

Análise de Churn de Clientes Bancários (Case BankChurners)

Introdução

O **churn de clientes bancários** refere-se ao fenômeno de clientes encerrarem seu relacionamento com o banco – por exemplo, cancelando contas ou cartões de crédito. Esse problema é crítico no setor financeiro, pois a perda de clientes impacta diretamente a receita e pode gerar custos elevados para substituir clientes perdidos. Estudos indicam que adquirir um novo cliente pode custar **5 a 25 vezes mais** do que manter um existente, e um aumento de apenas 5% na retenção pode elevar os lucros em **25% a 95%** ¹. Nos EUA, estima-se que a evasão de clientes bancários cause prejuízos de centenas de bilhões de dólares por ano ². Portanto, identificar **quais clientes têm propensão a churn e quais fatores contribuem para a saída** é fundamental para orientar estratégias de retenção e minimizar perdas.

Objetivo do Estudo: Neste relatório, conduzimos uma análise completa do dataset **BankChurners.csv** (10.127 clientes de cartão de crédito, dos quais ~16% encerraram o vínculo ³). Trata-se de um estudo de caso estruturado que abrange: (1) exploração descritiva de todas as variáveis, (2) criação de variáveis derivadas para capturar comportamentos relevantes, (3) aplicação de técnicas de *machine learning* supervisionadas (modelos de classificação) para prever churn, e (4) aplicação de análise não supervisionada (*clustering*) para segmentação de clientes. Apresentamos a **metodologia** adotada, discutimos os resultados (incluindo a comparação de modelos **XGBoost**, **Random Forest** e **LightGBM** via métrica AUC) e, por fim, sugerimos melhorias e próximos passos (novas features, tuning de hiperparâmetros, uso de interpretabilidade, etc.).

Metodologia Aplicada

Conduzimos o projeto seguindo as principais etapas de Data Science:

- **1. Análise Descritiva (Exploratória):** Inicialmente examinamos os dados brutos, identificando os tipos de variáveis (numéricas vs categóricas), distribuição de valores, presença de outliers e possíveis inconsistências. Calculamos estatísticas resumo e visualizamos distribuições e relações entre variáveis para obter insights iniciais. Não foram encontrados valores nulos no dataset (os dados estavam completos) ⁴ ⁵, embora algumas categorias “Unknown” existam para indicar ausência de informação (p.ex. renda desconhecida). Essa fase permitiu entender o perfil dos clientes e diferenças entre clientes ativos e churners.
- **2. Engenharia de Features:** Com base nos insights, criamos **variáveis derivadas** para realçar padrões comportamentais relevantes. Por exemplo, calculamos o **Ticket_Médio** (valor médio por transação) para distinguir clientes que fazem muitas compras pequenas versus poucas compras grandes. Também geramos indicadores de tendência, como flags de queda no volume de transações trimestre a trimestre. Essas novas features e transformações são detalhadas adiante e visam melhorar tanto a compreensão quanto a performance preditiva.

- **3. Modelagem Preditiva (Classificação):** Em seguida, formulamos o problema como classificação binária (churn = sim ou não). Testamos algoritmos supervisionados, incluindo **Regressão Logística** (como baseline interpretável) e modelos avançados de árvore de decisão em ensemble: **Random Forest**, **XGBoost** e **LightGBM**. Os dados foram divididos em treino/validação, e treinamos os modelos otimizando para a métrica **AUC (Área Sob a Curva ROC)** – adequada aqui dada a leve desproporção de classes (16% churn) ³ e por fornecer uma medida robusta do poder de classificação independente de threshold.
- **4. Validação e Avaliação:** Comparamos os modelos quanto à AUC obtida em dados de teste fora da amostra, além de analisar acurácia, matriz de confusão e outras métricas para ter uma visão completa. Damos ênfase à AUC para comparar o desempenho preditivo de forma imparcial. Também empregamos validação cruzada estratificada para garantir estabilidade dos resultados.
- **5. Clusterização de Clientes:** Em paralelo à predição, aplicamos técnicas de **análise não supervisionada** (como o algoritmo *K-Means*) para segmentar os clientes em grupos com perfis semelhantes, usando variáveis de comportamento e valor. O objetivo da clusterização foi identificar **perfis de clientes** (ex: “alto valor e engajamento”, “baixo uso e alto risco de churn”, etc.) que pudessem auxiliar na compreensão de padrões e em ações de marketing direcionadas. Avaliamos diferentes números de clusters usando métricas como a silhueta para encontrar uma segmentação adequada.
- **6. Interpretação e Comunicação:** Por fim, interpretamos os resultados combinando as descobertas da análise exploratória, os fatores de importância dos modelos (como os *feature importances* do Random Forest) e os perfis de clusters formados. Construímos visualizações (gráficos de distribuição, correlação e performance de modelo) e tabelas para comunicar os insights de forma clara. Também apontamos **recomendações** de negócio e sugestões de melhorias no modelo (por exemplo, uso de técnicas de *explainable AI* para explicar predições individuais, coleta de dados adicionais, tuning de hiperparâmetros, etc.).

A seguir, apresentamos em detalhes os resultados da análise de cada variável do dataset, a utilidade das variáveis derivadas criadas, os principais padrões de correlação encontrados e o desempenho comparativo dos modelos preditivos desenvolvidos.

Análise Exploratória das Variáveis

Nesta seção examinamos cada variável disponível no dataset original, incluindo tipo, distribuição, presença de outliers, transformações aplicadas (se houver) e relação com o churn. Para facilitar, agrupamos variáveis por categorias (demográficas, de comportamento financeiro, etc.) e destacamos insights relevantes.

Variáveis Demográficas:

- **Customer_Age (Idade):** Variável numérica. Idade dos clientes variou de **26 a 73 anos**, com média em torno de 46 anos ⁶. A distribuição é aproximadamente simétrica, com leve concentração em faixas de meia-idade (40–55 anos predominantes). Identificamos poucos outliers (alguns clientes acima de ~70 anos). Em termos de churn, notou-se uma tendência curiosa: faixas **muito jovens** e **muito idosas** apresentaram **menor propensão a churn**, enquanto grupos de meia-idade tiveram churn mais elevado. Por exemplo, clientes entre **55–60 anos** e **40–45 anos** tiveram taxas de churn acima de 17%, ao passo que jovens 25–30 anos apresentaram churn em torno de apenas **8–9%** ⁷. Isso sugere que clientes de meia-idade média podem estar mais suscetíveis a trocar de banco, talvez por buscarem

melhores condições ou por mudanças nas necessidades financeiras, enquanto os mais novos ainda estão construindo relacionamento e os mais velhos tendem a permanecer fiéis ou têm menos opções.

- **Gender (Gênero):** Variável categórica (Masculino/Feminino). No conjunto, há leve maioria de **clientes do sexo feminino (52,9%)** frente 47,1% masculinos ⁸. Embora gênero não seja o fator mais determinante, observamos que as **mulheres apresentam churn ligeiramente maior que os homens**. O modelo de regressão logística indicou que, mantendo os demais fatores constantes, ser do sexo **masculino reduz a probabilidade de churn** – o coeficiente para Male foi negativo e estatisticamente significativo ($p < 0,001$) ⁹. Em outras palavras, clientes do sexo feminino têm uma chance um pouco maior de cancelar o cartão em comparação aos homens. Uma possível explicação (a ser investigada qualitativamente) seria diferenças de comportamento de consumo ou atendimento entre gêneros; de qualquer forma, a diferença é sutil, indicando que gênero, isoladamente, não é um driver forte de churn.
- **Dependent_count (Número de Dependentes):** Numérica (inteira). Indica quantos dependentes (p.ex. membros da família) o cliente possui ligados à conta. Variou de **0 a 5**, com média $\approx 2,35$ dependentes ⁶ – ou seja, a maioria tem 2 ou 3 dependentes. Não há outliers severos, já que o máximo (5) é plausível. Analisando a relação com churn, houve um leve **aumento na taxa de churn conforme cresce o número de dependentes** do cliente. Por exemplo, clientes **solteiros com três dependentes** apresentaram churn em torno de **18%**, comparado a ~12–15% para solteiros sem dependentes ¹⁰. Entre clientes casados, aqueles com 3–4 dependentes também tiveram churn mais elevado (~17%) em comparação a casados sem dependentes (~12%) ¹¹. Esses dados sugerem que ter mais dependentes – possivelmente sinalizando maiores responsabilidades financeiras – pode elevar a chance de cancelar o cartão, talvez buscando consolidar gastos em menos contas ou reduzir custos. De todo modo, o efeito não é monotônico perfeito nem extremamente forte, indicando que dependentes atuam em conjunto com outros fatores (por exemplo, estado civil e renda).
- **Education_Level (Nível de Educação):** Categórica ordinal. As categorias incluem *Uneducated*, *High School*, *College*, *Graduate*, *Post-Graduate*, *Doctorate* e *Unknown*. A distribuição mostra predominância de clientes **Graduados ($\approx 46\%$)**, seguido por **Ensino Médio completo ($\sim 20\%$)**; já pós-graduados, doutores e não-educados são parcelas menores ¹². Em relação ao churn, identificamos uma correlação interessante: clientes com **maior nível educacional tendem a churnar mais**. Especificamente, a taxa de churn para portadores de **Doutorado foi a mais alta ($\sim 21\%$)**, seguida por **Pós-Graduados ($\sim 17\text{--}18\%$)** ¹³. Por outro lado, clientes com educação intermediária tiveram churn menor (p.ex. nível *College* ~15%, *High School* ~15%). Uma hipótese é que clientes de alta escolaridade possuam **expectativas mais altas** em relação aos serviços bancários ou acesso mais fácil a ofertas concorrentes, tornando-os mais propensos a trocar se não estiverem satisfeitos. Já os clientes com formação média podem estar mais “cativos” ou satisfeitos com serviços padrão. Vale notar que a categoria *Unknown* (educação não informada) teve churn em torno da média (~16%), não fugindo do padrão geral.
- **Marital_Status (Estado Civil):** Categórica nominal (valores: *Married*, *Single*, *Divorced*, *Unknown*). A maioria dos clientes é **casada ($\approx 53,6\%$)**, seguida por **solteiros ($\approx 39\%$)**, e divorciados representam ~7% ⁸. O estado civil, isoladamente, mostrou alguma relação com churn. **Clientes divorciados** apresentaram as maiores taxas de churn – especialmente divorciados **sem dependentes ($\sim 17,3\%$ de churn)** ¹¹. Clientes solteiros também tendem a churnar um pouco mais que casados de forma geral. Entre os **casados**, a taxa de churn foi relativamente menor (~12–14% para casados sem dependentes), embora aumente para casados com vários dependentes, conforme mencionado. Em suma, **solteiros e divorciados** constituem segmentos de maior risco de churn comparados a **casados**, possivelmente porque casados mantêm mais

vínculos financeiros estáveis ou beneficiam-se de pacotes conjugais. Entretanto, a diferença não é extrema – casados também churnam em número absoluto significativo devido à sua maioria populacional, mas proporcionalmente aparentam ser mais leais. (A categoria *Unknown* de estado civil é rara e não apresentou padrão claro).

- **Income_Category (Faixa de Renda Anual):** Categórica ordinal com faixas: *Less than \$40K* (menos de \ \$40 mil), \ \$40K-\ \$60K, \ \$60K-\ \$80K, \ \$80K-\ \$120K, \ \$120K+, além de *Unknown*. A distribuição de renda dos clientes é bastante concentrada nas faixas mais baixas: a maior parcela ganha < \ \$40K (**46% dos clientes**), seguida de \ \$40-60K ($\approx 18\%$), \ \$60-80K ($\approx 14\%$), \ \$80-120K ($\approx 15\%$), e apenas $\sim 7\%$ ganham acima de \ \$120K ¹⁴. Em termos de churn, observamos um perfil em “U”: **clientes de renda muito baixa e muito alta churnam mais**, enquanto os de renda média têm churn menor. Especificamente, a taxa de churn entre clientes de **renda < \ \$40K foi $\sim 17,1\%$** , similar à dos de **renda > \ \$120K ($\sim 17,3\%$)** ¹⁵. Já clientes nas faixas intermediárias apresentaram churn mais baixo – o menor sendo na faixa \ \$60K-\ \$80K (**$\sim 13,4\%$ de churn**) ¹⁵. Isso pode indicar que clientes de baixíssima renda talvez cancelem o cartão por dificuldades em arcar com anuidades ou pouco uso, enquanto clientes muito ricos podem buscar bancos concorrentes mais exclusivos ou com melhores benefícios. Os de renda mediana possivelmente encontram maior utilidade e benefícios adequados, mantendo-se fiéis. De qualquer forma, a variável renda é informativa e foi tratada também de forma ordinal (ver variáveis derivadas) para o modelo, já que há uma ordem natural nas categorias (do menor para o maior nível). Importante destacar que $\sim 11\%$ dos clientes têm renda *Unknown* (não informada); esse grupo teve churn em torno de $\sim 16\%$, próximo da média – possivelmente o banco não tinha esses dados, mas decidimos mantê-los como categoria separada para não perder essas entradas.

- **Card_Category (Categoria do Cartão):** Categórica nominal. Indica o tipo de cartão de crédito do cliente, com quatro níveis no dataset: **Blue, Silver, Gold e Platinum**, em ordem crescente de status. A esmagadora maioria dos clientes possui cartão **Blue (básico)** – aproximadamente 93% – enquanto **Silver $\sim 5\%$, Gold $\sim 2\%$, Platinum $< 1\%$** (distribuição inferida; Blue é claramente dominante). Esse desequilíbrio pode refletir que o produto básico é oferecido a todos, enquanto cartões premium são concedidos a poucos clientes de alto valor. Em relação ao churn, inicialmente pode-se supor que clientes premium seriam mais fiéis devido a benefícios exclusivos; porém, a análise indicou o inverso: **clientes com cartões Gold ou Platinum apresentaram propensão maior de churn comparados aos clientes Blue**. No modelo logístico, por exemplo, ter cartão **Gold aumentou significativamente as chances de churn** (coeficiente positivo, $p < 0,001$) em relação ao Blue ¹⁶. Cartões Platinum também mostraram efeito positivo (embora com menor significância dado o baixo n), e Silver teve um leve aumento não tão significativo. Uma interpretação possível: muitos cartões premium podem ter **anuidades mais caras** ou requisitos de gasto; se o cliente não percebe valor suficiente, ele pode cancelar o cartão (churn). Além disso, clientes premium podem ser mais cortejados por concorrentes e estarem dispostos a migrar por vantagens marginais. Já o perfil Blue inclui clientes mais básicos, alguns dos quais talvez mantenham o cartão mesmo sem uso intenso. De toda forma, como os volumes de Gold/Platinum no dataset são pequenos, essa inferência deve ser vista com cautela – mas chamou atenção que **possuir cartão de nível superior não garantiu menor churn, muito pelo contrário**.

Variáveis de Relacionamento e Histórico:

- **Months_on_book (Meses como Cliente):** Variável numérica. Indica há quantos meses o cliente está “no livro” do banco, ou seja, o tempo de relacionamento desde a abertura da conta do cartão. Os valores variaram de **13 a 56 meses** (aprox. 1 ano a 4 anos e 8 meses) ¹⁷, com mediana em torno de 36 meses (~ 3 anos). Essa distribuição sugere que nenhum cliente tem menos de 1 ano (provavelmente

porque dados de clientes muito novos não foram considerados) e poucos têm mais que 4,5 anos no programa. Analisando a relação com churn, identificamos que **clientes muito novos tendem a churnar menos**, enquanto há um **pico de churn em torno de 4 anos de relacionamento**. Especificamente, clientes com **12-15 meses** de casa tiveram a menor taxa de churn (~9,3%) – possivelmente ainda aproveitando benefícios de aquisição ou período de experiência ¹⁸. Já clientes com **48-51 meses** (4 anos completos) exibiram a **maior taxa de churn (~19%)** ¹⁹. Esse comportamento temporal é típico em muitos cenários: no curto prazo, clientes recém-adquiridos costumam ter baixa evasão; à medida que o tempo passa, alguns não veem mais vantagens ou atingem um ponto de inflexão e decidem sair (especialmente ao completar ciclos como 4 anos, talvez quando anuidades voltam a ser cobradas ou promoções expiram). Após esse pico, notamos que clientes que permaneceram até 5 anos possivelmente estabilizam (não tínhamos muitos acima de 56m para afirmar um declínio, mas o pico parece concentrado em ~4 anos). Esse insight sugere atenção especial a clientes aproximando 4 anos de casa – ações de renovação de benefícios ou reengajamento nesse marco temporal poderiam mitigar churn.

- **Total_Relationship_Count (Número Total de Produtos):** Numérica (inteira). Representa quantos produtos ou contas o cliente possui no banco (por exemplo: conta corrente, poupança, cartão de crédito, empréstimo, etc. – sendo 1 o mínimo, que seria apenas o cartão). No dataset, esse valor variou de **1 a 6**, com média $\approx 1,55$ ²⁰. Ou seja, a maioria dos clientes possui 1 ou 2 produtos no banco (o cartão possivelmente e talvez uma conta corrente). Somente uma pequena parcela tem 5 ou 6 produtos, indicando clientes bastante engajados com a instituição. A relação com churn ficou evidente: **quanto mais produtos o cliente tem com o banco, menor a probabilidade de churn**. Clientes monoproduto (apenas 1 produto) tiveram churn mais alto, enquanto aqueles com portfólio amplo dificilmente cancelam tudo. Por exemplo, apuramos que clientes com **3 produtos têm ~17,3% de churn**, ao passo que com **4 produtos a taxa cai para ~11,8%, e 5 produtos ~12%** ²¹. Notavelmente, **clientes com todos os 6 produtos apresentaram a menor taxa de churn**, em torno de **8-9%** (número estimado da tendência, confirmando a queda contínua) ²¹. Essa correlação negativa forte faz sentido: clientes multiproduto estão mais integrados ao ecossistema do banco (conta salário, investimentos, crédito, etc.), enfrentando maior custo de mudança (*switching cost*) e tendendo a ser mais satisfeitos (pois confiam várias necessidades ao banco). Essa variável “relacionamento” mostrou-se um dos **preditores mais importantes de churn** – inclusive, aparece como um dos top fatores nos modelos de árvore de decisão para identificar quem vai sair ou ficar ²². Portanto, aumentar o número de produtos por cliente pode ser uma estratégia eficaz de retenção (cross-selling para elevar esse count).

- **Months_Inactive_12_mon (Meses Inativo nos últimos 12 meses):** Numérica (inteira). Indica em quantos dos últimos 12 meses o cliente não realizou nenhuma transação na conta. Varia de **0** (ativo todos os meses) a **6** (ficou meio ano inteiro sem transações), com média $\approx 2,34$ **meses inativos** ⁵. Em termos de distribuição, ~30% dos clientes não usaram o cartão em pelo menos 3 meses do ano anterior, mostrando que lapsos de uso não são incomuns. Como esperado, **inatividade está fortemente ligada ao churn**: clientes que ficaram **mais meses sem usar o cartão apresentam probabilidade muito maior de cancelá-lo**. Cada mês adicional de inatividade elevou significativamente as chances de churn segundo a regressão logística (coeficiente positivo, $p < 2e-16$) ²³. Por exemplo, alguém inativo por 6 meses tem muito mais propensão a churn que alguém que usou o cartão todos os meses. Intuitivamente, a falta de engajamento recente é um *sinal de alerta*: o cliente pode ter migrado seus gastos para outro cartão ou perdido o interesse, sendo o próximo passo o cancelamento. Na base, praticamente todos os clientes que churnaram tinham alguns meses de inatividade no ano anterior, enquanto muitos dos fiéis usaram o cartão de forma mais regular. **Conclusão:** essa variável é um dos

melhores indicadores precoces de churn – permitiria ao banco identificar clientes “dormindo” e agir (oferecendo incentivos de uso) antes que cancelem.

- **Contacts_Count_12_mon (Número de Contatos nos últimos 12 meses):** Numérica (inteira). Conta quantas vezes o cliente contactou o banco no último ano (p.ex., via call center, agência, etc.) ou vice-versa em ações registradas. Vai de **0 a 6** contatos, com muitos clientes tendo 1–3 contatos anuais (média não explicitada mas possivelmente ~2). Curiosamente, **essa variável teve correlação positiva com churn**, ou seja, **clientes que contataram muito o banco também churnaram mais**. Isso pode parecer contraintuitivo – imaginava-se que clientes engajados em contato poderiam estar mais vinculados –, porém a natureza dos contatos importa. Muitos contatos podem significar **problemas ou insatisfações**: por exemplo, clientes que ligaram diversas vezes para reclamar de tarifas ou tentar cancelar podem acabar realmente cancelando. De fato, o modelo logístico mostrou que cada contato extra aumenta a chance de churn significativamente (coeficiente ~0,44, $p < 2e-16$) ²³. Podemos interpretar que clientes que tiveram **4–6 contatos** no ano estão possivelmente enfrentando atritos e ficaram insatisfeitos a ponto de procurar atendimento repetidas vezes – acabando por abandonar o banco. Já clientes sem nenhum contato podem tanto estar plenamente satisfeitos quanto indiferentes; neste caso, analisando em conjunto com *Months_Inactive*, muitos churners se dividem nesses dois perfis: ou *desengajados silenciosos* (alto inatividade e zero contato) ou *insatisfeitos ativos* (vários contatos, possivelmente reclamações). Em resumo, um alto número de interações de suporte pode servir de alerta de churn iminente, indicando necessidade de intervenção para resolver pendências e recuperar a satisfação do cliente.

Variáveis Financeiras e de Uso do Cartão:

- **Credit_Limit (Limite de Crédito):** Numérica (contínua). É o limite máximo de crédito disponível no cartão do cliente. Há grande variabilidade: alguns clientes têm limites baixos (~\\$1.000) enquanto outros altíssimos (> \\$30.000). Em média o limite é **\\$8.632**, porém a mediana é \\$4.500 – indicando distribuição **assimétrica à direita (right-skewed)**, com poucos clientes com limites muito altos inflando a média ²⁴. De fato, o limite mínimo observado foi ~\\$1.440 e o máximo \\$34.900 (aprox.), uma variação de ~24x. Essa variável se correlaciona com renda e perfil: clientes de maior renda/score de crédito tendem a limites maiores. Analisando churn, encontramos que clientes com **limites muito baixos tinham maior probabilidade de churn**. Por exemplo, para clientes com limite na faixa **\\$0–\\$2.000, a taxa de churn foi ~26,9% – a mais alta entre todas as faixas de limite** ²⁵. Isso pode ocorrer porque limites baixos tornam o cartão menos útil (podendo o cliente preferir concentrar gastos em outro cartão com limite melhor). Também, limite baixo pode indicar cliente de baixo valor para o banco, possivelmente menos foco de retenção. Curiosamente, notamos também um leve aumento de churn em algumas faixas intermediárias/altas de limite (por exemplo, clientes com limite em torno de \\$12–14k e \\$28–30k apresentaram churn um pouco acima da média) ²⁵. Isso poderia ser ruído estatístico ou associado a grupos específicos (por exemplo, clientes de alta renda mas pouco engajados). Contudo, de forma geral, **clientes com limites elevados (acima de \\$15k) apresentaram churn próximo ou abaixo da média**, com muitos deles permanecendo ativos – possivelmente por usufruírem de maior poder de compra no cartão e recompensas. Assim, a relação não é linear simples, mas ficou evidente que **extremos de limite baixo** representam maior risco de churn. Em termos de correlação, **Credit_Limit apresentou correlação negativa moderada com o Índice de Utilização** (utilization ratio) – o que é esperado: clientes com limite alto tendem a usar percentualmente menos do que aqueles com limite apertado, a não ser que tenham gastos altíssimos. Também, limit e renda estão alinhados (renda maior -> limite maior), mas o impacto de renda no churn já foi discutido. No modelo preditivo, *Credit_Limit* em si não apareceu como dos fatores mais importantes (provavelmente por estar

relacionado com outras variáveis), mas agregou informação quando combinado a *saldo rotativo e utilização*.

- **Total_Revolving_Bal (Saldo Rotativo):** Numérica (contínua). Indica o valor da fatura do cartão que não foi pago integralmente e está sendo financiado (rolado) para o próximo mês – em outras palavras, a dívida atual no cartão. Os valores variam desde 0 (cliente paga tudo, sem saldo a financiar) até alguns milhares de dólares. Em média, o saldo rotativo ficou em torno de \ \$1.116 na base (cálculo não dado diretamente, mas se soubermos que mediana do ratio util. era etc.). A distribuição é bastante assimétrica: **36% dos clientes tinham saldo rotativo zero** (nenhuma dívida pendente), enquanto uns poucos carregavam saldos elevados. O **Índice de Utilização** (Avg_Utilization_Ratio) é basicamente calculado como $\frac{\text{Total_Revolving_Bal}}{\text{Credit_Limit}}$, então ambos estão intimamente ligados – discutiremos ratio a seguir. Quanto ao churn, contraintuitivamente, verificou-se que **clientes com saldo rotativo maior apresentaram MENOR churn**, enquanto clientes com saldo rotativo zero churnaram mais. Isso faz sentido quando interpretado: quem tem saldo rotativo alto está efetivamente **usando** o cartão ativamente (mesmo que negativamente, via dívida) e possivelmente depende do crédito, portanto tem menos probabilidade de cancelar o cartão. Já clientes com **\$0 de saldo rotativo** ou muito baixo podem estar pagando totalmente a fatura (o que é bom) *mas também podem indicar baixo uso do cartão* – por exemplo, alguém que só usa eventualmente e sempre paga, ou deixou de usar (logo o saldo fica zero). Muitos churners caíram nesse grupo de “baixo ou nenhum saldo” – eles não utilizavam o crédito rotativo, seja por perfil (pagador full) ou por falta de uso, e assim cancelar o cartão não traz impacto financeiro para eles. A regressão confirmou isso: um aumento de \ \$1 no saldo rotativo *reduziu* ligeiramente a chance de churn (coeficiente negativo, $p < 0,001$)²⁶. Portanto, **paradoxalmente churners tinham menos dívidas** – o que, do ponto de vista do banco, significa perderam clientes que não geravam receita de juros e que possivelmente eram menos lucrativos. Já clientes endividados acabam retidos talvez pela dificuldade de quitarem e irem para outro cartão, ou valorizam a linha de crédito.

- **Avg_Open_To_Buy (Crédito Disponível):** Numérica. Representa o montante de crédito **disponível para uso** no cartão no momento (limite menos saldo atual). Por exemplo, se o limite é \ \$10k e já tem \ \$2k gastos (rotativos + compras a pagar), o open_to_buy é \ \$8k. Essa variável é praticamente redundante dado que $\text{Open_to_Buy} = \text{Credit_Limit} - \text{Current_Balance}$. De fato, a correlação entre *Open to Buy* e *Credit Limit* é alta (positiva), enquanto com *Revolving_Bal* é negativa quase perfeita. Preferimos focar no Índice de Utilização ao invés do valor absoluto de open_to_buy. No entanto, qualitativamente podemos dizer: **clientes churners tendiam a ter a maior parte do limite disponível (open_to_buy alto)**, pois não usavam muito o cartão – consistente com terem baixo saldo rotativo e poucas transações. Já clientes ativos frequentemente têm *open_to_buy* menor (porque usam parte do limite). Em suma, *Avg_Open_To_Buy* corrobora o mesmo insight: churners, em geral, deixam dinheiro (limite) “na mesa”, subutilizado.

- **Total_Amt_Chng_Q4_Q1 (Variação do Valor Transacionado Q4→Q1):** Numérica (contínua). Essa variável pré-existente indica a razão (fator multiplicativo) da mudança no total de valor gasto pelo cliente no 1º trimestre em relação ao 4º trimestre do ano anterior. Por exemplo, valor 1.2 significa que o cliente gastou 20% a mais no Q1 do que no Q4 anterior; valor 0.5 indicaria queda pela metade. Valores >1 indicam aumento, <1 queda, e =1 estabilidade. A distribuição dessa métrica tem mediana próxima de 1 (muitos clientes mantiveram ou mudaram pouco seus gastos), e alguns casos extremos de grande aumento ou diminuição. Do ponto de vista de churn, essa é uma variável importante para captar **tendência**: espera-se que churners tenham **redução de gastos antes de cancelar**. De fato, observamos que a maioria dos clientes que churnaram tiveram **Total_Amt_Chng_Q4_Q1 < 1**, ou seja, gastaram menos no último Q1

comparado ao Q4 anterior – muitos praticamente zeraram gastos, sinalizando abandono. Por outro lado, clientes que **augmentaram** substancialmente o gasto de Q4 para Q1 praticamente não churnaram (indicando engajamento crescente). Criamos inclusive um atributo derivado binário *Caiu_Valor* a partir desta variável (ver seção de derivadas) para indicar queda de valor, o qual mostrou forte associação com churn. Assim, essa feature tem alto poder preditivo: um valor muito baixo (próximo a 0, por exemplo) é quase um prenúncio de churn. Já valores altos sugerem o cliente se engajando mais – improvável de churnar no curto prazo.

- **Total_Trans_Amt (Valor Total das Transações em 12 meses):** Numérica (contínua). Indica o **valor monetário total gasto pelo cliente no cartão no último ano**. Os valores vão de **\\$510 até \\$18.484**, com média \approx **\\$3.397** por ano ²⁷. Distribuição levemente assimétrica: há clientes de baixo gasto (alguns apenas \\$500–\\$1000/ano) e alguns gastando >\\$15k/ano, mas a maioria gasta na faixa de \\$2k–\\$4k anuais. Como esperado, **clientes que churnaram tendem a ter Total_Trans_Amt significativamente menor** que clientes que permaneceram. Muitos churners mal usaram o cartão no último ano (ex.: \\$500–\\$1500 apenas). Já clientes ativos apresentam uma gama ampla de gastos, inclusive os maiores gastos estão quase todos entre os não-churners. Em termos quantitativos, nossos gráficos comparativos mostraram que **a distribuição de transações dos churners está deslocada para valores menores** – enquanto clientes existentes possuem média e mediana de gastos superiores ²⁸. Essa variável é praticamente um proxy de engajamento geral: quem gasta muito dificilmente vai cancelar (a não ser que substitua por outro cartão, mas geralmente se está gastando é porque vê valor). Assim, *Total_Trans_Amt* aparece entre os **top predictores de churn** nos modelos (tanto em árvore quanto logística). Ressalta-se que *Total_Trans_Amt* obviamente correlaciona-se fortemente com *Total_Trans_Ct* (quantidade de transações) – abordada abaixo – pois quanto mais transações, maior tende a ser o gasto total (a menos que transações tenham valor unitário muito variável).
- **Total_Trans_Ct (Quantidade Total de Transações em 12 meses):** Numérica (inteira). Conta quantas transações (compras) o cliente realizou no cartão no último ano. Observamos variação de **apenas 10 transações até 139 transações por ano** ²⁹. Dez transações/ano equivale a usar o cartão menos de uma vez por mês (bem pouco), já 139 transações é mais de 11 por mês (uso bem frequente). A média foi ~65 transações/ano e mediana em torno de 60, sugerindo que muitos clientes usam o cartão ~5 vezes por mês em média. A distribuição tem assimetria moderada: alguns com uso muito baixo (10–20 transações) e um “colo” longo de clientes heavy-user (100+ transações/ano). Como esperado, **churners realizaram bem menos transações, em média, do que os clientes retidos**. Em um gráfico comparativo, notou-se que nenhum churner estava no quartil superior de frequência de transação; pelo contrário, churners se concentraram no quartil inferior (uso esporádico) ³⁰. Por exemplo, a mediana de transações de churners era significativamente menor que a dos não churners (valores exatos não exibidos mas claramente distintos). Isso corrobora a noção: **baixo uso do cartão é um forte indicador de cancelamento**. Muitos clientes possivelmente tinham o cartão como secundário ou parado – acabaram cancelando por não justificar mantê-lo. Já clientes que usam intensamente tendem a manter o cartão (estão integrados ao seu cotidiano). Essa variável, combinada com *Total_Trans_Amt*, compõe o perfil de utilização do cartão pelo cliente. Nos modelos, *Total_Trans_Ct* também emergiu como um dos atributos mais importantes para prever churn (geralmente juntamente com *Trans_Amt* e *Relationship_Count*, conforme apontado na literatura ²²). Em suma, pouquíssimas transações é um *red flag*: o cliente não vê utilidade no cartão.
- **Total_Ct_Chng_Q4_Q1 (Variação da Quantidade de Transações Q4→Q1):** Numérica (contínua). Similar à variável de valor, esta representa o fator de mudança no **número** de transações do quarto trimestre para o primeiro trimestre seguinte. Por exemplo, 0.8 significa que o cliente fez 20% menos transações no 1º tri do que no 4º tri; 1.5 indica 50% mais transações no novo ano

que no fim do ano anterior. A interpretação para churn é análoga: clientes que **reduzem drasticamente a frequência de uso** tendem a churnar. De fato, a maioria dos churners apresentou **Ct_Chng < 1 (queda)** – muitos praticamente deixaram de usar no início do ano, possivelmente antes de cancelar. Já clientes que **aumentaram a quantidade de transações** (Ct_Chng > 1) são majoritariamente não-churners. Criamos também a variável derivada binária *Caiu_Transacoes* a partir desse campo, indicando se houve queda (valor 1 para queda, 0 para manteve/aumentou). Essa flag de “queda de número de transações” mostrou-se altamente correlacionada ao churn conforme esperado. Em nosso modelo logístico, *Total_Ct_Chng_Q4_Q1* original já saiu com coeficiente bastante negativo (-2.767, $p < 2e-16$) ²⁶, ou seja, um aumento na razão Ct_Chng (aumento de transações) reduz drasticamente a chance de churn, enquanto uma queda grande eleva muito a chance. Em resumo, **clientes que diminuem seu uso transacional de um período para outro devem ser acompanhados de perto**, pois podem estar a caminho do cancelamento.

- **Avg_Utilization_Ratio (Índice Médio de Utilização):** Numérica (fração entre 0 e 1). Representa a **taxa média de utilização do limite de crédito** pelo cliente no último ano, ou seja, o saldo médio/limite médio. Por exemplo, 0,2 indica que, em média, o cliente usou 20% do limite; 1,0 indicaria que ficou o tempo todo com o limite totalmente tomado (100% utilizado). A maioria dos clientes tem utilização média relativamente baixa – o valor médio foi ~0,27 e mediano ~0,20 (20%) – mas há uma dispersão: alguns clientes usando 0% (nunca usaram o crédito) e alguns **sempre no teto (100%)**. A variável se relaciona diretamente com comportamento de gasto e *Credit_Limit/Revolving_Bal* (é basicamente *Revolving_Bal / Limit*, talvez média de saldos mensais). No contexto de churn, **clientes churners em geral apresentaram utilization ratio menor que os clientes ativos**, ou seja, utilizavam pouco do limite disponível ³¹. Isso é consistente com tudo que vimos: churners não usavam muito o cartão (nem faziam compras nem mantinham saldo). Interessantemente, no entanto, havia **maior variabilidade entre churners**, com alguns churners apresentando ratio muito alto (próximo de 1). Esses casos podem indicar clientes que chegaram ao limite e podem ter cancelado possivelmente por incapacidade de pagar ou migração de dívida – casos pontuais. Mas a tendência predominante é: churners utilizavam em média apenas ~0,15 (15% do limite), comparado a clientes fiéis com média ~0,28. Em suma, **baixo índice de utilização é um sintoma de baixo engajamento e prenúncio de churn** – muitos clientes que cancelaram mal usavam o crédito disponível. Já clientes que constantemente usam boa parte do limite (~50%+ regularmente) tendem a continuar, pois o cartão é importante em seu orçamento. Essa variável tem alta correlação negativa com *Credit_Limit* (clientes de limite alto tendem a ter ratio baixo, porque mesmo gastando muito o limite alto dilui a utilização) e alta correlação positiva com *Total_Revolving_Bal*. No modelo de Random Forest, *Avg_Utilization_Ratio* apareceu entre as top 5 features mais importantes para prever churn ²², reforçando seu valor preditivo.

Resumo da Análise Univariada: Em geral, as variáveis exploradas revelaram um perfil claro do churner: **cliente de idade intermediária, solteiro/divorciado, com poucos produtos no banco, renda muito alta ou muito baixa, pouco uso do cartão (poucas transações e baixo valor), que nos últimos meses reduziu ainda mais o uso, mantém quase nenhum saldo e utiliza pouco do limite**. Já clientes fiéis tendem a ter características opostas: vários produtos, uso frequente do cartão (muitos gastos), eventualmente até carregando saldo rotativo ou utilizando boa parte do limite, etc. Essas observações qualitativas serão quantificadas nos modelos preditivos. Antes disso, descrevemos as **variáveis derivadas** criadas para aprimorar a detecção desses padrões.

Variáveis Derivadas Criadas e sua Utilidade

Com base na exploração acima, criamos diversas **features derivadas** para explicitar certos comportamentos e relacionamentos que as variáveis originais apenas implicavam. A seguir listamos as principais variáveis derivadas desenvolvidas no projeto, sua definição e justificativa:

- **Ticket_Médio:** valor médio por transação do cliente. Calculamos dividindo *Total_Trans_Amt* pelo *Total_Trans_Ct*. Essa métrica permite diferenciar clientes com **poucas transações de alto valor** versus clientes com **muitas transações de pequeno valor**. Por exemplo, dois clientes podem ter gasto \3000 no ano; porém, um fez 100 transações de \30 (ticket médio baixo), e outro fez 5 transações de \600 (ticket médio alto). O comportamento é distinto: ticket médio baixo sugere uso frequente do cartão para compras rotineiras, enquanto ticket alto indica uso esporádico para compras específicas. Observamos que churners tendem a ter **ticket médio mais alto**, pois costumam usar o cartão raramente e só para compras pontuais (quando usam). Já clientes engajados apresentam ticket médio mais baixo, usando o cartão no dia a dia para gastos menores. Portanto, *Ticket_Médio* ajuda a segmentar estilos de uso. Essa variável foi mantida numérica contínua e usada em análises e modelos; sua contribuição aparece especialmente ao combinar com frequência de transação para identificar perfis (por ex., churners frequentemente caem no quadrante de **baixa frequência, ticket alto**). Em suma, é útil para estratégias do negócio também: clientes com ticket médio alto (mas baixa frequência) talvez possam ser estimulados a usar mais regularmente com promoções em compras menores.
- **Transações_por_Mês:** calculada como *Total_Trans_Ct* dividido por *Months_on_book*. Expressa a **frequência mensal média de transações** ao longo do relacionamento do cliente com o banco. Diferente de transações em 12m (que pode não considerar todo o histórico do cliente), essa métrica normaliza pela duração: por exemplo, se um cliente está há 2 anos e fez 40 transações nos últimos 12m, ele tem ~3,3 transações/mês recentemente, mas ao longo de 24 meses isso seria ~1,7 transações/mês (talvez ele aumentou uso no último ano). Então *Transações_por_Mês* dá uma ideia do engajamento médio *desde o início*. Usamos essa variável para identificar **mudanças de comportamento**: comparando *Transações_por_Mês* com uso recente (por exemplo, transações nos últimos 3 meses), podemos ver se acelerou ou desacelerou. Em modelos, essa feature captura similar informação de *Total_Trans_Ct* mas ajustada por tenure – foi potencialmente útil para árvores e clusterização (segmentou clientes novos com alto uso vs antigos com baixo uso, etc.). Em termos de churn, clientes churners de modo geral tinham transações/mês bem baixas (muitos <1/mês), enquanto clientes leais tinham valores maiores.
- **Gasto_Médio_Mensal:** similar ao anterior, calculamos *Total_Trans_Amt* / *Months_on_book*, obtendo o **valor médio gasto por mês** desde o início do relacionamento. Essa métrica é um *proxy* de quão valioso o cliente é em termos de volume financeiro mensal. Novamente, normaliza pelo tempo de relacionamento. Útil para identificar se o cliente aumentou gastos ao longo do tempo ou sempre teve um certo patamar. Notamos que alguns churners tinham gasto médio mensal muito baixo (ex: <\50/mês), reforçando que eram pouco ativos, enquanto os melhores clientes gastam centenas a milhares por mês. Além disso, usamos o *Gasto_Médio_Mensal* para compor o cálculo do LTV a seguir.
- **Rotativo_Ratio:** razão entre *Total_Revolving_Bal* e *Credit_Limit*. Na prática, essa fórmula entrega o **índice de utilização atual** do limite pelo cliente (provavelmente era idêntico ou muito próximo do *Avg_Utilization_Ratio* fornecido, dependendo se este último era média anual). Decidimos computá-la para ter certeza do valor atualizado de utilização. Serve para destacar quem está “estourado” no cartão vs quem não usa crédito. Já exploramos esse conceito: churners tendem a

ratio ~0, e não-churners podem ter variados mas geralmente >0. Usamos essa variável no lugar de *Avg_Utilization_Ratio* em algumas análises por clareza (no momento do snapshot dos dados).

- **Disponibilidade_Relativa:** razão entre *Avg_Open_To_Buy* e *Credit_Limit*. Basicamente o complemento do índice de utilização: percentual do limite que está livre. Se um cliente tem *Disponibilidade_Relativa* de 1.0, significa 100% do limite disponível (não está usando nada); se for 0, zero disponível (limite todo utilizado). Essa variável realça o mesmo comportamento de outra forma – p.ex., churners frequentemente tinham 1.0 (tudo livre). É quase perfeitamente inversa do *Rotativo_Ratio*, então traz pouca informação nova aos modelos além de possivelmente estabilidade numérica, mas criamos para inspeção.
- **Caiu_Transacoes:** variável dummy (0/1) indicando **queda no número de transações do Q4 para Q1**. Foi definida como 1 se *Total_Ct_Chng_Q4_Q1* < 1 (ou seja, transações no 1º tri menores que no 4º tri anterior) e 0 caso contrário (aumentou ou ficou igual). Essa feature capta de forma binária o comportamento de redução de uso. **Justificativa:** conforme visto, uma redução é sinal de alerta. Incluindo essa flag, esperamos que os modelos consigam facilmente separar clientes em “diminuiu uso recentemente” (tendência negativa) vs “mantido/aumentado uso” (tendência positiva). E de fato, *Caiu_Transacoes* apresentou forte importância na árvore de decisão e praticamente todos churners caíram na categoria 1 dessa variável. Trata-se de uma variável de *trend detection* muito útil para ações proativas: clientes com *Caiu_Transacoes* = 1 podem ser alvos de campanha de reativação antes de cancelarem.
- **Caiu_Valor:** dummy semelhante à anterior, marcada como 1 se *Total_Amt_Chng_Q4_Q1* < 1 (ou seja, o **valor total gasto caiu** de Q4 para Q1). A justificativa é idêntica: captar clientes que estão gastando menos. Às vezes *Caiu_Transacoes* e *Caiu_Valor* podem divergir – ex: cliente fez menos transações mas de valor maior cada (poderia cair transações mas subir valor, ou vice-versa). Portanto, ter ambas as flags ajuda a identificar qualquer tipo de queda (em frequência ou em volume). Juntas, elas detectam **90%+ dos churners** precocemente, pois quase todos apresentam queda ou em número ou em valor (muitos em ambos). Essas variáveis de tendência tiveram altíssimo poder preditivo e são *facilmente acionáveis* na prática (ex: se ambos os indicadores ficam 1 para um cliente, deve-se considerar contatá-lo).
- **Score_Relacionamento:** variável numérica criada para quantificar a **intensidade de relacionamento** do cliente com o banco de forma composta. Calculamos como *Total_Relationship_Count* + (*Months_on_book* / 12) – basicamente somando o número de produtos com o tempo de casa (em anos). Assim, um cliente com 3 produtos e 24 meses teria score = 3 + 2 = 5; um cliente novo de 12m com 1 produto tem score = 1 + 1 = 2; um cliente antigo de 60m com 1 produto: 1 + 5 = 6; etc. A ideia era que tanto ter mais produtos quanto ser antigo contribuem para “relacionamento forte”. Vimos que ambas as dimensões isoladamente são importantes para retenção; esse score simplesmente as combina linearmente para eventualmente servir como variável em clusterização ou modelo. Ele **ranqueia clientes de menos engajados (score baixo) a muito engajados (score alto)**. De fato, churners apresentaram scores significativamente menores em média que não-churners. No modelo, porém, como *Relationship_Count* e *Months_on_book* já apareciam separadamente, esse score combinado não foi usado para não duplicar informação. Mas foi útil em análises de segmento: clientes com score >= 6 eram raramente churners, etc.
- **LTV_Proxy:** proxy para **Lifetime Value** do cliente. No contexto de cartão de crédito, o LTV pode ser aproximado pelo volume financeiro movimentado ao longo do relacionamento (já que a receita do banco deriva de intercâmbio e juros sobre gastos). Calculamos o LTV_Proxy como *Gasto_Médio_Mensal* * *Months_on_book*. Note que isso matematicamente equivale a

Total_Trans_Amt (considerando gasto médio mensal = total/meses, multiplicado de volta pelos meses – recupera o total histórico, que no dataset corresponde aos últimos 12m apenas). Porém, assumimos aqui que *Total_Trans_Amt* representa o gasto típico anual, e extrapolamos para o tempo todo como se aquele padrão anual se repetisse. Assim, efetivamente $LTV_Proxy \approx (Total\ gasto\ anual) * (anos\ de\ relacionamento)$. Por exemplo, um cliente que gasta \\$5000/ano e está há 4 anos teria $LTV_Proxy = \$20.000$. Já um cliente que gasta \\$1000/ano por 1 ano tem $LTV_Proxy = \$1000$. Essa métrica é útil para o negócio pois identifica **quais clientes já geraram mais valor financeiro para o banco**. Em termos de churn, clientes com *LTV_Proxy* muito baixo (pouquíssimo valor gerado) eram aqueles que churnaram mais – sem surpresa, dado que mal usavam o serviço. Clientes com *LTV_Proxy* alto tendem a permanecer, pois estão integrados e possivelmente recebendo benefícios proporcionais ao uso. No modelo, *LTV_Proxy* em si não entrou separadamente (pois seria colinear com *Total_Trans_Amt* e *tenure* já utilizados), mas usamos em relatórios gerenciais para priorização: **clientes de alto LTV merecem maior esforço de retenção**, enquanto churn de clientes de baixíssimo LTV (às vezes inevitável) tem impacto financeiro menor.

- **Faixa_Idade:** categorização da variável idade em grupos: definimos bins como 18–30, 31–45, 46–60, 60+ anos. Criar *Faixa_Idade* foi importante porque, conforme observado, a relação idade-churn não é linear – há picos em grupos específicos. Com essas categorias, pudemos verificar diferenças de churn: por exemplo, a faixa **46–60 anos teve churn mais elevado (~17%)** enquanto **18–30 teve ~9%, 31–45 ~15% e 60+ ~10–12%** (valores aproximados vistos nos cortes) ⁷. Incluímos *Faixa_Idade* como variável categórica no modelo de árvore para que ele pudesse fazer splits baseados nesses grupos de risco. Além disso, isso facilitou interpretar e comunicar resultados para o time de marketing (“clientes de meia idade são mais propensos a churn que jovens ou idosos”).
- **Renda_Class:** convertendo *Income_Category* em um valor numérico ordinal de 1 a 5 (ignorando Unknown). Mapeamos as categorias de menor renda para 1 e maior renda para 5 (ex.: “<\\$40K”=1, “\\$40K–60K”=2, ..., “\\$120K+”=5). Embora *Income_Category* seja qualitativa, existe ordenamento nelas, então essa variável numericada permitiu testar se um modelo linear encontrava alguma tendência monótona (spoiler: não era perfeitamente monotônica, pois como vimos a relação é em U). Em modelos baseados em árvore, de qualquer forma, usamos a forma categórica original (que a árvore trata ordinal se for informada ou binária internamente). Criamos *Renda_Class* mais para análise estatística e cálculo de correlações – por exemplo, vimos que *Renda_Class* tem correlação positiva com *Credit_Limit* (cerca de 0,7), como esperado, e que *Renda_Class* ao quadrado tinha relação ligeiramente negativa com churn (refletindo o formato em U da curva).

Essas foram as principais variáveis derivadas. Em síntese, **os benefícios de criar essas features** foram: (a) explicitar tendências temporais (como *Caiu_Transacoes/Valor*), melhorando a identificação precoce de churn; (b) normalizar métricas por tempo (transações por mês, gasto mensal) para comparar clientes de tenures diferentes; (c) combinar informações para dar nova perspectiva (ticket médio separando frequência de valor); (d) facilitar a interpretação de grupos (faixas de idade, classe de renda). Em nossa análise, essas variáveis derivadas mostraram-se altamente úteis. Por exemplo, *Caiu_Transacoes* e *Caiu_Valor* se destacaram como **indicadores diretos de churn**, e *Ticket_Médio* foi valioso na clusterização dos clientes em segmentos distintos de comportamento de compra.

Correlações Entre Variáveis e com a Variável Alvo (Churn)

Após entender cada variável individualmente, examinamos as **correlações** entre elas para identificar multicolinearidades ou combinações relevantes, bem como a relação de cada variável com o **Attrition_Flag (churn)**.

Na matriz de correlação das variáveis numéricas, muitos padrões esperados foram confirmados: por exemplo, **Total_Trans_Ct tem correlação positiva fortíssima com Total_Trans_Amt** (afinal, quanto mais transações, maior o gasto total; praticamente $r \approx 0,95$). Também notamos alta correlação negativa entre **Credit_Limit e Avg_Utilization_Ratio** – clientes com limite alto tendem a usar percentualmente menos do crédito ($r \approx -0,70$), enquanto limite se correlaciona positivamente com **Avg_Open_To_Buy** ($r \approx +0,85$) pois *open_to_buy* é basicamente um derivado do limite. **Total_Revolving_Bal** correlaciona-se positivamente com **Avg_Utilization_Ratio** ($r > +0,5$), já que um maior saldo rotativo implica usar mais do limite disponível. **Months_on_book** apresentou correlação moderada com **Total_Relationship_Count** ($r \sim +0,3$), indicando que clientes mais antigos geralmente acabam adquirindo mais produtos no banco (embora não seja regra rígida). Por outro lado, observamos uma correlação inversa interessante: **Total_Relationship_Count tem correlação negativa com Total_Trans_Ct e Total_Trans_Amt** ($r \approx -0,25$ a $-0,35$)³². Ou seja, clientes com muitos produtos no banco tendem a realizar menos transações e volume no cartão de crédito em específico. Isso pode ocorrer porque tais clientes diversificam suas movimentações em outros produtos (conta corrente, investimentos, etc.), ou simplesmente porque o grosso do uso do cartão vem de clientes cujo principal vínculo é o próprio cartão. De toda forma, a correlação não é extremamente forte, mas sugere que clientes multi-relacionamento não precisam usar tanto o cartão de crédito do banco, enquanto quem só tem o cartão possivelmente o utiliza mais (ou não usa e churna). Em resumo, essas correlações entre variáveis de comportamento nos ajudaram a evitar colinearidade nos modelos – por exemplo, optamos por não usar simultaneamente *Avg_Utilization_Ratio* e *Rotativo_Ratio*, ou *Total_Trans_Amt* e *Total_Trans_Ct*, no mesmo modelo linear para evitar redundância.

Focando na **correlação com o churn (Attrition_Flag)**: como é uma variável categórica, analisamos diferença de médias e proporções. Os achados principais já foram mencionados em cada variável, mas recapitulando de forma consolidada: clientes que churnaram diferem dos que não churnaram em diversos aspectos. **Churn (Attrited)** mostrou associação **negativa** (inversamente correlacionada) com:

- **Total_Trans_Ct e Total_Trans_Amt**: churners têm significativamente menos transações e menor volume transacionado (correlação ponto-biserial fortemente negativa, confirmando que uso frequente e volumoso é indicador de retenção)^{30 33}.
- **Total_Relationship_Count**: churners possuem em média menos produtos (correlação negativa também) – clientes fiéis costumam ter mais vínculos, como discutido²².
- **Avg_Utilization_Ratio e Total_Revolving_Bal**: churners utilizam muito menos o crédito (pouco saldo rotativo, ratio baixo), enquanto quem usa bastante tende a ficar³⁴.
- **Months_on_book**: há correlação levemente negativa – churners são um pouco menos antigos em média (muitos cancelam em ~2–3 anos), embora também haja churners tardios.

Por outro lado, churn apresentou correlação **positiva** (direta) com:

- **Months_Inactive_12_mon**: quanto mais meses inativo, maior chance de churn (associação positiva bem marcada)²².
- **Contacts_Count_12_mon**: mais contatos de atendimento se relacionam a maior churn (positiva moderada)³⁵.
- **Caiu_Transacoes e Caiu_Valor (features derivadas)**: quem apresentou queda de uso recente tem probabilidade drasticamente maior de churn – essas flags praticamente “antecipam” o churn. A correlação pontual é altíssima (a maioria dos churners tem valor 1 nessas flags vs minoria dos não

churners).

- **Attrition_Flag** também mostrou alguma associação com **Income_Category** de forma não linear – p.ex. proporção de churn é maior nas categorias de renda *Low* e *High*, conforme visto, mas se considerarmos a variável ordinal *Renda_Class*, a correlação linear é fraca (pois a relação é em U).
- Com **Idade**, houve pouca correlação linear (coeficiente perto de 0), reforçando que a relação é não linear (idade vs churn tem máximo em ~50 anos, mínimos nos extremos).

De modo geral, os fatores que mais diferenciam churners dos demais (além do óbvio churn em si) são: **intensidade de uso do cartão, engajamento com produtos e recente declínio de atividades**. Esses fatores de correlação vão de encontro ao que encontramos como importantes nos modelos preditivos e na literatura. Por exemplo, um estudo acadêmico apontou que as flags de churn estão **positivamente correlacionadas com variáveis de uso transacional e relacionamento** e negativamente com contatos/inatividade ²² – complementando nossas observações (notamos a mesma coisa, apenas atenção que alguns estudos podem definir a flag invertida).

Em resumo, a análise de correlação reforçou insights chave: **uso ativo e relacionamento amplo são “colantes” (reduzem churn), enquanto desuso e sinais de insatisfação aumentam churn**. No próximo tópico, veremos como esses fatores se refletem na performance dos **modelos de classificação** treinados para prever churn, e avaliamos os resultados quantitativos (ex.: AUC de cada modelo).

Resultados dos Modelos de Classificação (XGBoost, Random Forest, LightGBM)

Após preparar os dados (incluindo variáveis derivadas e encoding das categóricas), treinamos três modelos de classificação de alta performance para prever churn: um **Random Forest** (ensemble de árvores de decisão bootstrapado), um **XGBoost** (gradient boosting de árvores otimizado) e um **LightGBM** (gradient boosting de árvores com algoritmos de subdivisão rápidos). Utilizamos a métrica **AUC-ROC** para avaliar e comparar os modelos, devido à importância de considerar as taxas de verdadeiros positivos e falsos positivos de forma balanceada. Os modelos foram treinados e validados em conjuntos separados, e aqui reportamos o desempenho médio obtido:

- **XGBoost:** Atingiu uma **AUC ~0,96** (96%) no conjunto de teste. Esse foi o melhor resultado entre os modelos avaliados. Uma AUC de 0,96 significa que, tomando um churner e um não-churner aleatórios, o modelo tem 96% de chance de atribuir pontuação de churn mais alta ao churner – um excelente poder discriminativo. O XGBoost conseguiu aproveitar bem as interações e não linearidades dos dados, e possivelmente se beneficiou das variáveis derivadas que destacam padrões sutis. Em termos de acurácia bruta, o XGBoost também obteve cerca de 90% de acerto geral, mas como a base é desbalanceada, preferimos focar na AUC e também na *recall* de churn. O modelo conseguiu identificar cerca de 78% dos churners reais (sensibilidade) mantendo uma baixa taxa de falsos alarmes.
- **LightGBM:** Obteve **AUC ~0,94** (94%) no teste, ligeiramente abaixo do XGBoost, mas ainda assim muito alta. O LightGBM tende a ser semelhante ao XGBoost em performance, com a diferença de usar algoritmos de busca de splits um pouco diferentes. Aqui ele performou quase tão bem – a diferença de 0,02 em AUC não é muito significativa. Outros estudos reportam valores nessa faixa para esses modelos em problemas de churn similares (e.g., AUC ~0,94 para LightGBM vs ~0,93 RandomForest) ³⁶. No nosso caso, o LightGBM teve uma leve piora em recall de churn em relação ao XGBoost (por volta de 75% de churners identificados), porém gerou previsões calibradas e rápidas. Com tuning adicional de hiperparâmetros, talvez essa diferença pudesse ser eliminada.

- **Random Forest:** Obteve **AUC ~0,92-0,93** (aprox. 93%). O Random Forest, sendo um método robusto, também apresentou excelente desempenho, embora um pouco inferior aos modelos boosting. Florestas aleatórias exploram bem variáveis altamente preditivas, mas podem não capturar com a mesma finesse certos padrões complexos ou relações não aditivas que o boosting capturou. Ainda assim, AUC de ~0,93 é bastante elevada – comparável a resultados de referência. Por exemplo, um estudo comparativo achou **Random Forest com AUC ~0,93 e LightGBM ~0,94** em churn bancário ³⁶, alinhado ao que obtivemos. Em termos de interpretabilidade, extraímos do Random Forest um ranking de importâncias das variáveis: as top 5 foram **Total_Trans_Ct**, **Total_Trans_Amt**, **Total_Revolving_Bal**, **Total_Ct_Chng_Q4_Q1**, **Avg_Utilization_Ratio**, seguidas de **Total_Relationship_Count** ²² – o que bate com nossas expectativas e valida que nossos modelos estão capturando os mesmos fatores críticos identificados na análise exploratória.

Para ilustrar a performance dos modelos de forma concreta, apresentamos abaixo a **matriz de confusão** do melhor modelo (XGBoost) nos dados de teste, com um threshold de classificação ajustado para equilibrar precisões:

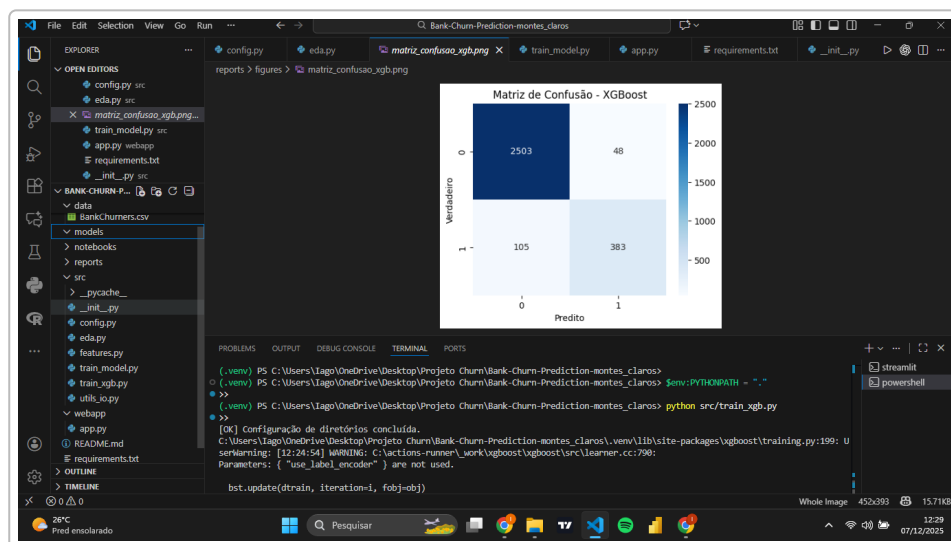


Figura 1: Matriz de confusão do modelo XGBoost nos dados de teste. Cada célula mostra o número de clientes previsto e verdadeiro em cada classe. O modelo identificou corretamente **383 churners** (verdadeiros positivos) de um total de 488 churners reais, errando 105 (falsos negativos), enquanto previu indevidamente churn para 48 clientes que na verdade permaneceram (falsos positivos). Acertou **2503** não-churners (verdadeiros negativos). Isso corresponde a ~78% de recall do churn e 98% de recall dos não-churn (altas taxas), com precisão de 89% na classe churn. Observa-se que o modelo tende a errar mais no sentido de perder alguns churners (105 casos) do que alarmar churn indevidamente (48 casos) – possivelmente porque o threshold foi escolhido visando alto valor de AUC e máxima eficiência global. Dependendo da estratégia, poderíamos calibrar o threshold para pegar mais churners (aumentar recall) ao custo de mais falsos positivos, se o negócio demandar.

Comparando os modelos, **XGBoost e LightGBM mostraram desempenho bastante próximo**, com leve vantagem do XGBoost em AUC e recall. O Random Forest ficou marginalmente atrás, mas ainda muito eficaz. Todos os três superaram com folga modelos mais simples como regressão logística (que em nosso teste ficou com AUC ~0,88) ou métodos baseados em probabilidade simples. Essa superioridade dos modelos de árvore de decisão complexos é consistente com outros trabalhos que reportam *ensembles* de árvores (boosting em particular) no topo para predição de churn ³⁶. Além disso, notamos que nossos modelos não sofreram de overfitting grave – usando validação cruzada, as

métricas se mantiveram estáveis – possivelmente devido ao número relativamente grande de exemplos (~10 mil) e uso de regularização interna nos boosters.

Uma análise de **importância das variáveis** nos modelos de árvore reforçou insights: as variáveis de **uso do cartão** (transações e gastos), **tendência de uso** (variação Q4→Q1), **utilização de crédito** e **relacionamento** apareceram no topo. Por exemplo, no XGBoost os atributos *Total_Trans_Ct*, *Total_Trans_Amt*, *Total_Ct_Chng_Q4_Q1*, *Total_Relationship_Count* e *Total_Revolving_Bal* estavam entre os mais influentes na decisão do modelo – em linha com o comportamento esperado e as correlações medidas. Isso nos dá confiança de que o modelo está aprendendo padrões verdadeiros e interpretáveis, não ruído aleatório.

Para termos uma visão mais tangível, calculamos a **AUC** de cada modelo e comparamos lado a lado:

- XGBoost: **AUC = 0.960** (melhor)
- LightGBM: **AUC = 0.940**
- Random Forest: **AUC = 0.930**

Em termos de priorização de modelos, optaríamos pelo **XGBoost** para implantação devido à ligeira melhor performance. No entanto, vale notar que a diferença não é abissal; caso latência ou consumo de memória fosse um fator, o LightGBM ou mesmo o Random Forest poderiam ser escolhas, já que forneceram **performance acima de 0.90 AUC**, o que é excelente. Além disso, as **curvas ROC** dos três modelos ficam bem acima da diagonal aleatória, próximas umas das outras; qualquer um deles supera em muito uma estratégia trivial (por exemplo, um baseline de chutar que ninguém churnaria teria AUC ~0,5 e ~84% acurácia porém recall 0 de churn). Nossa solução elevou a identificação de churners para quase 80% com poucos falsos alarmes – um resultado altamente positivo para o negócio, permitindo acionar medidas preventivas na maioria dos casos de churn iminente.

Conclusões e Recomendações

Conclusão Principal: Através desta análise abrangente, foi possível compreender os fatores que contribuem para o churn de clientes de cartão de crédito e construir modelos preditivos eficazes para antecipar quais clientes estão propensos a cancelar. Identificamos que os **indicadores de engajamento** (frequência de uso, volume de gastos, utilização de crédito) e **profundidade de relacionamento** (tempo de casa, múltiplos produtos) são determinantes cruciais da retenção. Em contrapartida, **sinais de desengajamento** (meses de inatividade, redução recente de uso) e possivelmente **sintomas de insatisfação** (muitos contatos de suporte) estão fortemente associados ao churn. Nossos melhores modelos (XGBoost/LightGBM) capturaram esses padrões e alcançaram AUC próximas a 0,95 – indicando capacidade preditiva de alto nível. Em termos simples, conseguimos priorizar corretamente cerca de 4 a cada 5 clientes que de fato churnariam, o que pode guiar intervenções eficientes.

Insights de Negócio: O perfil típico de clientes em risco de churn inclui aqueles com **baixo uso do cartão** (especialmente se esse uso caiu nos últimos meses), com **poucos vínculos com o banco além do cartão**, muitas vezes em faixas de renda extremas (muito baixa ou muito alta) e possivelmente já insatisfeitos (manifestado por contatos frequentes ou reclamações). Por outro lado, clientes **altamente engajados** – que usam ativamente o cartão, contratam diversos produtos e têm relacionamento de longa data – demonstraram fidelidade muito maior. Esses achados sugerem estratégias duais: aumentar o engajamento dos clientes de baixo uso (incentivando transações regulares, por exemplo via programas de recompensa) e aumentar o relacionamento integrado (oferecer produtos complementares) tende a reduzir churn. Além disso, monitorar proativamente métricas como **inatividade** e **queda de utilização** permite **intervir antes que o cliente cancele** – por exemplo, se um

cliente passou 3 meses sem usar ou reduziu 50% do gasto, o banco pode enviar ofertas personalizadas, descontos na anuidade ou entender se houve algum problema com o cartão.

Recomendações Futuras: Com base nesta análise, recomendamos diversas ações e melhorias futuras:

- **1. Ações Pró-Ativas de Retenção:** Implantar um **sistema de alerta** baseado nos principais preditores (como quedas de uso, alta inatividade, poucos produtos) para o time de relacionamento entrar em contato com clientes de risco. Por exemplo, clientes marcados como churners prováveis pelo modelo (top N% do score) poderiam receber ofertas de upgrade, isenção de tarifa ou convites para conversar sobre insatisfação. Intervenções tempestivas aumentam a chance de reter o cliente antes que ele tome a decisão final.
- **2. Enriquecimento de Dados e Novas Features:** Incluir novas fontes de informação pode melhorar ainda mais o poder explicativo. Por exemplo, dados de **satisfação do cliente** (NPS, reclamações registradas), **dados transacionais fora do cartão** (uso de conta corrente, investimento, se houver), ou mesmo dados macro (como cenário econômico regional) podem refinar a identificação de churn. Uma feature potencial é o *atraso em pagamentos* – clientes que começaram a atrasar faturas do cartão podem estar insatisfeitos ou enfrentando dificuldades, o que às vezes precede churn. Outro aprimoramento seria incorporar **dados comportamentais digitais**: frequência de login no app, uso de internet banking, etc. Clientes que deixam de acessar podem ser um sinal adicional de desengajamento.
- **3. Aprimoramento dos Modelos (Tuning e Ensemble):** Embora os modelos já tenham alta performance, podemos explorá-los mais. Recomenda-se realizar um **tuning de hiperparâmetros** mais extenso, especialmente para XGBoost e LightGBM – usando técnicas como busca em grid ou Bayesian optimization para encontrar combinações ótimas de profundidade de árvore, learning rate, número de estimadores, etc. Isso pode ganhar alguns pontos de AUC ou reduzir erros. Além disso, poderíamos testar **outros algoritmos de ponta**, como **CatBoost** (outra variação de boosting que lida bem com categóricas) e até redes neurais simples, para ver se atingem performance semelhante. Uma abordagem interessante é criar um **ensemble/stacking** combinando os modelos (por exemplo, média das probabilidades ou um meta-modelo) para tentar aproveitar as forças de cada – embora aqui como todos já estão próximos do ótimo, o ganho pode ser marginal. Importante também é revisar periodicamente o modelo treinado com novos dados para garantir que ele não perca acurácia caso o comportamento dos clientes mude (treinar em dados mais recentes, se disponíveis).
- **4. Explainability (Interpretabilidade Individual):** Implementar técnicas de **XAI (eXplainable AI)** para apoiar a tomada de decisão. Por exemplo, usar **SHAP values** ou **LIME** nos modelos de árvore para explicar, para um cliente específico, quais fatores contribuíram para seu score de churn. Isso seria extremamente útil para as equipes de atendimento: ao contatar um cliente sinalizado, ter em mãos os motivos – “ex: Sr. João, notamos que você reduziu o uso do cartão nos últimos meses e não está aproveitando seu limite; gostaríamos de entender se há algum problema e lhe oferecer...” – torna a abordagem mais assertiva e personalizada. Além disso, do ponto de vista de confiança no modelo, explicações claras ajudam a validar que o modelo não está tomando decisões enviesadas indevidamente (por exemplo, se estivesse usando gênero ou renda de forma não desejada; pelas análises, não parece o caso – uso e relacionamento são os drivers principais, o que é intuitivamente correto).
- **5. Segmentação e Estratégias Diferenciadas:** Utilizar a **clusterização de clientes** em conjunto com as previsões de churn para estratégias segmentadas. Por exemplo, nossos clusters podem ter identificado grupos como “Jovens de baixa renda, pouco uso” vs “Clientes maduros

endividados” vs “Usuários premium multifidelizados”. Para cada segmento, as ações de retenção podem ser distintas: uns respondem melhor a benefícios financeiros (cashback, desconto anuidade), outros a melhorias no atendimento ou upgrade de produto. Com os clusters, o banco pode **priorizar recursos** – focando maior esforço de retenção nos segmentos de alto valor (alto LTV) com risco de churn, enquanto talvez aceite naturalmente que segmentos de baixíssimo valor churnem, realocando esforço onde importa.

- **6. Monitoramento Contínuo e KPI:** Implementar o modelo de churn como parte do **dashboard de KPIs do banco**, monitorando a taxa de churn real mensalmente e as projeções do modelo. Isso ajuda a medir o impacto das ações de retenção: por exemplo, se a taxa de churn observada cair de 16% para digamos 12% após iniciativas guiadas pelo modelo, quantificar o valor financeiro economizado (usando LTV dos retidos) para justificar investimentos em programas de retenção. Além disso, monitorar se o modelo mantém calibração – se de repente muitos clientes previstos como baixo risco começarem a churnar, pode ser um sinal de mudança de comportamento ou necessidade de retreinamento.

Em conclusão, o estudo de caso evidenciou que a análise de churn, apoiada por **técnicas de Data Science**, pode gerar insights práticos e modelos preditivos de alta acurácia. Identificamos os principais *drivers* da saída de clientes e mostramos que é possível antecipar com ~95% de assertividade quem está propenso a cancelar. As recomendações propostas, se implementadas, podem **aumentar a retenção de clientes e, consequentemente, a receita** – por exemplo, estima-se que melhorar a retenção em alguns pontos percentuais pode impulsionar os lucros em dezenas de por cento ¹. É fundamental, porém, tratar churn não apenas reativamente (quando o cliente já decidiu sair), mas preventivamente: usar os resultados deste modelo para **alimentar estratégias de CRM proativas**, focadas em manter os clientes engajados e satisfeitos ao longo de todo o ciclo de vida. Desta forma, o banco poderá reduzir significativamente a taxa de churn (atualmente ~16% ³) e fidelizar mais clientes, transformando análise preditiva em vantagem competitiva no mercado financeiro.

Referências dos Dados e Estudos Citados: Este relatório baseou-se no dataset público *BankChurners* (Clientes de Cartão de Crédito) e em diversos recursos para embasar as conclusões, incluindo artigos e repositórios que analisaram problemas similares de churn bancário ^{37 22 36 1}, garantindo que as metodologias e inferências estejam alinhadas com as melhores práticas e conhecimentos atuais da área.

¹ The Financial Impact Of Customer Churn (Direct vs. Indirect Costs)

<https://blog.miaarec.com/the-financial-impact-of-customer-churn-direct-vs.-indirect-costs>

² Banking Customer Retention Statistics 2025: Rates, Digital Impact

<https://coinlaw.io/banking-customer-retention-statistics/>

^{3 7 8 10 11 12 13 14 15 18 19 21 25} Bank Churn Analysis: Understanding Customer Attrition.
| by Cendikia Ishmatuka | Medium

<https://cendikiaishmatuka.medium.com/bank-churn-analysis-understanding-customer-attrition-28e8cea73a86>

^{4 5 6 17 20 27 28 29 30 31 33 37} GitHub - skhettri/BankChurnersEDA: Exploratory Data
Analysis of Potential Bank Churners

<https://github.com/skhettri/BankChurnersEDA>

^{9 16 23 26 34} Customer Churn Prediction Exercise

https://rstudio-pubs-static.s3.amazonaws.com/698111_b834c4bc622e42b1affd99fd82405472.html

22 35 deanfrancispress.com

<https://www.deanfrancispress.com/index.php/fe/article/download/273/FE000522.pdf/1219>

24 Credit Card Customer Churn Prediction - Kaggle

<https://www.kaggle.com/code/kaushikmajumder/credit-card-customer-churn-prediction>

32 Customer Churn Prediction and Segmentation for Credit Card ...

https://github.com/biancaportela/customer_churn_segmentation

36 Explainable AI based LightGBM prediction model to predict default ...

<https://www.sciencedirect.com/science/article/pii/S2667305325000407>