

NAME-PAPU SWAIN

PROJECT-02(Income Qualification in Latin America)

OBJECTIVE-Identify the level of income qualification needed for the families in Latin America.

```
import seaborn as sns
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")
```

In [2]:

```
train=pd.read_csv("train.csv",index_col=0)
test=pd.read_csv("test.csv",index_col=0)
```

In [3]:

```
train.shape, test.shape
```

Out[3]:

```
((9557, 142), (23856, 141))
```

In [4]:

```
train.head()
```

Out[4]:

	id	v2a1	hacdor	rooms	hacapo	v14a	refrig	v18q	v18q1	r4h1	r4h2	...	SQBescolari	SQBage	SQBhogar_total	SQBdeje	SQBhogar_nin	SQBovercrowding	SQBdependency	SQB
ID_279628684	190000.0	0	3	0	1	1	0	NaN	0	1	...	100	1849	1	100	0	1.000000		0.0	
ID_1c7f8646d	135000.0	0	4	0	1	1	1	1.0	0	1	...	144	4489	1	144	0	1.000000		64.0	
ID_68de5c1c4	NaN	0	8	0	1	1	0	NaN	0	0	...	121	8464	1	0	0	0.250000		64.0	
ID_d671d8b9c	180000.0	0	5	0	1	1	1	1.0	0	2	...	81	289	16	121	4	1.777778		1.0	
ID_0561d69b5	180000.0	0	5	0	1	1	1	1.0	0	2	...	121	1369	16	121	4	1.777778		1.0	

5 rows × 142 columns

In [5]:

```
test.head()
```

Out[5]:

	id	v2a1	hacdor	rooms	hacapo	v14a	refrig	v18q	v18q1	r4h1	r4h2	...	age	SQBescolari	SQBage	SQBhogar_total	SQBdeje	SQBhogar_nin	SQBovercrowding	SQBdependency	SQB
ID_216873615	NaN	0	5	0	1	1	0	NaN	1	1	...	4		0	16	9	0	1	2.25		0.2
ID_1c7f8846d	NaN	0	5	0	1	1	0	NaN	1	1	...	41		256	1681	9	0	1	2.25		0.2
ID_05442cfa	NaN	0	5	0	1	1	0	NaN	1	1	...	41		289	1681	9	0	1	2.25		0.2

ID_2f6873615	NaN	0	5	0	1	1	0	NaN	1	1	...	4	0	16	9	0	1	2.25	0.2	
ID_1c7f8646d	NaN	0	5	0	1	1	0	NaN	1	1	...	41	256	1681	9	0	1	2.25	0.6	
ID_a5442c2fa	NaN	0	5	0	1	1	0	NaN	1	1	...	41	289	1681	9	0	1	2.25	0.2	
ID_8d4326f7a	NaN	0	14	0	1	1	1	1.0	0	1	...	59	256	3481	1	256	0	1.00	0.6	
ID_862966799	175000.0	0	4	0	1	1	1	1.0	0	0	...	18	121	324	1	0	1	0.25	64.0	

5 rows × 141 columns

From the above,it's clear that Target column is the output variable.

```
In [6]: train.info()

<class 'pandas.core.frame.DataFrame'>
Index: 9557 entries, ID_279628684 to ID_a38c64491
Columns: 142 entries, v2a1 to Target
dtypes: float64(8), int64(130), object(4)
memory usage: 10.4+ MB

In [7]: test.info()

<class 'pandas.core.frame.DataFrame'>
Index: 23856 entries, v2_2f6873615 to ID_34754556f
Columns: 141 entries, v2a1 to agepo
dtypes: float64(8), int64(129), object(4)
memory usage: 25.8+ MB
```

Understand the type of data

```
In [50]: train.describe(include=[object])
```

From the above,it's clear that Target column is the output variable.

	unique	2988	31	22	22
	top	1d8a6d014	yes	no	no
	freq	13	2192	3762	6230

In [9]:

```
test.describe(include="object")
```

Out[9]:

	idhogar	dependency	edjefe	edjefa
count	23856	23856	23856	23856
unique	7352	35	22	22

Understand the type of data

```
In [10]: train["idhogar"].unique()

Out[10]: array(['21eb7fcc1', '9e5d7a658', '2c7317ea8', ..., 'a8eefc29',
      '1220b6fc6', 'd6c886aa3'], dtype=object)

In [11]: train["dependency"].value_counts()[0:10]

Out[11]:
yes      2192
no       1747
.5       1497
2         730
1.5       713
.33333334 598
.66666669 487
8         378
.25       260
3         236
Name: dependency, dtype: int64

In [12]: train["dependency"][train["dependency"].isin(["yes","no"])==False].value_counts().head(1)

Out[12]:
.5      1497
Name: dependency, dtype: int64

In [13]: train.dependency[train.dependency.isin(["yes","no"])==False].mode()

Out[13]:
0      .5
dtype: object
```

changed object datatype to float of dependency column

changed object datatype to float of dependency column

```
In [15]: train.edje[train.edje.isin(["yes","no"])==False].astype("float")
```

```
Out[15]:
```

	id	v2a1	hacdor	rooms	hacapo	v14a	refrig	v18q	v18q1	r4h1	r4h2	...	age	SQBescolari	SQBage	SQBhogar_total	SQBdeje	SQBhogar_nin	SQBovercrowding	SQBdepende
ID_279628684	10.0																			
ID_1c7f8646d	12.0																			
ID_d671d8b9c	11.0																			
ID_0561d69b5	11.0																			
ID_ec89ba7b7	11.0																			
ID_d45ae367d	9.0																			
ID_c94744e07	9.0																			
ID_887cd68f8	9.0																			
ID_c8f54dc01	9.0																			
ID_a38c64491	9.0																			


```
Name: edje, Length: 5672, dtype: float64
```

```
In [16]: train.edje[train.edje.isin(["yes","no"])==False].astype("float").median()
```

```
Out[16]:
```

```
7.0
```

```
In [17]: np.median(train.edje[train.edje.isin(["yes","no"])==False].astype("float"))
```

```
Out[17]:
```

```
7.0
```

```
In [18]: med=np.median(train.edje[train.edje.isin(["yes","no"])==False].astype("float"))  
train["edje"]=train["edje"].replace("yes",med).replace("no",0).astype("float")  
test["edje"]=test["edje"].replace("yes",med).replace("no",0).astype("float")
```

changed object datatype to float of edje column

```
In [16]: train.edje[train.edje.isin(["yes","no"])==False].astype("float").median()
Out[16]: 7.0

In [17]: np.median(train.edje[train.edje.isin(["yes","no"])==False].astype("float"))
Out[17]: 7.0

In [18]: med=np.median(train.edje[train.edje.isin(["yes","no"])==False].astype("float"))
train["edjea"]=train["edjea"].replace("yes",med).replace("no",0).astype("float")
test["edjea"]=test["edjea"].replace("yes",med).replace("no",0).astype("float")
```

Count how many null values are existing in columns

```
test["edjefa"] = test["edjefa"].replace("yes", med).replace("no", 0).astype("float")
```

changed object datatype to float of edjefa column

```
In [20]: train.describe(include="object")
```

```
Out[20]:
```

	idhogar
count	9557
unique	2988
top	id8a5d014

Training dataset

In [23]:	<pre>print("The percentage of missing values in rez_esc is ",(train["rez_esc"].isna().sum()/train.shape[0])) print("The percentage of missing values in v18q1 is ",(train["v18q1"].isna().sum()/train.shape[0])) print("The percentage of missing values in v2a1 is ",(train["v2a1"].isna().sum()/train.shape[0])) The percentage of missing values in rez_esc is 0.829498216595166 The percentage of missing values in v18q1 is 0.76822376989895922 The percentage of missing values in v2a1 is 0.717798472323951</pre>									
Out[23]:	The percentage of missing values in rez_esc is 0.829498216595166 The percentage of missing values in v18q1 is 0.76822376989895922 The percentage of missing values in v2a1 is 0.717798472323951									

Testing dataset

v18q1	10126
v2a1	17403
SQBmeaned	31
meaneduc	31
hogar_adul	0
parentesco10	0
parentesco11	0
parentesco12	0
idhogar	0
dtype: int64	0

Training dataset

```
In [25]: print("The percentage of missing values in rez_esc is ",(train["rez_esc"].isna().sum()/train.shape[0]))
```