



DETECÇÃO DE FRAUDE DE CARTÃO DE CRÉDITO

JOZIANI MOTA VIEIRA

- **Formação:**

- Graduada em Estatística – UFOP



[linkedin.com/in/joziani-mota](https://www.linkedin.com/in/joziani-mota)



github.com/Joziani



jozianivieira@outlookl.com



SOBRE O PROBLEMA

- Como, cada vez menos as pessoas estão utilizando dinheiro físico, e optando mais por utilizar cartões, é de suma importância que as empresas que disponibilizam crédito tenham uma forma de verificar a probabilidade de um cliente fraudar o cartão.
- Assim podemos criar modelos para ajudar essas empresas nessa busca.

SOBRE OS DADOS

- O conjunto de dados contém transações feitas com cartões de crédito em setembro de 2013 por titulares de cartões europeus.
- Variáveis:
 - **Class** é a variável de resposta do modelo, assumindo valor 1 em caso de fraude e 0 caso contrário;
 - de **V1 à V28** são os componentes principais obtidos através de PCA;
 - **Time** contém os segundos decorridos entre cada transação e a primeira transação no conjunto de dados;
 - **Amount** é o valor da transação.

ANÁLISE EXPLORATÓRIA

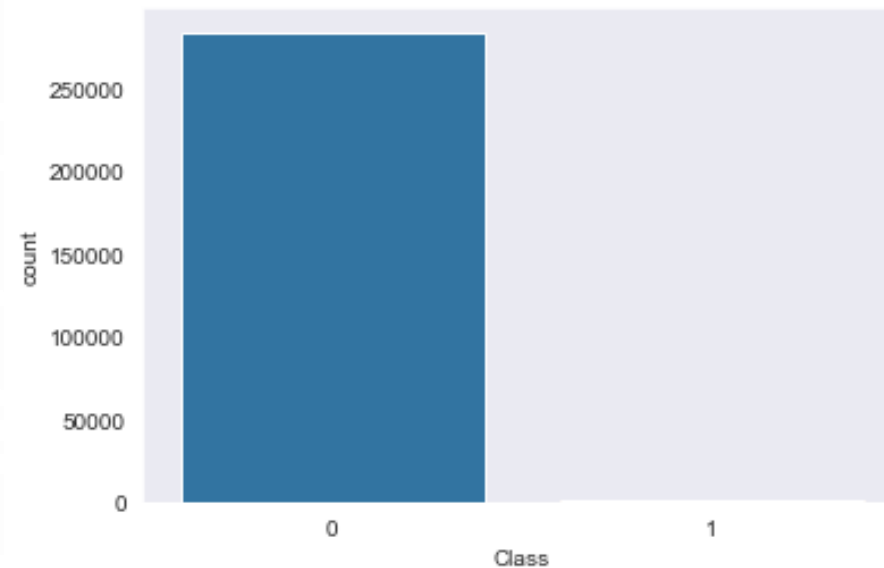
- Foi verificado, inicialmente, que a mostra não possuía dados faltantes.
- Após essa detecção, foi feita uma análise descritiva das variáveis.
 - Não foi incluídos os componentes principais para esta análise.

ANÁLISE DESCRITIVA

VARIÁVEL RESPOSTA

- Percebe-se que temos um desbalanceamento nos dados, já que em 99,83% não houve fraude.

Fraude	N	%
Não houve	284315	99,83
Houve	492	0,17



ANÁLISE DESCRITIVA

VARIÁVEIS NUMÉRICAS

- A média de segundos entre as transações e a primeira transação foi de 94.813,86, com desvio padrão de 47.488,15.
- A média do valor da transação foi de 88,35, com desvio padrão de 250,12.

Variáveis	N	Média	D.P.	Mínimo	1º Quartil	Mediana	3º Quartil	Máximo
Time	284807	94813,86	47488,15	0,00	54201,50	84692,00	139320,50	172792,00
Amount	284807	88,35	250,12	0,00	5,60	22,00	77,17	25691,16

ANÁLISE DE CORRELAÇÃO

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13	V14	V15	V16	V17	V18	V19	V20	V21	V22	V23	V24	V25	V26	V27	V28	Amount	Class
Time	1.00	0.12	-0.01	-0.42	-0.11	0.17	-0.06	0.08	-0.04	-0.01	0.03	-0.25	0.12	-0.07	-0.10	-0.18	0.01	-0.07	0.09	0.03	-0.05	0.04	0.14	0.05	-0.02	-0.23	-0.04	-0.01	-0.01	-0.01	-0.01
V1	0.12	1.00	0.00	-0.00	-0.00	0.00	-0.00	-0.00	-0.00	-0.00	0.00	0.00	0.00	-0.00	-0.00	0.00	0.00	-0.00	0.00	0.00	0.00	-0.00	-0.00	0.00	-0.00	-0.00	0.00	0.00	-0.23	-0.10	
V2	-0.01	0.00	1.00	0.00	-0.00	0.00	0.00	0.00	-0.00	0.00	-0.00	0.00	-0.00	0.00	-0.00	0.00	-0.00	0.00	-0.00	0.00	0.00	0.00	0.00	0.00	0.00	-0.00	0.00	-0.00	-0.53	0.09	
V3	-0.42	-0.00	0.00	1.00	0.00	-0.00	0.00	0.00	-0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	-0.00	0.00	-0.00	-0.00	0.00	-0.00	0.00	0.00	-0.21	-0.19	
V4	-0.11	-0.00	-0.00	0.00	1.00	-0.00	-0.00	-0.00	0.00	0.00	0.00	0.00	-0.00	0.00	0.00	0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	0.00	0.00	0.00	-0.00	0.00	-0.00	0.10	0.13
V5	0.17	0.00	0.00	-0.00	-0.00	1.00	0.00	0.00	0.00	0.00	-0.00	0.00	0.00	0.00	0.00	-0.00	0.00	0.00	0.00	-0.00	-0.00	-0.00	0.00	-0.00	-0.00	0.00	0.00	0.00	-0.00	-0.39	-0.09
V6	-0.06	-0.00	0.00	0.00	-0.00	0.00	1.00	0.00	-0.00	0.00	0.00	0.00	0.00	-0.00	0.00	-0.00	0.00	0.00	0.00	-0.00	-0.00	0.00	-0.00	0.00	-0.00	0.00	-0.00	0.00	0.22	-0.04	
V7	0.08	-0.00	0.00	0.00	-0.00	0.00	0.00	1.00	0.00	0.00	-0.00	0.00	-0.00	0.00	0.00	-0.00	0.00	0.00	0.00	-0.00	0.00	-0.00	-0.00	-0.00	0.00	-0.00	-0.00	-0.00	-0.00	0.40	-0.19
V8	-0.04	-0.00	-0.00	-0.00	0.00	0.00	-0.00	0.00	1.00	0.00	-0.00	0.00	0.00	-0.00	-0.00	0.00	-0.00	-0.00	-0.00	-0.00	0.00	0.00	0.00	0.00	-0.00	-0.00	-0.00	0.00	-0.00	-0.10	0.02
V9	-0.01	-0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	-0.00	0.00	-0.00	0.00	0.00	-0.00	-0.00	0.00	0.00	-0.00	-0.00	0.00	-0.00	-0.00	-0.00	0.00	-0.00	-0.00	0.00	-0.04	-0.10
V10	0.03	0.00	-0.00	0.00	0.00	-0.00	0.00	-0.00	-0.00	-0.00	1.00	-0.00	0.00	-0.00	0.00	0.00	0.00	0.00	0.00	0.00	-0.00	0.00	-0.00	0.00	-0.00	-0.00	-0.00	-0.00	0.00	-0.10	-0.22
V11	-0.25	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	-0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	-0.00	-0.00	-0.00	0.00	-0.00	0.00	-0.00	-0.00	-0.00	0.00	0.00	0.15
V12	0.12	0.00	-0.00	0.00	-0.00	0.00	0.00	-0.00	0.00	-0.00	0.00	0.00	1.00	-0.00	0.00	-0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	-0.00	-0.00	-0.00	0.00	-0.01	-0.26
V13	-0.07	-0.00	0.00	0.00	0.00	0.00	-0.00	0.00	-0.00	0.00	-0.00	0.00	-0.00	1.00	0.00	0.00	0.00	0.00	0.00	-0.00	0.00	0.00	0.00	-0.00	-0.00	-0.00	-0.00	0.00	0.00	0.01	-0.00
V14	-0.10	-0.00	-0.00	0.00	0.00	0.00	0.00	0.00	-0.00	0.00	0.00	0.00	0.00	0.00	1.00	-0.00	-0.00	0.00	0.00	0.00	-0.00	-0.00	0.00	0.00	0.00	-0.00	-0.00	0.00	0.00	0.03	-0.30
V15	-0.18	0.00	-0.00	0.00	0.00	-0.00	-0.00	-0.00	0.00	-0.00	0.00	0.00	-0.00	0.00	-0.00	1.00	0.00	0.00	0.00	-0.00	0.00	-0.00	-0.00	-0.00	0.00	0.00	-0.00	-0.00	-0.00	-0.00	-0.00
V16	0.01	0.00	0.00	0.00	-0.00	0.00	0.00	0.00	-0.00	-0.00	0.00	0.00	0.00	0.00	-0.00	0.00	1.00	0.00	-0.00	0.00	0.00	-0.00	-0.00	0.00	-0.00	-0.00	-0.00	0.00	0.00	-0.00	-0.20
V17	-0.07	-0.00	-0.00	0.00	-0.00	0.00	0.00	0.00	-0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	-0.00	-0.00	-0.00	-0.00	-0.00	0.00	-0.00	0.00	0.00	0.00	-0.00	0.01	-0.33
V18	0.09	0.00	0.00	0.00	-0.00	0.00	0.00	0.00	-0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	-0.00	-0.00	1.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	0.00	0.00	0.00	0.04	-0.11
V19	0.03	0.00	-0.00	0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	0.00	-0.00	0.00	-0.00	0.00	-0.00	0.00	-0.00	-0.00	1.00	0.00	0.00	-0.00	0.00	0.00	0.00	-0.00	-0.00	-0.06	0.03	0.00
V20	-0.05	0.00	0.00	-0.00	-0.00	-0.00	-0.00	0.00	0.00	0.00	-0.00	-0.00	0.00	0.00	-0.00	0.00	-0.00	-0.00	0.00	0.00	1.00	-0.00	0.00	0.00	0.00	0.00	-0.00	-0.00	-0.00	0.34	0.02
V21	0.04	-0.00	-0.00	0.00	-0.00	-0.00	0.00	-0.00	0.00	0.00	0.00	-0.00	0.00	0.00	-0.00	-0.00	-0.00	-0.00	-0.00	0.00	-0.00	1.00	0.00	0.00	0.00	-0.00	-0.00	0.00	0.11	0.04	
V22	0.14	-0.00	0.00	-0.00	-0.00	0.00	-0.00	-0.00	0.00	-0.00	-0.00	0.00	0.00	0.00	0.00	-0.00	-0.00	-0.00	-0.00	-0.00	0.00	0.00	1.00	-0.00	0.00	-0.00	-0.00	0.00	-0.06	0.00	
V23	0.05	0.00	0.00	-0.00	0.00	-0.00	0.00	-0.00	0.00	-0.00	0.00	-0.00	0.00	-0.00	0.00	0.00	0.00	0.00	-0.00	0.00	0.00	0.00	-0.00	1.00	0.00	-0.00	0.00	0.00	-0.11	-0.00	
V24	-0.02	-0.00	0.00	0.00	0.00	-0.00	-0.00	0.00	-0.00	-0.00	-0.00	0.00	0.00	-0.00	0.00	-0.00	-0.00	-0.00	-0.00	0.00	0.00	0.00	0.00	0.00	1.00	0.00	0.00	-0.00	0.01	-0.01	
V25	-0.23	-0.00	-0.00	-0.00	0.00	0.00	0.00	-0.00	-0.00	0.00	-0.00	-0.00	-0.00	-0.00	-0.00	0.00	-0.00	0.00	-0.00	0.00	0.00	-0.00	-0.00	0.00	1.00	0.00	-0.00	-0.00	-0.05	0.00	
V26	-0.04	-0.00	0.00	-0.00	-0.00	0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	0.00	0.00	-0.00	0.00	0.00	-0.00	-0.00	0.00	0.00	0.00	1.00	-0.00	-0.00	-0.00	0.00	
V27	-0.01	0.00	-0.00	0.00	0.00	0.00	-0.00	-0.00	0.00	-0.00	-0.00	-0.00	-0.00	-0.00	0.00	-0.00	0.00	0.00	0.00	-0.00	-0.00	-0.00	0.00	0.00	-0.00	-0.00	1.00	-0.00	0.03	0.02	
V28	-0.01	0.00	-0.00	0.00	-0.00	-0.00	0.00	-0.00	-0.00	0.00	0.00	-0.00	0.00	0.00	0.00	-0.00	0.00	-0.00	0.00	-0.00	-0.00	0.00	-0.00	0.00	-0.00	-0.00	-0.00	1.00	0.01	0.01	
Amount	-0.01	-0.23	-0.53	-0.21	0.10	-0.39	0.22	0.40	-0.10	-0.04	-0.10	0.00	-0.01	0.01	0.03	-0.00	-0.00	0.01	0.04	-0.06	0.34	0.11	-0.06	-0.11	0.01	-0.05	-0.00	0.03	0.01	1.00	0.01
Class	-0.01	-0.10	0.09	-0.19	0.13	-0.09	-0.04	-0.19	0.02	-0.10	-0.22	0.15	-0.26	-0.00	-0.30	-0.00	-0.20	-0.33	-0.11	0.03	0.02	0.04	0.00	-0.00	-0.01	0.00	0.00	0.02	0.01	0.01	1.00

- Observa-se que não há correlações fortes entre as variáveis

MÉTODO HOLDOUT

- Foi aplicado o método holdout para calculo de métricas do modelo.
- Que consiste em separar o banco de dados para treino e teste dos modelos.
- Neste projeto o banco de dados para treinar os modelos foi composto por 70% e o banco de dados para testar os modelos foi composto por 30% dos dados

DADOS DESBALANCEADOS

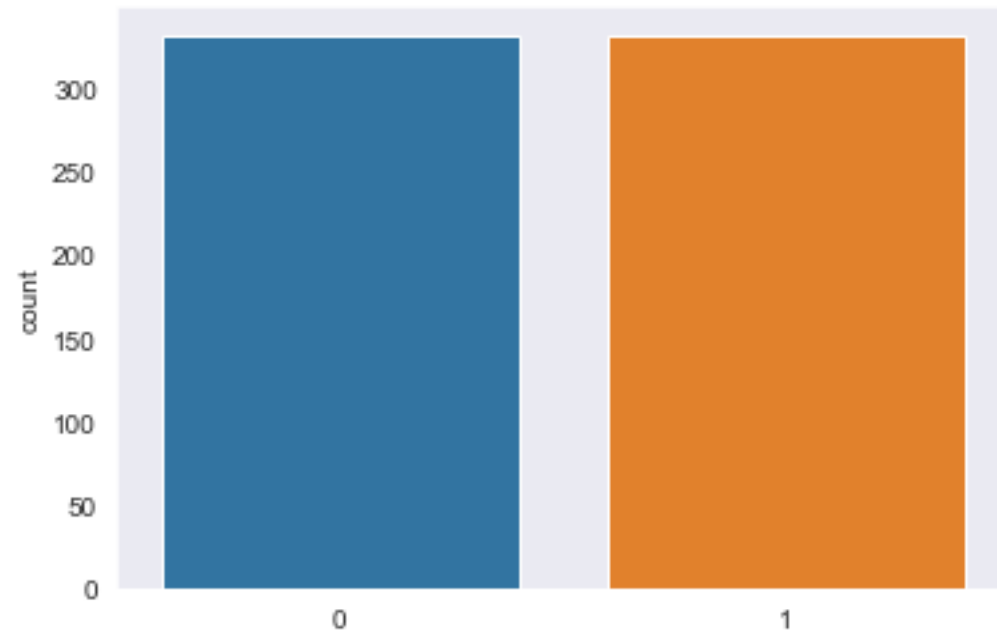
- Como foi observado, temos uma grande quantidade de dados sem fraude, em comparação à dados onde houve fraude.
- O que indica que os modelos irão responder muito bem para detectar quando não houver fraude, mas terá um desempenho inferior para detectar quando houver fraude.

UNDER-SAMPLING

- O método under-sampling consiste em reduzir o desbalanceamento dos dados, focando na classe majoritária, quando não houve fraude. Ou seja, elimina aleatoriamente entradas da classe desta classe.
- Foi optado pelo método under-sampling, já que o método over-sampling era inviável com uma diferença tão grande entre as classes.
- A desvantagem é que perdemos dados com esta abordagem, mas como o objetivo principal é prever fraudes, o método under-sampling se aplica melhor aqui.
- Será feita uma comparação entre utilizar o método e os dados originais.

UNDER-SAMPLING

- Depois de aplicar o método under-sampling, obtemos dados balanceados para a variável resposta:



MODELOS

- Optou-se por utilizar três algoritmos de machine learning, para fins de comparação, sendo eles:
 - Regressão Logística;
 - RandomForest;
 - Support Vector Machines - SVM.

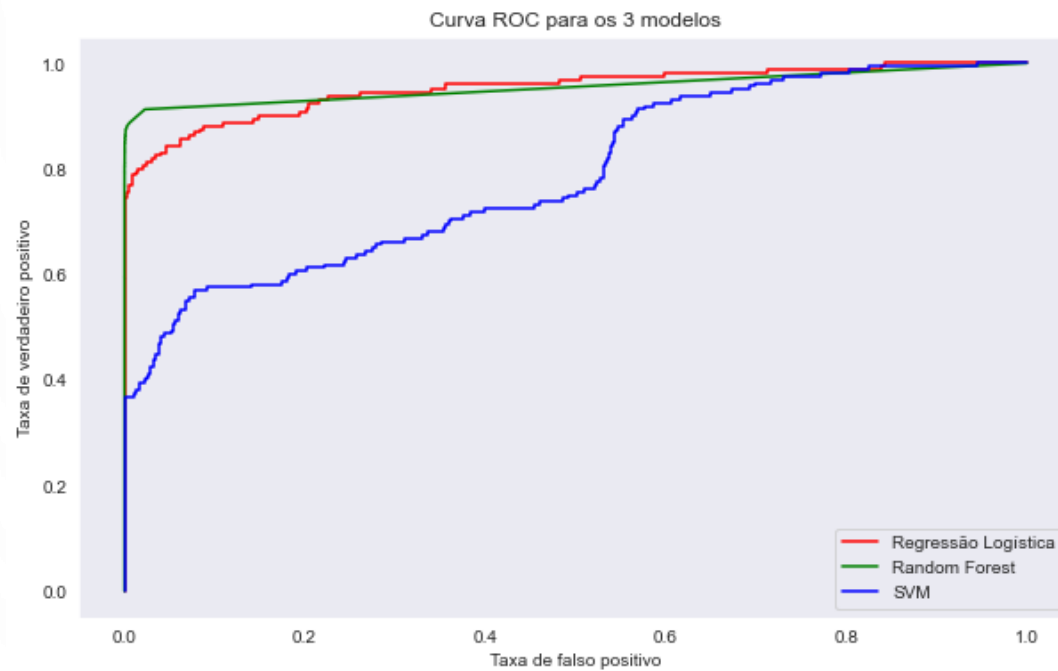
RESULTADOS

MÉTRICAS

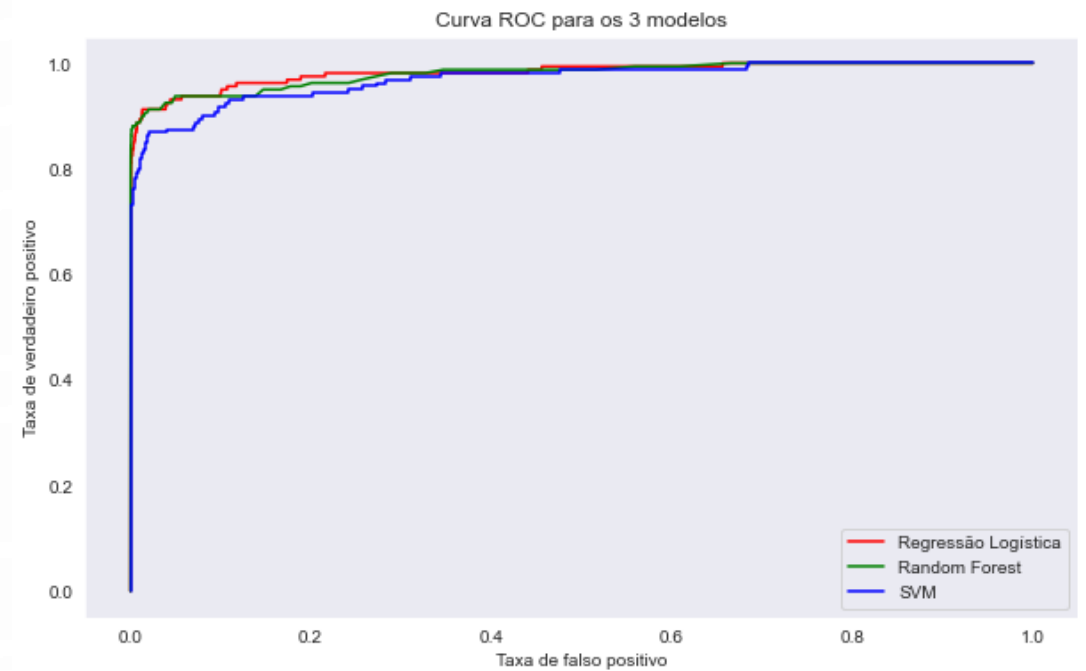
Métricas	Acurácia	Sensibilidade	Especificidade	AUC
Regressão Logística	1,00	1,00	0,61	0,95
Regressão - Under-sampling	0,96	0,96	0,92	0,98
RandomForest	1,00	1,00	0,81	0,96
RandomForest - Under-sampling	0,98	0,98	0,91	0,98
SVM	1,00	1,00	0,36	0,79
SVM - Under-sampling	0,99	0,99	0,80	0,97

RESULTADOS

CURVA ROC



Dados originais



Método under-sampling

CONCLUSÕES

- Vimos que os modelos com os dados originais tiveram uma acurácia ótima, mas se olharmos, por exemplo, para a especificidade não se saíram tão bem.
- Já tivemos uma melhora na especificidade com a utilização under-sampling.
- Por esse motivo, é importante, quando verificar ajuste e desempenho de modelos, olhar para mais de uma métrica.
- Como a regressão logística teve um desempenho bem próximo aos outros modelos, é interessante à escolher, já que com esse modelo tem-se mais informações, não apenas se terá fraude ou não.

A decorative graphic on the left side of the slide, consisting of a network of white lines and small circles on a blue gradient background, resembling a circuit board or a stylized tree structure.

OBRIGADA!