

# Atividade 1 Aprendizado Supervisionado

## Relatório de Atividade

Julio Vinicius Amaral Oliveira  
Matrícula: 230537  
Disciplina: MO432  
Professor: Marcos M. Raimundo

### SUMÁRIO

<b>I</b>	<b>Introdução</b>	<b>1</b>
<b>II</b>	<b>Metodologia (Parte Prática)</b>	<b>1</b>
II-A	Análise Exploratório de Dados (EDA) e Pré-processamento . . .	1
II-B	Abordagem 1: Modelagem da Classe Normal (Detecção de Anomalia) . . . . .	2
II-C	Abordagem 2: Classificação Supervisionada . . . . .	3
II-D	Uso de LGBM . . . . .	4
<b>III</b>	<b>Conclusão</b>	<b>4</b>
<b>IV</b>	<b>Respostas Teórico-Conceituais (Parte 2)</b>	<b>5</b>
IV-A	Exercício 1 . . . . .	5
IV-B	Exercício 2: Verossimilhança, entropia cruzada e regressão logística . . . . .	6
IV-C	Exercício 3: Decomposição viés-variância . . . . .	6
IV-D	Exercício 4: Avaliação em dados desbalanceados . . . . .	7
<b>V</b>	<b>Referências</b>	<b>8</b>
	<b>Apêndice</b>	<b>8</b>

### I. INTRODUÇÃO

- O objetivo do trabalho é aplicar técnicas básicas de aprendizado supervisionado para detectar fraudes em transações financeiras. O trabalho foi dividido em duas partes: a primeira parte vamos implementar dois modelos de detecção de fraude, enquanto a segunda parte resolveremos algumas questões teóricas relacionadas ao aprendizado supervisionado.
- O dataset utilizado é o Credit Card Fraud Detection. Ele possui transações realizadas com cartões de crédito em setembro de 2013. O conjunto de dados tem a maioria de suas features anonimizadas, com exceção de algumas variáveis como: tempo, valor da transação e a variável que indica se a transação é fraude ou não. Uma característica bem importante desse dataset é que ele é muito desbalanceado, com apenas 0.172% das transações sendo fraudes.

### II. METODOLOGIA (PARTE PRÁTICA)

#### *A. Análise Exploratório de Dados (EDA) e Pré-processamento*

Como dito anteriormente, o dataset é bastante desbalanceado, o primeiro passo realizado foi verificar a quantidade de fraudes e não fraudes utilizando o método `value_counts()` do pandas. Obtendo o seguinte resultado:

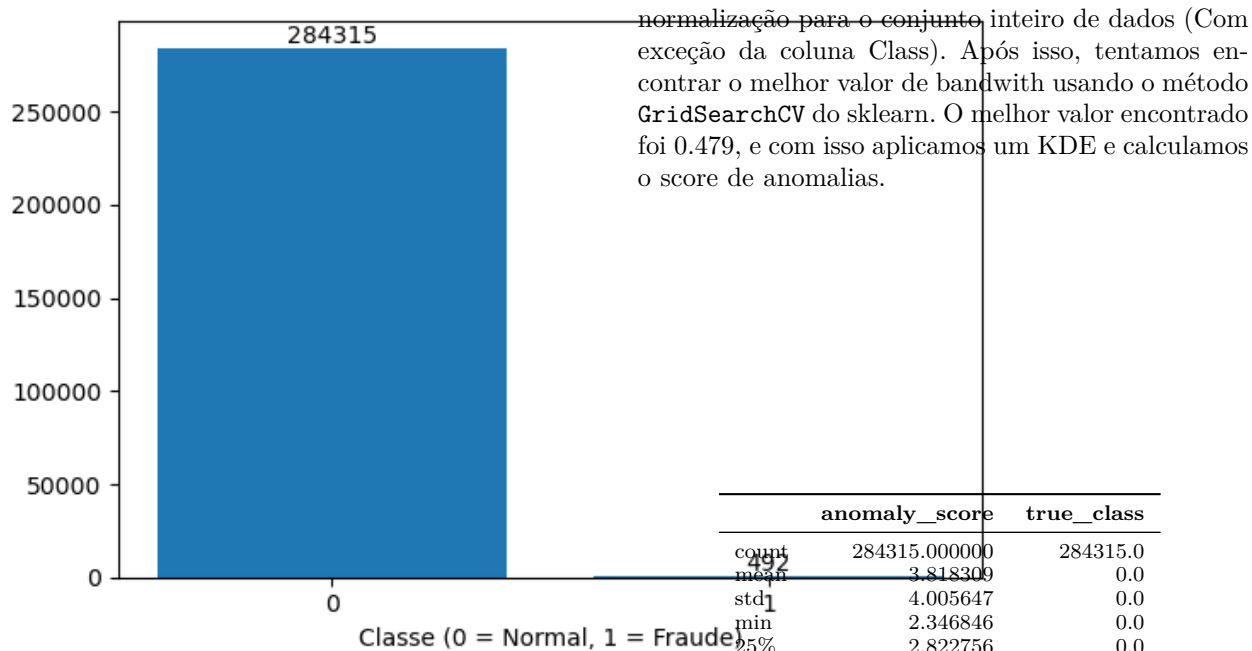


Figura 1. Proporção de fraudes e não fraudes no dataset

Após isso, procuramos entender como as features se relacionavam com a coluna de Class, que era nossa coluna alvo. Então usamos o método `corr()` do `pandas`, buscando a correlação entre Class e o restante das features, obtendo o seguinte resultado para as features com maior correlação absoluta:

Feature	Correlação com Class
V17	-0.326481
V14	-0.302544
V12	-0.260593
V10	-0.216883
V16	-0.196539
V3	-0.192961
V7	-0.187257

Tabela I  
CORRELAÇÃO DAS FEATURES COM CLASS

### B. Abordagem 1: Modelagem da Classe Normal (Detecção de Anomalia)

Para a primeira abordagem utilizamos o modelo de detecção de anomalias. Primeiramente, separamos apenas as 3 features mais relevantes para o modelo, esse número foi escolhido com base na performance do modelo. Foi percebido que ao diminuir cada vez mais o número de features tanto a medida de AUC-PR quanto a matriz de confusão resultavam em valores melhores. Após isso, separamos 2000 transações normais, e com o `StandardScaler` aplicamos uma normalização para esses dados. Depois aplicamos a

	anomaly_score	true_class
count	284315.000000	284315.0
mean	3.818309	0.0
std	4.005647	0.0
min	2.346846	0.0
25%	2.822756	0.0
50%	3.311860	0.0
75%	4.210312	0.0
max	557.518645	0.0

Tabela II  
TRANSAÇÕES NORMAIS

	anomaly_score	true_class
count	492.000000	492.0
mean	162.766740	1.0
std	250.679879	0.0
min	2.780611	1.0
25%	26.711496	1.0
50%	68.714390	1.0
75%	156.108730	1.0
max	1377.747682	1.0

Tabela III  
TRANSAÇÕES FRAUDULENTAS

Após isso, tentamos encontrar o melhor limiar para separar as transações, utilizamos o método `precision_recall_curve` do `sklearn`, e encontramos o melhor limiar para maximizar `f1` sendo 34.9. Com isso, conseguimos os seguintes resultados após plotar a matriz de confusão:

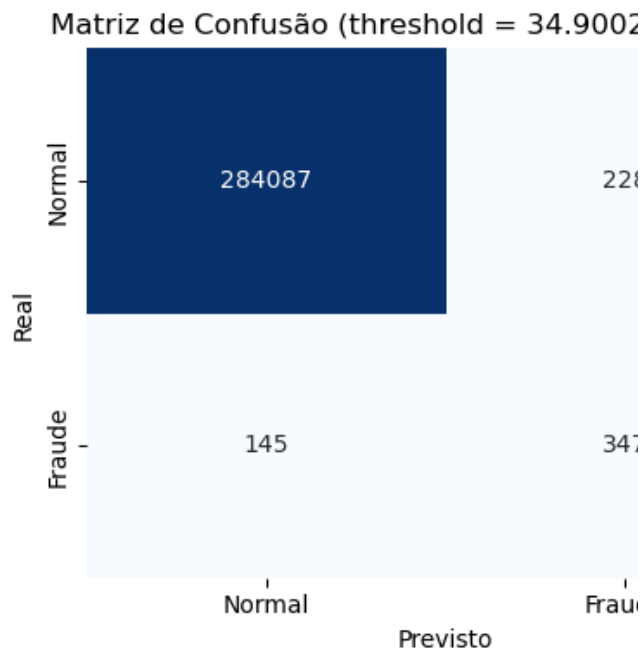


Figura 2. Matriz de confusão

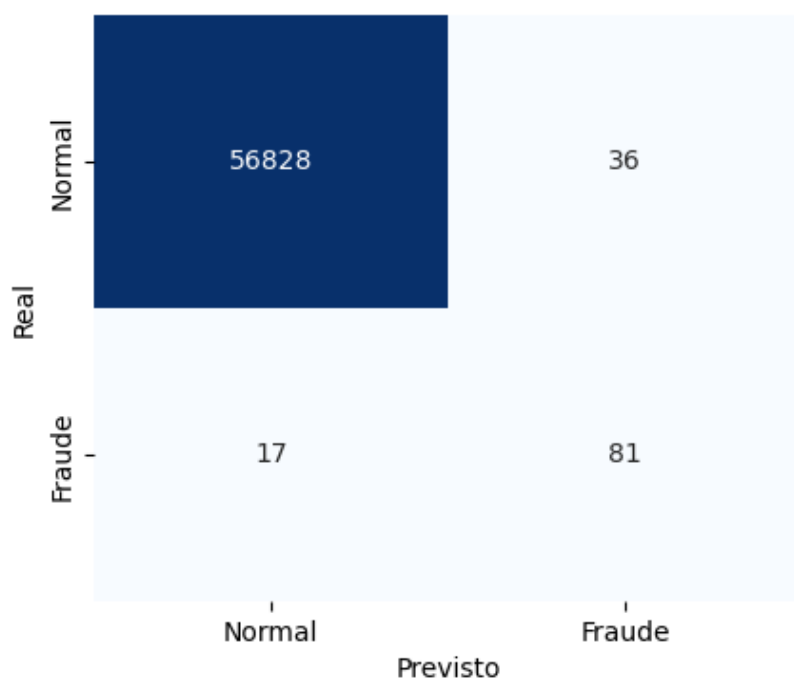


Figura 3. Matriz de confusão regressão logística

Segundo o próprio repositório do Kaggle, a matriz de confusão não é o melhor indicador de desempenho quando temos classes desbalanceadas e nesses casos é recomendado o uso da área abaixo da curva de precision-recall. Cujo resultado foi de 0.623

### C. Abordagem 2: Classificação Supervisionada

Nessa segunda abordagem utilizamos um modelo de regressão logística, utilizamos o método `train_test_split` do `sklearn` para separar 80% dos dados para treino e 20% para teste, isso foi realizado nos dados já normalizado com a ajuda do `StandartScaler`. Após isso, aplicamos o modelo de regressão logística e aplicamos o modelo nos dados de testes, obtendo a seguinte AUC-PR: 0.742, que foi um valor muito maior do que o obtido na primeira abordagem. Para ilustrar também o resultado, plotamos a matriz de confusão:

Para melhorar ainda mais o resultado, resolvemos encontrar um `threshold` utilizando o método `precision_recall_curve` do `sklearn`, e ao fazer isso encontramos um valor melhor de recall, entretanto o valor da precision diminuiu um pouco. Mas dado o contexto que o modelo deveria ser usado, ter um recall alto é melhor pois evita perdas financeiras para a empresa que utiliza o modelo, o valores obtidos foram

	Classe 0	Classe 1	Acurácia	Média Macro	Média Classe 0	Classe 0	Classe 1	Acurácia	Média Macro	Média Classe 0
Precision	0.999	0.829	0.999	0.914	Precision	0.999	0.692	0.999	0.846	0.999
Recall	1.0	0.643	0.999	0.821	Recall	0.999	0.827	0.999	0.913	0.999
F1-score	1.0	0.724	0.999	0.862	F1-score	0.999	0.753	0.999	0.877	0.999
Support	56864	98	0.999	56962	Support	56864	98	56962	56962	56962

Tabela IV  
MÉTRICAS DE AVALIAÇÃO DO MODELO PARA AS CLASSES 0 (NORMAL) E 1 (FRAUDE)

Tabela V  
MÉTRICAS DE AVALIAÇÃO DO MODELO (VALORES APROXIMADOS PARA 3 CASAS DECIMAIS)

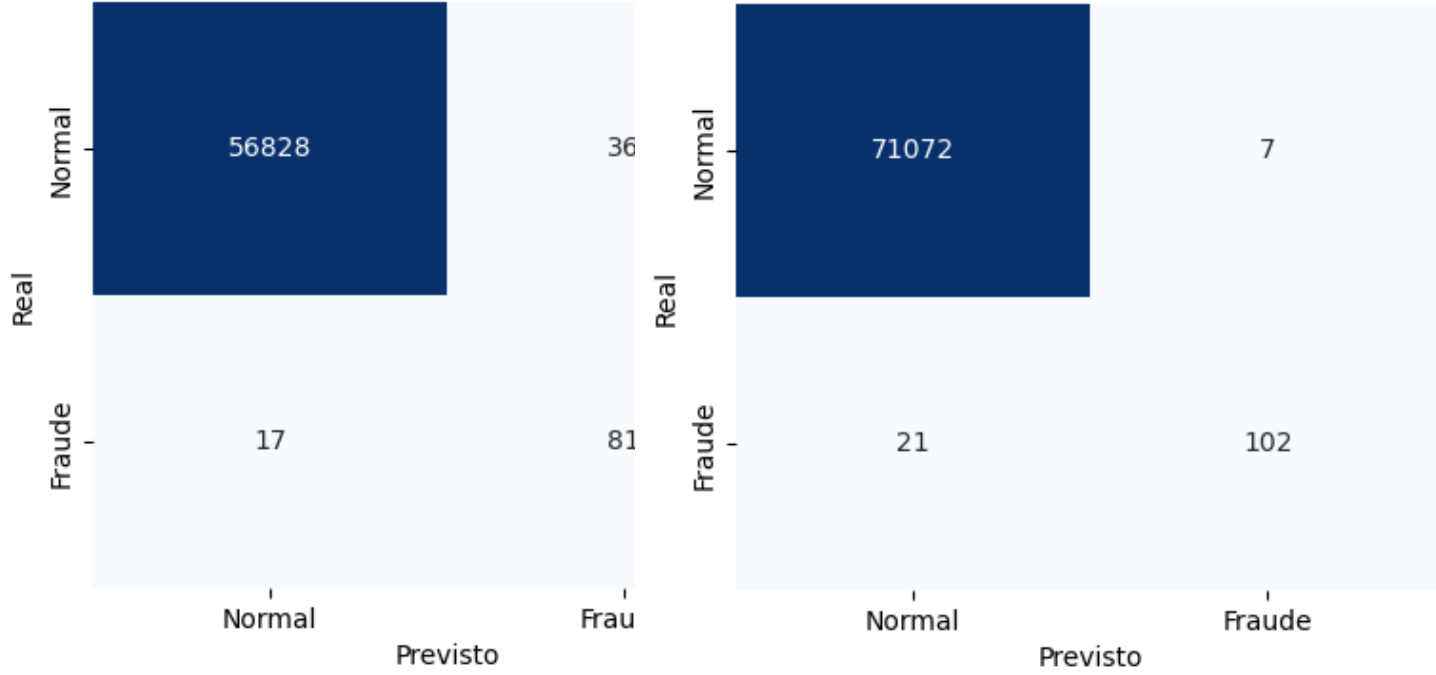


Figura 4. Matriz de confusão regressão logística

#### D. Uso de LGBM

Ao utilizar o modelo LGBM, obtivemos um resultado muito superior aos anteriores. Calculamos os scores utilizando o `predict_proba` e ajustamos o threshold para maximizar o F1-score. Os thresholds, precision e recall foram encontrados utilizando o método `precision_recall_curve` e calculamos o F1 através da fórmula:

$$F1 = \frac{2 \cdot (precision \cdot recall)}{precision + recall} \quad (1)$$

Seguem os resultados de precisão, recall e F1-score:

	0	1	Accuracy	macro avg	weighted avg
Precision	1.000	0.936	1.000	0.968	1.000
Recall	1.000	0.829	1.000	0.915	1.000
F1-score	1.000	0.879	1.000	0.940	1.000
Support	71079	123	1.000	71202	71202

Tabela VI  
MÉTRICAS DE AVALIAÇÃO DO MODELO LGBM

Também calculamos a AUC-PR, que foi de 0.8585. A utilização do LGBM trouxe uma distribuição muito heterogênea entre os scores, o que pode ser visto na figura 5. Enquanto os scores de dados normais se concentram em torno de 0, os scores de dados fraudulentos se concentram em torno de 1, mostrando que o LGBM conseguiu separar bem as classes.

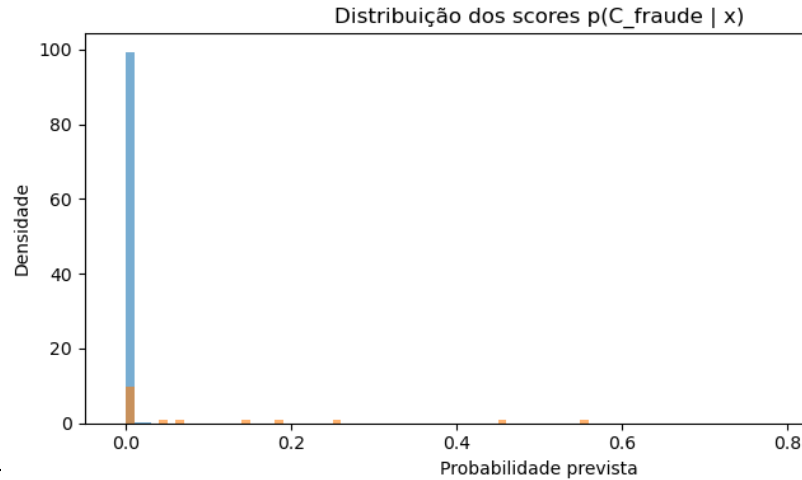


Figura 5. Distribuição dos scores do modelo LGBM

### III. CONCLUSÃO

A primeira abordagem, que utilizou o modelo de detecção de anomalias, obteve um resultado de AUC-

PR de 0.623, enquanto a segunda abordagem, que utilizou o modelo de regressão logística, obteve um resultado de AUC-PR de 0.742. Por fim, a terceira abordagem, que utilizou o modelo LGBM, obteve um resultado de AUC-PR de 0.8585. Isso mostra que o modelo LGBM foi o mais eficaz para detectar fraudes no dataset. Deve-se levar em conta também, que na primeira abordagem, o resultado só foi tão alto, devido ao fato de que reduzimos consideravelmente a quantidade de features, pois em testes prévios, ao utilizar todas as features, os resultados eram muito ruins, com AUC-PR próximos de 0.1. A abordagem 2 pareceu melhor, tanto em questão de resultado quanto em questão de performance, já que para realizar a abordagem 1, tivemos que separar um conjunto inferior de dados, utilizar menos features para chegar em um resultado minimamente aceitável. Já na abordagem 2, utilizamos todas as features e mesmo assim chegamos em um resultado muito bom. Quanto a escolha do limiar, utilizamos o método `precision_recall_curve` para encontrar o F1 com maior valor, em um cenário real, talvez o ideal fosse utilizar um limiar que favorecesse ainda mais o recall, dado que estamos lidando com fraudes que podem significar perda financeira para uma empresa, mas ao mesmo tempo é importante considerar que um modelo que tem uma precisão baixa pode causar problemas tanto para clientes que possam ter transações bloqueadas, causando insatisfação quanto para um time que precisará investigar essas transações.

#### IV. RESPOSTAS TEÓRICO-CONCEITUAIS (PARTE 2)

##### A. Exercício 1

##### (a) Fronteira de Decisão com Covariâncias Diferentes e Priors Iguais

Sabemos que a probabilidade posterior de uma classe  $C_k$  dado um vetor de features  $\mathbf{x}$  é dada pelo Teorema de Bayes:

$$p(C_k|\mathbf{x}) = \frac{p(\mathbf{x}|C_k)p(C_k)}{p(\mathbf{x})}.$$

A fronteira de decisão é quando  $p(C_1|\mathbf{x}) = p(C_2|\mathbf{x})$ . Segundo o enunciado,  $p(C_1) = p(C_2) = 0.5$ .

Sabemos que:

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{1}{(2\pi)^{D/2}|\boldsymbol{\Sigma}_k|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right)$$

Igualando  $p(\mathbf{x}|C_1) = p(\mathbf{x}|C_2)$

$$\frac{p(C_1|\mathbf{x})p(\mathbf{x})}{p(C_1)} = \frac{p(C_2|\mathbf{x})p(\mathbf{x})}{p(C_2)}$$

Ao simplificar, obtemos:

$$p(\mathbf{x}|C_1) = p(\mathbf{x}|C_2)$$

Portanto, igualando a expressão obtida acima e aplicando o logaritmo natural:

$$\begin{aligned} & -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) - \frac{1}{2} \ln |\boldsymbol{\Sigma}_1| \\ & = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_2^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) - \frac{1}{2} \ln |\boldsymbol{\Sigma}_2|. \end{aligned}$$

Dividindo por  $-1/2$  e rearranjando, a equação fica:

$$(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) + \ln |\boldsymbol{\Sigma}_1| = (\mathbf{x} - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_2^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) + \ln |\boldsymbol{\Sigma}_2|.$$

Como  $\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$  é escalar,  $\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} = (\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})^T = \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}$ , expandindo os termos e aplicando essa propriedade, obtemos:

$$\mathbf{x}^T \boldsymbol{\Sigma}_1^{-1} \mathbf{x} - 2\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}_1^{-1} \mathbf{x} + \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \ln |\boldsymbol{\Sigma}_1| = \mathbf{x}^T \boldsymbol{\Sigma}_2^{-1} \mathbf{x} - 2\boldsymbol{\mu}_2^T \boldsymbol{\Sigma}_2^{-1} \mathbf{x} + \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2 + \ln |\boldsymbol{\Sigma}_2|$$

Reorganizando os termos, obtemos:

$$\mathbf{x}^T (\boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\Sigma}_2^{-1}) \mathbf{x} - 2(\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}_2^{-1}) \mathbf{x} + (\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2) + \ln |\boldsymbol{\Sigma}_1| - \ln |\boldsymbol{\Sigma}_2| = 0$$

Essa equação tem um formato de  $\mathbf{x}^T A \mathbf{x} + b^T \mathbf{x} + c = 0$  que é uma superfície quadrática.

##### (b) Efeito da Alteração do Prior com Covariâncias Compartilhadas

Como os priors são  $p(C_1) = \pi$  e  $p(C_2) = 1 - \pi$ , tomando a condição para a fronteira de decisão  $p(\mathbf{x}|C_1)p(C_1) = p(\mathbf{x}|C_2)p(C_2)$  temos que:

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma})\pi = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma})(1 - \pi).$$

Aplicando o logaritmo natural nos dois lados:

$$\begin{aligned} & -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) - \frac{1}{2} \ln |\boldsymbol{\Sigma}| + \ln \pi \\ & = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) - \frac{1}{2} \ln |\boldsymbol{\Sigma}| + \ln(1 - \pi). \end{aligned}$$

Expandindo os termos, temos:

$$-\frac{1}{2}(\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - 2\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1) + \ln \pi = -\frac{1}{2}(\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - 2\boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2) + \ln(1 - \pi)$$

Simplificando e rearranjando, temos:

$$(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2}(\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2) + \ln\left(\frac{\pi}{1 - \pi}\right) = 0.$$

Essa expressão tem o formato de uma equação:

$$\omega^T \mathbf{x} + \omega_0 = 0$$

$$\omega_0 = -\frac{1}{2}(\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2) + \ln\left(\frac{\pi}{1 - \pi}\right)$$

$\ln\left(\frac{\pi}{1 - \pi}\right)$  é um deslocamento escalar na equação linear. Alterar o prior  $\pi$  resulta no deslocamento paralelo da fronteira de decisão.

- Se  $\pi > 0.5$  ( $p(C_1) > p(C_2)$ ), implica que  $\ln\left(\frac{\pi}{1-\pi}\right) > 0$ . Para satisfazer a equação,  $(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} \mathbf{x}$  precisa ser menor, fazendo com que a fronteira se expanda para a região de  $C_1$ .
- Se  $\pi < 0.5$  (ou seja,  $p(C_1) < p(C_2)$ ), então  $\ln\left(\frac{\pi}{1-\pi}\right) < 0$ , e a fronteira se desloca na direção oposta, expandindo a região de  $C_2$ .
- Se  $\pi = 0.5$ , o termo logarítmico é zero, e a fronteira não se desloca para nenhum dos lados, sendo definida pelas médias e covariância.

(c) *Uso da Distância de Mahalanobis para Classificação:*

A distância de Mahalanobis pode ser usada para classificação comparando a distância de um ponto  $\mathbf{x}$  para os centros das classes ( $\boldsymbol{\mu}_1$  e  $\boldsymbol{\mu}_2$ ), ponderadas pelas suas respectivas matrizes de covariância. A regra de classificação seria classificar  $\mathbf{x}$  como pertencente à classe  $C_1$  para a qual a distância de Mahalanobis quadrática  $\Delta_1^2$  é mínima. Ou seja:

- Se  $\Delta_1^2 \leq \Delta_2^2$ , classificar  $\mathbf{x}$  como  $C_1$  (Normal).
- Se  $\Delta_1^2 > \Delta_2^2$ , classificar  $\mathbf{x}$  como  $C_2$  (Fraude).

Premissas Implícitas:

- As classes devem seguir uma distribuição Gaussiana multivariada.
- Os parâmetros  $\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2$  são conhecidos ou foram bem estimados.
- Priors iguais para as classes  $p(C_1) = p(C_2)$ .

B. *Exercício 2: Verossimilhança, entropia cruzada e regressão logística*

(A) MOSTRAR QUE MAXIMIZAR A FUNÇÃO DE VEROSSIMILHANÇA É EQUIVALENTE A MINIMIZAR A ENTROPIA CRUZADA.

Para um modelo de regressão logística com alvos binários:

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n}$$

Sabendo que maximizar a função é o mesmo que maximizar o logaritmo natural, podemos aplicar o logaritmo natural para que a função se pareça mais com o nosso objetivo:

$$\ln p(\mathbf{t}|\mathbf{w}) = \sum_{n=1}^N \ln(y_n^{t_n} (1 - y_n)^{1-t_n}).$$

Usando as propriedades do logaritmo:

$$\ln p(\mathbf{t}|\mathbf{w}) = \sum_{n=1}^N [t_n \ln y_n + (1 - t_n) \ln(1 - y_n)].$$

A função de erro de entropia cruzada é definida por:

$$E(\mathbf{w}) = - \sum_{n=1}^N [t_n \ln y_n + (1 - t_n) \ln(1 - y_n)].$$

Comparando as duas equações, vemos que  $E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w})$ . Portanto, maximizar a verossimilhança é equivalente a minimizar a entropia cruzada.

(B) DERIVANDO  $\nabla E(\mathbf{w})$

Lembrando que  $y_n = \sigma(\mathbf{a}_n)$ , com  $\mathbf{a}_n = \mathbf{w}^T \boldsymbol{\phi}_n$ , e que

$$\frac{d\sigma(a)}{da} = \sigma(a)(1 - \sigma(a)),$$

a função de erro é:

$$E(\mathbf{w}) = - \sum_{n=1}^N [t_n \ln \sigma(\mathbf{w}^T \boldsymbol{\phi}_n) + (1 - t_n) \ln(1 - \sigma(\mathbf{w}^T \boldsymbol{\phi}_n))].$$

O gradiente de  $E(\mathbf{w})$  é:

$$\begin{aligned} \nabla E(\mathbf{w}) &= - \sum_{n=1}^N \left[ t_n \frac{1}{y_n} y_n (1 - y_n) \boldsymbol{\phi}_n + (1 - t_n) \frac{1}{1 - y_n} (-y_n (1 - y_n)) \boldsymbol{\phi}_n \right] \\ &= - \sum_{n=1}^N [t_n (1 - y_n) - (1 - t_n) y_n] \boldsymbol{\phi}_n \\ &= \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_n. \end{aligned}$$

(C) EXPLICAÇÃO DO PROBLEMA DE OVERFITTING

O overfitting ocorre quando os dados são linearmente separáveis. Nesse caso, existe um  $\mathbf{w}$  tal que  $\mathbf{w}^T \boldsymbol{\phi}_n$  separa perfeitamente as classes, fazendo com que  $\|\mathbf{w}\| \rightarrow \infty$ .

Para evitar isso, usa-se regularização  $L_2$ :

$$E_{\text{regularizado}}(\mathbf{w}) = - \sum_{n=1}^N [t_n \ln y_n + (1 - t_n) \ln(1 - y_n)] + \frac{\lambda}{2} \|\mathbf{w}\|^2,$$

O termo  $\frac{\lambda}{2} \|\mathbf{w}\|^2$  penaliza pesos grandes, fazendo com que o modelo faça um overfitting.

C. *Exercício 3: Decomposição viés-variância*

(a) *Explicação (bias)<sup>2</sup>, variance e noise*

a) *Bias<sup>2</sup>*: O viés quadrático mede o erro sistemático do modelo, ou seja, o quanto em média o modelo se difere do valor real. Um modelo com um viés alto não consegue entender a complexidade dos dados, levando a um underfitting.

b) *Variance*: A variância mede o quanto as nossas previsões para um dado ponto variam usando diferentes conjuntos de treino. Um modelo que tem alta variância é muito sensível às variações nos dados de treino, fazendo com que o modelo não consiga performar bem em dados que ele não conhece, levando a um overfitting.

c) *Noise*: O ruído é a parte do erro que não é possível de ser reduzida, pois ele é uma característica inerente aos dados.

(b) *Explicação* ( $bias$ )<sup>2</sup>, *variance* e *noise*

D. *Exercício 4: Avaliação em dados desbalanceados*

(a) *Acurácia padrão em datasets desbalanceados*

Em dados desbalanceados, a acurácia padrão pode dar uma falsa sensação que um modelo está muito bom. Já que um classificador simples que sempre classifica as transações como a classe majoritária pode ter uma alta acurácia pois ele irá classificar corretamente na grande maioria das vezes. Isso acontece porque a acurácia é definida como:

$$\text{Acurácia} = \frac{TP + TN}{TP + TN + FP + FN}$$

Em datasets desbalanceados, a classe majoritária tem um peso muito grande sobre essa métrica. Sendo assim, o classificador não cumpriria seu objetivo que é identificar fraudes, porém ainda sim teria uma acurácia alta.

(b) *Precisão, Recall e F1-Score*

**Precisão** Parcela das transações classificadas como fraude e que de fato são fraude:

$$\text{Precision} = \frac{TP}{TP + FP}$$

Em que:

- TP: quantidade de fraudes classificadas corretamente.
- FP: quantidade de transações normais classificadas como fraude.

**Recall** Parcela de todas as transações que de fato eram fraudes que foram classificadas corretamente:

$$\text{Recall} = \frac{TP}{TP + FN}$$

Em que:

- FN: quantidade de fraudes que foram classificadas como transações normais.

**F1-Score** Métrica harmônica entre precisão e recall:

$$F_1 = 2 \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

a) *Quando priorizar Recall sobre Precisão:*

- Em situações em que perder uma fraude gera prejuízo financeiro para a empresa que utiliza o classificador ou em cenários em que uma classificação errada pode gerar muito prejuízo para o utilizador, como em resultado de doenças.

b) *Quando priorizar Precisão sobre Recall:*

- Quando o custo de falsos positivos é alto, seja por implicar em esforço manual, gerar insatisfação ou diminuir a confiança do cliente nos serviços oferecidos pela empresa que se utiliza do classificador. Além disso, para que a precisão tenha mais relevância que o recall, é necessário que o impacto de um falso negativo não seja tão alto para o negócio. Por exemplo, em um sistema de identificação de câncer, um falso negativo poderia levar a um diagnóstico mais tardio da doença.

(c) *Curva Precisão-Recall e AUC-PR*

O procedimento para calcular a curva PR seria:

- 1) Ordenar todos os exemplos em ordem decrescente dos seus scores  $s_n$ .
- 2) Para cada limiar  $\tau$ :
  - Classificar  $x_n$  como fraude se  $s_n \geq \tau$ , caso contrário, classificar como normal.
  - Calcular

$$\text{Precision}(\tau) = \frac{\#\{n : s_n \geq \tau, t_n = 1\}}{\#\{n : s_n \geq \tau\}}, \quad \text{Recall}(\tau) = \frac{\#\{n : s_n \geq \tau, t_n = 1\}}{\#\{n : t_n = 1\}}$$

- 3) Plotar cada ponto do par  $(\text{Recall}(\tau), \text{Precision}(\tau))$  e formar a curva.

A área Sob a Curva PR é dada por

$$\text{AUC-PR} = \int_0^1 \text{Precision}(r) dr,$$

em que  $\text{Precision}(r)$  é a precisão correspondente ao recall  $r$ .

c) *Por que AUC-PR é preferível à ROC em cenários muito desbalanceados:* A curva ROC avalia a taxa de verdadeiros positivos e falsos positivos, em casos desbalanceados a taxa de falsos positivo

$$FPR = \frac{FP}{FP + TN}$$

Permanece baixa pois a quantidade de verdadeiro negativos é muito alta, então essa métrica mascara um modelo ruim. Enquanto a AUC-PR avalia diretamente os dados de fraude, ou seja, não temos um denominador muito alto que mascara o valor da métrica.

V. REFERÊNCIAS

APÊNDICE

Feature	Correlação com Class
V17	-0.326481
V14	-0.302544
V12	-0.260593
V10	-0.216883
V16	-0.196539
V3	-0.192961
V7	-0.187257
V11	0.154876
V4	0.133447
V18	-0.111485
V1	-0.101347
V9	-0.097733
V5	-0.094974
V2	0.091289
V6	-0.043643
V21	0.040413
V19	0.034783
V20	0.020090
V8	0.019875
V27	0.017580
Time	-0.012323
V28	0.009536
V24	-0.007221
Amount	0.005632
V13	-0.004570
V26	0.004455
V15	-0.004223
V25	0.003308
V23	-0.002685
V22	0.000805

Tabela VII

CORRELAÇÃO DAS FEATURES COM A VARIÁVEL CLASS