

# Capítulo 1

## Regresión Logística

Si un banco le da un préstamo a una persona ¿qué tan probable es que esa persona le pague? ¿Cómo modelamos esta probabilidad? De esto, en muchas ocasiones, se encarga la regresión logística.

Hagamos un ejemplo concreto, vamos trabajar con una base de datos que llamamos `base_cred` y la leemos python cómo sigue.

Primero instalamos la biblioteca que nos permite leer el archivo e inmediatamente llamamos a la base

```
1 import pandas as pd
2 base_cred=pd.read_csv('incumplimiento.csv')
```

y después de eso podemos imprimir la base escribiendo

```
1 base_cred
```

pero esto puede resultar muy engorroso a la hora de visualizar la base, así que basta ver las primeras 5 entradas,

```
1 base_cred.head(5)
```

	Folio	Incumplimiento	Estudiante	Saldo	Ingreso
2	1	No	No	729.526495	44361.625074
3	2	No	Yes	817.180407	12106.134700
4	3	No	No	1073.549164	31767.138947
5	4	No	No	529.250605	35704.493935
6	5	No	No	785.655883	38463.495879

Ahora, sí queremos ver las últimas entradas de la base, escribimos `tail` en lugar de `head`.

También podemos ver la información general de nuestra base

```
1 base_cred.info()
```

```

1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 10000 entries, 0 to 9999
3 Data columns (total 5 columns):
4 Folio                10000 non-null int64
5 Incumplimiento       10000 non-null object
6 Estudiante           10000 non-null object
7 Saldo                10000 non-null float64
8 Ingreso              10000 non-null float64
9 dtypes: float64(2), int64(1), object(2)
10 memory usage: 390.7+ KB

```

Después del breve resumen de la base que hemos importado al código, ¿cómo accedemos a cada una de las variables? Dos maneras comunes

```
1 base_cred.NombredelaVariable
```

```
1 base_cred['NombredelaVariable']
```

y así podemos hacer un análisis descriptivo de la gráfica. Antes debemos de instalar las siguientes bibliotecas

```

1 from matplotlib import*
2 from pylab import*

```

Intuitivamente ¿que variables creen que influyan en la probabilidad de incumplimiento? ¿el Ingreso? ¿el Saldo?, sería bueno graficar Ingreso vs Saldo y ver, a partir allí quienes incumplieron. Para esto necesitamos acceder a los datos *créditos que incumplieron e identificarlos* y *créditos que cumplieron e identificarlos*

```

1 base_cred[base_cred.Incumplimiento=='No']
2 base_cred[base_cred.Incumplimiento=='Yes']

```

Ya que los tenemos identificados, ahora seleccionamos Saldo de los que incumplieron y Saldo de los que cumplieron. Hacemos lo mismo con Ingresos. Esto nos va a facilitar “clasificar” a los cumplidores de los incumplidores en relación con el saldo y el ingreso.

```

1 s1=base_cred[base_cred.Incumplimiento=='No'].Saldo
2 i1=base_cred[base_cred.Incumplimiento=='No'].Ingreso
3 s2=base_cred[base_cred.Incumplimiento=='Yes'].Saldo
4 i2=base_cred[base_cred.Incumplimiento=='Yes'].Ingreso

```

y luego graficamos

```

1 s1=base_cred[base_cred.Incumplimiento=='No'].Saldo
2 i1=base_cred[base_cred.Incumplimiento=='No'].Ingreso
3 s2=base_cred[base_cred.Incumplimiento=='Yes'].Saldo
4 i2=base_cred[base_cred.Incumplimiento=='Yes'].Ingreso
5
6 figsize(16,8)
7 scatter(s1,i1,marker='o',s=40,edgecolors='g',facecolors='g',alpha=0.2,label='Bien
8 portados')
9 scatter(s2,i2,marker='o',s=40,edgecolors='red',facecolors='red',alpha=0.5,label='
10 Mal portados')
11 legend()
12 xlabel('Saldo')
13 ylabel('Ingreso')

```

y lo que debemos obtener es la figura 1.1 ¿qué observamos en esta gráfica? Nos sugiere que el Saldo del crédito tiene un peso importante en la probabilidad de incumplimiento. Así que, tal vez debamos de incluirlo como variable explicativa, de este modo un primer objetivo es modelar

$$Pr(\text{Incumplimiento} = Si | \text{Saldo})$$

Haciendo una identificación adecuada, vamos a trabajar con el modelo

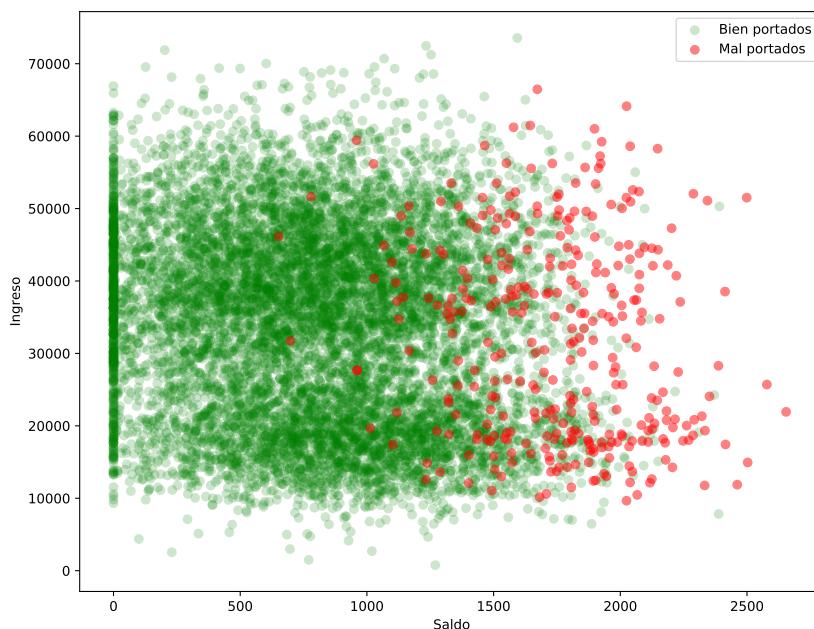


Figura 1.1: De verde aparecen las personas que incumplieron y de rojo las que no

$$p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}} \quad (1.1)$$

La figura 1.1 nos sugiere que una de las variables que influye en la probabilidad de incumplimiento es Saldo, entonces lo que sigue ahora es graficar

Saldo vs Incumplimiento

Y esto lo hacemos con el siguiente código,

```
1 scatter(base_cred.Saldo, base_cred.Incumplimiento, )
2 xlabel('Saldo')
3 ylabel('Incumplimiento')
```

y nos devuelve la figura 1.2, qué como vemos, esta a su vez nos sugiere que debemos ajustarle una función logística. Esto lo hacemos como sigue (siempre teniendo en cuenta quien es la variable de respuesta y quien es la variable explicativa). Hacemos el ajuste,

```
1 import sklearn.linear_model as skl_lm
2 base_cred['In_num']=base_cred.Incumplimiento.factorize()[0]
3 y=base_cred.In_num#Variable de Respuesta
4 x_1 = base_cred.Saldo.values.reshape(-1,1)#Variable Explicativa
```

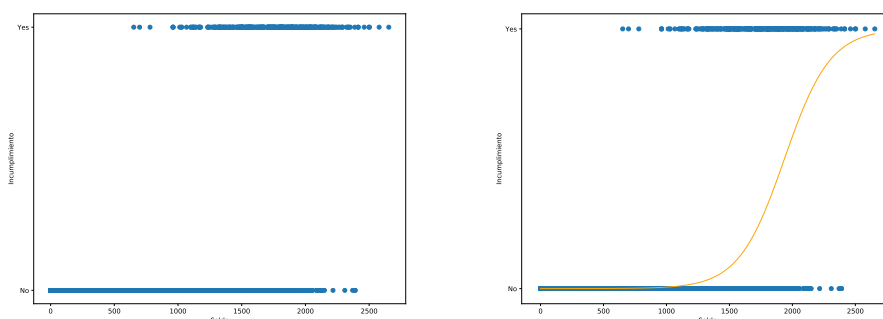


Figura 1.2: Con y sin curva ajustada

```

5 x_test = np.arange(base_cred.Saldo.min(), base_cred.Saldo.max()).reshape(-1,1)
6 clf = skl_lm.LogisticRegression(solver='newton-cg')
7 clf.fit(x_1,y)
8 prob=clf.predict_proba(x_test)

```

y luego graficamos la curva ajustada,

```

1 plot(prob[:,1],color='orange')
2 scatter(base_cred.Saldo,base_cred.Incumplimiento,)
3 xlabel('Saldo')
4 ylabel('Incumplimiento')

```

El resultado de la curva ajustada lo podemos ver en la parte de la derecha de la figura 1.2. Y entonces muchos se preguntarán ¿Es un buen ajuste? Po's quien sabe, de eso precisamente trata la tarea.

## Tarea 2(Regresión Logística): Primera parte

**Problema 1.** *La idea de este ejercicio es ver si la curva naranja ajusta bien en el siguiente sentido.*

**Algoritmo 1.** ¿Qué tal el ajuste?

1. De la base original, calcule la proporción de incumplimiento (si etiquetamos con 1 si incumplio y con 0 si no, entonces lo que estoy pidiendo es simplemente la proporción de 1's.).
2. Calcule  $\hat{\beta}_0$  y  $\hat{\beta}_1$  usando python con clf (esto ya lo hicimos en clases).
3. Para cada crédito, calcule

$$f(\hat{\beta}_0 + \hat{\beta}_1 \text{Saldo}_i),$$

donde  $f(x) = \frac{e^x}{1+e^x}$ .

4. Luego defina,  $\hat{Y}_i = 1$  si  $f(\hat{\beta}_0 + \hat{\beta}_1 \text{Saldo}_i) > 0,5$ . Y calcule la proporción de 1's y compare con la proporción original.

5. Luego calcule la proporción en donde sí coincidieron los 1's, la proporción en donde sí coincidieron los 0's, la proporción en donde en lugar de 1's salieron 0, y por último la proporción en donde en lugar de 0's salieron 1's. Explique detalladamente sus conclusiones.

**Problema 2.** *Luego, considere el modelo en donde tomamos como variables explicativas: Saldo, Estudiante e Ingreso y nuevamente vea el “ajuste” cómo en el algoritmo [1](#).*