

Credit Card Fraud Detection



Jéssica Ramos

Junho de 2021

O problema

- Classificar transações de cartão de crédito como de alto ou baixo risco de fraude.
- Composição da base:
 - Transações realizadas em Setembro de 2013;
 - 284.807 transações, das quais 492 são fraudes (0.172% das transações);
 - 28 componentes principais gerados com as variáveis originais (não disponibilizadas);
 - Valor da transação, momento da transação e identificação se a transação é fraudulenta.
- **Proposta:** Construção de modelo preditivo de classificação que atribua uma probabilidade de fraude para cada transação.

Estratégia de análise

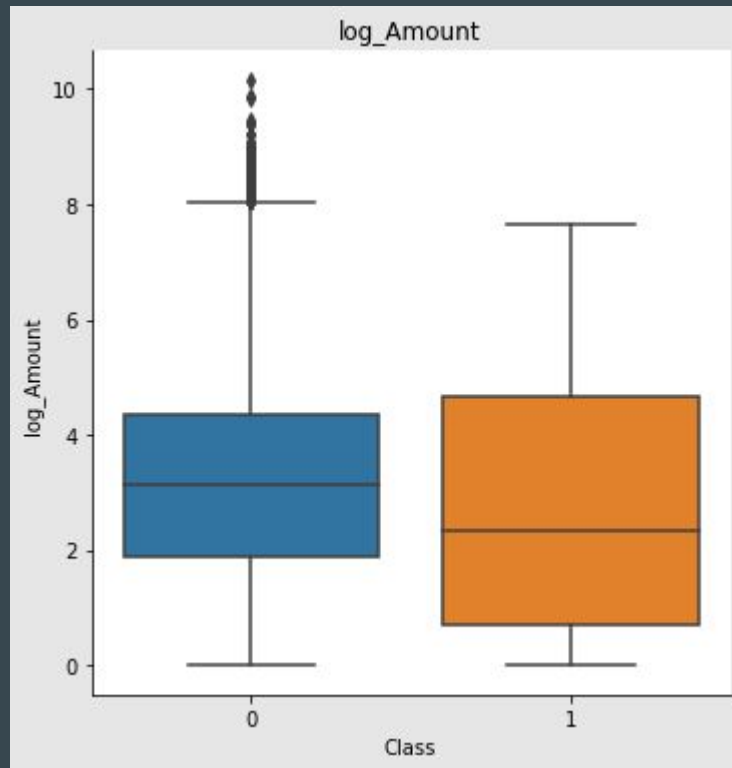
- Análise exploratória
 - Preparação da base
 - Treino/Teste
 - Desbalanceamento
 - Ajuste de algoritmos
 - Modelo Logístico
 - Random Forest
 - Gradient Boosting
 - Resultados finais
-

Análise Exploratória

A maior parte das variáveis na base não tem interpretação direta, pois são resultado de PCA. A vantagem é que são não correlacionadas.

Foi incluída a variável 'Amount' transformada pela função log, devido à distribuição muito assimétrica da variável original.

Em muitos casos não há separação clara quando a análise é bivariada.

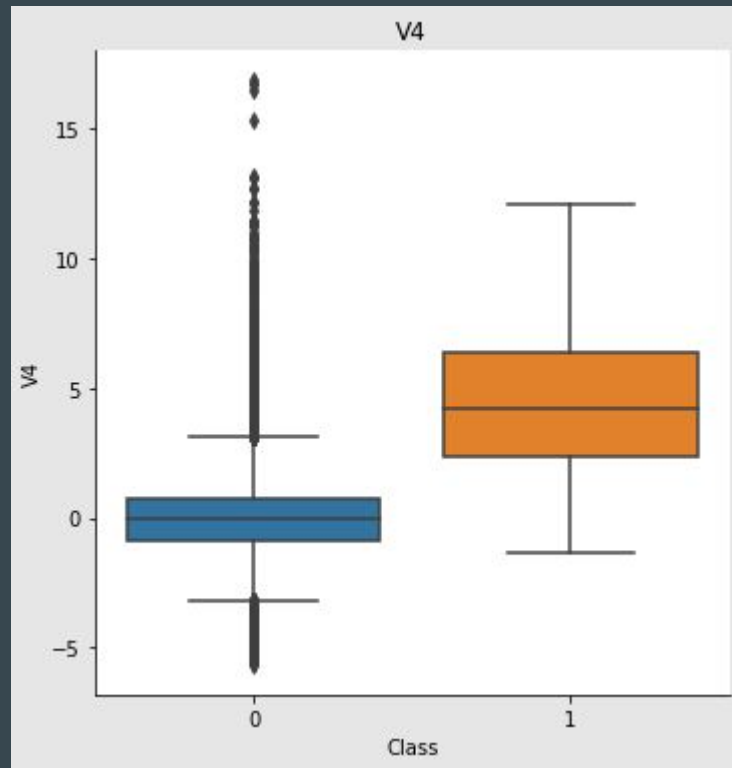


Análise Exploratória

A maior parte das variáveis na base não tem interpretação direta, pois são resultado de PCA. A vantagem é que são não correlacionadas.

Foi incluída a variável 'Amount' transformada pela função log, devido à distribuição muito assimétrica da variável original.

Em muitos casos não há separação clara quando a análise é bivariada.



Preparação da Base

Separação Treino/Teste

A base foi separada com base na coluna 'Time': as transações mais recentes foram colocadas na base de treino (out-of-time).

Treino:

199.368 observações (70%)

0,193% de fraudes

Teste:

85.439 observações (30%)

0,126% de fraudes

Desbalanceamento

O desbalanceamento pode prejudicar o desempenho do modelo na classe positiva.

Foi utilizado o método **SMOTE** para criar observações sintéticas na base de treino.

Composição final:

397.968 observações (+99.6%)

50% de fraudes

Ajuste de Algoritmos

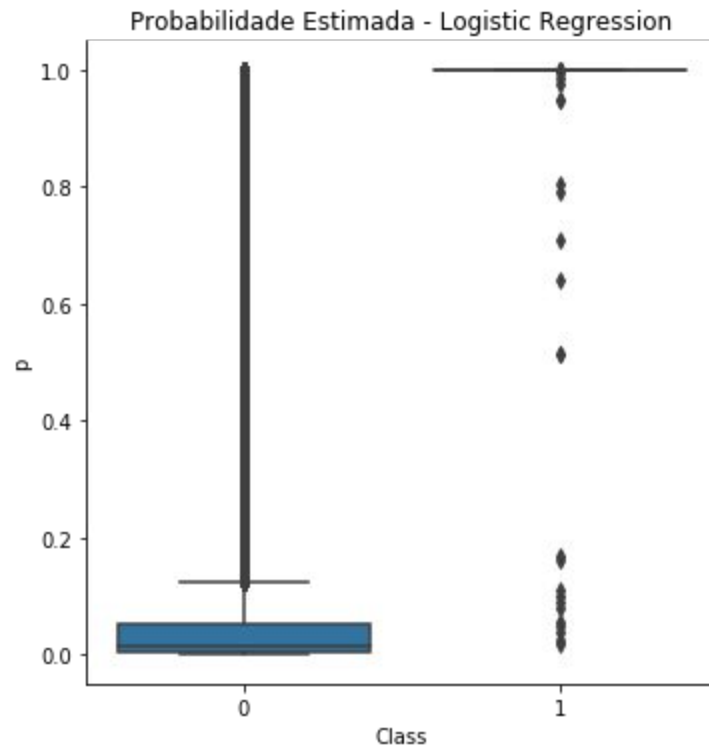
Modelo Logístico

Regularização L2

Lambda = 1

Threshold = 0,5

Otimização de lambda por cross-validation mostrou pouca variação no recall, então foi mantido o valor padrão.



Ajuste de Algoritmos

Modelo Logístico

Regularização L2

Lambda = 1

Threshold = 0,5

Otimização de lambda por cross-validation mostrou pouca variação no recall.

	Treino	Teste
Recall	92,8%	89,8%
Specificity	97,8%	97,4%
Precision	97,6%	4,1%

Resultado: Um modelo bem simples e o corte usual de threshold deu resultados muito bons.

Ajuste de Algoritmos

Random Forest

Random Grid Search usando uma base de validação (20% da base de treino).

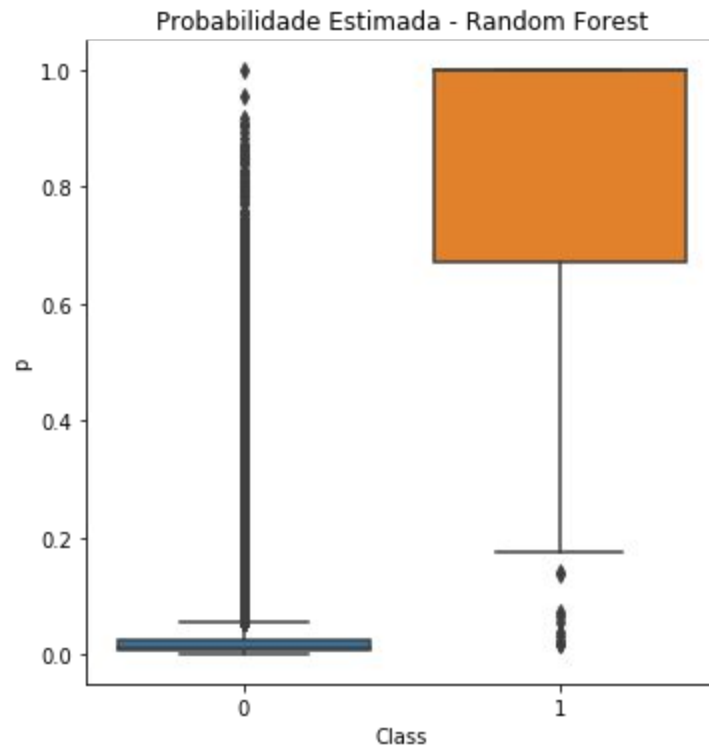
Estimadores: 500

Profundidade Máxima: 8

Mínimo de obs. por folha: 5

Perc. de features: 70%

Threshold: 0,2



Ajuste de Algoritmos

Random Forest

Random Grid Search usando uma base de validação (20% da base de treino).

Estimadores: 500

Profundidade Máxima: 8

Mínimo de obs. por folha: 5

Perc. de features: 70%

Threshold: 0,2

	Treino	Teste
Recall	99,3%	88,0%
Specificity	97,1%	97,5%
Precision	97,1%	4,3%

Resultado: Com um nível de specificity semelhante ao modelo logístico, o recall é menor em 1,8 p.p.

Ajuste de Algoritmos

Gradient Boosting

Random Grid Search usando uma base de validação (20% da base de treino).

Taxa de aprendizado: 0,05

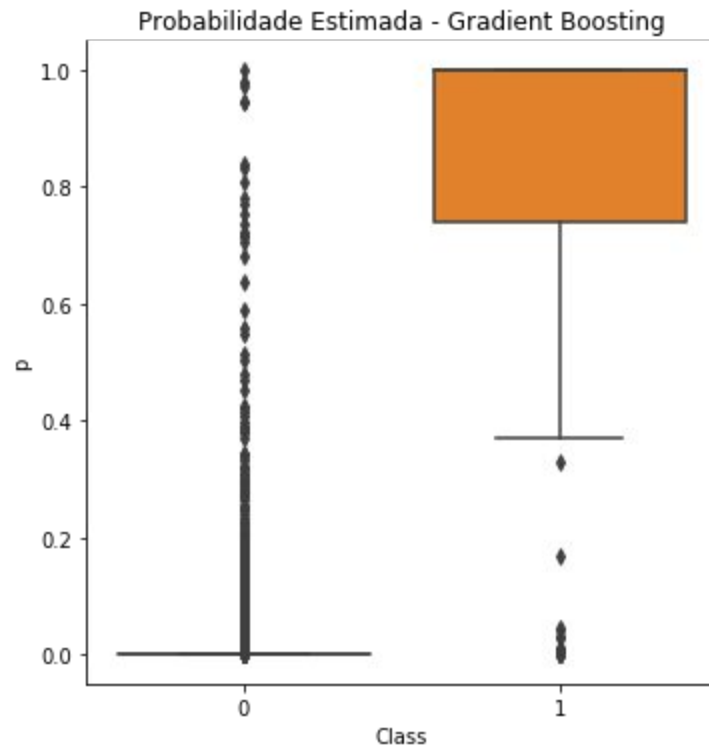
Estimadores: 300

Profundidade Máxima: 7

Mínimo de obs. por folha: 6

Perc. de features: 90%

Threshold: 0,005



Ajuste de Algoritmos

Gradient Boosting

Random Grid Search usando uma base de validação (20% da base de treino).

Taxa de aprendizado: 0,05

Estimadores: 300

Profundidade Máxima: 7

Mínimo de obs. por folha: 6

Perc. de features: 90%

Threshold: 0,005

	Treino	Teste
Recall	100,0%	87,0%
Specificity	97,7%	97,8%
Precision	97,7%	4,9%

Resultado: O desempenho do recall é bem menor em teste do que em trein O resultado não supera o modelo logístico.

Resultados Finais

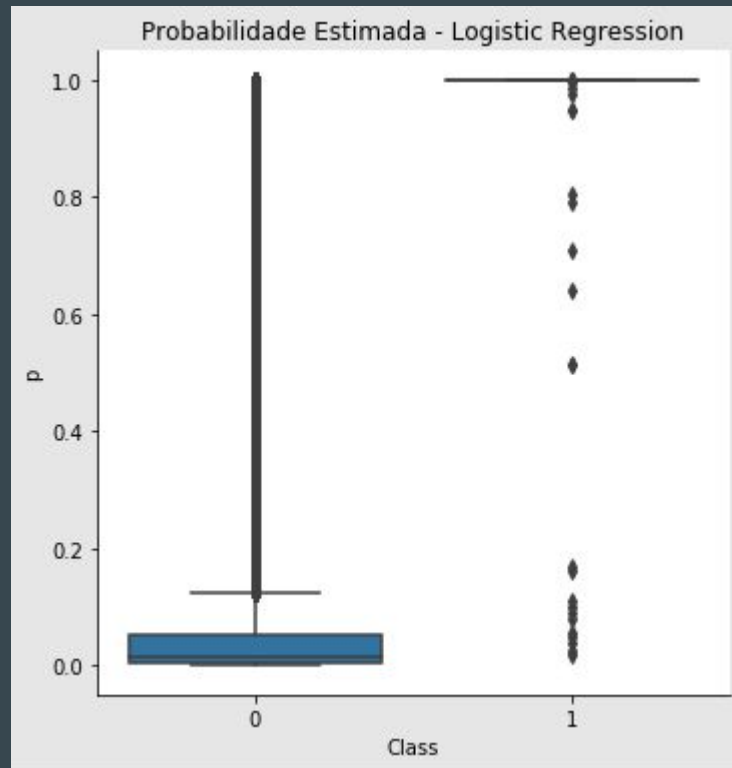
O **modelo logístico** apresentou a melhor separação na base de teste.

“Alto risco” = $P(\text{Fraude}) > 0,5$

“Baixo risco” = $P(\text{Fraude}) \leq 0,5$

2,7% das transações são de alto risco

Sugestão: Usar a classificação de alto risco para fazer uma segunda autenticação com o usuário, pois a precisão é baixa (4,1%).



Obrigada!