

Proyecto final HIGI: base de datos del INE

Daniela Jiménez Pano
NUA:427292

17 de junio de 2021

Resumen

En el presente documento se reportará el desarrollo y resultados del proyecto final para la asignatura de herramientas informáticas y gestión de la información, el cual solicitó el análisis estadístico de dos bases de datos, una que encontráramos los alumnos por cuenta propia y otra que nos indicara la profesora. Estos estudios fueron realizados en python, herramienta nueva que aprendimos el último parcial del semestre. Sin embargo, este documento mostrará solamente el análisis de la base de datos que indicó la maestra de la asignatura.

1. Introducción

Para el proyecto final de la materia mencionada, se pidió analizar dos bases de datos, los primeros datos obtenidos de la página del Instituto Nacional Electoral, estos llamados: Estadística de Padrón Electoral y Lista Nominal de Electores [1]. Un análisis estadístico se basa en recopilar, explorar y presentar cantidades de datos en los que se descubren ciertos patrones y tendencias. [2] El objetivo del análisis de este conjunto de datos fue obtener la predicción del número de personas dentro de la lista nominal 2 meses después de los últimos datos obtenidos, además de encontrar cómo afectaba el incremento (tasa de crecimiento/variación) de este número de personas para cada una de las secciones de los municipios del estado de Guanajuato.

2. Desarrollo

2.1. Primer análisis: Predicción del número de casillas y lista nominal.

Para esta actividad se llevaron a cabo ciertos pasos, lo primero a realizar fue descargar el conjunto de datos desde el mes de septiembre de 2019 hasta diciembre de 2020. Este análisis ya había sido realizado con anterioridad con ayuda de la herramienta Excel, sin embargo, la manera de trabajar en python no es tan cómodo trabajar debido a que sus comandos no son tan accesibles como en excel. Como los documentos venían organizados por estado además de involucrar los municipios y secciones de estos, se tenían que filtrar cada uno de los documentos para que pudieran organizarse solamente por el estado de nuestro interés, es decir, Guanajuato, para obtener la información específica de este:

```
1 data3=pd.read_csv(files[2],usecols=[0,1,2,3,9])
2 data3_GTO=data3[data3['ENTIDAD']==11]
3
```

Donde en cada uno de los datos con la terminación ".GTO" eran los datos a utilizar para llevar a cabo el análisis anteriormente mencionado, en cada uno de estos datos se eligieron las columnas que nos brindaran la información específica a utilizar, en este caso, la última columna definía el total de la lista nominal de cada municipio. Además, como se muestra en el código, la palabra 'ENTIDAD' fue de suma importancia para poder filtrar los datos y el número 11 especificaba a Guanajuato.

Y cada uno de los documentos fueron agrupados por municipio, de manera que nos ayudara a obtener un mejor acomodo de los datos, ya que también la asignación así lo pedía, con eso se utilizó lo siguiente en donde al final se imprimía cada dataframe para observar el resultado:

```

1  LISTA_NAL_MPO=data1_GTO[1:].groupby(['MUNICIPIO']).sum()
2  LISTA_NAL_MPO
3

```

Al terminar de filtrar cada documento, lo siguiente a realizar fue agrupar los documentos en orden debido a que los nombres de los documentos, como eran nombrados numericamente, no se acomodaban del menor al mayor, por lo tanto, fue necesario ejecutar un código que nos permitiera llevar a cabo eso y desplegarlos los documentos en el acomodados en el orden mencionado:

```

1  date=[]
2  date_=[]
3  files_=[]
4
5  for i,file in enumerate(files):
6      date.append(re.findall(r'\d+',file)[0])
7
8
9  temp=sorted(range(len(date)), key=date.__getitem__)
10
11 for i in temp:
12     date_.append(date[i])
13     print(date[i],files[i])
14     files_.append(files[i])
15

```

A continuación, se juntaron todos los dataframes creados a partir de cada uno de los documentos en una tabla para poder observar de una manera más clara cómo es que había un cambio en el aumento o decremento de cada uno de los meses con relación a la lista nominal y así ayudarnos a poder realizar la predicción deseada, es decir, cuántas personas estarían en la lista nominal para 2 meses después, en este caso para febrero de 2021.

Figura 1: Tabla (dataframe) donde se acomodaron los datos de todos los documentos de la base de datos

Además, con ayuda de la librería matplotlib. se añadió una gráfica para poder ver el comportamiento de ciertos municipios con respecto a la lista nominal, en este caso se pudo notar que los datos reflejaban un comportamiento lineal, por lo que sí era factible realizar un análisis correcto por mínimos cuadrados.

Ahora bien, para poder hacer la regresión lineal que se ha mencionado a lo largo del documento, se tuvo que convertir el dataframe en un array, esto para poder aplicar un ciclo for que permitiera

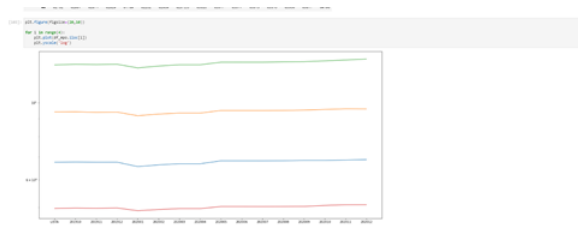


Figura 2: Graficación del comportamiento de la lista nominal por cada municipio.

al programa poder comprender qué datos estaba tomando en cuenta para poder hacer el análisis estadístico necesario:

```
1 municipios=np.asarray(df_mpo)
```

Cabe mencionar que la variable *dfmpo* muestra relación al dataframe creado a partir de los datos de los municipios. Teniendo esto, como el proyecto nos solicitaba observar el crecimiento (tasa de crecimiento/variación) de las personas en la lista nominal por cada una de las secciones de cada uno de los municipios de Guanajuato. Por lo que se decidió obtener la regresión utilizando un código que nos permitiera obtener la predicción en un año después por cada uno de los 46 municipios:

```
1 fits=[]
2 prediction_lnal=[]
3
4 for i in range(len(municipios)):
5     xx=np.arange(len(municipios[i]))
6     ma, ba = np. polyfit(xx, municipios[i],1,w=municipios[i])
7     fits.append([ma,ba])
8     pred=ma*(xx[-12]+12)+ba
9
10    prediction_lnal.append(pred)
```

Donde los datos que muestran la variable municipios hacen una referencia (un poco obvia) a los datos de los municipios de cada documento. Después de haber ingresado el código anterior, se imprimió la predicción como una columna más dentro del dataframe que contenía la información de todos los documentos, permitiéndonos así visualizar la cantidad de personas en la lista nominal predicha.

	LISTA	201610	201611	201612	201701	201702	201703	201704	201705	201706	201707	201708	201709	201710	201711	201712	Predicción UNIA ASES DESPUE
MUNICIPIO																	
1	67427	67441	67463	67463	68008	68770	68750.0	68837	68934	68952	68993	69204	69341	69422	69511	69533	6.95333e+04
2	34940	34939	35025	35070	35172	35270	35340.0	35407	35492	35559	35639	35705	35765	35814	35854	35894	5.85934e+04
3	12446	12557	12771	12979	13209	13461	13634.0	13824	14032	14257	14500	14752	15009	15269	15540	15809	1.52690e+05
4	49054	49197	49359	49519	49687	49851	49984.0	50150	50320	50497	50684	50867	51057	51250	51447	51648	5.16477e+04
5	67669	68120	67962	68219	68709	67953	68952.0	69125	69124	69114	69152	69207	69245	69272	69294	69314	7.01217e+04
6	4207	4225	4235	4255	4245	4212	4212.0	4244	4244	4241	4235	4245	4252	4261	4269	4281	4.29849e+03
7	21231	21247	21239	21262	21278	21291	21322.0	21357	21392	21425	21459	21494	21528	21564	21599	21634	2.17159e+05
8	31451	31444	31501	31453	31452	31466	31485.0	31490	31797	31797	31816	31816	32045	32045	32045	32045	3.22047e+04
9	60267	60364	60237	60255	60750	59944	59917.0	60051	60251	60269	60267	61071	61456	61479	61515	61555	6.15555e+04
10	5353	5375	5395	5415	5432	5456	5479.0	5496	5519	5542	5565	5588	5610	5632	5654	5676	5.62675e+03
11	14226	14274	14287	14289	14279	14271	14222.0	14210	14206	14204	14201	14201	14201	14201	14201	14201	1.42010e+04
12	23212	23248	23268	23268	23268	23268	23268.0	23268	23268	23268	23268	23268	23268	23268	23268	23268	2.32679e+04
13	15343	15346	15353	15361	15369	15369	15369.0	15353	15353	15346	15379	15389	15394	15394	15394	15394	1.53940e+04
14	113151	113190	113175	113176	113162	113167	113161.0	113160	113167	113164	113162	113163	113161	113161	113161	113161	1.13161e+05
15	11006	110674	110002	110004	110001	110119	110119.0	110119	110119	110119	110119	110119	110119	110119	110119	110119	1.10119e+05
16	10045	10085	10062	10075	10079	10087	10086.0	10086	10086	10086	10086	10086	10086	10086	10086	10086	1.00859e+04
17	423390	423391	423390	423390	423390	423390	423390.0	423390	423390	423390	423390	423390	423390	423390	423390	423390	4.23390e+05
18	28558	28551	28520	28547	28558	28570	28555.0	28555	28555	28555	28555	28555	28555	28555	28555	28555	2.85555e+04
19	41462	41462	41470	41473	41482	41487	41484.0	41484	41484	41484	41484	41484	41484	41484	41484	41484	4.14840e+04
20	115582	115584	115584	115584	115584	115584	115584.0	115582	115582	115582	115582	115582	115582	115582	115582	115582	1.15582e+05
21	43380	43405	43403	43403	43403	43404	43404.0	43404	43404	43404	43404	43404	43404	43404	43404	43404	4.34040e+04
22	18031	18030	18025	18025	18029	18030	18030.0	18030	18030	18030	18030	18030	18030	18030	18030	18030	1.80300e+04
23	124850	124840	124825	124835	124835	124835	124835.0	124835	124835	124835	124835	124835	124835	124835	124835	124835	1.24835e+05
24	10011	10011	10011	10011	10011	10011	10011.0	10011	10011	10011	10011	10011	10011	10011	10011	10011	1.00110e+04
25	54957	54954	54955	54955	54955	54954	54954.0	54954	54954	54954	54954	54954	54954	54954	54954	54954	5.49540e+04

Figura 3: Predicción de la lista nominal para febrero de 2021

Finalmente para obtener el número de casillas a necesitar para el mes de febrero del año 2021 se sumaron cada una de las predicciones divididas en 750 (que es el número mínimo para abrir una casilla) y con eso finalizar el análisis estadístico de esta base de datos.

3. Resultados y discusiones

Con los datos utilizados, se pudo llegar a que para el mes de febrero de 2021, es decir, unos cuantos meses después del último documento analizado, la cantidad de casillas a necesitar fue de 6038 casillas. Para obtener una vista más detallada de estos números, se adjunta la siguiente tabla:

Municipio	Número de casillas	Municipio	Número de casillas
1	91	14	155
2	128	15	191
3	177	16	23
4	68	17	579
5	94	18	40
6	6	19	60
7	516	20	1561
8	43	21	59
9	82	22	25
10	13	23	170
11	101	24	16
12	32	25	76
13	25	26	63

Cuadro 1: Predicción a Febrero de 2021: Primeros 23 municipios de Guanajuato.

Municipio	Número de casillas	Municipio	Número de casillas
27	287	40	18
28	113	41	67
29	38	42	156
30	112	43	21
31	127	44	64
32	83	45	11
33	123	46	85
34	6		
35	86		
36	9		
37	188		
38	14		
39	42		

Cuadro 2: Predicción a Febrero de 2021: Sigüientes 23 municipios de Guanajuato.

Ahora, haciendo una comparación a la misma asignación que se realizó con la herramienta de excel, las casillas que nos arrojó por municipio fueron de 6829, con esto, se puede observar que la diferencia entre los resultados no es tan diferente y podremos aceptar la cantidad que python nos brindó.

Desafortunadamente para poder predecir y realizar una regresión lineal al apartado de secciones, debido a que se intentó aplicar un código que se utilizó anteriormente en el análisis, sin embargo, marcaba un error en el que la librería matplotlib. no permitía hacer el ajuste de mínimos cuadrados con los datos que se indicaron, sin embargo, se asume a que esto sucedió debido a que los datos de las secciones eran demasiados (el código fue marcado como comentario dentro de mi análisis).

4. Conclusión

Gracias a esta actividad, en lo personal pude observar qué tan importante es realizar este tipo de análisis ya que se pueden hacer proyecciones futuras acerca de datos tan importantes como lo son las listas nominales del Instituto Nacional Electoral. Además de que las herramientas que nos brinda python para este tipo de proyectos serán de mucha ayuda para proyectos futuros además de que el trabajo en este programa se hace de una manera más rápida.

Con respecto a los resultados obtenidos, en un principio me causó conflicto el hecho de que los valores de las casillas para los municipios fueran diferentes e al análisis realizado en excel, sin embargo, recordando el trabajo que hice, me percaté acerca de que se usó una función en la que se redondeaban los valores de cualquier manera, por lo que provocó que el aumento de casillas fuera mayor a comparación

al presente proyecto. A pesar de esto, quedo satisfecha con los resultados ya que estos resultan bastante razonables.

Referencias

- [1] Griselda Rojas Peña. Estadística de padrón electoral y lista nominal de electores. url<https://www.ine.mx/transparencia/datos-abiertos//archivo/estadistica-padron-electoral-lista-nominal-electores>, 2020.
- [2] TechTarget. Análisis estadístico. url<https://searchdatacenter.techtarget.com/es/definicion/Analisis-estadistico>, s.f.