

CIENCIA Y ANALÍTICA DE DATOS

Reto

Profesor: María de la Paz Rico Fernández

Alumno: Juan Sebastián Ortega Briones A01794327

Equipo 13

18 de Noviembre del 2022

AGENDA

Datos

Limpieza

Análisis

Kmeans

Clasificación

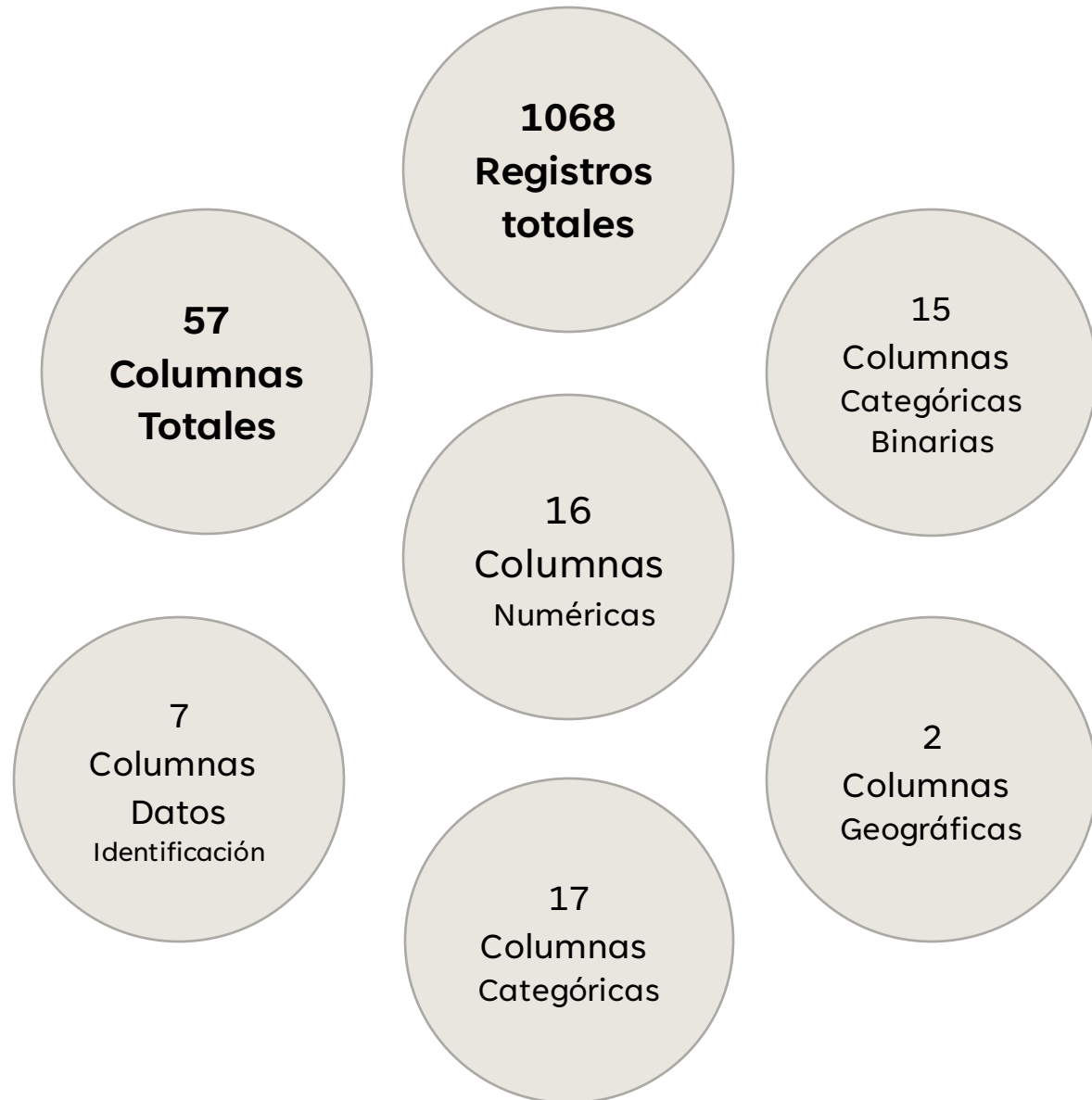
Resultados

Conclusiones

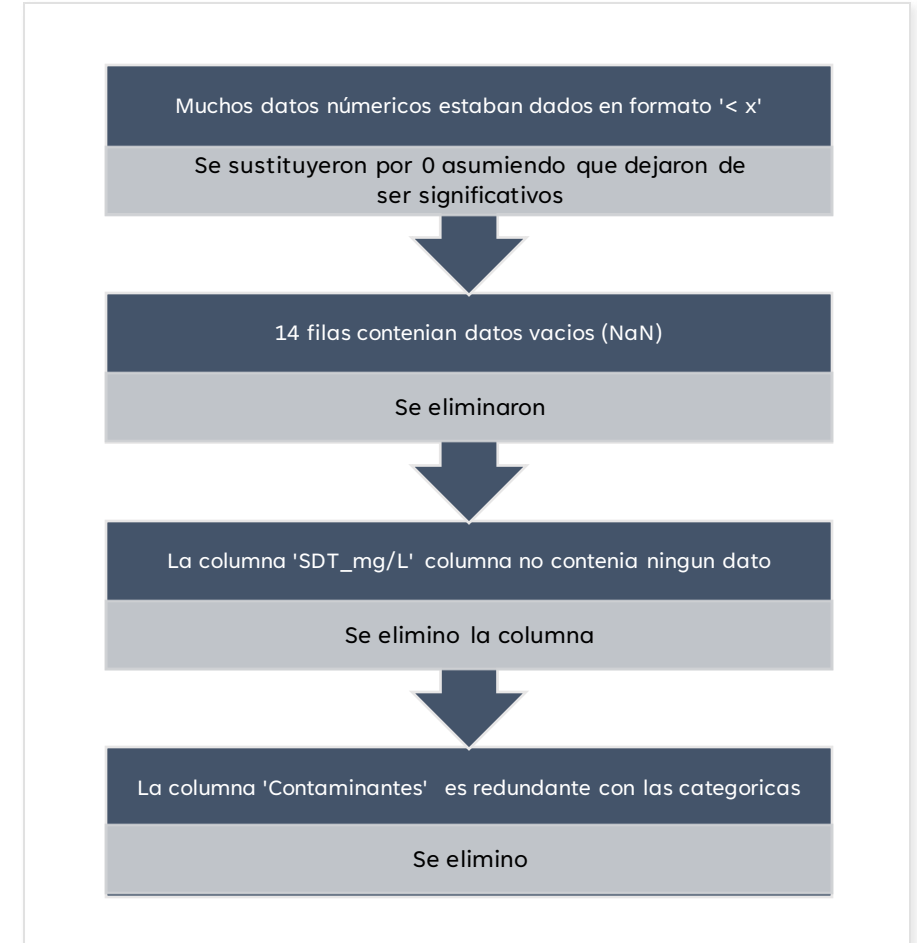
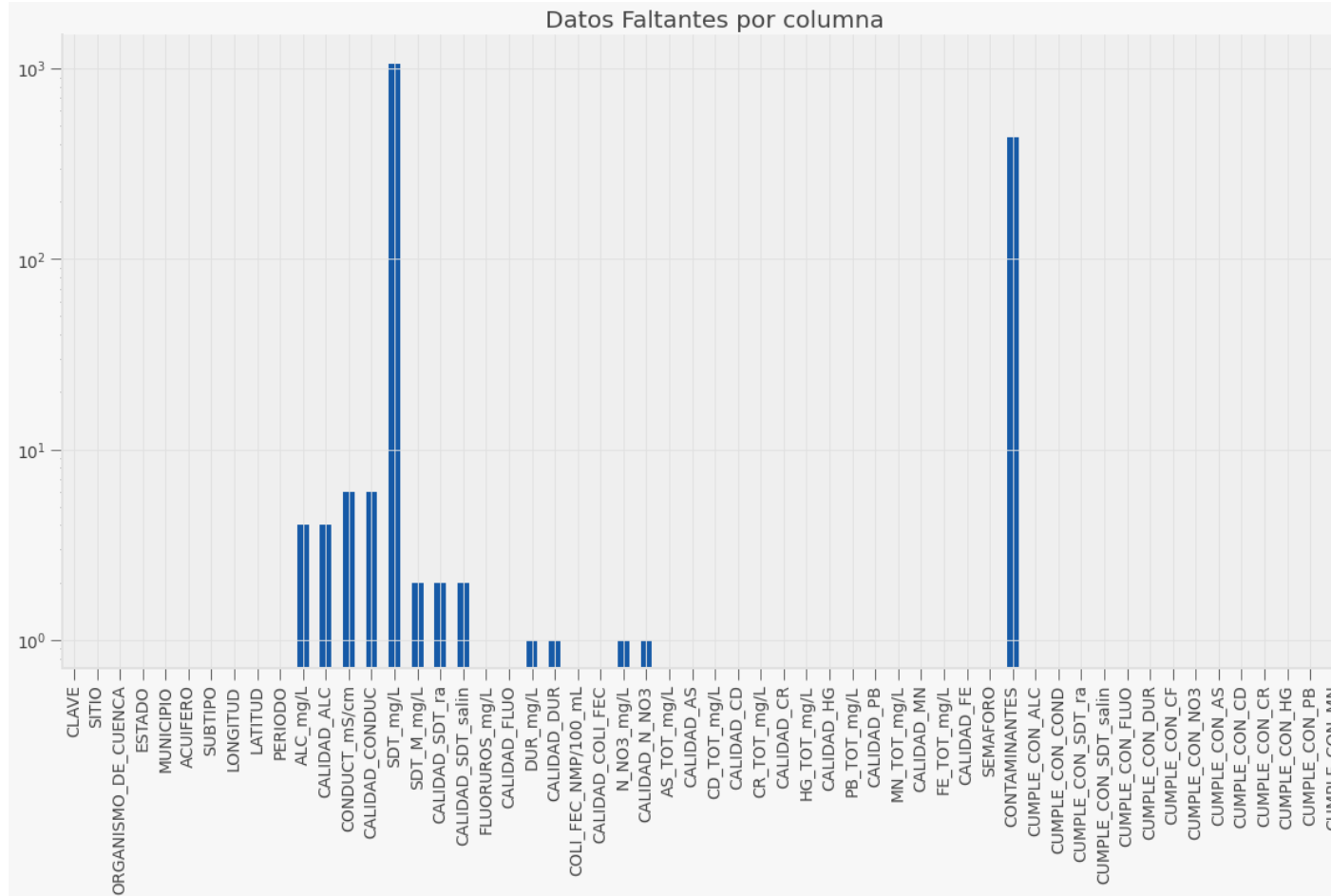


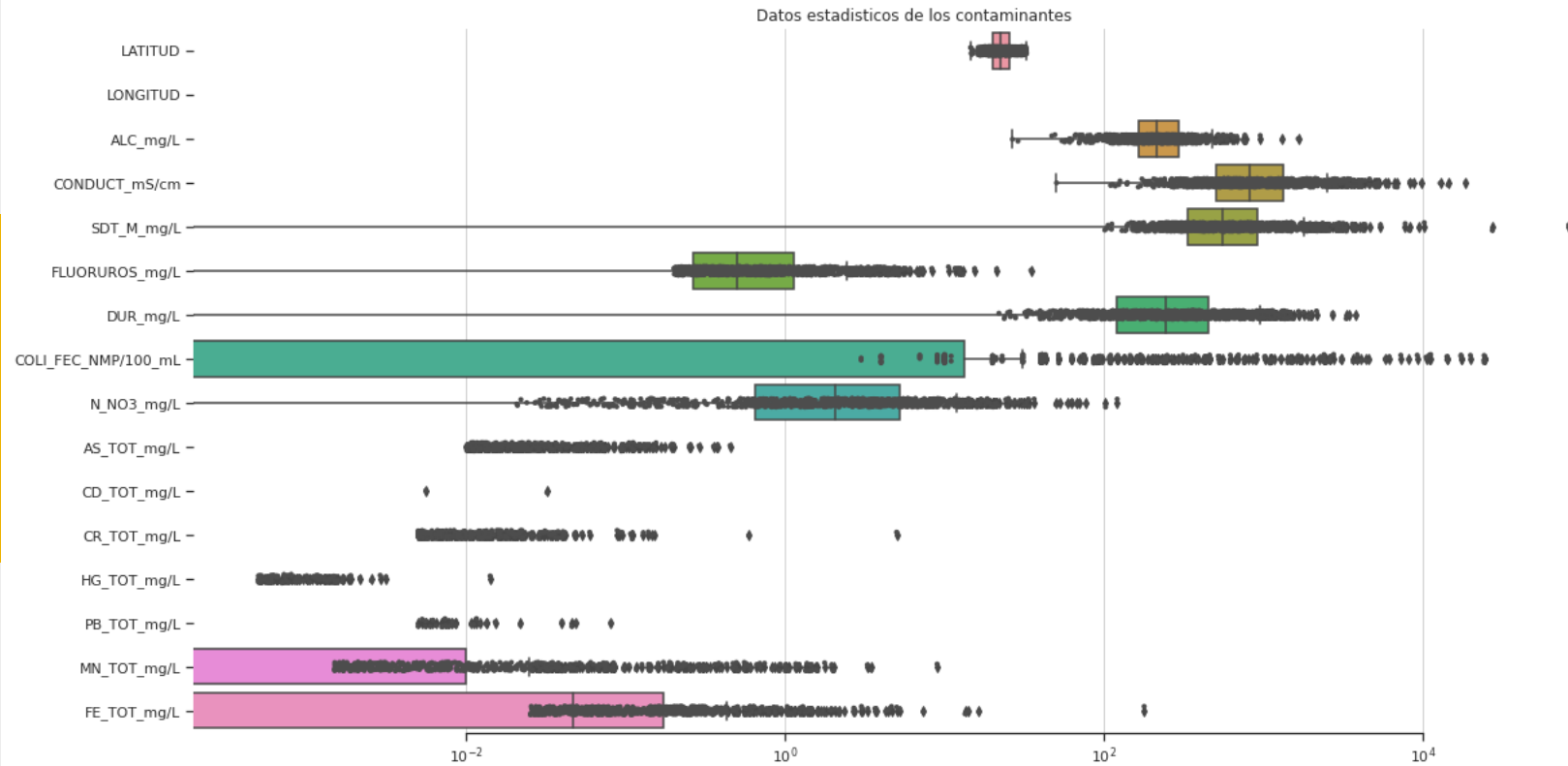
Los datos:

Datos de calidad del agua
subterránea de 5000
estaciones a nivel nacional
del año 2020



LIMPIEZA

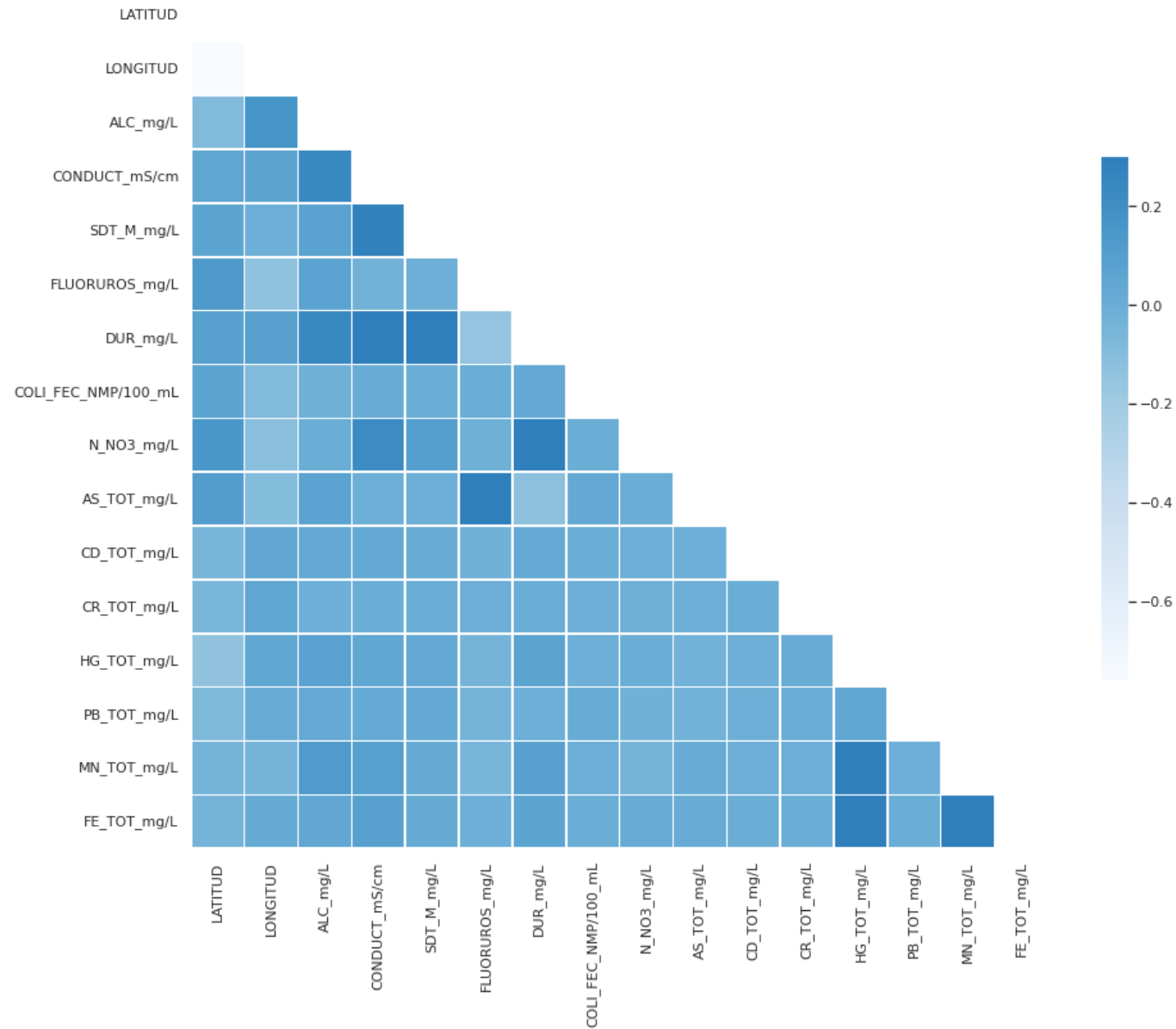




ANALISIS

Calidad del
agua subterránea 2022

Correlaciones de los datos



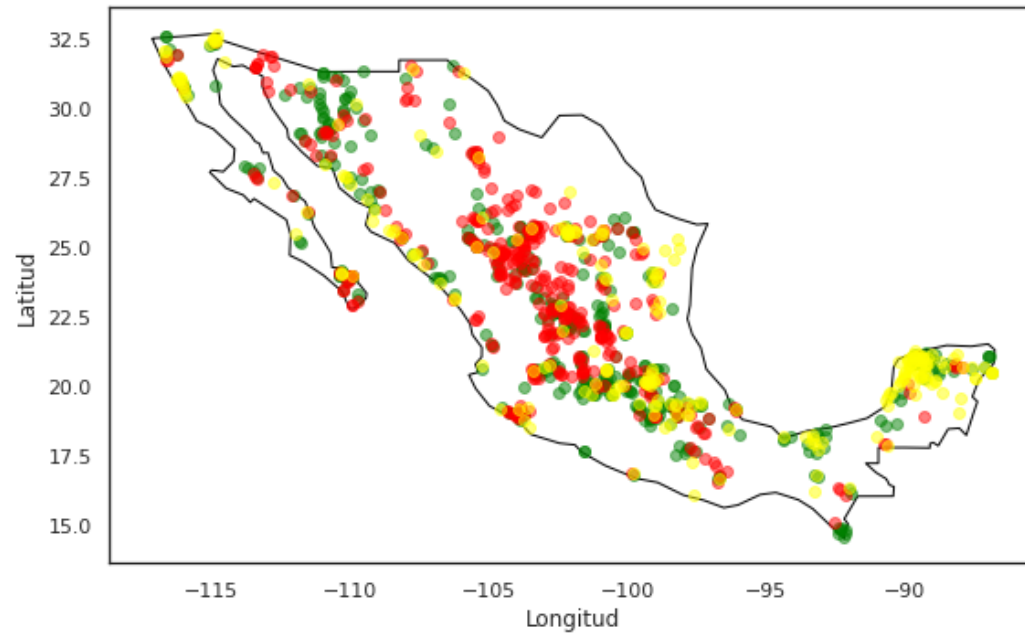
ANALISIS

Calidad del
agua subterránea 2022

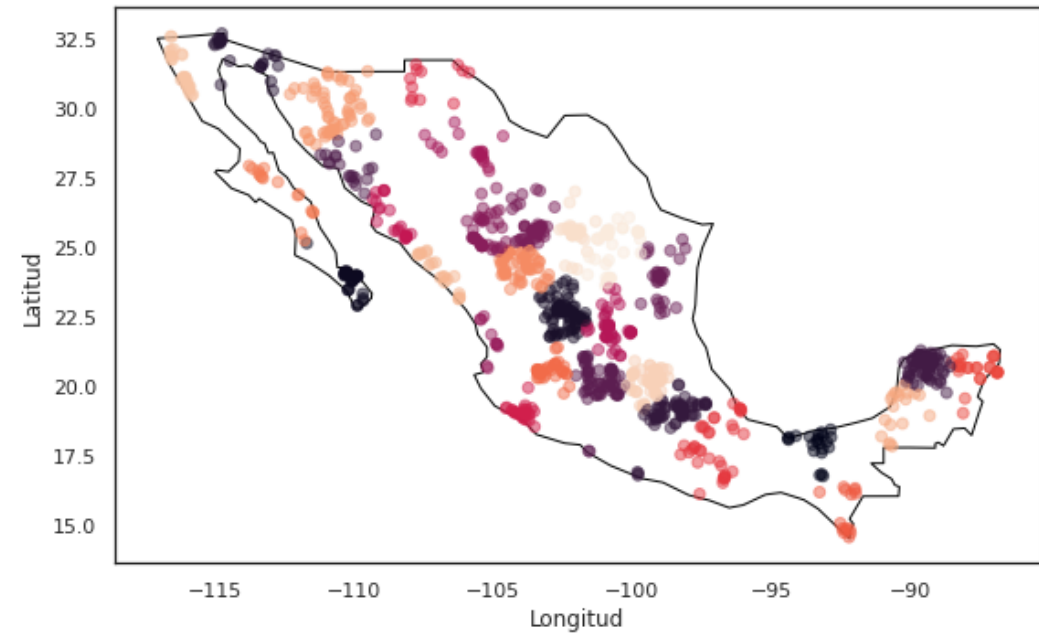
KMEANS

Comparativo de semáforo contra predicción de kmeans

Semáforo de contaminantes por estación de medición de agua



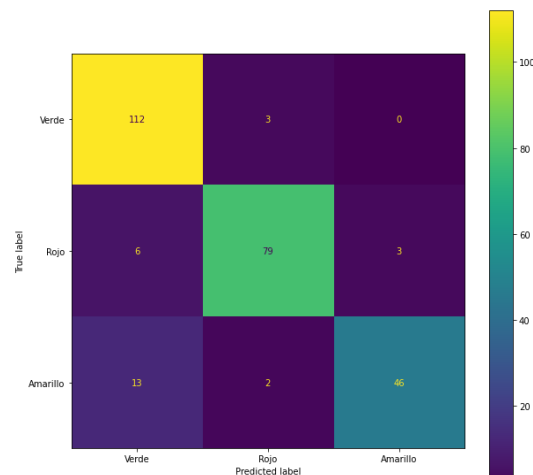
Agrupamiento por K-means



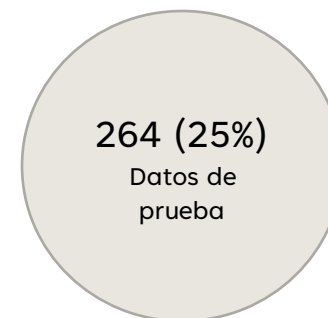
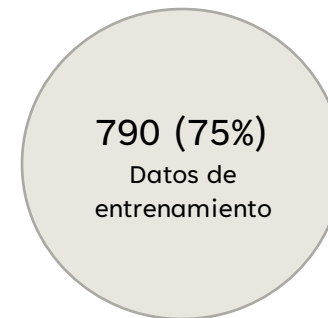
CLASIFICACIÓN

Después del análisis, se puede ver que el método más eficiente es Random Forest

Metodo	Eficiencia
LogisticRegression	0.84
RandomForestClassifier	0.97
SVC	0.60
VotingClassifier	0.90



Matrix de confusión



Variables más importantes

Porcentaje de importancia	Variable
43.42 %	DUR_mg/L
19.18 %	MN_TOT_mg/L
17.55 %	FLUORUROS_mg/L
12.18 %	FE_TOT_mg/L
4.37 %	N_NO3_mg/L
2.52 %	COLI_FEC_NMP/100_mL
0.74 %	AS_TOT_mg/L

RESULTADOS

Datos y análisis

La base de datos es muy buena, tiene muy pocos datos vacíos

La correlación más alta es entre dureza y fluoruros.

Es interesante la correlación entre el plomo y la latitud que es la segunda más alta.

kmeans

Usando kmeans con latitudes y longitudes me parece que no se observa una relación clara entre custers y calidad del agua, aunque si es clara en clusters y cercanía de estaciones de medición.

Clasificación

La predicción usando los métodos de clasificación más eficiente fue random forest y la menos eficiente fue SVC.

Los resultados no cambiaron al intercambiar los métodos de votación.

CONCLUSIONES

Usando el método Random Forest es posible determinar con gran eficiencia (>90%), la potabilidad del agua subterránea usando las 7 variables principales.

Para determinar clusters geográficos que determinen la calidad del agua es necesario usar estas variables en kmeans: DUR_mg/L, MN_TOT_mg/L, FLUORUROS_mg/L, FE_TOT_mg/L