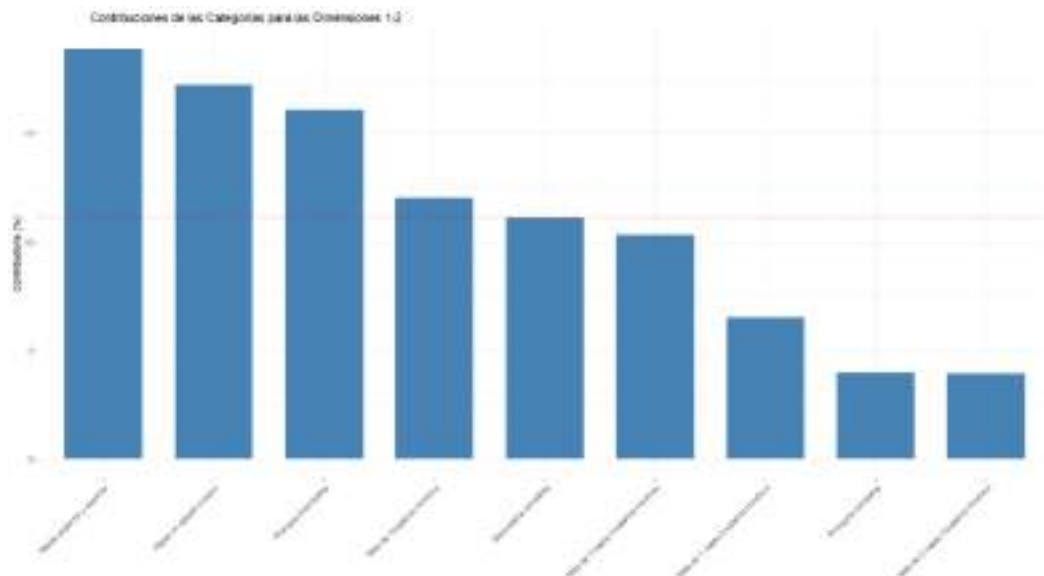




**XII Muestra Académica de Trabajos de Investigación de la Licenciatura en Administración**



Fuente: elaboración propia



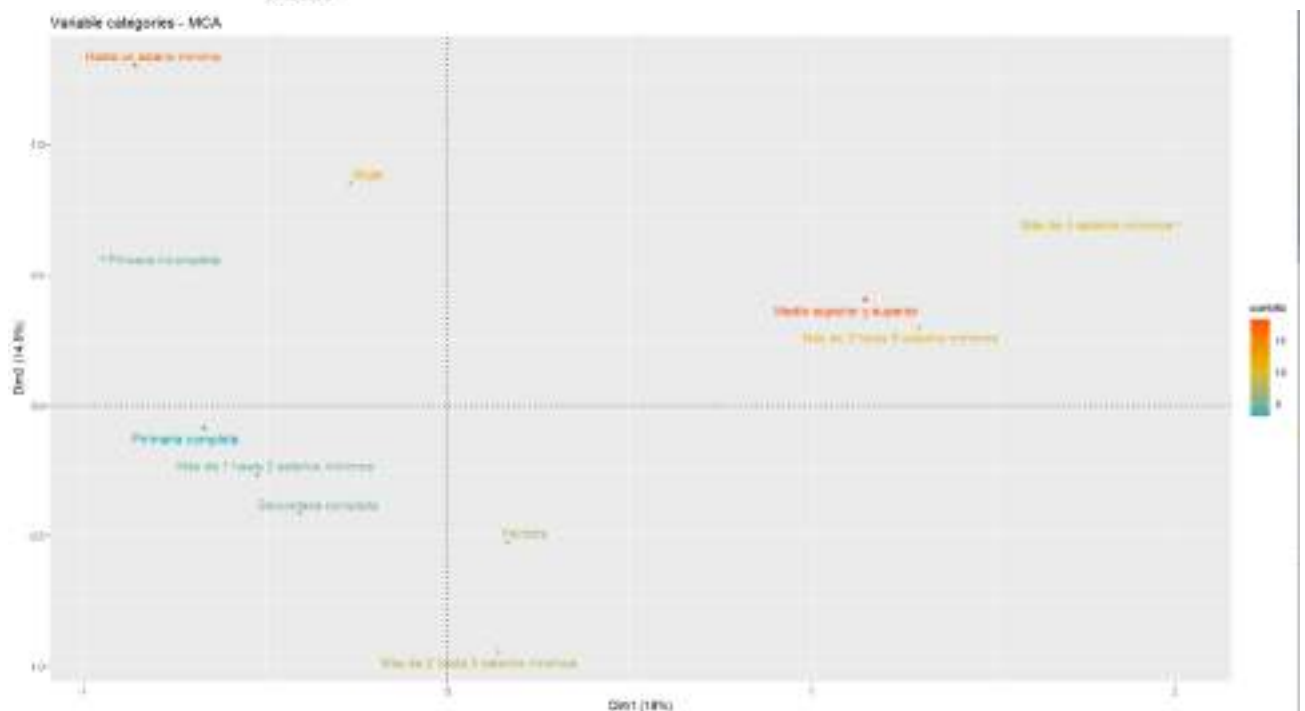
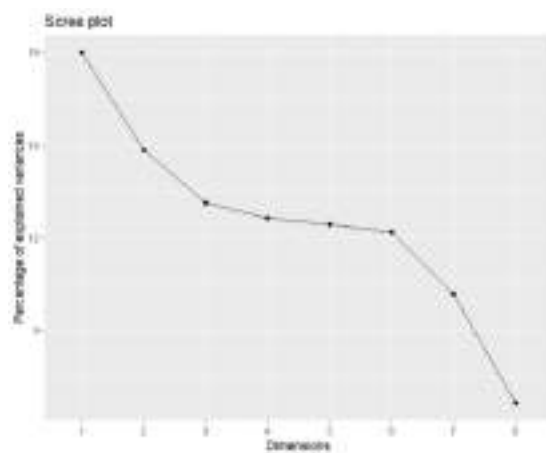
Fuente: elaboración propia



Si procedemos a agregar la variable sexo:

### MCA con variables SEXO INGRESOS y NIVEL EDUCATIVO

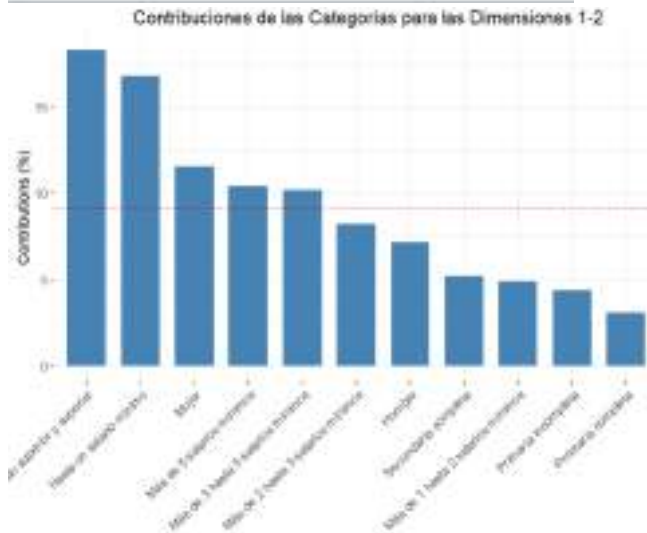
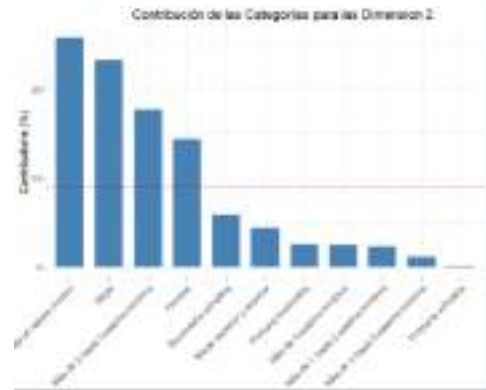
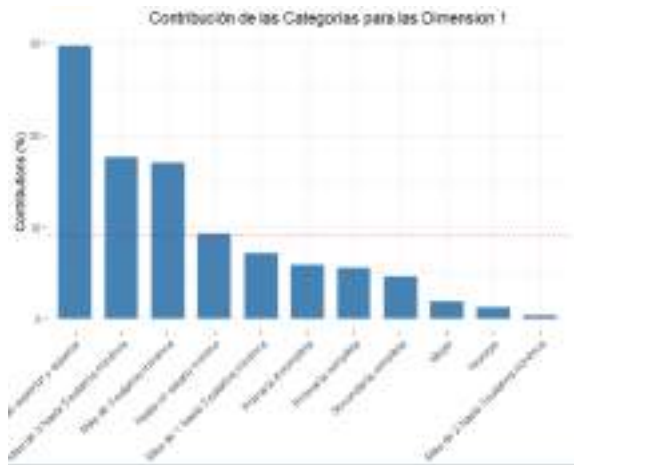
	Dim 1	Dim 2	Dim 3	Dim 4	Dim 5
sex	0.04401985	0.4472590	0.1633222	0.01442663	0.0178647
ing_salarios	0.73686660	0.5853750	0.3405291	0.43235279	0.6391154
niv_edu	0.65800100	0.1537156	0.5461953	0.56329740	0.3365933



Esta vez con la variable sexo junto a nivel de ingresos y nivel educativo, se aprecia que en la muestra analizada de la población, el sexo femenino se encuentra en el segundo cuadrante junto a salarios bajos y niveles bajos de educación, en cuanto al sexo masculino, se aprecia que se encuentra posicionado en un nivel intermedio mas cercano al nivel educativo “secundaria completa” y con salarios cercanos a la categoría superior a 2 salarios mínimos.



**XII Muestra Académica de Trabajos de Investigación de la Licenciatura en Administración**





## CLUSTER

Para el empleo de clústers se procedió a analizar los estados extremos en cuanto a los años de educación obtenidos previamente en los análisis descriptivos. Es decir, analizaremos el estado de México siendo este el que poseía mayores años de estudio y el estado de Sonora que contaba con los menores años de estudio.

En primera instancia se realizó la limpieza de aquellos datos atípicos que en este caso en base a las distancias de Mahalanobis se determinaron los atípicos como aquellos que superaban el valor de 13,27.

State	Mahalanobis Distance
[1] México	14.5799611
[10] México	16.8813084
[19] México	19.2466284
[127] México	13.2368727
[190] México	23.9191704
[208] México	25.5574382

Luego de la detección de 6 outliers, se procedió a eliminar los mismos.

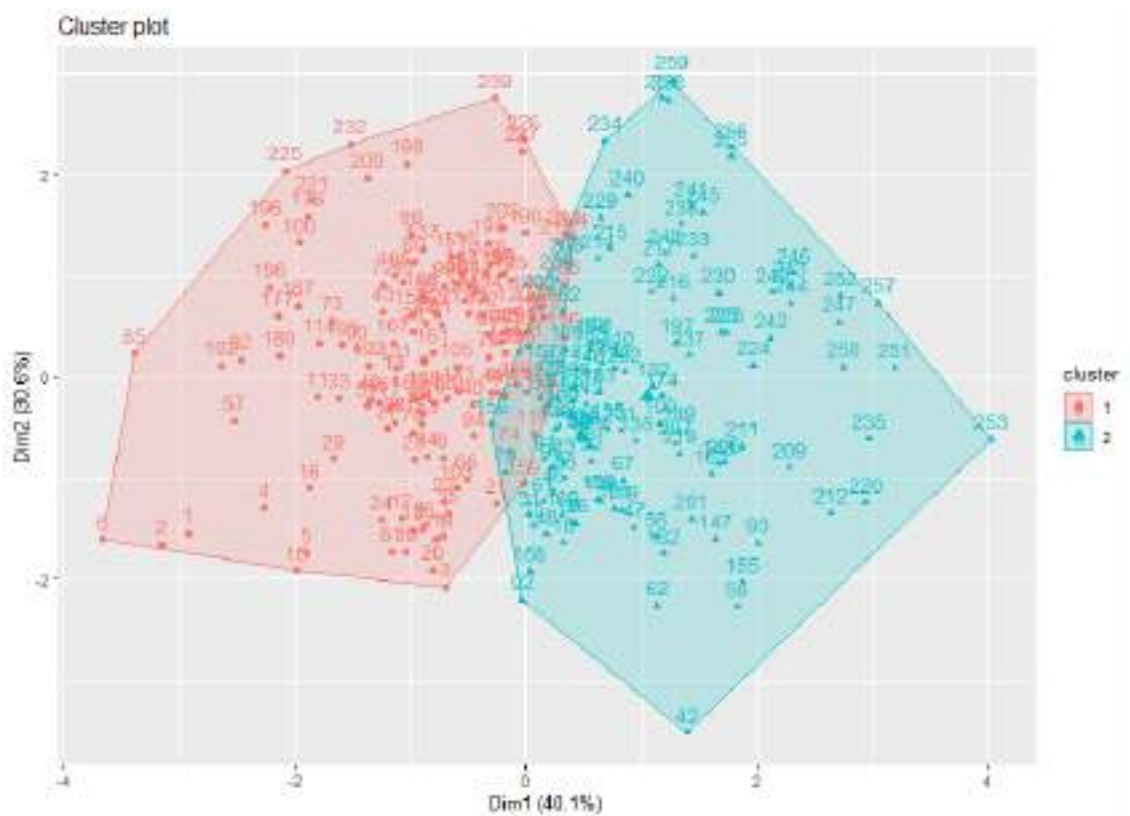
State	Year	Age	Income	Unemployed
1 México	23	9	86	5000
2 México	19	12	88	3200
3 México	38	16	89	4000
4 México	43	11	84	9200
5 México	42	9	84	5000
7 México	35	11	80	2500
8 México	30	9	80	6000
9 México	40	9	71	4100
10 México	32	8	71	3000
11 México	30	9	71	3470
12 México	19	9	71	5480
13 México	48	9	71	7740
14 México	41	9	71	3720
15 México	38	8	71	7740
16 México	32	9	71	4300
17 México	48	9	71	10000
18 México	19	9	71	8500
19 México	37	9	71	21200
20 México	43	12	70	9100
21 México	55	9	70	4300
22 México	38	9	70	7500
23 México	48	9	70	8800
24 México	35	9	70	5100
25 México	53	9	66	4300
26 México	55	9	66	3000



**XII Muestra Académica de Trabajos de Investigación de la Licenciatura en Administración**

En cuanto a la matriz de distancias se eligieron los números de cluster teniendo en cuenta los indicadores óptimos de número de clusters con método de distancia euclídea y empleando en el scrip: wss, gap\_stat, y silhouette. El mismo correspondía con el número óptimo de dos (2) conglomerados. De esta manera se logra una explicación dimensional del 70%

Clúster elaborado con el estado de México:



Fuente: elaboración propia

Se conformaron 2 grupos

Media de Cada Variable por Conglomerado, Utilizando los datos originales

cluster	edad	años_esc	hrsocup	ingreso_mensual
1	35.48571	14.657143	43.47819	12284.781
2	41.52174	8.088379	42.08938	5585.012

En este caso se aprecia que, de los 2 conglomerados, las variables más representativas son el ingreso junto a los años de estudio, dándose una correlación positiva entre las mismas, a su vez, el primer grupo representado en color rojo del lado izquierdo del gráfico, el análisis de los promedios es, para edad de 35 años, en años de escolaridad 15 años, horas de ocupación de 43, y un ingreso alrededor de los \$12.285 moneda local.



**XII Muestra Académica de Trabajos de Investigación de la Licenciatura en Administración**

Mientras que el segundo grupo conformado por los promedios son los siguientes, 42 años de edad, 8 años de escolaridad, 42 hs de ocupación laboral, y salarios en promedio de \$5585.

**Comparación con el estado de Sonora:**

El mismo procedimiento empleado en México se llevó a cabo con el estado Sonora a fin de apreciar los contrastes posibles dado que ambos representan los máximos y mínimos en el ranking de los años de estudio respectivamente.

Por lo que, siguiendo el mismo proceso mencionado anteriormente, la matriz de distancia de Mahalanobis arrojada con los datos atípicos es la siguiente:

```

> print(CO2)
 [1] 6.34636507 5.81679738 6.31995574 8.70763183 5.02321097 5.76506861 6.92811165 7.70536320
 [9] 5.02368772 4.20682531 4.65246252 7.56740041 2.30371747 2.48852100 2.34267923 4.32067376
 [17] 7.99644407 2.72437737 1.90221608 6.15568073 6.42349196 8.14088106 4.86844730 3.47631126
 [25] 3.83501325 4.04458007 1.46678868 2.42459746 0.66622349 3.56030419 3.82585010 3.62654023
 [33] 0.57589710 1.38417992 0.92547780 2.05201113 1.65744474 0.60192869 4.20194814 2.96691036
 [41] 0.65332336 2.98015628 1.04654293 17.67764448 2.36265027 2.07710156 2.02539556 2.24664640
 [49] 8.12968275 4.54464887 6.96997863 4.33479790 2.45487919 0.72740817 0.61597077 3.07688012
 [57] 8.76734766 6.42424515 11.01549691 3.50969774 4.69532970 1.29227866 2.54475017 2.69010436
 [65] 1.94342436 1.67385734 0.92677179 0.73915893 1.06372421 2.01206795 6.61606542 2.28129854
 [73] 1.07382826 0.79769164 2.62510457 0.32478906 0.31570096 0.26074369 3.06392068 0.55267425
 [81] 7.69490903 0.32130567 1.79450721 0.33917752 1.99155328 0.12947429 2.93389915 1.49737082
 [89] 0.54218545 0.24172868 1.47398496 0.35615750 0.03339049 0.20241554 1.04392720 0.99428572
 [97] 0.45780700 0.95015412 0.91949529 1.98171826 0.78613803 0.29358589 0.72827813 0.47897350
 [105] 3.86870792 0.47189929 0.48587555 1.41646470 1.57575042 2.42072466 3.85716312 2.85295953
 [113] 4.39649730 3.92084197 4.91811238 1.40441267 1.39380701 1.17207537 5.08078235 7.09412901
 [121] 3.47603240 1.51177221 1.93477688 8.73964716 14.57479604 4.36576077 3.11520700 2.46425928
 [129] 2.33020342 15.73254015 2.79335958 10.52859717 1.37072700 0.95557769 1.86853852 1.96809504
 [137] 3.88124729 1.45332552 2.90059054 4.96404189 4.83089203 6.57843627 4.88857389 5.97503055
 [145] 5.85840414 3.27189398 11.20809654 62.39855960 8.14845510 4.44954909 5.07780385 4.48939995
 [153] 8.10449414 7.40649029 7.28185374 0.20781137 7.51727586
  
```

Eliminación de los 4 outliers detectados:

edad	anios_esc	hrsocup	ingreso_mensual
31	17	96	23000
32	12	84	17200
66	6	84	5000
74	9	84	4600
30	12	84	8000
23	17	84	7740
30	12	84	17200
29	16	84	10500



**XII Muestra Académica de Trabajos de Investigación de la Licenciatura en Administración**

Clúster elaborado con el estado de Sonora:



Media de Cada Variable por Conglomerado, Utilizando los datos originales

cluster	edad	anios_esc	hrsocup	ingreso_mensual
1	51.22059	5.970588	43.92647	4340.118
2	29.95506	10.494382	43.76405	5513.292

En esta ocasión las variables más representativas analizando Sonora, son la edad y los años de estudio.

El primer clúster representado en Rojo en el lado derecho del gráfico se encuentra agrupando aquellos individuos que poseen en promedio una edad de 51 años con 6 años de estudio promedio, 44hs ocupadas e ingresos por \$4340 moneda local.

Mientras que el segundo grupo color celeste ubicado en el lado izquierdo del gráfico posee un promedio de edad de 30 años quienes presentan en promedio estudios por más de 10 años, horas de ocupación similares al primer grupo e ingresos promedio de \$5513 los cuales implican en promedio un aumento en salarios mensuales del 30% respecto al primer grupo.

Podríamos especular con que, en la muestra estudiada de Sonora, las personas catalogadas como más jóvenes poseen en promedio el doble de años de estudio y perciben consecuentemente un 30% más en ingresos que el grupo con edades promedio de 51años.



## REGRESIÓN LINEAL MÚLTIPLE

Para la regresión lineal múltiple se emplearon las variables Ingresos, Años de escolaridad, Edad, Horas de ocupación.

Se tomó como variable dependiente a Ingresos dado que es la que queremos explicar y el resto de variables mencionadas como explicativas del fenómeno ingresos.

Antes de aplicar la regresión se llevó a cabo la detección y eliminación de datos atípicos para que no interfieran en el resultado final con interpretaciones erróneas.

Se empleó la herramienta de análisis de datos de Excel obteniéndose lo siguiente:

<b>Estadísticas de la regresión</b>	
Coeficiente de correlación múltiple	<b>0,852678793</b>
Coeficiente de determinación R <sup>2</sup>	<b>0,727061124</b>
R <sup>2</sup> ajustado	<b>0,726910497</b>
Error típico	<b>4750,991529</b>
Observaciones	<b>10266</b>

Al estar ante regresión múltiple el  $R^2$  más significativo para las interpretaciones es el ajustado dado que penaliza el hecho de agregar más variables al modelo. En este sentido, podemos decir que las variables años de estudio, edad y horas de ocupación explican a la variable ingresos en un 73%.

Con respecto al desvío típico significa que al querer realizar estimaciones/predicciones tendremos un error en el resultado esperado de  $\pm$  \$4751.

ANÁLISIS DE VARIANZA					
Descripción	Grados de libertad	Suma de cuadrados	Producto de la cuadrada	F	Valor crítico de F
Regresión	3	6,1700E+11	2,05697E+11	9112,941829	0
Residuos	10263	2,31656E+11	22571900,51		
Total	10266	8,48745E+11			

Con el Valor F elevado en este caso con un valor de 9113, estaríamos en condiciones de decir que el modelo es adecuado, como también al observar el valor crítico de F igual a cero.

(Este estadístico F indica si todos los coeficientes de la regresión, conjuntamente, son distintos de cero. Es decir, indica si el coeficiente que acompaña a “Años de educación”, “edad” y el coeficiente que acompaña a “horas de ocupación” son distintos de cero, que es lo mismo que decir que son conjuntamente significativos. En general, queremos que el estadístico F sea lo más grande posible, o también que el “Valor crítico de F” sea lo más pequeño posible. En este caso, tenemos que el estadístico F elevado y su valor crítico es cero, por lo que sí podemos decir que los coeficientes son conjuntamente significativos.)





**XII Muestra Académica de Trabajos de Investigación de la Licenciatura en  
Administración**

Descripción	Coefficientes	Error típico	Estadístico t	Probabilidad	Inferior 95%	Superior 95%	Inferior 95,0%	Superior 95,0%
Intercepción	0							
edad	16,28	2,45	6,64	0,00	11,47	21,09	11,47	21,09
esc	491,57	9,00	54,61	0,00	473,93	509,21	473,93	509,21
hrsocup	45,27	2,35	19,30	0,00	40,67	49,87	40,67	49,87

Vemos que la variable que más repercute en el resultado final del ingreso mensual es los años de estudio (dado que posee el valor de su coeficiente varias veces superior frente al resto de variables explicativas), seguido de horas de ocupación y finalmente la edad del individuo.

En cuanto al **error típico**, el cual nos indica la variación de la estimación en los coeficientes posee valores acordes para buenas estimaciones en relación a los coeficientes.

A su vez los **estadísticos t** nos indican buena representatividad, como también el **probabilidad** o también conocido como **valor-p** se busca que los mismos sean cercanos a cero y en este caso tenemos los valores iguales a cero. Con respecto a los intervalos de confianza presentan un rango bajo por lo que son bastante aceptables.

Por todo lo anterior nos encontramos en condiciones de aceptar los resultados obtenidos por la regresión lineal múltiple como adecuados para realizar estimaciones.

La ecuación de regresión lineal múltiple quedaría como sigue:

$$Y = 16 * (\text{edad}) + 491,57 * (\text{años de escolaridad}) + 45,27 * (\text{hs de ocupación})$$



## APÉNDICE

### El Scrip del trabajo se presenta a continuación

```
setwd("C:/Users/Martin/Desktop/R/Archivos/3.- GGPlot 2")

library(ggplot2)

library(readxl)

library(tidyr)

library(dplyr)

library(plotly)

library(ggplotly)

enoe<-read_xlsx("mu_enoe.xlsx")

#MAPA DE PUNTOS DE DISPERSION INGRESOS VS AÑOS DE ESTUDIO VS TIPO DE EMPLEO

ggplot(data = enoe) +

  geom_point(mapping = aes(x = anios_esc, y = ingreso_mensual, color=tipo_empleo))

#CLASIFICACIÓN SEGÚN SEXO, TIPO DE EMPLEO, AÑOS DE ESCOLARIDAD E INGRESO MENSUAL

ggplotly(

ggplot(data=enoe)+

  geom_point(mapping = aes(x =anios_esc, y =ingreso_mensual))+

  facet_grid(tipo_empleo~sex))

#LINEA DE TENDENCIA EXPLONIENDO LOS AÑOS DE ESCOLARIDAD CON LOS INGRESOS MENSUALES

ggplot(data=enoe, mapping = aes(x =anios_esc, y =ingreso_mensual))+

  geom_point(mapping = aes(color=niv_edu), show.legend = TRUE)+

  geom_smooth()

#PROPORCIÓN DE HOMBRES VS MUJERES EN LA MUESTRA ANALIZADA

ggplot(data=enoe)+

  geom_bar(mapping = aes(x=sex, y=..prop.., group=1, ))

#conteo de hombres y mujeres de la muestra

ggplot(data=enoe)+

  geom_bar(mapping = aes(x=sex, fill=sex))

#conteo de Individuos por sexo clasificado según su nivel de estudio

ggplot(data=enoe)+

  geom_bar(mapping = aes(x=sex, fill=niv_edu))

ggplot(data=enoe, mapping = aes(x=sex, fill=niv_edu,)))+

  geom_bar( position = "dodge")+

  labs(title="Observaciones por sexo y nivel educativo", x="Sexo", y="-

Observaciones")

#idem anterior pero por porporciones
```



**XII Muestra Académica de Trabajos de Investigación de la Licenciatura en  
Administración**

```
ggplot(data=enoe, mapping = aes(x=sex, fill=niv_edu))+  
  
  geom_bar( position = "fill")  
  
library(ggplot2)  
  
library(tidyselect)  
  
library(htmltools)  
  
library(plotly)  
  
library(yaml)  
  
#Generar un plotly basado en un boxplot  
  
ggplotly(  
  
  ggplot(enoe,aes(x=estado,y=anios_esc,fill=estado))+  
  
  geom_boxplot(color='black',show.legend = T))  
  
#ANÁLISIS DE CORRESPONDENCIA#####  
  
library(ca)  
  
library("ggplot2")  
  
library("factoextra")  
  
library("FactoMineR")  
  
library("gridExtra")  
  
enoe<-read.delim("clipboard", dec=",")  
  
for (i in 1:ncol(enoe)) enoe[,i]=as.factor(enoe[,i])  
  
str(enoe)  
  
#Análisis descriptivo  
  
ggplotly( ggplot(enoe,aes(x=estado))+ geom_bar(fill= "#DDB4EB"))  
  
ggplotly(ggplot(enoe,aes(x=sex))+ geom_bar(fill= "#FFD4A5"))  
  
ggplotly(ggplot(enoe,aes(x=edad))+ geom_bar(fill= "#41894A"))  
  
ggplotly(ggplot(enoe,aes(x=pos_ocu))+ geom_bar(fill= "#43857A"))  
  
ggplotly(ggplot(enoe,aes(x=ing_salarios))+ geom_bar(fill= "#FFEC75"))  
  
ggplotly(ggplot(enoe,aes(x=niv_edu))+ geom_bar(fill= "#FFEC30"))  
  
ggplotly( ggplot(enoe,aes(x=num_trabajos))+ geom_bar(fill= "#DDB4EB"))  
  
ggplotly(ggplot(enoe,aes(x=tipo_empleo))+ geom_bar(fill= "#41894A"))  
  
enoe <- MCA(enoe, graph = FALSE)  
  
print(enoe)  
  
print(enoe$var)  
  
fviz_screplot(enoe,geom="line")+  
  
  theme_grey()  
  
fviz_mca_var(enoe, col.var = "blue", addEllipses = FALSE, repel = TRUE) + theme_minimal()  
  
#Contribuciones al eje 1
```



**XII Muestra Académica de Trabajos de Investigación de la Licenciatura en Administración**

```
fviz_contrib(enoef, choice = "var", axes = 1, top = 15)+labs(title = " Contribución de las Categorías para las Dimension 1")
```

```
#Contribuciones al eje 2
```

```
fviz_contrib(enoef, choice = "var", axes = 2, top = 15)+labs(title = " Contribución de las Categorías para las Dimension 2")
```

```
#Contribuciones a ambos ejes
```

```
fviz_contrib(enoef, choice = "var", axes = 1:2, top = 15)+labs(title = " Contribuciones de las Categorías para las Dimensiones 1-2")
```

```
#Gráfico de calor por contribuciones
```

```
fviz_mca_var(enoef, col.var = "contrib",
```

```
  gradient.cols = c("#00AFBB", "#E7B800", "#FC4E07"),
```

```
  ggtheme = theme_grey()
```

```
  , repel = TRUE)
```

```
#boxplot
```

```
ggplot(enoef,aes(x="Descripción",y="Total"))+
```

```
  geom_boxplot(color='black',show.legend = F)
```

```
boxplot(enoef)
```

```
getwd()
```

```
library(cluster)
```

```
clusplot(wine.stand, k.means.fit$cluster,
```

```
  main='2D representation of the Cluster solution',
```

```
  color=TRUE, shade=TRUE,
```

```
  labels=2, lines=0)
```

```
# CLUSTER#####
```

```
# Cargamos el data set
```

```
enoefcluster<-read.delim("clipboard", dec=",")
```

```
for (i in 1:ncol(enoef)) enoef[,i]=as.factor(enoef[,i])
```

```
head(enoefcluster)
```

```
str(enoefcluster)
```

```
library(NbClust)
```

```
library(cluster)
```

```
library(ggplot2)
```

```
library(factoextra)
```

```
#detectar outliers
```

```
mean<-colMeans(enoefcluster)
```

```
Sx<-cov(enoefcluster)
```

```
D2<-mahalanobis(enoefcluster, mean, Sx, inverted = FALSE)
```

```
print(D2)
```



## XII Muestra Académica de Trabajos de Investigación de la Licenciatura en Administración

```
pchisq(D2, df= 6, lower.tail= FALSE)

qchisq(.99, df= 6)

# Como existen diferencias en las escalas de las variables, debemos escalar las mismas

enocluster.esc <- scale(enocluster, center = TRUE, scale = TRUE)

#Obtenemos la matriz de distancias

dist <- hclust(d = dist(x = enocluster.esc, method = "euclidean"), method = "complete")

res.dist <- get_dist(enocluster.esc, stand = TRUE, method = "euclidean")

fviz_dist(res.dist, gradient = list(low = "#00AFBB", mid = "white", high = "#FC4E07"))

matriz.dis.euclid<-dist(enocluster.esc,method="euclidean",diag=TRUE)

round(print(matriz.dis.euclid),2)

#Analizamos el n?mero ?ptimo de clusters

library(ggplot2)

library(factoextra)

fviz_nbclust(x = enocluster.esc, FUNcluster = kmeans, method = "wss", k.max = 15,

            diss = get_dist(enocluster.esc, method = "euclidean"), nstart = 50)

fviz_nbclust(x = enocluster.esc, FUNcluster = kmeans, method = "gap_stat", k.max = 15,

            diss = get_dist(enocluster.esc, method = "euclidean"), nstart = 50)

fviz_nbclust(x = enocluster.esc, FUNcluster = kmeans, method = "silhouette", k.max = 15,

            diss = get_dist(enocluster.esc, method = "euclidean"), nstart = 50)

#Pedimos los indicadores para un análisis jerárquico con distancia euclídea y método de ward

library(NbClust)

res<-NbClust(enocluster.esc, distance = "euclidean", min.nc=2, max.nc=8, method = "ward.D2", index = "alllong")

res$All.index

res$Best.nc

res$Best.partition

fviz_nbclust(res)

#Realizamos un dendrograma

fviz_dend(x = dist, k = 6, cex = 0.6) + geom_hline(yintercept = 4, linetype = "dashed") + labs(title = "clustering", subtitle = "Distancia eucl?dea, K=4")

fviz_dend(dist, cex = 0.8, k=6,

          rect = TRUE,

          k_colors = "jco",

          rect_border = "jco",

          rect_fill = TRUE,

          horiz = FALSE)

#Probamos otro tipo de dendrograma

library(igraph)
```



## XII Muestra Académica de Trabajos de Investigación de la Licenciatura en Administración

```
fviz_dend(dist, cex = 0.8, lwd = 0.8, k = 4,

  rect = TRUE,

  k_colors = "jco",

  rect_border = "jco",

  rect_fill = TRUE,

  type = "phylogenic")

fviz_dend(dist, cex = 0.8, lwd = 0.8, k = 4,

  rect = TRUE,

  k_colors = "jco",

  rect_border = "jco",

  rect_fill = TRUE,

  type = "circular")

# Seleccionamos al azar K Centroides en la base de datos normalizada.

# Para obtener los mismos resultados cuantas veces se reproduzca tenemos que

# usar una semilla de números aleatorios

# para K = 4 y asignaciones aleatorias nstart = 50

set.seed(123)

km.res <- kmeans(enocluster.esc, 6, nstart = 25)

print(km.res)

fviz_cluster(km.res, data = enocluster.esc, ellipse.type = "convex")

#Otra visualización

fviz_cluster(object = km.res, data = enocluster.esc, show.clust.cent = TRUE,

  ellipse.type = "euclid", star.plot = TRUE, repel = TRUE) +

labs(title = "Resultados clustering K-means") +

theme_bw() +

theme(legend.position = "none")

# Calculamos la media de cada variable por cluster utilizando los datos originaleslibrary(dplyr)

A<-aggregate(enocluster, by = list(km.res$cluster), FUN = mean)

print(A)

#Otra manera

library(kableExtra)

aggregate(enocluster, by = list(cluster = km.res$cluster), mean) %>%

kable(caption = "Media de Cada Variable por Conglomerado, Utilizando los datos originales" ,

  align = "c",

  digits = 6) %>%

kable_classic_2(html_font = "sans-serif",
```



```
lightable_options = c("hover", "striped")) %>%  
  
row_spec(0,  
  
  bold = T,  
  
  color = "white",  
  
  background = "#219B6D")  
  
# Validación interna de los clusters  
  
library("clValid")  
  
intern <- clValid(enocluster_esc, nClust = 2:6,  
  
  clMethods = c("hierarchical", "kmeans", "pam", "clara"),  
  
  validation = "internal")  
  
# Summary  
  
summary(intern)  
  
plot(intern)  
  
library(psych)  
  
anova1 <- aov(Murder ~ km.res$cluster, data = enocluster)  
  
anova2 <- aov(Assault ~ km.res$cluster, data = enocluster)  
  
anova3 <- aov(UrbanPop ~ km.res$cluster, data = enocluster)  
  
anova4 <- aov(Rape ~ km.res$cluster, data = enocluster)  
  
summary(anova1)  
  
summary(anova2)  
  
summary(anova3)  
  
summary(anova4)  
  
#Añadimos el número de grupo a los datos originales  
  
enoclustercluster <- cbind(enocluster, km.res$cluster)  
  
print(enoclustercluster)
```

## BIBLIOGRAFÍA

Aldás, J., & Uriel, E. (2017). *Análisis multivariante aplicado con R*. Madrid: Parafino SA.

García, J., Molina, J., & Bustamante, ä. (2018). *Ciencia de datos, Técnicas analíticas y aprendizaje estadístico en un enfoque práctico*. Bogotá: Alfaomega.

Hair, Anderson, Tatham, & Black. (1999). *Análisis multivariante*. Madrid: Prentice Hall, Quinta edición.

Lind, D., Marchal, W., & Wathen, S. (2012). *Estadística aplicada a los negocios y la economía 15° Edición*. México: Mc Graw Hill.