



**Universidade de Brasília
Departamento de Estatística**

Interpretação de redes neurais

Davi Guerra Alves

Projeto apresentado para o Departamento de Estatística da Universidade de Brasília como parte dos requisitos necessários para obtenção do grau de Bacharel em Estatística.

**Brasília
2023**

Davi Guerra Alves

Interpretação de redes neurais

Orientador(a): Thais Carvalho Valadares Rodrigues

Projeto apresentado para o Departamento de Estatística da Universidade de Brasília como parte dos requisitos necessários para obtenção do grau de Bacharel em Estatística.

**Brasília
2023**

Resumo

O objetivo central deste trabalho foi de buscar métodos para interpretar um modelo de redes neurais. Concomitantemente a essa interpretação, realizou-se uma comparação dos resultados desse modelo com um modelo estatístico convencional, a regressão logística. Utilizaram-se dados relacionados a empréstimos, obtidos da plataforma *Kaggle*. A variável de estudo analisada foi a "Condição do empréstimo", que classifica os empréstimos dos clientes como "Bom" ou "Ruim".

A técnica adotada para interpretar os modelos de redes neurais foi o SHAP. Os resultados comparativos entre os modelos revelam que o modelo de rede neural produziu previsões mais heterogêneas, enquanto o modelo logístico ficou limitado pelo desbalanceamento dos dados, rotulando a maior parte dos casos baseado na classe majoritária.

Palavras-chaves: SHAP, redes neurais, regressão logística, interpretação, empréstimos.

Lista de Tabelas

1	Matriz de confusão	25
2	Número de colunas antes e depois da preparação dos dados	30
3	Divisão dos dados para a modelagem dos modelos de regressão logística e redes neurais	30
4	Número de observação em cada categoria da variável resposta	36
5	Valores do coeficiente de Pearson entre as covariáveis e a variável resposta	41
6	Valores do coeficiente de contingência entre as covariáveis e a variável resposta	41
7	Estimativa dos coeficientes do modelo logístico e o erro padrão associado .	42
8	Report do modelo logístico	43
9	Report da rede neural	45
10	Número de parâmetros nos modelos	51
11	Comparação dos resultados de Falsos Positivos e Falsos Negativos	52
12	Comparação dos resultados da Regressão logística e Rede neural. Em verde, o modelo que obteve o melhor resultado na respectiva métrica.	52
13	Relação entre os coeficientes da regressão logística com os valores absolutos de shap	53
14	Tempo de predição (em ms) de cada modelo, em uma amostra com 50 observações.	54
15	Tempo de execução para realizar a interpretação das variáveis	54

Lista de Figuras

1	Neurônio da Rede Neural	12
2	Tipos de Função de Ativação.	13
3	Rede Neural com uma camada oculta	14
4	Arquitetura padrão de um rede neural <i>feedforward</i>	14
5	Comportamentos dos pesos em relação à função de perda. O ponto w_A representa um ponto local mínimo e w_B representa um ponto global mínimo.	16
6	Representação do método do gradiente descendente para a estimação de um parâmetro.	17
7	Ganho do jogador 3 em relação à todas as permutações de jogadores.	20
8	Relação entre permutações e coalizões.	21
9	Cálculo de f_S , sendo S o conjunto de variáveis X_1, X_3, X_4 , dentre as observações de um conjunto de dados.	23
10	Fluxograma da implementação do modelo logístico	31
11	Fluxograma da implementação da rede neural	32
12	Arquitetura de rede neural inicial	33
13	Fluxograma da escolha dos modelos de redes neurais estimados	33
14	Arquitetura da rede neural depois do processo de tunelagem	34
15	Variáveis explicativas em relação à condição do empréstimo	37
16	Variáveis explicativas em relação à condição do empréstimo	38
17	Variáveis explicativas em relação à condição do empréstimo	39
18	Variáveis explicativas em relação à condição do empréstimo	40
19	Matrix de confusão do modelo logístico	43
20	Matrix de confusão da rede neural	45
21	Média absoluta dos valores de shap	46
22	Valor de SHAP para uma observação do tipo Verdadeiro Positivo	47
23	Valor de SHAP para uma observação do tipo Verdadeiro Negativo	48
24	Valor de SHAP para uma observação do tipo Falso Positivo	49
25	Valor de SHAP para uma observação do tipo Falso Negativo	49

26	Valores de shap para as 80 observações utilizadas	50
27	Gráfico de força em uma observação	50
28	Gráfico de força para múltiplas observações	51

Sumário

1 Introdução	9
2 Referencial teórico	10
2.1 Regressão logística	10
2.1.1 Interpretação	11
2.2 Redes neurais artificiais	12
2.2.1 Neurônio	12
2.2.2 Arquitetura	13
2.2.3 <i>Forward propagation</i>	15
2.2.4 Função de perda	16
2.2.5 <i>Backpropagation</i>	16
2.3 SHAP	18
2.3.1 Valores de Shapley	18
2.3.2 <i>Shapley Additive Explanations</i>	22
2.4 Medidas de associação	23
2.4.1 Coeficiente de Pearson	24
2.4.2 Coeficiente de contingência	24
2.5 Métricas de avaliação	24
2.5.1 Matriz de confusão	25
2.5.2 Acurácia	25
2.5.3 Precisão	26
2.5.4 Recall	26
2.5.5 F1-score	26
3 Metodologia	28
3.1 Conjunto de dados	28
3.1.1 Variáveis	28
3.1.2 Limpeza dos dados	30
3.2 Modelagem dos dados	31
3.2.1 Regressão logística	31

3.2.2 Rede neural	32
3.3 Interpretação dos modelos.	34
4 Resultados	36
4.1 Análise descritiva.	36
4.1.1 Condição do empréstimo	36
4.1.2 Relação entre as covariáveis e a variável resposta	37
4.2 Regressão logística	42
4.3 Rede neural	44
4.4 Interpretação da rede neural	46
4.5 Benchmark entre regressão logística e redes neurais	51
4.5.1 Complexidade da arquitetura	51
4.5.2 Resultado dos modelos	52
4.5.3 Tempo de execução	54
5 Conclusão	56
6 Anexo.	58

1 Introdução

As redes neurais são modelos matemáticos que, unidos às técnicas computacionais, visam tentar reproduzir o funcionamento da estrutura neural presente no ser humano, buscando, assim, realizar tarefas complexas, como o reconhecimento de padrões, identificação de imagens, processamento de linguagem natural, etc. No entanto, apesar de sua eficácia em muitas aplicações, as redes neurais podem ficar muito complexas conforme sua arquitetura cresce, sendo consideradas como "caixas pretas", devido à sua complexidade e falta de transparência.

Por isso, a interpretação de redes neurais é uma área cada vez mais essencial, pois busca entender como esses modelos tomam decisões e quais fatores influenciam suas saídas. Entender o porquê uma rede neural tomou tal decisão é importante em diversas áreas, como a área da saúde, em diagnósticos médicos, e na área bancária, analisando um risco de crédito.

Uma das técnicas mais promissoras para a interpretação de redes neurais é o SHAP (Shapley Additive Explanations), que foi introduzido em 2017 (LUNDBERG; LEE, 2017). O SHAP é uma técnica de interpretação que fornece explicações locais e globais para as saídas da rede neural. Ele é baseado no conceito matemático de valor de Shapley (SHAPLEY, 1953), que atribui uma contribuição de importância para cada recurso de entrada na saída da rede neural.

Portanto, esse trabalho tem como objetivo principal explorar a técnica SHAP para a interpretação de redes neurais e sua aplicação em diversas áreas, pois, ao compreender como as redes neurais funcionam, e quais são os recursos mais importantes para suas decisões, será possível tornar o método mais confiável e transparentes para os usuários.

Ao optar por empregar um modelo estatístico tradicional, como a regressão logística, que já possui um conjunto de ferramentas interpretativas consolidadas, este estudo também se concentra em uma análise comparativa entre o modelo de redes neurais e o modelo logístico. Essa abordagem visa avaliar tanto os resultados quanto as interpretações geradas por ambos os modelos, proporcionando uma compreensão mais profunda de como essas abordagens se comportam em um cenário comum. Essa comparação não apenas lança luz sobre as diferenças de desempenho, mas também destaca as diferenças interpretativas distintas entre os dois modelos, enriquecendo assim a compreensão sobre a escolha adequada de modelos em diferentes contextos.

2 Referencial teórico

2.1 Regressão logística

A regressão logística é um método estatístico utilizado para modelar a probabilidade de uma variável dependente categórica. É comumente utilizada para problemas de classificação binária, onde a variável dependente possui apenas duas categorias, como sim/não, positivo/negativo, 0/1 (HOSMER, 2013).

O cálculo da regressão logística é baseado na probabilidade da variável aleatória Y ser igual a 1, onde Y é uma variável aleatória com distribuição Bernouli, com parâmetro p de sucesso, cuja fórmula é dada por:

$$P(Y = 1|x_1, x_2, \dots, x_k) = \frac{1}{1 + e^{-(\beta_0 + x_1\beta_1 + x_2\beta_2 + \dots + x_k\beta_k)}} \quad (2.1.1)$$

onde cada variável explicativa (x_1, x_2, \dots, x_k) tem um parâmetro β correspondente, influenciando o resultado de Y .

A estimação dos coeficientes $(b_0, b_1, b_2, \dots, b_k)$ na regressão logística é geralmente realizada por meio do método da máxima verossimilhança. O objetivo é encontrar os valores dos coeficientes que maximizam a função de verossimilhança, representando a probabilidade de observar os dados observados dado o modelo.

A função de verossimilhança (L) para a regressão logística é dada pelo produto das probabilidades condicionais de observar os eventos (valores da variável dependente) dados os valores das variáveis independentes. Para facilitar o cálculo, geralmente trabalhamos com o logaritmo natural da função de verossimilhança, conhecido como log-verossimilhança(l).

A log-verossimilhança para a regressão logística é:

$$l(\beta) = \sum_{i=1}^N [y_i \beta^T x_i - \log(1 + e^{\beta^T x_i})] \quad (2.1.2)$$

- N é o número total de observações.
- y_i é a variável dependente binária da i -ésima observação (0 ou 1).
- p_i é a probabilidade predita de $Y = 1$ para a i -ésima observação, dada pela função logística.

A ideia é encontrar os valores de $(b_0, b_1, b_2, \dots, b_k)$ que maximizam essa função.

Isso geralmente é feito usando métodos computacionais, como o algoritmo de otimização Newton-Raphson ou o Gradiente Descendente.

2.1.1 Interpretação

Para se interpretar o modelo logístico é utilizado a Razão de Chances, que calcula a razão entre duas chances, sendo a chance de evento definida como a probabilidade do sucesso do evento sobre a probabilidade do fracasso.

$$\text{Chance} = \frac{P(Y = 1|X = x)}{P(Y = 0|X = x)} \quad (2.1.3)$$

Para calcular a Razão de Chances, inicialmente, é determinada a chance de ocorrência de um evento específico, designada como "Chance de X". Em seguida, é calculada a chance desse mesmo evento ocorrer quando a variável correspondente é incrementada em uma unidade, referida como "Chance de X + 1". Dessa forma, a fórmula para a Razão de Chances é expressa por:

$$\text{Razão de Chances} = \frac{\text{Chance de } X + 1}{\text{Chance de } X} \quad (2.1.4)$$

Esse mesmo resultado pode ser obtido quando se calcula a exponencial dos coeficientes do modelo logístico.

$$RC = \exp(\beta_k) \quad (2.1.5)$$

Logo, para o crescimento de 1 unidade em X_k , variável associada ao β_K , a chance do evento ocorrer aumento em β_k unidades, considerando as demais variáveis constantes.

Um RC igual a 1 indica que a variável independente não tem efeito no resultado de Y(nenhuma associação). Um RC maior que 1 sugere uma associação positiva, enquanto um RC menor que 1 sugere uma associação negativa.

Outra medida interpretativa é a função log odds ou logíto. Ela é uma função que calcula o log da razão do evento acontecer e dele não acontecer, cuja formulação é dada por:

$$\text{logit}(P(Y = 1)) = \log \left(\frac{P(Y = 1)}{1 - P(Y = 1)} \right) \quad (2.1.6)$$

onde esse resultado nada mais é do que $\beta_0 + X_1\beta_1 + X_2\beta_2 + \dots + X_k\beta_k$.

Portanto, a utilidade de se analisar o log-odds é justamente uma ponte entre

olhar coeficiente e a probabilidade final, pois um coeficiente positivo indica que o aumento na variável está associado a um aumento nas log-odds (e, portanto, na probabilidade), enquanto um coeficiente negativo está associado a uma diminuição nas log-odds (e na probabilidade).

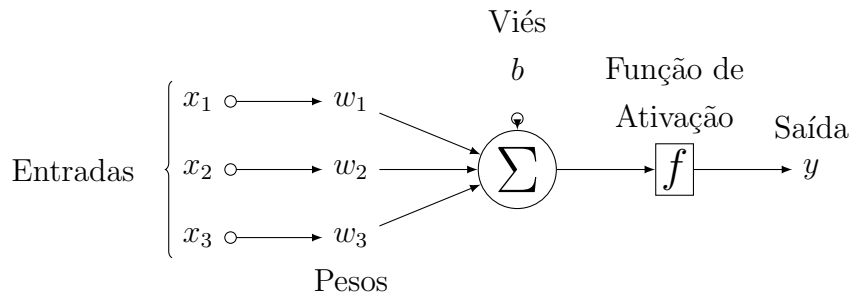
2.2 Redes neurais artificiais

Redes Neurais Artificiais (ou *Deep Learning*) é uma técnica preditiva presente no campo de Inteligência Artificial. As redes neurais tem sido amplamente utilizadas devido ao seu alto poder preditivo e também à flexibilidade de se aplicar esse método em diversos contextos, permitindo ser um modelo com menos restrições que os modelos tradicionais estatísticos.

2.2.1 Neurônio

Uma rede neural tem esse nome devido à tentativa de se reproduzir o comportamento do cérebro humano. Sua arquitetura é composta por um conjunto de unidades denominadas neurônios, e cada neurônio é responsável por receber informações, fazer o tratamento do que foi recebido, e repassar o resultado disso para frente. A Figura 2.2.1 ilustra a estrutura de 1 neurônio. Quando as informações x_i entram no neurônio, acontece primeiramente um processo onde é ponderada cada informação que foi recebida, os chamados **pesos**. Logo em seguida ocorre a soma dessa ponderação. Feito isso, é realizado mais um processo de soma, agora adicionando uma informação própria daquele neurônio nesse resultado. Essa informação é chamada de **bias** (ou Viés). Antes desse resultado ser repassado para outro neurônio, ele passa por uma função que vai definir a natureza daquela informação, chamada de **função de ativação**, retornando assim uma saída y .

Figura 1: Neurônio da Rede Neural

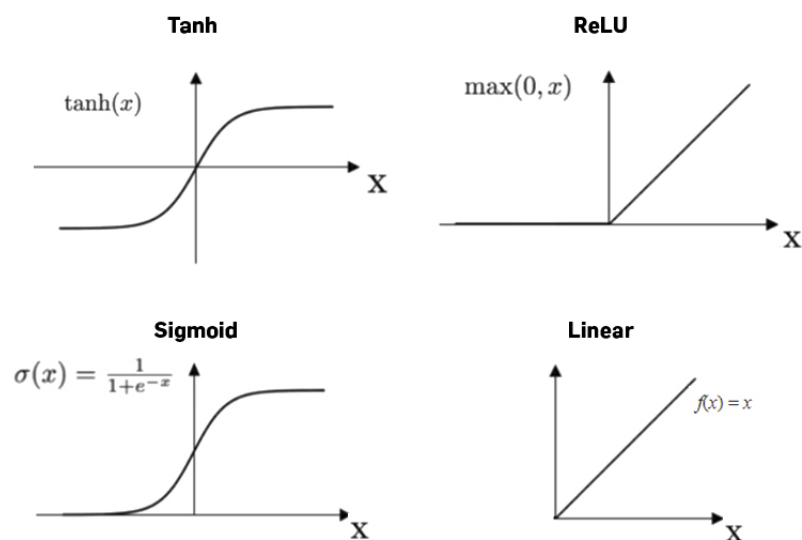


A Figura 1 pode ser representada matematicamente da seguinte maneira:

$$y = f\left(\beta + \sum_{i=1}^{d_x} w_i x_i\right) \quad (2.2.1)$$

onde d_x é o número de entradas.

Figura 2: Tipos de Função de Ativação.



Fonte: <https://machine-learning.paperspace.com/wiki/activation-function>

A Figura 2 mostra alguns tipos de funções de ativação que um neurônio pode ser atribuído. Note como a maioria dessas funções restringe o valor de x (nesse caso o valor calculado no neurônio), no caso da função Sigmoide e da Tangente hiperbólica (\tanh) limitando o valor do neurônio em um intervalo, a ReLU, que é bastante utilizada, desconsidera os valores negativos e existe a função Linear que basicamente só vai repassar a informação do neurônio para frente.

2.2.2 Arquitetura

Uma rede neural é estruturada em camadas formadas por um conjunto de neurônios. Conforme ilustrado na Figura 3, temos as camadas de entrada, as camadas ocultas e a camada de saída. A camada de entrada é o ponto de partida da rede neural, pois é onde as informações das variáveis entram. Logo em seguida encontram-se as camadas ocultas, que são as principais responsáveis por criar redes mais complexas, pois o número de camadas e o número de neurônio dentro dessas camadas podem ser moldados ou adicionados

dependendo do objetivo empregado pela rede, conforme ilustrado na Figura 4. E por fim existe a camada de saída que contém o(s) valor(es) predito(s) pela rede.

Figura 3: Rede Neural com uma camada oculta

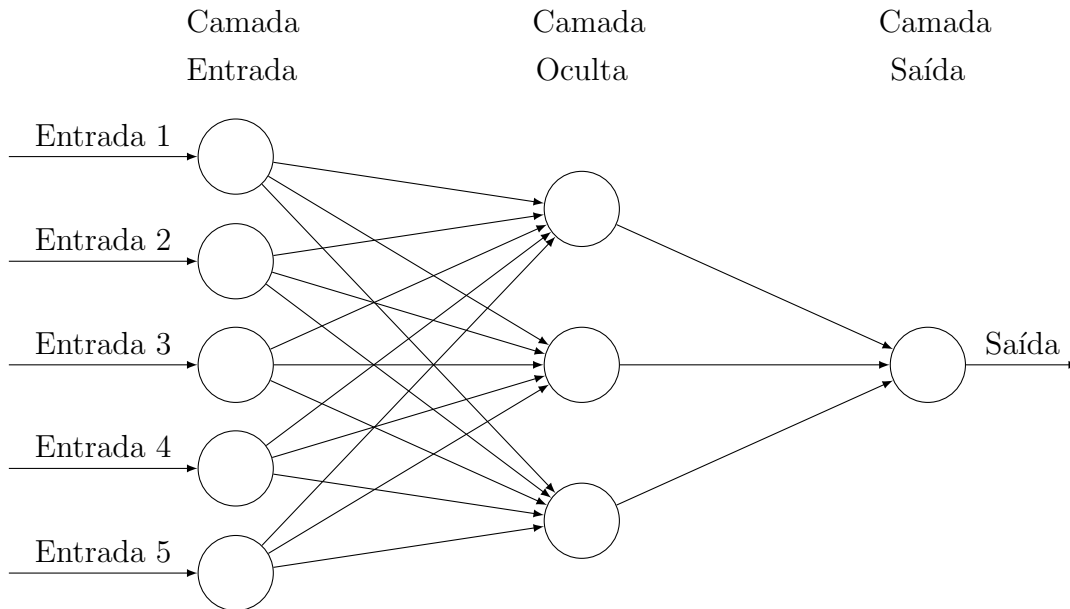
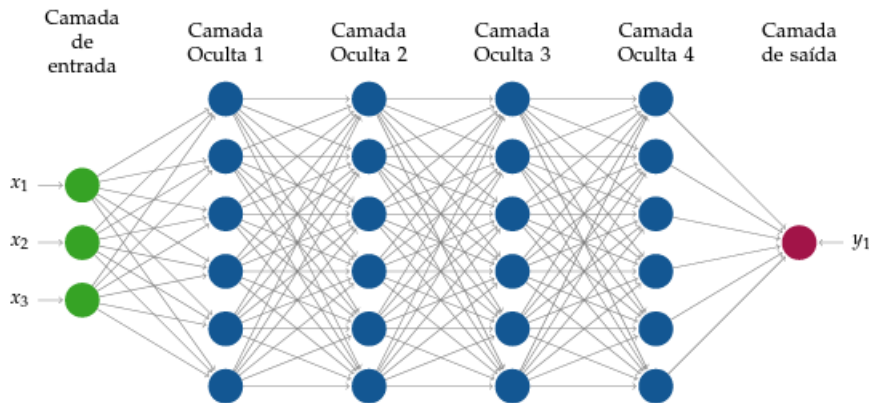


Figura 4: Arquitetura padrão de um rede neural *feedforward*



Fonte: (IZBICKI; SANTOS, 2020)

Note que as Figuras 3 e 4 evidenciam um potencial muito grande de crescimento da rede e naturalmente esse aumento pode acabar gerando um custo computacional elevado quando a rede estiver em treinamento.

2.2.3 Forward propagation

O processo *Forward propagation* (ou propagação direta) é o responsável por transmitir as informações, desde a camada de entrada, passando pelas camadas ocultas, até chegar na camada de saída. O *Forward propagation* utiliza da generalização a Equação 2.2.1 para cada neurônio presente nas camadas internas da rede neural. Por isso temos que, para cada j -ésimo neurônio, da camada l :

$$z_j^{(l)} = b_j^{(l)} + \sum_{i=1}^{d_{(l-1)}} w_{ij} a_i^{(l-1)}$$

onde:

- w_{ij} é o peso associado à conexão entre o neurônio i na camada $l - 1$ e o neurônio j na camada l ;
- $a_i^{(l-1)}$ é a saída do neurônio i na camada anterior ($l - 1$);
- $b_j^{(l)}$ é o viés (bias) associado ao neurônio j na camada l .

Logo em seguida é aplicada uma função de ativação g em $z_i^{(l)}$ que vai ser a responsável por gerar o resultado final $a_i^{(l)}$, do i -ésimo neurônio na l -ésima camada.

$$a_i^{(l)} = g(z_i^{(l)})$$

Esse processo vai ser realizado camada a camada, sequencialmente. Logo, supondo que uma rede neural tenha l camadas ocultas e, cada camada contendo d_l neurônios, considerando também w_{ij} como o peso presente no i -ésimo neurônio com a j -ésima saída na camada seguinte ($l + 1$), onde $l = 0, \dots, H$. Temos que o resultado final da propagação é igual a:

$$f(\mathbf{x}) = \mathbf{a}^{H+1} = g(b_j^{(H+1)} + \sum_{i=1}^{d_H} w_{ij} a_i^H) \quad (2.2.2)$$

Note que a previsão da rede vem diretamente do resultado obtido da última camada oculta, e esse depende da camada que o antecede e assim sucessivamente até chegar na camada de entrada.

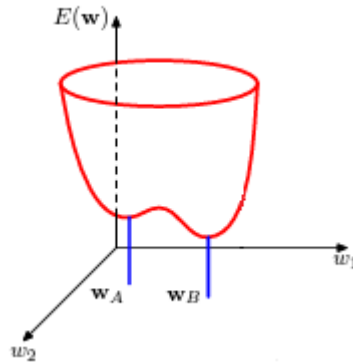
2.2.4 Função de perda

Para se obter informações sobre o desempenho do modelo, é escolhida uma função de perda. Uma função bastante utilizada é a do erro do quadrático médio:

$$EQM(f) = \frac{1}{n} \sum_{k=1}^n (f(\mathbf{x}_k) - y_k)^2$$

Essa função é uma indicadora do quão longe, em média, os valores preditos estão distantes dos valores reais. Note que o resultado da função f depende exclusivamente dos parâmetros da rede (viés e pesos), por isso, se essa função de perda tende a 0, significa que os parâmetros dessa rede alcançaram um ponto mínimo global. Entretanto, devido a complexidade desse modelo, acabam-se escolhidos pontos locais mínimos, que, dependendo do contexto, acabam satisfazendo o objetivo. A Figura 5 ilustra esse comportamento:

Figura 5: Comportamentos dos pesos em relação à função de perda. O ponto w_A representa um ponto local mínimo e w_B representa um ponto global mínimo.



Fonte: (BISHOP, 2006)

2.2.5 Backpropagation

Como a função de perda está relacionada com os parâmetros (θ) da rede, para se minimizar a função de perda $R(\theta)$, é necessário encontrar os valores de θ que resolvam esse problema de otimização. Para fazer isso, é necessário calcular o gradiente de $R(\theta)$ em relação à θ (JAMES et al., 2013),

$$\nabla R(\theta) = \frac{\partial R(\theta)}{\partial \theta} \quad (2.2.3)$$

A rede neural, durante todo o treinamento, aplica esse processo do cálculo do gradiente de $R(\theta)$ em relação à θ . Esse é um processo iterativo, com o objetivo de mudar

o valor de θ afim de conseguir minimizar a função de perda. Com isso a Equação 2.2.3 pode ser descrita nesse processo iterativo como:

$$\nabla R(\theta^m) = \left. \frac{\partial R(\theta)}{\partial \theta} \right|_{\theta=\theta^m},$$

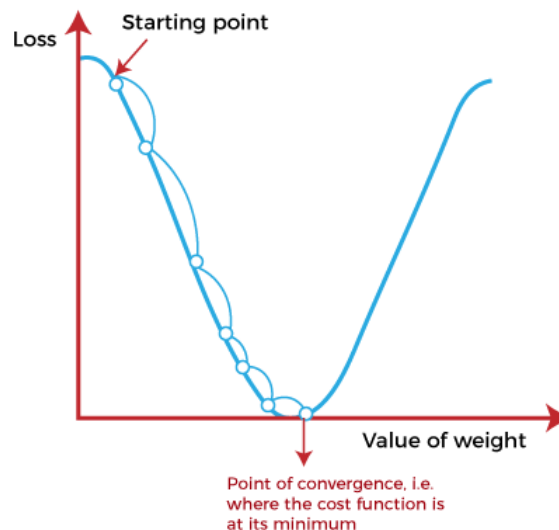
onde $\theta = \theta_m$ significa que o cálculo do gradiente está sendo realizado na iteração m .

E, para conseguir atualizar esse θ , conforme é calculado o gradiente durante as iterações, é utilizada a técnica de gradiente descendente, que pode ser descrita como:

$$\theta^{m+1} \leftarrow \theta^m - \lambda \frac{\partial R(\theta^m)}{\partial \theta^m}$$

sendo λ o parâmetro que vai definir a magnitude de influência da derivada $\frac{\partial R(\theta^m)}{\partial \theta^m}$ em θ^m .

Figura 6: Representação do método do gradiente descendente para a estimação de um parâmetro.



Fonte: <https://www.javatpoint.com/gradient-descent-in-machine-learning>

A Figura 6 demonstra o processo do gradiente descendente. Os parâmetros são iniciados com algum valor e, conforme ocorre os processos iterativos de aprendizado, o parâmetro converge para um mínimo da função de perda. Note que a distância entre cada ponto é definida pelo λ ou taxa de aprendizado.

Todo esse processo é realizado em cada parâmetro que existe na rede neural. Assim como as informações das variáveis são passadas camada a camada, saindo da camada de entrada, passando pelas camadas ocultas e chegando na camada de saída, visto anteriormente como *Forward propagation*, a informação do resultado da rede na função de perda é passada de forma contrária. O gradiente de cada parâmetro é calculado primeiro

nas camadas mais próximas da saída, e essa informação é repassada para trás, chegando até os parâmetros próximos aos da camada de entrada. Esse processo é chamado de *Backpropagation* (WERBOS, 1974).

2.3 SHAP

A estrutura de uma rede neural, por mais que proporcione bons resultados, mostra uma deficiência na parte interpretativa. Conhecida por ser uma "caixa-preta" pelo fato de sua estrutura ser muito complexa, existe a necessidade de se entender as previsões feitas. Para isso, existem técnicas que abordam o tema de interpretação de modelos de redes neurais e dentro delas existe a técnica SHAP, que através dela é possível entender como as variáveis de entrada influenciam as previsões do modelo, fornecendo resultados sobre sua lógica e permitindo uma explicação clara e confiável. Isso contribui para a transparência, confiabilidade e aceitação dos modelos, além de auxiliar na detecção de vieses e discriminação.

2.3.1 Valores de Shapley

Os valores de Shapley foram desenvolvidos por Lloyd Shapley (SHAPLEY, 1953) no contexto da teoria de jogos, e essa técnica ganhou força na área de inteligência artificial pela sua capacidade de conseguir interpretar modelos preditivos tidos como "caixa-preta".

No método criado por Shapley, existiam uma quantidade de jogadores que exerciam juntos determinada atividade, e o intuito era observar o ganho que um jogador (ou um conjunto de jogadores), obtinha ao ser adicionado para realizar a mesma tarefa, sem a presença do restante do grupo.

Podemos definir \mathbf{F} como o conjunto dos jogadores disponíveis para se realizar a tarefa, logo $\mathbf{F} = \{1, 2, \dots, \mathbf{M}\}$, onde \mathbf{M} é o número total de jogadores. Definindo \mathbf{S} como uma coligação do conjunto \mathbf{F} ($\mathbf{S} \subseteq \mathbf{F}$), temos, por exemplo, as seguintes possibilidades de \mathbf{S} , quando \mathbf{M} é igual a 3:

$$\{\{\emptyset\}, \{1\}, \{2\}, \{3\}, \{1, 2\}, \{1, 3\}, \{2, 3\}, \{1, 2, 3\}\}$$

Existe também uma função que vai mapear um conjunto de valores e retornar um número real, chamada de ν . Com isso, o retorno de $\nu(\mathbf{S})$ é um número real que pode ser definido como o "trabalho da coligação \mathbf{S} " ou o "trabalho dos jogadores presentes no conjunto \mathbf{S} ". Esse valor é equivalente ao total ganho que os jogadores podem obter caso trabalhem juntos em uma determinada coligação.

Para calcular o ganho ao adicionar i -ésimo jogador em uma tarefa, pode-se calcular o ganho quando é adicionada aquela jogador na coligação menos a coligação sem a adição daquele jogador, ficando da seguinte maneira:

$$\nu(\mathbf{S} \cup \{i\}) - \nu(\mathbf{S}) \quad (2.3.1)$$

No exemplo acima, para calcular o efeito da coligação $\{3\}$, pode-se realizar o seguinte processo:

$$\text{Contribuição de } \{3\} = \nu(\{1, 2, 3\}) - \nu(\{1, 2\})$$

Mas suponha que as variáveis (ou jogadores) $\{2\}$ e $\{3\}$ sejam extremamente semelhantes. Quando é calculado o ganho após inserir $\{2\}$ na coligação $\{1,2\}$, é possível notar um aumento substancial, mas quando é adicionado $\{3\}$ na coligação $\{1,2,3\}$, o ganho obtido é muito pouco. Como as variáveis exercem um papel parecido, o ganho maior ficou sujeito à variável que foi adicionada primeiro na coligação, não necessariamente porque uma é mais importante que a outra.

Por isso, para calcular o real ganho da variável $\{i\}$, é necessário testar todas as permutações de \mathbf{F} (conjunto de jogadores) e obter a contribuição de $\{i\}$ em cada uma delas, para então fazer a média dessas contribuições. Por exemplo, definindo $\mathbf{F} = \{1,2,3,4\}$, suponha que estamos interessados em calcular a contribuição de $\{3\}$, logo, podemos obter a seguinte permutação de \mathbf{F} :

$$[3, 1, 2, 4]$$

Calculando a contribuição de $\{3\}$, temos:

$$\nu(\{3\}) - \nu(\emptyset)$$

Outra permutação poderia ser :

$$[2, 4, 3, 1]$$

Calculando a contribuição de $\{3\}$, nessa permutação temos:

$$\nu(\text{coligação de } [2, 4, 3]) - \text{coligação de } [2, 4])$$

Uma observação deve ser feita: a função ν considera a coligação como argumento,

não a permutação. A coligação é um conjunto, com isso a ordem dos elementos não importa, mas a permutação é uma coleção ordenada de elementos. Na permutação do tipo $[3,1,2,4]$, 3 é a primeira variável adicionada e 4 é a última. Por isso, para cada permutação a ordem dos elementos pode mudar a contribuição do total ganho, contudo o total ganho da permutação somente depende dos elementos, não da ordem. Logo:

$$\nu(\text{coligação de } [3, 1, 2, 4]) = \nu(\{1, 4, 2, 3\})$$

Sendo assim, para cada permutação \mathbf{P} , é preciso primeiro calcular o ganho da coligação das variáveis que foram adicionadas antes de $\{i\}$, e esse conjunto pode ser chamado de coligação \mathbf{S} . Feito isso, agora é preciso calcular o ganho das coligações que são formadas ao adicionar $\{i\}$ em \mathbf{S} , e podemos chamar isso de $\mathbf{S} \cup \{i\}$. Com isso, a contribuição da variável $\{i\}$, denotada por ϕ_i , é:

$$\phi_i = \frac{1}{|\mathbf{F}|!} \sum_{\mathbf{P}} [\nu(\mathbf{S} \cup \{i\}) - \nu(\mathbf{S})] \quad (2.3.2)$$

O número total de permutações de \mathbf{F} é $|\mathbf{F}|!$. Logo, podemos dividir a soma das contribuições por $|\mathbf{F}|!$ para encontrar o valor esperado de contribuição de $\{i\}$. A Figura 7 mostra como é feito esse calculo para um determinado jogador $\{i\}$.

Figura 7: Ganho do jogador 3 em relação à todas as permutações de jogadores.

	\mathbf{P}	$\nu(\mathbf{S} \cup \{i\}) - \nu(\mathbf{S})$	$i=3$
$ \mathbf{F} !$	$[1, 2, 3, 4, 5]$	$\nu(\{1, 2, 3\}) - \nu(\{1, 2\})$	
	$[2, 1, 3, 4, 5]$	$\nu(\{1, 2, 3\}) - \nu(\{1, 2\})$	
	$[3, 1, 2, 4, 5]$	$\nu(\{3\})$	
	\dots	\dots	
	$[1, 2, 4, 5, 3]$	$\nu(\{1, 2, 3, 4, 5\}) - \nu(\{1, 2, 4, 5\})$	

$$\phi_i = \frac{1}{|\mathbf{F}|!} \sum_{\mathbf{P}} (\nu(\mathbf{S} \cup \{i\}) - \nu(\mathbf{S}))$$

Fonte: <https://towardsdatascience.com/introduction-to-shap-values-and-their-application-in-machine-learning-8003718e6827>

É possível perceber que algumas permutações possuem a mesma contribuição, desde que suas coligações $\mathbf{S} \cup \{i\}$ e \mathbf{S} sejam as mesmas. Com isso, para reduzir o processo do cálculo de contribuição de cada permutação, pode-se identificar quantas vezes a

permutação gerada vai resultar em uma contribuição que seja igual a outra.

Para fazer isso, é necessário descobrir quantas permutações podem ser formadas de cada coligação. Podemos definir $\mathbf{F} - \{i\}$ como o conjunto de todas as variáveis excluindo a variável $\{i\}$, e \mathbf{S} como uma das coligações de $\mathbf{F} - \{i\}$ ($\mathbf{S} \subseteq \mathbf{F} - \{i\}$).

Logo, para cada coligação \mathbf{S} temos $|\mathbf{S}|!$ possíveis permutações, que corresponde às possibilidades de variáveis e suas respectivas ordens antes de adicionar a variável $\{i\}$.

Tendo os conjuntos $\mathbf{S} \cup \{i\}$ e \mathbf{S} definidos, resta agora achar as possíveis permutações das variáveis restantes. E para saber o valor restante é preciso calcular o tamanho do conjunto gerado por: $\mathbf{F} - (\mathbf{S} \cup \{i\} + 1)$ que basicamente é o que resta das variáveis para completar o conjunto \mathbf{F} .

A Figura 8 mostra o que acontece quando se escolhe o jogador i , nesse caso $i = 3$. Note que, na linha das coligações, é definido as possíveis coligações de \mathbf{S} , que seria as permutações dos jogadores 1 e 2, temos a coligação de um único elemento na coluna $\{i\}$, que sempre vai ser o próprio elemento, em seguida a coligação dos jogadores restantes. Na linha das permutações é definida todas as possíveis permutações para \mathbf{S} , para $\{i\}$ e para $\mathbf{F} - \mathbf{S} - \{i\}$. E na última linha é representado o tamanho do conjunto formado pela permutação/coligação descrita anteriormente.

Figura 8: Relação entre permutações e coalizões.

Coalitions	S	+	$\{i\}$	+	$F-S-\{i\}$	=	F
	$\{1, 2\}$	+	$\{3\}$	+	$\{4, 5\}$	=	$\{1, 2, 3, 4, 5\}$
Permutations					$[4, 5]$	=	$[1, 2, 3, 4, 5]$
	$[1, 2]$				$[5, 4]$	=	$[1, 2, 3, 5, 4]$
		+	$[3]$	+			
	$[2, 1]$				$[4, 5]$	=	$[2, 1, 3, 4, 5]$
					$[5, 4]$	=	$[2, 1, 3, 5, 4]$
Number of Permutations	$ S !$	+	1	+	$(F - S -1)!$	=	$ S !(F - S -1)!$

Fonte: <https://towardsdatascience.com/introduction-to-shap-values-and-their-application-in-machine-learning-8003718e6827>

Com isso, podemos reescrever a Equação 2.3.2 da seguinte maneira:

$$\phi_i = \sum_{\mathbf{S} \subseteq \mathbf{F} - \{i\}} \frac{|\mathbf{S}|!(|\mathbf{F}| - |\mathbf{S}| - 1)!}{|\mathbf{F}|!} [\nu(\mathbf{S} \cup \{i\}) - \nu(\mathbf{S})],$$

onde ϕ_i é o valor de shapley para a variável $\{i\}$.

2.3.2 *Shapley Additive Explanations*

Fazendo a relação do valor de Shapley para o SHAP (*Shapley Additive Explanations*), temos que a função característica ν é equivalente à função $f(x)$ responsável por fazer as predições. E os valores de SHAP são calculados a partir das observações que entram no modelo. Com isso, a fórmula do valor de SHAP, para cada conjunto de observação e variável especificada, se dá por:

$$\phi_i(f, \mathbf{x}) = \sum_{\mathbf{S} \subseteq \mathbf{F} - \{i\}} \frac{|\mathbf{S}|!(|\mathbf{F}| - |\mathbf{S}| - 1)!}{|\mathbf{F}|!} [f_{\mathbf{S} \cup \{i\}}(\mathbf{x}_{\mathbf{S} \cup \{i\}}) - f_{\mathbf{S}}(\mathbf{x}_{\mathbf{S}})] \quad (2.3.3)$$

Perceba que $f_{\mathbf{S}}(\mathbf{x}_{\mathbf{S}})$ representa o resultado do modelo com somente as variáveis que estão na coligação \mathbf{S} , algo que na realidade não é permitido na maioria dos modelos. Por isso, uma aproximação desse resultado é a seguinte:

$$f_{\mathbf{S}}(\mathbf{x}_{\mathbf{S}}) \approx E[f(\mathbf{x}|\mathbf{x}_{\mathbf{S}})] \approx \frac{1}{k} \sum_{i=1}^k f(\mathbf{x}_{\mathbf{S}}^{(i)}, \mathbf{x}_{\mathbf{S}}) \quad (2.3.4)$$

Para determinar os valores das variáveis que não pertencem ao conjunto $\mathbf{x}_{\mathbf{S}}$, a Equação 2.3.4 esclarece que para isso o valor esperado dessas variáveis é calculado. Isso resulta na criação de uma estimativa para $f(\mathbf{x})$, levando em consideração apenas as variáveis presentes em \mathbf{S} .

Figura 9: Cálculo de f_S , sendo S o conjunto de variáveis X_1, X_3, X_4 , dentre as observações de um conjunto de dados.

$$\mathbf{x} = \{x_1, x_2, x_3, x_4, x_5\} \quad \mathbf{x}_S = \{x_1, x_3, x_4\} \quad \mathbf{x}_{\bar{S}} = \{x_2, x_5\}$$

X_1	X_2	X_3	X_4	X_5
$x_1^{(1)}$	$x_2^{(1)}$	$x_3^{(1)}$	$x_4^{(1)}$	$x_5^{(1)}$
$x_1^{(2)}$	$x_2^{(2)}$	$x_3^{(2)}$	$x_4^{(2)}$	$x_5^{(2)}$
...
$x_1^{(k)}$	$x_2^{(k)}$	$x_3^{(k)}$	$x_4^{(k)}$	$x_5^{(k)}$

$$\begin{aligned}
 & \frac{f(\mathbf{x}_{\bar{S}}^{(i)}, \mathbf{x}_S)}{f(x_1, x_2^{(1)}, x_3, x_4, x_5^{(1)})} \\
 & \frac{f(x_1, x_2^{(2)}, x_3, x_4, x_5^{(2)})}{\dots} \\
 & + \frac{f(x_1, x_2^{(k)}, x_3, x_4, x_5^{(k)})}{\sum_{i=1}^k f(\mathbf{x}_{\bar{S}}^{(i)}, \mathbf{x}_S)}
 \end{aligned}$$

$$f_S(\mathbf{x}_S) \approx E[f(\mathbf{x})|\mathbf{x}_S] \approx \frac{1}{k} \sum_{i=1}^k f(\mathbf{x}_{\bar{S}}^{(i)}, \mathbf{x}_S)$$

Fonte: <https://towardsdatascience.com/introduction-to-shap-values-and-their-application-in-machine-learning-8003718e6827>

Ao analisar a Figura 9, tem-se que as variáveis X_1, X_3 e X_4 representam o conjunto S . E escolhendo um valor x_1, x_3 e x_4 dessas variáveis, respectivamente, é definido o conjunto \mathbf{x}_S . Para obter $f_S(\mathbf{x}_S)$, que é calcular o resultado do modelo somente com as variáveis presentes em \mathbf{x}_S , calcula-se f para cada observação do conjunto de dados, travando x_1, x_3 e x_4 na função e utilizando o valor das variáveis complementares, que neste caso são os valores de X_2 e X_5 , em suas respectivas observações. Feito isso, é calculada média desses valores que corresponde justamente com o resultado da função $f_S(\mathbf{x}_S)$.

Com $f_S(\mathbf{x}_S)$ estimado, é possível calcular a Equação 2.3.3. Os dados que a técnica SHAP utiliza são divididos em 2 tipos: as observações necessárias para a estimação de $f_S(\mathbf{x}_S)$ e as observações em que se deseja calcular o valor de SHAP.

2.4 Medidas de associação

A utilização de medidas associativas, como o coeficiente de Pearson para variáveis contínuas e o coeficiente de contingência para variáveis categóricas, desempenha um papel crucial na análise estatística, proporcionando a visualização das relações entre diferentes variáveis. Essas medidas quantificam a força e a natureza das associações, permitindo a identificação de padrões e tendências nos dados (WITTE; WITTE, 2017).

2.4.1 Coeficiente de Pearson

O coeficiente de Pearson, também conhecido como correlação de Pearson, é uma medida estatística que avalia a relação linear entre duas variáveis contínuas.

$$\rho = \frac{cov(X, Y)}{\sqrt{var(X)var(Y)}} \quad (2.4.1)$$

Este coeficiente varia de -1 a 1, onde -1 indica uma relação linear negativa perfeita, 1 indica uma relação linear positiva perfeita, e 0 indica ausência de relação linear.

2.4.2 Coeficiente de contingência

O coeficiente de contingência é uma medida estatística utilizada para avaliar a associação entre duas variáveis categóricas em uma tabela de contingência. Ele é especialmente útil quando se trabalha com dados categóricos, fornecendo uma indicação da força e direção da associação entre as variáveis.

$$C = \sqrt{\frac{\chi^2}{N + \chi^2}} \quad (2.4.2)$$

O coeficiente de contingência varia de 0 a 1, onde 0 indica nenhuma associação e 1 indica uma associação perfeita. Valores mais próximos de 1 indicam uma forte relação entre as variáveis, enquanto valores mais próximos de 0 sugerem independência.

2.5 Métricas de avaliação

A avaliação do desempenho dos modelos é fundamental para obter insights sobre sua eficácia nas previsões ou classificações. A utilização de estratégias que resumem o desempenho por meio de métricas específicas é crucial nesse processo. A análise dessas métricas proporciona uma compreensão mais aprofundada do modelo, permitindo identificar pontos fortes e áreas de melhoria. Essa avaliação não apenas valida a qualidade das previsões, mas também orienta os próximos passos na pesquisa, direcionando ajustes necessários no modelo ou indicando caminhos para refinamento. Dessa forma, a escolha e interpretação adequadas das métricas são passos essenciais para uma avaliação informada e um progresso significativo na pesquisa (GERON, 2019).

2.5.1 Matriz de confusão

Uma matriz de confusão é uma tabela usada para avaliar o desempenho de um modelo de classificação. Seu papel é de expor os resultados das predições do modelo quando comparadas com os valores reais.

A matriz de confusão organiza as previsões do modelo em quatro categorias, comumente chamadas de Verdadeiro Positivo (VP), Falso Positivo (FP), Verdadeiro Negativo (VN) e Falso Negativo (FN). Essas categorias são definidas da seguinte maneira:

- Verdadeiro Positivo (VP): Exemplos que foram corretamente classificadas como pertencentes à classe positiva.
- Falso Positivo (FP): Exemplos que foram erroneamente classificadas como pertencentes à classe positiva, quando na verdade pertencem à classe negativa.
- Verdadeiro Negativo (VN): Exemplos que foram corretamente classificadas como pertencentes à classe negativa.
- Falso Negativo (FN): Exemplos que foram erroneamente classificadas como pertencentes à classe negativa, quando na verdade pertencem à classe positiva.

		Previsão	
		Negativo	Positivo
Real	Negativo	Verdadeiro Negativo (VN)	Falso Positivo (FP)
	Positivo	Falso Negativo (FN)	Verdadeiro Positivo (VP)

Tabela 1: Matriz de confusão

2.5.2 Acurácia

A acurácia é a proporção de predições corretas feitas por um modelo em relação ao número total de predições. A fórmula básica para calcular a acurácia é dada por:

$$\text{Acurácia} = \frac{\text{Número de predições corretas}}{\text{Número total de predições}} = \frac{VP + VN}{VP + VN + FP + FN} \quad (2.5.1)$$

Essa métrica fornece uma visão geral do desempenho do modelo, indicando a porcentagem de instâncias corretamente classificadas. No entanto, a acurácia pode ser enganosa em casos onde as classes não estão balanceadas. Em situações desse tipo, um modelo que prevê sempre a classe majoritária pode ter uma acurácia alta, mesmo que não seja eficaz.

2.5.3 Precisão

A precisão é definida como a proporção de exemplos classificados corretamente como positivos em relação ao total de exemplos classificadas como positivas (verdadeiras positivas mais falsos positivos).

$$\text{Precisão} = \frac{VP}{VP + FP} \quad (2.5.2)$$

A precisão é particularmente útil quando os falsos positivos são mais problemáticos ou custosos em comparação com os falsos negativos. Por exemplo, em um sistema de detecção de spam, classificar erroneamente um e-mail legítimo como spam (falso positivo) pode ser mais prejudicial do que deixar passar um e-mail de spam (falso negativo).

2.5.4 Recall

O recall, também conhecido como sensibilidade, é outra métrica utilizada no contexto de classificação, focada em capturar a proporção de exemplos positivos que foram corretamente identificadas pelo modelo em relação ao total de exemplos positivos existentes.

$$\text{Recall} = \frac{VP}{VP + FN} \quad (2.5.3)$$

O recall é especialmente útil quando os falsos negativos (exemplos positivos não identificadas pelo modelo) são mais críticos ou custosos do que os falsos positivos. Por exemplo, em um sistema de detecção de fraudes, é crucial identificar todas as transações fraudulentas, mesmo que isso signifique aceitar algumas transações normais erroneamente classificadas como fraudulentas.

2.5.5 F1-score

O F1-score é uma métrica de avaliação que combina as métricas de precisão e recall em um único valor.

$$F1 = 2 \cdot \frac{\text{Precisão} \cdot \text{Recall}}{\text{Precisão} + \text{Recall}} \quad (2.5.4)$$

O F1-score é a média harmônica entre a precisão e o recall. A média harmônica é utilizada porque penaliza extremos, sendo particularmente sensível a baixos valores em qualquer uma das métricas.

O F1-score varia de 0 a 1, onde 1 indica o melhor desempenho possível, equilibrando tanto a precisão quanto o recall. Essa métrica é particularmente útil quando há um desequilíbrio significativo entre as classes, pois é menos sensível a grandes quantidades de verdadeiros negativos.

3 Metodologia

A metodologia adotada nesta pesquisa se fundamenta na análise do conjunto de dados *"Loan Data for Dummy"*, retirado da plataforma *Kaggle*, visando a compreensão e modelagem de padrões associados a operações de empréstimos. Dois métodos distintos, Regressão Logística e Redes Neurais, serão empregados para investigar as relações existentes nos dados e aprimorar as previsões. A implementação desses modelos será realizada utilizando tanto a linguagem de programação R quanto Python. Além disso, a técnica SHAP (*Shapley Additive Explanations*) será integrada para proporcionar uma interpretação aprofundada do modelo de redes neurais, ampliando a transparência nas decisões preditivas.

3.1 Conjunto de dados

O banco de dados *"Loan Data for Dummy"* é uma base de dados do Kaggle, projetada para simular informações relacionadas a operações de empréstimos. Desenvolvido para fins educacionais e de pesquisa, esse conjunto tem sua origem de um modelo de banco *"peer to peer"* sediado na Irlanda, no qual o banco disponibiliza recursos a potenciais clientes, obtendo lucros com base no risco que assume. Os dados disponíveis no Kaggle representam uma versão fictícia de uma situação real, com a maior parte dos dados manipulados ou criados sinteticamente para preservar as informações dos clientes originais.

A variável que será foco do estudo é a "Condição do empréstimo". Através dessa variável, é possível discernir se um empréstimo foi classificado como "bom" ou "ruim", proporcionando uma avaliação da qualidade e risco associados a cada transação. No contexto deste conjunto de dados, presume-se que a "Condição do empréstimo" seja uma variável binária, onde, por exemplo, "0" poderia indicar um empréstimo em boas condições e "1" indicaria o contrário. A compreensão aprofundada dessa variável é essencial para a construção e interpretação adequada dos modelos subsequentes, como a regressão logística e redes neurais, permitindo uma análise mais precisa e informada do risco associado aos empréstimos.

3.1.1 Variáveis

A base de dados é composta por 30 variáveis, incluindo a variável resposta, e existem 887379 observações. Para se avaliar a variável "Condição do empréstimo" foi utilizada algumas variáveis presentes na base de dados, como:

1. **Tempo de emprego:** Representa o tempo de emprego do solicitante expresso numericamente. Um valor de 5 indicaria que o indivíduo está empregado há 5 anos.
2. **Tipo de residência:** Indica o status de moradia do solicitante, como proprietário, inquilino ou outra forma de ocupação residencial.
3. **Renda anual:** Reflete a renda anual do solicitante, uma medida crucial para avaliar a capacidade de pagamento do empréstimo. Pode ser expressa numericamente, por exemplo, 50,000.
4. **Valor do empréstimo:** Representa o valor do empréstimo solicitado pelo requerente, geralmente expresso em termos monetários, como 10,000.
5. **Prazo:** Indica o prazo do empréstimo, especificando o período de tempo durante o qual o empréstimo deve ser reembolsado. Pode ser, por exemplo, 36 meses.
6. **Tipo de aplicação:** Refere-se ao tipo de aplicação, indicando se é uma aplicação individual ou conjunta.
7. **Finalidade:** Descreve a finalidade do empréstimo, como consolidação de dívidas, compra de casa, educação, entre outros.
8. **Tipo do juros:** Indica a natureza dos pagamentos de juros, se são fixos ou variáveis.
9. **Taxa de juros:** Representa a taxa de juros associada ao empréstimo, geralmente expressa como uma porcentagem, como 10.
10. **Grau:** Refere-se à classificação de risco do tomador de empréstimo atribuída pela instituição financeira, como A, B, C, etc.
11. **DTI:** Significa "Debt-to-Income" (Dívida-para-Renda) e representa a proporção entre as dívidas mensais e a renda mensal do requerente, proporcionando uma medida da capacidade de pagamento.
12. **Valor bruto pago:** Representa o valor total pago, incluindo o principal e os juros, ao final do empréstimo.
13. **Valor líquido pago:** Indica o total de principal (quantia inicial do empréstimo) recuperado até o momento.
14. **Valor recuperado:** Representa o valor recuperado em caso de inadimplência ou perda.
15. **Parcelas:** Refere-se à parcela mensal que o requerente do empréstimo deve pagar, incluindo tanto o principal quanto os juros.
16. **Região:** Indica a região geográfica associada ao requerente do empréstimo.

3.1.2 Limpeza dos dados

Para diminuir a complexidade da base de dados, as variáveis passaram por 3 critérios de avaliação antes de serem utilizadas nos modelos:

1. Identificação de variáveis que não impactariam o resultado do modelo;
2. Comparação de variáveis que gerem a mesma informação;
3. Extração de variáveis presentes em apenas uma das categorias da variável resposta.

O item 1. destaca-se a variável "ID" como independente da variável resposta, atuando unicamente como identificador do cliente, sem exercer influência no resultado final do modelo.

O item 2. se refere à situação da base de dados em que o autor realizou uma rotulação numérica de variáveis já categorizadas, como por exemplo: "Tipo de juros" e "Tipo de juros Cat", onde na primeira variável tem as opções "Juros simples" e "Juro compostos" e na segunda variável o autor associa os números "1" e "2" respectivamente a essas variáveis.

O item 3. esclarece variáveis que desempenham funções em apenas uma das categorias da variável resposta. Como é o caso da variável "Recuperações totais" onde a mesma é presente apenas no caso do cliente ter sido inadimplente, se relacionando com a categoria "Empréstimo ruim" da variável resposta. Para o caso de "Empréstimo bom" os valores da variável estão zerados.

Com isso a base de dados ficou da seguinte maneira:

Base de dados	Número de Colunas
Antes da extração	30
Depois da extração	18

Tabela 2: Número de colunas antes e depois da preparação dos dados

Com o tratamento dos dados realizado, foi feita a separação dos dados para a modelagem dos dois modelos.

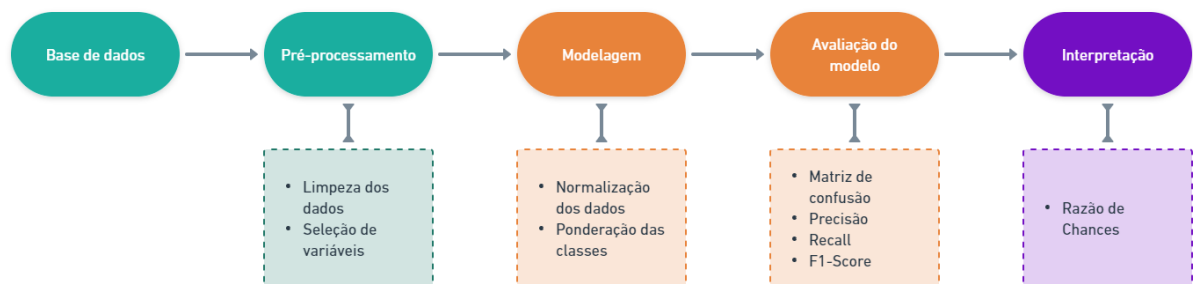
Tipo de dado	% dos dados originais
Treino	70%
Validação	10%
Teste	20%

Tabela 3: Divisão dos dados para a modelagem dos modelos de regressão logística e redes neurais

3.2 Modelagem dos dados

3.2.1 Regressão logística

Figura 10: Fluxograma da implementação do modelo logístico



Fonte: Autoria própria

A metodologia para a implementação da regressão logística incluiu várias etapas essenciais para garantir a robustez e eficácia do modelo. Inicialmente, foi realizada uma cuidadosa etapa de pré-processamento, que envolveu a limpeza e tratamento das variáveis. Durante essa fase, foram identificados e tratados possíveis valores ausentes, outliers e erros nos dados, contribuindo para a qualidade do conjunto de dados.

Outro aspecto crítico foi a categorização adequada das variáveis, quando aplicável. Isso incluiu a transformação de variáveis categóricas em formatos adequados para análise estatística, garantindo que todas as variáveis estivessem em uma forma consistente para o modelo de regressão logística.

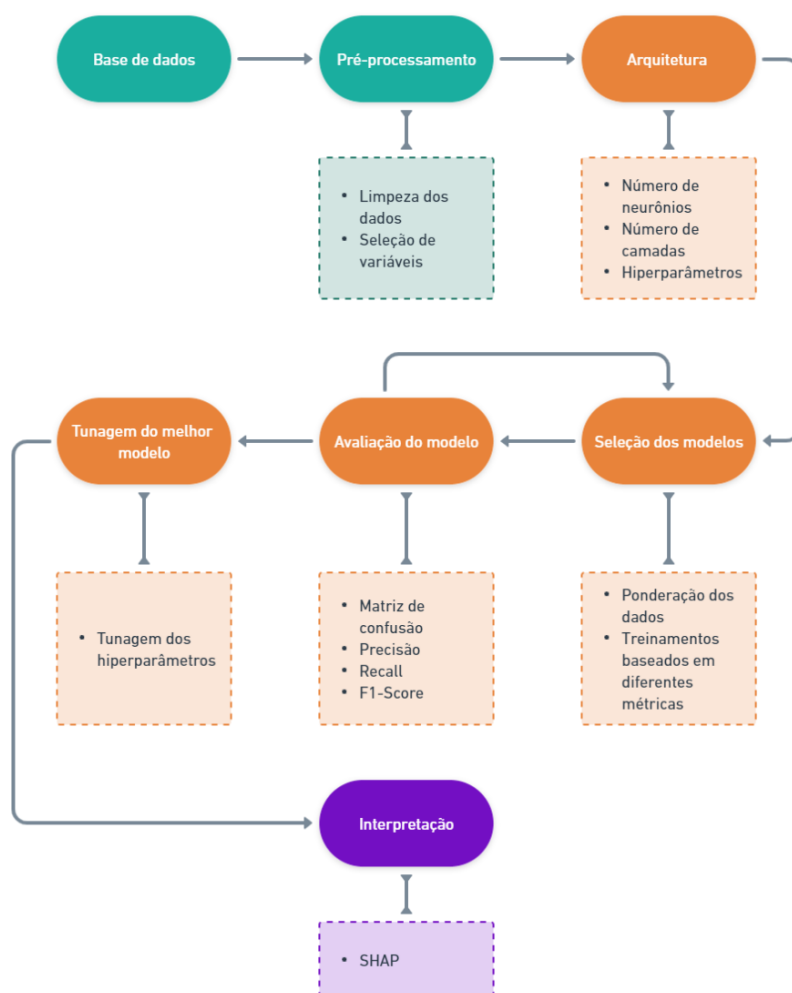
Para garantir que as variáveis contribuíssem igualmente para o modelo, os dados foram normalizados, usando uma padronização com uma escala com média zero e desvio padrão unitário.

Dada a natureza desbalanceada do conjunto de dados, onde as classes "Empréstimo bom" e "Empréstimo ruim" têm proporções significativamente diferentes, foi aplicada a ponderação das classes. Isso foi realizado durante a fase de treinamento do modelo, atribuindo pesos diferentes às classes para compensar o desequilíbrio e garantir que o modelo não fosse enviesado em direção à classe majoritária.

Tendo o modelo definido, foi feita uma avaliação das predições do modelo no conjunto de teste, usando métricas apropriadas para problemas de classificação, como precisão, recall, F1-score e matriz de confusão. Essas métricas permitiram uma compreensão abrangente do desempenho do modelo, especialmente no que diz respeito à capacidade de prever corretamente os casos de "Empréstimo ruim" e "Empréstimo bom".

3.2.2 Rede neural

Figura 11: Fluxograma da implementação da rede neural



Fonte: Autoria própria

A implementação da metodologia para a construção e ajuste de redes neurais envolveu diversas etapas cruciais para garantir a eficácia do modelo. Inicialmente, foi realizada uma fase de pré-processamento, que consistiu na limpeza e tratamento das variáveis. Durante esse estágio, foram tratados valores ausentes, outliers e possíveis erros nos dados, contribuindo para a qualidade do conjunto de dados. Essa etapa foi a mesma apresentada pelo modelo logístico.

A escolha da arquitetura inicial foi um passo significativo. Foi necessário definir o número de camadas ocultas, a quantidade de neurônios em cada camada e a função de ativação a ser utilizada. Essas escolhas iniciais foram baseadas tanto em conhecimentos

prévios do problema quanto em experimentações para encontrar a configuração que melhor se adequava aos dados. E para isso foi utilizado uma rede neural com a seguinte arquitetura:

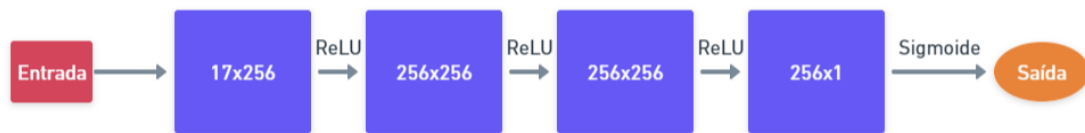
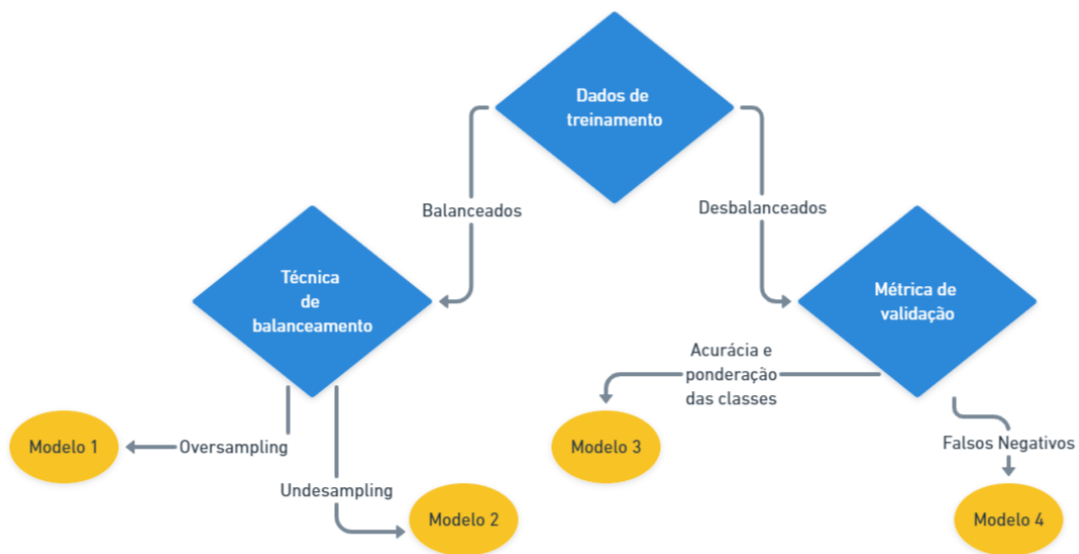


Figura 12: Arquitetura de rede neural inicial

Diversos modelos de redes neurais foram desenvolvidos e treinados em contextos variados, considerando abordagens distintas, tais como:

- Utilização de dados equilibrados para realizar os treinamentos;
- Modelagem dos modelos baseadas em métricas de treinamento específicas;
- Ponderação das classes da variável resposta.

Figura 13: Fluxograma da escolha dos modelos de redes neurais estimados



Fonte: Autoria própria

Para realizar o treinamento dos modelos com dados equilibrados, foram empregados dois processos distintos: *upsampling* e *oversampling*. No âmbito da modelagem, em que se priorizaram métricas específicas de treinamento, além da acurácia, métrica padrão, foram incorporados o Recall do modelo no conjunto de validação e o número de Falsos Negativos. Este último é particularmente crucial em cenários de empréstimos, sendo considerado o erro mais relevante. Por fim, foi utilizada também um modelo que aplicasse

ponderação nas classes, buscando assim equilibrar a classe majoritária. Essas abordagens visaram explorar diferentes aspectos do treinamento para encontrar a configuração mais eficaz, adaptada às situações do problema em análise.

Cada modelo foi treinado utilizando o conjunto de treinamento e validado para avaliar seu desempenho. A métrica de desempenho, geralmente relacionada à precisão, recall e F1-score, foi usada para comparar e selecionar os melhores modelos.

O modelo mais promissor foi então submetido a uma etapa de tunagem de hiperparâmetros. Ajustes finos nos hiperparâmetros, como taxa de aprendizado, número de épocas de treinamento e tamanho do lote, foram realizados para otimizar o desempenho do modelo. A arquitetura final do modelo ficou como:

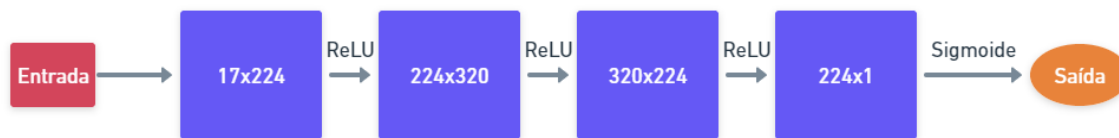


Figura 14: Arquitetura da rede neural depois do processo de tunelagem

Sendo o modelo final escolhido aquele treinado com as classes ponderadas e baseado na métrica da acurácia.

A avaliação final do modelo ocorreu no conjunto de validação e teste. Essa etapa foi crucial para garantir que o modelo não apenas se ajustasse bem aos dados de treinamento, mas também generalizasse de maneira eficaz para novos dados. As métricas de desempenho foram novamente utilizadas para avaliar a capacidade do modelo de fazer previsões precisas e úteis.

3.3 Interpretação dos modelos

Definindo os modelos logístico e o de redes neurais, foi feita a interpretação de ambos. Realizar a interpretação dos modelos é uma etapa crucial na análise de dados, mostrando o funcionamento e as relações das variáveis no contexto do problema em questão.

No caso do modelo logístico, a interpretação se concentrou nos parâmetros estimados para cada variável. Esses coeficientes fornecem uma medida da magnitude e direção da influência de cada variável na predição da variável resposta. Além disso, a interpretação envolveu a análise da Razão de Chances (RC), que expressa como a chance de o evento ocorrer se torna multiplicativamente maior ou menor com a mudança em uma unidade na variável explicativa.

Para a interpretação da rede neural, utilizou-se a técnica SHAP (*SHapley Additive exPlanations*). Essa abordagem proporciona uma compreensão mais profunda ao atribuir a contribuição de cada variável para a saída do modelo em nível individual. Com o auxílio de gráficos SHAP, pôde-se observar como cada variável influencia as predições e identificar padrões de comportamento em diferentes cenários.

4 Resultados

Esta seção inicia explorando a relação entre as variáveis explicativas e a variável resposta. Os resultados derivados tanto do modelo logístico quanto da rede neural serão detalhadamente apresentados, destacando a ênfase na interpretação de ambos os modelos. Além disso, será conduzido um benchmark comparativo entre as duas abordagens, proporcionando uma análise crítica de suas performances.

4.1 Análise descritiva

4.1.1 Condição do empréstimo

A variável "Condição do empréstimo" é a variável resposta desse estudo, como foi definido anteriormente. Com isso temos o seguinte comportamento dessa variável:

Condição do empréstimo	Número de observações	Frequência relativa
Empréstimo bom	819950	92,4%
Empréstimo ruim	67429	7,59%

Tabela 4: Número de observação em cada categoria da variável resposta

A Tabela 4 mostra a distribuição da variável "Condição do empréstimo". Uma variável composta majoritariamente por observações do tipo "Empréstimo bom", onde a mesma está presente em mais de 90% das observações na base de dados, mostrando que a cada 12 empréstimos rotulados como "bons", existe 1 rotulado como "ruim".

4.1.2 Relação entre as covariáveis e a variável resposta

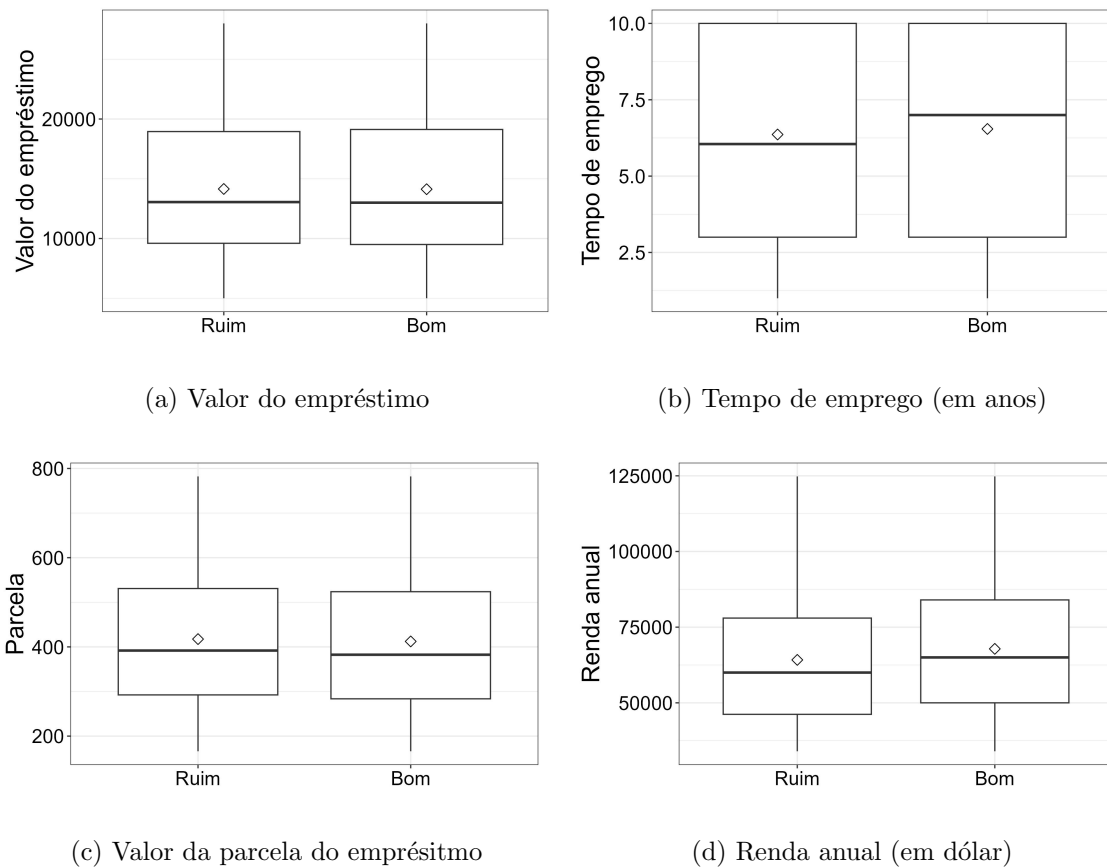
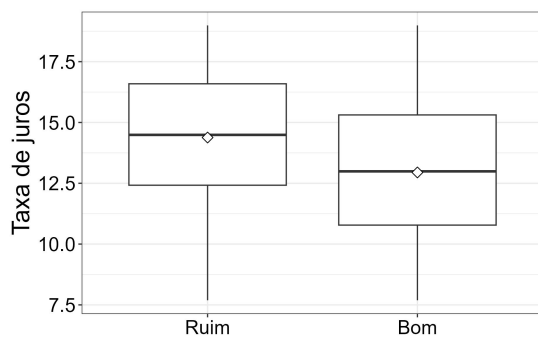
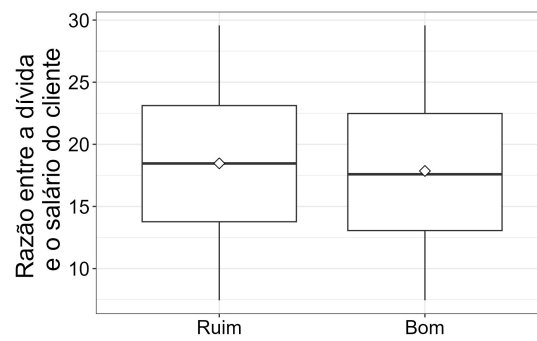


Figura 15: Variáveis explicativas em relação à condição do empréstimo

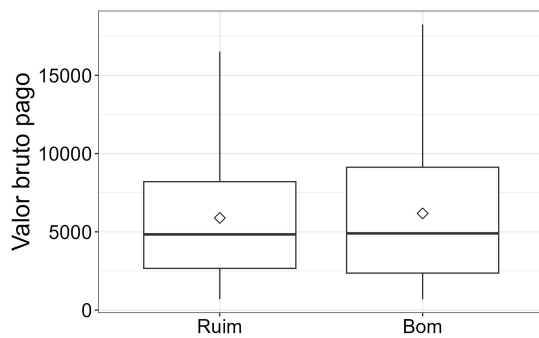
O comportamento da variável resposta nas Figuras 15a e 15c demonstrou semelhanças, onde, em ambos os casos, não foi evidenciada uma clara diferença entre o valor do empréstimo e o valor da parcela em relação às categorias da variável resposta. A Figura 15b também apresenta um comportamento semelhante entre as classes "Empréstimo ruim" e "Empréstimo bom", mas com um detalhe: a mediana do tempo de trabalho dos clientes rotulados como "Empréstimo ruim" foi inferior em comparação ao outro caso. Por fim, a Figura 15d indica que clientes com uma renda anual elevada tendem a ser categorizados como "Empréstimo bom".



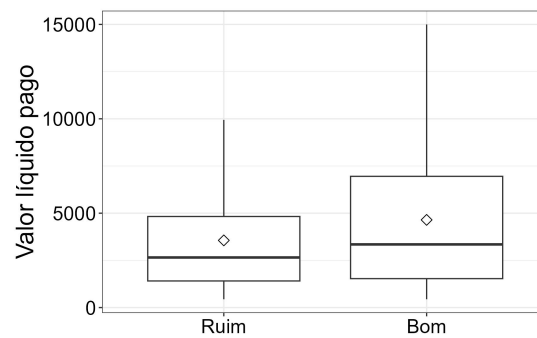
(a) Taxa de juros do empréstimo



(b) Razão entre a dívida e o salário do cliente



(c) Valor bruto do empréstimo pago



(d) Valor líquido do empréstimo pago

Figura 16: Variáveis explicativas em relação à condição do empréstimo

A Figura 16a evidencia uma relação significativa entre taxas de juros elevadas e empréstimos considerados ruins. A Figura 16b complementa a informação fornecida pela Figura 15d, indicando que clientes com renda mais elevada tendem a cumprir adequadamente com seus pagamentos. As Figuras 16c e 16d seguem padrões semelhantes, sugerindo que clientes que quitaram a maior parte do empréstimo são frequentemente rotulados como bons pagadores.

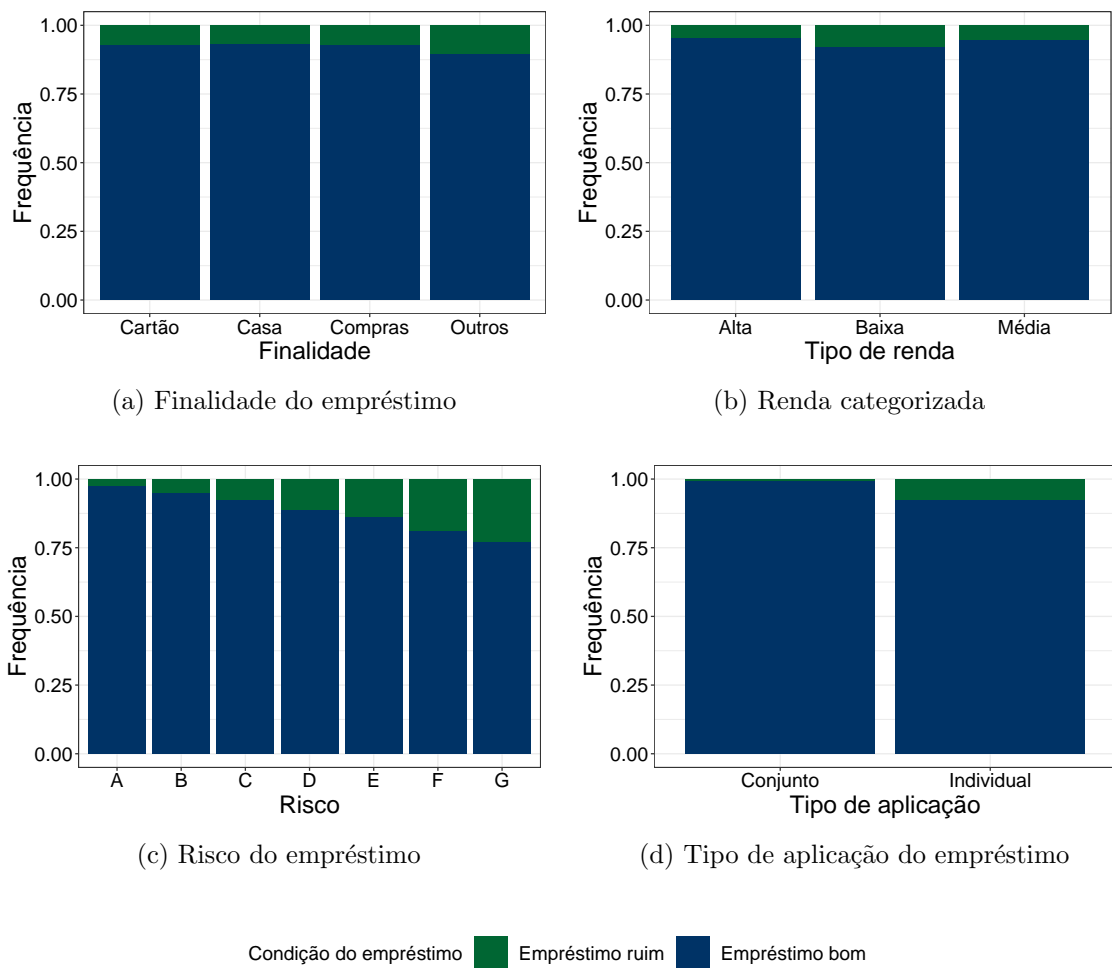


Figura 17: Variáveis explicativas em relação à condição do empréstimo

A Figura 17a ilustra que as categorias da variável "Finalidade" seguem a proporção natural da condição do empréstimo, conforme indicado na Tabela de Condição do Empréstimo. Na Figura 17b, as categorias "Alta" e "Média" exibem proporções menores de empréstimos ruins em comparação com a categoria "Baixa", que apresenta uma proporção de quase 10% de empréstimos ruins. A Figura 17c revela um padrão de "cascata", indicando que à medida que o risco do empréstimo aumenta, a proporção de empréstimos ruins nas últimas categorias também aumenta, sendo a categoria G a mais afetada, com quase 25% de empréstimos classificados como ruins. Na Figura 17d, a categoria "Empréstimo conjunto" não registrou observações de empréstimos ruins, concentrando a maioria desses empréstimos na categoria "Empréstimo individual".

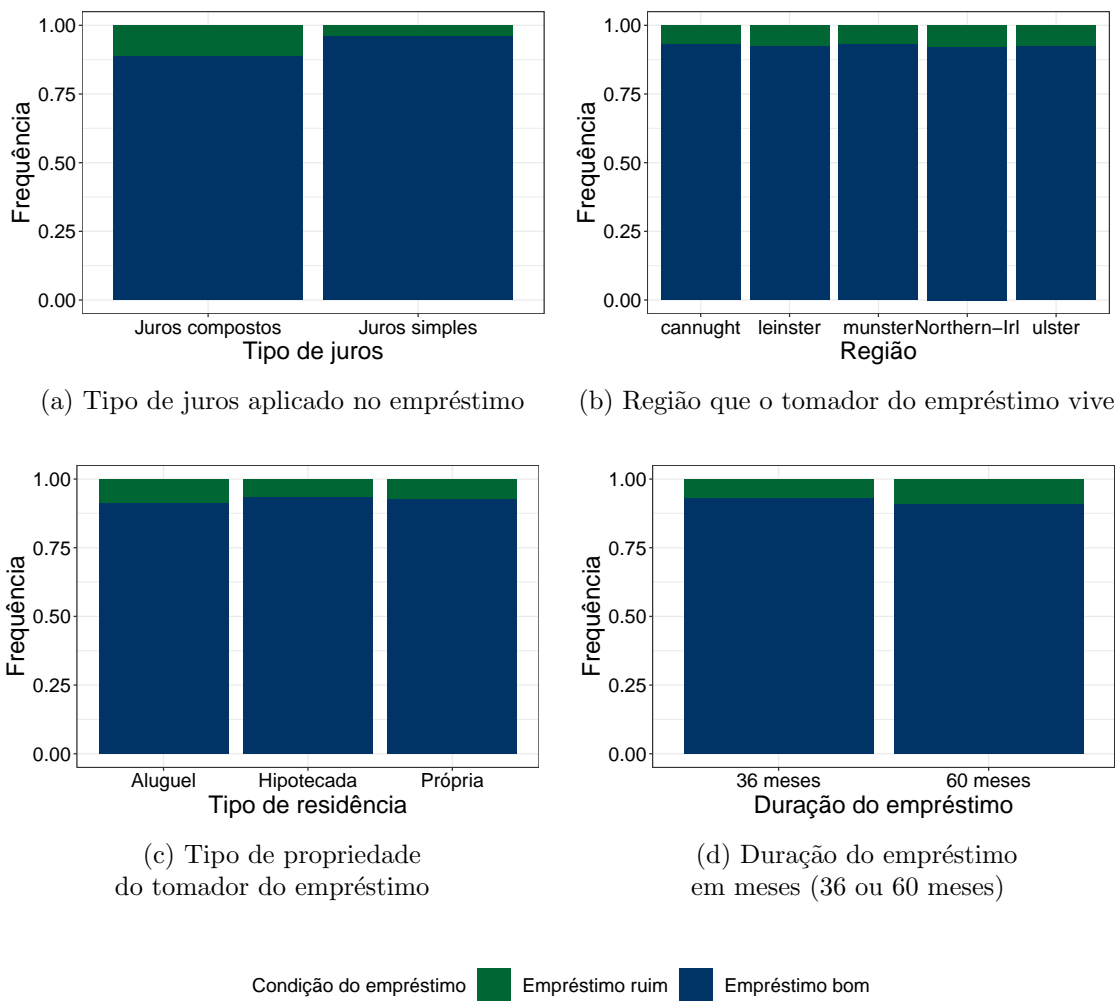


Figura 18: Variáveis explicativas em relação à condição do empréstimo

Na Figura 18, o gráfico 18a evidencia que empréstimos obtidos sob juros compostos possuem uma proporção mais elevada de rotulações ruins em comparação com empréstimos sob juros simples. As Figuras 18b e 18c destacam uma proporção natural refletida pela distribuição das categorias da variável resposta, conforme apresentado na Tabela 4. Já a Figura 18d revela uma proporção mais significativa de empréstimos ruins quando estes tendem a demorar mais para serem pagos.

Covariáveis	Coefficiente de correlação
Tempo de trabalho	-0.02
Renda anual	-0.03
Valor do empréstimo	0.00
Taxa de juros	0.18
DTI	0.01
Valor bruto pago	-0.04
Valor líquido pago	-0.10
Parcela	-0.01
Duração do empréstimo	0.01

Tabela 5: Valores do coeficiente de Pearson entre as covariáveis e a variável resposta

A partir da análise da Tabela 5, nota-se que as correlações entre as variáveis explicativas e a variável resposta são de baixa magnitude. Os coeficientes calculados indicam uma relação linear fraca ou inexistente entre essas variáveis. Esses resultados sugerem que outros fatores ou relações não lineares podem estar desempenhando um papel mais significativo na explicação da variabilidade na variável resposta.

Covariáveis	Coefficiente de contingência
Tipo de residência	0.04
Tipo de aplicação	0.01
Finalidade	0.03
Tipo de juros	0.14
Risco	0.15
Região	0.01
Prazo	0.04
Renda	0.04

Tabela 6: Valores do coeficiente de contingência entre as covariáveis e a variável resposta

Ao analisar a Tabela 6, nota-se que a maioria dos coeficientes de contingência entre as covariáveis e a variável resposta são próximos de zero. Destaca-se que a variável "Risco" exibe o maior valor de associação, atingindo 0.15. Entretanto, é importante ressaltar que esse valor ainda é relativamente baixo. Os coeficientes sugerem, em geral, uma falta de associação significativa entre as covariáveis mencionadas e a variável resposta.

4.2 Regressão logística

Como visto na Equação 2.1.1, a fórmula do modelo logístico é dada por:

$$P(Y = 1|x_1, x_2, \dots, x_k) = \frac{1}{1 + e^{-(\beta_0 + x_1\beta_1 + x_2\beta_2 + \dots + x_k\beta_k)}} \quad (4.2.1)$$

E com base nela os seguintes resultados foram encontrados:

Covariáveis	Coefficientes	Erro padrão
Valor líquido pago	-4.733	0.034
Valor bruto pago	3.321	0.028
Tipo de aplicação	-1.848	0.106
Taxa de juros	1.412	1.412
Valor do empréstimo	-1.406	0.039
Risco	-1.049	0.014
Tipo de juros	-0.462	-0.462
Prazo	-0.203	0.028
Renda anual	-0.195	0.010
DTI	-0.151	0.011
Renda categorizada	-0.111	0.015
Tempo de trabalho	-0.061	0.009
Região	0.038	0.003
Duração do empréstimo	0.033	0.009
Tipo de residência	0.019	0.003
Finalidade	0.019	0.002
Parcela	0.005	0.000

Tabela 7: Estimativa dos coeficientes do modelo logístico e o erro padrão associado

Ao analisar os resultados apresentados na Tabela 7, fica evidente que as variáveis "Valor líquido pago" e "Valor bruto pago" exercem uma influência significativa no valor final de $P(Y = 1)$. Essas duas variáveis estão diretamente associadas à quantia do empréstimo que o cliente já quitou, indicando sua relevância na predição do resultado. Ao calcular a Razão de chances dessas duas variáveis, temos que:

- "Valor líquido pago": apresenta um RC de 0.008, o que sugere que, mantendo todas as outras variáveis constantes, a chance de o empréstimo ser classificado como bom é 125 vezes maior do que ser classificado como ruim.

- "Valor bruto pago": exibe um RC de 27.68, indicando que, ao manter todas as outras variáveis constantes, a chance de o empréstimo ser classificado como ruim é 27 vezes maior do que ser classificado como bom.

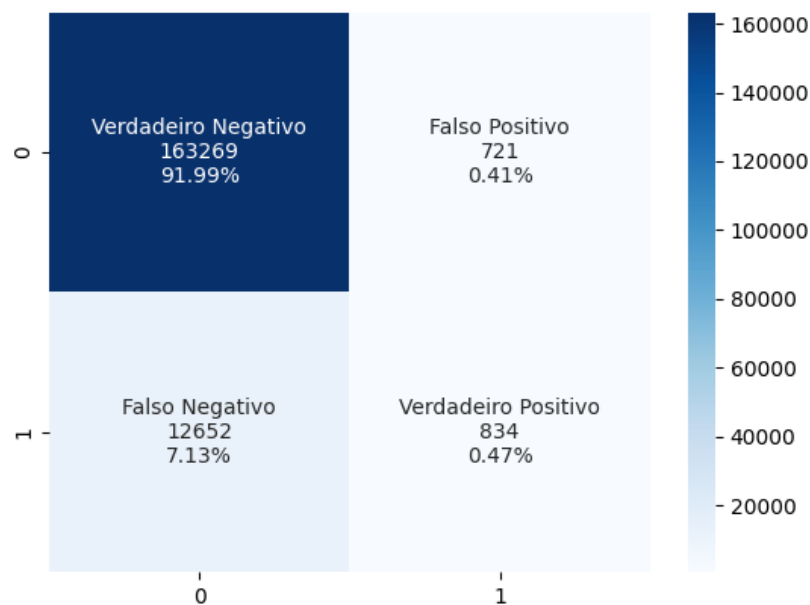


Figura 19: Matrix de confusão do modelo logístico

Visualizando os dados da Figura 19 é possível analisar os resultados do modelo logístico no conjunto de teste. O conjunto de teste apresenta uma distribuição da variável resposta de com mais de 92% dos casos como um empréstimo bom, e o restante como o empréstimo ruim.

	Precisão	Recall	F1-Score	Tamanho da amostra
0	0.928	0.995	0.961	163990
1	0.536	0.062	0.111	13486
Média macro	0.732	0.528	0.535	177476
Média ponderada	0.898	0.925	0.896	177476
Acurácia	0.925			

Tabela 8: Report do modelo logístico

Com base nos dados apresentados na Tabela 8 e na Figura 19, observamos que o modelo exibe uma acurácia elevada. Ele é capaz de fazer previsões precisas na maioria dos casos, alcançando uma taxa de 92,46% de classificações corretas no conjunto de teste. No entanto, é crucial destacar que essa elevada acurácia é influenciada pela proporção

significativa de casos onde o empréstimo é rotulado como "bom", presente em mais de 92% dos dados de teste. Como resultado, o modelo tende a classificar uma parte considerável dos dados como "0", refletindo a influência dessa distribuição desigual na estimação dos parâmetros do modelo logístico.

Ao examinarmos a precisão do modelo, observamos uma taxa de acerto de 73% nas previsões em comparação com as rótulos reais do conjunto de teste. É importante ressaltar a notável precisão na categoria "Empréstimo bom", atingindo quase 92%. No entanto, vale destacar que esse valor elevado está correlacionado ao desequilíbrio nos dados, onde a classe "Empréstimo bom" é predominante.

Ao avaliar o recall do modelo logístico, observamos, em média, valores mais baixos em comparação com a precisão. O recall médio é de 52,8%, indicando que, ao analisar as porcentagens das rótulos reais, o modelo conseguiu acertar um pouco mais da metade delas. Esse desempenho é atribuído ao alto número de falsos negativos no modelo, visto que, ao considerar o total de "Empréstimos ruins" (13.486), o modelo acertou apenas 834 desses casos.

O F1-score acaba refletindo a real situação do modelo, pois ele balanceia os bons resultados apresentados pela precisão com os resultados ruins do recall. O F1-score médio apresentado foi de 52,57%.

A avaliação global do modelo logístico revela um viés significativo, amplificado pelo desequilíbrio nos dados. Embora o modelo tenha alcançado uma taxa geral de acerto de 92%, sua incapacidade de distinguir adequadamente entre "Empréstimos bons" e "Empréstimos ruins" é evidente. Este desempenho inferior sugere limitações na capacidade do modelo de generalizar e discriminar efetivamente entre as categorias, indicando a necessidade de refinamentos ou considerações adicionais para melhorar sua robustez.

4.3 Rede neural

Esse seção vai abordar os resultados obtidos pelo modelo de rede neural.

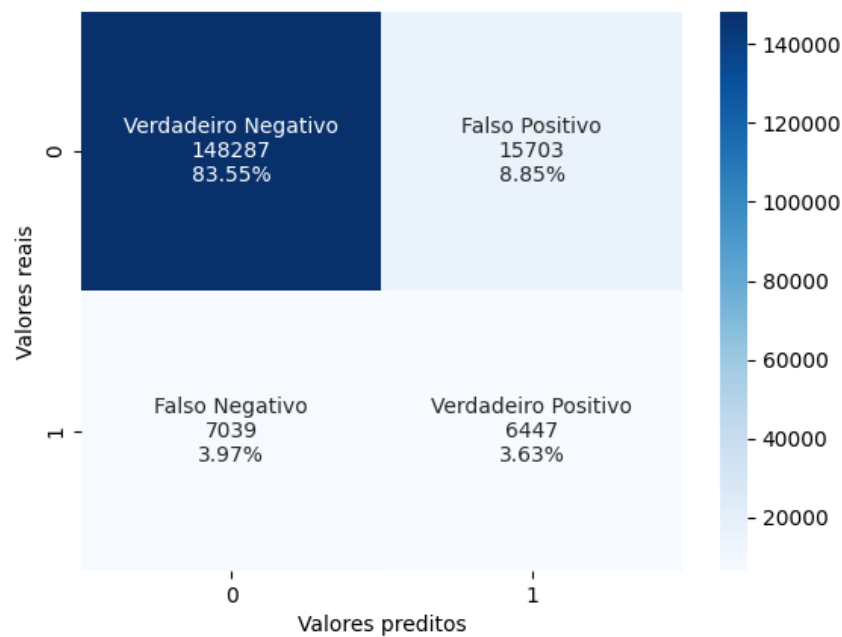


Figura 20: Matrix de confusão da rede neural

A Figura 20 estima que aproximadamente 12% da amostra de teste estimada pelo modelo consiste em empréstimos ruins, uma estimativa ligeiramente distante da distribuição real apresentada na Tabela 4. No entanto, esse resultado sugere que o modelo conseguiu generalizar seus resultados, quando comparado com o modelo logístico. Por outro lado, o número de empréstimos classificados como "0" diminuiu, representando mais de 84% da amostra de teste. Esse ajuste parece ter aumentado o número de Falsos Positivos, provavelmente devido à tentativa do modelo em variabilizar os resultados.

	Precisão	Recall	F1-Score	Tamanho da amostra
0	0.954	0.904	0.929	163990
1	0.291	0.478	0.362	13486
Média macro	0.622	0.691	0.645	177476
Média ponderada	0.904	0.872	0.886	177476
Acurácia	0.872			

Tabela 9: Report da rede neural

Os resultados apresentados na Tabela 9 indicam um desempenho satisfatório do modelo, especialmente em termos de acurácia e métricas avaliadas para a classe "0". Contudo, ao comparar esses resultados com as métricas da classe "1", percebe-se que o modelo ainda é influenciado pelo elevado número de observações na classe "0". Entretanto, uma análise mais detalhada do Recall da classe "1" revela que o modelo conseguiu reduzir

significativamente o número de Falsos Negativos, identificando corretamente quase metade dos empréstimos considerados ruins na base original. Essa melhoria no Recall da classe "1" resultou em uma precisão inferior para essa classe, refletindo que menos de 30% das previsões do modelo foram corretas nesse contexto, como evidenciado na Tabela 9.

Em termos gerais, as decisões relacionadas à arquitetura do modelo, seus hiperparâmetros, estratégia de treinamento e outros fatores contribuíram para que o modelo de redes neurais realizasse previsões de alta qualidade, não se limitando apenas ao desbalanceamento dos dados. Os resultados obtidos pelo modelo de redes neurais foram satisfatórios devido à eficácia da arquitetura da rede. Portanto, buscar aprimorar ainda mais essa arquitetura pode ser uma abordagem promissora na busca por resultados ainda melhores.

4.4 Interpretação da rede neural

Após definir o modelo de rede neural e examinar seus resultados, esta seção aborda a interpretação do modelo. Para isso, foram construídos gráficos com base nos resultados do SHAP, destacando as variáveis de maior importância no resultado final. Utilizando uma amostra de 80 observações, o valor de SHAP foi calculado para cada observação, permitindo uma análise tanto individual quanto conjunta dessa amostra.

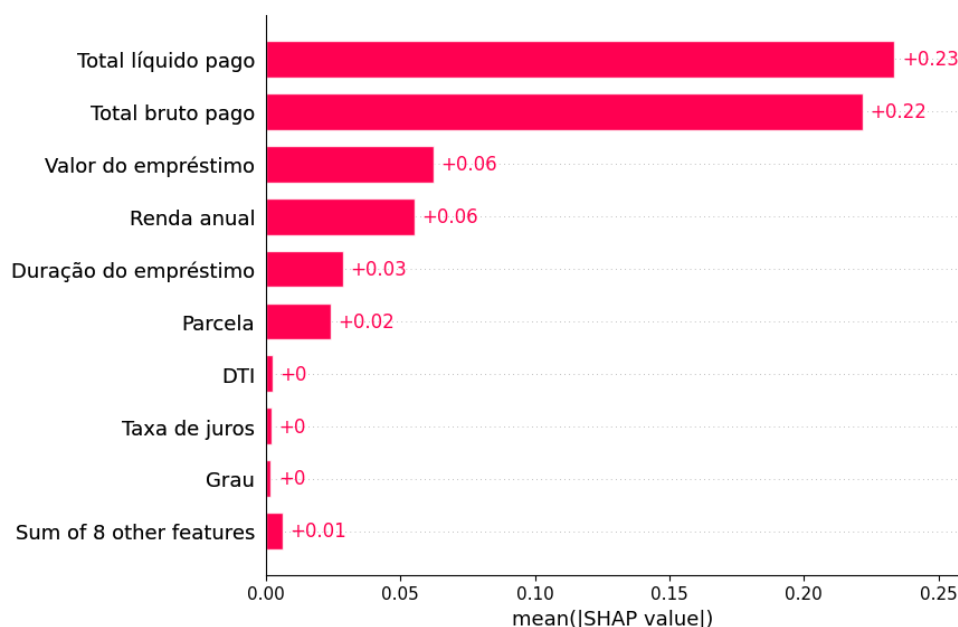


Figura 21: Média absoluta dos valores de shap

O Gráfico 21 exibe a média absoluta dos valores de SHAP para cada variável nas 80 observações consideradas. Essa representação oferece insights sobre quais variáveis o modelo considerou mais relevantes durante as predições, destacando a magnitude da

contribuição de cada variável. Vale notar que, por se tratar de valores absolutos, o gráfico não proporciona informações sobre se a contribuição de cada variável é positiva ou negativa. Entretanto, os próximos gráficos irão elucidar essa questão ao detalhar a contribuição específica de cada variável.

Ao analisar cada variável no gráfico, destaca-se que "Total líquido pago" e "Total bruto pago" são as que mais contribuem para o resultado final do modelo. O gráfico enfatiza a relevância de nove variáveis, omitindo o restante devido à sua baixa contribuição. Ao observar as variáveis omitidas, percebe-se que, somadas, suas contribuições aproximam-se de 0.01, evidenciando a sua baixa influência no resultado final do modelo de rede neural.

Os próximos gráficos proporcionam a visualização dos valores SHAP para cada variável em observações individuais. O gráfico também revela os valores específicos utilizados em cada variável. A apresentação de quatro figuras distintas tem como objetivo compreender como o valor SHAP é atribuído em quatro cenários distintos que surgem em problemas de classificação binária.

Os 2 primeiros gráficos mostram casos onde o modelo acertou suas predições, tanto para casos de "Empréstimos ruins" quanto para "Empréstimos bons".

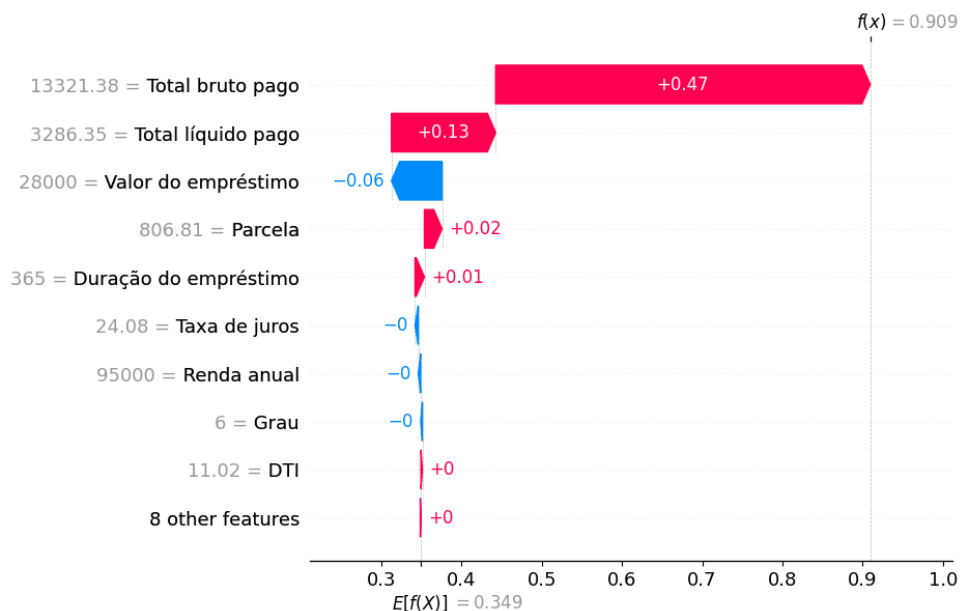


Figura 22: Valor de SHAP para uma observação do tipo Verdadeiro Positivo

Na Figura 22, cenário em que o modelo acertou a classificação de um empréstimo como ruim, observa-se que as variáveis que mais contribuíram foram as mesmas observadas na Figura 21, comportamento esse que vai prevalecer nos demais casos. Ao analisar os valores dessas duas variáveis, nota-se que, para esse cliente específico, ainda resta um montante significativo do empréstimo a ser pago, totalizando quase 90% do valor líquido

pendente e quase 50% do valor bruto pendente. Uma diferença muito grande entre o valor do empréstimo e o que falta a ser pago pode ser um dos fatores que está ocasionando a rotulação de observações como "Empréstimos ruins".

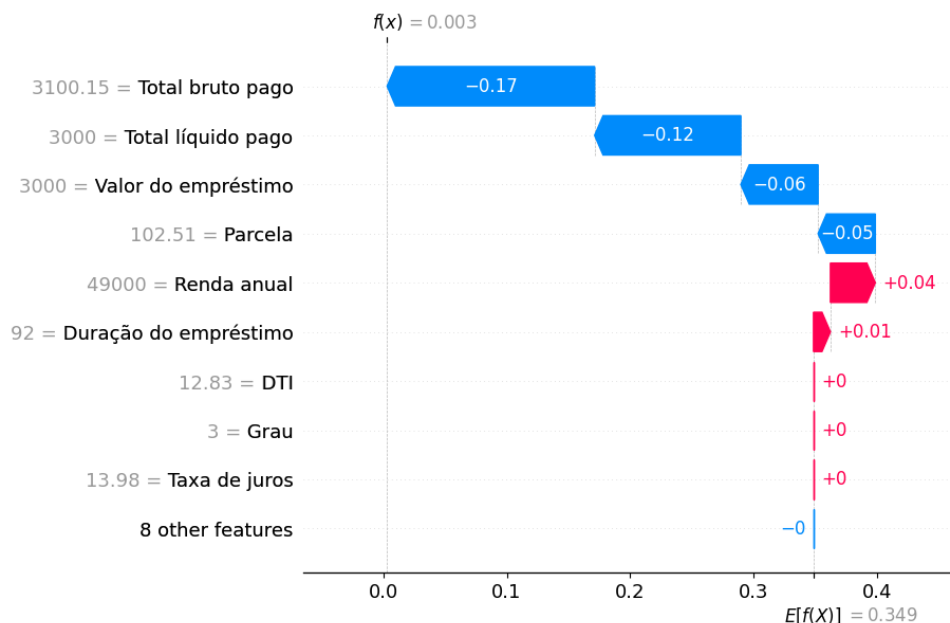


Figura 23: Valor de SHAP para uma observação do tipo Verdadeiro Negativo

Ao analisar a Figura 23, que representa o cenário em que o modelo acertou a rotulação de um empréstimo bom, nota-se um valor final bem próximo de 0. Indícios de que o modelo teve mais confiança ao realizar essa predição. Ao analisar os valores de cada variável, é possível perceber que o cliente já está finalizando ou finalizou o empréstimo, dado que as variáveis de pagamento do empréstimo chegaram no valor real do empréstimo.

Analisando as Figuras 22 e 23, é possível observar a flexibilidade do modelo em lidar com diferentes valores da mesma variável. Quando o modelo encontra valores que indicam um empréstimo ruim, atribui valores positivos para a contribuição das variáveis, aproximando-as de 1. Da mesma forma, para empréstimos bons, o modelo adiciona uma contribuição negativa, direcionando o resultado para 0. Isso demonstra como o modelo responde de forma dinâmica às variações nos valores das variáveis.

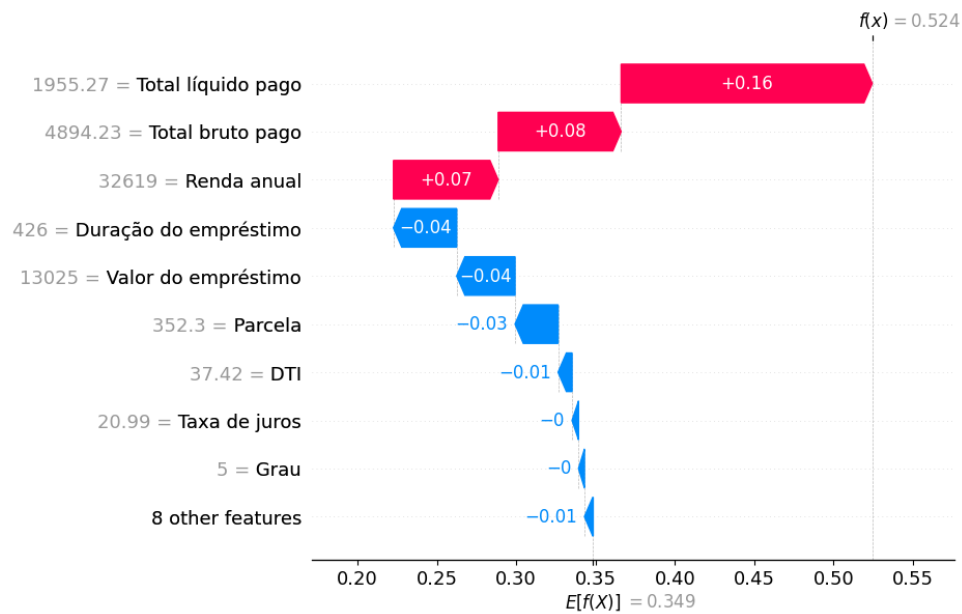


Figura 24: Valor de SHAP para uma observação do tipo Falso Positivo

Ao analisar um exemplo em que o modelo erroneamente classifica um "Empréstimo ruim", conforme ilustrado na Figura 24, destacam-se as características já discutidas nas interpretações anteriores. O resultado do modelo foi muito próximo do limiar de decisão, que é 0,5, indicando que o modelo teve alguma indecisão ao realizar a predição desse caso.

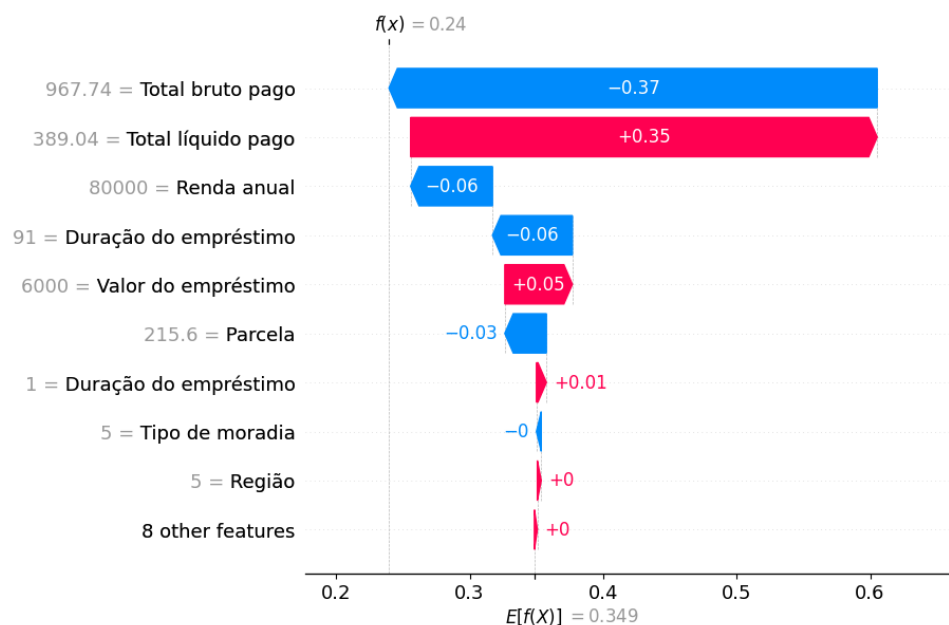


Figura 25: Valor de SHAP para uma observação do tipo Falso Negativo

Observando a Figura 25, que retrata um cenário em que o modelo classifica de forma equivocada um "Empréstimo bom", é possível perceber um comportamento de divergência entre as variáveis que mais contribuem com o resultado do modelo. Esse

padrão de divergência não foi observado nos casos anteriores, e pode ter influenciado o modelo a realizar uma predição incorreta.

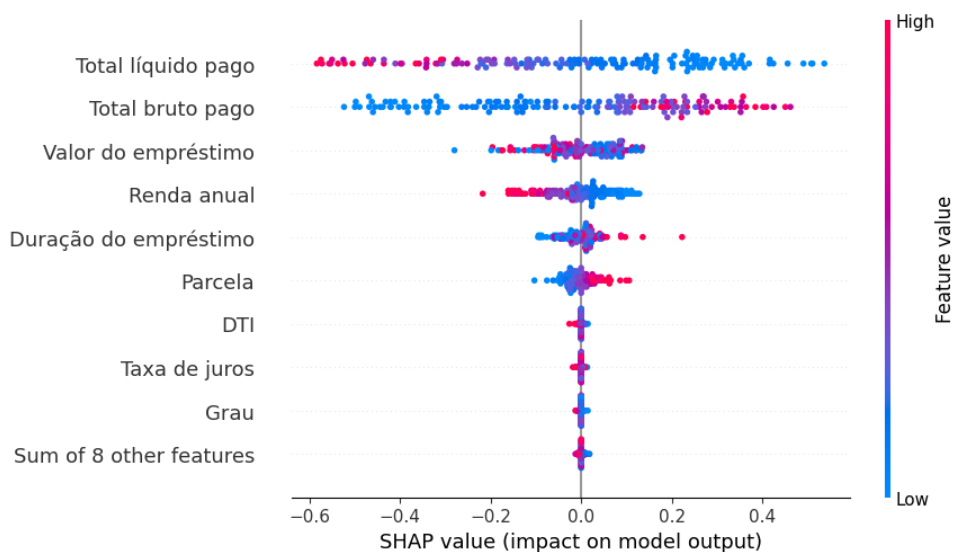


Figura 26: Valores de shap para as 80 observações utilizadas

A Figura 26 apresenta o comportamento do valor de SHAP para cada variável entre as 80 observações utilizadas no experimento. No gráfico, o eixo x refere-se à distribuição do valor de SHAP, enquanto o eixo y representa cada variável. O eixo das cores ilustra a distribuição dos valores da variável, sendo que cores mais quentes indicam valores maiores e cores mais frias indicam valores menores.

Ao analisar as variáveis que mais influenciam o resultado do modelo, observa-se uma distribuição mais dispersa no eixo x. Ao focar nas duas variáveis de maior contribuição, nota-se um comportamento heterogêneo dos valores em comparação com os valores de SHAP. Para a primeira variável, elevados valores tendem a impactar negativamente no resultado do modelo, enquanto a segunda variável apresenta o comportamento oposto, com valores mais baixos contribuindo negativamente no resultado do modelo.

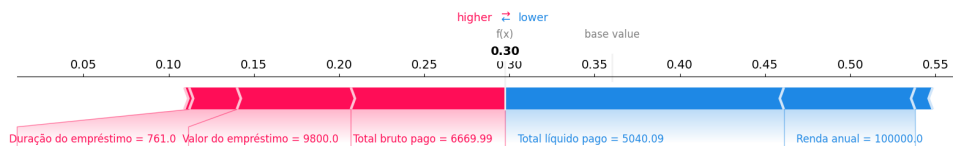


Figura 27: Gráfico de força em uma observação

A Figura 27 ilustra a "força" de contribuição de cada variável no modelo. Este gráfico proporciona uma visão clara das variáveis que tiveram impacto positivo e negativo no resultado final. Cada barra representa a intensidade da contribuição de uma variável

específica, sendo que aquelas com maior influência concentram-se no centro, enquanto as de menor influência ficam nas extremidades. A figura é centrada no valor final predito pelo modelo, que, neste caso, é observado como 0.3, indicando um empréstimo classificado como bom.

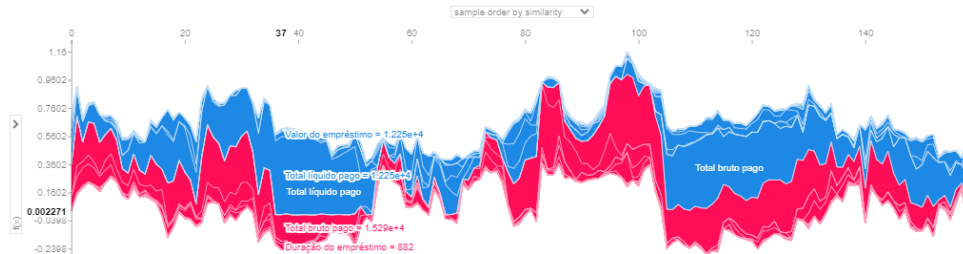


Figura 28: Gráfico de força para múltiplas observações

A Figura 28 é uma generalização da Figura 27, considerando as 80 observações. Este gráfico é um recorte obtido de uma visualização dinâmica gerada pelo pacote SHAP. Devido à natureza dinâmica do gráfico original, ao transformá-lo em uma imagem estática, parte de sua capacidade de representação é limitada. No entanto, é possível observar alguns padrões nesse recorte. Por exemplo, cores mais quentes indicam contribuições positivas quando o modelo faz previsões próximas de 1, com tons de vermelho representando adições das variáveis ao resultado. Da mesma forma, cores mais frias indicam contribuições negativas quando o modelo faz previsões próximas de 0.

4.5 Benchmark entre regressão logística e redes neurais

Ao obter interpretações do modelo de redes neurais, abre-se a possibilidade de realizar comparações entre os modelos de regressão logística e redes neurais no âmbito interpretativo. Até recentemente, essa capacidade de interpretação estava exclusivamente associada ao modelo logístico em comparação com o modelo de redes neurais. Essa análise comparativa se soma às comparações já realizadas entre os modelos, englobando aspectos como arquitetura, tempo de treinamento, tempo de predição e os resultados gerados por ambos os modelos. Os resultados abaixo evidenciam essas comparações.

4.5.1 Complexidade da arquitetura

Modelo	Número de parâmetros
Regressão Logística	18
Rede Neural	151233

Tabela 10: Número de parâmetros nos modelos

Ao examinar a Tabela 10, destaca-se a significativa disparidade na complexidade entre os modelos. Enquanto o modelo logístico possui apenas um parâmetro para cada covariável, a rede neural apresenta mais de 8400 parâmetros para cada parâmetro do modelo logístico, sendo esses distribuídos nos neurônios e camadas da rede.

4.5.2 Resultado dos modelos

Métricas	Regressão logística	Rede neural
Falsos Positivos	721	15703
Falsos Negativos	12652	7039

Tabela 11: Comparação dos resultados de Falsos Positivos e Falsos Negativos

Ao analisar os casos em que os modelos cometeram erros, conforme apresentado na Tabela 11, é evidente que, em geral, o modelo logístico cometeu menos erros do que a rede neural. O modelo logístico registrou um total de pouco mais de 13 mil classificações incorretas, enquanto a rede neural teve mais de 18 mil erros. No entanto, é importante destacar que a rede neural cometeu menos erros nos casos de Falsos Negativos. Esse cenário é significativo no contexto de empréstimos, pois rotular erroneamente um possível cliente inadimplente como alguém que cumprirá com o empréstimo pode ter implicações mais graves.

Os resultados a seguir vão identificar melhor como foi o processo preditivo dos dois modelos no conjunto de teste.

Métricas	Regressão logística	Rede neural
Precisão (Classe 0)	0.928081	0.954
Precisão (Classe 1)	0.536334	0.291
Recall (Classe 0)	0.995603	0.904
Recall (Classe 1)	0.0618419	0.478
F1-Score (Classe 0)	0.960657	0.929
F1-Score (Classe 1)	0.110897	0.362
Acurácia	0.924649	0.872

Tabela 12: Comparação dos resultados da Regressão logística e Rede neural. Em verde, o modelo que obteve o melhor resultado na respectiva métrica.

Ao examinar a Tabela 12, fica claro que a regressão logística superou a rede neural em quatro categorias. No entanto, é crucial compreender que esse desempenho destacado da regressão logística decorre da sua falta de generalização. Este modelo, ao ser treinado

com dados desbalanceados, enfrenta limitações devido à sua arquitetura menos complexa. Isso ressalta a necessidade de equilibrar a simplicidade do modelo com a capacidade de generalização para obter resultados mais confiáveis.

Ao avaliar as métricas críticas no contexto de empréstimos, como o Recall da classe 1 e, conseqüentemente, o F1-score dessa classe, destaca-se que a rede neural supera a regressão logística. Notavelmente, o recall na rede neural é 40% superior, evidenciando uma performance mais robusta no que diz respeito à identificação de empréstimos ruins.

Covariáveis	Coefficientes da regressão logística	Médias absoluta dos valores de SHAP
Total líquido pago	-4.733	0.23363
Total bruto pago	3.321	0.22619
Valor do empréstimo	-1.406	0.06298
Renda anual	-0.195	0.06132
Tempo de trabalho	-0.061	0.03105
Parcela	0.005	0.02586
DTI	-0.151	0.00242
Taxa de juros	1.412	0.00232
Risco	-1.049	0.00169
Duração do empréstimo	0.033	0.00166
Tipo de moradia	0.019	0.00117
Finalidade	0.019	0.00094
Tipo de juros	-0.462	0.00077
Região	0.038	0.00071
Prazo	-0.203	0.00068
Renda categorizada	-0.111	0.00058
Tipo da aplicação	1.848	0

Tabela 13: Relação entre os coeficientes da regressão logística com os valores absolutos de shap

A Tabela 13 realiza uma comparação entre os coeficientes do modelo logístico e a média dos valores absolutos de SHAP, sendo esse último calculado com base em uma amostra composta por 80 observações. A média dos valores de SHAP proporciona uma medida da magnitude das contribuições das variáveis, apresentando, de certa forma, uma analogia com os coeficientes da regressão logística.

Ambos os modelos exibiram semelhança nas duas variáveis que mais impactam o resultado final. No entanto, ao classificar as demais variáveis em ordem de importância, observa-se uma divergência significativa no grau de influência que esses valores exercem.

4.5.3 Tempo de execução

Examinar o tempo de predição é crucial para alcançar resultados mais rápidos, sendo em alguns contextos um fator de extrema relevância, quase tão significativo quanto as métricas de desempenho do modelo.

Estatísticas	Regressão logística	Rede neural
Mínimo	0.0	52.067995
Quartil 25	0.0	53.327155
Média	0.3335619	59.446688
Mediana	0.0	56.855202
Quartil 75	0.88143349	66.278052
Máximo	1.50370598	92.03124
Variância	0.28385463	59.008917
Desvio padrão	0.5327801	7.681726

Tabela 14: Tempo de predição (em ms) de cada modelo, em uma amostra com 50 observações.

A Tabela 14 destaca a extrema eficiência do modelo logístico em comparação com o modelo de redes neurais. O modelo logístico demonstrou um tempo de predição mais curto, com a predição mais demorada levando apenas 1,5 milissegundos, sendo mais de 60 vezes mais rápido do que a predição mais demorada do modelo de redes neurais. Essa tendência é consistente em relação às outras métricas de desempenho.

Observações	Tempo de execução
1	317s
80	7.07hrs

Tabela 15: Tempo de execução para realizar a interpretação das variáveis

Ao contrário do modelo logístico, onde não há um tempo de execução associado à interpretação, pois a interpretação é derivada dos coeficientes estimados, a situação é diferente no modelo de redes neurais. A interpretação de uma única observação demandou mais de 5 minutos, e o cálculo das interpretações para as 80 observações utilizadas nos resultados anteriores exigiu mais de 7 horas. Essa diferença substancial de tempo destaca a complexidade computacional envolvida na interpretação de modelos mais intrincados, como redes neurais.

5 Conclusão

A presente pesquisa buscou ampliar o entendimento sobre a interpretabilidade de modelos de redes neurais, explorando métodos que permitam elucidar as predições desses modelos complexos. A comparação sistemática entre um modelo de redes neurais e a tradicional regressão logística foi central para a análise, devido à natureza interpretativa já conhecida do modelo logístico.

A utilização da técnica SHAP (SHapley Additive exPlanations) para a interpretação do modelo de redes neurais revelou-se uma ferramenta poderosa, permitindo uma compreensão mais profunda das contribuições de cada variável nas predições do modelo. Este método proporcionou uma visão holística, destacando as características individuais que mais influenciaram nas decisões do modelo.

Ao confrontar os resultados preditivos da rede neural com a regressão logística, observou-se que o modelo de redes neurais apresentou predições mais heterogêneas, enquanto a regressão logística, por sua simplicidade, ficou refém do desbalanceamento dos dados.

Essa pesquisa contribui para a discussão em torno da interpretabilidade em inteligência artificial e fornecendo resultados valiosos para a aplicação prática desses modelos em contextos onde a transparência é essencial. A busca contínua por métodos interpretativos robustos é vital para a implementação responsável e eficaz de modelos de aprendizado de máquina em diversos domínios.

Referências

- BISHOP, C. M. *Pattern Recognition and Machine Learning*. [S.l.]: Springer, 2006.
- GERON, A. *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*. [S.l.]: O'Reilly Media, 2019. ISBN 978-1492032649.
- HOSMER, D. W. *Applied Logistic Regression*. [S.l.]: John Wiley, 2013.
- IZBICKI, R.; SANTOS, T. M. dos. *Aprendizado de máquina: uma abordagem estatística*. [S.l.: s.n.], 2020. ISBN 978-65-00-02410-4.
- JAMES, G. et al. *An Introduction to Statistical Learning with Applications in R*. New York: Springer, 2013.
- LUNDBERG, S.; LEE, S.-I. A unified approach to interpreting model predictions. In: *Advances in neural information processing systems*. [S.l.: s.n.], 2017. v. 30, p. 4765–4774.
- SHAPLEY, L. S. A value for n-person games. *Contributions to the Theory of Games*, v. 2, p. 307–317, 1953.
- WERBOS, P. J. *Beyond regression: New tools for prediction and analysis in the behavioral sciences*. Tese (Doutorado) — Harvard University, 1974.
- WITTE, R. S.; WITTE, J. S. *Statistics*. 11th. ed. New York: Wiley, 2017. ISBN 978-0471557613.

6 Anexo