

Reporte: Exploración en aguas superficiales

EQUIPO 104

Integrantes:

Eddie Guadalupe Elorza Ruiz | **A01793547**

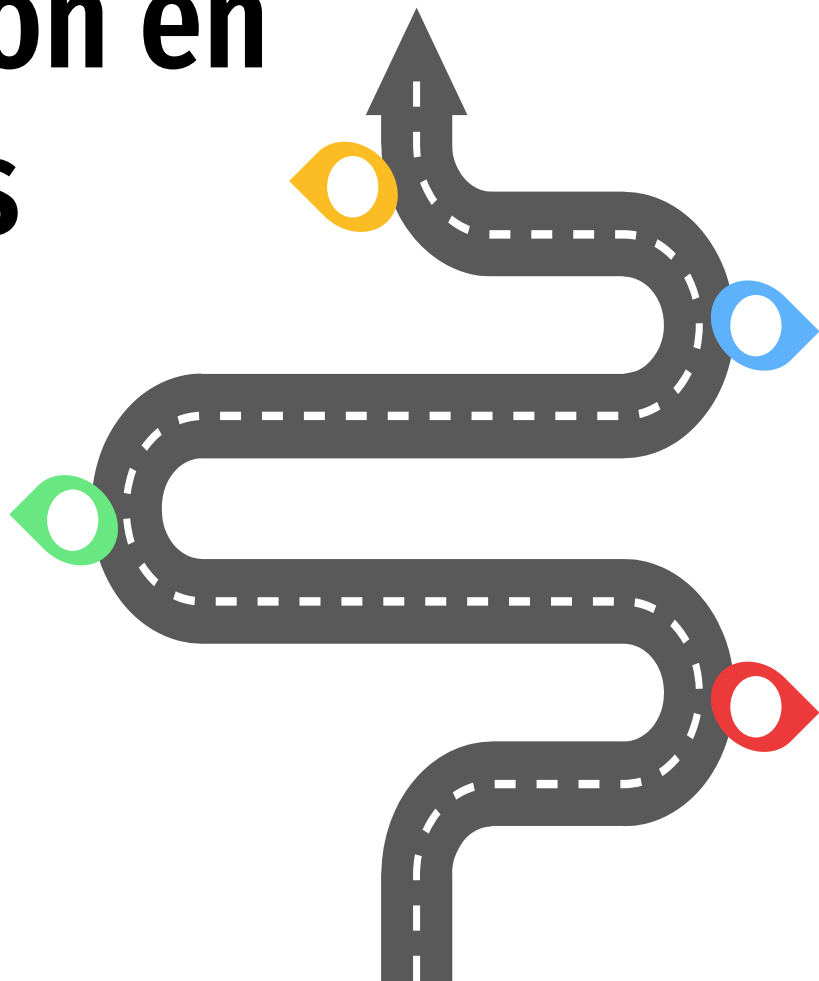
Yves Turley Macias Vargas. | **A00813752**

Materia: Ciencia y analítica de datos (Gpo 10)

Profesor Titular: PhD. María de la Paz Rico Fdz

Profesor Tutor: Victoria Guerrero Orozco

18 de noviembre 2022

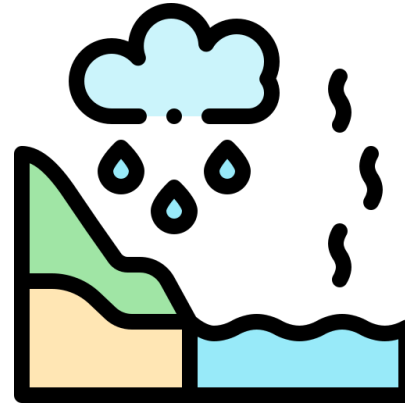
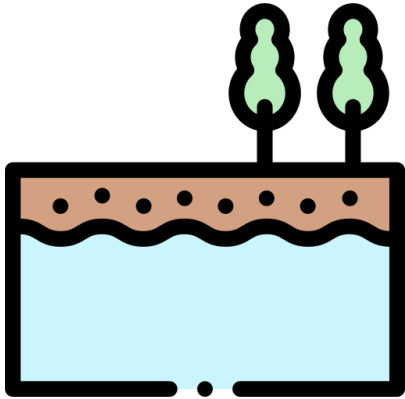


Metodología



Selección de Dataset

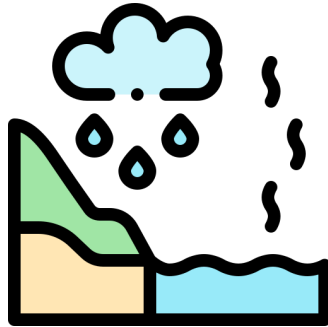
La primera tarea asignada fue la de seleccion de un dataset. Se nos dieron dos opciones. La de cuerpos de agua superficiales y cuerpos de agua subterrneos en México.



Justificación de selección de “Dataset”

Decidimos optar por el de cuerpos de agua superficiales en México.

El dataset es un campo de datos pero con mucha oportunidad para aplicar diferentes técnicas de limpieza y poner en practica métodos aprendidos a lo largo del curso. Se puede visualizar espacios en blanco en la imagen.



	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W		
	CLAVE	SITIO	ORGANISMA	ESTADO	MUNICIPIO	CUENCA	CUERPO	D TIPO	SUBTIPO	LONGITUD	LATITUD	PERIODO	DBO_mg/l	CALIDAD	DOQ_mg/l	CALIDAD	SS_T_mg/l	CALIDAD	COLI_FEC	CALIDAD	E_COLI_N	CALIDAD	ENTEROC	CF	
1	DLBAJ08	PRESA EL LERMA	SA AGUASAL	RINCÓN D RIO	SAN PRESA EL L	SENTICO	PRESA			-102.339	22.2473	2020		6 Buena cali	54.08	Contaminu	13.75	Excelente	1162	Contaminu	98	Excelente			
2	DLBAJ09	LOS CABO	PENINSUL	BAJA CALI	LOS CABO	SAN JOSE	OCEANO	P	COSTERO	OCEANO	P	2020					<10	Excelente					20	Ex	
3	DLBAJ10	LOS CABO	PENINSUL	BAJA CALI	LOS CABO	SAN LUCA	OCEANO	F	COSTERO	OCEANO	F	2020					<10	Excelente					<3	Ex	
4	DLBAJ102	LOS CABO	PENINSUL	BAJA CALI	LOS CABO	SAN LUCA	BAHIA	SAN COSTERO	BAHIA	-109.886	22.89609	2020					13.9667	Excelente					<3	Ex	
5	DLBAJ103	LOS CABO	PENINSUL	BAJA CALI	LOS CABO	SAN LUCA	BAHIA	SAN COSTERO	BAHIA	-109.897	22.87694	2020					<10	Excelente					30	Ex	
6	DLBAJ104	BAHIA CAE	PENINSUL	BAJA CALI	LOS CABO	SAN LUCA	BAHIA	SAN COSTERO	BAHIA	-109.903	22.88	2020					22.0667	Excelente					<3	Ex	
7	DLBAJ105	LOS CABO	PENINSUL	BAJA CALI	LOS CABO	SAN LUCA	BAHIA	SAN COSTERO	BAHIA	-109.905	22.8831	2020					13.9667	Excelente					90	Ex	
8	DLBAJ106	LAGUNA T	PENINSUL	BAJA CALI	LA PAZ	TODOS SA	TODOS SA	COSTERO	LAGUNA	-110.239	23.393	2020					57.85	Buena calidad					402	Cc	
9	DLBAJ109	MANANTIL	PENINSUL	BAJA CALI	LA PAZ	TODOS SA	TODOS SA	LITOTIC	ARRIOYO	-110.224	23.45805	2020	<2	Excelente	<10	Excelente	33.9	Buena cali	3873	Contaminu	512	Buena calidad			
10	DLBAJ112	AGUA CAL	PENINSUL	BAJA CALI	LOS CABO	SANTO	SANTO	SANTIAGO	LITOTIC	ARRIOYO	-109.808	23.43995	2020	<2	Excelente	<10	Excelente	25.6	Buena cali	189	Buena cali	<3	Excelente		
11	DLBAJ120	BOCA DE L	PENINSUL	BAJA CALI	LOS CABO	SAN JOSE	ISAN JOSE	LITOTIC	H	ARRIOYO	-109.826	23.39128	2020	<2	4.26	Buena cali	27.98	Acceptable	14	Excelente	1408	Contaminu	84	Excelente	
12	DLBAJ121	HUMEDAL	PENINSUL	BAJA CALI	LA PAZ	LAS POICIT	HUMEDAL	LITOTIC	H	ARRIOYO	-110.952	24.50289	2020	<2	Excelente	<10	Excelente	<10	Excelente	15531	Fuertemer	538	Buena calidad		
13	DLBAJ123	OASIS DE L	PENINSUL	BAJA CALI	LA PAZ	LAS POICIT	CUERPO	D	LITOTIC	H	ARRIOYO	-111.003	24.46953	2020	<2	Excelente	<10	Excelente	<10	Excelente	10	Excelente	<3	Excelente	
14	DLBAJ124	MANANTIL	PENINSUL	BAJA CALI	LOS CABO	SAN JOSE	ISAN JOSE	LITOTIC	H	CANAL	-109.779	23.337	2020	<2	6.4	Acceptable	<10	Excelente	<10	Excelente	24196	Fuertemer	14136	Fuertemente contamina	
15	DLBAJ126	MANANTIL	PENINSUL	BAJA CALI	LA PAZ	SAN BARTI	SAN BARTI	LITOTIC	ARRIOYO	-109.845	23.72546	2020	<2	Excelente	<10	Excelente	<10	Excelente		218	Acceptable	<3	Excelente		
16	DLBAJ132	HUMEDAL	PENINSUL	BAJA CALI	COMONDI	SANTO	DO	HUMEDAL	LITOTIC	H	ARRIOYO	-111.805	26.06218	2020	<2	Excelente	<10	Excelente	<10	Excelente	663	Acceptable	74	Excelente	
17	DLBAJ133	ESTERO EL	PENINSUL	BAJA CALI	MULEGE	SAN IGNAI	ESTERO	COSTERO	ESTERO	-113.462	26.80955	2020					<10	Excelente					<3	Ex	
18	DLBAJ134	ESTERO SA	PENINSUL	BAJA CALI	COMONDI	SANTO	DO	ESTERO	COSTERO	ESTERO	-112.093	25.69695	2020					19.5333	Excelente				<3	Ex	
19	DLBAJ135	HUMEDAL	PENINSUL	BAJA CALI	COMONDI	SANTO	DO	HUMEDAL	LITOTIC	H	ARRIOYO	-111.833	26.03686	2020	<2	Excelente	<10	Excelente	10.1	Excelente	1408	Contaminu	368	Buena calidad	
20	DLBAJ138	PARQUE N	PENINSUL	BAJA CALI	LORETO	LORETO	GOLFO	DE	COSTERO	BAHIA	-111.235	25.7233	2020				<10	Excelente					<3	Ex	
21	DLBAJ139	BAHIA DE	PENINSUL	BAJA CALI	LORETO	LORETO	GOLFO	DE	COSTERO	BAHIA	-111.256	25.74527	2020				<10	Excelente					<3	Ex	
22	DLBAJ140	PLAYA JU	PENINSUL	BAJA CALI	LORETO	LORETO	GOLFO	DE	COSTERO	BAHIA	-111.304	25.80818	2020				<10	Excelente					<3	Ex	
23	DLBAJ141	PARQUE N	PENINSUL	BAJA CALI	LORETO	LORETO	GOLFO	DE	COSTERO	BAHIA	-111.34	25.8758	2020				<10	Excelente					<3	Ex	
24	DLBAJ143	BAHIA DE	PENINSUL	BAJA CALI	LORETO	LORETO	GOLFO	DE	COSTERO	BAHIA	-111.331	25.8363	2020					16.8571	Excelente				<3	Ex	
25	DLBAJ144	BAHIA LOF	PENINSUL	BAJA CALI	LORETO	LORETO	GOLFO	DE	COSTERO	BAHIA	-111.335	26.00486	2020					13.8667	Excelente				<3	Ex	

Datos de calidad del agua de si

<

Imagen 1. Dataset elegido, referencia visual

Limpieza

Se corrieron métodos para identificar espacios vacíos dentro de la base de datos. Así como para identificar el tipo de elemento para trabajar. La cantidad de elementos faltantes dentro de los atributos.

Técnicas que se implementaron:

- Se eliminaron columnas y filas que no agregaban valor al análisis
- Se cambiaron formatos para poder hacer las interpretaciones adecuadas

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 3493 entries, 0 to 3492
Data columns (total 55 columns):
#   Column                Non-Null Count  Dtype
---  -
0   CLAVE                  3493 non-null  object
1   SITIO                  3493 non-null  object
2   ORGANISMO_DE_CUENCA   3493 non-null  object
3   ESTADO                 3493 non-null  object
4   MUNICIPIO              3493 non-null  object
5   CUENCA                 3492 non-null  object
6   CUERPO DE AGUA        3479 non-null  object
7   TIPO                   3493 non-null  object
8   SUBTIPO                3479 non-null  object
9   LONGITUD               3493 non-null  float64
10  LATITUD                3493 non-null  float64
11  PERIODO                3493 non-null  int64
```

Imagen 2. Formato de los datos.

```
[ ] #Revisamos si tenemos valores faltantes
df.isnull().values.any()
```

True

```
#El resultado anterior nos marco en efect
#En nuestra primer etapa de validación de
#decidimos primero ubicar las columnas de
df.isna().any()
```

```
CLAVE                False
SITIO                False
ORGANISMO_DE_CUENCA  False
ESTADO              False
MUNICIPIO            False
CUENCA               True
CUERPO DE AGUA       True
TIPO                 False
SUBTIPO              True
LONGITUD             False
-----
```

Imagen 3. Presencia de elementos faltantes

Análisis

Definimos los atributos con los que vamos a trabajar.

1.- Nos limitamos a trabajar con variables numéricas y graficamos su comportamiento.

2 Revisamos que los atributos nichos con los que trabajamos tengan una correlación directa o una correlación indirecta alta. Nos percatamos de una correlación directa alta y nos permite seguir con el análisis.

3. Todo el análisis parece bien realizado, es necesario identificar si tenemos “outliers” en una presencia significativa como para afectar nuestro análisis. En la siguiente grafica vemos que si tenemos pero son minoría como para poder ser despreciados en las interpretaciones, sin la necesidad de una normalización extra.

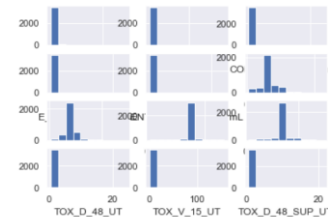


Imagen 4. Elementos nicho y su comportamiento

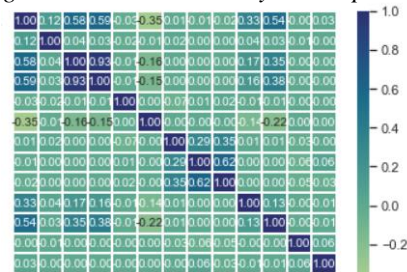


Imagen 5. Confirmación de correlacion

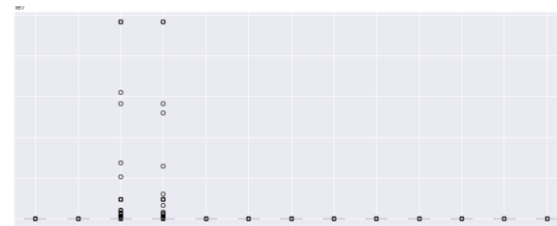


Imagen 6. Mapa de comportamiento.

01

02

03

Clasificación

La clasificación y ubicación de los centroides fue revisada con un método de matriz de confusión y se determinó la precisión de este. Una precisión superior al 99%. Por ende podemos decir que los datos son muy conclusivos en la siguiente imagen.

El método utilizado para mostrar lo siguiente es denominado "K-means".

El método agrupa las formas de agua superficiales en una categoría de semáforos. La verde siento la mejor condición, amarilla poco recomendada y la roja no recomendada.

A su vez podemos ubicar los centroides marcados por puntos azules en la grafica.

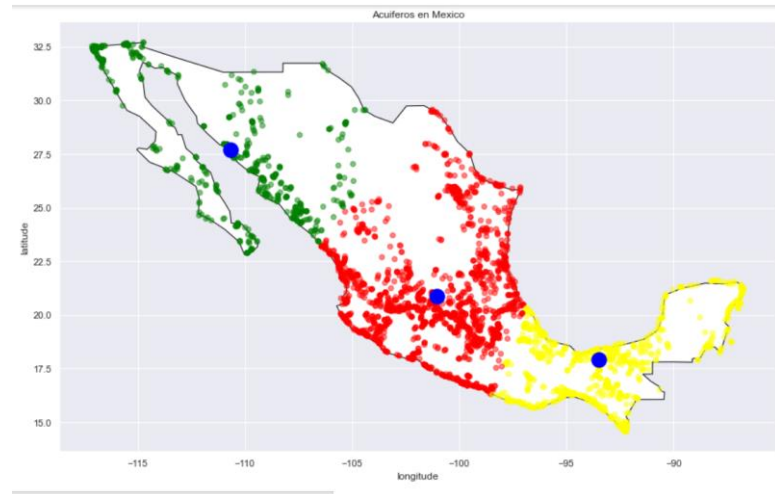


Imagen 7. Mapa de México usando el método "K-means"

01

02

03

04

Resultados

Es necesario analizar las siguientes graficas para tener una mejor perspectiva. Los mantos acuíferos superficiales fueron analizados de manera individual y podemos ver que tras tomar en cuenta presencias de metales pesados, cargaso bacterianas, así como otros aspectos nocivos para la salud humana podemos identificar en las costas de nuestro país una carga notoria de agua “buena” para nuestra humanidad pero si utilizamos un método de ML para hacer este análisis y otro método para corrobórrar su efectividad podemos ver que hay una fuerte tendencia en la zona nor-oeste del país a una calidad de agua superior, denotando que la zona centro y zona nor-este un “no-recomendado” y finalmente en la zona sur y la península una presencia de agua “poco recomendada”.

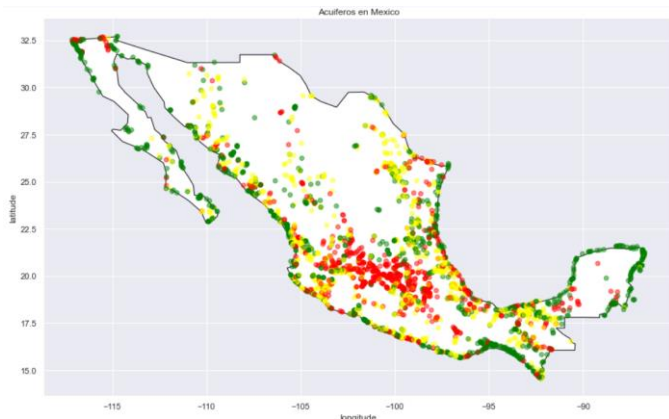


Imagen 8. Mapa de México de la calidad de los mantos acuíferos

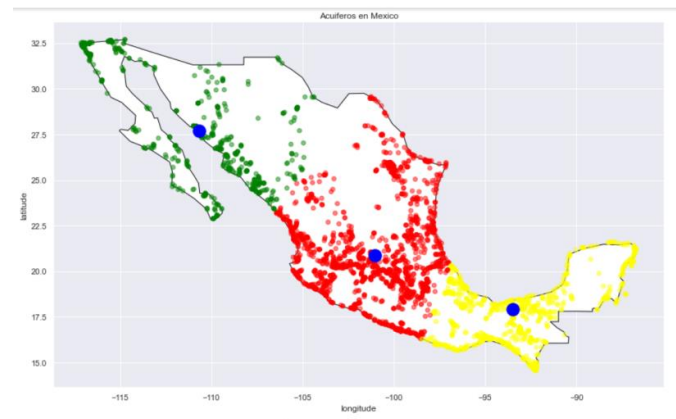


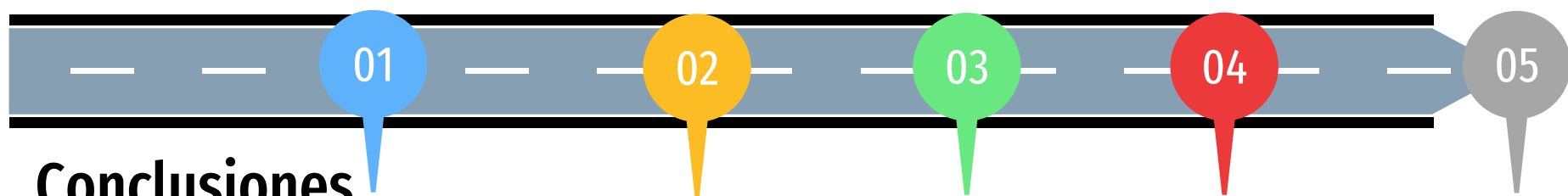
Imagen 9. Mapa de México usando el método "K-means"



Conclusiones

Las conclusiones de nuestro muestreo son las siguientes:

- 1.- La zona norte-oeste tiene mayor calidad de agua, centro y nor-este un “no recomendada” y la zona sur y península “poco recomendada”.
- 2.- Este ejercicio no permite acciones conclusivas, un análisis de ciencia de datos debe estar basado en problemáticas del negocio, la única indicación adecuada que se puede entregar es que puede servir para las estrategias de organismos gubernamentales en sus planes de acción.
- 3.- El ejercicio y los métodos utilizados son de gran poder y permiten una perspectiva global de la geografía del país en temas de mantos acuíferos superficiales, la comparación de estos hallazgos con las cuencas de agua entregaría mayor valor.



Conclusiones

Las conclusiones de nuestro muestreo son las siguientes:

4.- El análisis implementado requirió de la necesidad de borrar una gran cantidad de atributos debido al manejo del banco de datos, una propuesta directa del equipo de ciencia de datos es que el equipo de recolección de información no tiene un “modern data warehouse”. Estos análisis pueden entregar mayor valor al tener la infraestructura técnica adecuada.

5.- Es necesario hacer estos análisis comparativos a lo largo de 20 o 30 años para ver una evolución en los tratamientos y acciones de la sociedad junto con gobierno, o en su defecto ver los efectos nocivos.

Proyecto

Actividad de cierre de la materia de Ciencia de datos
Maestría en Inteligencia artificial aplicada
ITESM