

CoderHouse

Curso Data Science

Informe del Proyecto Final

**Análisis socioeducativo de los habitantes de la Ciudad de
Buenos Aires**

Profesor: Damian Dapuetto

Tutor: Héctor Alonso

Grupo de Trabajo:

Lucia Buzzeo, Lucia Hukovsky,
Jose Saint German, Juan Martín Carini

4 de septiembre de 2022



Índice

1. Presentación del problema y fuente de información	3
1.1. Presentación del problema	3
1.2. Definición de la fuente de información	3
2. Pregunta y objetivos de investigación	3
3. Orden de trabajo	3
4. Introducción a las variables: Análisis exploratorio de los datos	3
4.1. Análisis univariado	6
4.1.1. Género y edad	6
4.1.2. Comuna	6
4.1.3. Ingreso familiar per capita	6
4.1.4. Años de escolaridad	7
4.1.5. Máximo nivel educativo (Target)	8
4.2. Análisis bivariado	8
4.2.1. Comparación entre variables numéricas	9
4.2.2. Comparación de variables categóricas con numéricas	9
4.2.3. Variable numéricas con comuna	11
4.3. Análisis multivariado	12
5. Modelos analíticos	15

1. Presentación del problema y fuente de información

1.1. Presentación del problema

Nos es de gran de interés vivir en una comunidad con políticas públicas eficaces que mejoren las condiciones de vida de las personas. En este sentido, hemos decidido analizar los diferentes ejes que en nuestro país se rigen por políticas publicas. Al respecto, encontramos una gran limitación en el eje de educación al reconocer que su acceso dista de ser equitativo. Este aspecto no nos resultó una novedad, sin embargo, nos dio el pie para comenzar una investigación que permita dar una explicación teórica a la problemática. En concreto, nos ha permitido conocer mejor la situación educativa actual de CABA y descubrir las principales variables que afectan el nivel educativo.

El análisis realizado en el marco del presente proyecto podría establecer una base de requerimientos que permitan generar políticas públicas efectivas, no solo en el ámbito educativo, sino en el económico, cultural, social y geográfico, entre otros.

1.2. Definición de la fuente de información

Para trabajar esta problemática, hemos decidido recurrir a la [datos abiertos](#) del Gobierno de la Ciudad de Buenos Aires para el año 2019. El mismo está disponible en la base de [Encuesta Anual de Hogares](#) del GCBA.

Esta encuesta contiene información demográfica, social, económica, educativa y de salud de 14319 habitantes de la Ciudad, la cual es una muestra representativa que permite obtener un vistazo de la población de la Ciudad.

2. Pregunta y objetivos de investigación

Nuestro objetivo principal es descubrir las principales variables intervinientes en el nivel máximo educativo alcanzado por la población de la Ciudad Autónoma de Buenos Aires (CABA).

De este objetivo principal se desprenden los siguientes sub-objetivos:

- Determinar si la ubicación geográfica del encuestado es determinante para alcanzar ciertos niveles educativos. De este objetivo se desprende determinar la relación entre el nivel educativo y la comuna del encuestado, así como la relación entre la misma variable y el hecho de que el encuestado habite en una villa de emergencia.
- Establecer la fuerza con la que el nivel socio-económico afecta la variable target.
- Explorar la relación del target con otras variables, como el sexo del encuestado, la cantidad de hijos, la afiliación de salud o la edad.

3. Orden de trabajo

Este trabajo estará dividido en 3 partes:

1. **Introducción a las variables del problema:** Se hará un análisis de las variables en donde se buscará conocer su performance dentro del dataset y su potencial significanos para la pregunta que buscamos responder. A la vez, queremos ver cómo las variables interactúan entre si. Esta parte es lo que se conoce como análisis univariado, bivariado y multivariado,
2. **Modelos analíticos:** En esta sección se llevarán a cabo diversos modelos analíticos y algoritmos que nos servirán para acercarnos a la respuesta a nuestra pregunta de investigación,
3. **Conclusión:** Haremos conclusiones finales sobre nuestros hallazgos. Además, discutiremos posibles limitaciones que tuviera y plantearemos futuras líneas de análisis a partir del análisis presente.

4. Introducción a las variables: Análisis exploratorio de los datos

Una ves cargado el dataset con el que vamos a trabajar, miramos sus variable, el tipo que son y si tienen nulls:

RangeIndex: 14319 entries, 0 to 14318

Data columns (total 31 columns):

#	Column	Non-Null Count	Dtype
0	id	14319 non-null	int64
1	nhogar	14319 non-null	int64
2	miembro	14319 non-null	int64
3	comuna	14319 non-null	int64
4	dominio	14319 non-null	object
5	edad	14319 non-null	int64
6	sexo	14319 non-null	object
7	parentesco_jefe	14319 non-null	object
8	situacion_conyugal	14318 non-null	object
9	num_miembro_padre	14319 non-null	object
10	num_miembro_madre	14319 non-null	object
11	estado_ocupacional	14319 non-null	object
12	cat_ocupacional	14319 non-null	object
13	calidad_ingresos_lab	14319 non-null	object
14	ingreso_total_lab	14319 non-null	int64
15	calidad_ingresos_no_lab	14319 non-null	object
16	ingreso_total_no_lab	14319 non-null	int64
17	calidad_ingresos_totales	14319 non-null	object
18	ingresos_totales	14319 non-null	int64
19	calidad_ingresos_familiares	14319 non-null	object
20	ingresos_familiares	14319 non-null	int64
21	ingreso_per_capita_familiar	14319 non-null	int64
22	estado_educativo	14319 non-null	object
23	sector_educativo	14316 non-null	object
24	nivel_actual	14319 non-null	object
25	nivel_max_educativo	13265 non-null	object
26	años_escolaridad	14257 non-null	object
27	lugar_nacimiento	14318 non-null	object
28	afiliacion_salud	14315 non-null	object
29	hijos_nacidos_vivos	6535 non-null	object
30	cantidad_hijos_nac_vivos	14319 non-null	object

dtypes: int64(10), object(21)
memory usage: 3.4+ MB

Ahora, en base a los datos arrojados por la tabla de arriba, generamos diversas transformaciones de variables, así como la creación de la variable “Target”, pues es la que usaremos para todo el análisis:

- creamos el target para “nivel_max_educativo”,
- remplazamos los valores de “años_escolaridad” para que todos sean numéricos,
- la variable “cantidad_hijos_nac_vivos” se puede pasar a numérica si se toma “no corresponde” como NAN,
- hay determinadas variables (comuna, id, hogar y miembro) que están como numéricas pero deberían ser categóricas,
- por otro lado, variables como sexo y dominio pueden pasarse a numérico mediante one hot encoding,
- generamos la variable “target” como copia de “Target” para tener ambas versiones,
- pasamos la variable “Target” a one hot encoding,
- y por último renombramos algunas variables para que sean más cortas.

Luego, con el dataset ya acomodado, comenzamos analizándolo en su conjunto. Miramos las nuevas modificaciones en las variable, el tipo que son y si tienen nulls:

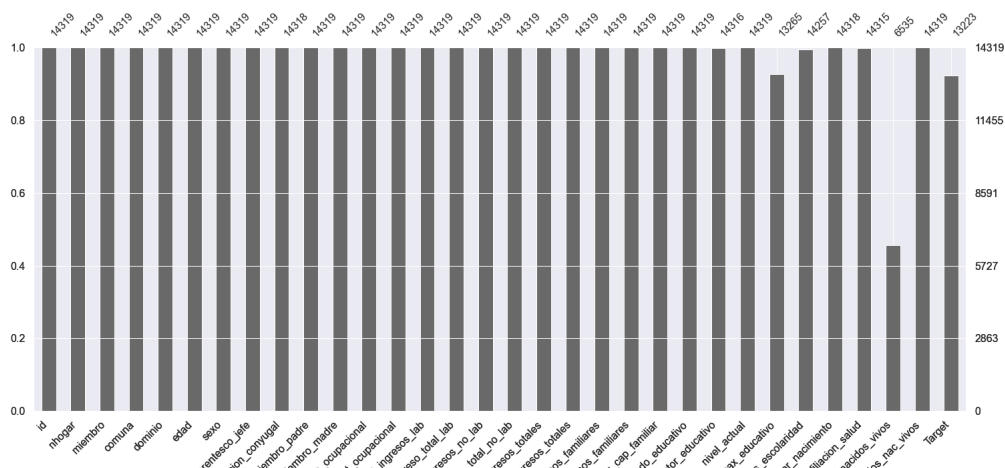
RangeIndex: 14319 entries, 0 to 14318

Data columns (total 36 columns):

#	Column	Non-Null Count	Dtype
0	id	14319 non-null	object
1	nhogar	14319 non-null	object
2	miembro	14319 non-null	object
3	comuna	14319 non-null	object
4	edad	14319 non-null	int64
5	parentesco_jefe	14319 non-null	object
6	situacion_conyugal	14318 non-null	object
7	num_miembro_padre	14319 non-null	object
8	num_miembro_madre	14319 non-null	object
9	estado_ocupacional	14319 non-null	object
10	cat_ocupacional	14319 non-null	object
11	calidad_ingresos_lab	14319 non-null	object
12	ingreso_total_lab	14319 non-null	int64
13	calidad_ingresos_no_lab	14319 non-null	object
14	ingreso_total_no_lab	14319 non-null	int64
15	calidad_ingresos_totales	14319 non-null	object
16	ingresos_totales	14319 non-null	int64
17	calidad_ingresos_familiares	14319 non-null	object
18	ingresos_familiares	14319 non-null	int64
19	ing_per_cap_familiar	14319 non-null	int64
20	estado_educativo	14319 non-null	object
21	sector_educativo	14316 non-null	object
22	nivel_actual	14319 non-null	object
23	nivel_max_educativo	13265 non-null	object
24	años_escolaridad	14257 non-null	float64
25	lugar_nacimiento	14318 non-null	object
26	afiliacion_salud	14315 non-null	object
27	hijos_nacidos_vivos	6535 non-null	object
28	cant_hijos_nac_vivos	14319 non-null	int64
29	sexo_Varon	14319 non-null	uint8
30	dominio_villas	14319 non-null	uint8
31	target	13223 non-null	object
32	Target_inicial	14319 non-null	uint8
33	Target_prim_completo	14319 non-null	uint8
34	Target_sec_completo	14319 non-null	uint8
35	Target_superior	14319 non-null	uint8

dtypes: float64(1), int64(7), object(22), uint8(6)

memory usage: 3.4+ MB

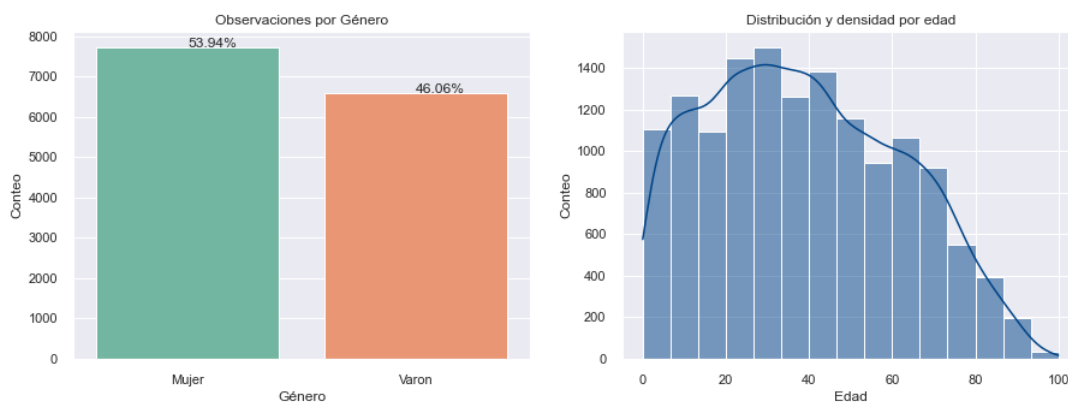


Detectamos que nuestra variable target tiene 1054 valores nulos. Es importante tener este dato presente cuando querramos correr un algoritmo de clasificación.

4.1. Análisis univariado

4.1.1. Género y edad

Comenzamos con un pantallazo general sobre las primeras cualidades de los datos, como muestra representativa para la EPH, sobre quiénes son los ciudadanos representado en el dataset.

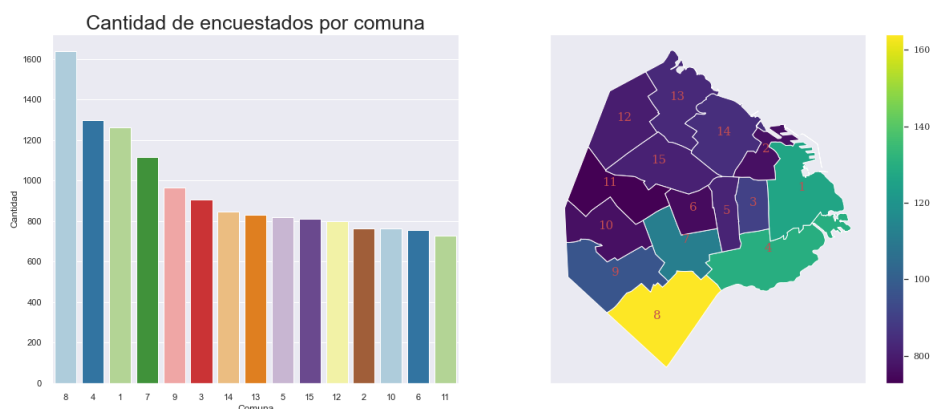


En la variable género los datos parecen equilibrados en las categorías. Para el caso de la variable 'edad', la distribución se asemeja a la de una normal.

4.1.2. Comuna

Seguimos observando la variable "comuna". En la misma se muestra la comuna de la Ciudad de Buenos Aires del entrevistado, de manera de tener una ubicación geográfica. Consideramos importante revisar esta variable ya que tenemos como hipótesis que el nivel educativo alcanzado puede estar dependiendo de la zona geográfica de la ciudad en la que se encuentra el entrevistado.

Para esto vamos a generar un mapa, así que utilizaremos el mapa de comunas de la Ciudad de Buenos Aires, transformamos las variables que vamos a usar para joinear el mapa con la base de manera que coincidan, transformamos la base para contabilizar la frecuencia con la que aparece cada comuna en la base. Y por último unimos ambos datasets y generamos una nueva variable con las coordenadas para poder agregar etiquetas en el centro geográfico de cada comuna, que nos da como resultado los siguientes graficos:

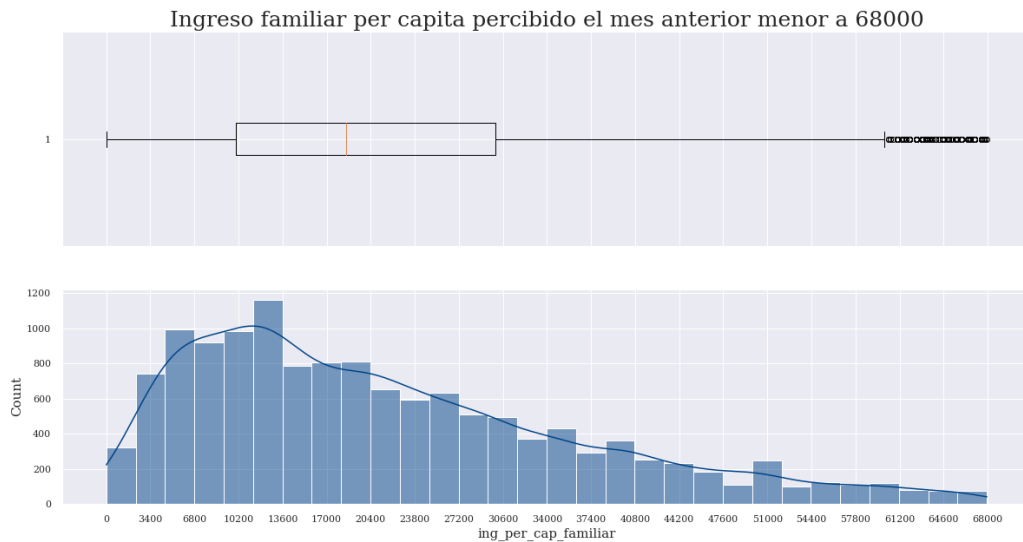


Observando ambos gráficos vemos que las comunas 1,4,7 y 8 tienen mayor cantidad de casos. Queda por verse si en posteriores análisis es necesario abordar esta diferencia para evitar sesgos. Para eso, será necesario tomar en cuenta el porcentaje de la población total de cada comuna.

4.1.3. Ingreso familiar per capita

Ahora probamos con observar los ingresos familiares. Creemos que puede ser un indicador interesante del nivel educativo.

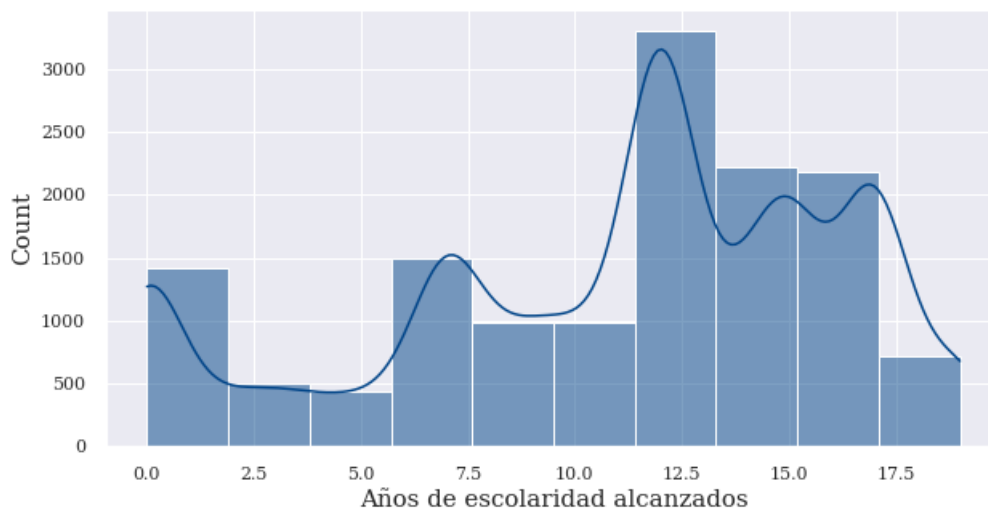
Para esto, armamos una función para graficar y jugar con el nivel del filtrado de la variable y obtener un histograma que permita apreciar mejor la distribución de la variable sin tantos outliers. Probamos graficando con el máximo de la variable:



Y como hay muchos outliers que impiden ver la distribución correctamente, los quitamos de los gráficos: De este forma vemos que, aún removiendo los outliers, la distribución sigue sesgada.

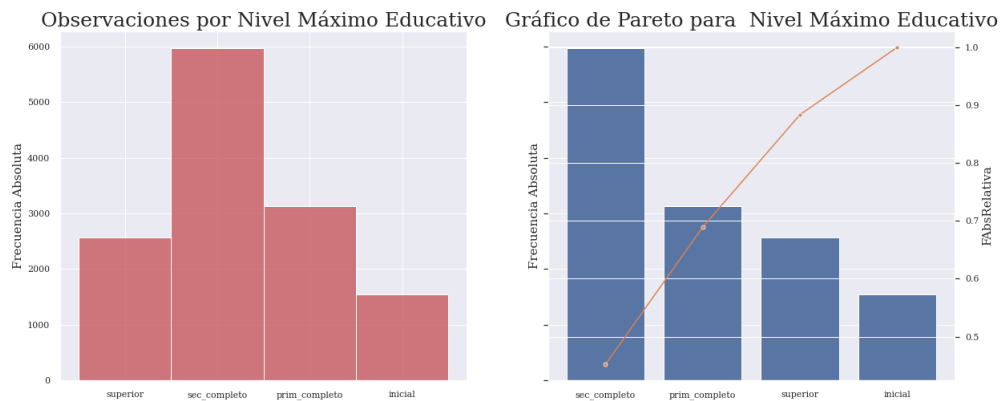
4.1.4. Años de escolaridad

Analizamos mediante un gráfico de barras los años de escolaridad alcanzados por los encuestados:



A simple vista se observan tres "picos": en el valor mínimo, alrededor del 7.5 y alrededor del 12.5. Podemos inferir que estos tres casos corresponden a no tener estudios, solo haber transcurrido el primario y haber transcurrido hasta la educación secundaria, respectivamente.

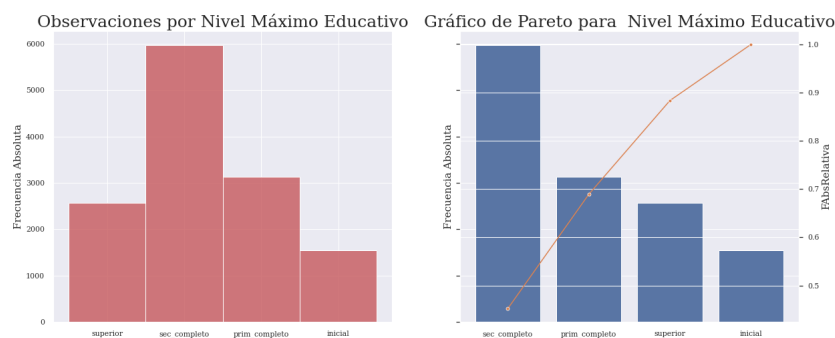
4.1.5. Máximo nivel educativo (Target)



Podemos observar que el nivel máximo educativo más alcanzado es el secundario completo, seguido por el primario. Contrario de lo que habíamos intuido anteriormente, el nivel superior quedó en tercer lugar. Adicionalmente, el nivel secundario y primario explican casi el 77 % de los datos.

4.2. Análisis bivariado

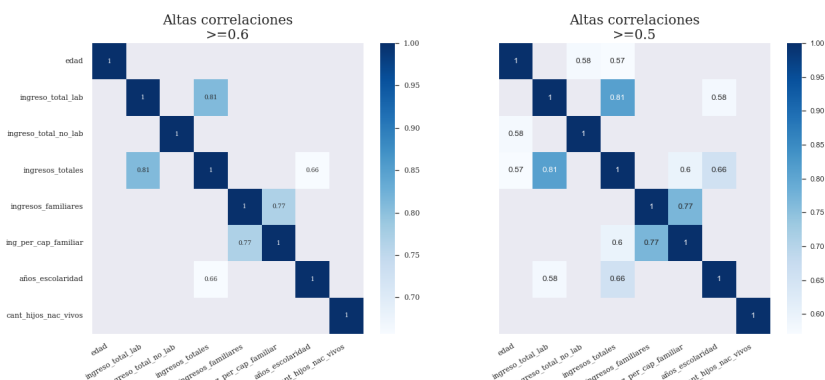
Para comenzar el análisis bivariado del problema, realizamos diferentes heatmaps para ver si algo nos llama la atención entre las variables numéricas.



A simple vista, no se observan fuertes correlaciones.

Podemos notar que la variable .años_escolaridad correlaciona moderadamente bien con variables relacionadas al ingreso.

La principal correlación positiva es .años_escolaridad con ingreso familiar per cápita (.ing_per_cap_familiar"), lo cual hace sentido teórico.



Aquí, vemos que los años de escolaridad alcanzados por los entrevistados tienen algo relación (66 %) con la variable ingresos_totales".

Por último corremos una tabla de correlación y filtramos las de valores más altos

	Variable_1	Variable_2	corr_value
2	ingreso_total_lab	ingresos_totales	0.8028054238319398
6	ingresos_familiares	ing_per_cap_familiar	0.7616082762779849
4	ingresos_totales	ing_per_cap_familiar	0.6214687193744823
5	ingresos_totales	años_escolaridad	0.602779948043472
1	edad	ingreso_total_no_lab	0.5734168774116969
7	años_escolaridad	Target	0.571440140244098
3	ingreso_total_lab	años_escolaridad	0.5389383108306696

Conclusiones:

- Como es esperable, hay alta correlación entre las variables relacionadas al ingreso.
- A su vez, encontramos una alta correlación (66 %) entre los ingresos y los años de escolaridad.
- También observamos una relación positiva entre la edad y los ingresos totales.

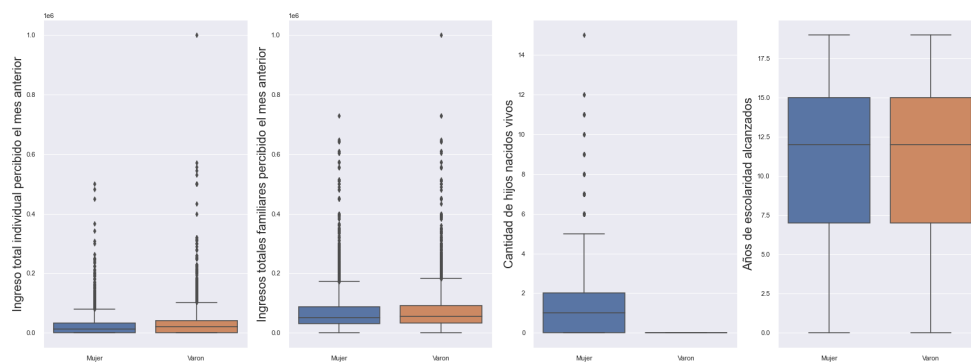
4.2.1. Comparación entre variables numéricas



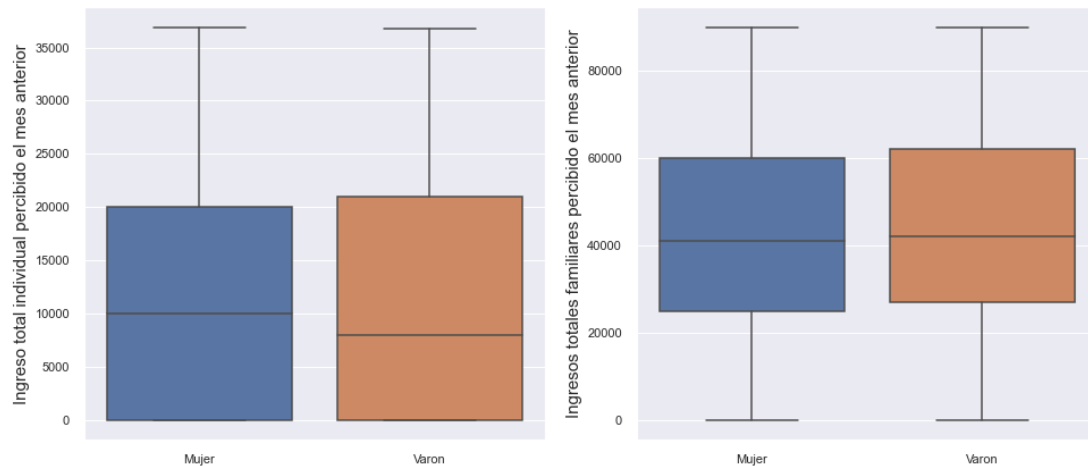
Se puede ver que desde los 30 años en adelante el ingreso total de la persona se corresponde con el ingreso familiar. Por ende suele haber un unico ingreso fuerte por grupo familiar.

4.2.2. Comparación de variables categóricas con numéricas

Adicionalmente, vamos a comparar algunas variables con nuestro target, comenzando con los ingresos totales.

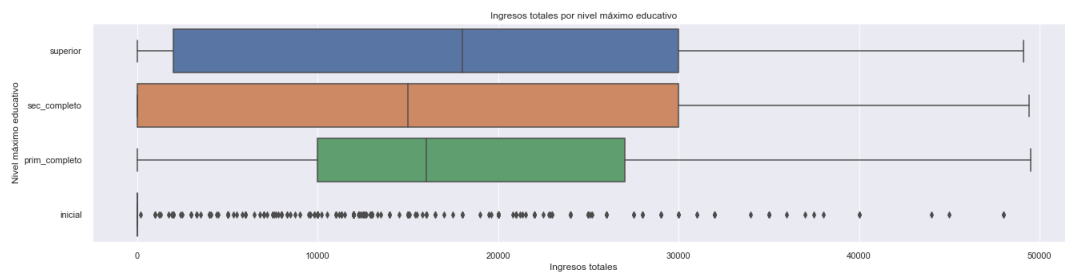


Probemos quitando outliers, a excepciónn de la cantidad de hijos nacidos vivos (puesto que no arrojará ningún dato nuevo) y de años de escolaridad (que no tiene outliers)

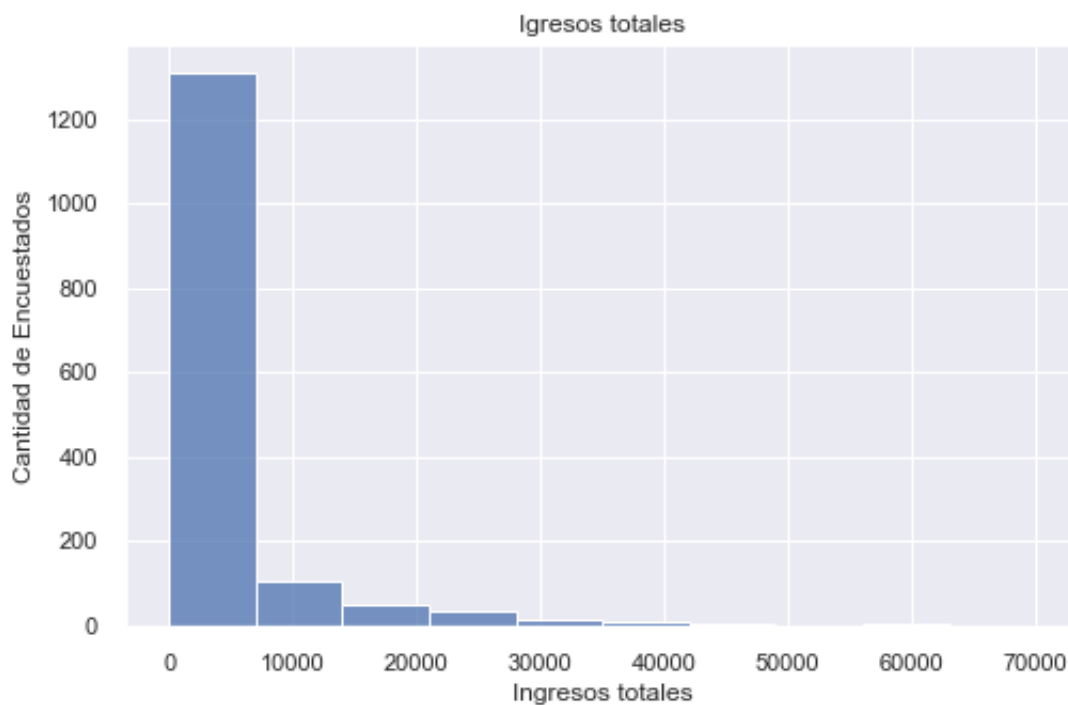


Por parte de las variables de ingreso, no parece haber nada disruptivo. La distribución por ingreso y años de escolaridad pareciera ocurrir pero no en un orden lineal.

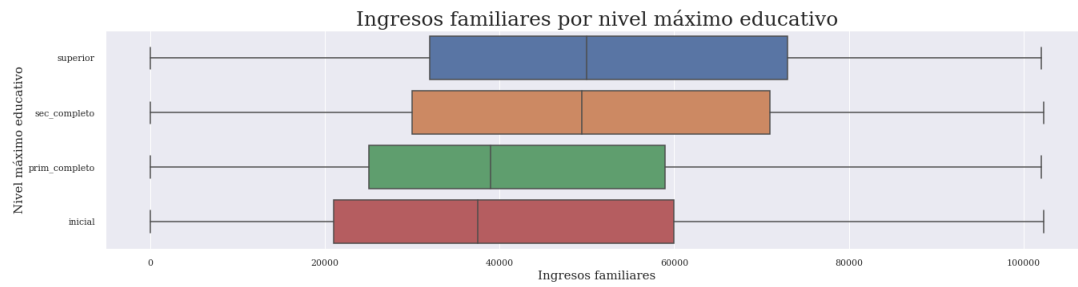
Llama la atención la variable "sexo": por algún motivo, todos los encuestados hombres figuran sin hijos nacidos vivos. Alternativamente, se podría investigar la metodología de la encuesta para ver si hay alguna respuesta. Adicionalmente, los hombres parecieran tener ingresos totales y familiares mayores que las mujeres, pero no pareciera que haya distribuciones desiguales en los años de escolaridad.



Parece que para el nivel inicial la remoción de outliers en otra categoría sigue siendo insuficiente para mostrar la distribución real de la variable. Echemos un vistazo a los valores de esta categoría.

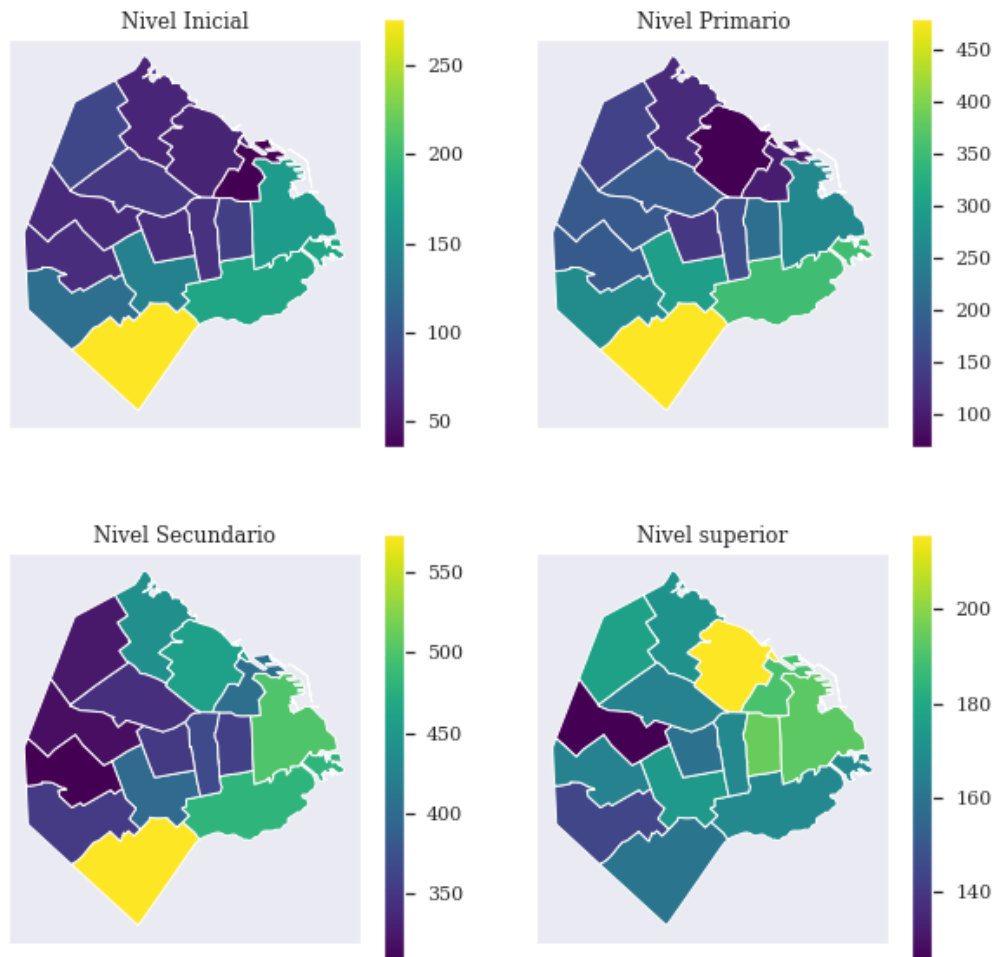


Logicamente, la enorme mayoría de los ingresos para el valor inicial dan 0, puesto que incluye a personas que en ese momento estaban cursando su educación inicial, por lo que tenían entre 2 y 6 años.

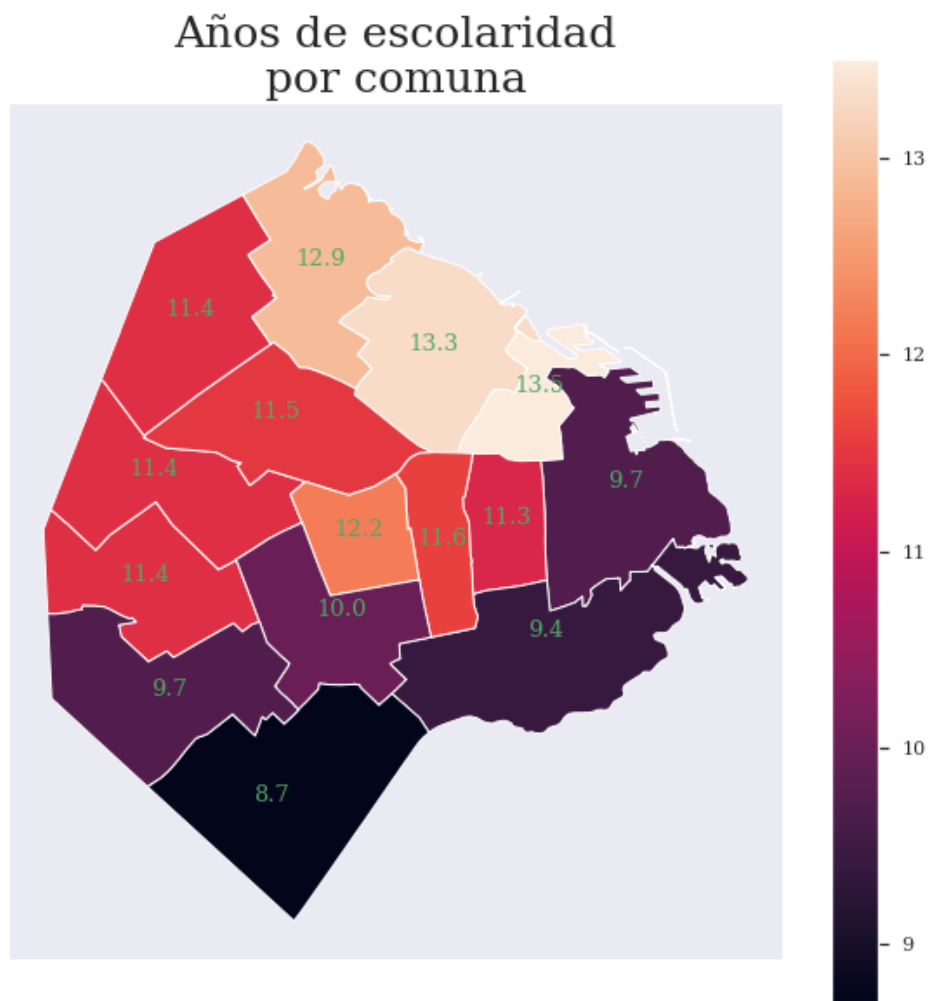


En definitiva, se observa un corrimiento de los valores centrales (dentro de la caja) hacia la izquierda a medida que aumenta el nivel educativo.

4.2.3. Variable numéricas con comuna



Se observa que en el sur de la ciudad hay mayor cantidad de encuestados con niveles de inicial, primario y secundario completo, mientras que el norte (particularmente el barrio de Palermo) tiene mayor cantidad de personas con estudios superiores. En menor medida también las comunas del este (comunmente llamado el "centro" de la ciudad) destacan por la cantidad de encuestados con nivel superior.

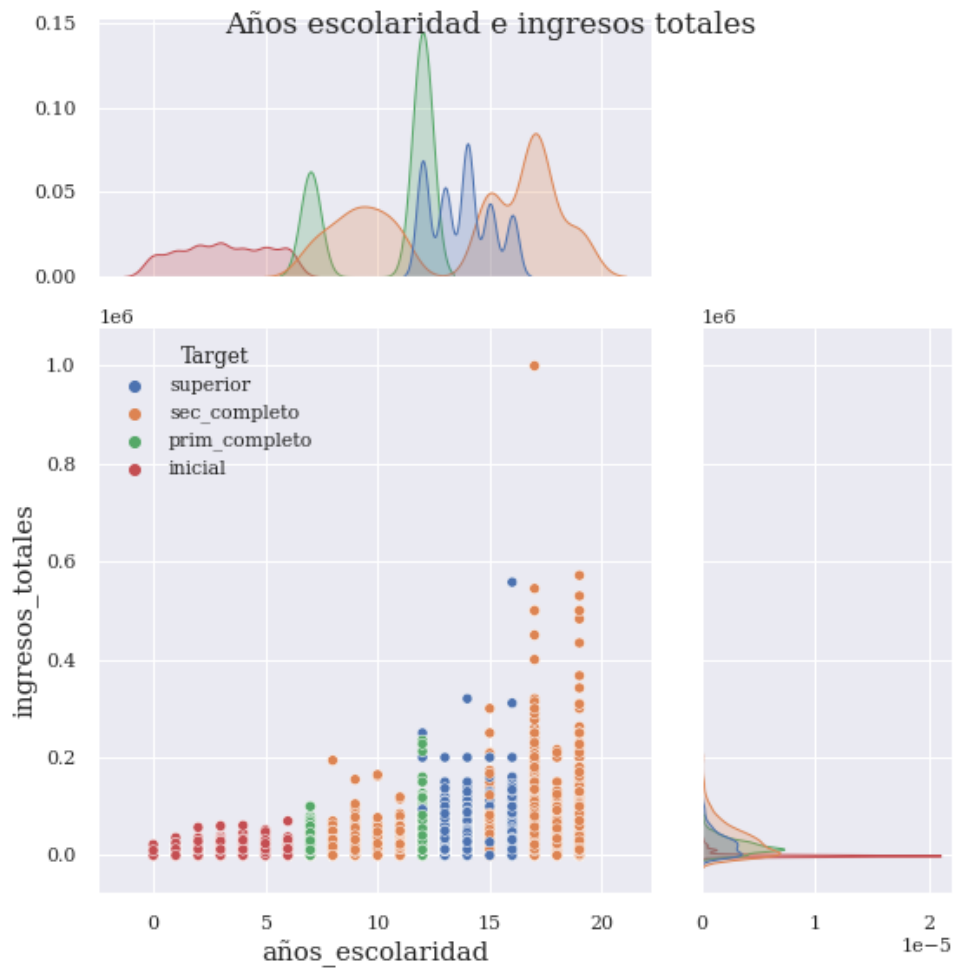


Lo que podemos observar en los últimos dos gráficos es una clara división geográfica del nivel educativo:

- Las comunas del norte son las que tienen mayor nivel educativo.
- Las comunas del centro tienen niveles medios.
- Las comunas del sur (con la comuna 6 en el centro de la ciudad como outlier) y la comuna 1 en el este son las que tienen niveles más bajos.

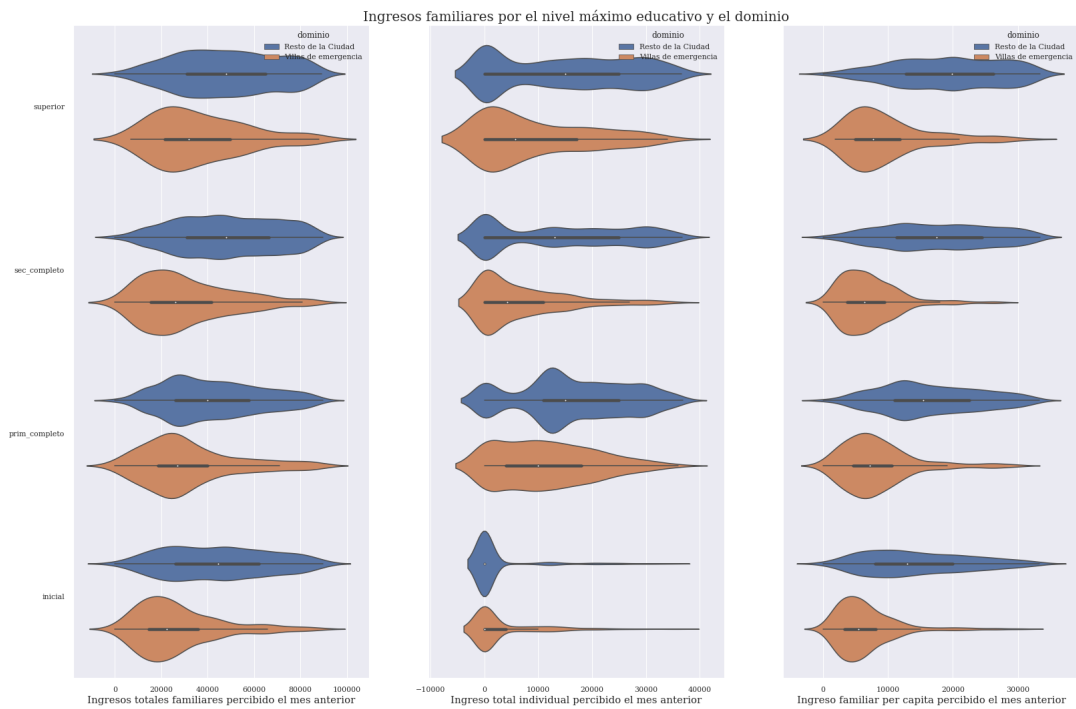
4.3. Análisis multivariado

Probamos de cruzar años de escolaridad, nivel máximo educativo y los ingresos totales.

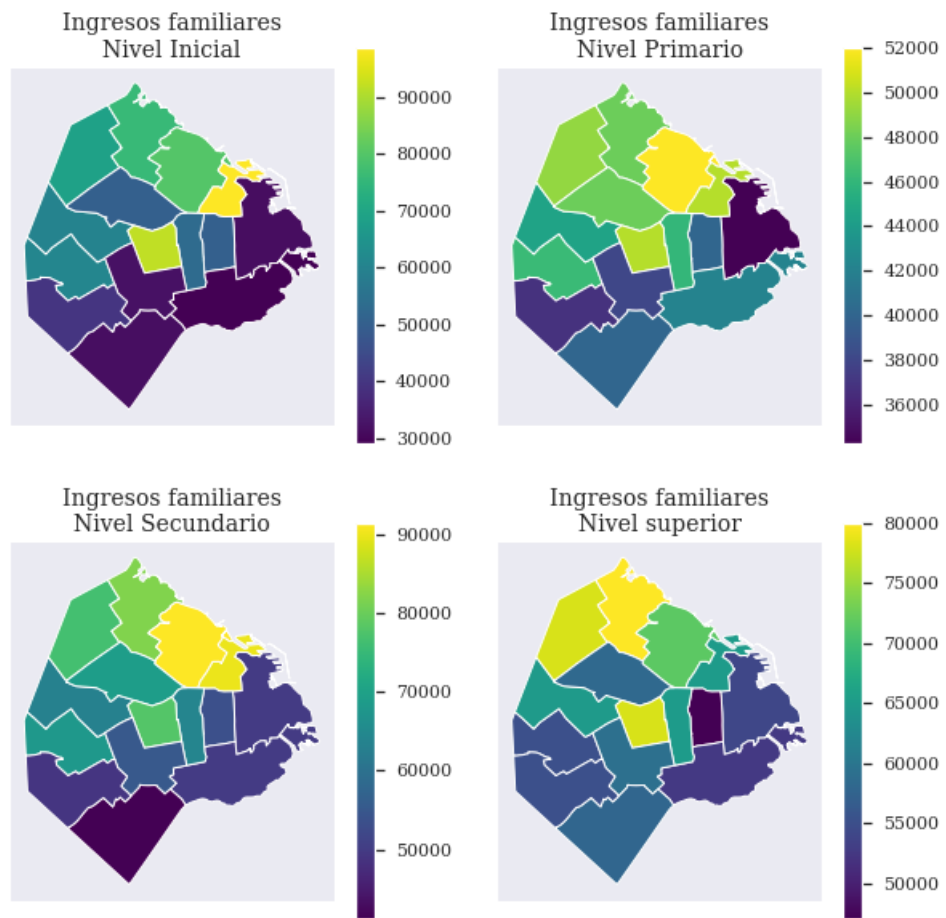


Conclusiones de la visualización:

- Hasta los 6 años, como era esperable, todos los casos llegan al nivel inicial.
- Vemos dos años en que aparece el primario completo: 7 y 12 años. Estimamos que se debe a la división entre los que comenzaron su educación en la primaria y los que comenzaron en el nivel inicial.
- A partir de los 12 años vemos un aumento consistente de los ingresos totales.



Aquí obtuvimos un descubrimiento interesante: no importa el nivel máximo educativo, los casos que no provienen de villas de emergencia (dominio=“villas_de_emergencia”) obtienen en promedio ingresos más altos en todos los niveles educativos. El alcanzar estudios superiores no parece homogeneizar ambos conjuntos. Esto se puede observar en el segundo gráfico, ya que el violín naranja acumula mayor cantidad de casos hacia la derecha, en comparación con los violines azules que tienen una mayor distribución.



Aquí podemos observar que a medida que avanza el nivel educativo máximo se atenúan levemente las diferencias de ingresos familiares entre comunas. Queda pendiente cruzar estos datos con la edad, para saber si el hecho de incluir a menores de edad está sesgando los valores para nivel inicial, primario y secundario.

5. Modelos analíticos