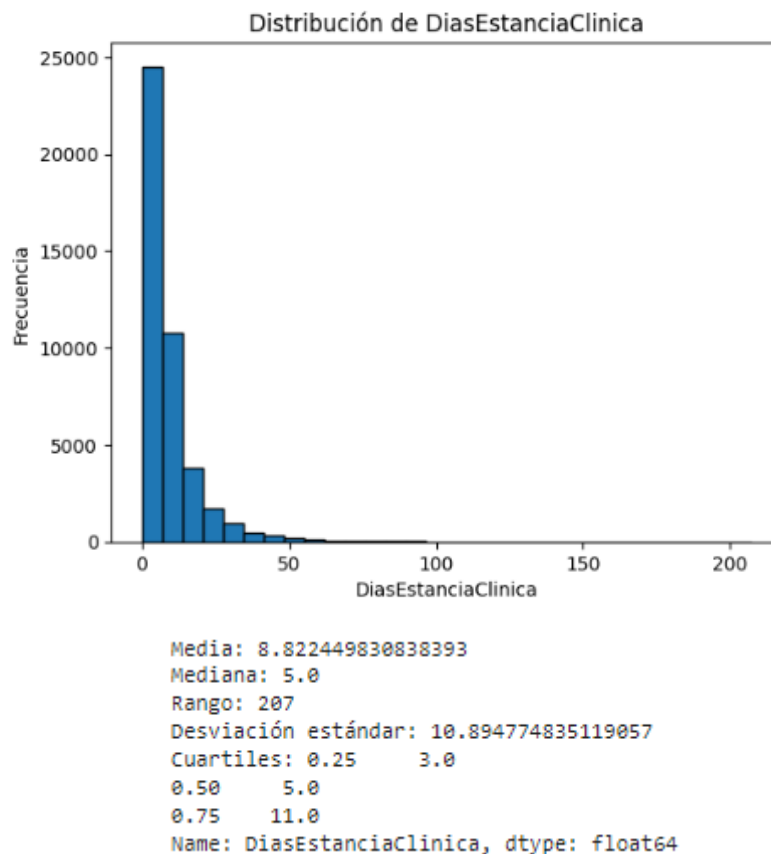


Análisis exploratorio de datos

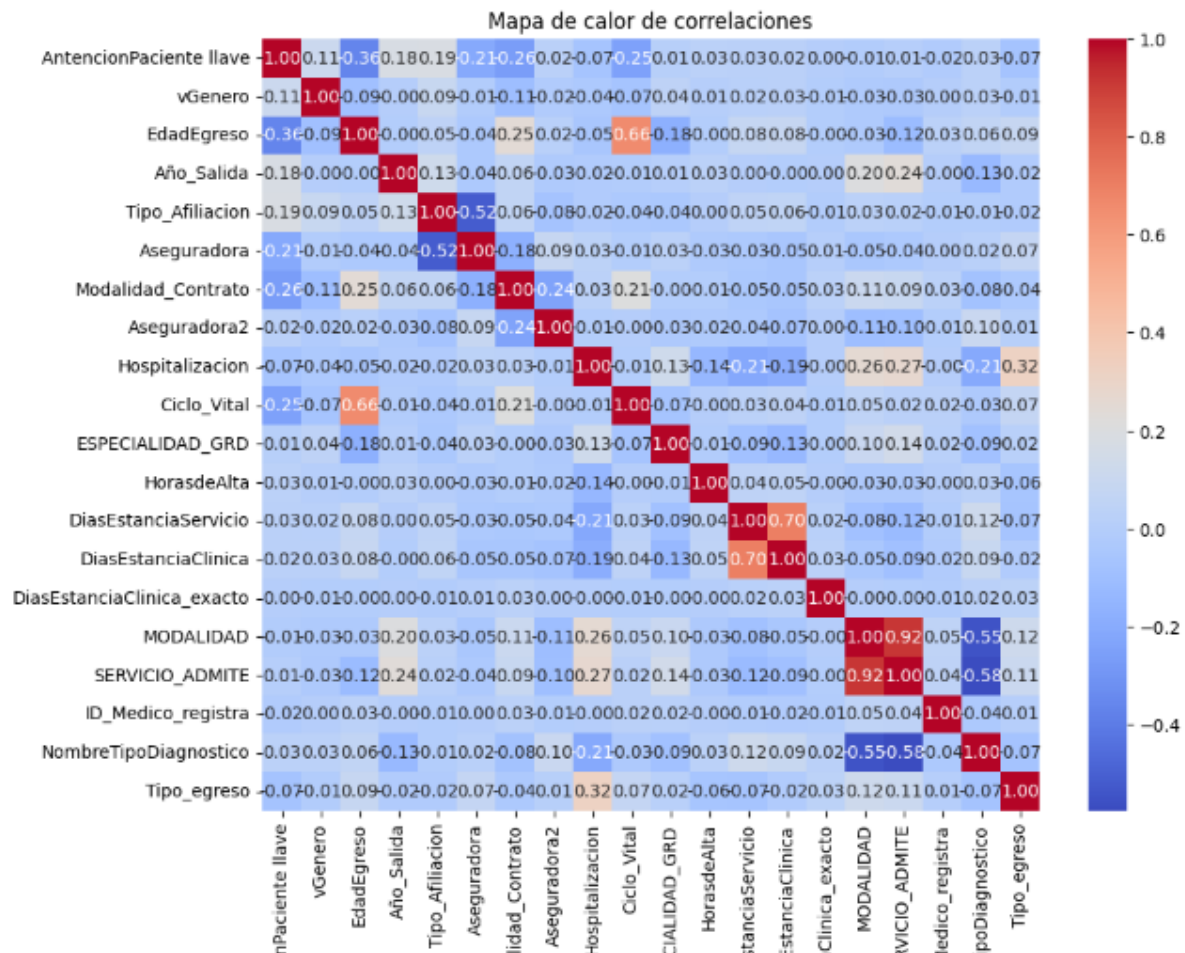
Inicialmente se comenzó por hacer un análisis con respecto a la cantidad y tipo de datos que se tiene en el dataset, dicho procedimiento se hizo mediante la función .info a partir de esta se encontró que se cuenta con un dataset de 73 columnas y 43154 filas. En el dataset se cuentan variables numéricas y categóricas que describen datos de pacientes ingresados a determinada institución hospitalaria.

A partir del data set presentado se busca tratar de predecir el tiempo de estancia de un paciente en la unidad hospitalaria por lo que dicha variable será nuestro propósito de predicción. Por lo que se procede a realizar un diagrama de distribución además de calcular la mediana, la media, el rango, la desviación estándar y los cuartiles para así tener una mejor comprensión del comportamiento de la variable. Se encuentran los siguientes resultados:



A partir de la información podemos determinar que la variable días de estancia de servicio cuenta con una distribución de asimetría positiva en la cual la mayor cantidad de valores se concentran en los valores más bajos y a medida que aumenta el valor de la variable disminuye su frecuencia. Además sabemos que el tiempo de estancia media es de 8.82 días aproximadamente y la mediana es 5.0 aproximadamente, dichos valores nos dan un indicativo de la magnitud que debemos aspirar encontrar una vez se apliquen modelos para predecir dicha variable.

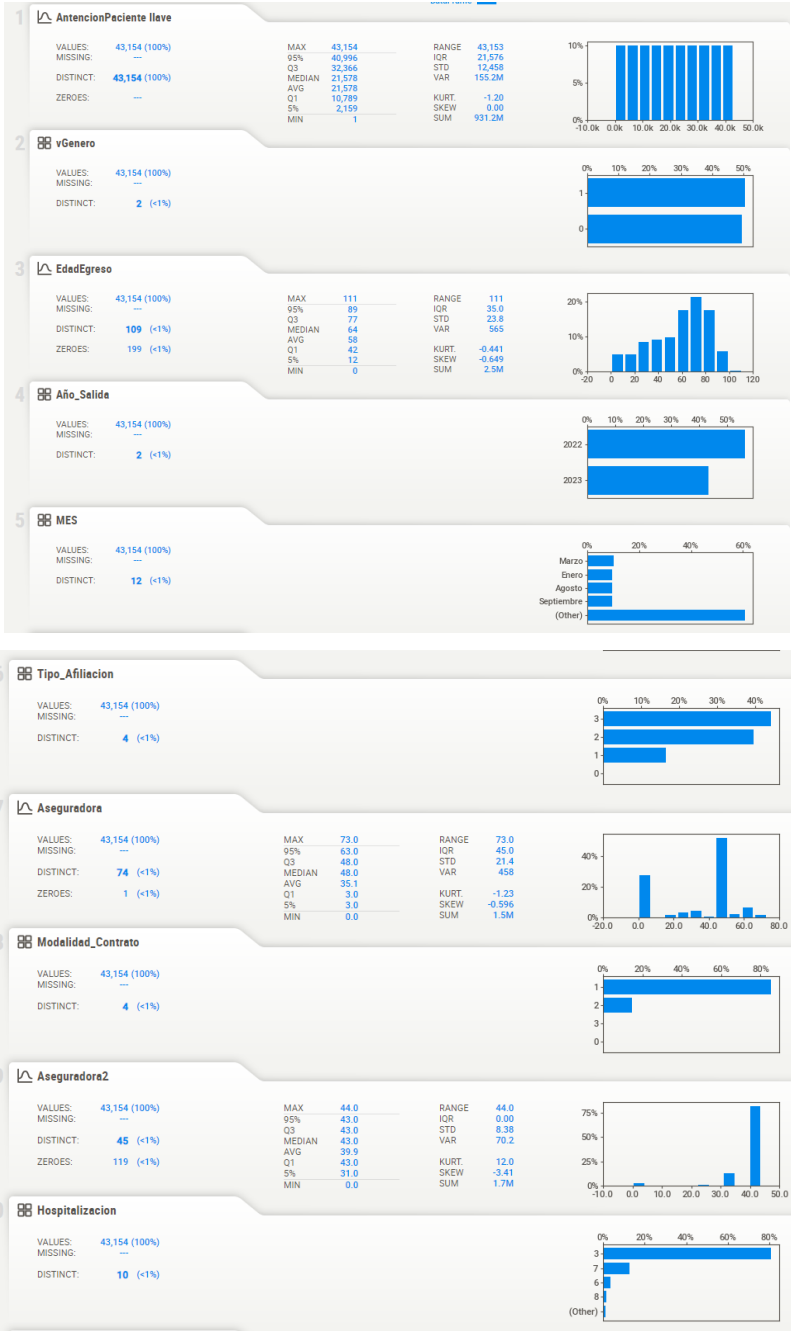
Una vez analizada nuestra variable de salida se procede a hacer un mapa de calor con las correlaciones entre cada variable numérica entregada en el data set, obteniéndose el siguiente resultado:

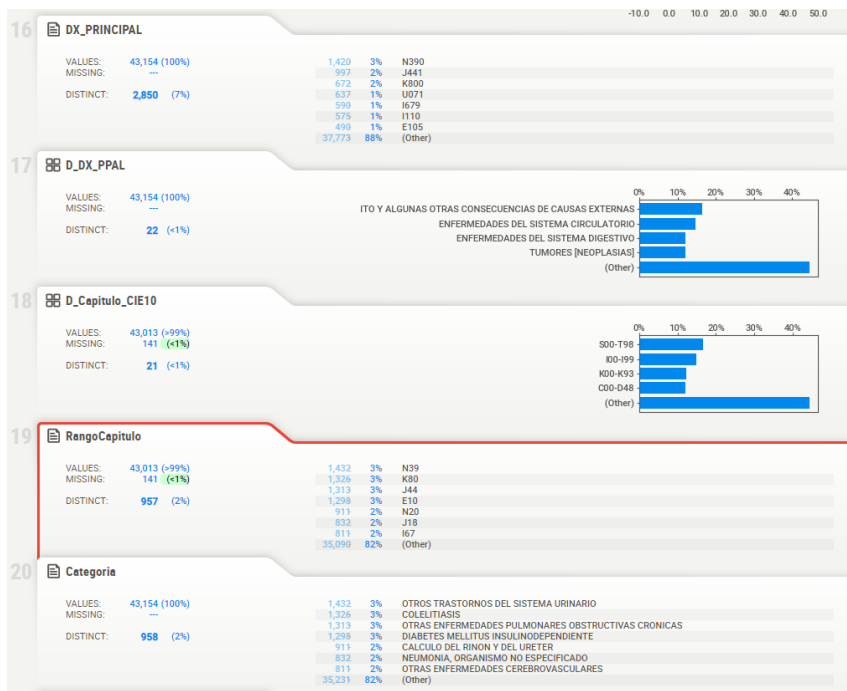
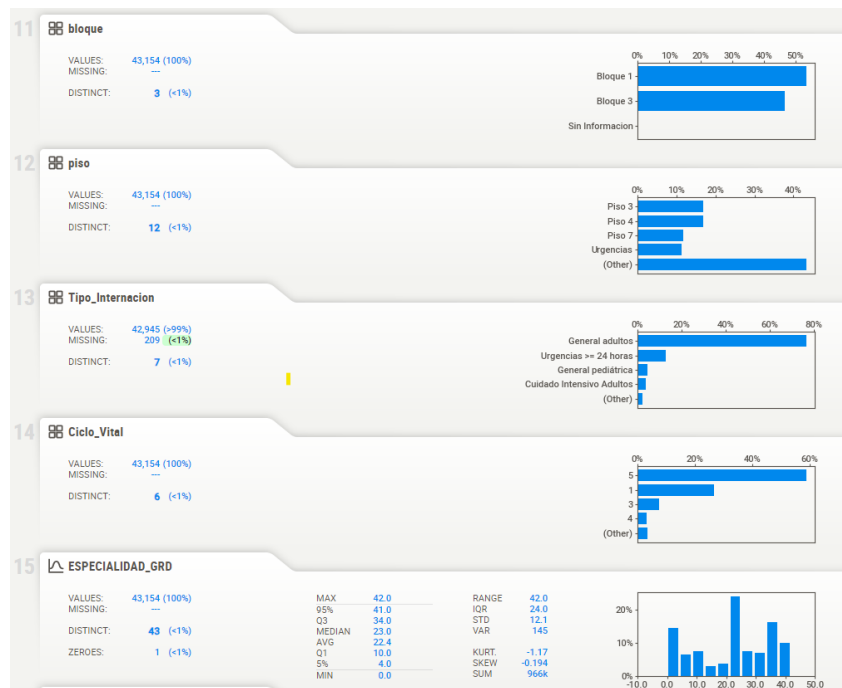


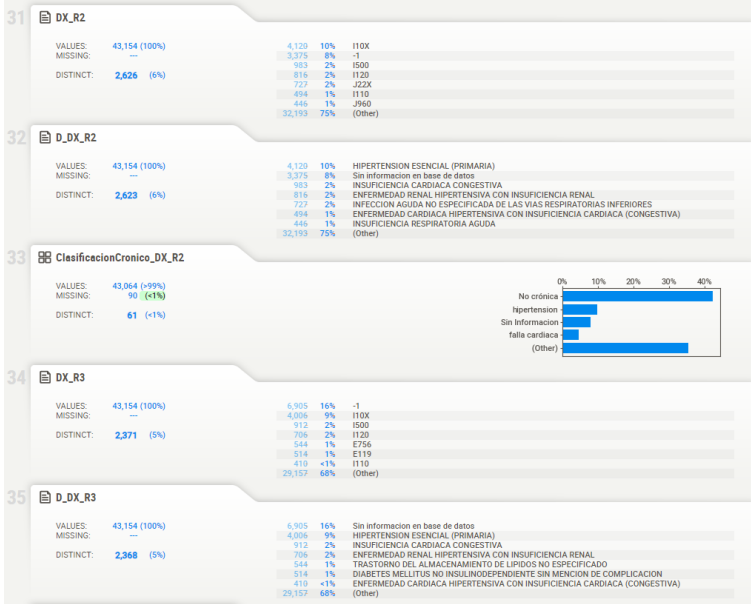
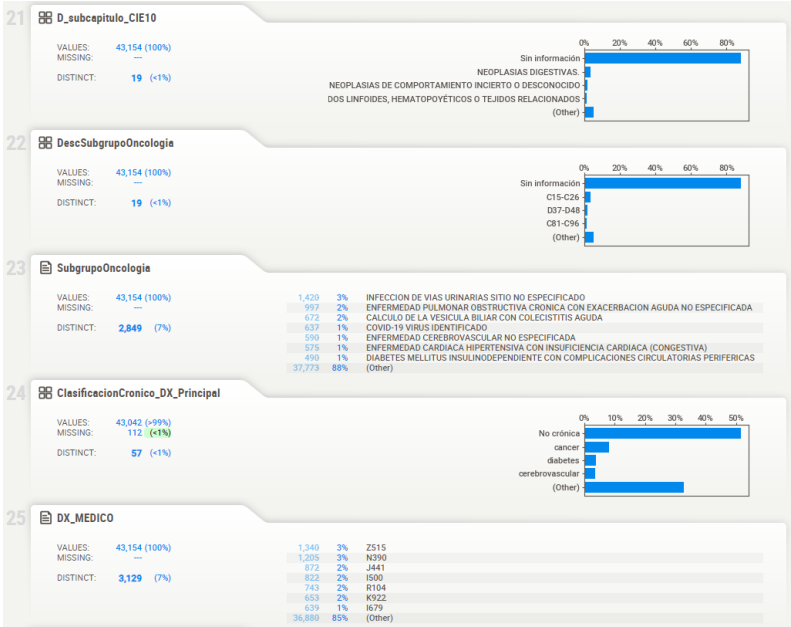
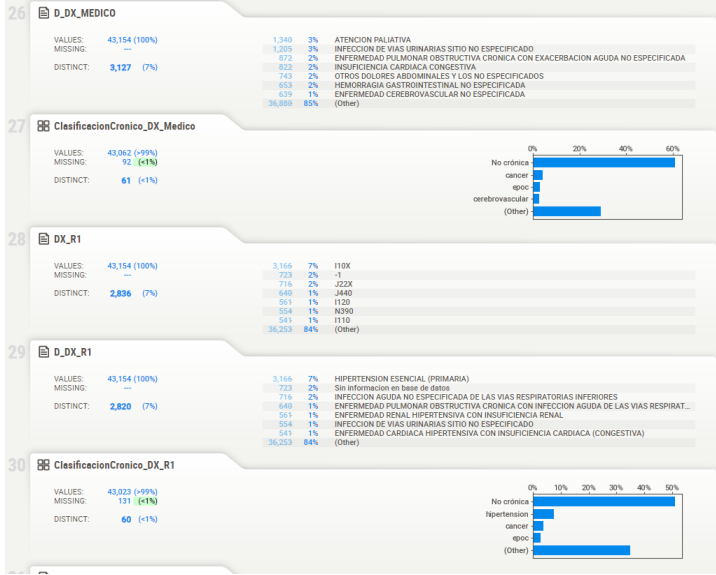
Se puede decir que dos variables tienen una correlación positiva perfecta si el coeficiente es 1 esto es que si una variable incrementa la otra también lo hace,, si el coeficiente es -1 indica que las variables tienen una correlación negativa, esto implica que cuando una incrementa la otra disminuye y viceversa. Cuando el coeficiente es 0 implica que no existe ninguna correlación entre las variables.

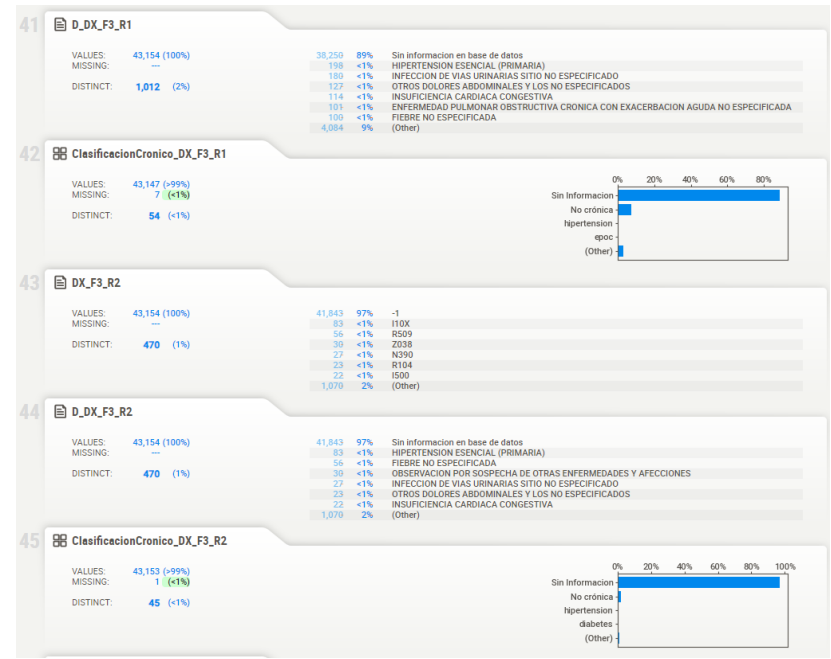
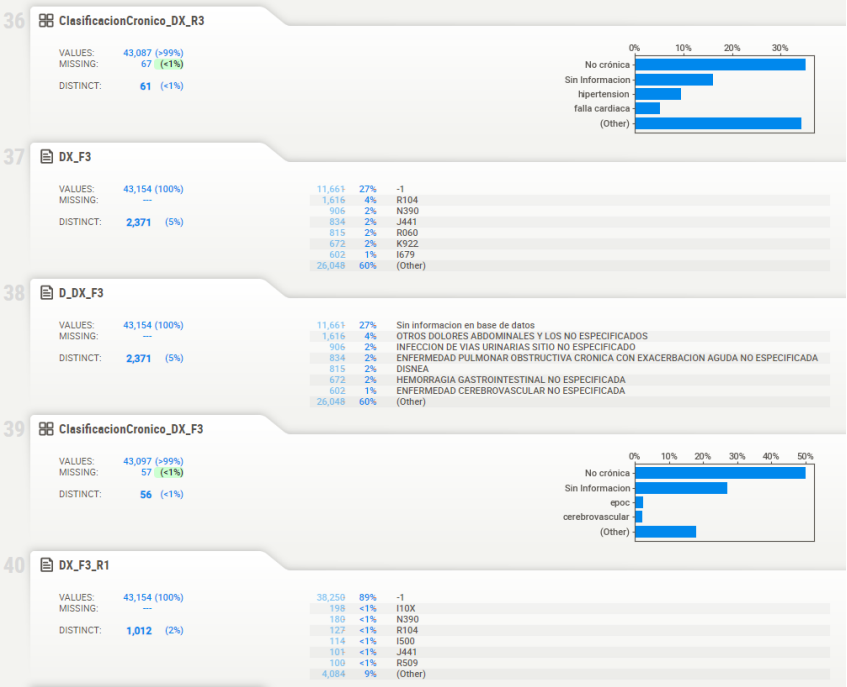
Lo que nos interesa principalmente es encontrar los coeficientes de correlación entre cada variable numérica con la variable días de estancia clínica para así comenzar a descartar variables que no nos aportan información relevante a la predicción de nuestra variable de interés.

Luego, haciendo uso del software sweetviz se analizaron cada una de las variables presentadas en el data set. De esta forma se encontró cuales variables presentaban gran cantidad de datos faltantes o información poco aportante a nuestra variable de salida. Por lo que se eliminaron todas las variables con que tuvieran mas del 80% de los datos como información faltante:









46

DX_F3_R3

VALUES: 43,154 (100%)
MISSING: ---
DISTINCT: 401 (<1%)

42,242 98%
96 <1%
48 <1%
19 <1%
15 <1%
15 <1%
12 <1%
707 2%

-1
Z038
I10X
R060
R509
R104
I500
(Other)

47

D_DX_F3_R3

VALUES: 43,154 (100%)
MISSING: ---
DISTINCT: 401 (<1%)

42,242 98%
96 <1%
48 <1%
19 <1%
15 <1%
15 <1%
12 <1%
707 2%

Sin informacion en base de datos
OBSERVACION POR SOSPECHA DE OTRAS ENFERMEDADES Y AFECCIONES
HIPERTENSION ESENCIAL (PRIMARIA)
DISNEA
FIEBRE NO ESPECIFICADA
OTROS DOLORES ABDOMINALES Y LOS NO ESPECIFICADOS
INSUFICIENCIA CARDIACA CONGESTIVA
(Other)

48

ClasificacionCronico_DX_F3_R3

VALUES: 43,152 (>99%)
MISSING: 2 (<1%)
DISTINCT: 44 (<1%)

0% 20% 40% 60% 80% 100%

Sin informacion
No crónica
hipertension
epoc
(Other)

49

DX_pre_cx

VALUES: 43,154 (100%)
MISSING: ---
DISTINCT: 1,461 (3%)

29,123 67%
482 1%
449 1%
295 <1%
247 <1%
239 <1%
204 <1%
12,115 28%

-1
N201
K800
S721
N40X
I702
N200
(Other)

50

D_DX_pre_cx

VALUES: 43,154 (100%)
MISSING: ---
DISTINCT: 1,461 (3%)

29,123 67%
482 1%
449 1%
295 <1%
247 <1%
239 <1%
204 <1%
12,115 28%

Sin informacion en base de datos
CALCULO DEL URETER
CALCULO DE LA VESICULA BILIAR CON COLECISTITIS AGUDA
FRACTURA PERTROCANTERIANA
HIPERPLASIA DE LA PROSTATA
ATEROSCLEROSIS DE LAS ARTERIAS DE LOS MIEMBROS
CALCULO DEL RINON
(Other)

51

ClasificacionCronico_DX_pre_cx

VALUES: 43,123 (>99%)
MISSING: 31 (<1%)
DISTINCT: 44 (<1%)

0% 20% 40% 60%

Sin informacion
No crónica
cancer
vascular periférica
(Other)

52

DX_pos_cx

VALUES: 43,154 (100%)
MISSING: ---
DISTINCT: 1,492 (3%)

29,123 67%
442 1%
411 <1%
298 <1%
252 <1%
241 <1%
226 <1%
12,161 28%

-1
N201
K800
S721
N40X
I702
N200
(Other)

53

D_DX_pos_cx

VALUES: 43,154 (100%)
MISSING: ---
DISTINCT: 1,492 (3%)

29,123 67%
442 1%
411 <1%
298 <1%
252 <1%
241 <1%
226 <1%
12,161 28%

Sin informacion en base de datos
CALCULO DEL URETER
CALCULO DE LA VESICULA BILIAR CON COLECISTITIS AGUDA
FRACTURA PERTROCANTERIANA
HIPERPLASIA DE LA PROSTATA
ATEROSCLEROSIS DE LAS ARTERIAS DE LOS MIEMBROS
CALCULO DEL RINON
(Other)

54

ClasificacionCronico_DX_pos_cx

VALUES: 43,126 (>99%)
MISSING: 28 (<1%)
DISTINCT: 43 (<1%)

0% 20% 40% 60%

Sin informacion
No crónica
cancer
otras genitourinarias
(Other)

55

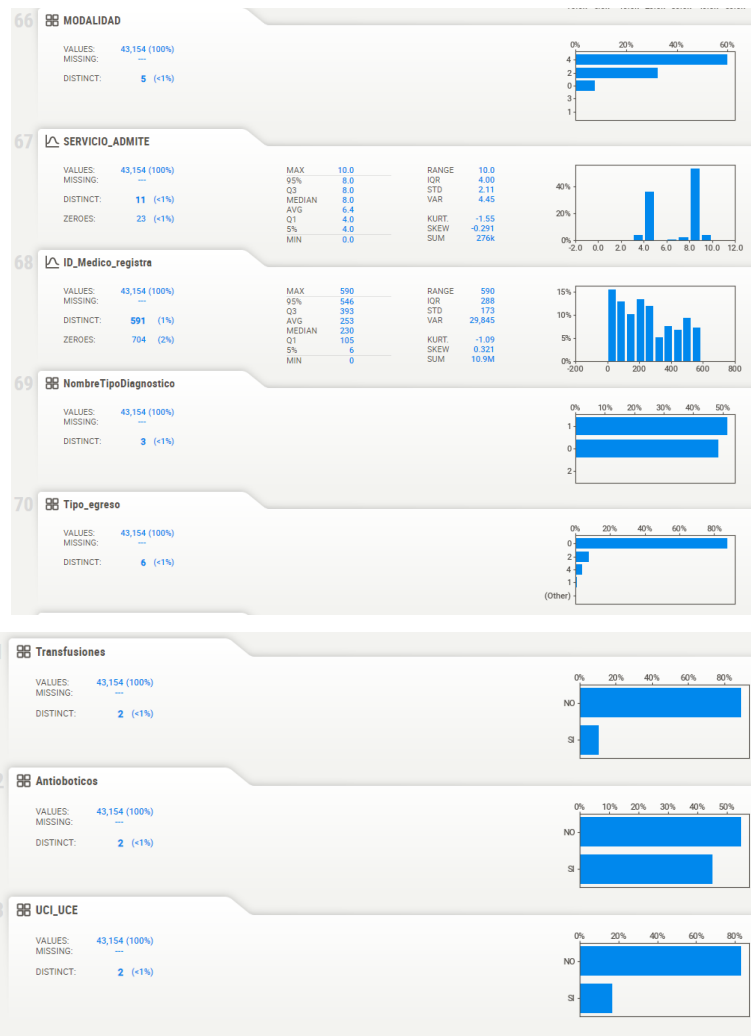
DX_MUERTE

VALUES: 43,154 (100%)
MISSING: ---
DISTINCT: 576 (1%)

39,736 92%
303 <1%
192 <1%
174 <1%
129 <1%
128 <1%
104 <1%
2,388 6%

-1
J440
U071
I132
I120
C349
I219
(Other)

56	<div><div>D_DX_MUERTE</div><div><div><div>VALUES: 43,154 (100%)</div><div>MISSING: ---</div><div>DISTINCT: 575 (1%)</div></div><div><div>39,736 92%</div><div>303 <1%</div><div>192 <1%</div><div>174 <1%</div><div>129 <1%</div><div>128 <1%</div><div>104 <1%</div><div>2,388 6%</div></div><div><div>Sin informacion en base de datos</div><div>ENFERMEDAD PULMONAR OBSTRUCTIVA CRONICA CON INFECCION AGUDA DE LAS VIAS RESPIRAT...</div><div>COVID-19 VIRUS IDENTIFICADO</div><div>ENFERMEDAD CARDIORRENAL HIPERTENSIVA CON INSUFICIENCIA CARDIACA Y RENAL (CONGESTIV...</div><div>ENFERMEDAD RENAL HIPERTENSIVA CON INSUFICIENCIA RENAL</div><div>TUMOR MALIGNO DE LOS BRONQUIOS O DEL PULMON PARTE NO ESPECIFICADA</div><div>INFARTO AGUDO DEL MIOCARDIO SIN OTRA ESPECIFICACION</div><div>(Other)</div></div></div></div>
57	<div><div>D_Subcapitulo_DX_Muerte</div><div><div><div>VALUES: 43,154 (100%)</div><div>MISSING: ---</div><div>DISTINCT: 296 (<1%)</div></div><div><div>39,736 92%</div><div>306 <1%</div><div>194 <1%</div><div>181 <1%</div><div>146 <1%</div><div>129 <1%</div><div>118 <1%</div><div>2,344 5%</div></div><div><div>Sin informacion en base de datos</div><div>OTRAS ENFERMEADES PULMONARES OBSTRUCTIVAS CRONICAS</div><div>USO DE EMERGENCIA DE U07</div><div>ENFERMEDAD CARDIORRENAL HIPERTENSIVA</div><div>TUMOR MALIGNO DE LOS BRONQUIOS Y DEL PULMON</div><div>ENFERMEDAD RENAL HIPERTENSIVA</div><div>INFARTO AGUDO DEL MIOCARDIO</div><div>(Other)</div></div></div></div>
58	<div><div>Subcapitulo_DX_Muerte</div><div><div><div>VALUES: 43,144 (>99%)</div><div>MISSING: 10 (<1%)</div><div>DISTINCT: 296 (<1%)</div></div><div><div>39,736 92%</div><div>306 <1%</div><div>194 <1%</div><div>181 <1%</div><div>146 <1%</div><div>129 <1%</div><div>118 <1%</div><div>2,344 5%</div></div><div><div>-1</div><div>J44</div><div>U07</div><div>I13</div><div>C34</div><div>I12</div><div>I21</div><div>(Other)</div></div></div></div>
59	<div><div>capitulo_DX_Muerte</div><div><div><div>VALUES: 43,144 (>99%)</div><div>MISSING: 10 (<1%)</div><div>DISTINCT: 16 (<1%)</div></div><div><div>0%20%40%60%80%</div><div>-1</div><div>I00-I99</div><div>C00-D48</div><div>J00-J99</div><div>(Other)</div></div></div></div>
60	<div><div>D_capitulo_DX_Muerte</div><div><div><div>VALUES: 43,154 (100%)</div><div>MISSING: ---</div><div>DISTINCT: 17 (<1%)</div></div><div><div>0%20%40%60%80%</div><div>Sin informacion en base de datos</div><div>ENFERMEADES DEL SISTEMA CIRCULATORIO</div><div>TUMORES (NEOPLASIAS)</div><div>ENFERMEADES DEL SISTEMA RESPIRATORIO</div><div>(Other)</div></div></div></div>
61	<div><div>ClasificacionCronico_DX_MUERTE</div><div><div><div>VALUES: 43,143 (>99%)</div><div>MISSING: 11 (<1%)</div><div>DISTINCT: 38 (<1%)</div></div><div><div>0%20%40%60%80%</div><div>Sin Informacion</div><div>cancer</div><div>No cronica</div><div>cerebrovascular</div><div>(Other)</div></div></div></div>
62	<div><div>HoresdeAlta</div><div><div><div>VALUES: 43,154 (100%)</div><div>MISSING: ---</div><div>DISTINCT: 41 (<1%)</div><div>ZEROS: 3,154 (7%)</div></div><div><div>MAX 11,448</div><div>95% 120</div><div>Q3 72</div><div>AVG 50</div><div>MEDIAN 48</div><div>Q1 24</div><div>5% 0</div><div>MIN -8,376</div></div><div><div>RANGE 19,824</div><div>IQR 48.0</div><div>STD 78.5</div><div>VAR 6,168</div><div>KURT. 13,394</div><div>SKEW 43.8</div><div>SUM 2.1M</div></div><div><div>0%50%100%</div><div>-10.0k-5.0k0.0k5.0k10.0k15.0k</div></div></div></div>
63	<div><div>DiasEstanciaServicio</div><div><div><div>VALUES: 43,154 (100%)</div><div>MISSING: ---</div><div>DISTINCT: 84 (<1%)</div><div>ZEROS: 1,485 (3%)</div></div><div><div>MAX 122</div><div>95% 18</div><div>Q3 7</div><div>AVG 6</div><div>MEDIAN 4</div><div>Q1 2</div><div>5% 1</div><div>MIN 0</div></div><div><div>RANGE 122</div><div>IQR 5.00</div><div>STD 6.52</div><div>VAR 42.6</div><div>KURT. 22.7</div><div>SKEW 3.53</div><div>SUM 245k</div></div><div><div>0%25%50%75%</div><div>-250255075100125150</div></div></div></div>
64	<div><div>DiasEstanciaClinica</div><div><div><div>VALUES: 43,154 (100%)</div><div>MISSING: ---</div><div>DISTINCT: 134 (<1%)</div><div>ZEROS: 331 (<1%)</div></div><div><div>MAX 207</div><div>95% 28</div><div>Q3 11</div><div>AVG 9</div><div>MEDIAN 5</div><div>Q1 3</div><div>5% 1</div><div>MIN 0</div></div><div><div>RANGE 207</div><div>IQR 8.00</div><div>STD 10.9</div><div>VAR 119</div><div>KURT. 34.2</div><div>SKEW 4.30</div><div>SUM 381k</div></div><div><div>0%25%50%75%</div><div>-50050100150200250</div></div></div></div>
65	<div><div>DiasEstanciaClinica_exacto</div><div><div><div>VALUES: 43,154 (100%)</div><div>MISSING: ---</div><div>DISTINCT: 43,122 (>99%)</div><div>ZEROS: ---</div></div><div><div>MAX 45,069</div><div>95% 28</div><div>Q3 11</div><div>AVG 11</div><div>MEDIAN 6</div><div>Q1 3</div><div>5% 1</div><div>MIN 0</div></div><div><div>RANGE 45,069</div><div>IQR 7.94</div><div>STD 306</div><div>VAR 93,627</div><div>KURT. 21,521</div><div>SKEW 147</div><div>SUM 472k</div></div><div><div>0%50%100%</div><div>-10.0k0.0k10.0k20.0k30.0k40.0k50.0k</div></div></div></div>



A partir del mapa de calor con correlaciones y los diagramas de frecuencia se opta por eliminar las siguientes variables puesto que no aportan información valiosa a la predicción de nuestra variable de salida ya sea por datos faltantes o falta de correlación:

'AtencionPaciente llave', 'Año_Salida', 'MES','Aseguradora', 'Aseguradora2', 'Modalidad_Contrato', 'DX_F3_R3', 'DX_F3_R3', 'bloque', 'piso', 'RangoCapitulo', 'Categoria', 'D_subcapitulo_CIE10','DescSubgrupoOncologia', 'SubgrupoOncologia', 'DX_MEDICO', 'D_DX_MEDICO','D_Subcapitulo_DX_Muerte','Subcapitulo_DX_Muerte','ClasificacionCronico_DX_F3_R2','capitulo_DX_Muerte','D_capitulo_DX_Muerte','ClasificacionCronico_DX_F3_R1','ClasificacionCronico_DX_MUERTE', 'DX_R1', 'D_DX_R2', 'D_DX_R3', 'DX_F3_R3', 'D_DX_F3_R2', 'DiasEstanciaClinica_exacto', 'D_DX_F3', 'DX_F3_R1','D_DX_F3_R1', 'SERVICIO_ADMITE', 'ID_Medico_registra', 'NombreTipoDiagnostico', 'D_DX_pre_cx', 'DX_pre_cx', 'vGenero', 'Tipo_egreso'

Modificaciones a la base de datos

Enfocandonos en el diagnostico principal del paciente como una de las variables determinantes para la predicción se procede a eliminar todas las variables relacionada con diagnosticos secundarios y se obtiene el siguiente data set, con las siguientes variables:

#	Column	Non-Null Count	Dtype
0	EdadEgreso	43154 non-null	float64
1	Tipo_Afiliacion	43154 non-null	int64
2	Hospitalizacion	43154 non-null	int64
3	Tipo_Internacion	42945 non-null	object
4	Ciclo_Vital	43154 non-null	int64
5	ESPECIALIDAD_GRD	43154 non-null	int64
6	DX_PRINCIPAL	43154 non-null	object
7	HorasdeAlta	43154 non-null	int64
8	DiasEstanciaServicio	43154 non-null	int64
9	DiasEstanciaClinica	43154 non-null	int64
10	MODALIDAD	43154 non-null	int64
11	Transfusiones	43154 non-null	int64
12	Antibioticos	43154 non-null	int64
13	UCI_UCE	43154 non-null	int64

Luego se hizo una codificación de las variables categóricas y normalización de las variables numéricas para que así la base de datos se adapte mejor a un modelo predictivo. Luego de dichas modificaciones se obtuvo una base de datos con 2869 columnas y el mismo numero de filas. Se plantearon los siguientes modelos:

Random forest para la variable días de estancia clínica:

```
Error Cuadrático Medio: 56.572474062210226
Coeficiente de Determinación (R^2): 0.5642520704092626
Mean Absolute Error (MAE): 2.747564384932129
```

Random forest para la variable días de estancia servicio:

```
Error Cuadrático Medio: 16.457387860388785
Coeficiente de Determinación (R^2): 0.6284483985817577

Mean Absolute Error (MAE): 1.8653939599580893
```

Según los resultados obtenidos mediante el random forest se encontró una predicción para los días de estancia clínica con un error de mas o menos 2.75 días y para los días de estancia de servicio se encontró un error de mas o menos 1.87 días. .

Árbol de decisiones para la variable días de estancia clínica:

Error Cuadrático Medio: 77.53201959345512
Mean Absolute Error (MAE): 3.2090333294712856

Árbol de decisiones para la variable días de estancia servicio:

Error Cuadrático Medio: 28.399808828640946
Mean Absolute Error (MAE): 2.205711968485691

Regresión lineal para la variable días de estancia clínica:

Error Cuadrático Medio: 2.761431616846068e+18
Coeficiente de Determinación (R^2): -2.1269851278271376e+16
Mean Absolute Error (MAE): 2.6900858113605524

Regresión lineal para la variable días de estancia servicio:

Error Cuadrático Medio: 22.4502364330495
Coeficiente de Determinación (R^2): 0.4931503486653164
Mean Absolute Error (MAE): 1.8653939599580893

sdads

Red Neuronal

270/270 [=====] - 1s 5ms/step - loss: 25.7859 - mae: 2.4953
Error Cuadrático Medio: 25.78586769104004
Error absoluto medio: 2.495349407196045

Conclusiones

Días de Estancia Clínica:

El modelo de Árbol de Decisiones tiene un Error Cuadrático Medio (MSE) de aproximadamente 77.53 y un MAE de alrededor de 3.21, lo que indica un rendimiento moderado.

El modelo de Random Forest muestra un MSE de alrededor de 56.57 y un MAE de alrededor de 2.75, lo que indica un mejor rendimiento en comparación con el Árbol de Decisiones.

El modelo de Regresión Lineal tiene un MSE extremadamente alto, lo que sugiere un mal ajuste del modelo a los datos en este caso.

Días de Estancia Servicio:

El modelo de Árbol de Decisiones tiene un MSE de aproximadamente 28.40 y un MAE de alrededor de 2.21, lo que muestra un rendimiento moderado.

El modelo de Random Forest muestra un MSE de alrededor de 16.46 y un MAE de aproximadamente 1.87, lo que indica un mejor rendimiento en comparación con el Árbol de Decisiones.

El modelo de Regresión Lineal tiene un MSE de alrededor de 22.45 y un MAE de aproximadamente 1.87, lo que indica un rendimiento razonable para la Regresión Lineal en este caso.

Red Neuronal:

La red neuronal tiene un MSE de aproximadamente 25.79 y un MAE de alrededor de 2.50. Si bien no es el mejor rendimiento en términos de MAE, la red neuronal es capaz de capturar relaciones más complejas en los datos.

En general, el modelo de Random Forest parece funcionar bien en ambos casos (Días de Estancia Clínica y Días de Estancia Servicio), proporcionando predicciones precisas con MAE bajos. La Regresión Lineal y el Árbol de Decisiones también dan resultados aceptables en algunos casos, pero el modelo de Red Neuronal puede ser una opción interesante para capturar relaciones más complejas si estás dispuesto a invertir en su entrenamiento y ajuste de hiperparámetros.

Se podrían seguir ajustando los hiperparámetros para mejorar aún más las predicciones.