

Construcción de Modelos de Credit Scoring con R

Nov 2022

FRANCISCO J. RODRÍGUEZ ARAGÓN

SCORE MANAGER at ONEY

Ph D in Statistic by Cordoba Univ.

Associate Professional Risk Manager Certified by APRM

Operational Risk Manager Certified by APRM



Summary

1. ¿Qué es Oney?

2. ¿Qué es un modelo de credit scoring?

3. Pasos esenciales en la construcción de modelo de CS

4. DEMO de la librería *scorecard*

INTRODUCCIÓN

¿Qué es Oney?

¿Qué es Oney?

Especialista en medios de pago y financiación para el retail y el consumidor final

- 12 países
- 8M de clientes
- 35 años acompañando a los retailers
- +600 partners en Europa

Tarjetas

Pago Aplazado 3x 4x

E-financiación 6x 10x 12x

Crédito al consumo

Seguros





PARTE 1

¿Qué es un modelo de Credit Scoring?

one



oney | 5

Los modelos más utilizados en la actualidad

- En esencia los modelos de Credit Scoring no son más que una reformulación de un modelo de regression logística subyacente
- Son modelos poco conocidos por los Data Scientist actuales a pesar de su gran aplicación actual
- Los modelos de credit scoring están en:
 - Cada compra que se realiza con una tarjeta de crédito, aunque aquí poco a poco están llegando (que no imponiéndose) los modelos de Machine Learning
 - En la concesión de hipotecas y en general en todo tipo de préstamos al consumo ya que el regulador bancario prácticamente obliga a las entidades a usar este tipo de modelos en sus créditos

En España, las últimas cifras del **Banco de España** reflejan esta tendencia al alza. Por un lado, el saldo de créditos al consumo alcanzó casi los 188.000 millones de euros en junio, **un 5% más que a principios de año**. Si bien hay un pico cíclico en esta fecha, por todos los créditos que se piden en junio para contratar vacaciones, este auge se apoya en otras dos cifras.

188.000 mill. De euros al menos gestionados por estos modelos en el sector revolving!!!!

RIESGOS

Alarma bancaria por un 'boom' de créditos al consumo para combatir la inflación

La financiación al consumo crece un 5% y el saldo en tarjetas 'revolving' alcanza máximos de dos años y medio. Mientras, la morosidad de las financieras toca el pico desde 2016

Por **Jorge Zuloaga**

03/08/2022 - 05:00

Un modelo “Tarjeta de Puntuación”

Variable	Tramo	Tasa mora	WOE	Score	Beta	IV	SCR
Pasivo	-100	31,46%	-0,797	36	-0,4032	0,2775	10,99%
	100-800	17,90%	-0,053	45			
	800-2000	18,23%	0,305	49			
	2000-3700	11,10%	0,504	51			
	3700+	1,11%	0,992	57			
Triad	-650	22,74%	-0,353	39	-0,6724	0,2667	35,91%
	650-675	15,21%	0,136	48			
	675-700	9,11%	0,724	59			
	700+	1,60%	1,456	74			
Ratio_AC_PAS	1	34,18%	-0,921	36	-0,3718	0,1871	5,16%
	2	22,34%	-0,330	42			
	3	16,89%	0,017	46			
	4	12,98%	0,537	49			
Nacionalidad	1	27,48%	-0,335	33	-0,7051	0,1345	18,66%
	2	22,37%	-0,110	39			
	3	18,75%	-0,110	41			
	4	16,50%	0,046	45			
	5	13,25%	0,303	53			
	6	7,50%	0,937	58			
Incidencias	CON INCIDENCIAS	27,40%	-0,601	34	-0,5171	0,1149	6,56%
	SIN INCIDENCIAS	14,53%	0,196	48			
Antigüedad	-11	19,61%	-0,165	41	-0,6715	0,0802	9,92%
	11-29	18,84%	-0,116	43			
	29-53	14,21%	0,222	50			
	53+	9,60%	0,666	58			
Profesion	1	22,78%	-0,355	40	-0,4904	0,0432	2,92%
	2	19,20%	-0,139	43			
	3	16,19%	0,068	46			
	4	14,78%	0,176	48			
	5	9,95%	0,627	54			
Antig_Emp	-12	19,36%	-0,149	43	-0,6211	0,0419	3,75%
	12-36	15,28%	0,137	48			
	36+	11,80%	0,436	53			
Est_Civil	1	19,52%	-0,159	42	-0,7968	0,0351	4,80%
	2	17,41%	-0,019	45			
	3	14,32%	0,213	50			
	4	12,33%	0,386	54			
Provincia	1	20,26%	-0,206	42	-0,6165	0,0319	3,13%
	2	18,05%	-0,063	44			
	3	16,53%	0,044	46			
	4	13,64%	0,270	50			
	5	11,45%	0,469	54			

En los scoring las variables deben ser trameadas ¿Cómo se consigue esto?

Un modelo de Credit Scoring o de tarjeta de puntuación es un algoritmo donde a cada variable de entre una colección, se les asocia, en función de su valor, una puntuación denominada score parcial. Las variables son tomadas en general sobre un individuo, empresa, transacción, etc y por tanto la suma de los valores de todas las variables consideradas sobre dicho elemento de la población es lo que se denomina, scoring de dicho elemento

Fuente de la scorecard: <https://docplayer.es/49254028-Desarrollo-y-validacion-de-modelo-de-scoring-de-admision-para-tarjetas-de-credito-con-metodologia-de-inferencia-de-denegados.html>

En los scoring se generan puntuaciones que pueden ser algebraicamente sumadas para cada cliente ¿Cómo se consigue esto?

PARTE 2

Pasos esenciales en la construcción de un modelo de CS



¿Cómo se estima la probabilidad de Impago?

- Como en todo modelo, el paso más difícil es la construcción de una tabla que contengan *variables explicativas* junto con una *variable objetivo* o *target*
- Se trata de predecir la probabilidad de impago condicionada al valor de unas *variables explicativas*

$$P(Y = 1 | v_1; v_2; \dots; v_n) = ?$$

- Mientras que, en una regresión logística, el anterior número para un determinado cliente se estima del siguiente modo:

$$P(Y = 1 | v_1; \dots; v_n) = \frac{e^{a_0 + a_1 v_1 + \dots + a_n v_n}}{1 + e^{a_0 + a_1 v_1 + \dots + a_n v_n}}$$

- En un Credit Scoring, dicho valor se estima del siguiente modo (bajo las condiciones de construcción que se verán en la demo, ambas formulaciones son equivalentes):

$$Score = a + b \left[\frac{P}{1 - P} \right] = a + b[a_0 + a_1 v_1 + \dots + a_n v_n]$$

Es habitual tomar el valor de a igual a: 487.123 y el de b igual a: 28.8539. Aunque pueden tomarse otros como en el ejemplo siguiente (600; 28.85)

Un modelo muy interpretable

- La construcción de este tipo de modelos requiere un trameado previo de las variables: La transformación Weight of Evidence que convierte todas las variables en tramos de puntuación continua de modo que al final se saben cosas como:

$$\begin{aligned} \text{If } Age \in [18; 35] \quad \text{Points} &= 600 + 28.85 \cdot (3.45 + 5.45 \cdot -0.29) = 654.30 \\ \text{If } Age \in [36; 45] \quad \text{Points} &= 600 + 28.85 \cdot (3.45 + 5.45 \cdot 0.27) = 742.29 \\ \text{If } Age \in [46; Inf] \quad \text{Points} &= 600 + 28.85 \cdot (3.45 + 5.45 \cdot 0) = 699.53 \end{aligned}$$

***** Importancia de variables Voting_Classifier *****

Variables	Logit	Random_Forest
log_providerA_score	0.376384	0.366804
provider2_score	0.002200	0.231568
order_amount	0.000355	0.108292
num_installments_initial	0.038441	0.099995
dum_cellular_ip_address	0.368038	0.051138
dum_wifi_ip_userType	-0.355291	0.023141
dum_traveler_ip_address	0.000000	0.000000
pm_card_expiration_year	-0.000051	0.016783
provider1_score	-0.231073	0.044099
dum_hosting_ip_address	0.000000	0.000000
dum_master_method_card_type	0.011679	0.005051
dum_extranjero_pm_card_country_code	0.000000	0.000070
num_sdad	-0.005074	0.040019
device_cookies_enabled	0.020754	0.003667
dum_anex_method_card_type	0.000000	0.000000
dum_business_ip_address	0.000000	0.000126
customer_visit	-0.020313	0.000013
dum_satellite_ip_userType	0.000000	0.000000
ip_reputation	-0.481478	0.000235

- ¿Por qué un CS frente a un Random Forest?

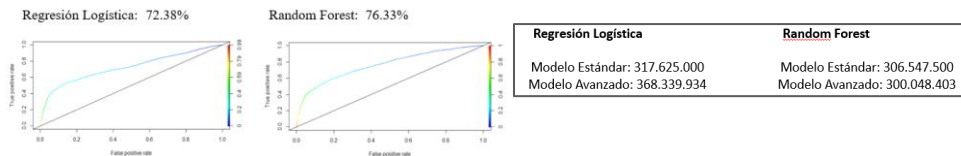
EBA DISCUSSION PAPER ON MACHINE LEARNING FOR IRB MODELS

EBA/DP/2021/04

11 NOVEMBER 2021

En este caso particular se observa como un Random Forest mejora la predictividad desde un 0.7238 de una regresión logística hasta un 0.7633, es decir, algo más de un 5.45%

Fuente: <https://github.com/FJROAR/Ejemplo-Blog-RWA>



PARTE 3

DEMO DE LA LIBRERÍA scorecard



Desarrollo de la demo CS con R

- La librería *scorecard* fue desarrollada para *R* y actualmente Python ha copiado su funcionalidad también, por lo que existe en los 2 lenguajes. A pesar de la importancia de estos modelos, esta librería es bastante reciente y se puso disponible alrededor del 2019. Su uso en R ahorra muchas horas de trabajo en la construcción de estos modelos
- Paso 1: Lectura de la información

```
#https://cran.r-project.org/web/packages/scorecard/vignettes/demo.html  
  
library(data.table)  
library(dplyr)  
library(scorecard)  
  
df <- fread("data/creditcard.csv", sep = ",")  
  
names(df)
```

Data	
df	284807 obs. of 31 variables

Desarrollo de la demo CS con R

- Paso 2: Limpieza de datos (adicional) y análisis de las variables (por tiempo, se elude este pasoe ya que los datos estaban muy limpios, en la realidad aquí debe usarse todo el tiempo necesario)
- Paso 3: Separate sample (hasta aquí, no hay nada distinto respecto a un modelo de Machine Learning actual)

```
df_list <- split_df(df, y = "Class", ratios = c(0.6, 0.4), seed = 30)
label_list <- lapply(df_list, function(x) x$Class)

df_train <- df_list$train
df_test <- df_list$test
```

Desarrollo de la demo CS con R

- Paso 4: TRAMEADO DE VARIABLES. En muchos modelos de ML este paso no se da y no es obligatorio, en cambio en los modelos CS es obligatorio darlo

```
bins_train <- woebin(df_train, y = "Class")

#Ejemplo de un bin:
#bins_train$Time
```

variable	bin	count	count_distr	neg	pos	posprob	woe	bin_iv	total_iv	breaks	is_special	values
Time	[-Inf,30000)	11395	0.06682736	11346	49	0.0043001316	0.923304704	9.339604e-02	0.3027489	30000	FALSE	
Time	[30000,40000)	12553	0.07361859	12543	10	0.0007966223	-0.766227983	3.021969e-02	0.3027489	40000	FALSE	
Time	[40000,85000)	61728	0.36201133	61623	105	0.0017010109	-0.006725162	1.631823e-05	0.3027489	85000	FALSE	
Time	[85000,110000)	8813	0.05168491	8772	41	0.0046522183	1.002356901	8.908773e-02	0.3027489	110000	FALSE	
Time	[110000,125000)	14296	0.08384062	14288	8	0.0005595971	-1.119628817	6.330402e-02	0.3027489	125000	FALSE	
Time	[125000, Inf)	61729	0.36201720	61650	79	0.0012797875	-0.291675712	2.672512e-02	0.3027489	Inf	FALSE	

- Paso 5: Construcción de un modelo de regresión logística subyacente con las variables trameadas (training del modelo)

```
dt_woe_list = lapply(df_list, function(x) woebin_ply(x, bins_train))
m1 = glm( Class ~ ., family = binomial(), data = dt_woe_list$train)

summary(m1)
```

```
> summary(m1)

Call:
glm(formula = Class ~ ., family = binomial(), data = dt_woe_list$train)

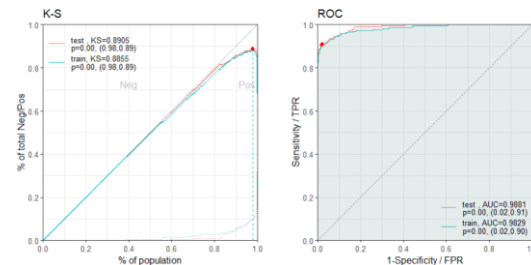
Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.4489  -0.0153  -0.0084  -0.0051   4.6185

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept) -5.117497    0.131948  -38.784  < 2e-16 ***
Time_woe     -0.328063    0.189628  -1.730  0.083624 .
V1_woe        0.221128    0.119524   1.850  0.064303 .
V2_woe       -0.326117    0.093742  -3.479  0.000504 ***
```


Desarrollo de la demo CS con R

- Paso 6: Evaluación del modelo (habitual en ML también)

```
pred_list = lapply(dt_woe_list, function(x) predict(m1, x, type='response'))
perf = perf_eva(pred = pred_list, label = label_list, show_plot = c("ks", "roc"),
               pred_desc = F)
```



- Paso 7: TARJETA DE PUNTUACIÓN. Este paso es exclusivo de estos modelos de CS

```
card <- scorecard(bins_train, m1)
```

variable	bin	count	count_distr	neg	pos	posprob	woe	bin_iv	total_iv	breaks	is_special_values	points
Time	[-Inf,30000)	11395	0.06682736	11346	49	0.0043001316	0.923304704	9.339604e-02	0.3027489	30000	FALSE	22
Time	[30000,40000)	12553	0.07361859	12543	10	0.0007966223	-0.766227983	3.021969e-02	0.3027489	40000	FALSE	-18
Time	[40000,85000)	61728	0.36201133	61623	105	0.0017010109	-0.006725162	1.631823e-05	0.3027489	85000	FALSE	0
Time	[85000,110000)	8813	0.05168491	8772	41	0.0046522183	1.002356901	8.908773e-02	0.3027489	110000	FALSE	24
Time	[110000,125000)	14296	0.08384062	14288	8	0.0005595971	-1.119628817	6.330402e-02	0.3027489	125000	FALSE	-26
Time	[125000, Inf)	61729	0.36201720	61650	79	0.0012797875	-0.291675712	2.672512e-02	0.3027489	Inf	FALSE	-7

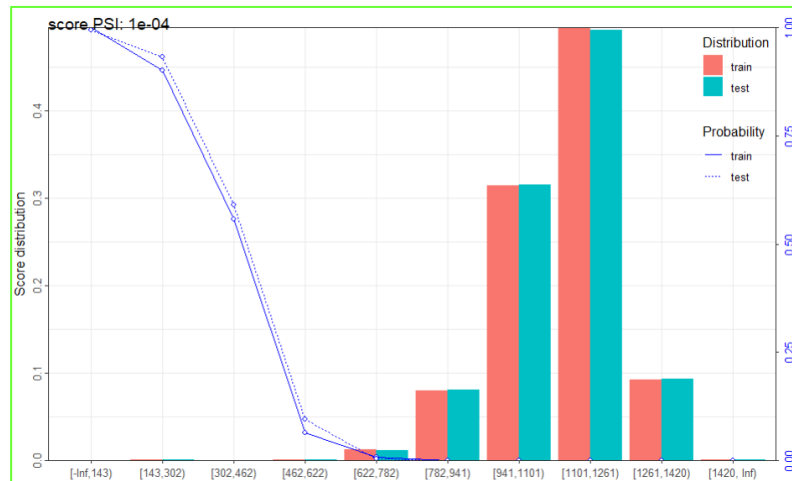
Desarrollo de la demo CS con R

- Paso 8: Aplicación de la tarjeta a conjunto de datos y evaluaciones adicionales como la estabilidad poblacional training vs test

```
# Obtain Credit Scores
score_train <- scorecard_ply(df_train, card)
score_test  <- scorecard_ply(df_test, card)

score_list = lapply(df_list, function(x) scorecard_ply(x, card))

# Analyze the PSI
perf_psi(score = score_list, label = label_list)
```



CONCLUSIONES FINALES

- La librería *scorecard* permite a día de hoy completar el ciclo de modelización para CS ahorrando mucho tiempo de programación al estadístico
- Es una librería muy flexible y optimizada para cubrir con eficiencia todas las necesidades de los CSs
- Se comporta bien ante conjuntos con elevado número de registros y variables, aunque en general en estos modelos los datasets suelen ser inferiores a 100.000 registros
- Tanto los tramos como algunos parámetros adicionales para construir tarjetas en distintas escalas, son modificables de modo trivial en el código evitando muchos errores operacionales en la construcción
- Su uso en R resulta muy sencillo y permite generar rápidamente modelos que se pueden poner en producción con extrema facilidad

A photograph of a woman in a light green shirt and blue jeans pushing two young girls with curly hair on a skateboard. The girls are wearing a pink and white striped shirt and a green shirt respectively. They are all smiling and looking towards the camera. The background is a blurred outdoor setting with trees and a fence.

**¡Muchas Gracias por su
atención!**

oney
YOUR MONEY YOUR WAY