



UNR Universidad
Nacional de Rosario

LICENCIATURA EN ESTADÍSTICA

DETECCIÓN DE FRAUDE

Un análisis mediante modelos lineales generalizados

Autores: Franco Santini - Nicolas Gamboa - Andrés Roncaglia

Docentes: Boggio Gabriela - Harvey Guillermina - Costa Victorio

2024

Tabla de contenidos

Introducción	1
Variables:	1
Análisis descriptivo	3
Modelado	8
Análisis de residuos	12
Evaluación de la componente sistemática	12
Comprobación de la distribución propuesta	13
Interpretaciones	13
Capacidad predictiva	14

Introducción

El fraude con tarjetas de crédito es una de las principales amenazas que sufren los bancos. Con el auge de la tecnología las transacciones digitales facilitaron los traspasos de dinero y los medios de pago electrónicos son algo de cada día, pero junto con las ventajas también vinieron las consecuencias, y es que los métodos de fraude se han vuelto más sofisticados, generando pérdidas significativas a los bancos y afectando la confianza de los usuarios. Actividades como el uso no autorizado de tarjetas, la clonación de datos y transacciones fraudulentas requieren el desarrollo de tecnologías avanzadas para la detección temprana y la prevención.

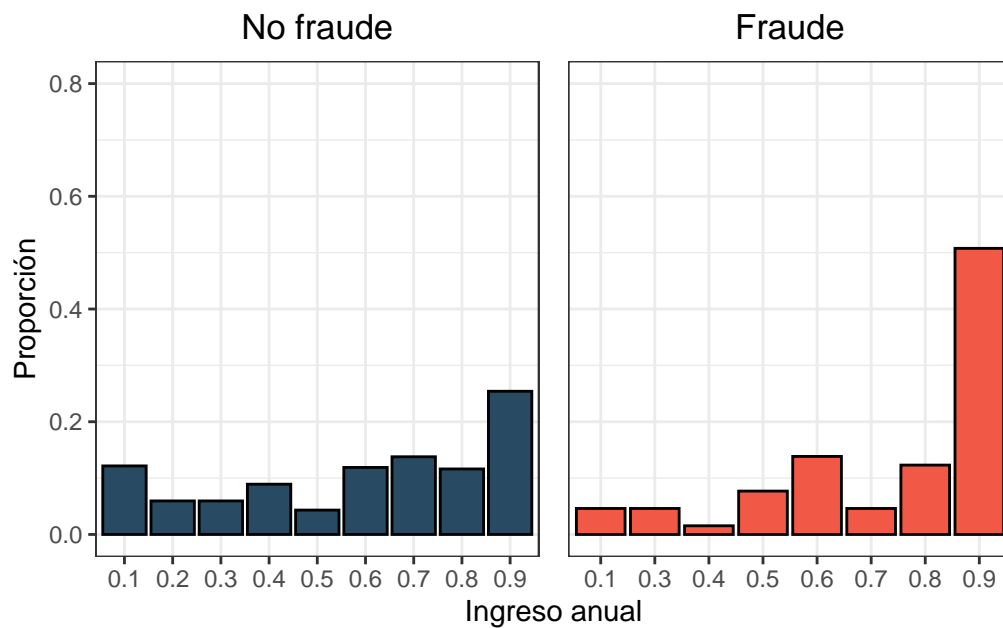
Variables:

- `fraud_bool`: Indicadora de si la transacción fue fraude o no
- `income`: Ingreso anual en cuantiles
- `name_email_similarity`: Similitud del nombre en el email y el nombre del solicitante
- `prev_address_months_count`: Es el número de meses que la persona estuvo viviendo en su locacion anterior
- `current_address_months_count`: Es el número de meses que la persona estuvo viviendo en su locacion actual
- `customer_age`: Edad del cliente en décadas
- `days_since_request`: Días desde la solicitud
- `intended_balcon_amount`: Valor de la transacción inicial para aplicar al credito
- `payment_type`: Tipo del plan de pago
- `zip_count_4w`: Número de aplicaciones con el mismo código postal en las últimas 4 semanas
- `velocity_6h`: Es la velocidad del total de solicitudes de transferencias de la tarjeta en las últimas 6 horas
- `velocity_24h`: Es la velocidad del total de solicitudes de transferencias de la tarjeta en las últimas 24 horas
- `velocity_4w`: Es la velocidad del total de solicitudes de transferencias de la tarjeta en las últimas 4 semanas
- `bank_branch_count_8w`: Número total de solicitudes en la seleccionada rama del banco en las últimas 8 semanas
- `date_of_birth_distinct_emails_4w`: Número de emails de aplicantes con la misma fecha de nacimiento en las últimas 4 semanas
- `employment_status`: Estado de empleo del solicitante
- `credit_risk_score`: Score de riesgo de la aplicación
- `email_is_free`: Tipo del dominio del email del aplicante (email pago o gratis)

- `housing_status`: Estado residencial del aplicante
- `phone_home_valid`: Validez del telefono fijo provisto
- `phone_mobile_valid`: Validez del telefono movil provisto
- `bank_months_count`: Antigüedad de la cuenta anterior en meses
- `has_other_cards`: Indicador de si la persona tiene otra tarjeta en el mismo banco
- `proposed_credit_limit`: Crédito limite propuesto por el aplicante
- `foreign_request`: Indicadora de si la solicitud fue hecha en el mismo pais que el banco
- `source`: Fuente online de la aplicación (Internet / app movil)
- `session_length_in_minutes`: Tiempo de la sesion en la pagina del banco en minutos
- `device_os`: Sistema operativo del dispositivo que hizo la solicitud
- `keep_alive_session`: Indicadora de si el solicitante decidió mantener la sesión iniciada al ingresar
- `device_distinct_emails_8w`: Número de emails distintos en la página del banco desde el mismo dispositivo usado en las últimas 8 semanas
- `device_fraud_count`: Número de solicitudes fraudulentas desde el dispositivo utilizado
- `month`: Mes en el que fue realizada la solicitud

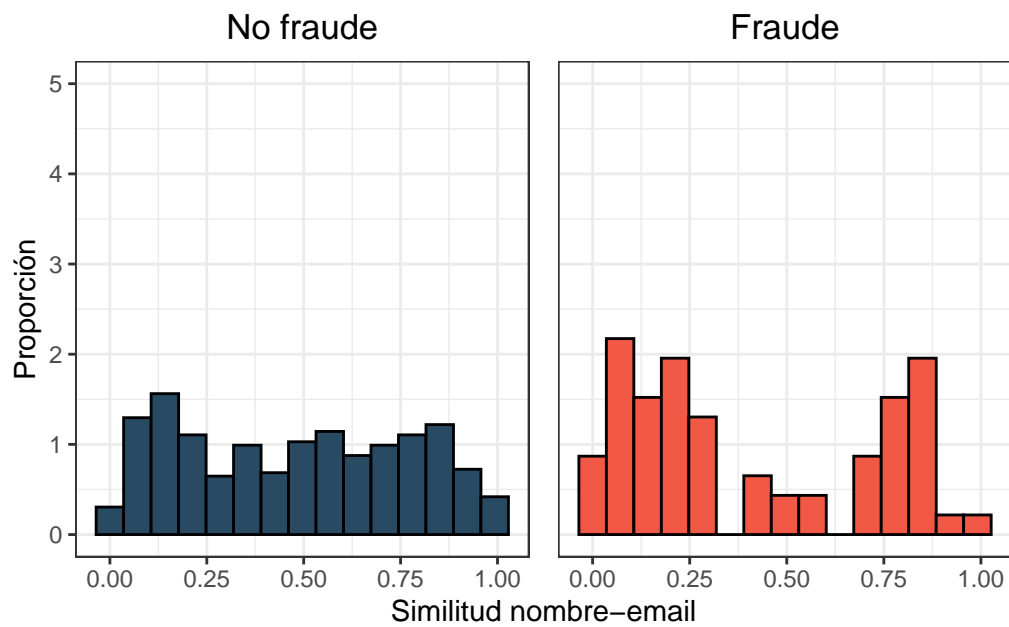
Análisis descriptivo

Figura 1: Distribución del ingreso anual según si la transacción es fraudulenta o no



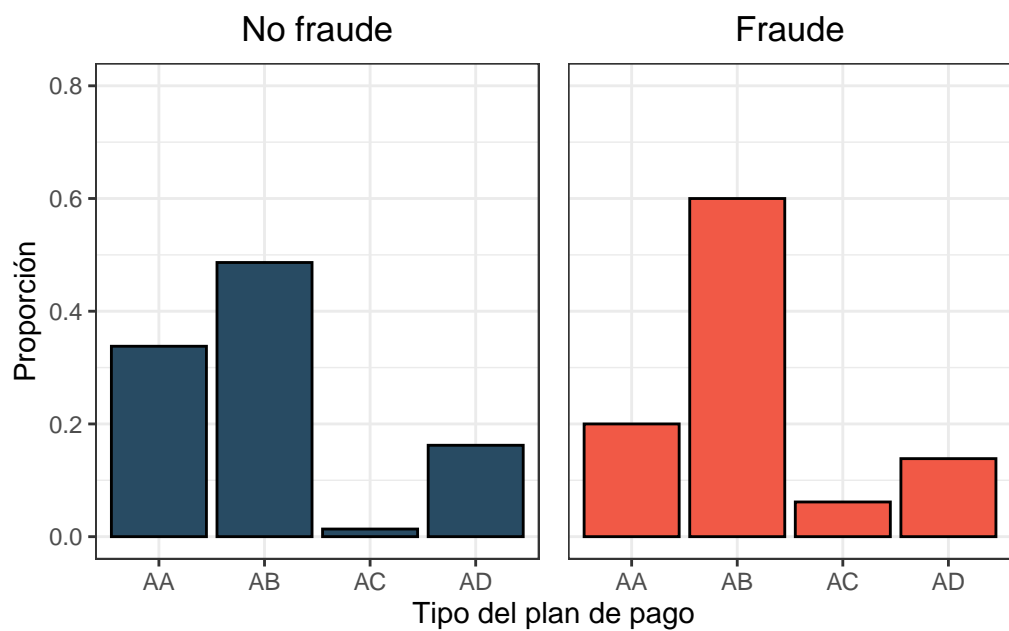
Se puede observar que las personas que cometieron fraude tienden a tener un ingreso anual registrado mayor. La distribución tiene una mayor asimetría a la izquierda.

Figura 2: Distribución del índice de similitud entre en nombre del solicitante y el nombre en el email según si la transacción es fraudulenta o no



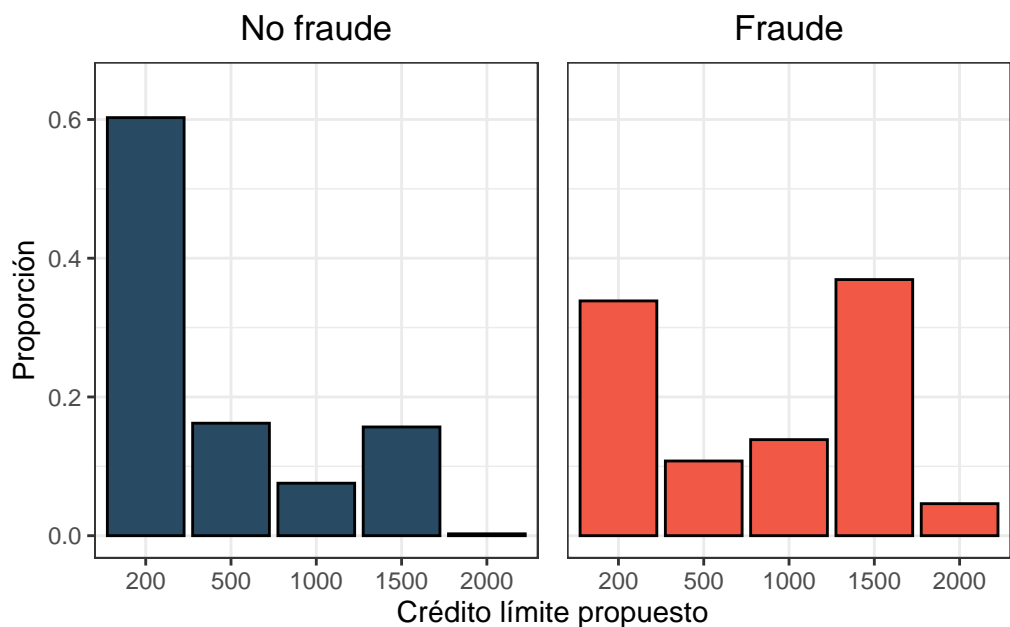
Las transacciones realizadas con emails no muy similares al nombre real de la persona parecen ser más propensas a ser fraudulentas.

Figura 3: Proporción del tipo de pago según si la transacción es fraudulenta o no



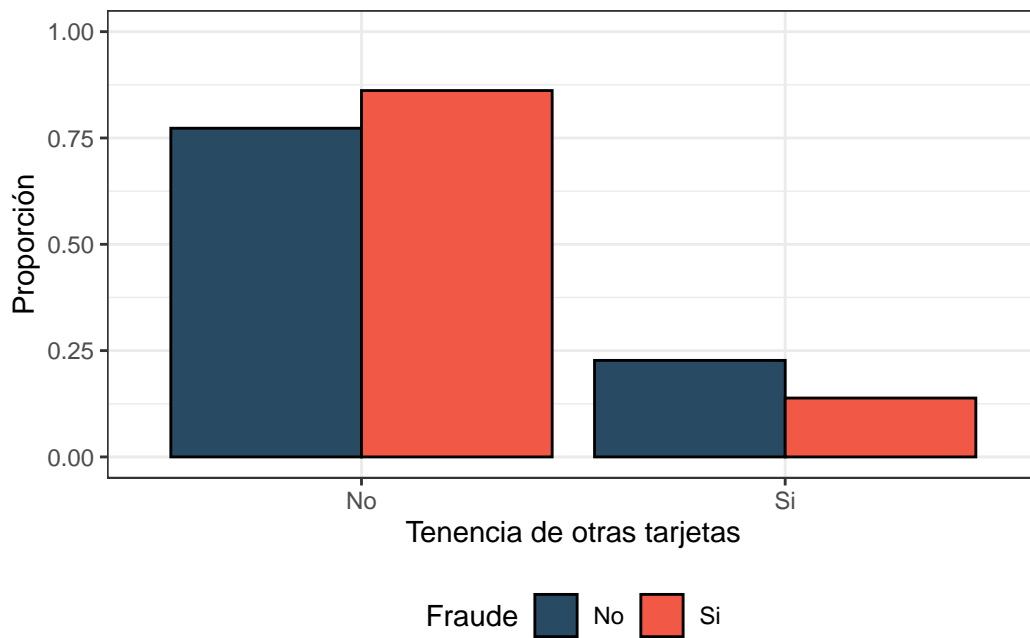
En general, las personas que cometen fraude parecen preferir los métodos de pago “AB” y “AC” por encima del resto, al contrario de las personas que operan de manera legítima que prefieren de igual manera los tipos de pago “AA”, “AB” y “AC”. Se puede notar también que la forma de pago “AE” no es muy popular.

Figura 4: Distribución del límite crediticio propuesto por el solicitante según si la transacción es fraudulenta o no



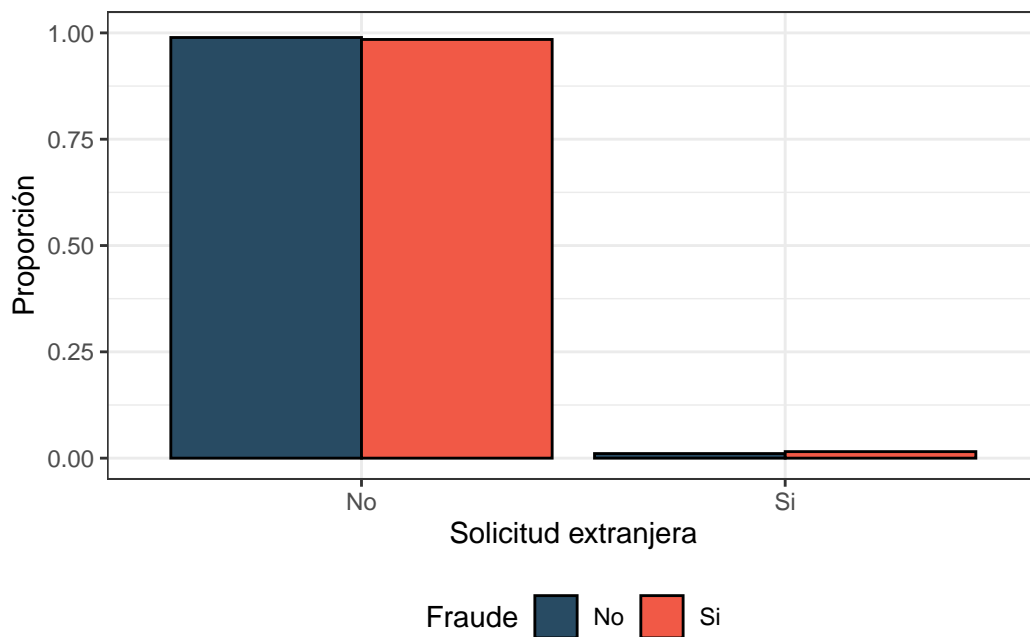
Se puede destacar en este gráfico que las personas que cometen fraude son ligeramente más propensas a pedir créditos más altos.

Figura 5: Proporción de la tenencia de otra tarjeta en el mismo banco según si la transacción es fraudulenta o no



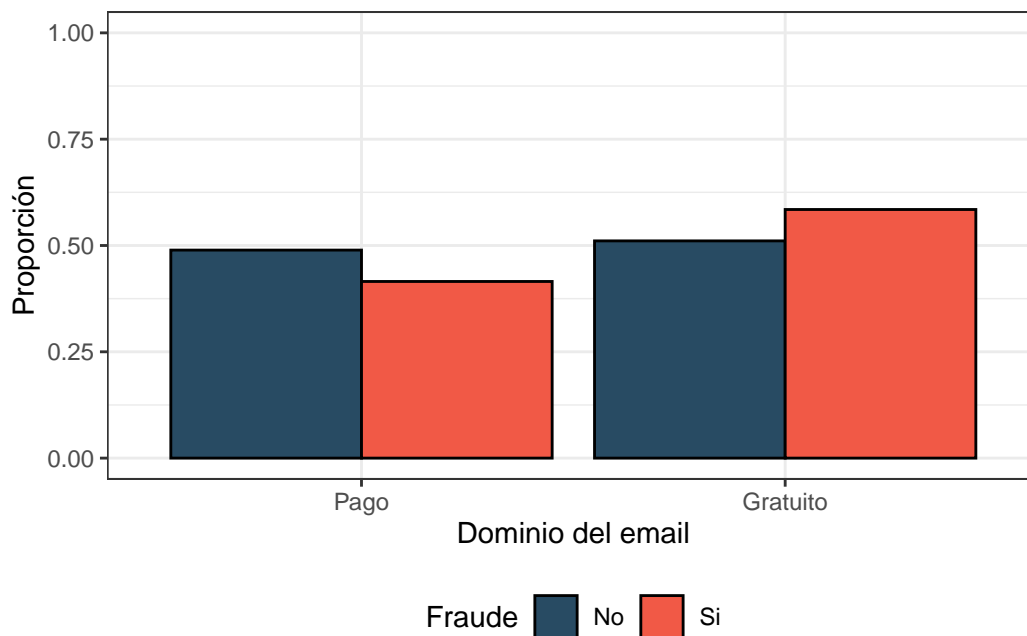
En cuanto a la tenencia de otra tarjeta en el mismo banco, suele ser común no poseer otra, sin embargo las personas que operan de forma legal se inclinan a tener más de una tarjeta un poco más que aquellos que cometen fraude.

Figura 6: Proporción de la locación de la solicitud según si la transacción es fraudulenta o no



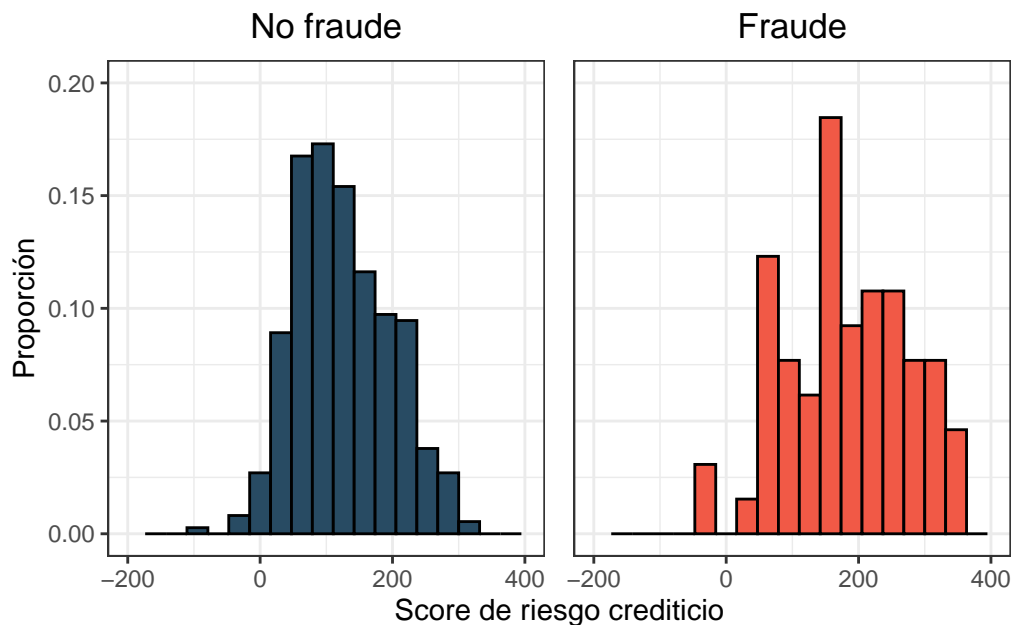
Se puede observar que las personas que cometen fraude, parecen hacer más solicitudes del exterior que las personas que no cometen fraude, aunque la diferencia parece ser sutil.

Figura 7: Proporción del tipo de dominio del email según si la transacción es fraudulenta o no



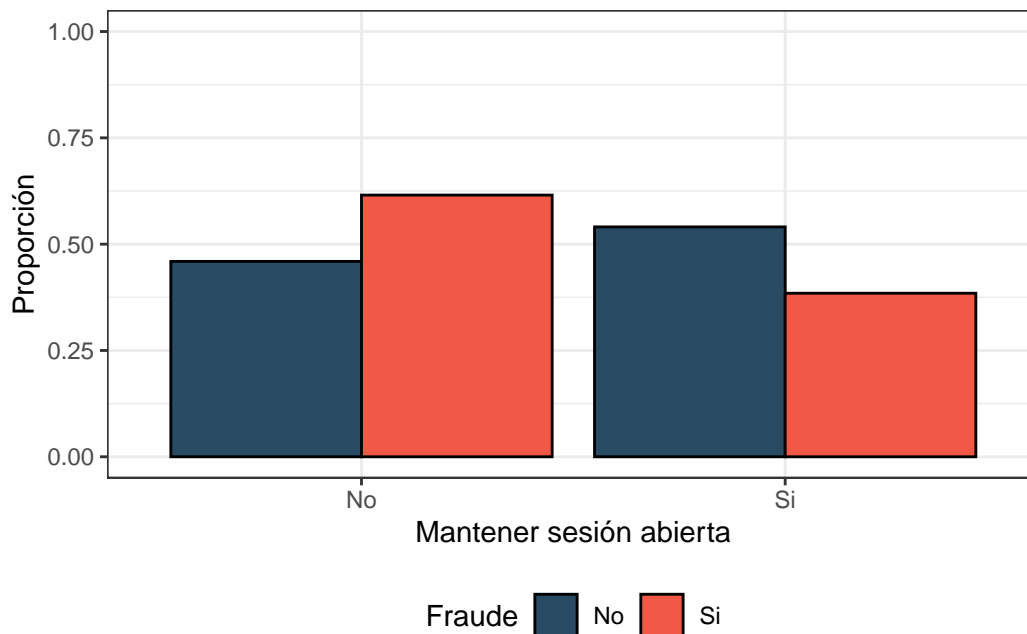
También se puede destacar que las operaciones fraudulentas parecen ser más comunes cuando el dominio del email del solicitante es gratuito que cuando es pago.

Figura 8: Distribución del score de riesgo interno según si la transacción es fraudulenta o no



La distribución del score de riesgo para las personas que cometen fraude es simétrica y centrada alrededor de 200, mientras que la distribución del score de riesgo para las personas que no cometen fraude parece ser más asimétrica y tener una media menor.

Figura 9: Proporción de opciones de inicio de sesión según si la transacción es fraudulenta o no



Por lo general, cuando las transacciones son fraudulentas la persona decide no mantener la sesión abierta en la cuenta del banco en mayor proporción que cuando las transacciones son legítimas.

Modelado

Teniendo todo esto en cuenta buscamos un modelo que ajuste bien a nuestros datos, para esto primero realizamos una seleccion de variables paso a paso, obteniendo el siguiente modelo:

```
Call: glm(formula = fraud_bool ~ income + customer_age + housing_status +
  phone_home_valid + phone_mobile_valid + has_other_cards +
  device_os + device_distinct_emails_8w + proposed_credit_limit_cat,
  family = binomial(link = "logit"), data = data, na.action = na.omit)
```

Coefficients:

(Intercept)	income0.2
-6.41801	-14.96080
income0.3	income0.4
0.12022	-0.67524
income0.5	income0.6

	1.49373		0.58133
	income0.7		income0.8
	-0.30468		1.13962
	income0.9		customer_age
	1.50709		0.02969
	housing_statusBB		housing_statusBC
	-0.96459		-1.37565
	housing_statusBD		housing_statusBE
	-17.39331		-1.96381
	housing_statusBF		phone_home_valid1
	-16.38542		-1.39200
	phone_mobile_valid1		has_other_cards1
	-0.83981		-1.19154
	device_osmacintosh		device_osother
	1.07758		-0.79738
	device_oswindows		device_osx11
	0.76089		-16.19908
	device_distinct_emails_8w	proposed_credit_limit_cat500	
	4.23761		-0.27256
proposed_credit_limit_cat1000	proposed_credit_limit_cat1500		
	1.31513		1.11132
proposed_credit_limit_cat2000			
	3.41638		

Degrees of Freedom: 434 Total (i.e. Null); 408 Residual
Null Deviance: 366.9
Residual Deviance: 226.2 AIC: 280.2

Call:
glm(formula = fraud_bool ~ income + housing_status + phone_home_valid +
phone_mobile_valid + has_other_cards + proposed_credit_limit_cat +
device_os + device_distinct_emails_8w, family = binomial(link = "logit"),
data = data, na.action = na.omit)

Deviance Residuals:

Min	1Q	Median	3Q	Max
-1.88538	-0.44638	-0.22239	-0.05847	2.94083

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-5.2595	1.3811	-3.808	0.000140 ***
income0.2	-15.0770	1270.3069	-0.012	0.990530
income0.3	0.1086	1.0413	0.104	0.916962
income0.4	-0.5542	1.2671	-0.437	0.661876
income0.5	1.5185	1.0481	1.449	0.147377
income0.6	0.6439	0.8506	0.757	0.449065
income0.7	-0.3238	1.0025	-0.323	0.746706

income0.8	1.2991	0.8432	1.541	0.123393	
income0.9	1.5979	0.7363	2.170	0.029993	*
housing_statusBB	-1.0958	0.4930	-2.223	0.026223	*
housing_statusBC	-1.5991	0.4445	-3.598	0.000321	***
housing_statusBD	-17.9029	1624.7479	-0.011	0.991208	
housing_statusBE	-2.4476	0.7678	-3.188	0.001433	**
housing_statusBF	-17.1950	4390.1860	-0.004	0.996875	
phone_home_valid1	-1.2150	0.4202	-2.892	0.003832	**
phone_mobile_valid1	-0.8868	0.5036	-1.761	0.078261	.
has_other_cards1	-1.2152	0.5013	-2.424	0.015356	*
proposed_credit_limit_cat500	-0.2293	0.5509	-0.416	0.677291	
proposed_credit_limit_cat1000	1.2763	0.5748	2.220	0.026387	*
proposed_credit_limit_cat1500	1.0929	0.4403	2.482	0.013052	*
proposed_credit_limit_cat2000	3.4409	1.4713	2.339	0.019355	*
device_osmacintosh	1.1927	0.6765	1.763	0.077886	.
device_osother	-0.9154	0.5223	-1.753	0.079649	.
device_oswindows	0.7807	0.4343	1.797	0.072264	.
device_osx11	-16.3382	3604.6664	-0.005	0.996384	
device_distinct_emails_8w	4.2669	1.0295	4.145	3.4e-05	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 366.89 on 434 degrees of freedom
Residual deviance: 229.89 on 409 degrees of freedom
AIC: 281.89

Number of Fisher Scoring iterations: 17

\$Logit(_i) = _0 \$

Se procede a comprobar el enlace:

Analysis of Deviance Table

Model 1: fraud_bool ~ income + housing_status + phone_home_valid + phone_mobile_valid +
has_other_cards + proposed_credit_limit_cat + device_os +
device_distinct_emails_8w

Model 2: fraud_bool ~ income + housing_status + phone_home_valid + phone_mobile_valid +
has_other_cards + proposed_credit_limit_cat + device_os +
device_distinct_emails_8w + pred.2.logit

	Resid. Df	Resid. Dev	Df	Deviance	Pr(>Chi)
1	409	229.89			
2	408	229.60	1	0.28314	0.5947

Analysis of Deviance Table

```

Model 1: fraud_bool ~ income + housing_status + phone_home_valid + phone_mobile_valid +
  has_other_cards + proposed_credit_limit_cat + device_os +
  device_distinct_emails_8w
Model 2: fraud_bool ~ income + housing_status + phone_home_valid + phone_mobile_valid +
  has_other_cards + proposed_credit_limit_cat + device_os +
  device_distinct_emails_8w + pred.2.probit
Resid. Df Resid. Dev Df Deviance Pr(>Chi)
1      409      230.04
2      408      229.54  1  0.49868  0.4801

```

Analysis of Deviance Table

```

Model 1: fraud_bool ~ income + housing_status + phone_home_valid + phone_mobile_valid +
  has_other_cards + proposed_credit_limit_cat + device_os +
  device_distinct_emails_8w
Model 2: fraud_bool ~ income + housing_status + phone_home_valid + phone_mobile_valid +
  has_other_cards + proposed_credit_limit_cat + device_os +
  device_distinct_emails_8w + pred.2.cloglog
Resid. Df Resid. Dev Df Deviance Pr(>Chi)
1      409      228.43
2      408      228.42  1 0.011481  0.9147

```

Los tres enlaces son apropiados, por lo tanto nos quedamos con el enlace logit por su facilidad en la interpretación.

The Hosmer-Lemeshow goodness-of-fit test

Group	Size	Observed	Expected
1	44	0	0.006753789
2	44	0	0.207370063
3	44	1	0.593732358
4	45	0	1.127851542
5	45	3	1.858645408
6	44	3	3.020927674
7	44	8	5.382100961
8	44	7	9.149775769
9	44	16	16.283726011
10	37	27	27.369116838

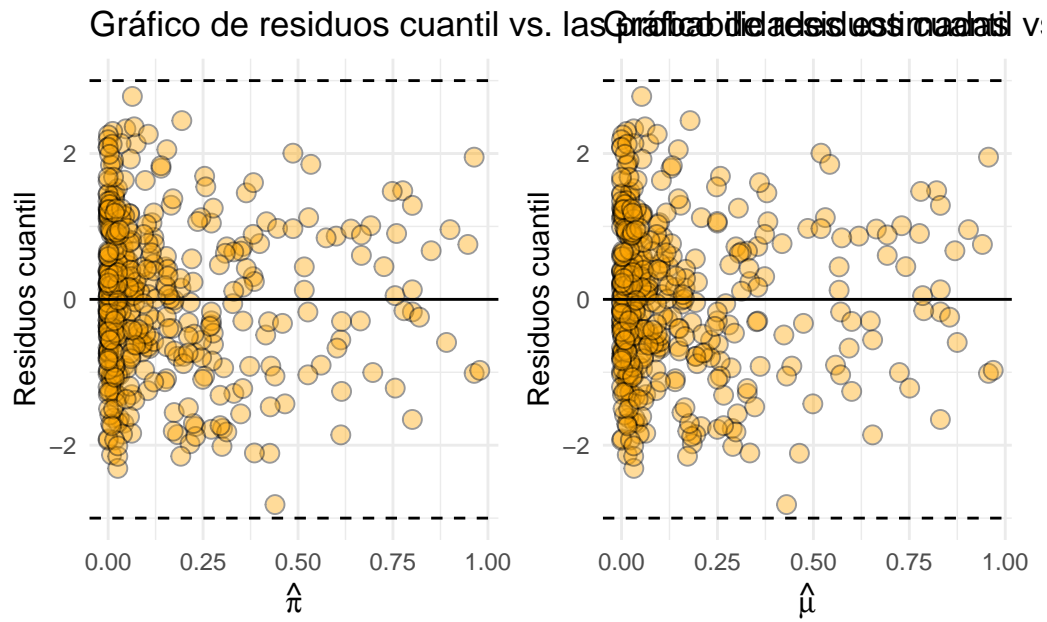
```

Statistic = 4.5005
degrees of freedom = 8
p-value = 0.80938

```

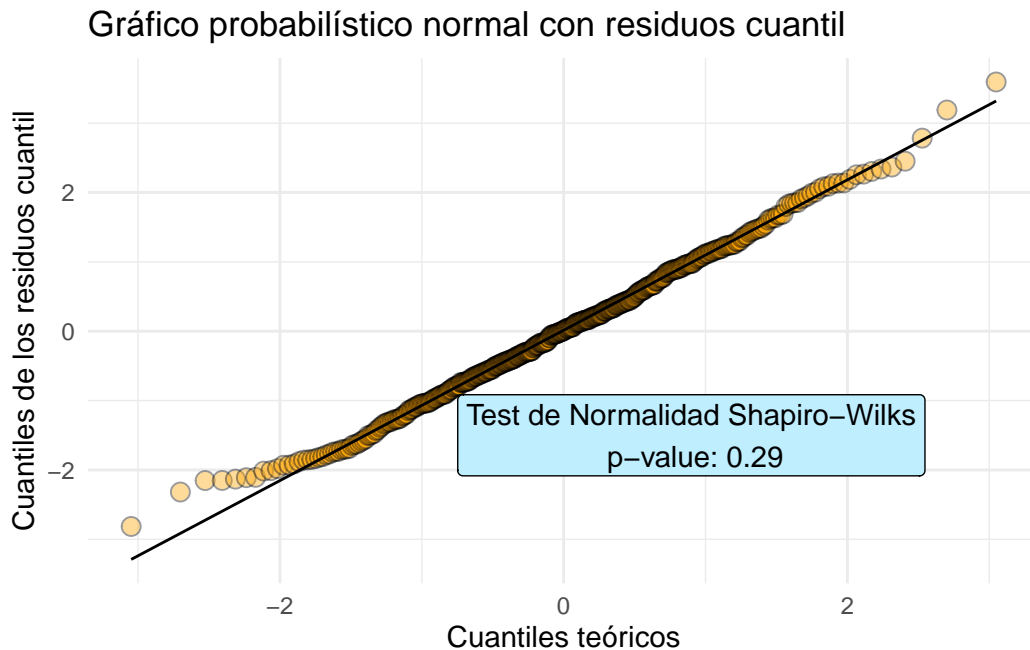
Analisis de residuos

Evaluación de la componente sistemática



Dado que no se ve ningún patrón y ningún punto se escapa de las bandas se puede decir que la componente sistemática seleccionada es adecuada.

Comprobación de la distribución propuesta



Viendo el gráfico y el test de Shapiro-Wilks para la normalidad de los errores, se puede concluir que la elección de la distribución de la variable es correcta.

Interpretaciones

La chance de que un cliente comita fraude dado que tiene otra tarjeta en el mismo banco es un 68% menor que un cliente que no tiene, cuando el resto de las variables permanecen constantes.

Del mismo modo, un cliente que proporciono un telefono fijo válido tiene una chance de cometer fraude un 71% menor que uno que proporciono uno no válido, cuando el resto de las variables permanecen constantes.

A medida que el crédito límite propuesto aumenta en mil unidades monetarias, la chance de que un cliente cometa fraude aumenta un 191%, manteniendo el resto de variables fijas. VARIFICAR SI ES LINEAL O NO

Las características de las personas más propensas a cometer fraude (con una probabilidad estimada del 99% aproximadamente) son:

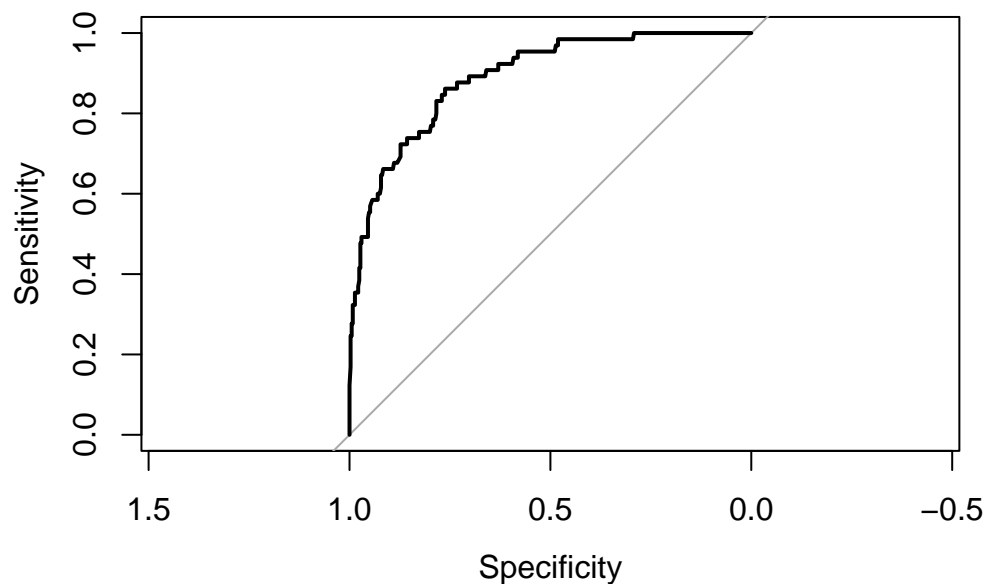
- Ingreso en el decil más alto
- El estado residencial del aplicante es “BA”
- El número de teléfono fijo proporcionado no es válido
- El número de celular proporcionado no es válido
- No tiene otra tarjeta en el mismo banco

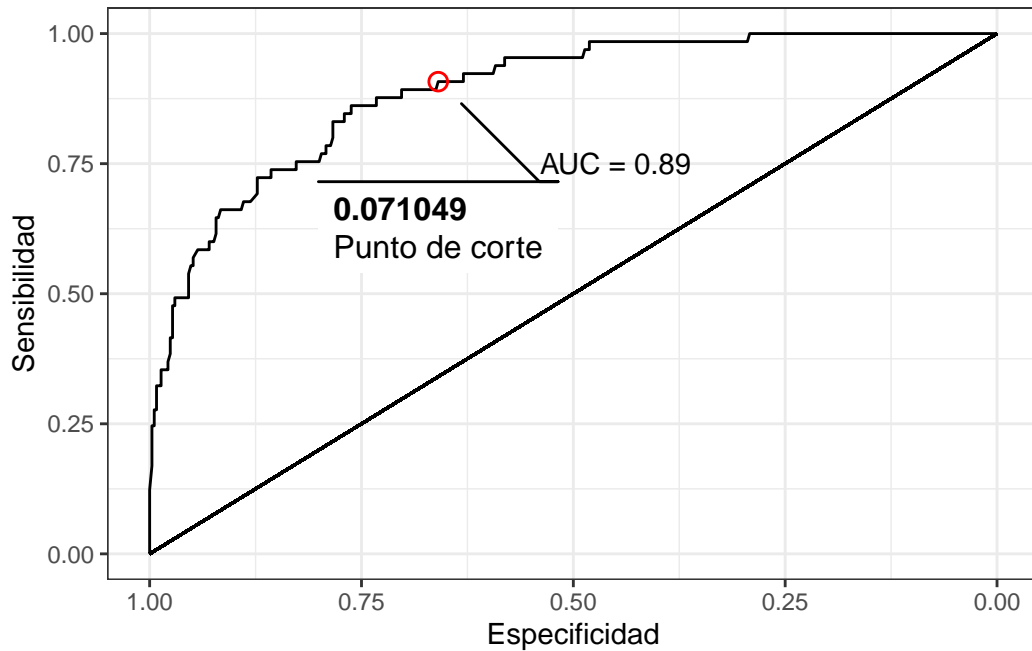
- El límite del crédito propuesto fue 2000 euros
- El sistema operativo usado fue macOS
- Se registraron 2 emails distintos en la página del banco desde el mismo dispositivo, en las últimas 8 semanas

Las características de las personas menos propensas a cometer fraude (con una probabilidad estimada de aproximadamente 0%) son:

- Ingreso en el segundo decil
- El estado residencial del aplicante es “BD”
- El número de teléfono fijo proporcionado es válido
- El número de celular proporcionado es válido
- Tiene otra tarjeta en el mismo banco
- El límite del crédito propuesto fue 500 euros
- El sistema operativo usado fue X11
- Se registró un solo email en la página del banco desde dispositivo utilizado, en las últimas 8 semanas

Capacidad predictiva





Se decidió utilizar el mejor punto de corte que maximice la especificidad garantizando una sensibilidad de al menos un 90%, esto para evitar el máximo número de fraudes posibles sin afectar a los clientes legítimos.

Confusion Matrix and Statistics

	Reference	
Prediction	0	1
0	132	22
1	66	84

Accuracy : 0.7105

95% CI : (0.656, 0.7609)

No Information Rate : 0.6513

P-Value [Acc > NIR] : 0.01659

Kappa : 0.4187

McNemar's Test P-Value : 4.566e-06

Sensitivity : 0.7925

Specificity : 0.6667

Pos Pred Value : 0.5600

Neg Pred Value : 0.8571

Prevalence : 0.3487

Detection Rate : 0.2763

Detection Prevalence : 0.4934

Balanced Accuracy : 0.7296

'Positive' Class : 1