

Curso de Big Data

Trabajo Practico 3

Integrantes del grupo: Sebastián Einstoss, Sofia Marinkovic y Pilar Ruiz Orrico

Parte 1: Análisis de la base de hogares y cálculo de pobreza

Ejercicio 1

Algunas variables de la base de hogar que consideramos que serían muy útiles para predecir la pobreza por ingresos y potencializar el tp2 serían las relacionadas con las cüestiones estructurales del hogar. Algunos ejemplos son:

- acceso a agua (IV6)
- tipo de suelo del hogar (IV3)
- si tiene o no baño en el hogar (IV6).

Existe cierta evidencia en la literatura que la pobreza estructural y la pobreza por ingreso tienen cierta correlación en el largo plazo, por eso la incorporación.

Ejercicio 2

	CODUSU	ANO4	TRIMESTRE	NRO_HOGAR	REALIZADA	REGION	MAS_500	
0	TQRMNOQUPHMKKUCDEJAH00701956		2021	1	1	1	44	N
1	TQRMNOPWXPWHMOKRCDEGLDF00701361	2021		1	1	1	41	N

2 rows × 88 columns

	CODUSU	ANO4	TRIMESTRE	NRO_HOGAR	REALIZADA	REGION	MAS_500	
15	TORMNOQXUJHMMUCDEJAH00693031	2021		1	1	1	1	S
21	TQRMNOSRPHLLNRCDEJAH00651171	2021		1	1	1	1	S
24	TQRMNOPVUHKLMLNCDEIAAD00655817	2021		1	1	1	1	S
36	TQRMNORYTHLOPMUCDEJAH00655933	2021		1	1	1	1	S
43	TQRMNOGPVHLMNNUCDEIAAD00655042	2021		1	1	1	1	S

5 rows × 88 columns

	CODUSU	ANO4	TRIMESTRE	NRO_HOGAR	COMPONENTE	H15	REGION	M
0	TQRMNOPPPHKLMLNCDEFAIH00646702	2021		1	1	1	1	43
1	TQRMNOPPRHKLMLNCDEFAIH00655104	2021		1	1	1	1	43

2 rows × 177 columns

	CODUSU	ANO4	TRIMESTRE	NRO_HOGAR	COMPONENTE	H15	REGION	
36799	TQRMNOPPRHKLMLNCDEIAAD00655703	2021		1	1	1	1	
36800	TQRMNOPPRHKLMLNCDEIAAD00655703	2021		1	1	2	1	
36801	TQRMNOPPPWHKMLNUCDEIAAD00655837	2021		1	1	1	1	
36802	TQRMNOPPPWHKMLNUCDEIAAD00655837	2021		1	1	2	1	
36803	TQRMNOPPPWHKMLNUCDEIAAD00655837	2021		1	1	3	1	

5 rows × 177 columns

Ejercicio 3

Longitud base hogares: 1474 ¡Longitud base individual: 4082

	CODUSU	ANO4	TRIMESTRE	NRO_HOGAR	COMPONENTE	H15	REGION	M
0	TQRMNOPPRHKLMLNCDEIAAD00655703	2021		1	1	1	1	1
1	TQRMNOPPRHKLMLNCDEIAAD00655703	2021		1	1	2	1	1
2	TQRMNOPPPWHKMLNUCDEIAAD00655837	2021		1	1	1	1	1
3	TQRMNOPPPWHKMLNUCDEIAAD00655837	2021		1	1	2	1	1
4	TQRMNOPPPWHKMLNUCDEIAAD00655837	2021		1	1	3	1	1

5 rows × 242 columns

Ejercicio 4

Funciones para limpiar la base de datos

Para limpiar la base de datos de observaciones sin sentido vamos a usar la función "drop()" de pandas. Para ello vamos a ponerle diferentes condiciones, si no las cumple se dropa la fila.

Asimismo para rellenar los missing values usaremos de pandas la función "fillna()" y completaremos con lo que consideremos acredo (Por ejemplo la mediana).

Para convertir variables categoricas en dummies usaremos del paquete numpy la función "where()".

Por último vamos a eliminar las variables sin sentido para el análisis. Para ello usaremos nuevamente la función "drop()" de pandas.

Ejercicio 5

En este inciso filtramos aquellos valores que no tiene sentido, como por ejemplo: Edades negativas, variables de ingreso que tomen valores negativos.

También reemplazamos los valores faltantes por la mediana y dropamos variables que no consideramos relevantes. En el código se aclara cuales son.

Finalmente para realizar el analisis convertimos las variables categoricas en dummies.

	CODUSU	ANO4	TRIMESTRE	NRO_HOGAR	COMPONENTE	H15	REGION	M
0	TQRMNOPPRHKLMLNCDEIAAD00655703	2021		1	1	1	1	1
1	TQRMNOPPRHKLMLNCDEIAAD00655703	2021		1	1	2	1	1

2 rows × 242 columns

	CODUSU	ANO4	TRIMESTRE	NRO_HOGAR	COMPONENTE	H15	REGION	M
0	TQRMNOPPRHKLMLNCDEIAAD00655703	2021		1	1	1	1	1
1	TQRMNOPPRHKLMLNCDEIAAD00655703	2021		1	1	2	1	1
2	TQRMNOPPPWHKMLNUCDEIAAD00655837	2021		1	1	2	1	1
3	TQRMNOPPPWHKMLNUCDEIAAD00655837	2021		1	1	3	1	1
4	TQRMNOPPPWHKMLNUCDEIAAD00655837	2021		1	1	4	1	1

5 rows × 242 columns

	CODUSU	NRO_HOGAR	REGION	MAS_500	AGLOMERADO	CH04	CH06	PP0
0	TQRMNOPPRHKLMLNCDEIAAD00655703	1	1	1		32	1	29
1	TQRMNOPPRHKLMLNCDEIAAD00655703	1	1	1		32	2	29

2 rows × 962 columns

Ejercicio 6

5 Variables relevantes de la encuesta de hogares

IV1[4] = Tipo de vivienda, pieza en hotel / pensión

IV8 = Baño/letrina

IV5 = Ingresos de subsidio o ayuda social(en dinero)del gobierno, iglesias

IV3[3] = Tipo de suelo, ladrillo suelto / tierra

IV6[2] = Acceso al agua, fuera de la vivienda pero dentro del terreno

	IV1_4	IV8	IV5	IV3_3	IV6_2
count	3308.000000	3308.000000	3308.000000	3308.000000	3308.000000
mean	0.003023	0.998186	0.189541	0.002116	0.013603
std	0.054907	0.042556	0.391997	0.045959	0.115855
min	0.000000	0.000000	0.000000	0.000000	0.000000
25%	0.000000	1.000000	0.000000	0.000000	0.000000
50%	0.000000	1.000000	0.000000	0.000000	0.000000
75%	0.000000	1.000000	0.000000	0.000000	0.000000
max	1.000000	1.000000	1.000000	1.000000	1.000000

Ejercicio 7

	Tabla de equivalencias de necesidades energéticas. Unidades de adulto equivalente, según sexo y edad	Unnamed: 1	Unnamed: 2	Unnamed: 3
0		NaN	NaN	NaN
1		NaN	NaN	NaN
2	Edad	Mujeres	Varones	NaN
3		NaN	NaN	NaN
4	Menor de 1 año	0.35	0.35	NaN

	Tabla de equivalencias de necesidades energéticas. Unidades de adulto equivalente, según sexo y edad	Unnamed: 1	Unnamed: 2
0		NaN	NaN
1		NaN	NaN
2	Edad	Mujeres	Varones
3		NaN	NaN
4	Menor de 1 año	0.35	0.35

	Tabla de equivalencias de necesidades energéticas. Unidades de adulto equivalente, según sexo y edad	Unnamed: 1	Unnamed: 2
2	Edad	Mujeres	Varones
4	Menor de 1 año	0.35	0.35
5	1 año	0.37	0.37
6	2 años	0.46	0.46
7	3 años	0.51	0.51
8	4 años	0.55	0.55

	Edad	sexo	value
0	Menor de 1 año	Mujeres	0.35
1	1 año	Mujeres	0.37
2	2 años	Mujeres	0.46
3	3 años	Mujeres	0.51
4	4 años	Mujeres	0.55

	CH06	adulto_equiv	CH04
0	0	0.35	2
1	1	0.37	2
2	2	0.46	2
3	3	0.51	2
4	4	0.55	2
5	5	0.60	2
6	6	0.64	2
7	7	0.66	2
8	8	0.68	2
9	9	0.69	2

	CODUSU	NRO_HOGAR	REGION	MAS_500	AGLOMERADO	CH04	CH06	PP0
0	TQRMNOPPRHKLMLNCDEIAAD00655703	1	1	1		32	1	29
1	TQRMNOPPRHKLMLNCDEIAAD00655703	1	1	1		32	2	29
2	TQRMNOPPPWHKMLNUCDEIAAD00655837	1	1	1		32	2	53
3	TQRMNOPPPWHKMLNUCDEIAAD00655837	1	1	1		32	2	22
4	TQRMNOPPPWHKMLNUCDEIAAD00655837	1	1	1		32	1	19

5 rows × 963 columns

	ad_equiv_hogar
CODUSU	
TQRMNOPPOHJMLQCDEJAH00702455	1.67
TQRMNOPPOHJMLQCDEJAH00698190	3.10
TQRMNOPPOHJONGCDEJAH00693114	1.76
TQRMNOPPOHJMOSCDCEJAH00656008	2.53
TQRMNOPPOHLMPPCCDEJAH00701610	3.92

	CODUSU	NRO_HOGAR	REGION	MAS_500	AGLOMERADO	CH04	CH06	PP0
0	TQRMNOPPRHKLMLNCDEIAAD00655703	1	1	1		32	1	29
1	TQRMNOPPRHKLMLNCDEIAAD00655703	1	1	1		32	2	29
2	TQRMNOPPPWHKMLNUCDEIAAD00655837	1	1	1		32	2	53
3	TQRMNOPPPWHKMLNUCDEIAAD00655837	1	1	1		32	1	27
4	TQRMNOPPPWHKMLNUCDEIAAD00655837	1	1	1		32	1	19

5 rows × 964 columns

Ejercicio 8

2474

821

0.24916540212443095

	CODUSU	NRO_HOGAR	REGION	MAS_500	AGLOMERADO	CH04	CH06	PP0
0	TQRMNOPPRHKLMLNCDEIAAD00655703	1	1	1		32	1	29
1	TQRMNOPPRHKLMLNCDEIAAD00655703	1	1	1		32	2	29
7	TQRMNOPPPXHLNLCDEIAAD00690221	1	1	1		32	1	66
8	TQRMNOPPPXHLNLCDEIAAD00690221	1	1	1		32	1	27
9	TQRMNOPPPXHLNLCDEIAAD00690221	1	1	1		32	2	65

5 rows × 965 columns

Ejercicio 9

	CODUSU	NRO_HOGAR	REGION	MAS_500	AGLOMERADO	CH04	CH06	PP0
0	TQRMNOPPRHKLMLNCDEIAAD00655703	1	1	1		32	1	29
1	TQRMNOPPRHKLMLNCDEIAAD00655703	1	1	1		32	2	29
7	TQRMNOPPPXHLNLCDEIAAD00690221	1	1	1		32	1	66
8	TQRMNOPPPXHLNLCDEIAAD00690221	1	1	1		32	1	27
9	TQRMNOPPPXHLNLCDEIAAD00690221	1	1	1		32	2	65

5 rows × 966 columns

	CODUSU	NRO_HOGAR	REGION	MAS_500	AGLOMERADO	CH04	CH06	PP0
0	TQRMNOPPRHKLMLNCDEIAAD00655703	1	1	1		32	1	29
1	TQRMNOPPRHKLMLNCDEIAAD00655703	1	1	1		32	2	29
7	TQRMNOPPPXHLNLCDEIAAD00690221	1	1	1		32	1	66
8	TQRMNOPPPXHLNLCDEIAAD00690221	1	1	1		32	1	27
9	TQRMNOPPPXHLNLCDEIAAD00690221	1	1	1		32	2	65

5 rows × 966 columns

Ejercicio 10

1868

680

```
pobre
0    2385188
1    1356105
Name: PONDIT, dtype: int64

36.24696060960743
```

La tasa de pobreza para hogares de GBA del primer trimestre 2021 se ubica en 36,25%.

Si bien da un resultado similar al que publicó el INDEC (37,1%) nunca nos va a dar exactamente lo mismo dado que no sabemos en que mes fueron tomados los datos de los diferentes hogares y por lo tanto estamos comparando los ingresos contra la canasta básica promedio del trimestre y no con la que le correspondería a cada familia (según el mes de relevamiento de los datos).

Parte 2: Construcción de funciones

Ejercicio 1: función evalua_metodo

Ejercicio 2: función cross_validation

Ejercicio 3: función evalua_config

Ejercicio 4: función evalua_multiples_metodos

Diccionario de métodos

valores que pueden tomar "penalty" y "C"

Parte 3: Clasificación y Regularización

Ejercicio 1

928

927

Ejercicio 2

	accuracy	auc	ecm	hiperparametro	modelo
0	0.794078	0.761558	0.205922	1.0	Regresión logística
1	0.775236	0.745101	0.224764	NaN	Análisis de Discriminante Lineal
2	0.736205	0.684876	0.263795	NaN	3 vecinos cercanos
3	0.804845	0.771605	0.195155	NaN	Arbol de decisión
4	0.810229	0.779627	0.189771	NaN	Support vector machines (SVM)
5	0.837147	0.800245	0.162853	NaN	Bagging
6	0.804845	0.753609	0.195155	NaN	Random Forests
7	0.826380	0.784949	0.173620	NaN	Boosting

Ejercicio 3

Para elegir el λ por validación cruzada hay que seguir una serie de pasos.

Para un λ dado se realizan los siguientes pasos:

En primer lugar partimos los datos de la base en k partes iguales. Luego, ajustamos el modelo dejando afuera una de las k partes. Computamos el error de predicción para los datos no utilizados. Este procedimiento lo repetimos para todos los k (en cada caso se deja una parte k afuera del ajuste del modelo). Luego, se promedian todos los errores de predicción.

Este procedimiento lo repetimos para los diferentes valores de λ.

Por último comparamos los errores de predicción promedio calculados con cada λ y nos quedamos con el λ que genera un menor resultado. Ese será nuestro λ elegido.

¿Por qué se deja afuera una parte de testeo?

Nuestro objetivo es medir el error de pronóstico fuera de la muestra. Para no separar arbitrariamente la base en entrenamiento y test se utiliza este método que permite separar la muestra en k partes e ir estimando con k-1 partes. Siempre hay una parte k que no se usa en la estimación porque es la que se va a usar luego para medir el error de predicción. De esta forma, se puede ver esto "fuera de los datos de entrenamiento" ya que se dejó una parte por fuera del ajuste de modelo. En concreto se estima todo k veces y en cada estimación se deja afuera una de las partes.

Ejercicio 4

Por un lado, si k es chico maximiza los datos para estimar pero es sensible a los valores particulares de la muestra.

Por el otro lado, si k es grande maximiza los datos para evaluar pero el modelo estimado es menos preciso.

En el caso particular de k=n el modelo se estima n veces con n-1 datos.

Como "regla general" se usa un k=5 o k=10

Ejercicio 5

```
[99999.99999999999,
10000.0,
1000.0,
100.0,
10.0,
1.0,
0.1,
0.01,
0.001,
0.0001,
1e-05]</
```