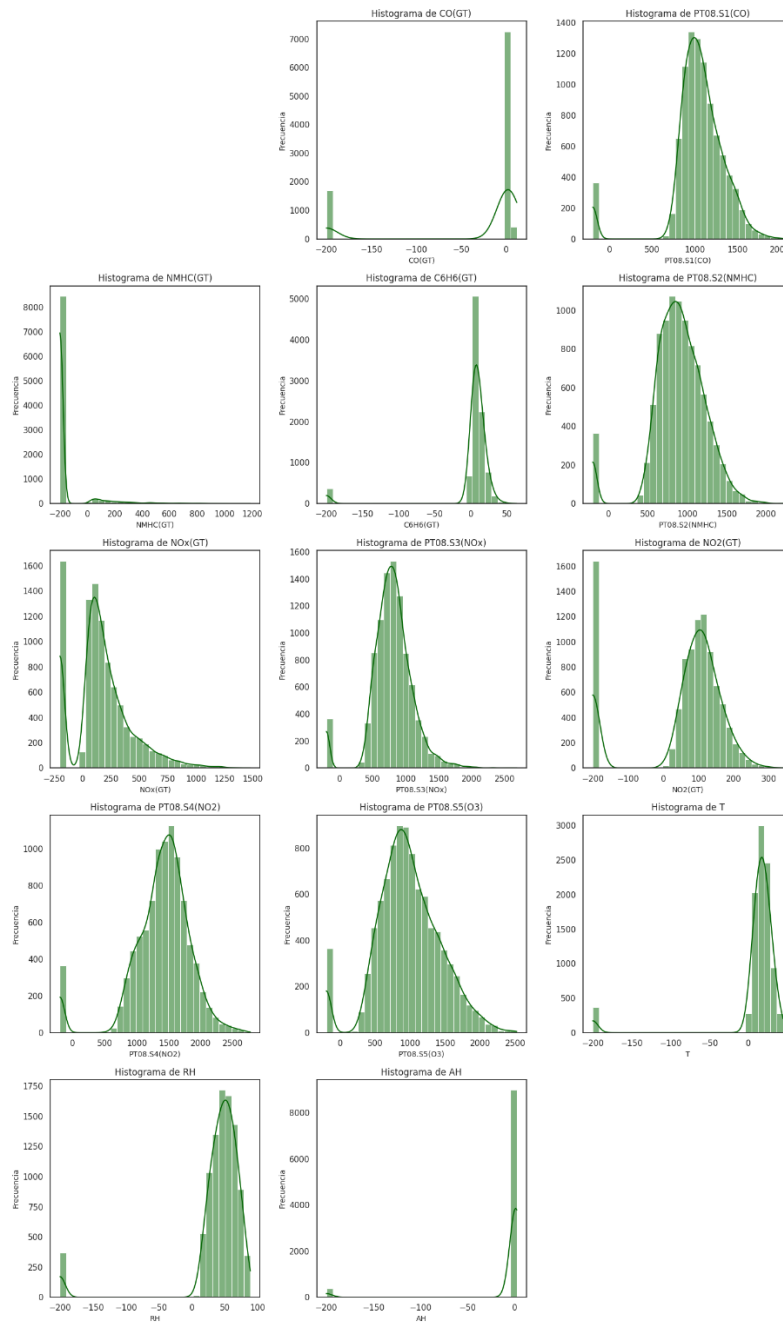


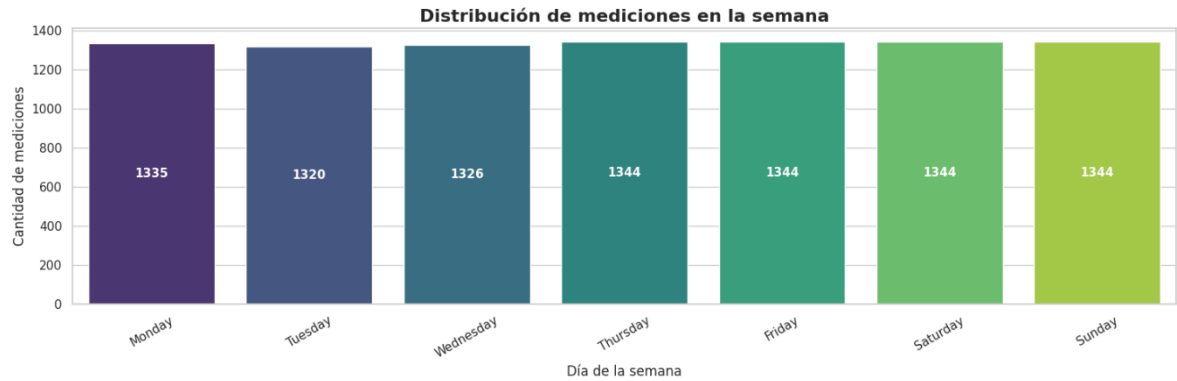
Análisis Exploratorio de Datos (EDA): Base de datos “Air Quality”

- ¿Qué patrones o tendencias observaste en los histogramas y gráficas de densidad (PDF)? ¿Alguna variable parece seguir una distribución normal?



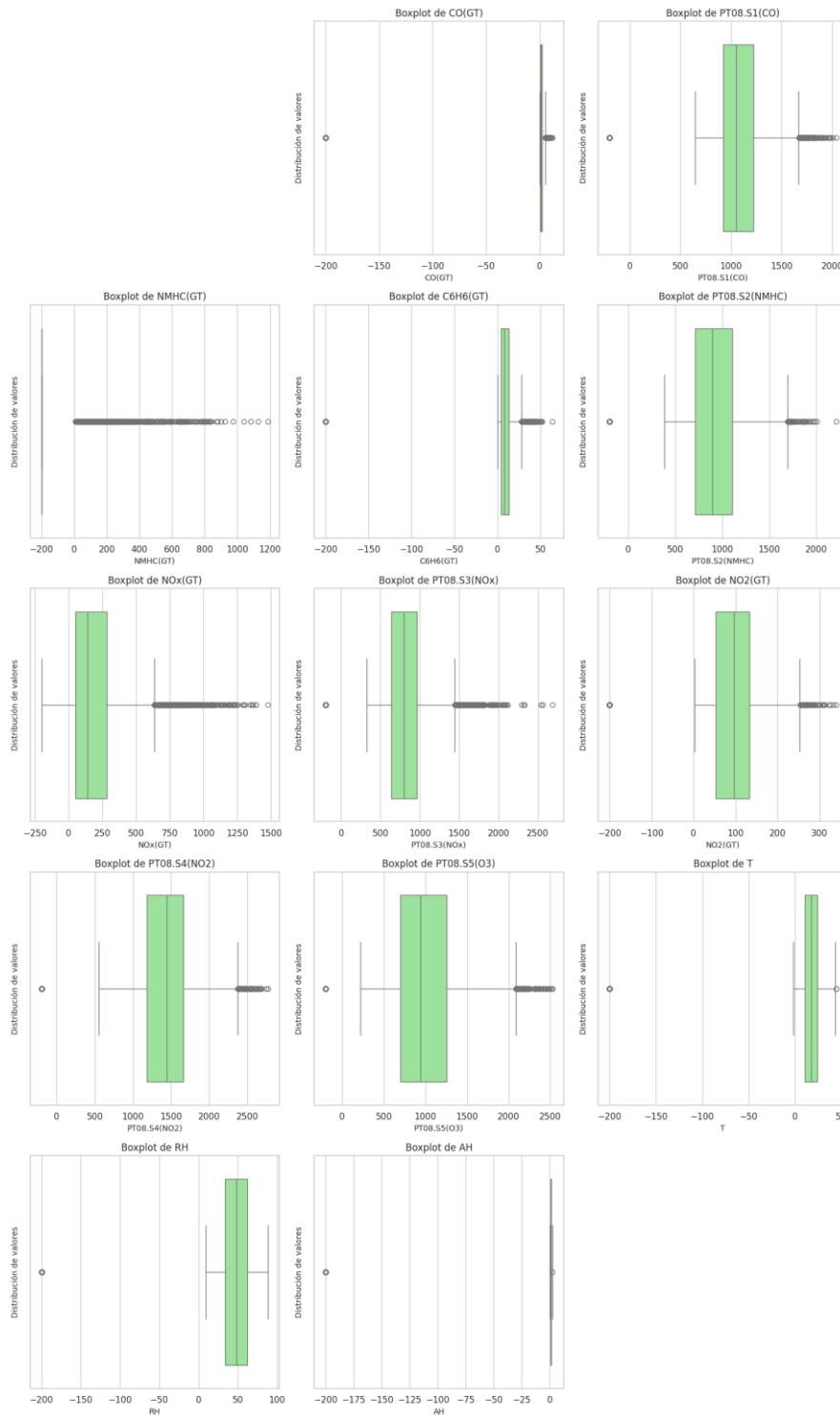
Todas las variables presentan distribuciones asimétricas, con un fuerte sesgo. Ninguna parece seguir una distribución normal definitivamente.

- *¿Qué información útil obtuviste de la gráfica de barras para el día de la semana?*



Tomando en cuenta que la información se compone por datos sobre la calidad del aire, se puede apreciar que los días que menos se contamina son el lunes, martes y miércoles, de ahí en adelante la contaminación que registró fue la misma. Probablemente la diferencia se deba a los días laborales.

- *¿Identificaste outliers en los boxplots? ¿Cómo podrían afectar estos outliers al modelo de regresión lineal?*

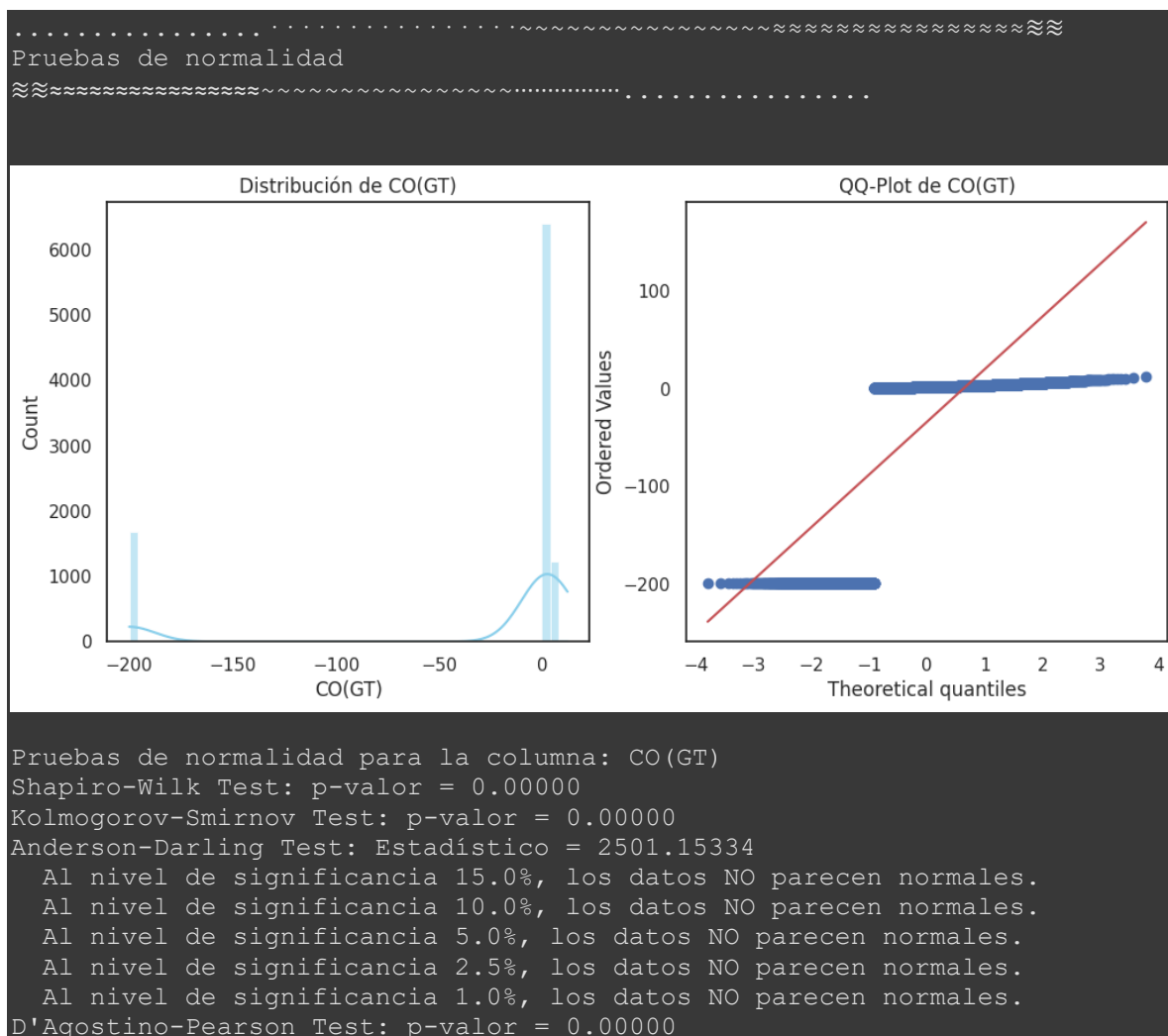


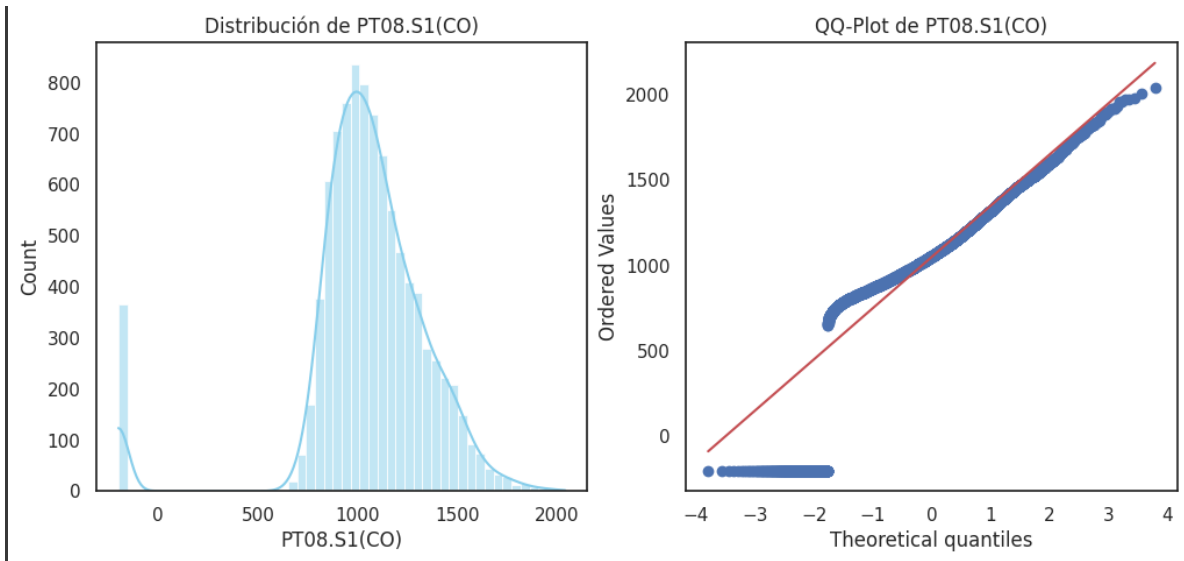
Se identificaron varios outliers, en casi todas las columnas. En algunos casos, hay valores que llaman la atención por encontrarse en un solo punto en gran cantidad, lo que muy probablemente indica un error, ya sea en la medición, en el registro u otro. Valores tan extremos afectan la media y varianza de las variables, lo que afectaría de forma

negativa a un modelo de regresión lineal. Deberían ser tratados antes de someterlos a dicho proceso.

Pruebas de Normalidad:

- ¿Qué conclusiones obtuviste de las pruebas de normalidad (Shapiro-Wilk, Anderson-Darling, Kolmogorov-Smirnov)? ¿Qué variables no siguen una distribución normal?





Pruebas de normalidad para la columna: PT08.S1(CO)

Shapiro-Wilk Test: p-valor = 0.00000

Kolmogorov-Smirnov Test: p-valor = 0.00000

Anderson-Darling Test: Estadístico = 355.41246

Al nivel de significancia 15.0%, los datos NO parecen normales.

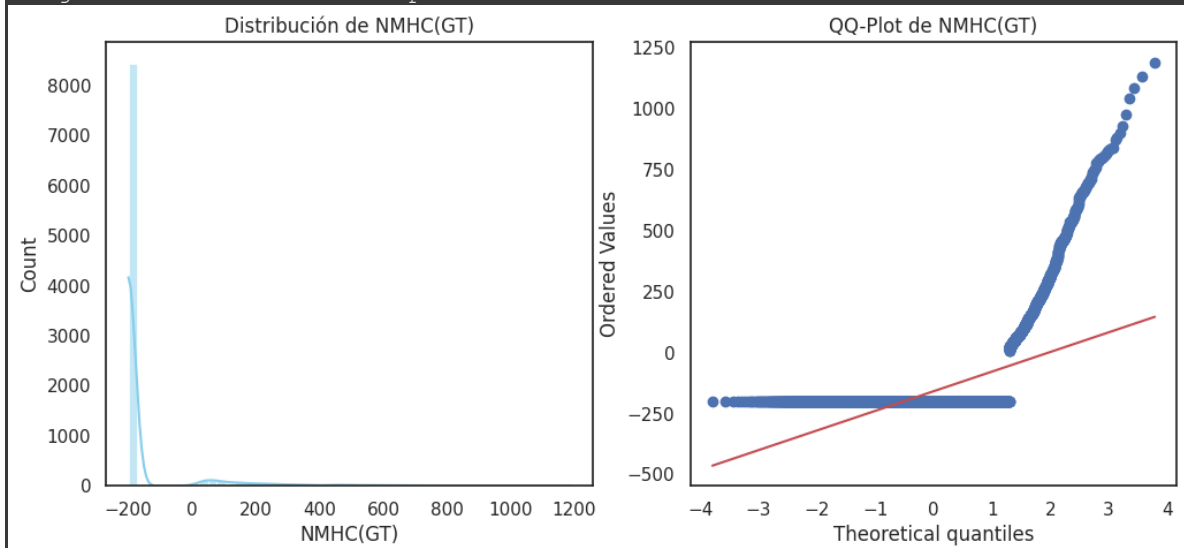
Al nivel de significancia 10.0%, los datos NO parecen normales.

Al nivel de significancia 5.0%, los datos NO parecen normales.

Al nivel de significancia 2.5%, los datos NO parecen normales.

Al nivel de significancia 1.0%, los datos NO parecen normales.

D'Agostino-Pearson Test: p-valor = 0.00000



Pruebas de normalidad para la columna: NMHC(GT)

Shapiro-Wilk Test: p-valor = 0.00000

Kolmogorov-Smirnov Test: p-valor = 0.00000

Anderson-Darling Test: Estadístico = 2820.94917

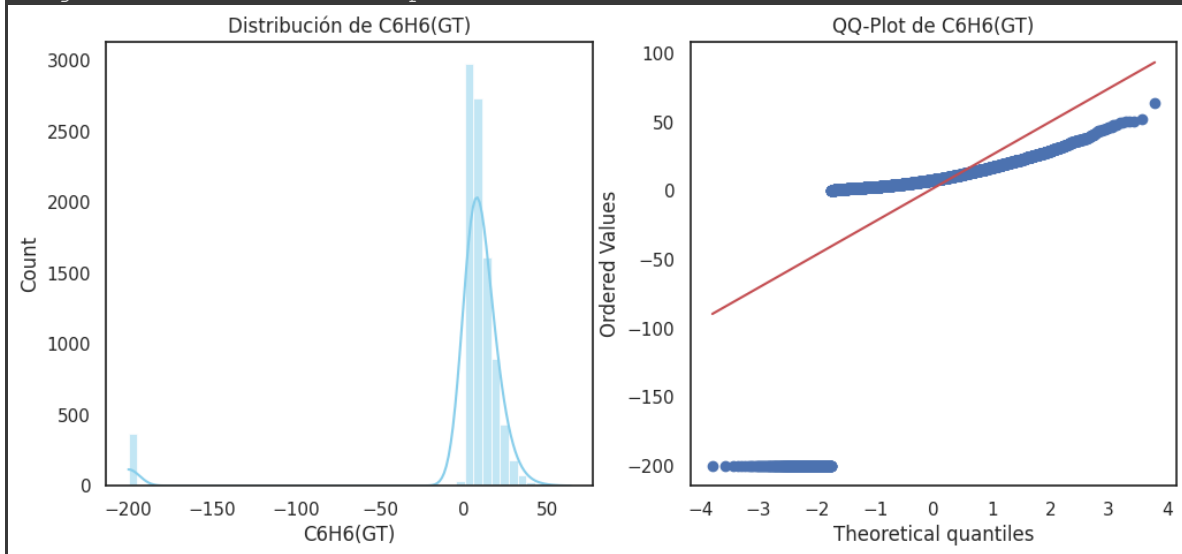
Al nivel de significancia 15.0%, los datos NO parecen normales.

Al nivel de significancia 10.0%, los datos NO parecen normales.

Al nivel de significancia 5.0%, los datos NO parecen normales.

Al nivel de significancia 2.5%, los datos NO parecen normales.

Al nivel de significancia 1.0%, los datos NO parecen normales.
D'Agostino-Pearson Test: p-valor = 0.00000



Pruebas de normalidad para la columna: C6H6(GT)

Shapiro-Wilk Test: p-valor = 0.00000

Kolmogorov-Smirnov Test: p-valor = 0.00000

Anderson-Darling Test: Estadístico = 2272.58484

Al nivel de significancia 15.0%, los datos NO parecen normales.

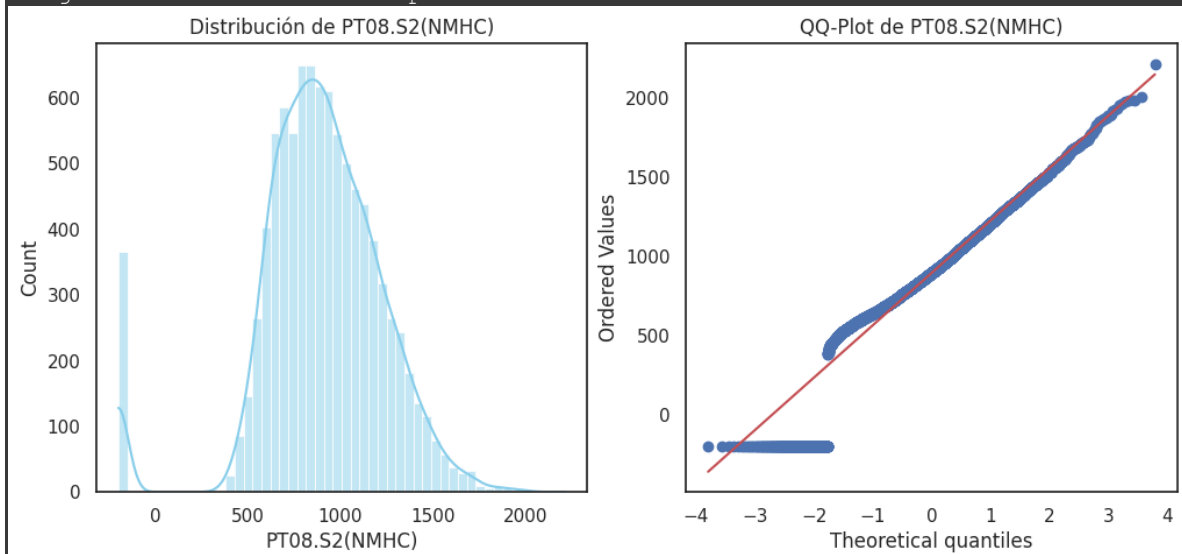
Al nivel de significancia 10.0%, los datos NO parecen normales.

Al nivel de significancia 5.0%, los datos NO parecen normales.

Al nivel de significancia 2.5%, los datos NO parecen normales.

Al nivel de significancia 1.0%, los datos NO parecen normales.

D'Agostino-Pearson Test: p-valor = 0.00000



Pruebas de normalidad para la columna: PT08.S2(NMHC)

Shapiro-Wilk Test: p-valor = 0.00000

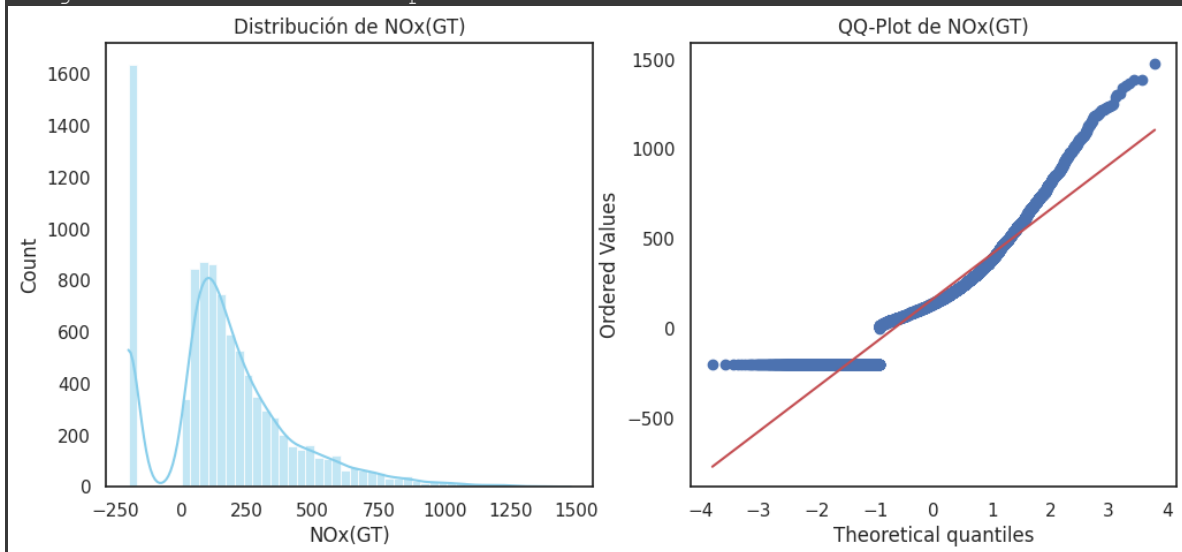
Kolmogorov-Smirnov Test: p-valor = 0.00000

Anderson-Darling Test: Estadístico = 119.20819

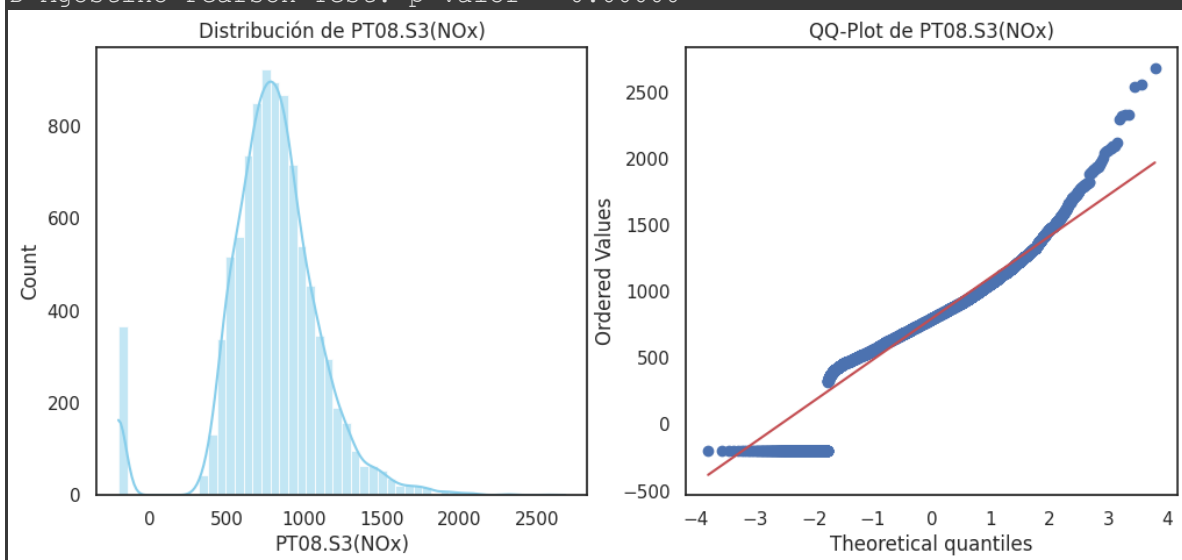
Al nivel de significancia 15.0%, los datos NO parecen normales.

Al nivel de significancia 10.0%, los datos NO parecen normales.

Al nivel de significancia 5.0%, los datos NO parecen normales.
 Al nivel de significancia 2.5%, los datos NO parecen normales.
 Al nivel de significancia 1.0%, los datos NO parecen normales.
 D'Agostino-Pearson Test: p-valor = 0.00000

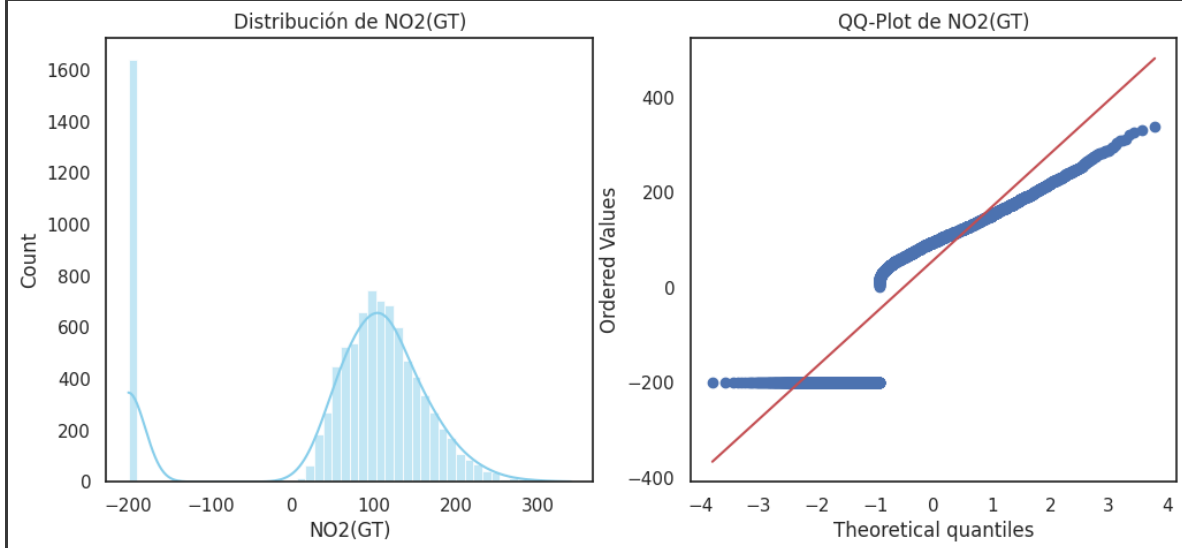


Pruebas de normalidad para la columna: NOx(GT)
 Shapiro-Wilk Test: p-valor = 0.00000
 Kolmogorov-Smirnov Test: p-valor = 0.00000
 Anderson-Darling Test: Estadístico = 188.79429
 Al nivel de significancia 15.0%, los datos NO parecen normales.
 Al nivel de significancia 10.0%, los datos NO parecen normales.
 Al nivel de significancia 5.0%, los datos NO parecen normales.
 Al nivel de significancia 2.5%, los datos NO parecen normales.
 Al nivel de significancia 1.0%, los datos NO parecen normales.
 D'Agostino-Pearson Test: p-valor = 0.00000

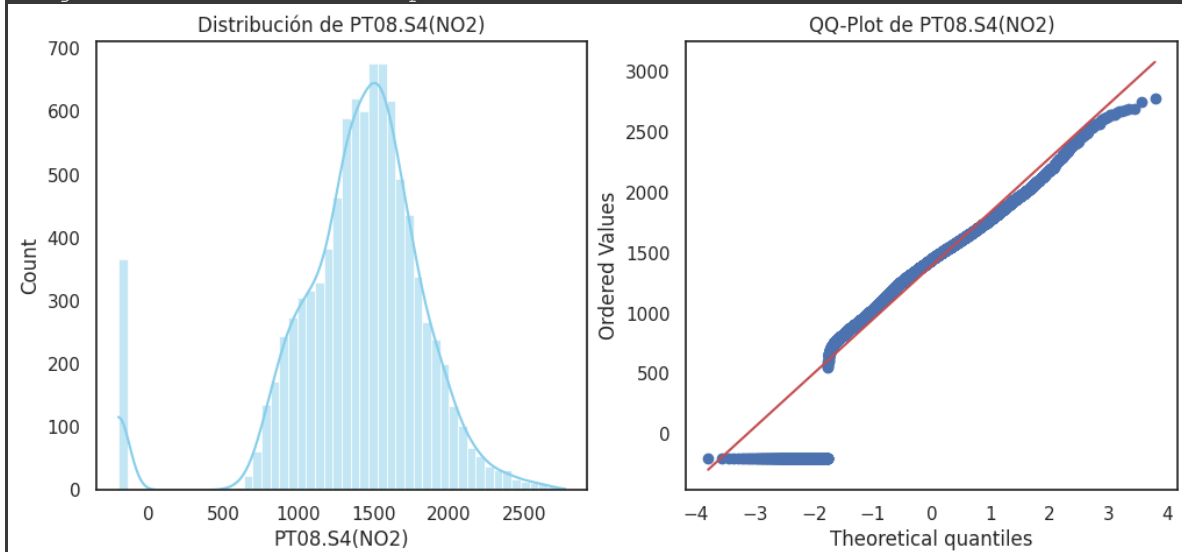


Pruebas de normalidad para la columna: PT08.S3(NOx)
 Shapiro-Wilk Test: p-valor = 0.00000
 Kolmogorov-Smirnov Test: p-valor = 0.00000
 Anderson-Darling Test: Estadístico = 159.03112

Al nivel de significancia 15.0%, los datos NO parecen normales.
 Al nivel de significancia 10.0%, los datos NO parecen normales.
 Al nivel de significancia 5.0%, los datos NO parecen normales.
 Al nivel de significancia 2.5%, los datos NO parecen normales.
 Al nivel de significancia 1.0%, los datos NO parecen normales.
 D'Agostino-Pearson Test: p-valor = 0.00000



Pruebas de normalidad para la columna: NO2 (GT)
 Shapiro-Wilk Test: p-valor = 0.00000
 Kolmogorov-Smirnov Test: p-valor = 0.00000
 Anderson-Darling Test: Estadístico = 894.66163
 Al nivel de significancia 15.0%, los datos NO parecen normales.
 Al nivel de significancia 10.0%, los datos NO parecen normales.
 Al nivel de significancia 5.0%, los datos NO parecen normales.
 Al nivel de significancia 2.5%, los datos NO parecen normales.
 Al nivel de significancia 1.0%, los datos NO parecen normales.
 D'Agostino-Pearson Test: p-valor = 0.00000



Pruebas de normalidad para la columna: PT08.S4 (NO2)
 Shapiro-Wilk Test: p-valor = 0.00000

Kolmogorov-Smirnov Test: p-valor = 0.00000

Anderson-Darling Test: Estadístico = 183.34158

Al nivel de significancia 15.0%, los datos NO parecen normales.

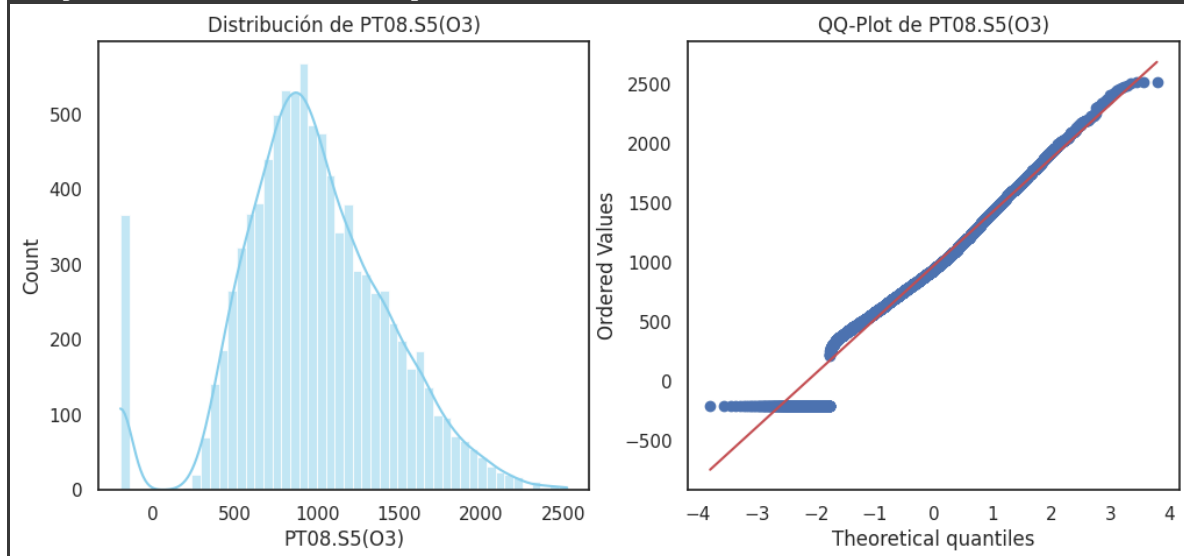
Al nivel de significancia 10.0%, los datos NO parecen normales.

Al nivel de significancia 5.0%, los datos NO parecen normales.

Al nivel de significancia 2.5%, los datos NO parecen normales.

Al nivel de significancia 1.0%, los datos NO parecen normales.

D'Agostino-Pearson Test: p-valor = 0.00000



Pruebas de normalidad para la columna: PT08.S5(O3)

Shapiro-Wilk Test: p-valor = 0.00000

Kolmogorov-Smirnov Test: p-valor = 0.00000

Anderson-Darling Test: Estadístico = 45.68234

Al nivel de significancia 15.0%, los datos NO parecen normales.

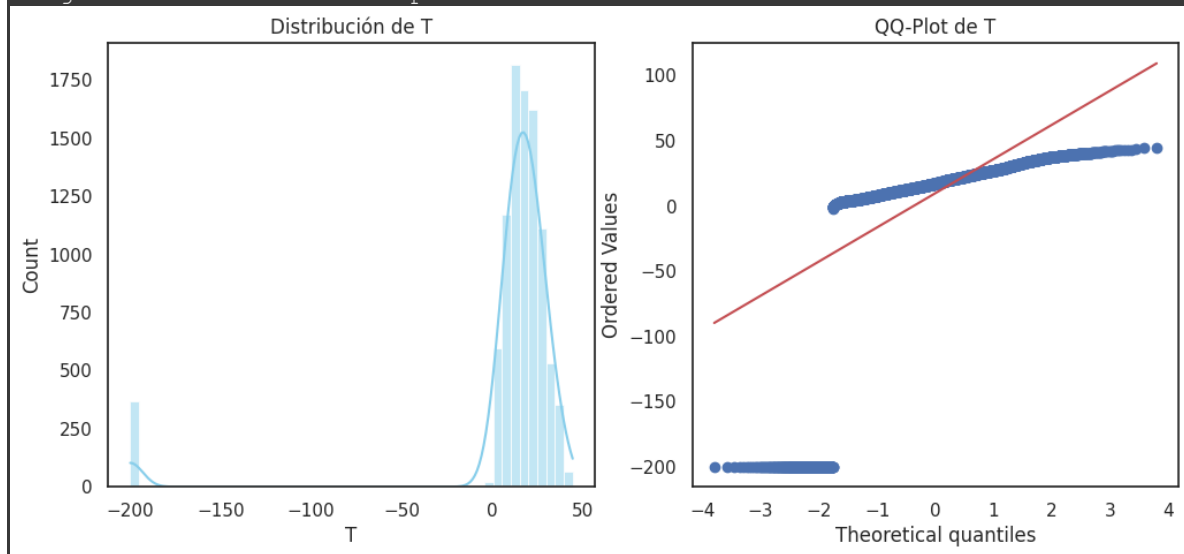
Al nivel de significancia 10.0%, los datos NO parecen normales.

Al nivel de significancia 5.0%, los datos NO parecen normales.

Al nivel de significancia 2.5%, los datos NO parecen normales.

Al nivel de significancia 1.0%, los datos NO parecen normales.

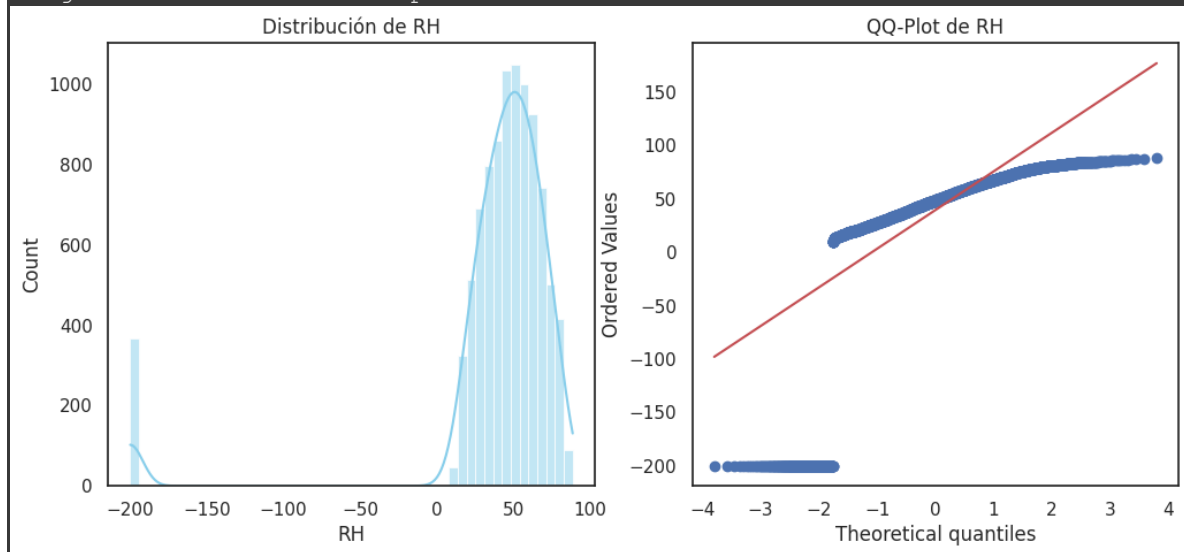
D'Agostino-Pearson Test: p-valor = 0.00000



```

Pruebas de normalidad para la columna: T
Shapiro-Wilk Test: p-valor = 0.00000
Kolmogorov-Smirnov Test: p-valor = 0.00000
Anderson-Darling Test: Estadístico = 2041.45858
  Al nivel de significancia 15.0%, los datos NO parecen normales.
  Al nivel de significancia 10.0%, los datos NO parecen normales.
  Al nivel de significancia 5.0%, los datos NO parecen normales.
  Al nivel de significancia 2.5%, los datos NO parecen normales.
  Al nivel de significancia 1.0%, los datos NO parecen normales.
D'Agostino-Pearson Test: p-valor = 0.00000

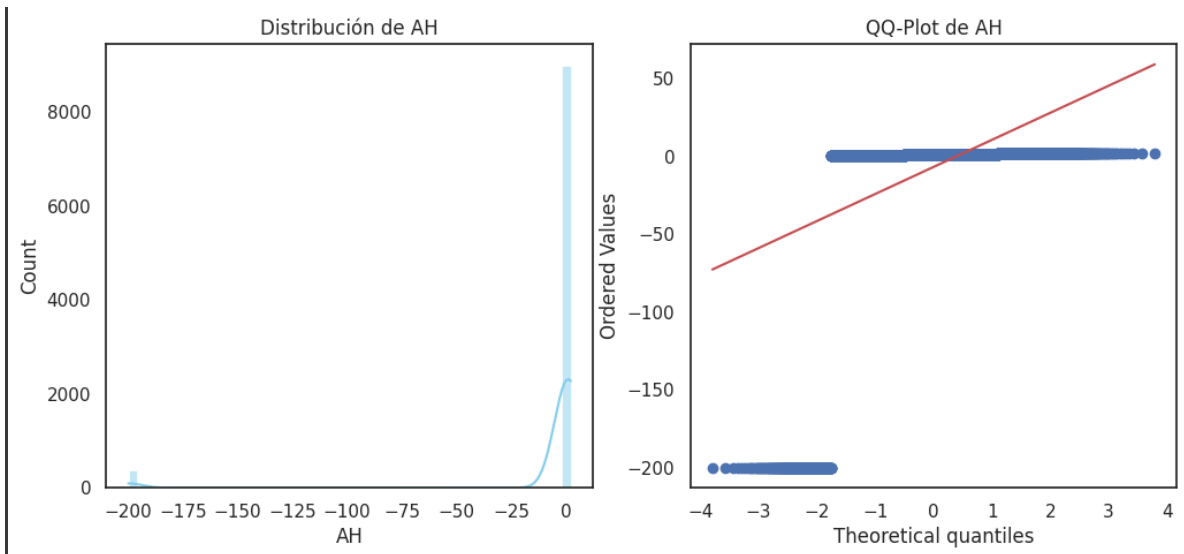
```



```

Pruebas de normalidad para la columna: RH
Shapiro-Wilk Test: p-valor = 0.00000
Kolmogorov-Smirnov Test: p-valor = 0.00000
Anderson-Darling Test: Estadístico = 1329.04460
  Al nivel de significancia 15.0%, los datos NO parecen normales.
  Al nivel de significancia 10.0%, los datos NO parecen normales.
  Al nivel de significancia 5.0%, los datos NO parecen normales.
  Al nivel de significancia 2.5%, los datos NO parecen normales.
  Al nivel de significancia 1.0%, los datos NO parecen normales.
D'Agostino-Pearson Test: p-valor = 0.00000

```



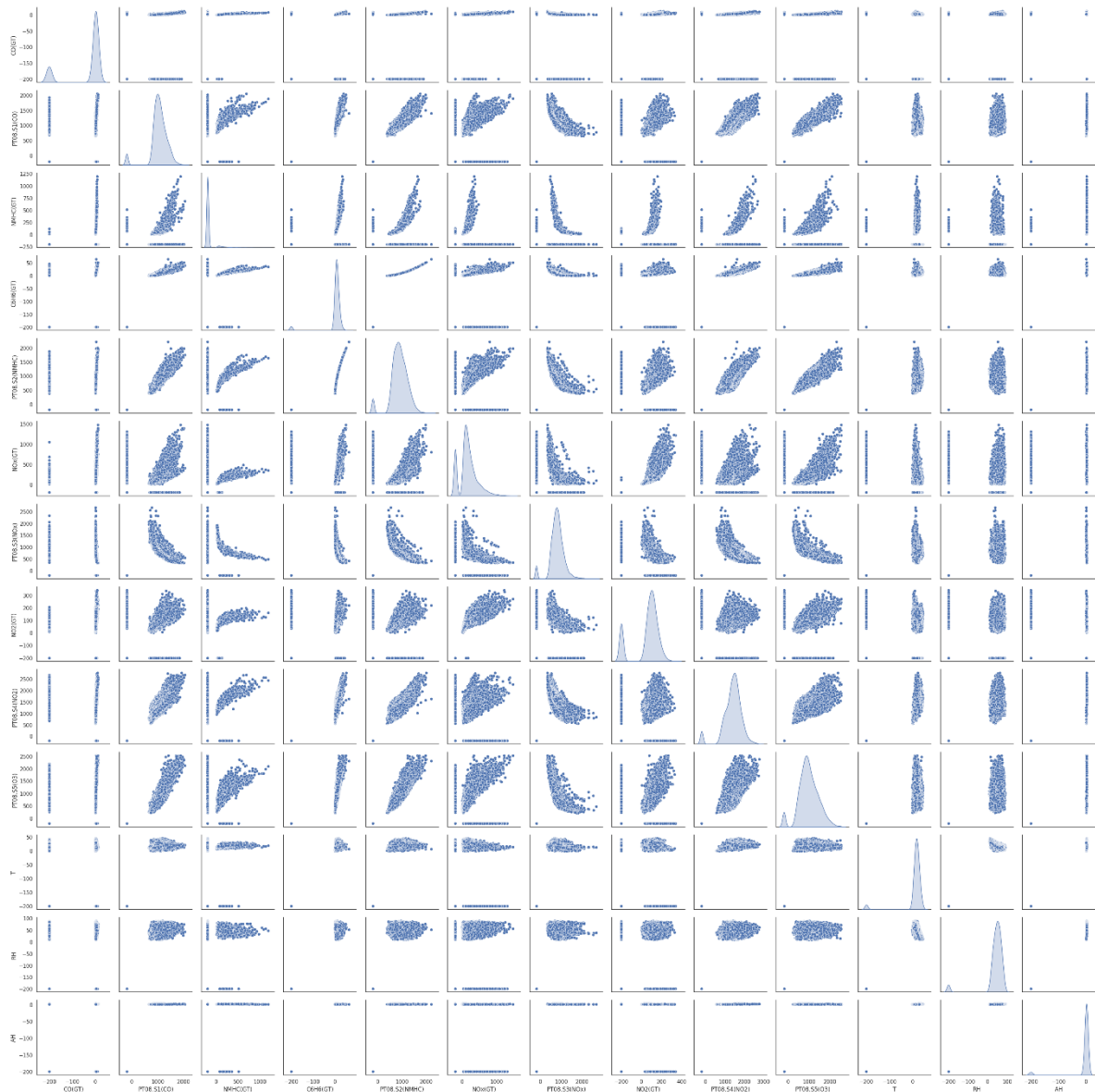
```

Pruebas de normalidad para la columna: AH
Shapiro-Wilk Test: p-valor = 0.00000
Kolmogorov-Smirnov Test: p-valor = 0.00000
Anderson-Darling Test: Estadístico = 3357.68062
  Al nivel de significancia 15.0%, los datos NO parecen normales.
  Al nivel de significancia 10.0%, los datos NO parecen normales.
  Al nivel de significancia 5.0%, los datos NO parecen normales.
  Al nivel de significancia 2.5%, los datos NO parecen normales.
  Al nivel de significancia 1.0%, los datos NO parecen normales.
D'Agostino-Pearson Test: p-valor = 0.00000

```

Ninguna columna tiene una distribución normal, tienen muchos valores atípicos y extremos, lo que hace que las pruebas no puedan calcular estadísticas. Eliminé varias columnas vacías, pero requiere imputación.

- *¿Cómo interpretas los QQplots? ¿Qué variables se desvían significativamente de la normalidad?*



Todas las variables se alejan de la normalidad, todas presentan desviaciones significativas con la línea diagonal. Mas que nada por los outliers que se disparan en algunas ocasiones.

Tratamiento de Datos Faltantes:

- ¿Qué estrategia utilizaste para manejar los datos faltantes? ¿Por qué elegiste esa estrategia?

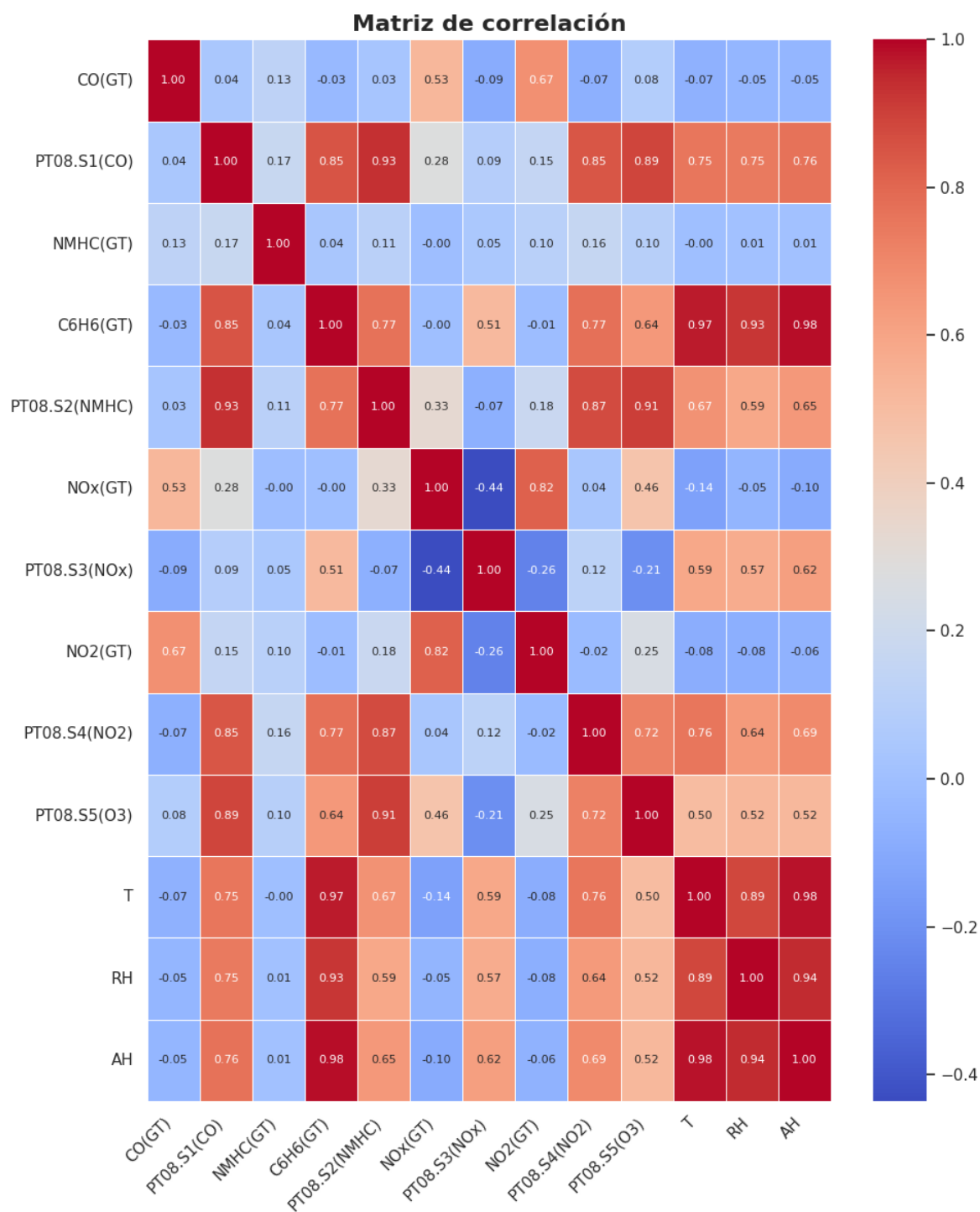
Se eliminaron 3 columnas del dataset original; “Unnamed: 15”, “Unnamed: 16” y “Time”, por no contener datos y no aportar nada.

- ¿Cómo cambió el EDA después de la imputación de datos? ¿Observaste diferencias significativas en las distribuciones de las variables?

Se aplicó una transformación logarítmica para reducir el sesgo y no hubo diferencia en la normalidad. Para ver cambios en las distribuciones se deben imputar los datos de otra forma.

Matriz de Correlación y Pairplot:

- *¿Qué relaciones lineales identificaste en la matriz de correlación y el pairplot? ¿Alguna variable tiene una correlación fuerte con la variable objetivo?*



Las columnas T, RH, AH, están fuertemente correlacionadas entre si, y con algunos contaminantes como C6H6(GT). La serie de sensores debería de relacionarse siempre con el valor que miden, ahí pueden estar afectando malas mediciones. Para encontrar las relaciones solo con la variable objetivo primero se debe definir cuál es, es decir que queremos saber.

- ¿Cómo podrías utilizar esta información para seleccionar características (features) en un modelo de regresión lineal?

Con esta información, se puede encontrar las variables que tienen una correlación más alta, y a su vez descartar las que no. Para que el modelo sea alimentado solo con información de calidad.