

# Universidad Nacional Autónoma de México

FACULTAD DE ESTUDIOS SUPERIORES ACATLÁN



## COVID-19 EN MÉXICO - PROYECTO FINAL

*Diplomado Ciencia de Datos - Módulo I*

Autor:

André Marx Puente Arévalo

Profesora:

Carla Paola Malerva Reséndiz

28 / Junio / 2020

# Índice

<b>1. Introducción</b>	<b>2</b>
1.1. Acotamiento del Problema . . . . .	2
<b>2. Análisis Exploratorio Previo al Procesamiento de Datos</b>	<b>3</b>
2.1. Diccionario de Datos . . . . .	3
2.2. Creación de la Tabla Analítica de Datos . . . . .	4
2.3. Estadísticas Descriptivas . . . . .	5
2.4. Visualización de datos . . . . .	6
<b>3. Calidad de Datos</b>	<b>9</b>
3.1. Variables Categóricas . . . . .	10
3.2. Variable Continua . . . . .	11
3.3. Variable Tipo Texto . . . . .	11
3.4. Variables Tipo Fecha . . . . .	11
3.5. Completitud . . . . .	11
3.6. Generación de la Variable Objetivo . . . . .	11
<b>4. Datos Anómalos</b>	<b>13</b>
<b>5. Análisis Descriptivo Post-Procesamiento de Datos</b>	<b>14</b>
5.1. Visualización de datos . . . . .	14
<b>6. Imputación de Valores Faltantes</b>	<b>16</b>
<b>7. Ingeniería de Variables</b>	<b>17</b>
7.1. Variables Categóricas . . . . .	17
7.2. Variable Objetivo . . . . .	17
<b>8. Reducción de Dimensiones</b>	<b>18</b>
<b>9. Conclusión</b>	<b>19</b>
<b>10. Apéndice</b>	<b>20</b>
10.1. Variables Restantes . . . . .	20
10.2. Imputación de Valores Faltantes - Restantes . . . . .	21
10.3. Visualización de Datos - Comparación . . . . .	21

# 1. Introducción

En la actualidad, el mundo entero ha sido doblegado por la aparición de un nuevo virus denominado coronavirus SARS-Cov-2, que provoca una enfermedad infecciosa llamada COVID-19. Sus orígenes provienen de China, posteriormente se extendió a todos los continentes del mundo, provocando una pandemia a nivel mundial.

La mayoría de las personas infectadas con el virus COVID-19 experimentarán una enfermedad respiratoria leve a moderada y se recuperarán sin necesidad de un tratamiento especial. Las personas mayores y aquellas con problemas médicos subyacentes como enfermedades cardiovasculares, diabetes, enfermedades respiratorias crónicas y cáncer tienen más probabilidades de desarrollar enfermedades graves.

Por el momento, la mejor manera de prevenir y ralentizar la transmisión es estar bien informado sobre el virus COVID-19, la enfermedad que causa y cómo se propaga, ya que hasta la fecha, sigue sin existir vacuna que nos ayude a contrarrestarla. El virus se transmite principalmente a través de gotitas de saliva o secreciones nasales cuando una persona infectada tose o estornuda.

Actualmente se han reportado aproximadamente 45 millones de casos confirmados alrededor del mundo, de los cuales, 1,200,000 han sido muertes confirmadas según la Organización Mundial de la Salud para inicios del mes de noviembre del 2020.

## 1.1. Acotamiento del Problema

Dado el contexto anterior, en el presente proyecto se plantea realizar limpieza y análisis estadístico de un conjunto de datos, el cual, es proporcionado por el gobierno mexicano, contiene los registros hasta el 31 de octubre del 2020 de casos diarios asociados a COVID-19 a nivel federal. Cuenta con contenido desagregado por sexo, edad, nacionalidad, padecimientos asociados, entre otros. En un principio se cuentan con 2,403,499 registros y 38 columnas.

Se seleccionó México como objeto de estudio porque es uno de los países que más afectaciones ha sufrido por el virus en muchos ámbitos como en la pérdida de vidas humanas o económicamente hablando y muchos otros, a su vez es el país en el que radico actualmente y este virus ha llegado a cambiar nuestro estilo de vida de una forma inimaginable.

Ahora bien, resulta de interés el poder resumir la información que contienen los datos con pocas cifras o gráficas, ¿cómo cambian los datos tras su procesamiento? ¿se puede medir el porcentaje de pacientes por cada estado? En el presente proyecto se pretende dar solución a estos problemas con los datos anteriormente mencionados, al igual crear una tabla con la información necesaria y preparada para desarrollar un modelo que sea capaz de predecir la clasificación final de un paciente en función de su edad, padecimientos médicos, lugar de residencia, entre otros. Se clasificará como negativo (1, 0), sospechoso (0, 1) o confirmado (0, 0).

## 2. Análisis Exploratorio Previo al Procesamiento de Datos

La finalidad del Análisis Exploratorio es examinar los datos previamente a la aplicación de cualquier técnica estadística. De esta forma el analista consigue un entendimiento básico de sus datos y de las relaciones existentes entre las variables analizadas. Proporciona métodos sencillos para organizar y preparar los datos, detectar fallos en el diseño y recogida de datos, tratamiento y evaluación de datos ausentes, identificación de casos atípicos y comprobación de los supuestos subyacentes en la mayor parte de las técnicas multivariantes.

### 2.1. Diccionario de Datos

Dentro del conjunto de datos brindado en el portal de “Datos Abiertos” se encuentran los siguientes datos que han sido tipificados entre:

- Variables Continuas (c)
- Variables Categóricas (v)
- Variables de Texto (t)
- Variables de Fecha (d)

Variable	Tipo	Descripción
FECHA_ACTUALIZACION	Fecha	La base de datos se alimenta diariamente, esta variable permite identificar la fecha de la última actualización
ID_REGISTRO	Texto	Cadena identificadora del caso.
ORIGEN	Categórica	La vigilancia centinela se realiza a través del sistema de unidades de salud monitoras de enfermedades respiratorias (USMER). Las USMER incluyen unidades médicas del primer, segundo o tercer nivel de atención y también participan como USMER las unidades de tercer nivel que por sus características contribuyen a ampliar el panorama de información epidemiológica, entre ellas las que cuentan con especialidad de neumología, infectología o pediatría.
SECTOR	Categórica	Identifica el tipo de institución del Sistema Nacional de Salud que brindó la atención.
ENTIDAD_UM	Categórica	Identifica la entidad donde se ubica la unidad médica que brindó la atención.
SEXO	Categórica	Identifica al sexo del paciente.
ENTIDAD_NAC	Categórica	Identifica la entidad de nacimiento del paciente.
ENTIDAD_RES	Categórica	Identifica la entidad de residencia del paciente.
MUNICIPIO_RES	Categórica	Identifica el municipio de residencia del paciente.

TIPO_PACIENTE	Categórica	Identifica el tipo de atención que recibió el paciente en la unidad. Se denomina como ambulatorio si regresó a su casa o se denomina como hospitalizado si fue ingresado a hospitalización.
FECHA_INGRESO	Fecha	Identifica la fecha de ingreso del paciente a la unidad de atención.
FECHA_SINTOMAS	Fecha	Identifica la fecha en que inició la sintomatología del paciente.
FECHA_DEF	Fecha	Identifica la fecha en que el paciente falleció.
INTUBADO	Categórica	Identifica si el paciente requirió de intubación.
NEUMONIA	Categórica	Identifica si al paciente se le diagnosticó con neumonía.
EDAD	Continua	Identifica la edad del paciente.
NACIONALIDAD	Categórica	Identifica si el paciente es mexicano o extranjero.
EMBARAZO	Categórica	Identifica si la paciente está embarazada.
HABLA_LENGUA_INDIG	Categórica	Identifica si el paciente habla lengua indígena.
CLASIFICACION_FINAL	Categórica	Identifica si el paciente es un caso de COVID-19 según el catálogo CLASIFICACION_FINAL.

**Nota:** En el "Apéndice", en la sección de "Variables Restantes" se encuentran el resto de variables con su respectivo nombre, tipificación y descripción.

Es preciso mencionar que se cuenta con un total de 38 variables de las cuales 32 son de tipo categóricas, una es de tipo continua, cuatro de tipo fecha y una de tipo texto. Todas estas variables cuentan con su respectivo catálogo, ya que, en la base de datos original los registros son números (en la mayoría de los casos), pero estos, no nos dicen nada si no sabemos qué significa.

## 2.2. Creación de la Tabla Analítica de Datos

En el caso del conjunto de datos de casos de COVID-19 que estoy analizando, en un principio no se presentan datos faltantes o registros en blancos, ni registros duplicados, por lo que procedí a realizar los respectivos cruces con los catálogos (obtenidos de la misma fuente que el conjunto de datos) para saber qué significan los valores numéricos contenidos en el conjunto de datos.

Analizando el conjunto de datos, me percaté que para realizar el cruce con el catálogo de los municipios es necesario generar otra variable de apoyo tanto en el catálogo como en mi conjunto de datos, esto porque en el catálogo de municipios, la clave se reinicia por cada estado, entonces lo que hice fue generar una nueva variable que contenga la información de ambas columnas, es decir, CLAVE\_ENTIDAD - CLAVE\_MUNICIPIO. Dicha variable recibe el nombre de CLAVE\_ENTIDAD\_MUNICIPIO.

## 2.3. Estadísticas Descriptivas

La estadística descriptiva es la rama de las matemáticas que recolecta, presenta y caracteriza un conjunto de datos con el fin de describir apropiadamente las diversas características de ese conjunto.

Dada la naturaleza de nuestras variables, sólo es posible obtener las estadísticas descriptivas de la variable "c\_EDAD", ya que, es la que presenta valores continuos. Sus estadísticas son:

	c_EDAD
<b>Conteo</b>	2,403,499
<b>Media</b>	41
<b>Desviación</b>	16
<b>Mínimo</b>	0
<b>25 %</b>	29
<b>50 %</b>	40
<b>75 %</b>	53
<b>Máximo</b>	120

Cuadro 1: Estadísticas descriptivas de la variable c\_EDAD



Figura 1: Box plot de la variable c\_EDAD

Como se puede observar en la Gráfica 1 y Figura 1, la media de la edad es 41, esto quiere decir, que en promedio las personas que han contraído el virus tienen 41 años, más en general se aprecia que los contagios se concentran en un intervalo de edad que es de 29 a 53 años. Siendo algo "benéfico" para la población, ya que, las personas mayores a 65 años son del sector vulnerable, lo que significa que tienen más probabilidades de que se agraven sus síntomas. Analizando más a profundidad la Figura 1, nótese que presentamos unos cuantos datos atípicos en la edad, son las personas más longevas que han sido víctima del virus.

Por otro lado me resulta de interés analizar las edades respecto al sexo del paciente para intentar llegar a conclusiones más contundentes.

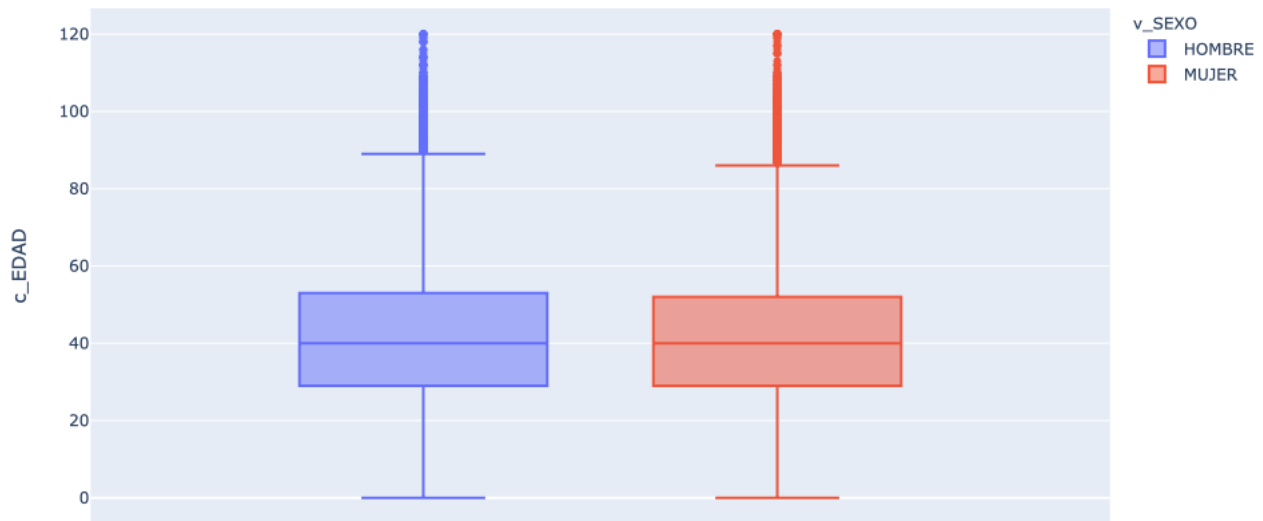


Figura 2: Box plot de la variable c\_EDAD respecto a cada sexo

Tras analizar la Figura 2, nos damos cuenta que no hay mucha variación si tomamos en cuenta el sexo de cada paciente, las principales diferencias es que la edad promedio en la que más contagiados hay en ambos sexos es de 40 años y en el caso particular de las mujeres el intervalo de edad en el que se concentran las pacientes es de 29 a 52 años.

## 2.4. Visualización de datos

Como persona que está viviendo actualmente la pandemia, me han surgido dudas como la siguiente: ¿se han atendido más hombres que mujeres por COVID-19? para dar solución a esta pregunta analicé la variable v\_SEXO.

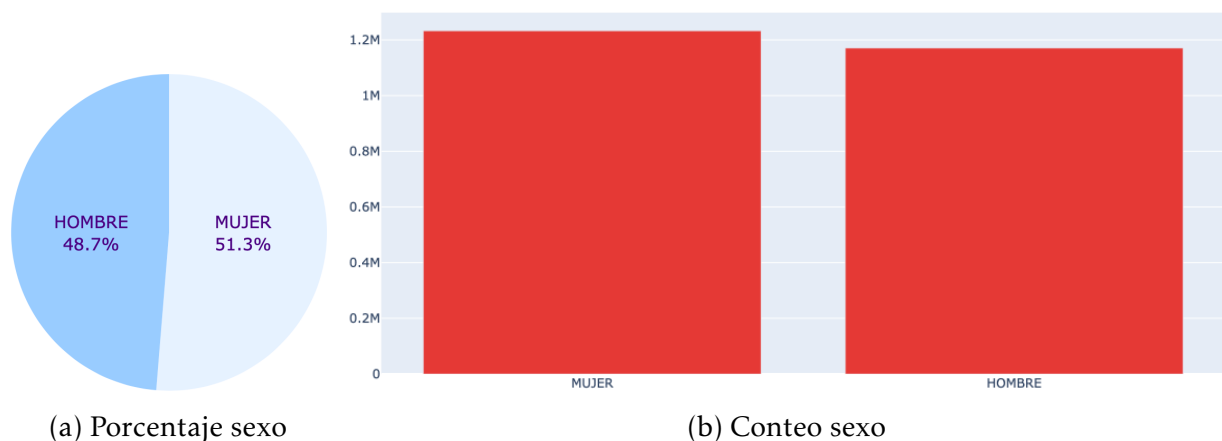


Figura 3: Variable v\_SEXO

Como se puede observar en la Figura 3, la proporción de hombres y mujeres es prácticamente la misma, pero las mujeres presentan un número mayor de casos, siendo este de 1,233,012 contra 1,170,487 casos de hombres.

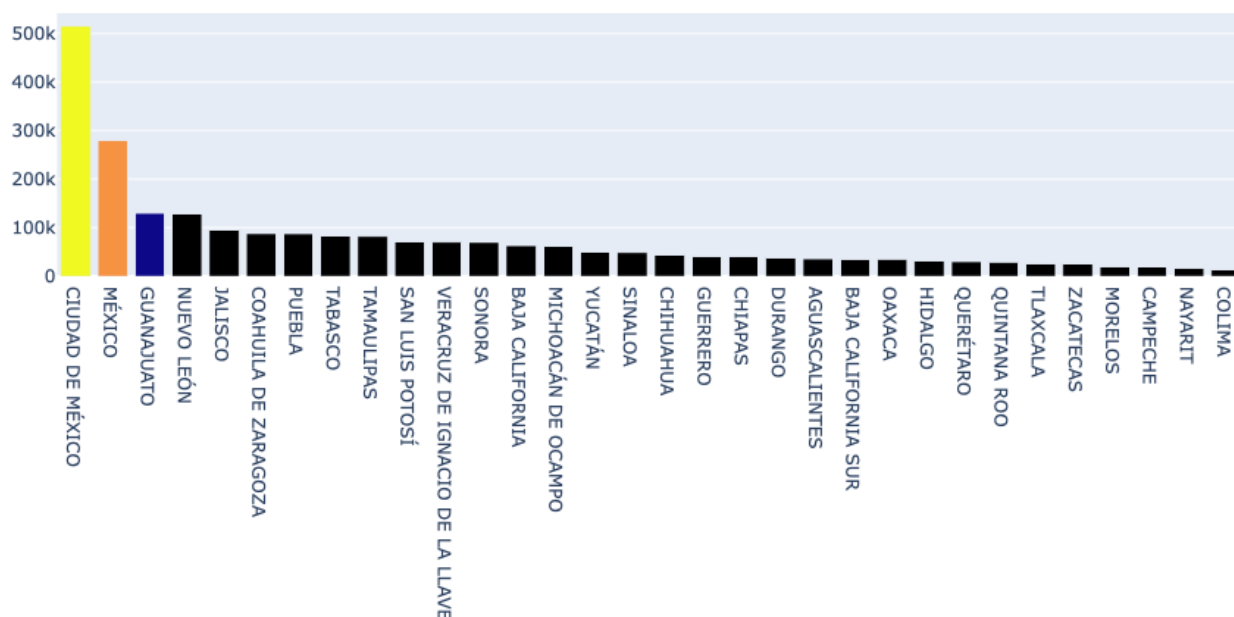


Figura 4: Número de pacientes por cada estado

En la Figura 4, se puede apreciar el número de pacientes que se han reportado desde que inició el registro de los pacientes (01/01/2020) hasta la última actualización de los datos (31/10/2020) siendo la Ciudad de México el estado de la república con más casos (esto es consecuencia de ser la ciudad más poblada del país), seguido del Estado de México y las ciudades más pobladas del país hasta el estado con menos registros de casos que es Colima.

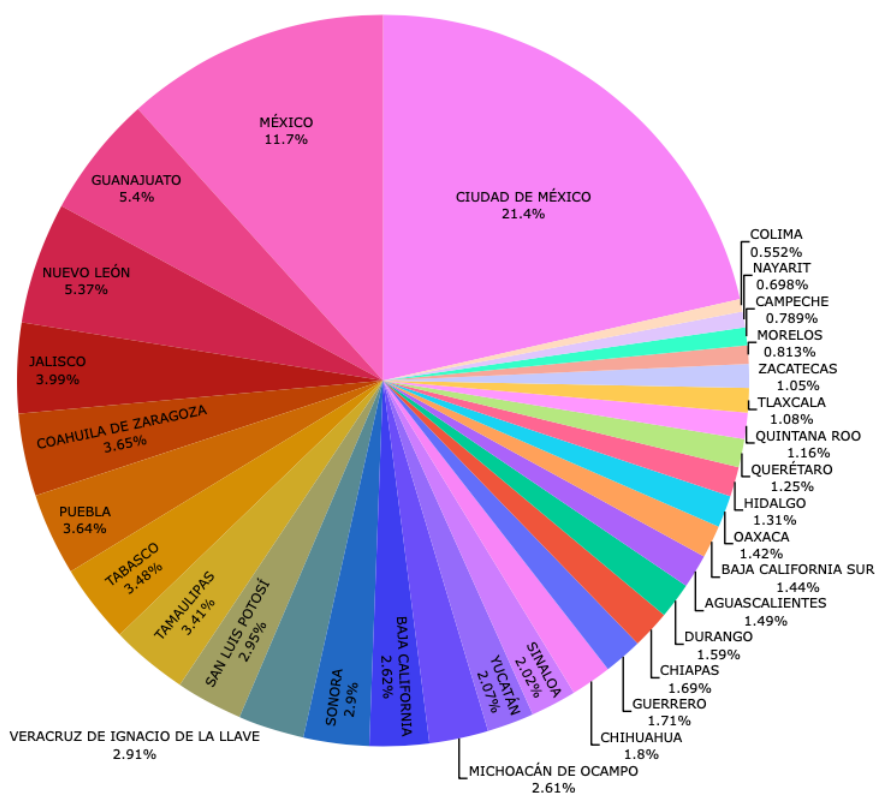


Figura 5: Porcentaje de pacientes por estado



Es preciso mencionar que los datos tienen coherencia como se muestra en la Figura 4 y en la Figura 5, ya que los que tienen mayor número de casos, son aquellos que tienen mayor número de habitantes y viceversa. Aunque, resulta bastante interesante estudiar el caso de Jalisco, ya que no se encuentra posicionado en el top 4, teniendo un mayor número de habitantes, alcanzando la cifra de **8,368,602** para 2020 según datos de la INEGI, cuando Nuevo León y Guanajuato cuenta con aproximadamente 6 millones de habitantes y tienen un mayor número de casos.

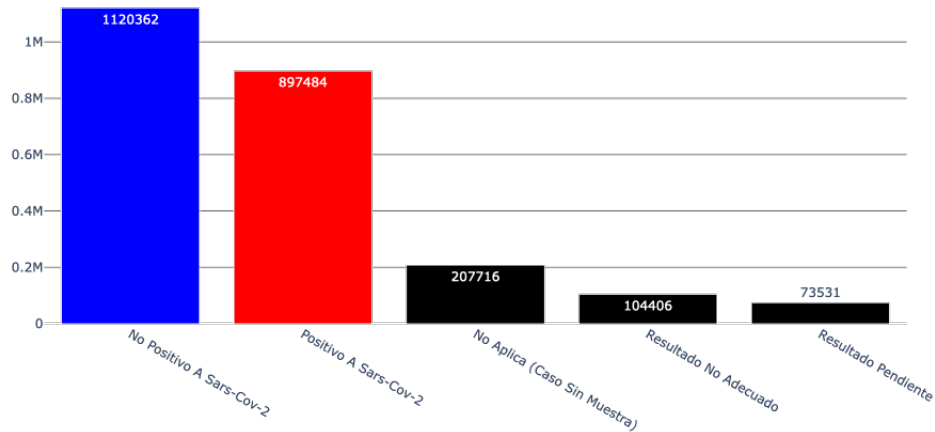


Figura 6: Número de resultados de laboratorio por cada tipo

Con la Figura 6 se pretende analizar el número de personas que se han realizado estudios para saber si tienen COVID-19. Por lo que se observa, la gran mayoría ha salido con un resultado favorable, es decir, no son positivos, pero un número considerable de personas, equivalente al **37%** aproximadamente, han resultado positivos. Esta figura igual permite detectar que con los datos que se cuenta hasta el momento, les hace falta pasar por un proceso de limpieza y de verificación de calidad, ya que, la barra correspondiente a "No aplica" corresponde a lo que serían datos faltantes en la muestra.

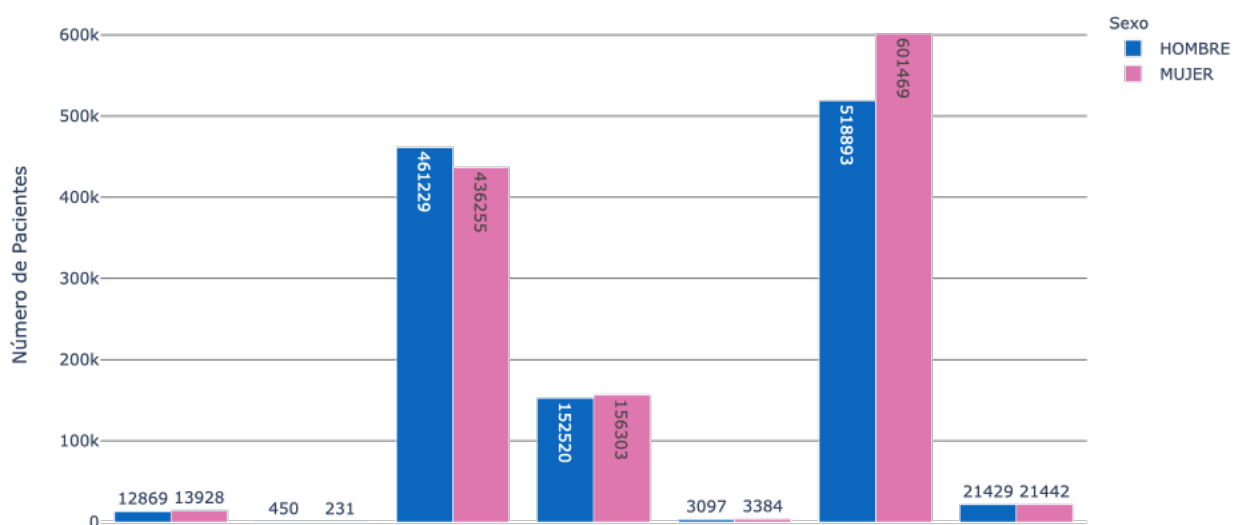


Figura 7: Total de clasificaciones finales por sexo

En la Figura 7, se pueden apreciar las clasificaciones que se tienen hasta el momento previo a la calidad de datos y el número de pacientes por cada sexo. Resultan interesantes

los siguientes datos, en los que se enlistan las clasificaciones finales en el orden que aparecen en el gráfico anterior (de izquierda a derecha):

Clasificación Final	% Hombres	% Mujeres
CASO DE COVID-19 CONFIRMADO POR ASOCIACIÓN CLÍNICA EPIDEMIOLÓGICA	1.1	1.13
CASO DE COVID-19 CONFIRMADO POR COMITÉ DE DICTAMINACIÓN	0.04	0.02
CASO DE SARS-COV-2 CONFIRMADO POR LABORATORIO	39.4	35.38
CASO SOSPECHOSO	13.03	12.68
INVÁLIDO POR LABORATORIO	0.26	0.27
NEGATIVO A SARS-COV-2 POR LABORATORIO	44.33	48.78
NO REALIZADO POR LABORATORIO	1.83	1.74

Cuadro 2: Porcentaje de clasificaciones finales por sexo

Como se puede apreciar en la Figura 7 y en el Cuadro 2, la mayoría de personas que asisten a atenciones médicas por COVID-19 resultan tener clasificaciones negativas y esto es debido a que muchas personas se sugestionan y sumado a lo anterior, los síntomas de esta enfermedad son muy parecidos a los de una gripa o resfriado común, por lo que las personas confunden la sintomatología por los de una enfermedad recurrente.

El tema de la sugestión es un delicado, ya que invade a muchas personas el pensamiento de poder tener COVID-19 y por temor realizan acciones que perjudican al resto, un ejemplo claro de esto es que al inicio de la cuarentena (al menos en el caso de México) la gente se alteró demasiado y empezó a hacer compras de pánico, agotando en muchos establecimientos el papel higiénico y haciendo que este aumentara drásticamente de precio por la escasez y la alta demanda. Por lo que el gobierno (no sólo el de México) han tomado iniciativas para evitar este tipo de acciones, ejemplo de estas fue el no promover el uso de cubrebocas al inicio de la pandemia, ya que si lo hacían se hubiera generado una alta demanda y escasez del producto, provocando que gente como los médicos y enfermeras que realmente los necesitaban en esos momentos, no tuvieran suficiente material.

**Nota:** En el "Apéndice", se encuentra una sección llamada "Visualización de Datos - Comparación" en donde se encuentran el resto de gráficas de las variables con su respectiva comparación (en caso de que hayan sufrido cambios).

### 3. Calidad de Datos

La calidad de datos le permite preparar y gestionar los datos, al tiempo que los pone a disposición de toda su organización. Los datos de alta calidad permiten a los sistemas estratégicos integrar todos los datos relacionados para proporcionar una visión completa de la organización y las interrelaciones dentro de la misma.

La calidad de los datos es una característica esencial que determina la confiabilidad de la toma de decisiones.

Se comenzó con el proceso de calidad de datos verificando que no existieran registros duplicados tanto de manera general como en la variable correspondiente al ID de identificación de paciente, se concluyó que **no** se encontraba ningún registro o paciente duplicado.

### 3.1. Variables Categóricas

Continuando con el proceso, se analizaron las variables dependiendo de su tipificado. Dado que la mayoría de variables eran categóricas, inicié analizando este tipo de variables. Se verificó que cada una de estas tomara valores dentro de su dominio únicamente y se encontró que muchas variables contenían categorías correspondientes a valores que deberían ser nulos, tales se muestran en la siguiente tabla:

Variables	Categoría
v_SECTOR, v_ENTIDAD_NAC, v_INTUBADO, v_NEUMONIA, v_MIGRANTE	no especificado
v_EMBARAZO, v_DIABETES, v_EPOC, v_ASMA, v_HIPERTENSION, v_OBESIDAD	se ignora
v_PAIS_ORIGEN, v_PAIS_NACIONALIDAD	se desconoce

Cuadro 3: Variables y su categoría correspondiente a nulos

Es importante mencionar que las variables mencionadas anteriormente sólo son un subconjunto de aquellas que contenían estos datos, ya que al final se descubrió que 22 de las 32 variables categóricas guardaban datos correspondientes a valores nulos en alguna de estas tres categorías, por lo que tras identificarlos se cambiaron a valores tipo NaN.

Por otro lado, la variable correspondiente a la entidad de residencia del paciente (v\_ENTIDAD\_RES) fueron modificadas sus categorías, esto pensando en la mejor manera de aprovechar la información que contiene, ya que como es una variable categórica, contiene los 32 estados que hay en el país. El problema radica en que cuando se vaya a transformar la información en datos numéricos, se generarían 31 variables nuevas, para no generar tantas variables, se agruparon en cinco regiones, tales que definió el gobierno mexicano al diseñar su estrategia de seguridad a nivel nacional. Las regiones son:

- **noroeste:** baja california, baja california sur, chihuahua, sinaloa, sonora.
- **noreste:** coahuila de zaragoza, durango, nuevo leon, san luis potosi, tamaulipas.
- **occidente:** aguascalientes, colima, guanajuato, jalisco, michoacan de ocampo, nayarit, queretaro, zacatecas.
- **centro:** ciudad de mexico, mexico, guerrero, hidalgo, morelos, puebla, tlaxcala.
- **sureste:** campeche, chiapas, oaxaca, quintana roo, tabasco, veracruz de ignacio de la llave, yucatan.

De esta manera, pasamos de tener que crear 31 variables a sólo **cuatro**, guardando la mayor cantidad de información posible.

### 3.2. Variable Continua

Se analizó los valores que tomaba la variable correspondiente a la edad y como se muestra anteriormente, no presenta valores nulo, pero sí contiene outliers, tales que en un proceso más adelante serán modificados.

### 3.3. Variable Tipo Texto

En este caso no se realizó ninguna análisis extra a la variable correspondiente al ID de indentificación de los pacientes, ya que todos su datos ya venían homologados y sin duplicados.

### 3.4. Variables Tipo Fecha

Para poder hacer uso de la información contenida en las variables correspondientes a fechas (de última actualización, de ingreso del paciente, en que se presentaron síntomas y de defunción) se tuvo que realizar una transformación en el tipo de dato, ya que, en un principio estas variables contenían puros datos tipo string los cuales no permitían el manejo adecuado de la información, por lo que se transformaron a valores tipo "dateTime".

La transformación realizada en cada una de las cuatro variables, provocó que la correspondiente a la fecha de defunción (d\_FECHA\_DEF) tuviera una enorme cantidad de valores tipo NaT, es decir, valores tipo nulos. Esto porque la mayoría de pacientes que son atendidos (correspondientes al 95 % aproximadamente del total de pacientes) no fallecen.

### 3.5. Completitud

El término "Completitud" se refiere a cuando todos los campos y registros están dentro del conjunto de datos, no existen espacios en blanco.

Tras haber verificado la calidad de datos en todos los tipos de variables, es importante verificar como ha cambiado su completitud, ya que, aquellas que **no** pasen el umbral de completitud del 80% serán eliminadas, pues ya tendrían un número alto de valores nulos y si se intentan imputar (término que se explicará más adelante) podría cambiar mucho su distribución.

Analizando el Cuadro 4, es notorio que no todas las variables tienen una completitud del 100% como se creía en un principio, en realidad sólo 10 cuentan con todos sus datos. Nótese que hay dos variables que no alcanzan el umbral establecido para la completitud, pero solo se eliminó "v\_MIGRANTE", dado que la variable de tipo fecha contiene información que es interesante con los pocos datos que posee.

### 3.6. Generación de la Variable Objetivo

La variable objetivo es la dependiente, es decir, es la variable que se intentará predecir o modelar en función del resto de variables independientes.

Variables	Compleitud
v_MIGRANTE	0.397 %
v_HABLA LENGUA INDIG	96.287 %
d_FECHA_DEF	5.44 %
v_OTRO_CASO	87.941 %
v_NEUMONIA	99.238 %
v_ENTIDAD_NAC	99.525 %
v_EMBARAZO	99.617 %
v_DIABETES	99.663 %
v_TABAQUISMO	99.673 %
v_HIPERTENSION	99.684 %

Cuadro 4: Compleitud de las variables más afectadas

En el caso del presente proyecto, se pretende generar una tabla lista para poder modelar con sus datos, y se fijó como variable objetivo la que contiene la clasificación final de los pacientes atendidos por COVID-19, pero como se puede ver a continuación, en un principio dicha variable cuenta con las siguientes categorías:

Clasificación Final	Número de pacientes
CASO DE COVID-19 CONFIRMADO POR ASOCIACIÓN CLÍNICA EPIDEMIOLÓGICA	26,675
CASO DE COVID-19 CONFIRMADO POR COMITÉ DE DICTAMINACIÓN	658
CASO DE SARS-COV-2 CONFIRMADO POR LABORATORIO	894,337
CASO SOSPECHOSO	307,925
INVÁLIDO POR LABORATORIO	6,463
NEGATIVO A SARS-COV-2 POR LABORATORIO	1,117,027
NO REALIZADO POR LABORATORIO	42,730

Cuadro 5: Número de pacientes por clasificaciones originales

Muchas de las clasificaciones anteriores resultan redundantes, ya que, según la descripción de dicha variable se pueden agrupar de la siguiente manera:

- **confirmado:** CASO DE SARS-COV-2 CONFIRMADO POR LABORATORIO, CASO DE COVID-19 CONFIRMADO POR ASOCIACIÓN CLÍNICA EPIDEMIOLÓGICA, CASO DE COVID-19 CONFIRMADO POR COMITÉ DE DICTAMINACIÓN.
- **sospechoso:** CASO SOSPECHOSO, NO REALIZADO POR LABORATORIO, INVÁLIDO POR LABORATORIO.
- **negativo:** NEGATIVO A SARS-COV-2 POR LABORATORIO.

Por lo que la variable objetivo "tgt\_CLASIFICACION\_FINAL" queda de la siguiente manera:

Clasificación Final	Número de pacientes
confirmado	921,670
sospechoso	357,118
negativo	1,117,027

Cuadro 6: Número de pacientes por cada calificación

De esta manera no se cuenta con categorías redundantes, ni tampoco se perdió información relevante para el momento de obtener resultados o conclusiones de los datos.

## 4. Datos Anómalos

Los datos anómalos u outliers (en inglés) son aquellos que tienen características diferentes de la multitud, por lo que resulta importante tratarlos, ya que si no, al modelar con los datos podrían provocar resultados erróneos.

Como en el proceso de calidad de datos ya se había revisado que las variables categóricas y las de tipo fecha tuvieran datos dentro de la naturaleza de su dominio, la única que resta por detectar outliers es en la continua (c\_EDAD), por lo que se le aplicó una función que mediante tres métodos de detección de datos atípicos los cuales son:

- Inter Quantile Range (IQR)
- Z-Score
- Percentiles

La función encuentra los índices de los datos atípicos que al menos hayan sido detectados por dos de los métodos anteriores. En el caso particular de la variable correspondiente a la edad se detectaron **7,684**, lo que representa al **0.32%** del total de los registros, por lo que se procedió a eliminar estas observaciones anómalas.

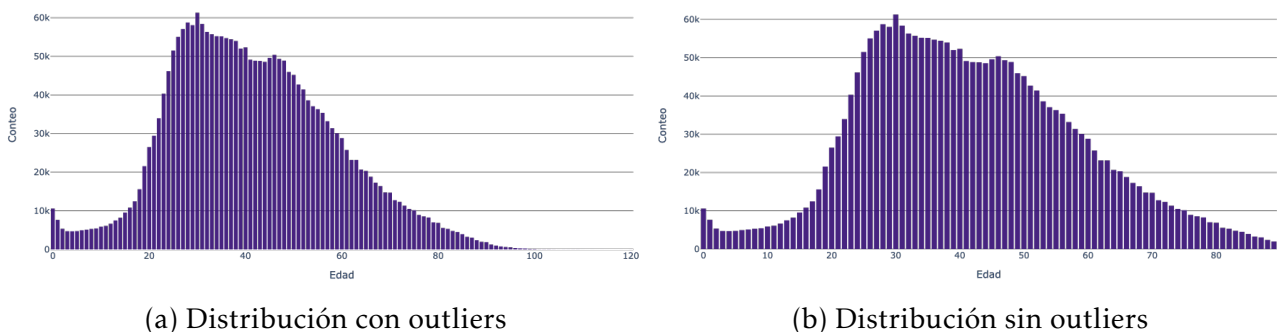


Figura 8: Comparación de la variable c\_EDAD

La Figura 8, muestra el cambio que sufrió la variable edad al detectar y eliminar sus outliers, en un principio, se contaba con una edad máxima de **120**, tras el tratamiento se cuenta con una edad máxima de **89**.

## 5. Análisis Descriptivo Post-Procesamiento de Datos

### 5.1. Visualización de datos

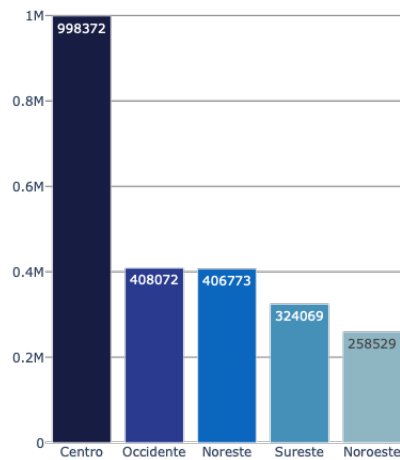


Figura 9: Número de pacientes por regiones

La Figura 9 es el contraste de la Figura 4, porque muestra las nuevas categorías que se crearon para la variable `v_ENTIDAD_RES`, en la que como se puede apreciar, no se perdió demasiada información, dado que si bien no es tan específica como en un principio, nos sigue permitiendo dimensionar la cantidad de pacientes que hay en los estados. Analizando la gráfica se llega a las siguientes conclusiones:

- El 42 % de los pacientes atendidos por COVID-19 provienen de la región centro.
- El número de pacientes de la región occidente y noreste es equivalente al 17% cada una.
- La región sureste representa el 13 % de los pacientes totales.
- La región noroeste representa el 11 %.

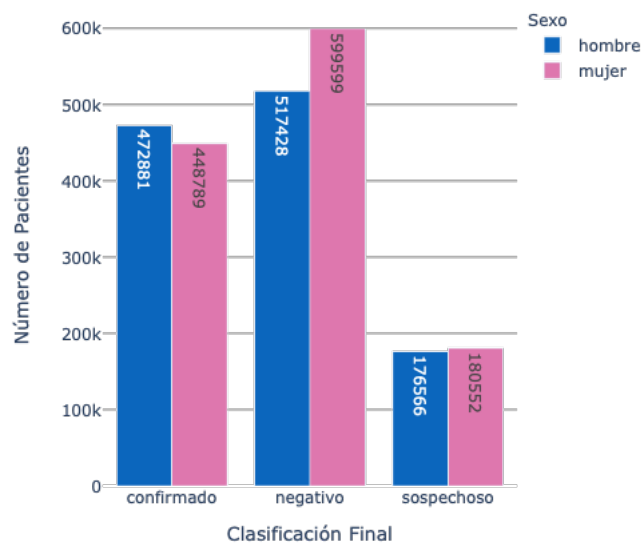


Figura 10: Clasificaciones finales por sexo

Mediante la Figura 10 se puede contrastar la Figura 7, ya que muestra los cambios que se realizaron a la variable objetivo para reducir el número de categorías que eran redundantes. Analizando la gráfica se tiene que el **36 %** de las pacientes ha sido clasificada como caso confirmado, el **49 %** como negativo y el **15 %** como caso sospechoso. En cuanto a los pacientes el **40 %** son casos confirmados, **44 %** casos negativos y **16 %** casos sospechosos.

A partir de lo que se lleva del tratamiento de los datos, es posible generar gráficas utilizando las variables de tiempo, pues en procesos anteriores, fueron tipificadas de manera correcta para realizar estos análisis.

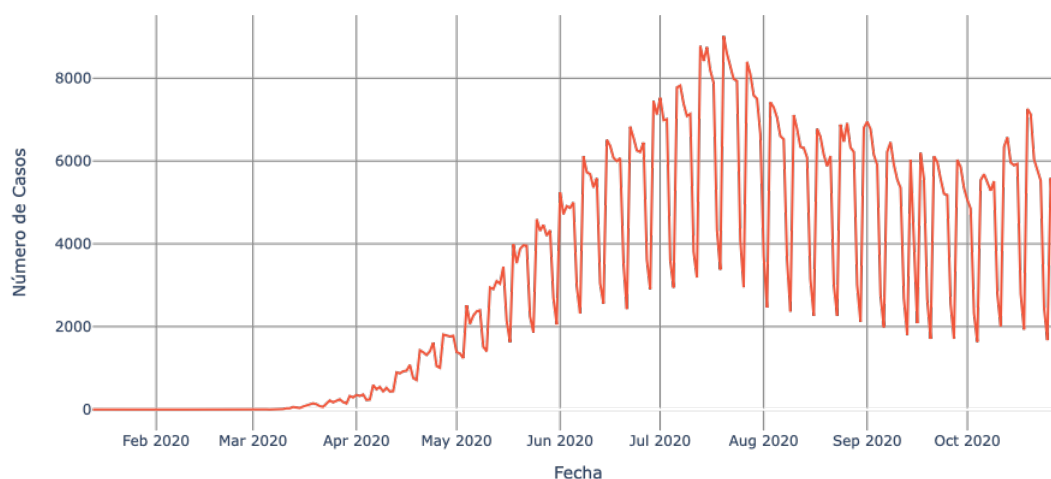


Figura 11: Número de casos confirmados diario

En la Figura 11, se puede visualizar cómo se va comportando el número de casos confirmados durante el paso de los días. Es importante destacar que los datos comienzan desde el 1º de enero hasta el 31 de octubre del 2020. Se logran llegar a las siguientes conclusiones:

- El 20 de julio es el pico más alto con **9,051** casos confirmados.
- El 31 de octubre se alcanza el pico más bajo con **296** casos confirmados.
- En promedio al día se tienen **3,168** nuevos casos confirmados.

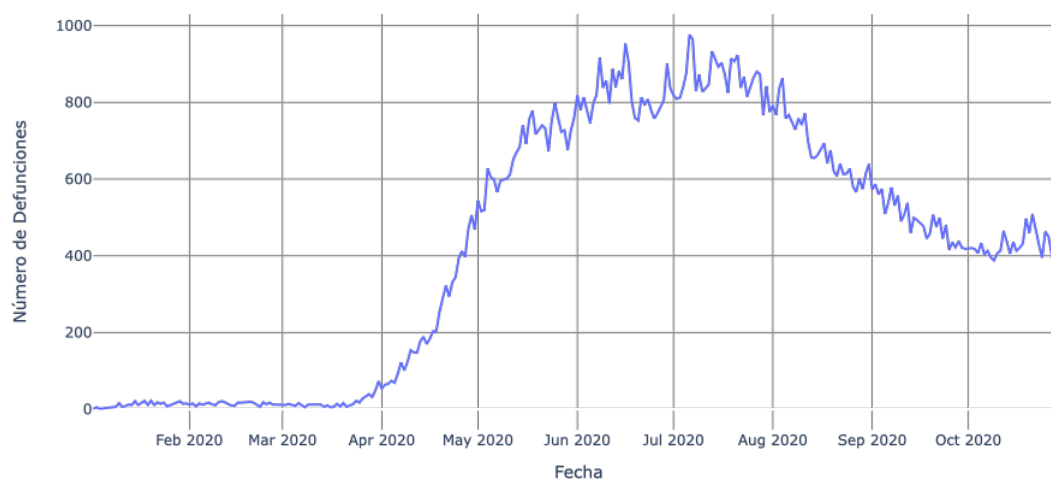


Figura 12: Número de fallecimientos diarios



Por otra parte, en la Figura 11 se puede apreciar el comportamiento del número de fallecimientos de pacientes que fueron atendidos por COVID-19 con el paso de los días y se pueden resaltar los siguientes aspectos:

- El 6 de julio se alcanzó el pico más alto con **985** defunciones.
- El 31 de octubre se alcanza el pico más bajo con **1** defunción.
- En promedio al día se tienen **431** defunciones.

## 6. Imputación de Valores Faltantes

La imputación de valores faltantes consiste en métodos que permiten aumentar la completitud de las variables, logrando que alcancen el 100 %.

Antes de proceder a imputar valores, es necesario conocer los datos para detectar el método a utilizar, en el caso presente se trabaja con una serie de tiempo pero dado que se imputaron valores a variables categóricas, no es necesario verificar si presenta tendencia o no. Partiendo de lo anterior se decidió imputar los valores haciendo uso de la moda, es decir, tomando el valor o categoría que más se repite.

A partir de este punto, se dividió el conjunto de datos en dos subconjuntos: prueba y entrenamiento. El primero contiene el **80 %** y el segundo el **20 %** restante. Es importante aclarar que los registros de cada subconjunto son seleccionados por orden, es decir, para el caso del conjunto de entrenamiento se seleccionan los registros considerando el orden de la fecha de ingreso del paciente y hasta el registro en el que se acumule el 80% de la tabla original. Se realiza de esta manera, ya que la moda que se obtiene es sobre el conjunto del entrenamiento y es valor que se le imputa a ambos conjuntos.

Resultan importantes estos subconjuntos porque cuando se vaya a modelar con los datos, se entrenará al modelo con el conjunto de entrenamiento y se evaluará el desempeño de predicción del modelo con el conjunto de prueba.

<b>Variables</b>	<b>Moda</b>
v_SECTOR	ssa
v_ENTIDAD_NAC	ciudad de mexico
v_INTUBADO	no aplica
v_NEUMONIA	no
v_EMBARAZO	no
v_HABLA LENGUA INDIG	no
v_DIABETES	no
v_EPOC	no
v_ASMA	no
v_INMUSUPR	no
v_HIPERTENSION	no

Cuadro 7: Variables y su moda

**Nota:** El resto de variables y su respectiva moda se encuentra en el "Apéndice" en la sección de "Imputación de Valores Faltantes - Restantes".

El Cuadro 7 muestra un subconjunto de las variables que fueron completadas y el valor que les fue imputado.

## 7. Ingeniería de Variables

La ingeniería de variables consiste en generar nuevas variables con las que ya se cuenta en un principio con el propósito de conservar la mayor cantidad de información posible o extraer nueva información.

### 7.1. Variables Categóricas

Este proceso consistió en transformar todas las variables categóricas, tales que contienen valores tipo "string", a variables que sólo contengan valores de tipo **numérico**. Para realizar el objetivo anterior, se crearon **variables dummies** (variables que sólo contienen 1 si ocurre el evento que describen ó 0 si no ocurre el evento que describen) de casi todas las variables categóricas menos de: v\_ENTIDAD\_UM, v\_ENTIDAD\_NAC, v\_RESULTADO\_LAB, v\_PAIS\_ORIGEN, v\_MUNICIPIO\_RES, v\_PAIS\_NACIONALIDAD. Debido a que estas seis tendrán un proceso diferente que será explicado más adelante.

En resumen, se pasó de tener 36 variables a terminar con un total de **51**, de las cuales 10 no fueron modificadas (las seis mencionadas anteriormente y las cuatro tipo fecha), las 26 restantes fueron reemplazadas por 41 variables nuevas que contienen la información de las variables en forma numérica.

### 7.2. Variable Objetivo

Debido a que la variable objetivo es de tipo categórica, para poder usarla en un modelo estadístico, es necesario transformar sus valores a tipo numérico, por lo que nuevamente se generaron variables dummies. En este caso, se pasó de tener la variable "tgt\_CLASIFICACION\_FINAL" únicamente a tener dos variables objetivo "tgt\_CLASIFICACION\_FINAL\_negativo" y "tgt\_CLASIFICACION\_FINAL\_sospechoso", tal que la información queda expresada como:

- **confirmado:** (0, 0)
- **sospechoso:** (0, 1)
- **negativo:** (1, 0)

Correspondiendo a la primer entrada del vector como el valor que toma la variable "tgt\_CLASIFICACION\_FINAL\_negativo" en el registro  $i$  y la segunda entrada del vector al valor de la variable "tgt\_CLASIFICACION\_FINAL\_sospechoso" en el registro  $i$ .

## 8. Reducción de Dimensiones

La reducción de dimensiones tiene como objetivo que mediante la proyección de los datos a un subespacio de menor dimensión se plantea captar la “**esencia**” de los datos. En otras palabras consiste en reducir el número de variables independientes conservando la mayor cantidad de información de los datos.

Tras haber creado variables dummies, el conjunto de datos aumentó significativamente en cuanto al número de variables, por lo que se aplicaron diferentes métodos que ayudaran a medir su desempeño en conjunto e individual.

Recuérdese que para este punto sólo hemos eliminado la variable `v_MIGRANTE` porque tras aplicarle el método de **relación de valor perdido**, no alcanzó el umbral de completitud.

Se eliminaron las variables `v_ENTIDAD_UM`, `v_ENTIDAD_NAC`, `v_PAIS_ORIGEN`, `v_PAIS_NACIONALIDAD` y `v_MUNICIPIO_RES` debido a que no aportan mucha información que sea útil para la variable objetivo, y a su vez, la información relevante de estas variables se encuentra almacenada en las variables `v_ENTIDAD_RES` y `v_NACIONALIDAD`.

Por otro lado, no aporta mayor información la variable `v_RESULTADO_LAB` por lo que fue eliminada, ya que la variable objetivo contiene prácticamente la misma información y esta es la importante. A su vez, la variable de texto que contiene el ID que identifica a cada paciente (`t_ID_REGISTRO`) fue eliminada porque no aporta información que sea de utilidad para la variable objetivo. En cuanto a las variables de tipo fecha, hay una que no aporta información a la variable dependiente y esta es `d_FECHA_ACTUALIZACION` que indica el último día que se actualizaron los datos (en este caso es en el día en el que se descargaron) y es constante para todos los registros.

Ahora se aplicará un método llamado **filtro de alta correlación**, el cual consiste en eliminar las variables que resulten redundantes, pues contienen información que otras variables guardan.

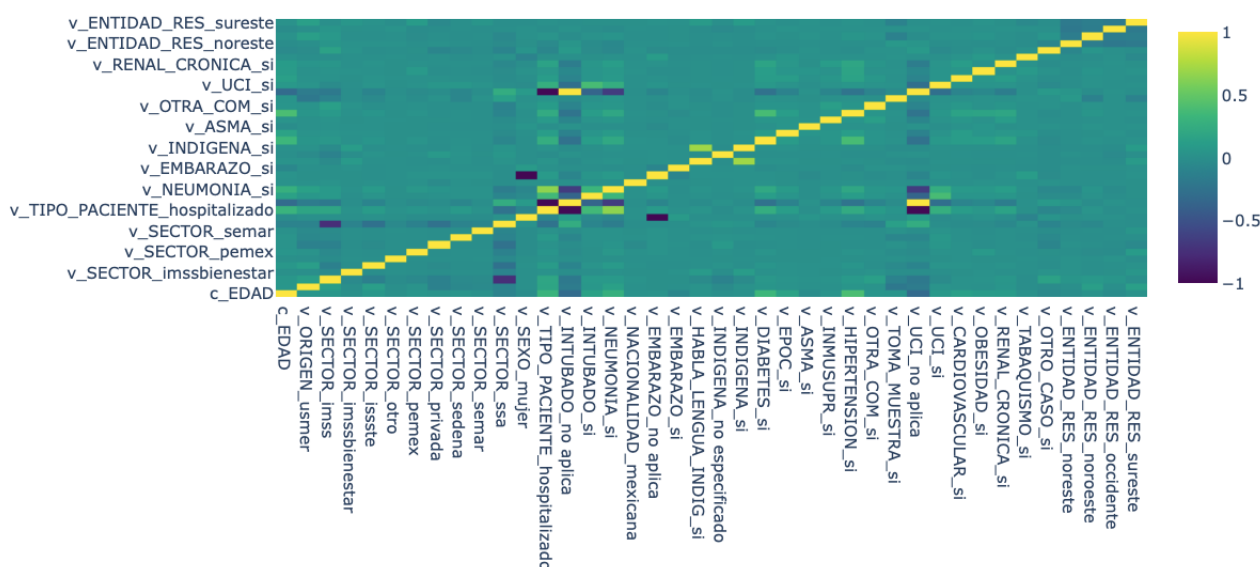


Figura 13: Correlaciones entre variable independientes

Analizando la Figura 13, se eliminaron **seis** variables que tienen su correlación mayor a 0.6 o menos a -0.6 con alguna otra variable.

- Se eliminó la variable v\_SECTOR\_ssa pues presentaba una alta correlación negativa con la variable v\_SECTOR\_imss, esto quiere decir, al aumentar el número de pacientes del imss, disminuyen los pacientes de la Secretaria de Salubridad y Asistencia de manera proporcional.
- Se eliminó la variable v\_EMBARAZO\_no aplica, dado que tiene una alta correlación negativa con la variable v\_SEXO\_mujer.
- Se eliminaron las variables v\_INTUBADO\_no aplica y v\_UCI\_no aplica por tener una correlación alta negativa con la variable v\_TIPO\_PACIENTE\_hospitalizado. Por otro lado, también presenta una alta correlación con esta última la variable v\_NEUMONIA\_si, es decir, si aumentan los pacientes de tipo hospitalizados igual aumentará el número de pacientes que tiene neumonía.
- Se eliminó la variable v\_HABLA LENGUA\_INDIG\_si dado que tiene una alta correlación con la variable v\_INDIGENA\_si.

Con el método anterior, se midió el desempeño de las variables independientes por parejas, así que, se pretende analizar ahora su desempeño de todas en conjunto, por lo que se procedió a medir su nivel de **multicolinealidad**, para esto se usa una métrica llamada Factor de Inflación de la Varianza (por sus siglas en inglés VIF).

Variable	VIF
v_NACIONALIDAD_mexicana	22.211
v_TOMA_MUESTRA_si	10.821

Cuadro 8: Variables y su respectivo VIF

Para este método se consideró un umbral para el VIF de 10, es decir, aquellas variables que obtuvieran un VIF mayor que el umbral establecido serían eliminadas, ya que la información que almacenan se puede obtener mediante la combinación lineal del resto de variables.

En resumen de esta sección, se logró reducir de 51 variables a **35**, es decir, se eliminó un total de 17 (considerando la que fue eliminada tras revisar la completitud).

## 9. Conclusión

En conclusión, se logró cumplir con el objetivo principal, el cual fue crear una tabla analítica de datos, donde su contenido ya estuviera preparado para poder realizar modelos estadísticos sobre ellos. Se cumplió tras haber transformado a valores numéricos todos los registros y conservar variables que fueran de utilidad para predecir la variable objetivo. Al igual que se separó de manera adecuada la información en conjunto de entrenamiento y de prueba.

Por otra parte, durante este proceso se logró extraer y presentar información que resultara de interés para el tema, tales como el porcentaje de pacientes por cada estado, identificación del estado con más pacientes, día en el que se presentaron más casos confirmados y día en el que se presentó el mayor número de defunciones, entre muchos otros.

Finalmente, también se logró visualizar los cambios realizados durante el procesamiento de datos y se explicó la razón y el medio por el que se realizaron.

## 10. Apéndice

### 10.1. Variables Restantes

En esta sección encontramos el resto de variables no mencionadas anteriormente pero que resulta también de interés su análisis para tener información y resultados completos.

Variable	Tipo	Descripción
INDIGENA	Categórica	Identifica si el paciente se autoidentifica como una persona indígena.
DIABETES	Categórica	Identifica si el paciente tiene un diagnóstico de diabetes.
EPOC	Categórica	Identifica si el paciente tiene un diagnóstico de EPOC.
ASMA	Categórica	Identifica si el paciente tiene un diagnóstico de asma.
INMUSUPR	Categórica	Identifica si el paciente presenta inmunosupresión.
HIPERTENSION	Categórica	Identifica si el paciente tiene un diagnóstico de hipertensión.
OTRAS_COM	Categórica	Identifica si el paciente tiene diagnóstico de otras enfermedades.
CARDIOVASCULAR	Categórica	Identifica si el paciente tiene un diagnóstico de enfermedades cardiovasculares.
OBESIDAD	Categórica	Identifica si el paciente tiene diagnóstico de obesidad.
RENAL_CRONICA	Categórica	Identifica si el paciente tiene diagnóstico de insuficiencia renal crónica.
TABAQUISMO	Categórica	Identifica si el paciente tiene hábito de tabaquismo.
OTRO_CASO	Categórica	Identifica si el paciente tuvo contacto con algún otro caso diagnosticado con SARS CoV-2.
TOMA_MUESTRA	Categórica	Identifica si al paciente se le tomó muestra.
MIGRANTE	Categórica	Identifica si el paciente es una persona migrante.
UCI	Categórica	Identifica si el paciente requirió ingresar a una Unidad de Cuidados Intensivos.

RESULTADO.LAB	Categórica	Identifica el resultado del análisis de la muestra reportado por el laboratorio de la Red Nacional de Laboratorios de Vigilancia Epidemiológica (INDRE, LESP y LAVE) y laboratorios privados avalados por el INDRE cuyos resultados son registrados en SISVER.
PAIS_NACIONALIDAD	Categórica	Identifica la nacionalidad del paciente.
PAIS_ORIGEN	Categórica	Identifica el país del que partió el paciente rumbo a México.

## 10.2. Imputación de Valores Faltantes - Restantes

Variables	Moda
v_OTRA_COM	no
v_UCI	no aplica
v_CARDIOVASCULAR	no
v_OBESIDAD	no
v_RENAL_CRONICA	no
v_TABAQUISMO	no
v_OTRO_CASO	si
v_PAIS_ORIGEN	no aplica
v_PAIS_NACIONALIDAD	maxico
v_MUNICIPIO_RES	iztapalapa

## 10.3. Visualización de Datos - Comparación

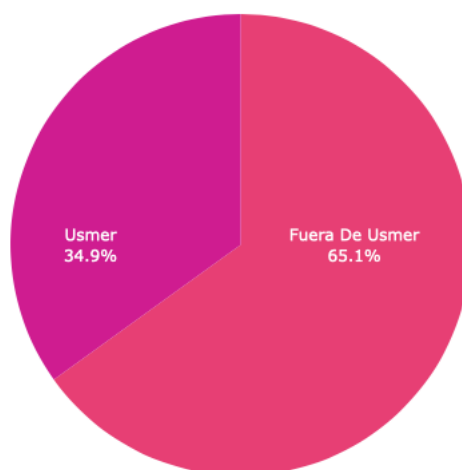


Figura 14: Porcentaje del origen de los pacientes

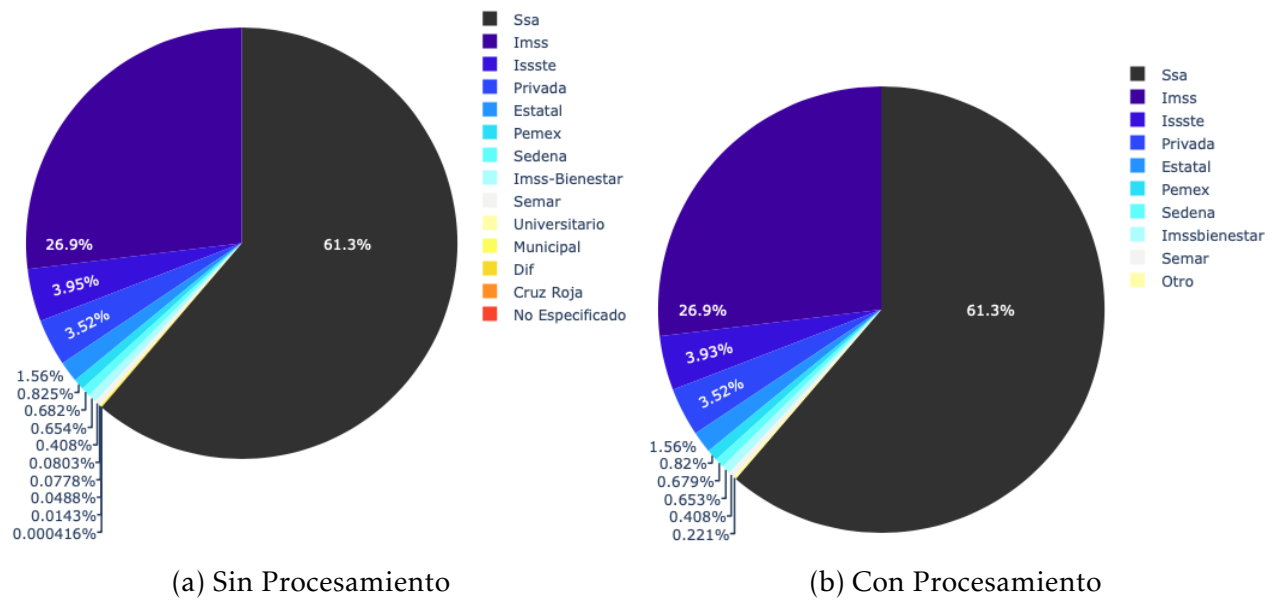


Figura 15: Comparación de la variable v\_SECTOR

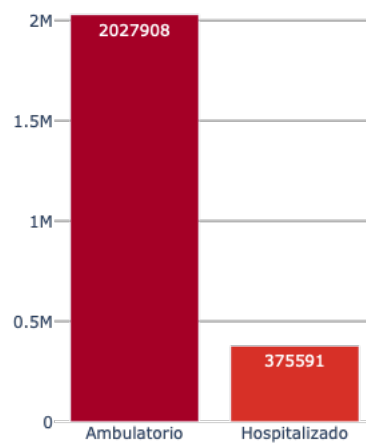


Figura 16: Número de pacientes por cada tipo

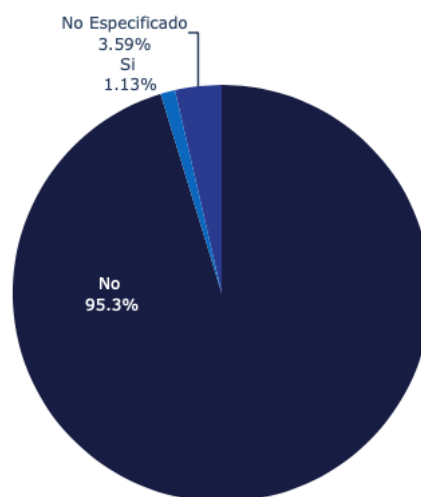


Figura 17: Porcentaje de pacientes que son indígenas

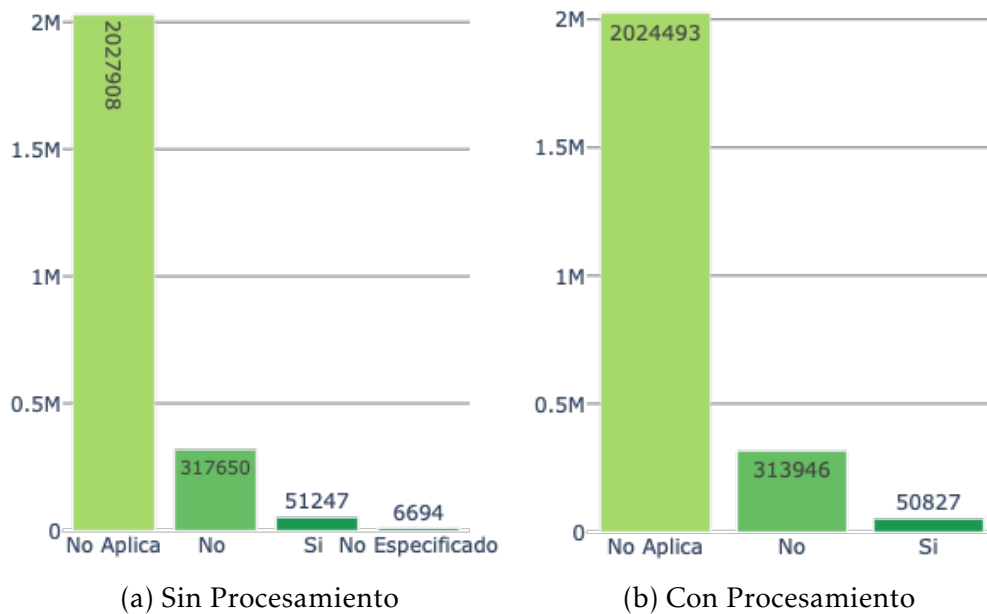


Figura 18: Comparación de la variable v\_INTUBADO

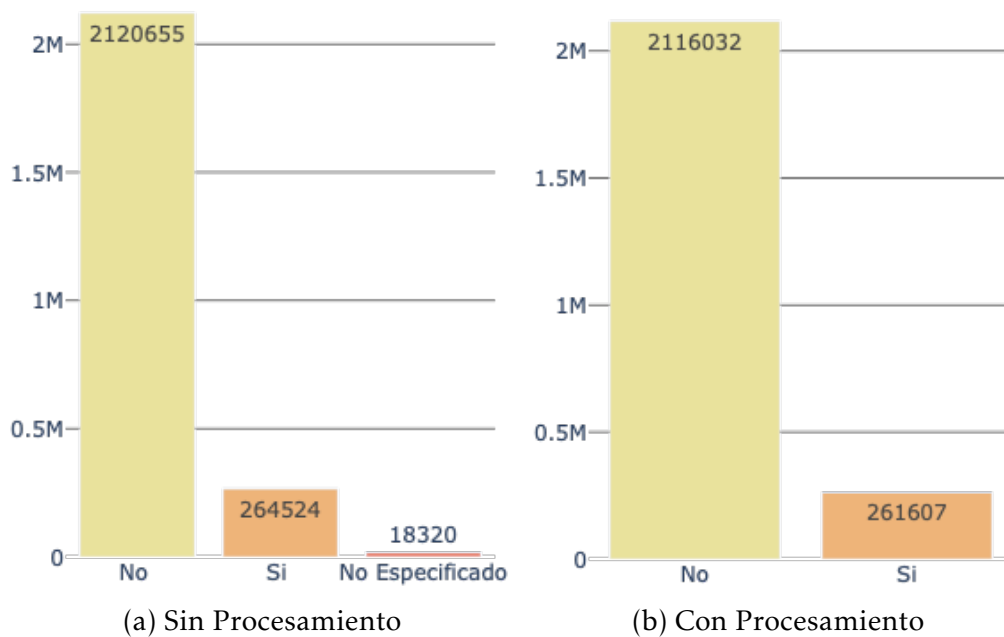


Figura 19: Comparación de la variable v\_NEUMONIA



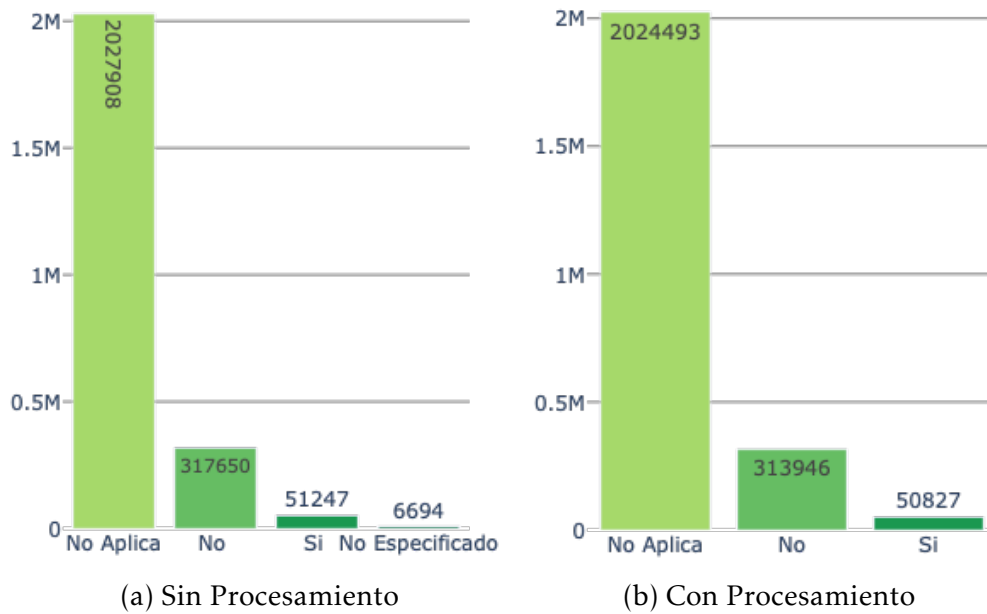


Figura 20: Comparación de la variable v\_INTUBADO

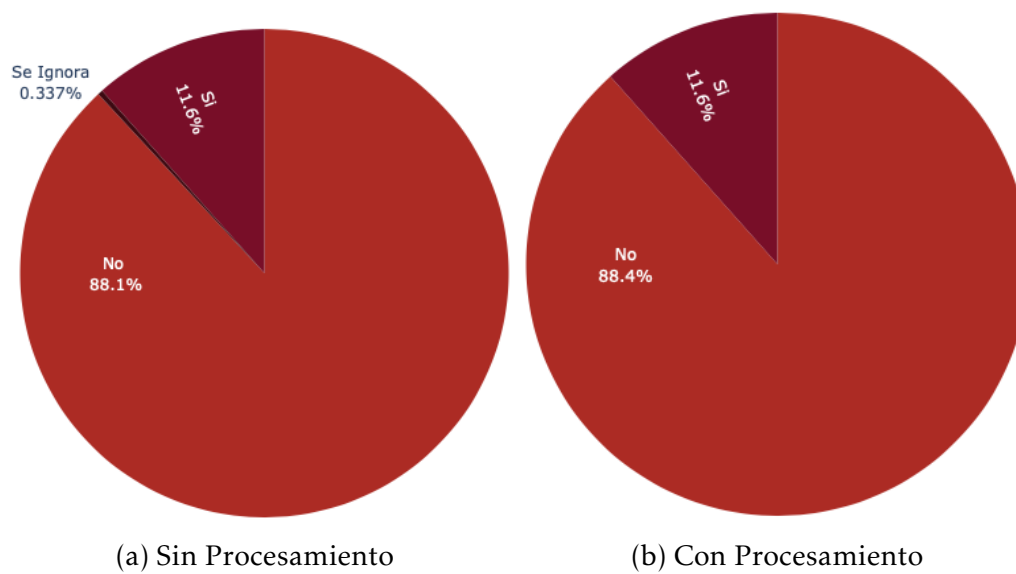


Figura 21: Comparación de la variable v\_DIABETES

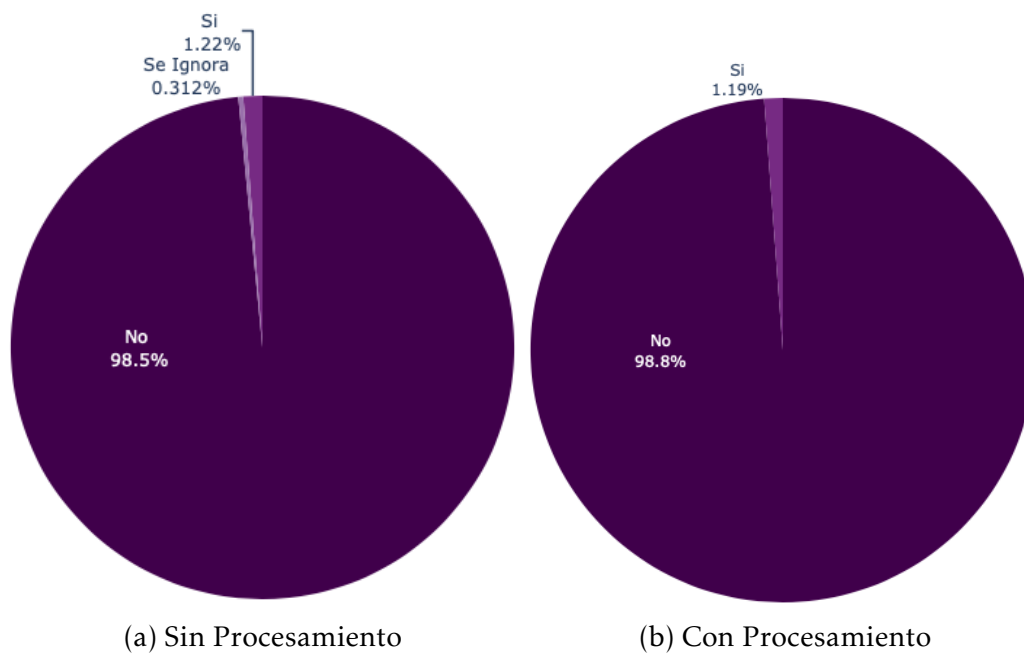


Figura 22: Comparación de la variable  $v\_EPOC$

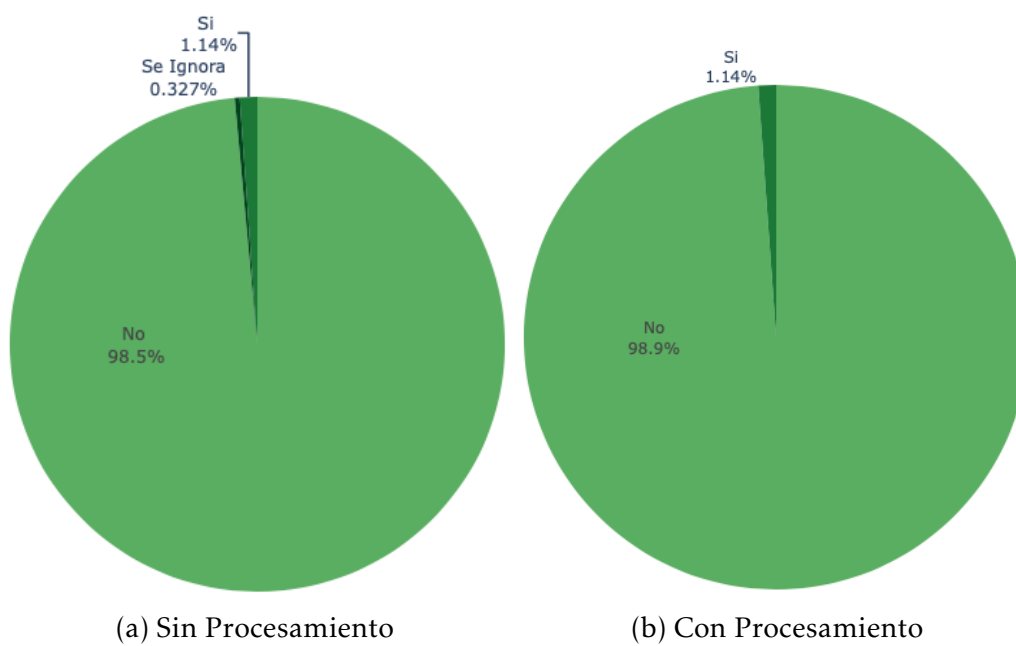


Figura 23: Comparación de la variable  $v\_INMUSUPR$

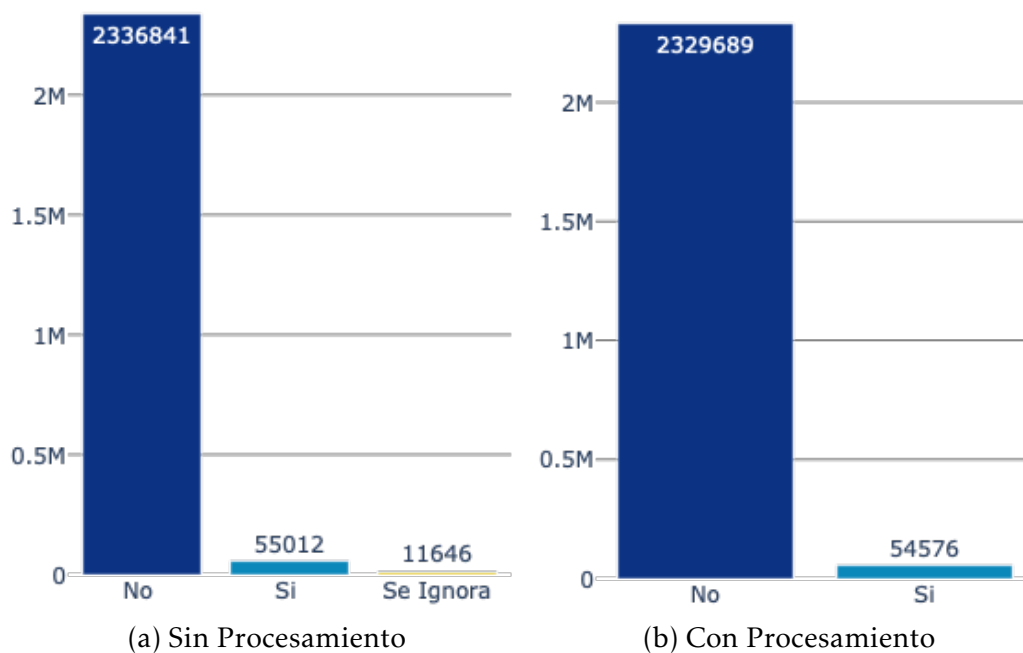


Figura 24: Comparación de la variable `v_OTRA_COM`

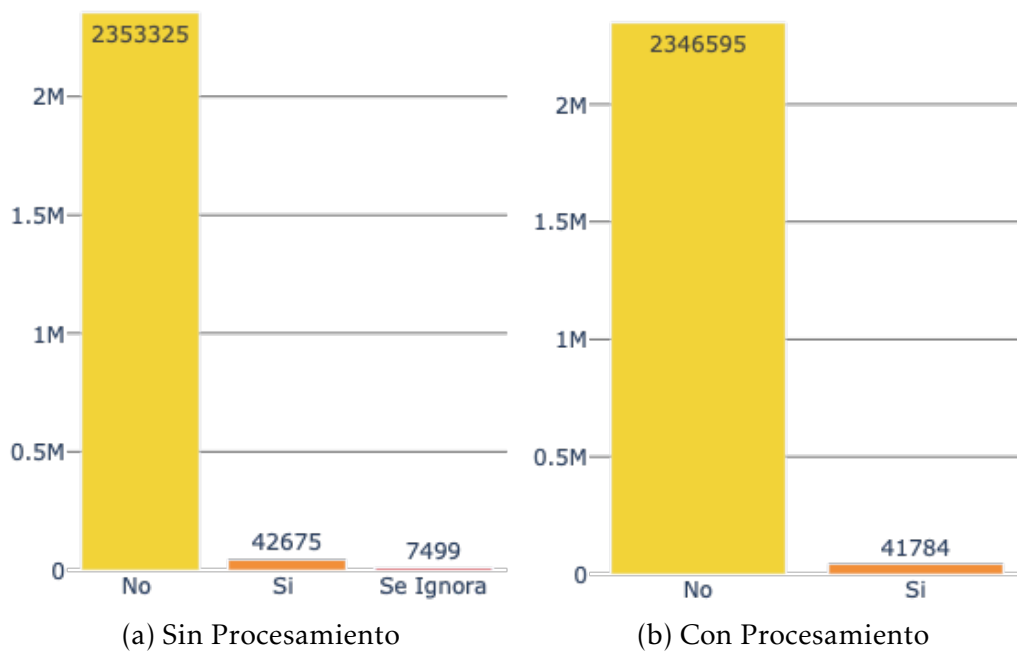


Figura 25: Comparación de la variable `v_CARDIOVASCULAR`

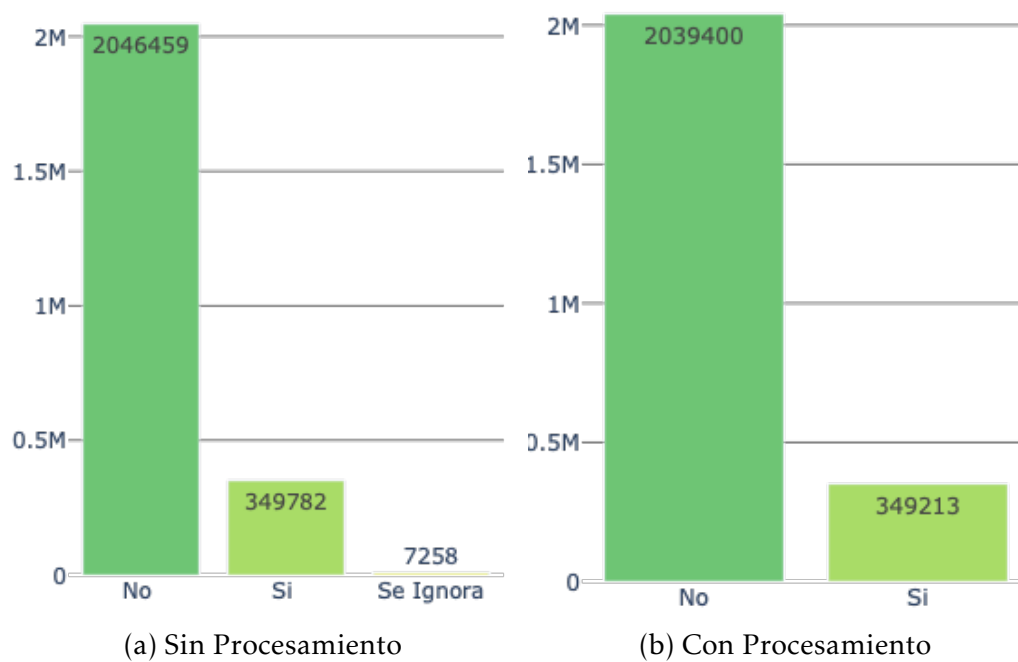


Figura 26: Comparación de la variable `v_OBESIDAD`

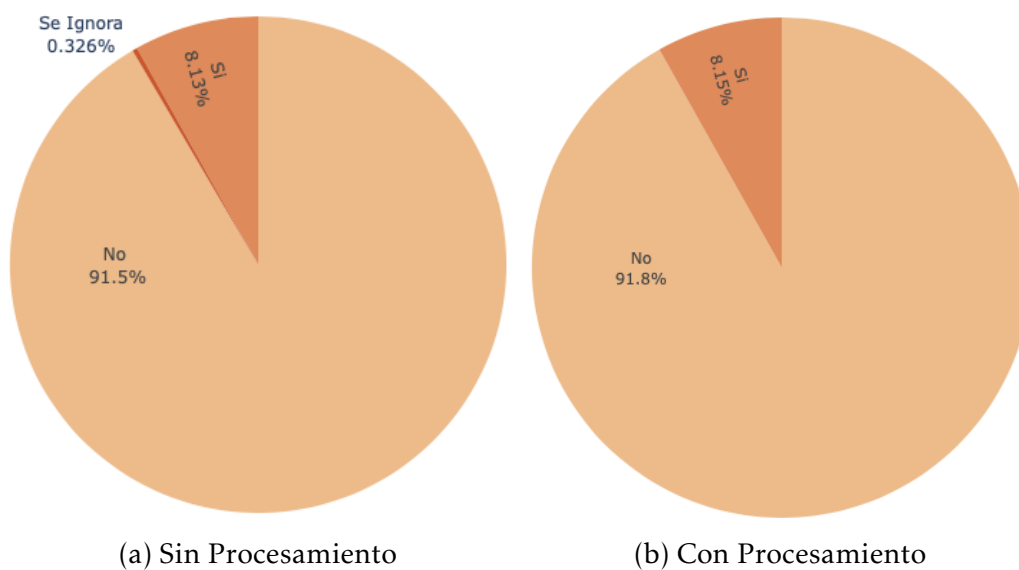


Figura 27: Comparación de la variable `v_TABAQUISMO`

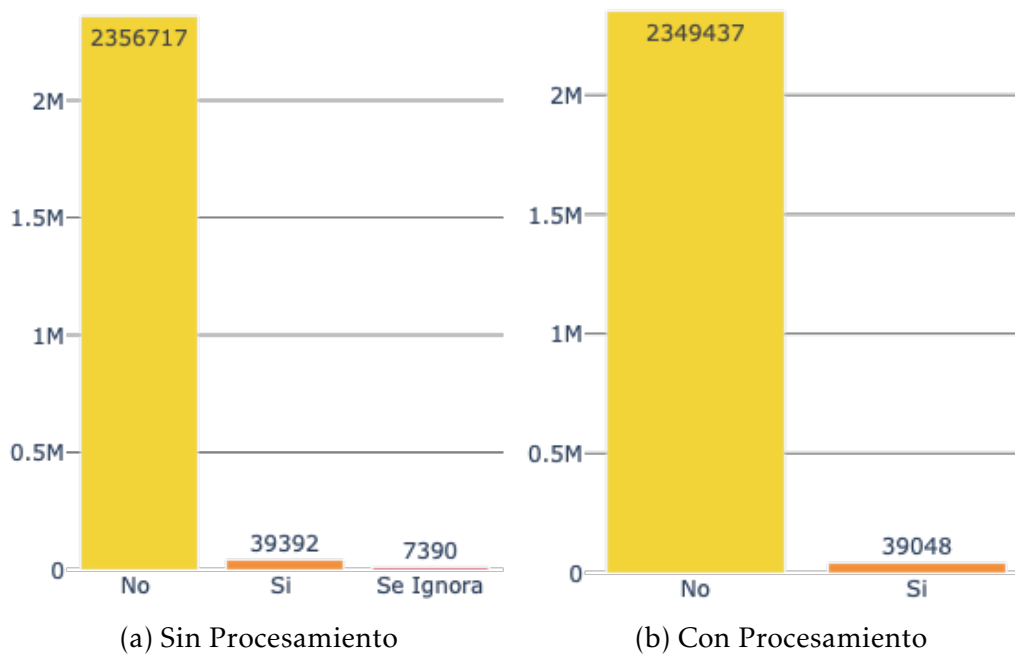


Figura 28: Comparación de la variable `v_RENAL_CRONICA`

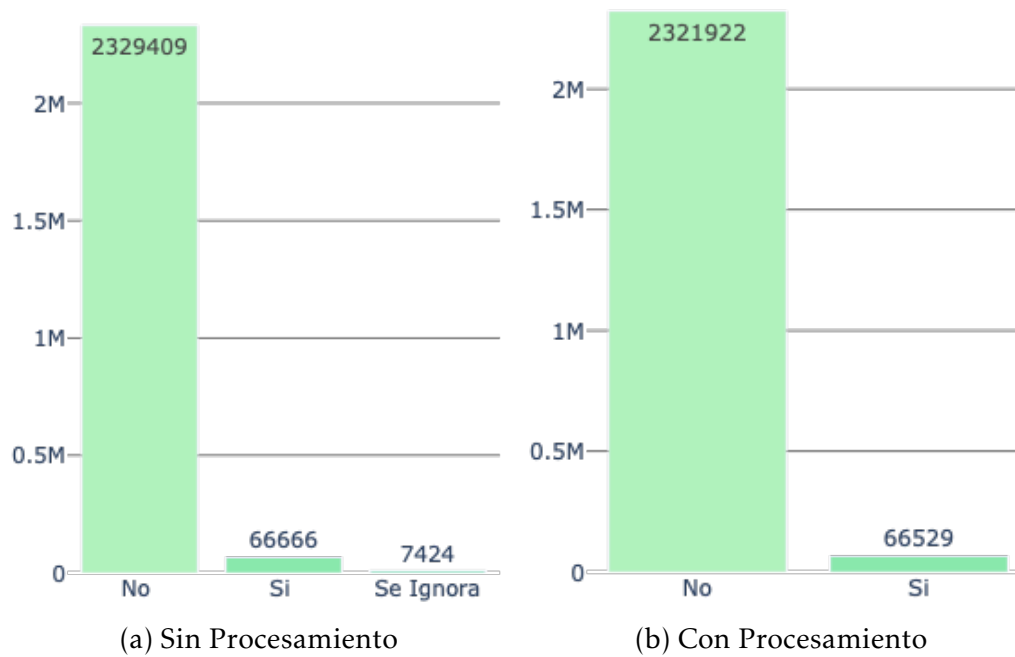


Figura 29: Comparación de la variable `v_ASMA`

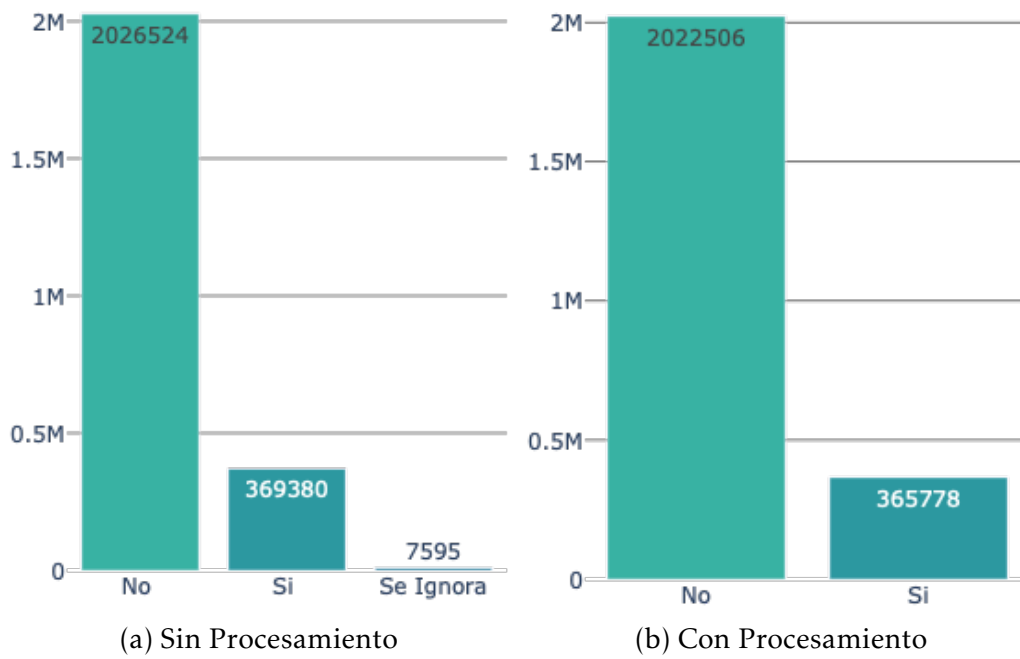


Figura 30: Comparación de la variable `v_HIPERTENSION`

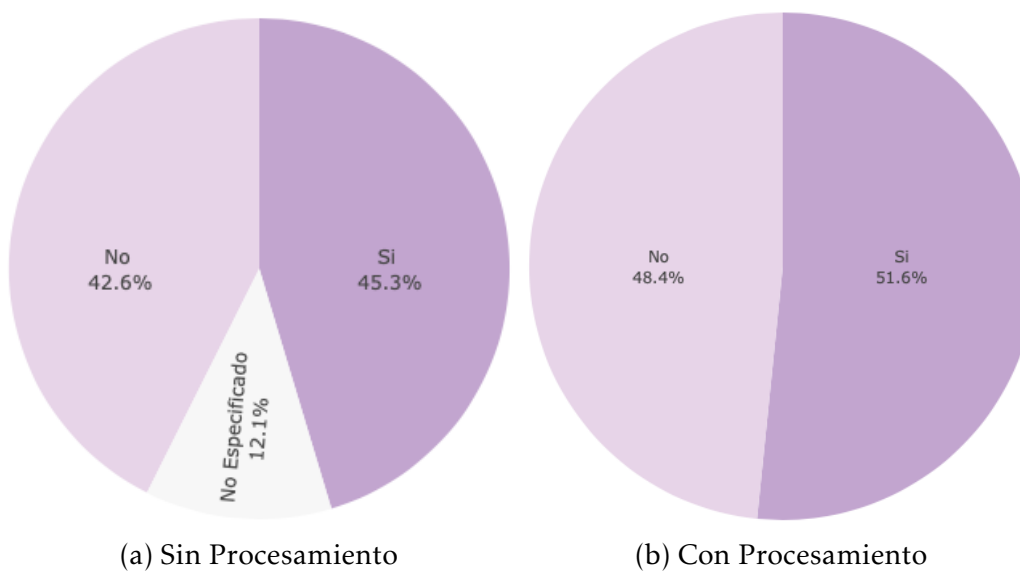


Figura 31: Comparación de la variable `v_OTRO_CASO`

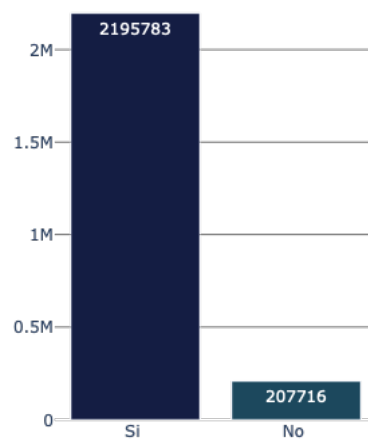


Figura 32: Numero de pacientes con que se les tomó muestra

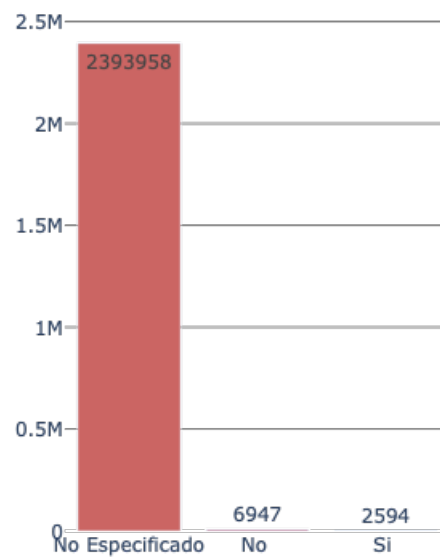


Figura 33: Numero de pacientes que son migrantes

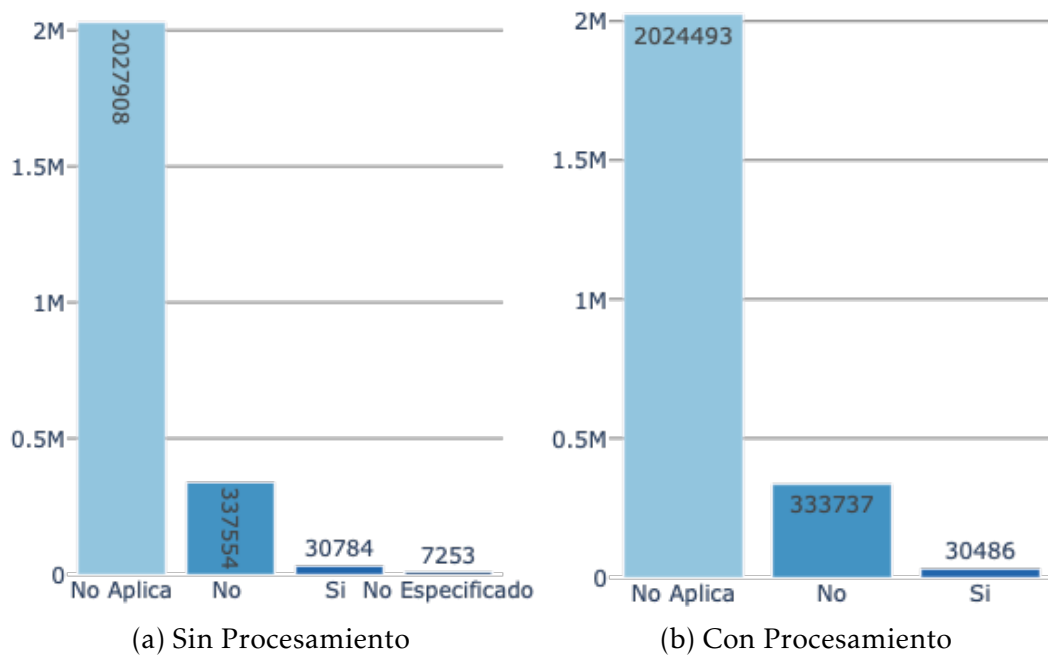


Figura 34: Comparación de la variable v\_UCI