

**Instituto Tecnológico y de Estudios Superiores de Monterrey**

**Campus Guadalajara**



**Inteligencia artificial avanzada para la ciencia de datos I  
(Gpo 101)**

**M1.2 Datos Faltantes y Outliers**

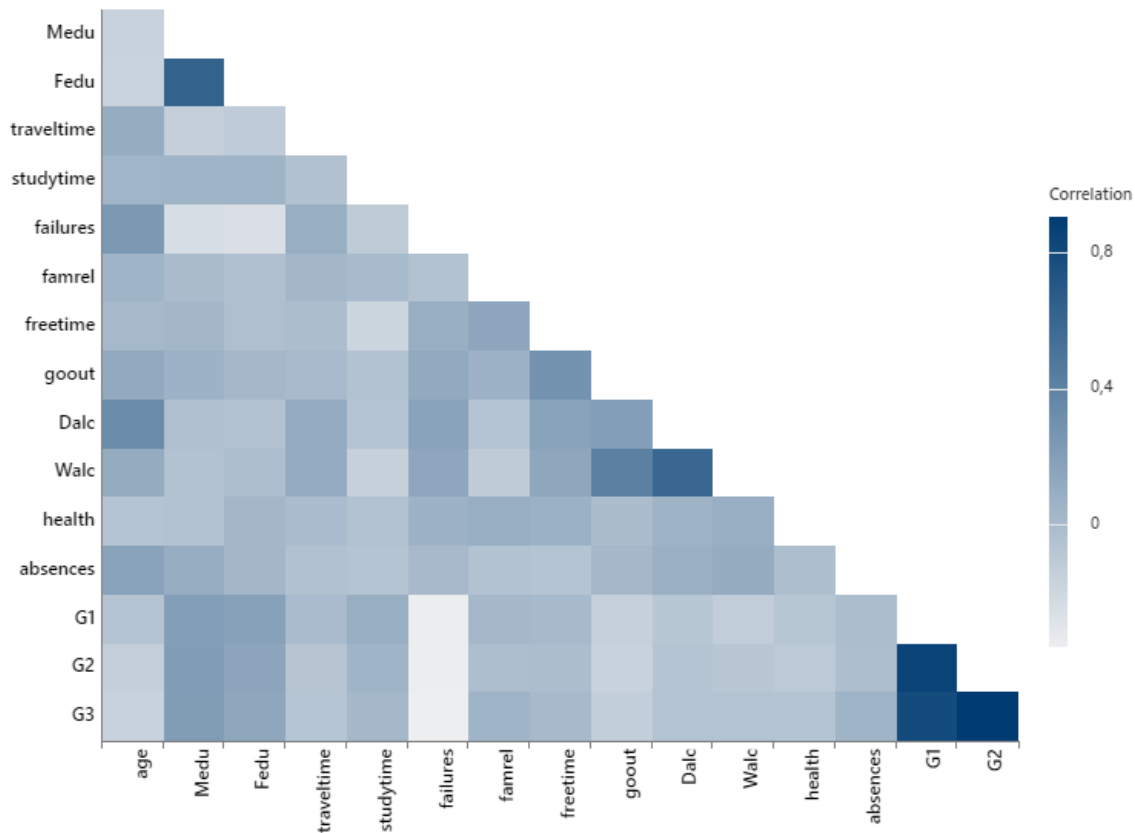
**Samuel Padilla Esqueda**

**| A01641383**

**Agosto 2024**



Walc	0,598														
health	0,057	0,092													
absences	0,077	0,117	-0,020												
G1	-0,080	-0,126	-0,073	-0,011											
G2	-0,060	-0,085	-0,098	-0,018	0,852										
G3	-0,057	-0,052	-0,061	0,049	0,801	0,905									



Podemos observar que los valores de correlación de **traveltime** con el resto de las variables de muy pequeño, su correlación mas grande es 0,112 con la variable age, pero es un valor tan cercano a 0 que implica una correlación muy débil. Ya que no depende de otra variable, descartamos el mecanismo MAR, por lo que puede ser que aplique el modelo MCAR o NMAR.

La variable **Absences** tiene un caso parecido, donde la mayor correlación se observa con la variable age, con un valor de 0,173. Aunque es un valor muy cercano a 0, también es significativamente mayor que el resto de correlaciones por lo que creo que si puede haber una correlación, aunque no sea muy fuerte, de esto modo descartamos el método MAR, ya que hay evidencia de que Absences depende de otra variable, es posible que se trate de un mecanismo MCAR o MAR.

**3. Obtener estadísticas descriptivas de los datos (histograma, media, desviación estándar, mediana, moda, etc).**

Usando las fórmulas de Excel podemos obtener estas estadísticas (sin imputación):

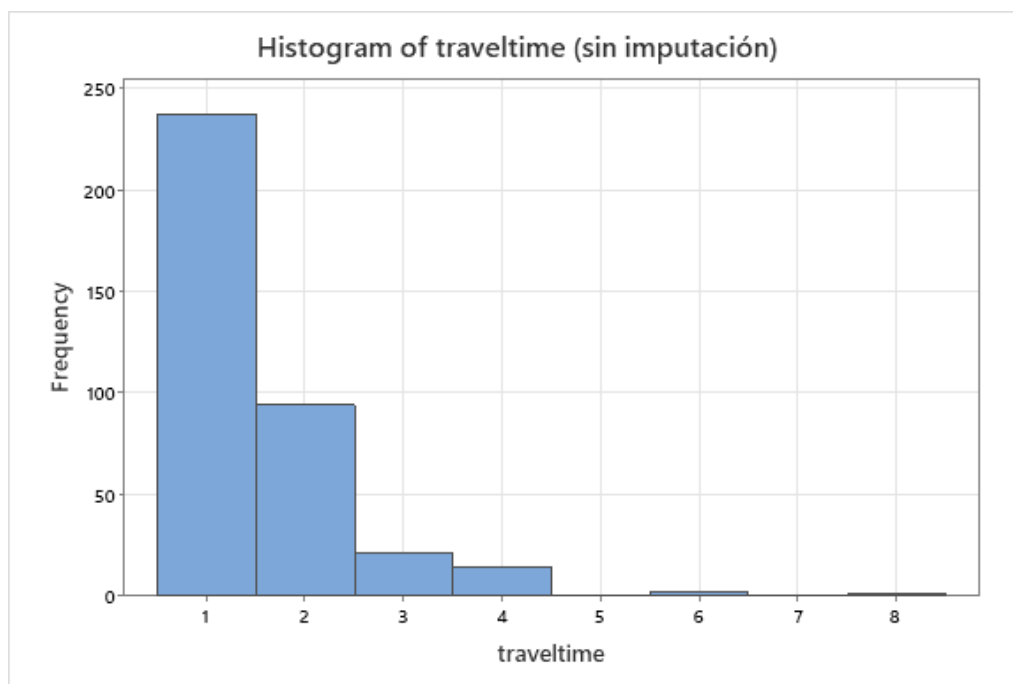
Traveltime:

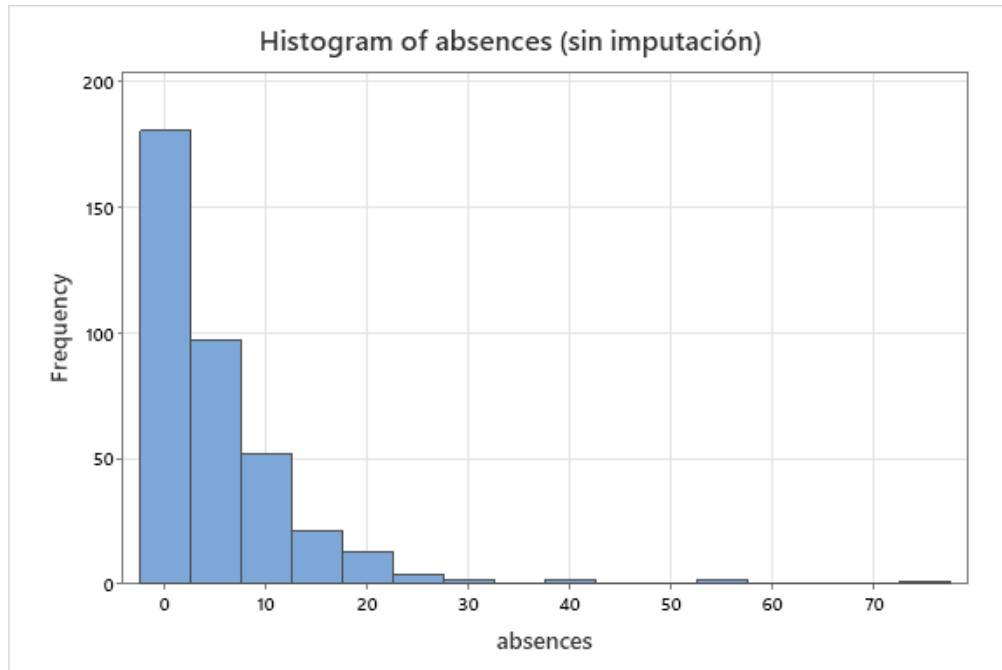
%datos faltantes	7.05
Mec datos faltantes	MCAR o NMAR
Media	1.528455285
desviación estándar	0.9028204892
mediana	1
moda	1

Absences:

%datos faltantes	5.61
Mec datos faltantes	MCAR o MAR
Media	5.542780749
desviación estándar	8.089117488
mediana	3.5
moda	0

Usando minitap obtenemos los siguientes histogramas:





#### 4. Utilizar el método de imputación adecuado para cada una de las variables con datos faltantes.

##### ◦ Imputación Simple: Media, Mediana, Moda

Implementamos imputación simple con moda, para poner los datos faltantes en el mismo lugar donde están la mayoría de los datos y evitar que se altere demasiado la distribución de los datos.

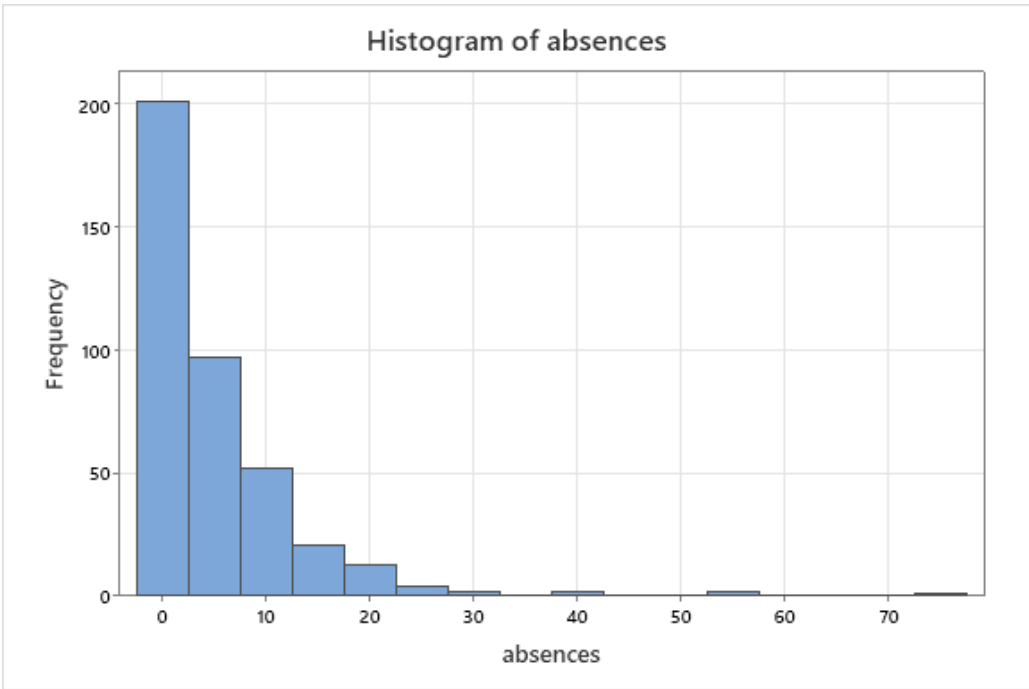
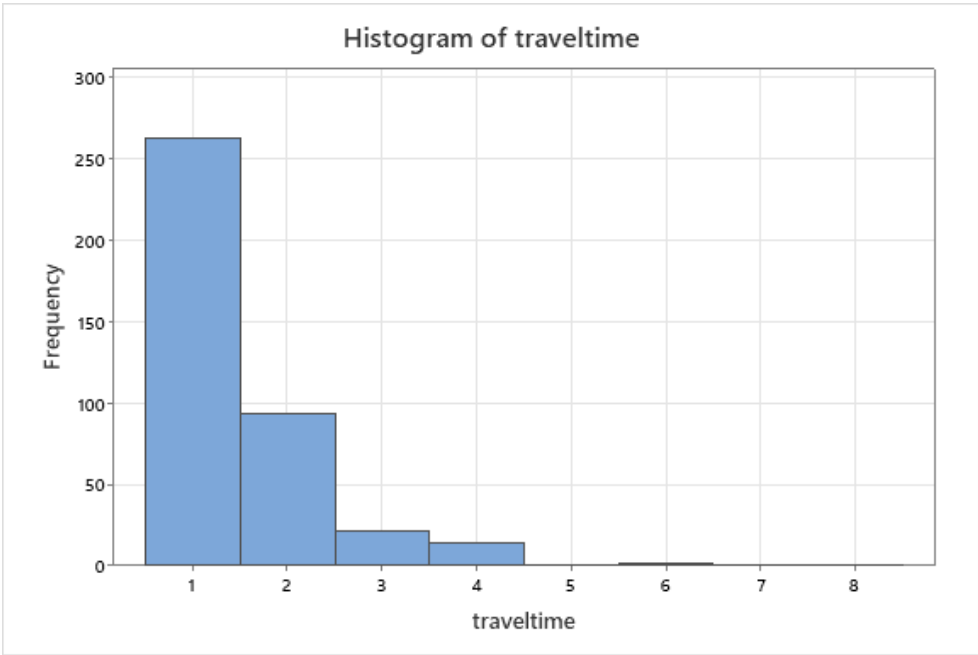
Reemplazamos en Excel los valores NA con el valor de la moda de cada respectiva variable. Obtenemos estas nuevas estadísticas.

Travelttime:

%datos faltantes	0
Mec datos faltantes	MCAR o NMAR
Media	1.493670886
desviación estándar	0.8823339165
mediana	1
moda	1

Absences:

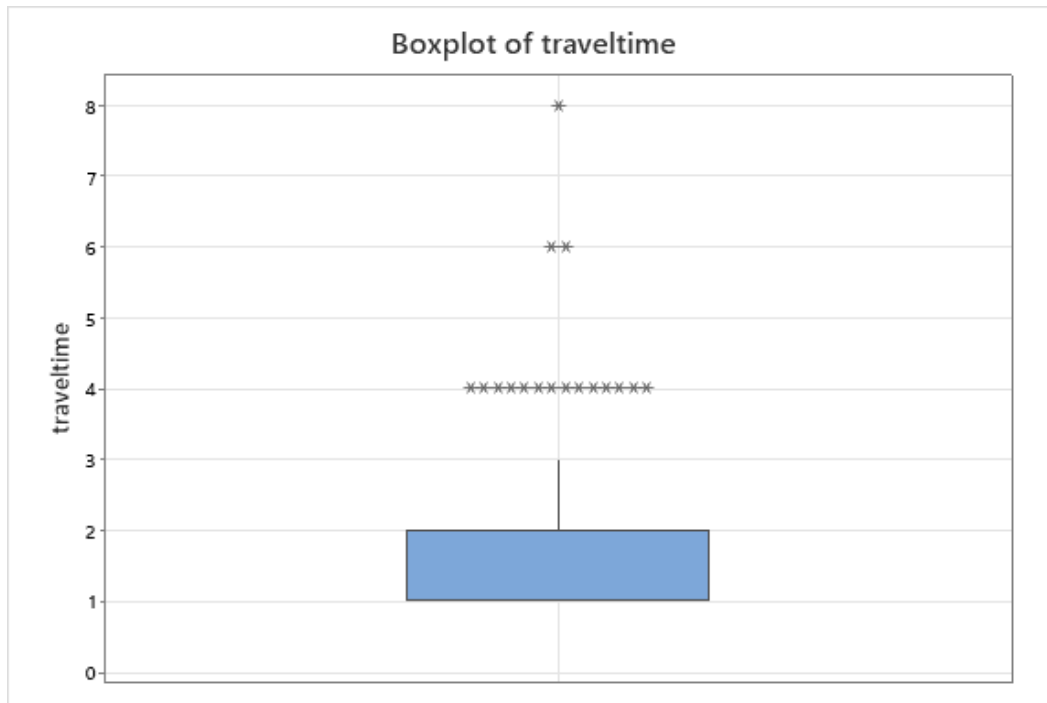
%datos faltantes	0
Mec datos faltantes	MCAR o MAR
Media	5.248101266
desviación estándar	7.968479476
mediana	2
moda	0



## 5. Realizar un boxplot e interpretarlo.

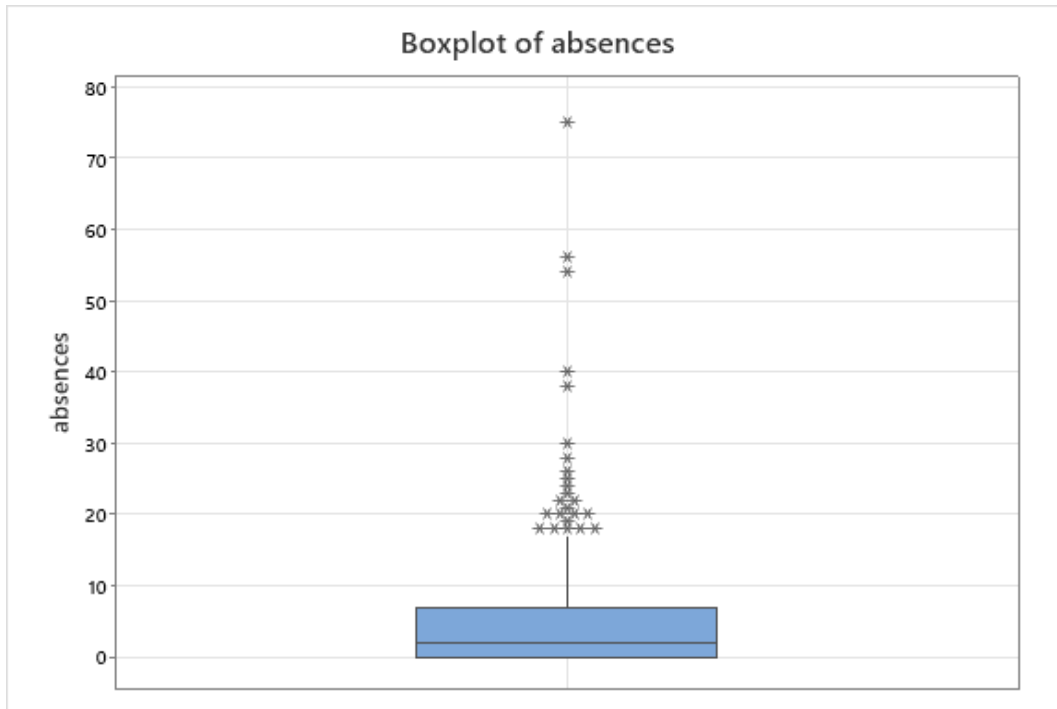
Usamos boxplot de los datos tras la imputación simple usando la moda.

Usamos los boxplots para identificar valores atípicos.



Observamos que el 50% de los datos están entre 1 y 2, dentro de la caja.

Hay una gran cantidad de valores atípicos que superan por mucho el promedio de los valores, la mayoría con un valor de 4.



Observamos que el 50% de los datos están entre 0 y 8-9, dentro de la caja. Hay una gran cantidad de valores atípicos que superan por mucho el promedio de los valores, llegando hasta a mas de 50. Observamos que la mediana se encuentra en la parte inferior de la caja.