



UNIVERSIDAD  
POLITÉCNICA  
DE YUCATÁN



# CLASSIFICATION FINAL REPORT

Data Mining

**PROFESSOR:**

Arturo Raymundo

**STUDENTS:**

Fatima Martínez Torres

Morales Zapata

Mario

**DATE:**

November 16th, 2020

**GROUP:**

Data 6A

---

# FINAL REPORT

---

## Classification

To begging with this report, it is essential for us to describe a little bit about the process needed in order to extract our data and the mental process we passed to decide on one type of classification.

### Dataset

First of all, we must say we decided to work on the Instituto Nacional de Estadística y Geografía (INEGI)'s data on Households and Housing Units, in its subsection called Household Income and Expenditure, because we thought we could extract more information and we discussed the fact that there are hundreds of thousands of houses in the Mexican territory, so we were expecting to have a good amount of information.

Downloading the information was not precisely intuitive.

From the INEGI's Home Page, we tried to download the Tabular Data. Nonetheless, there was no intuitive way to do so. Therefore, we explored the Related Programs tab which leaded us to the National Survey of Household Income and Expenditure (Encuesta Nacional de Ingresos y Gastos de los Hogares ENIGH). This survey has the purpose of providing a statistical overview of the behavior of household income and expenditure in terms of its amount, source and distribution; plus, the information it gives on the occupational and socio-demographic characteristics of the household members, as same as the infrastructure and equipment.

Once we were in this page, we used the Open Data option to download the latest data from 2018 in a CSV format. On the zip file downloaded, we ended up with several folders, each of which had inside (at least) five sub-folders with the needed information to get to know the data stored in those CSV files.

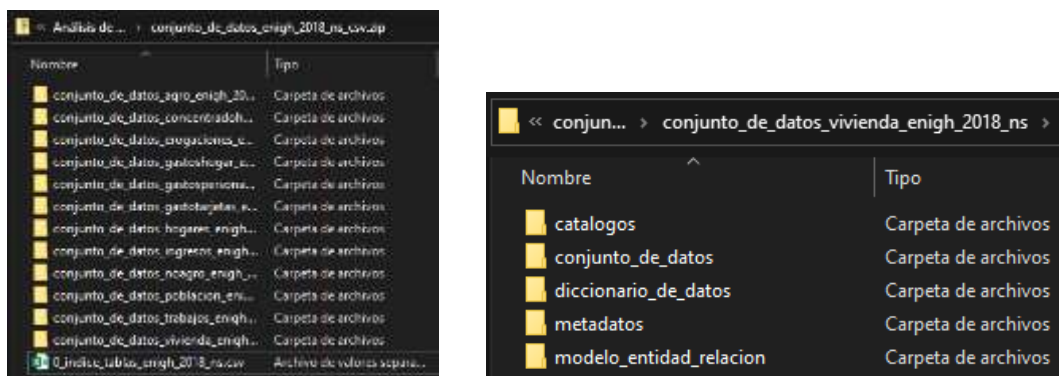


Figure 1 Overview of the data

(You can see what the folders contain in the [link](#))

Even though all of the files in there were house-related, we decided to work just on the one that was called *conjunto\_de\_datos\_viviendas\_enigh\_2018\_ns* which gives us the characteristics of the properties by the household members. That one only focuses on the physical characteristics of the house (type of floor, material, number of rooms, etc.).

That dataset contained more than fifty columns and more than seventy thousand rows, so it was a good one to work with. One of those fifty columns had the information of, so to speak, the socioeconomic status of the house with four different results:

- Bajo
- Medio bajo
- Medio alto
- Alto

## Methodology

We decided then to train a classification algorithm to identify (and classify of course) a house given its characteristics into one of these categories.

Now, following the idea of we having to classify an input into a category, a logistic regression was our best option; nonetheless, the output is not binary but polychotomous, so following that constrain, the most accurate model we could possibly use, is a multivariable logistic regression.

## Data Preprocessing

We started our data preprocessing by analyzing and looking for NaN and Null values, but we were not able to find any; nonetheless that did not mean that everything in the dataset was worth the shoot of using it into our model. Keeping that in mind, we dropped a lot of columns because the content of those columns did not make sense if you are trying to find the house's socioeconomic tag (let us remind you that this category does not refer to the people living in the house but to the physical characteristics of it, so for example, we did not care if that house has to pay a million dollars on electricity, because even the poorest house can consume that same energy if needed).

#	Column	Non-Null Count	Non-Null	Count	Type
0	Bajo	73405	non-null	int64	
1	Medio bajo	73405	non-null	int64	
2	Medio alto	73405	non-null	int64	
3	Alto	73405	non-null	int64	
4	...	...	...	...	...
44	...	73405	non-null	int64	
45	...	73405	non-null	int64	
46	...	73405	non-null	int64	
47	...	73405	non-null	int64	
48	...	73405	non-null	int64	
49	...	73405	non-null	int64	
50	...	73405	non-null	int64	
51	...	73405	non-null	int64	
52	...	73405	non-null	int64	
53	...	73405	non-null	int64	
54	...	73405	non-null	int64	
55	...	73405	non-null	int64	
56	...	73405	non-null	int64	
57	...	73405	non-null	int64	
58	...	73405	non-null	int64	
59	...	73405	non-null	int64	
60	...	73405	non-null	int64	
61	...	73405	non-null	int64	
62	...	73405	non-null	int64	
63	...	73405	non-null	int64	
64	...	73405	non-null	int64	
65	...	73405	non-null	int64	
66	...	73405	non-null	int64	
67	...	73405	non-null	int64	
68	...	73405	non-null	int64	
69	...	73405	non-null	int64	
70	...	73405	non-null	int64	
71	...	73405	non-null	int64	
72	...	73405	non-null	int64	
73	...	73405	non-null	int64	
74	...	73405	non-null	int64	
75	...	73405	non-null	int64	
76	...	73405	non-null	int64	
77	...	73405	non-null	int64	
78	...	73405	non-null	int64	
79	...	73405	non-null	int64	
80	...	73405	non-null	int64	
81	...	73405	non-null	int64	
82	...	73405	non-null	int64	
83	...	73405	non-null	int64	
84	...	73405	non-null	int64	
85	...	73405	non-null	int64	
86	...	73405	non-null	int64	
87	...	73405	non-null	int64	
88	...	73405	non-null	int64	
89	...	73405	non-null	int64	
90	...	73405	non-null	int64	
91	...	73405	non-null	int64	
92	...	73405	non-null	int64	
93	...	73405	non-null	int64	
94	...	73405	non-null	int64	
95	...	73405	non-null	int64	
96	...	73405	non-null	int64	
97	...	73405	non-null	int64	
98	...	73405	non-null	int64	
99	...	73405	non-null	int64	

Figure 2 Non-Null values

Once we stayed with just the columns that did make sense for our classification goal, we noticed that all the data was represented numerically, but everything was explained in the data dictionary provided by the INEGI. The thing with this fact, is that they were supposed to represent the information through numbers (check the dictionary of the data [here](#)), so most of the columns were type integer (int64), but some others were type object.

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 73405 entries, 0 to 73404
Data columns (total 43 columns):
 #   column              non-null count  dtype
0   mat_pared           73405 non-null  int64
1   mat_techos           73405 non-null  int64
2   mat_pisos           73405 non-null  object
3   cocina              73405 non-null  int64
4   cocina_dor          73405 non-null  object
5   cuart_dorm           73405 non-null  int64
6   num_cuarto           73405 non-null  int64
7   disp_agua           73405 non-null  int64
8   dotat_agua          73405 non-null  object
9   excusado            73405 non-null  int64
10  uso_compan           73405 non-null  object
11  aanti_agua          73405 non-null  object
12  boro_comp           73405 non-null  object
13  boro_excus          73405 non-null  object
14  boro_reged          73405 non-null  object
15  drenaje             73405 non-null  int64
16  disp_elect          73405 non-null  int64
17  focos_inca          73405 non-null  object
18  focos_ahor          73405 non-null  object
19  combustible          73405 non-null  object
20  estufa_chi          73405 non-null  object
21  eli_basura          73405 non-null  int64
22  renta               73405 non-null  object
23  estim_pago          73405 non-null  object
24  pago_viv            73405 non-null  object
25  viv_usada           73405 non-null  object
26  lavadero            73405 non-null  int64
27  fregadero           73405 non-null  int64
28  regadera            73405 non-null  int64
29  tinaco_azo          73405 non-null  int64
30  cisterna            73405 non-null  int64
31  pileta              73405 non-null  int64
32  calent_sol          73405 non-null  int64
33  calent_gas          73405 non-null  int64
34  medidor_luz         73405 non-null  int64
35  bomba_agua          73405 non-null  int64
36  tanque_gas          73405 non-null  int64
37  aire_acond          73405 non-null  int64
38  calefacc            73405 non-null  int64
39  tot_resid           73405 non-null  int64
40  tot_hog             73405 non-null  int64
41  tam_loc             73405 non-null  int64
42  est_socio           73405 non-null  int64
dtypes: int64(27), object(16)
memory usage: 24.1+ MB
```

Figure 3 Cleaned

Further analysis on those columns proved them to have errors in the information saved in there (for example, they were supposed to be completely numeric, but some of them also had "&" or empty spaces), giving us the reason why the data type were of the object type. Unfortunately, all rows had at least one mistake in some column, so instead of dropping all of our dataset, we decided to drop the columns that had at least one mistake. This could sound extreme, but let us consider the fact that the treatment required to replacing this errors must be exhaustive, because we would change at least one characteristic on the row in all the rows in our data... As this task is mainly focused on the classification process, we decided not to go deep into that part and just remaining with the columns that were saved correctly.

```
1 # Eliminamos todas las columnas que sean objetos
2 df.drop(objeto, axis=1,inplace=True)
3 df.head()
```

	mat_pared	mat_techos	cocina	cuart_dorm	num_cuarto	disp_agua	excusado	drenaje	disp_elect
0	8	10	1	4	6	1	1	1	
1	8	10	1	5	6	1	1	1	
2	8	10	1	2	7	1	1	1	
3	8	10	1	2	5	1	1	1	
4	8	10	1	3	4	1	1	1	

5 rows x 27 columns

Figure 4 Table with no object-datatype

## Classification system

After all the cleaning, we decided to have a first approach to our classification in order to see if all the data were ready to be classified, or at least strong enough.

Hence, we start with our training data for the Logistic Regression using as parameter to compare the *est\_socio* (estatus socioeconómico) column, and since the values were given by integers from 1 to 4, those were the data compared on our chosen Confusion Matrix, so that we can visualize the precision of the predicted data.

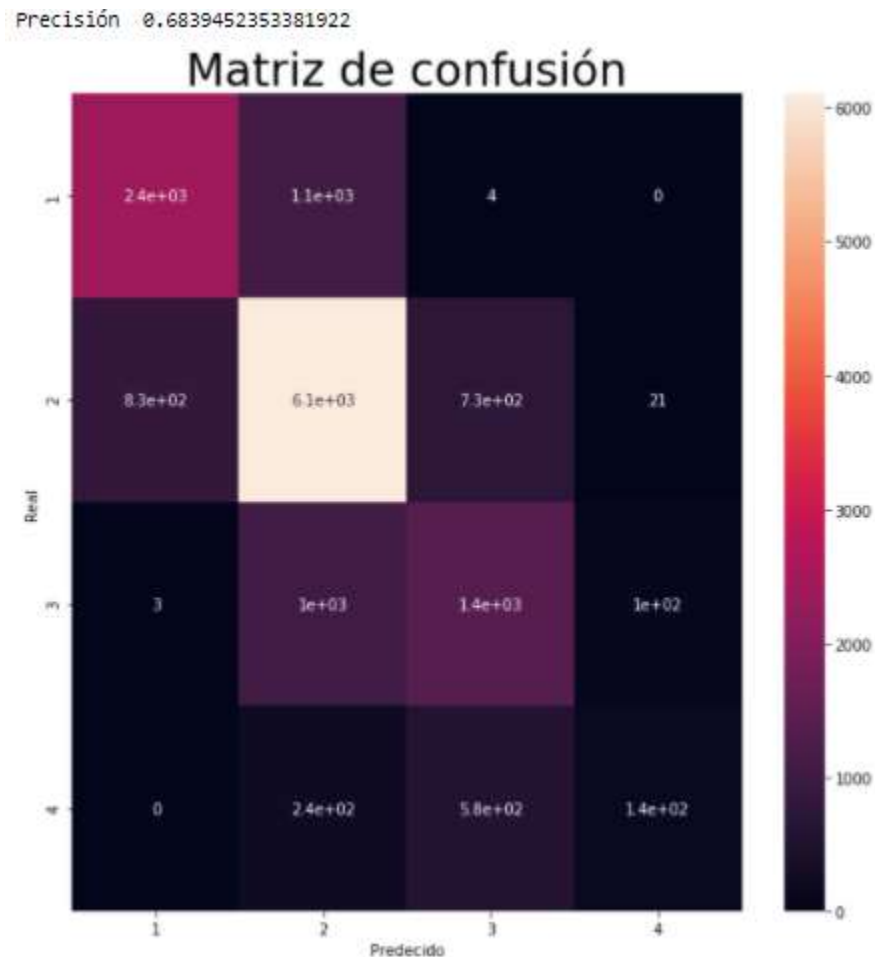


Figure 5 Confusion Matrix

We could notice, then, that the accuracy on our predicted data was not enough since there were thousands of data that showed so.

It is obvious that our model gives very poor predictions; We could not really trust it, therefore, we had to implement some modifications in our analysis in order to improve the accuracy of the model.

Most of the data found in the column represents data of nominal type, for example, in the column *viv\_type* represents the type of house on which the information was captured; the possible values within this column are the numeric values from 1 to 5 where:

1. Independent house
2. Apartment in building
3. Neighborhood housing
4. Living in the roof room
5. Local not built for room

Once knowing this information, it is useless for us to think of a logical order for this column just as suggested by its numerical values, we cannot think that  $1 < 2$  as it simply does not make sense and is a potential threat to our algorithm to apply later.

In that way, we decided to replace the nominal values on some columns, which might have show us important information on our analysis, into their real meaning. So that we could use the data as other type but numerical, since this type does not give a quantity but information.

However, since we got to the idea of working with the table as dummy values, we found out that those columns that shows Boolean data will be unnecessary to change their values of 1 or 2, so we kept those columns as we originally have it.

```
1 # Creo una lista de todas las columnas nominales. Por cuestiones de lo forms en que el INEGI guarde los datos, también
2 # tengo que incluir en la lista todos los nombres de las columnas que representen información binaria (sí, no)
3 nominales = ['num_cuarto', 'cuart_dorm', 'bano_com', 'bano_excus', 'bano_regad', 'focos_inca', 'focos_ahor',
4             'renta', 'estiu_pago', 'pago_viv', 'tot_resid', 'tot_hog', 'cocina', 'cocina_dor', 'excusado',
5             'uso_compen', 'estufa_chi', 'lavadero', 'fregadero', 'regadera', 'tineco_azo', 'cisterna', 'pileta',
6             'calent_sol', 'calent_gas', 'medidor_luz', 'bomba_agua', 'tanque_gas', 'aire_acond', 'calefact', 'viv_usada']
7
8 # creo una lista con todas las columnas del dataframe
9 nominales = df.columns.tolist()
10
11 # elimino todos los nombres de columnas nominales de la lista total de columnas
12 for element in nominales:
13     if element in nominales:
14         nominales.remove(element)
```

Figure 6 Omitting Boolean Columns

```
1 # Reemplazo todas las valores nominales por su significado real
2 for i in range(len(nominales)):
3     csv = 'Data\\'+nominales[i]+'.csv'
4     # Acceso al dataframe que contiene el diccionario de los datos de la respectiva columna
5     temp = pd.read_csv(csv, encoding = 'latin')
6     for index, row in temp.iterrows():
7         # Leo la combinación de valores de cada diccionario de datos
8         numero = temp[nominales[i]][index]
9         valor = temp['descripcion'][index]
10        df[nominales[i]].replace(numero, valor, inplace=True)
```

```
1 # Dataframe a trabajar
2 df.head(2)
```

	mat_pared	mat_techo	cocina	cuart_dorm	num_cuarto	disp_agua	excusado	drenaje	disp_elect	eli_basura	...	calent_gas	medidor_luz	bomba_agua
0	Tuboque, ladrillo, block, piedra, cantera, ceram...	Losa de concreto o viguetas con bovedilla	1	4	0	Agua entubada dentro de la vivienda	1	La red pública	Del servicio público	La tiran en un contenedor o depósito		1	1	1
1	Tuboque, ladrillo, block, piedra, cantera, ceram...	Losa de concreto o viguetas con bovedilla	1	5	0	Agua entubada dentro de la vivienda	1	La red pública	Del servicio público	La tiran en un contenedor o depósito		1	1	1

2 rows x 27 columns

Figure 7 Replacing nominal values

We finally got the data frame as needed to convert our nominal columns to a set of dummy columns, which, indeed, gave us a higher precision on our Confusion Matrix.

1	# Creo un dataframe con todas las características dummy
2	dummy = pd.get_dummies(df[nominales[:-1]])
3	#dummy_1 = pd.get_dummies(df[si_no])
4	# Elimino las columnas "originales" exceptuando columna de resultados. La razón es que crearía redundancia con los dummies
5	df.drop(nominales[:-1],axis=1,inplace=True)
6	#df.drop(si_no,axis=1,inplace=True)
7	# Hago un merge con mi dataframe original
8	df = df.merge(dummy,left_index=True,right_index=True)
9	#df = df.merge(dummy_1,left_index=True,right_index=True)

1	# Dataframe listo para utilizarse
2	df.head()

	cocina	cuart_dorm	num_cuarto	excusado	lavadero	fregadero	regadera	tinaco_azo	cisterna	pileta	...	eli_basura_La recoge un camión o carrito de basura	eli_basura_La tiran al río, lago o mar	eli_basura_La tiran en barranca o grieta	eli
0	1	4	6	1	1	1	1	1	1	2	...	0	0	0	0
1	1	5	6	1	1	1	1	1	1	2	...	0	0	0	0
2	1	2	7	1	1	1	1	1	1	2	...	0	0	0	0
3	1	2	5	1	1	1	1	1	1	2	...	0	0	0	0
4	1	3	4	1	1	1	1	1	1	1	...	0	0	0	0

5 rows x 67 columns

Figure 8 Data Frame as Dummy columns

Precisión 0.7006802721088435

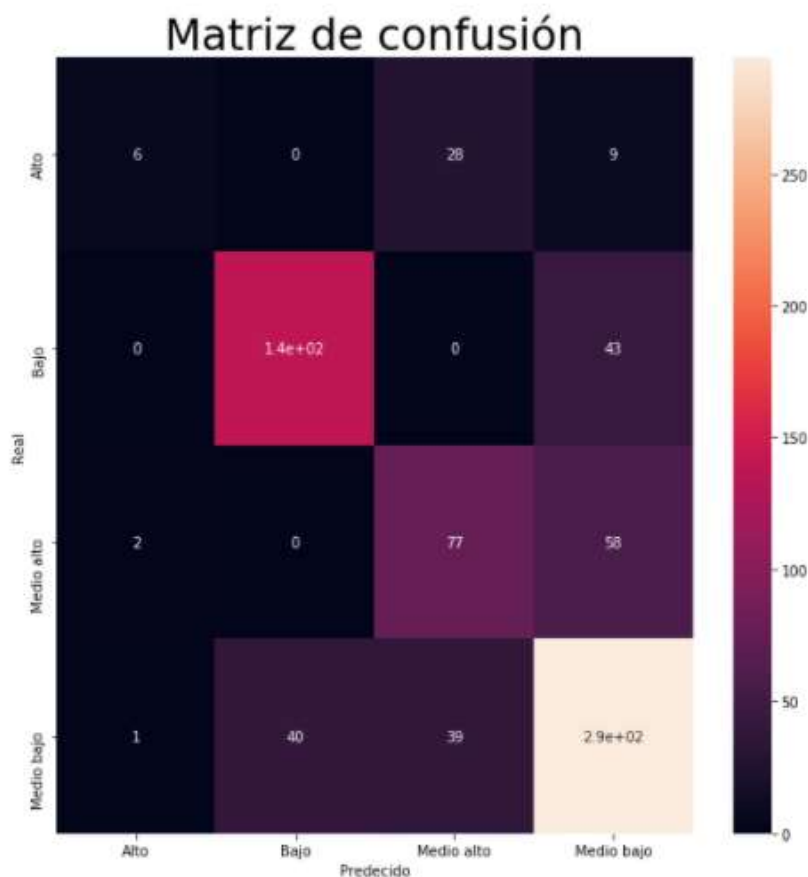


Figure 9 More accurate Confusion Matrix

## Conclusion

As a conclusion, we would like to focus on the fact that the accuracy of our model is not even close of being reliable; nonetheless we could indeed find some improvement between our first approach and the second one with the categorical data expressed as it supposed to.

Another thing to take into consideration is the fact that we kept less than half of the original number of columns, so even though our model does not perform the best way it could, it is not that bad given the amount of information we lost. With further treatment in the columns that were not saved in a right way, it is very likely that our model's performance could be really good.



## References

- INEGI. (2018). *Encuesta Nacional de Ingresos y Gastos de los Hogares (ENIGH). 2018 Nueva serie*. INEI. <https://www.inegi.org.mx/programas/enigh/nc/2018/#Documentacion>
- Martínez Torres, F., & Morales Zapata, M. A. (2020). *DataMining*. GitHub. <https://github.com/fatimamt/Data-Mining>