

# An investigation of challenges in the machine learning lifecycle and the importance of MLOps: A survey

Bruno Faustino Amorim  
Universidade Tecnológica Federal do Paraná  
Dois Vizinhos, Paraná, Brasil  
bamorim@alunos.utfpr.edu.br

Lincoln M. Costa  
Universidade Federal do Rio de Janeiro  
Rio de Janeiro, Rio de Janeiro, Brasil  
costa@cos.ufrj.br

Alinne C. Corrêa Souza  
Universidade Tecnológica Federal do Paraná  
Dois Vizinhos, Paraná, Brasil  
alinnesouza@utfpr.edu.br

Francisco Carlos M. Souza  
Universidade Tecnológica Federal do Paraná  
Dois Vizinhos, Paraná, Brasil  
franciscosouza@utfpr.edu.br

## ABSTRACT

Quite recently, considerable attention has been paid to developing artificial intelligence and data science areas. This has been driven by scientific advances and the growing number of software and services that are popularizing machine learning techniques and algorithms and driving people with less knowledge in areas such as statistics and mathematics to create their predictive models. As a result, the machine learning field is no longer only scientific and has aroused the interest of companies from different domains. These events led to the emergence of multiple tools such as Scikit-Learn, Tensorflow, Keras, Pycaret, and a vast number of cloud-based machine learning services that provide an acceleration in the development of predictive models at speeds never seen. However, many challenges remain in operationalizing and maintaining machine learning-centered products, making many business initiatives frustrated. In this scenario, practical experience shows that machine learning is only a slice of a more extensive set of practices and technologies necessary to build solutions in this area. In this paper, the main goal is to identify the challenges currently faced by data scientists in developing Machine Learning-centric products and how Machine Learning Operations can support overcoming them. For this purpose, a survey was conducted that collected answers from 66 Brazilian professionals in data science. From the challenges identified, the importance of Machine Learning Operations practices as an integrated part of the Machine Learning lifecycle was explored. Finally, this work contributes to filling the gap in Machine Learning Operations in daily activities involving data science and advancing this research field in Brazil.

## CCS CONCEPTS

• **Computing methodologies** → *Machine learning*.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*Computer on the Beach '23*, 30 de março a 01 de abril, 2023, Florianópolis, SC

© 2023 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## KEYWORDS

Machine Learning, MLOps

### ACM Reference Format:

Bruno Faustino Amorim, Alinne C. Corrêa Souza, Lincoln M. Costa, and Francisco Carlos M. Souza. 2023. An investigation of challenges in the machine learning lifecycle and the importance of MLOps: A survey. In *Proceedings of Computer on the Beach (Computer on the Beach '23)*. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUÇÃO

O crescente aumento no volume dos dados, somado aos avanços acadêmicos, direcionou empresas de diversos setores a investir em iniciativas envolvendo a resolução de problemas ou oferta de serviços utilizando Aprendizado de Máquina (AM) (do inglês, *Machine Learning - ML*). Além disso, a busca constante por inovação, resultado de mercados altamente competitivos, está tornando o uso de AM como uma peça fundamental. As iniciativas envolvendo AM são variadas, como na detecção de fraude [19], prevenção de doenças [20], bem como no controle de catástrofes naturais [13].

Apesar dos avanços dos algoritmos de AM, o sucesso de softwares baseados em AM não dependem somente de modelos preditivos altamente precisos, fazendo com que muitas expectativas não sejam alcançadas [8]. A maioria dos esforços acadêmicos, até o momento, concentra-se no desenvolvimento de técnicas e algoritmos de AM. Todavia, na prática, as empresas muitas vezes optam pela utilização de métodos mais simples para o desenvolvimento dos modelos, motivados pela complexidade de implementação, manutenção e por critérios de interpretabilidade [14], [6] [5]. O sucesso na criação de soluções baseadas em AM está intimamente relacionada em aplicar esses modelos como parte de ecossistemas de software.

As atividades do ciclo de vida de AM são de natureza exploratória, iterativa e requerem muita experimentação. Durante esse ciclo, existe um árduo trabalho entre a primeira fase de especificação dos requisitos e a disponibilização dos modelos preditivos para os usuários ou sistemas finais. Grande parte desse processo é centrado na manipulação de dados, já que esses dados podem sofrer diversos tipos de variações, seja por alterações nas aplicações que os geram, problemas de qualidade, ou mudança no comportamento de determinados eventos [2].

Fundamentado nisso, esse artigo visa explorar os desafios enfrentados pelos profissionais de ciência de dados durante o ciclo de

vida de AM, contemplando desde o desenvolvimento até a implantação de modelos em produção. A partir desse estudo, foi abordado como o MLOps pode otimizar o ciclo de vida de AM por meio de práticas, automação de processos e estrutura. Neste contexto, é importante destacar que práticas de engenharia de software podem ser aplicadas para o desenvolvimento, implantação, manutenção e monitoramento dos modelos como parte das aplicações empresariais, pois acelera as entregas de modelos, e consequentemente, otimiza o trabalho das equipes de ciência de dados.

Este tema de pesquisa tem sido investigado em trabalhos como [16], [18] e [22]. A literatura existente estuda a aplicação de MLOps nas atividades de AM com a finalidade de amenizar os desafios diários enfrentados pelos profissionais de ciência de dados [16]. No entanto, existe uma lacuna na abrangência dos aspectos práticos observados, já que na indústria a importância do MLOps ainda é vagamente explorada, fazendo com que muitas empresas enfrentem problemas em diversas fases do ciclo de vida de AM [12].

As principais contribuições deste estudo podem ser resumidas como: *i*) um *survey* para identificar possíveis desafios no ciclo de vida de AM; *ii*) síntese dos desafios identificados por meio de agrupamentos de atividades do ciclo de vida; e *iii*) a percepção sobre a área de MLOps na visão de profissionais da ciência de dados. Como principais observações, foi identificado que 90% dos participantes já realizam a implantação de modelos de AM em ambiente produtivo. Contudo, nota-se que os profissionais de ciência de dados estão empenhando tempo e esforço considerável com atividades além do desenvolvimento de modelos preditivos. Além disso, as operações de AM também carecem de profissionais de Engenharia.

O restante do artigo está estruturado da seguinte forma: a Seção 2 apresenta os conceitos relacionados à Aprendizagem de Máquina. Na Seção 3 são discutidos os trabalhos relacionados. Na Seção 4 é detalhado o processo de condução do *survey*. Na Seção 5 são discutidos os resultados alcançados, bem como a importância do MLOps. Por fim, na Seção 5 encerra este trabalho apresentando as conclusões.

## 2 ASPECTOS CONCEITUAIS

Nesta Seção são apresentados os conceitos relacionados ao Ciclo de vida de AM, bem como MLOps.

### 2.1 Ciclo de vida de Aprendizagem de Máquina

Com diferentes formas de representação, encontrou-se na literatura várias propostas de fluxos para o ciclo de vida de AM [2] [17] [3]. A construção de modelos de AM é um processo constituído por diversas etapas, onde durante a busca do modelo generalizável com melhor desempenho, várias dessas etapas são reiteradas, tornando o processo de desenvolvimento altamente iterativo e exploratório.

Como embasamento teórico foi utilizado um estudo de caso da Microsoft [2], que estruturou o ciclo de vida de AM em nove atividades, conforme pode ser visto na Figura 1. Para o nosso atual contexto, essas atividades foram agrupadas em dois conjuntos: 1) Desenvolvimento: contendo todas as etapas referentes ao desenvolvimento do modelo, desde o mapeamento de requisitos até a avaliação do modelo. 2) Implantação: contendo as etapas de implantação e monitoramento do modelo.

A primeira etapa é a análise dos requisitos do modelo, onde são definidos quais recursos podem ser implementados com AM, recursos úteis para um determinado produto existente ou novo e tipos de modelos mais apropriados para o problema [2]. Na coleta de dados os cientistas obtêm os dados necessários para o trabalho em questão. Nessa fase, surgem os desafios, principalmente relacionados as dimensões de variedade e volumetria dos dados [17]. Durante a preparação e transformação, é realizada a limpeza e transformação dos dados brutos para análise <sup>1</sup>. A fase de engenharia de recursos é utilizada para selecionar e transformar variáveis. Isso envolve a criação, transformação, extração e seleção de recursos <sup>2</sup>. Na fase de treinamento do modelo, os modelos escolhidos são treinados e ajustados sobre os dados transformados [2]. A avaliação do modelo é responsável por determinar se o desempenho do modelo é satisfatório para atingir as metas definidas pelo negócio <sup>3</sup>. Em seguida, o modelo pode ser implantado e disponibilizado para os sistemas ou usuários finais [3].

### 2.2 Operações de Aprendizagem de Máquina (MLOps)

MLOps tem por objetivo unificar o desenvolvimento de sistemas de AM (ML) e a operação de sistemas de AM (Ops) <sup>4</sup>. Apesar dos cientistas de dados possuírem ferramentas e habilidades para implementação e treino de modelos preditivos, a implantação e sustentação desses modelos em ambiente produtivo requer outros conhecimentos técnicos fora do domínio de ciência de dados. O maior desafio está na criação de sistemas de AM operantes continuamente em produção [21] <sup>5</sup>.

Apesar do rápido crescimento no desenvolvimento dos modelos, as empresas possuem muitas dificuldades para operacionalizá-los, fazendo com que grande parte dos modelos criados nunca cheguem em produção [8] [23]. Como a criação de produtos ou serviços baseados em AM envolve outros componentes além dos modelos preditivos, cientistas de dados não são capazes de sozinhos realizar todo o trabalho contido em todo o ciclo de vida de AM, sendo necessária uma equipe multidisciplinar e conhecimentos especializados em outras áreas de conhecimento como Engenharia de Dados e Engenharia de Software.

Além disso, a construção de sistemas de AM apresenta um conjunto de particularidades que não são tradicionalmente vistas em engenharia de software. Contrariamente ao desenvolvimento de software tradicional, as aplicações de AM acrescentam bases de dados que mudam constantemente e afetam a maneira de desenvolver e manter esses tipos de sistemas.

## 3 TRABALHOS RELACIONADOS

Embora a literatura sobre MLOps seja ampla, alguns *surveys* fornecem uma visão das dificuldades durante os estágios do ciclo de vida de AM. O trabalho de Mäkinen et al. [16] aborda uma pesquisa

<sup>1</sup> <https://www.oracle.com/a/ocom/docs/data-science-lifecycle-ebook.pdf>

<sup>2</sup> [https://docs.aws.amazon.com/pt\\_br/wellarchitected/latest/machine-learning-lens/feature-engineering.html](https://docs.aws.amazon.com/pt_br/wellarchitected/latest/machine-learning-lens/feature-engineering.html)

<sup>3</sup> [https://docs.aws.amazon.com/pt\\_br/wellarchitected/latest/machine-learning-lens/feature-engineering.html](https://docs.aws.amazon.com/pt_br/wellarchitected/latest/machine-learning-lens/feature-engineering.html)

<sup>4</sup> <https://cloud.google.com/architecture/mlops-continuous-delivery-and-automation-pipelines-in-machine-learning>

<sup>5</sup> <https://cloud.google.com/architecture/mlops-continuous-delivery-and-automation-pipelines-in-machine-learning>

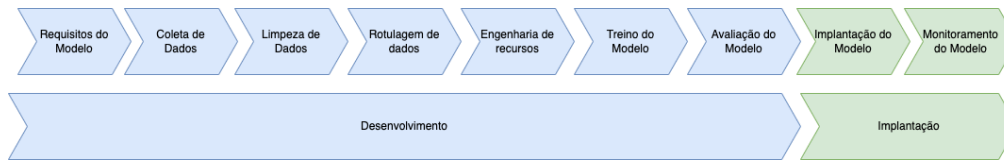


Figure 1: Ciclo de vida de Aprendizado de Máquina.

com 331 cientistas de dados, em que foi identificada a importância do MLOps no contexto das atividades diárias desses profissionais. Neste trabalho também foram levantados aspectos profissionais e empresariais dos participantes, além de abranger um número maior de elementos focados em todo o ciclo de vida de AM. Por fim, os resultados apontam que os cientistas de dados estão expandindo suas atuações para tarefas relacionadas a infraestrutura e implantação de modelos de AM.

Em [18] foi conduzida uma pesquisa que descreve obstáculos que surgem durante as fases do ciclo de vida de implantação de modelo de AM, considerando os estágios de gerenciamento de dados, aprendizado de modelo, verificação de modelo e implantação de modelo. Os principais obstáculos identificados estão relacionados à reutilização de código e anti-padrões de Engenharia de Software. É importante destacar que foram explorados aspectos como ética e segurança da informação, que devem ser considerados durante o desenvolvimento dos modelos.

No estudo de Souza [22] foi investigado a complexidade na adoção de ferramentas e práticas de MLOps. Como principal contribuição, por meio de experimentos, foram apresentadas como algumas ferramentas de código aberto podem ser aplicadas na implantação do MLOps.

O estudo de Kim et al. [11] apresenta uma pesquisa com 793 cientistas de dados da Microsoft investigando os problemas enfrentados, práticas recomendadas durante a construção de modelos de AM e o atual contexto de trabalho dos cientistas de dados. Os resultados apontam uma tendência desses profissionais atuarem em funções diretamente inseridas nas equipes de software. Além disso, o trabalho também mostra que funcionários contratados para outras posições na empresa assumiram tarefas de ciência de dados como parte do seu trabalho. O estudo é limitado ao contexto corporativo da Microsoft, logo algumas conclusões podem não ser representativas para empresas de outros setores e de diferentes outros portes.

Por fim, o trabalho de Jain et al. [9] destaca a importância do controle da qualidade dos dados em tarefas de AM, por meio de métricas de qualidade de dados como desequilíbrio de classes e homogeneidade dos dados que podem ser utilizadas e monitoradas como parte do ciclo de vida de AM.

Apesar da abrangência dos trabalhos apresentados em contemplar as dificuldades durante todos os estágios do ciclo de vida de AM, ainda existe uma lacuna na cobertura dos aspectos não técnicos. Entre elas, a percepção de valor vista pelas áreas de negócio, bem como a falta de informações sobre as composições e estruturas dos times dos profissionais de ciência de dados das empresas analisadas. Neste contexto, o presente artigo tem como objetivo preencher essa lacuna e contribuir para a literatura existente, trazendo uma visão holística que investiga os atuais desafios do ciclo de vida de AM. Portanto, foi conduzido um estudo exploratório referente às

adversidades práticas enfrentadas pelos profissionais atuantes na área de ciência de dados envolvendo a manipulação de dados.

## 4 AVALIAÇÃO EXPERIMENTAL

Esta seção detalha o *survey* conduzido para identificar os desafios enfrentados atualmente por profissionais que atuam na área de ciência de dados no desenvolvimento de produtos centrados em AM e como o MLOps pode ajudar a minimizá-los. O *survey* foi planejado seguiu o processo proposto por Kasunic [10] para o design efetivo de *surveys*. Além disso, foram utilizados direcionamentos descritos por Kitchenham e Pfleger [1].

### 4.1 Identificação dos objetivos de investigação e Questões de Pesquisa

O objetivo do *survey* consiste na identificação dos desafios no desenvolvimento de produtos utilizando AM e como a aplicação de MLOps pode minimizar esses desafios. Neste contexto, os objetivos do *survey* foram especificados segundo o modelo *Goal-Question-Metric (GQM)* proposto por Basili e Weiss [4], conforme pode ser visto a seguir:

*"Analisar desenvolvimento de produtos centrados em AM com o propósito de identificar os desafios enfrentados e como o MLOps pode minimizar esses desafios do ponto de vista de profissionais que atuam na área de ciência de dados no contexto de ciclo de vida de AM."*

A partir do objetivo, as seguintes Questões de Pesquisas (QPs) foram identificadas:

**QP<sub>1</sub>: Qual perfil dos profissionais atuantes na área de ciência de dados e da empresa na qual trabalham?** Nesta QP, buscou-se identificar informações relacionadas à formação, cargo, experiência no cargo ocupado, nível (júnior, pleno ou sênior) do cargo ocupado e o nível de conhecimento sobre MLOps. Além disso, foram coletadas informações relacionadas à empresa como o tipo (nacional ou multinacional), número de funcionários, ramo, como a empresa se encontra atualmente em relação ao desenvolvimento e execução de modelos e quais profissionais compõem os times de dados.

**QP<sub>2</sub>: Quais são as principais atividades desempenhadas pelos profissionais atuantes na área de ciência de dados nas empresas?** Esta QP busca identificar as atividades realizadas com mais frequência durante o desenvolvimento de modelos de AM, bem como a frequência dos tipos de dados manipulados. Além disso, investigou-se qual ambiente, linguagens e ferramentas são utilizadas para auxiliar o desenvolvimento de modelos de AM.

**QP<sub>3</sub>: Quais desafios são enfrentados durante o ciclo de vida de AM?** Esta QP visa identificar os principais desafios relacionados à manipulação dos dados, ao desenvolvimento dos modelos, a não

implantação de modelos em produção, bem como a implantação e o monitoramento. Além disso, nesta QP espera-se apresentar como as práticas de MLOps podem minimizar os desafios identificados no ciclo de vida.

## 4.2 Identificação do público-alvo e planejamento de amostragem

O público-alvo definido para participar do *survey* é composto por profissionais que atuam na área de ciência de dados que estejam trabalhando em empresas de diferentes ramos distribuídos geograficamente pelo Brasil. Após a definição do público-alvo, foi planejado o processo de obtenção de amostras. Nesta etapa considerou-se localizar profissionais que trabalham em qualquer nicho de negócio. Assim, a identificação desses profissionais foi realizada por meio de contato via e-mail e rede social.

## 4.3 Planejamento e escrita do questionário

O detalhamento do protocolo desenvolvido para a condução do *survey* pode ser visualizado no link <https://shre.ink/mNKt>. Além do protocolo, o questionário aplicado junto aos participantes pode ser acessado em: <https://bityli.com/fBqFghav>.

As 25 questões foram estruturadas em cinco seções: (i) contém uma apresentação do *survey* em que é descrito o objetivo do mesmo e o público-alvo; (ii) coleta informações que caracterizam os participantes; (iii) identifica informações relacionadas as atividades conduzidas pelos participantes; (iv) coleta informações a respeito dos desafios enfrentados pelos cientistas de dados durante o ciclo de vida de AM; e (v) obtém formas de como MLOps podem otimizar os processos de AM. O *survey* pode ser acessado em: <https://drive.google.com/file/d/1q8Ob2LIPKCY9Cf0lrbf2Nzm37iiOd7wF>.

## 4.4 Execução do *survey* piloto

Segundo Kasunic [10], é fundamental a condução de um *survey* piloto, pois é possível detectar possíveis problemas existentes no mesmo, verificar se as perguntas são compreensíveis, se as perguntas certas foram feitas para atingir o objetivo, e quanto tempo os participantes levam para completarem o questionário. Para avaliar o *survey* foram aplicadas quatro questões abertas propostas por Hauck et. al [7], sendo: (1) O questionário contém tudo que é esperado para contemplar o seu objetivo?; (2) O questionário contém quaisquer informações não desejáveis ou desnecessárias ao contexto e objetivo da pesquisa?; (3) Você conseguiu compreender adequadamente as perguntas?; e (4) Existe algum erro ou inconsistência no questionário?. Além destas questões também foram consideradas as opiniões de especialistas [15].

Nesse contexto, um grupo de quatro participantes distribuídos entre formados e não formados, estagiando e trabalhando, foram convidados por e-mail para participar do teste piloto. Esses participantes, foram selecionados pelos critérios de disponibilidade e proximidade, participaram do estudo piloto respondendo às questões propostas por Hauck et. al [7] e enviaram *feedbacks* sobre o *survey*.

Um grupo de profissionais foi convidado por e-mail para participar do estudo piloto. Esses profissionais foram escolhidos pelos critérios de disponibilidade e proximidade como o grupo onde esta pesquisa foi realizada. A avaliação dos profissionais foi positiva, com sugestões para: (i) reduzir o número de questões; (ii) deixar

algumas perguntas como não obrigatórias; e (iii) incluir pelo menos uma pergunta aberta.

## 4.5 Coleta e Análise dos Dados

O questionário foi divulgado para os profissionais da área de ciência de dados via e-mail e a rede social LinkedIn<sup>6</sup> e ficou disponível no período de 29 de Junho a 5 de Agosto de 2022. Após a coleta dos dados, os mesmos foram analisados com o objetivo de verificar a consistência e completude das respostas[1]. Em seguida, foram realizadas análises quantitativas e qualitativas conforme o tipo de perguntas utilizadas na pesquisa. Para a interpretação dos dados quantitativos foi utilizada estatística descritiva; e análise de discurso e visualização de dados para a análise qualitativa.

## 5 RESULTADOS

O *survey* obteve a contribuição de 66 profissionais atuantes na área de ciência de dados. Os resultados serão apresentados de acordo com as QPs, conforme apresentadas na Seção 4.

### 5.1 Perfil dos profissionais atuantes na área de ciência de dados e das empresas (QP<sub>1</sub>)

A pesquisa coletou respostas de profissionais de diferentes organizações, sendo 68% (45/66) empresas nacionais brasileiras e 31% (21/66) multinacionais. Com base na quantidade de funcionários, de acordo com o Sebrae<sup>7</sup>, a maioria dessas empresas são de grande porte, sendo que 27% (18/66) possuem de 100 a 500 funcionários, e 54% (36/66) possuem mais de 500 funcionários. Dessas empresas, 19% (13/66) foram classificadas como empresas de Software/Internet e 24% (16/66) de Serviços Financeiros.

Sobre os profissionais participantes, 45% (30/66) possuem ensino superior completo, 27% (18/66) são mestres, 16% (11/66) são especialistas, 7% (5/66) são doutores e 3% (2/66) possuem ensino superior incompleto.

Uma classificação de quatro níveis foi criada para identificar o nível de maturidade das empresas na utilização de Aprendizado de Máquina nas operações. Para isso, realizou-se a seguinte pergunta: “Qual das frases abaixo mais se aproxima do momento atual da sua empresa?”.

- **Nível 1:** a empresa está desenvolvendo seus primeiros modelos, porém ainda sem gerar valor para as frentes de negócio.
- **Nível 2:** as áreas de negócio já consomem os modelos, porém são executados localmente e não há monitoramento ou automações.
- **Nível 3:** a empresa executa seus modelos em ambiente produtivo em nuvem, com consumos via API ou *Batch*, porém ainda sem processos automatizados.
- **Nível 4:** a empresa executa seus modelos em ambiente produtivo em nuvem, e com processos como monitoramento, implantação e retreino automatizados.

A criação dessa classificação, em parte, foi inspirada em um trabalho realizado pelo Google<sup>8</sup>.

<sup>6</sup><https://br.linkedin.com/>

<sup>7</sup>[https://www.sebrae.com.br/Sebrae/Portal%20Sebrae/UFs/SP/Pesquisas/MPE\\_conceito\\_empregados.pdf](https://www.sebrae.com.br/Sebrae/Portal%20Sebrae/UFs/SP/Pesquisas/MPE_conceito_empregados.pdf)

<sup>8</sup><https://cloud.google.com/architecture/ml-ops-continuous-delivery-and-automation-pipelines-in-machine-learning>

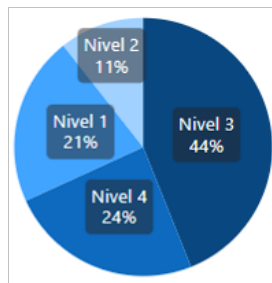


Figure 2: Níveis de Maturidade

A maioria das empresas está no nível 3, sendo empresas que já possuem modelos de AM em produção, porém ainda sem processos automatizados. Ademais, um percentual considerável de empresas está no nível 1, desenvolvendo seus primeiros modelos de AM, porém ainda sem a geração de valor para o negócio (Figura 2).

Os times de dados das empresas são compostos por diferentes perfis, observando que em sua maioria são 95% (63/66) formados por cientistas de dados, 83% (55/66) engenheiros de dados e 68% (45/66) analistas de BI (Figura 3). Apesar da maior parte das empresas serem de médio e grande porte, 37% (25/66) não possuem equipes dedicadas para a implantação de modelos de AM em produção.

## 5.2 Principais atividades desempenhadas pelos profissionais atuantes na área de ciência de dados nas empresas ( $QP_2$ )

Para atingir o objetivo dessa seção, foram elencadas perguntas para identificar os vários aspectos que contemplam todas as atividades de desenvolvimento no ciclo de vida de AM, realizando perguntas relacionadas a diferentes domínios, desde o ambiente utilizado para o desenvolvimento dos modelos, incluindo ferramentas e tecnologias, até as atividades mais frequentes e maiores problemas e dificuldades enfrentados durante essa fase.

O primeiro aspecto a ser considerado é que apesar do aumento na variedade das informações, causado com o advento do *Big Data*, nota-se que os times de ciência de dados raramente trabalham com dados não estruturados de imagem ou áudio. Suas atuações concentram-se praticamente sempre com dados estruturados e semi-estruturados (Figura 4).

Sobre as atividades mais frequentes, a "limpeza e transformação dos dados" é a atividade apontada como mais recorrente pelos participantes. Isso não é surpresa, já que essa atividade é conhecida uma das mais comuns e trabalhosas nessa área. Houve também um alto índice de respostas apontando para atividades de coleta de dados.

A implantação de novos modelos em produção mostra-se como atividade de alta frequência. Apesar disso, dos 57 participantes que afirmam que suas empresas possuem modelos em produção, 45% (26/57) deles responderam que suas empresas tem até no máximo 5 modelos em produção e 19% (11/57) deles tem de 6 a 10 modelos em produção. Dada a quantidade relativamente baixa de modelos em produção, suspeita-se que os poucos modelos em operação trazem uma sobrecarga operacional considerável para os times envolvidos.

Nesse cenário, o MLOps pode trazer ganhos significativos automatizando tarefas como implantação, retreino, e monitoramento dos modelos.

Outro ponto importante a se notar é que o desenvolvimento de novos modelos foi marcado por 57% (38/66) dos participantes como "Frequentemente" ou "Sempre". Relacionando essa informação a outros estudos como [23], suspeita-se que boa parte dos modelos desenvolvidos não chegam até a produção (Figura 5). Isso também pode estar relacionado diretamente a aspectos de dificuldades encontradas no contexto de dados, como será apresentado no decorrer do artigo.

A respeito das ferramentas mais utilizadas durante o desenvolvimento, 87% (58/66) dos participantes utilizam Scikit-learn<sup>9</sup>, 60% (40/66) XGBoost<sup>10</sup>. Apesar da baixa quantidade de manipulação de dados não estruturados, como áudio ou textos, nota-se um alto uso de ferramentas de Aprendizado Profundo como TensorFlow<sup>11</sup> e Keras<sup>12</sup>, onde 53% (35/66) utilizam TensorFlow e 48% (32/66) Keras. Com isso, questiona-se se não seria possível utilizar algoritmos e ferramentas menos complexas para atingir os mesmos resultados, talvez tornando o processo de modelagem mais simples. Sobre o Aprendizado de Máquina Automatizado (*AutoML*), vemos que 22% (15/66) dos participantes utilizam ferramentas como Pycaret<sup>13</sup>. Por fim, nota-se que a minoria, isto é, 30% (20/66) dos profissionais participantes, utiliza ferramentas e pacotes R<sup>14</sup>.

Quanto ao ambiente de trabalho, cerca de 63% (42/66) dos participantes desenvolvem seus modelos de AM em nuvem e 27% (18/66) em ambiente local ou *on-premise*, ou seja, diretamente na infraestrutura de servidores da empresa. Para o desenvolvimento dos modelos várias linguagens de programação podem ser aplicadas, sendo Python utilizada em 98% (65/66) e R 28% (19/66). Além disso, linguagens como Java, Scala ou *SaaS* raramente são utilizadas. As respostas indicam que boa parte dos times de ciência de dados dificilmente sofrem com a falta de ferramentas ou conhecimento técnico durante o desenvolvimento dos modelos.

## 5.3 Desafios enfrentados durante o ciclo de vida de AM ( $QP_3$ )

Os resultados demonstram que há uma grande dificuldade com aspectos de qualidade dos dados. Afirmarções de "dados necessários insuficiente ou incompletos", mostram-se em 81% (54/66) das vezes como uma das mais altas ocorrências quando somadas as 29 respostas marcadas como "Algumas vezes" e 25 respostas "Frequentemente". Nota-se também um alto índice de respostas afirmando dificuldades no acesso aos dados. O problema na qualidade dos dados se agrava quando é analisado o total marcado como "dados não confiáveis". Além disso, cerca de 53% (35/66) dos participantes relataram sofrer em algum nível com problemas no versionamento dos conjuntos de dados (*datasets*) (Figura 6).

Um alto número de participantes afirmou ter problemas com o rastreamento de experimentos e parâmetros durante a fase de desenvolvimento (Figura 7). Além disso, a fim de ampliar a abrangência

<sup>9</sup><https://scikit-learn.org/stable/>

<sup>10</sup><https://xgboost.readthedocs.io/en/stable/>

<sup>11</sup><https://www.tensorflow.org/>

<sup>12</sup><https://keras.io/>

<sup>13</sup><https://pycaret.org/>

<sup>14</sup><https://www.r-project.org/>

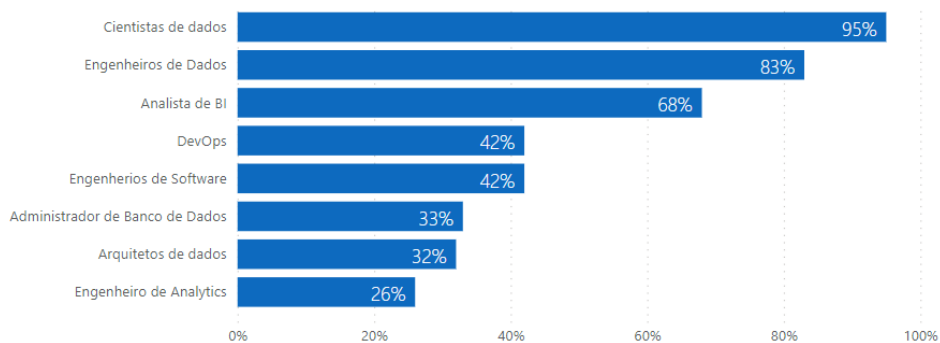


Figure 3: Composição dos times de dados.

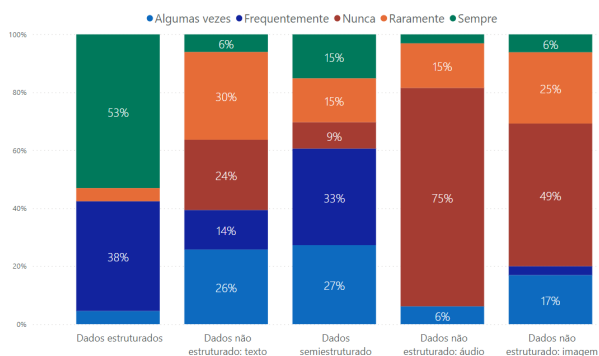


Figure 4: Tipos de dados frequentemente manipulados.

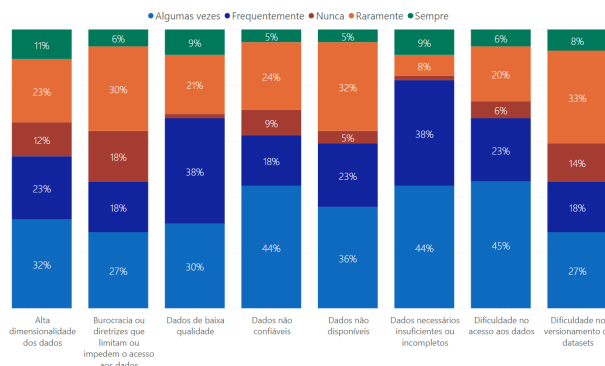


Figure 6: Problemas identificados na manipulação de dados.

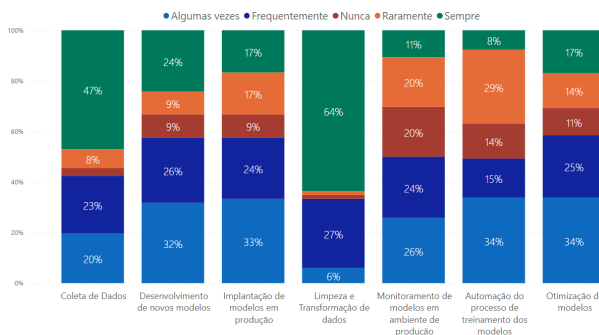


Figure 5: Atividades mais frequentes durante o desenvolvimento.

deste trabalho, cobrimos também aspectos não-técnicos. Grande parte das respostas demonstra que as áreas de negócio desconfiam do valor gerado com os modelos de AM (Figura 7). Com isso, a hipótese criada é que a falta de reprodutibilidade e falta de qualidade dos dados sejam um dos principais fatores que alimentam esse cenário de desconfiança. Outro importante aspecto é o alto nível de afirmações sobre expectativas irreais (Figura 7). A suspeita é que as áreas de negócio podem não compreender que a sua participação é importante para o processo de modelagem e definição de métricas. As expectativas também podem ser frustradas devido aos

problemas de ausência de dados necessários e dados de baixa qualidade. As organizações parecem não ter consciência de que modelos matemáticos são incapazes de gerar resultados positivos sem dados minimamente confiáveis. Muitas vezes as áreas de negócio criam a expectativa de que os modelos de AM são a saída para quase todos os problemas. Como resultado, isso pode gerar alta desconfiança.

Por fim, percebem-se poucas ocorrências apontando para a falta de recursos financeiros durante essa fase.

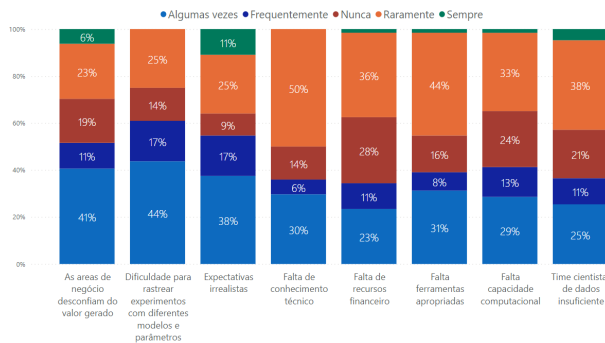


Figure 7: Dificuldades no desenvolvimento dos modelos de AM.

## 5.4 Implantação de modelos de AM

Neste ponto, analisou-se os 57 participantes que afirmaram realizar a implantação de modelos de AM em ambiente produtivo. Desses participantes, 43% (25/57) apontou não ter uma equipe dedicada para essa frente de trabalho.

Em geral, dois perfis profissionais ou mais atuam juntos durante a implantação, sendo que na maioria dos casos os cientistas e/ou engenheiros de dados são os responsáveis pelas implantações. Na minoria das implantações, observa-se a atuação em 23% (13/57) dos casos os Engenheiros de Aprendizado de Máquina, 16% (9/57) Engenheiros de Software e 16% (9/57) equipes de DevOps (Figura 8). Em poucos casos observam-se profissionais mais orientados a práticas de MLOps, como Engenheiros de Aprendizado de Máquina, atuando nas implantações. A estrutura de MLOps, pode ajudar os cientistas de dados a dedicar mais tempo em trabalhos de modelagem, do que implantando e orquestrando modelos de AM em produção.

A respeito das dificuldades enfrentadas durante a implantação e monitoramento dos modelos, 59% (34/57) dos participantes queixam-se de ter times insuficientes de engenharia (Engenharia de Dados, Engenharia de Aprendizado de Máquina ou DevOps) e 52% (30/57) dizem não ter conhecimento técnico suficiente para esse tipo de tarefa. Todavia, 43% (25/57) dos participantes dizem ter baixa ou nenhuma visibilidade sobre as métricas de modelos em produção. Outro fator importante é que boa parte, isto é 52% (30/57), se queixa da dificuldade em reproduzir na produção, os mesmos resultados obtidos durante a fase de desenvolvimento. Problemas com a visibilidade das métricas dos modelos em produção e a complexidade de retreino, mostram-se também como dificuldades enfrentadas pelos times de ciência de dados (Figura 9).

Analisou-se também o grupo de 13% (9/66) dos participantes que afirmaram que suas empresas não realizam implantação de modelo em produção. Nesse caso, as principais dificuldades relatadas foram a falta de conhecimento técnico e times de engenharia (Engenharia de Dados, Engenharia de Aprendizado de Máquina ou DevOps) não suficientes, seguido da falta de uma camada de consumo amigável como *dashboards* ou aplicações *web* e a falta de ferramentas adequadas.

Sobre o formato de consumo dos modelos, os resultados indicam que 47% (27/57) dos modelos em produção são entregues em formato de APIs, 14% (8/57) como aplicações para processamento em lote (*batch*), enquanto 10% (6/57) são entregues como *scripts* em formato *jupyter notebook*<sup>15</sup> para serem executados pelos usuários finais.

Observa-se que apesar da maioria das organizações já realizarem a implantação de modelos em produção, muitos desafios ainda são enfrentados. Com isso, há uma ampla gama de oportunidades para a aplicação de práticas de MLOps. Dificuldades envolvendo o rastreamento de experimentos, reprodutibilidade de resultados, versionamento de conjuntos de dados (*datasets*) e automações são questões intrínsecas ao MLOps. A aplicação de práticas de MLOps pode reduzir a sobrecarga operacional dos times de ciência de dados, automatizar tarefas e beneficiar a implantação e sustentação de modelos de AM em diferentes escalas.

Por fim, notou-se que as empresas normalmente possuem poucos modelos preditivos em produção. Levantou-se como hipótese que

isso se deve, em parte, por times de engenharia insuficientes e falta de conhecimento técnico. Além disso, a medida que novos modelos são colocados em produção, torna-se cada vez mais difícil mantê-los sem processos automatizados.

## 6 CONSIDERAÇÕES FINAIS

Neste artigo, os desafios enfrentados durante todo o ciclo de vida de AM foram explorados. Notou-se que os times de ciência de dados têm gasto tempo e esforço consideráveis com atividades além do desenvolvimento de modelos preditivos. Em contraste com outros trabalhos, os resultados deste artigo exibem os problemas enfrentados com questões não técnicas, falta de profissionais de engenharia e uma visão geral sobre o valor percebido pelos negócios.

Neste contexto, identificou-se que 86% (57/66) dos participantes já realizam a implantação de modelos em produção, porém a maioria ainda sem processos automatizados. Grande parte dessas implantações são realizadas pelos cientistas de dados em conjunto com times de engenharia. Como resultado, foi observado que há falta de recursos como times de engenharia e desafios técnicos e não-técnicos são percebidos durante a implantação.

De acordo com a pesquisa, 44% (26/66) das empresas encontram-se no Nível 3 de maturidade, tendo modelos de AM em produção, porém ainda sem processos automatizados. Uma parte considerável das empresas, isto é 21% (14/66), ainda estão no Nível 1, desenvolvendo seus primeiros modelos de AM, ainda não gerando valor para o negócio.

Sobre desafios não técnicos, nota-se dificuldades enfrentadas como a desconfiança das áreas de negócio sobre a geração de valor dos modelos de AM e expectativas não realistas das áreas de negócio. A implantação e sustentação de modelos em produção é dependente de outras áreas, profissionais e tecnologias que vão além dos processos e conhecimentos de modelagem. Nesse sentido, mostrou-se que práticas de MLOps podem ser utilizadas para superar desafios como reprodutibilidade de experimentos, automação de tarefas como implantação e retreino de modelos.

Os cientistas de dados dedicam boa parte do seu tempo trabalhando em tarefas como limpeza e transformação de dados. Nesse sentido, o MLOps tem um papel importante, podendo ajudar diretamente em todas as fases do ciclo de vida de AM. Em problemas apontados como os do caso da alta frequência de retreino, abordagens como treinamento contínuo podem ajudar na automação desse tipo de tarefa. A construção de *pipelines* de integração contínua e entrega contínua possibilitam a implantação de modelos de AM em produção de forma ágil e confiável.

Manter uma operação de modelos de AM requer uma estrutura de times de dados multifuncionais. Modelos de Aprendizado de Máquina apenas entregam seu real valor quando são mantidos de forma saudável em ambiente produtivo e disponível constantemente para consumo pelos usuários. Para que isso ocorra, apenas o desenvolvimento de modelos de AM não é o suficiente, sendo necessário uma estrutura capaz de manter tais modelos continuamente operantes.

Esta pesquisa contribui para o campo de AM, porque em primeiro lugar, aborda de forma ampla o ciclo de vida de Aprendizado de Máquina, discorrendo sobre questões como problemas relacionados a gestão de dados, ferramentas e tecnologias de desenvolvimento,

<sup>15</sup><https://jupyter.org/>

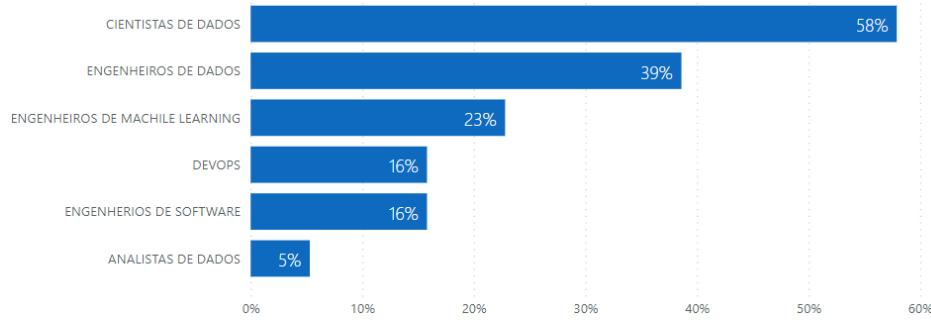


Figure 8: Profissionais que realizam a implantação dos modelos de AM.

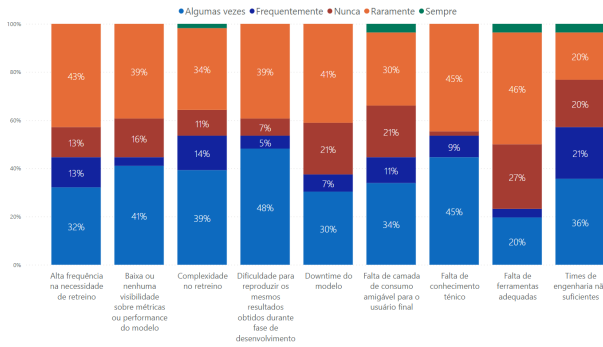


Figure 9: Desafios durante a fase de implantação.

profissionais responsáveis pela implantação, expectativas não realistas, falta de profissionais de engenharia e desconfiança das áreas de negócio na geração de valor, e em segundo lugar, explora como MLOps é importante para a otimização das operações diárias de Aprendizado de Máquina. Como limitações desse trabalho, a maior parte dos participantes são de empresas nacionais.

Como trabalho futuro, pretende-se: (i) desenvolver um estudo de caso sobre a experiência na implementação das práticas de MLOps citadas neste artigo; (ii) conduzir uma investigação aprofundada sobre os impactos da gestão de dados eficiente como parte do ciclo de vida de Aprendizado de Máquina.

REFERENCES

[1] Kitchenham B. A. and S. L. Pfleeger. 2008. Guide to Advanced Empirical Software Engineering. Springer London, London, Chapter Personal opinion surveys., 63–92.

[2] Saleema Amershi, Andrew Begel, Christian Bird, Robert DeLine, Harald Gall, Ece Kamar, Nachiappan Nagappan, Besmira Nushi, and Thomas Zimmermann. 2019. Software engineering for machine learning: A case study. In *2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*. IEEE, 291–300.

[3] Rob Ashmore, Radu Calinescu, and Colin Paterson. 2021. Assuring the Machine Learning Lifecycle: Desiderata, Methods, and Challenges. *ACM Comput. Surv.* 54, 5, Article 111 (may 2021), 39 pages. <https://doi.org/10.1145/3453444>

[4] V. Basili and D. Weiss. 1984. A methodology for collecting valid software engineering data. *IEEE Transactions on Software Engineering* 10, 6 (1984), 728–738.

[5] Or Biran and Courtenay Cotton. 2017. Explanation and justification in machine learning: A survey. In *IJCAI-17 workshop on explainable AI (XAI)*, Vol. 8. 8–13.

[6] Kristian Bondo Hansen. 2020. The virtue of simplicity: On machine learning models in algorithmic trading. *Big Data & Society* 7, 1 (2020), 2053951720926558.

[7] J. C. R. Hauck, C. G. Von Wangenheim, and A. Von Wangenheim. 2011. *Método de Aquisição de Conhecimento para Customização de Modelos de Capacidade/Maturidade de Processos de Software*. Relatório Técnico. INCoD, Florianópolis, SC.

[8] Lawrence E Hecht. 2019. Add it up: How long does a machine learning deployment take. *The New Stack. TheNewStack* 12 (2019).

[9] Abhinav Jain, Hima Patel, Lokesh Nagalapatti, Nitin Gupta, Sameep Mehta, Shanmukha Guttula, Shashank Mujumdar, Shazia Afzal, Ruhi Sharma Mittal, and Vitobha Munigala. 2020. Overview and Importance of Data Quality for Machine Learning Tasks. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (Virtual Event, CA, USA) (KDD '20)*. Association for Computing Machinery, New York, NY, USA, 3561–3562. <https://doi.org/10.1145/3394486.3406477>

[10] Mark Kasunic. 2005. *Designing an effective survey*. Pittsburgh, PA.: Carnegie Mellon University.

[11] Miryung Kim, Thomas Zimmermann, Robert DeLine, and Andrew Begel. 2017. Data scientists in software teams: State of the art and challenges. *IEEE Transactions on Software Engineering* 44, 11 (2017), 1024–1038.

[12] Dominik Kreuzberger, Niklas Köhl, and Sebastian Hirschl. 2022. Machine Learning Operations (MLOps): Overview, Definition, and Architecture. *arXiv preprint arXiv:2205.02302* (2022).

[13] Jake Lever and Rossella Arcucci. 2022. Towards Social Machine Learning for Natural Disasters. In *International Conference on Computational Science*. Springer, 756–769.

[14] Zachary C Lipton. 2018. The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery. *Queue* 16, 3 (2018), 31–57.

[15] M. S. Litwin. 1995. *How to Measure Survey Reliability and Validity*. SAGE Publication.

[16] Sasu Mäkinen, Henrik Skogström, Eero Laaksonen, and Tommi Mikkonen. 2021. Who needs MLOps: What data scientists seek to accomplish and how can MLOps help?. In *2021 IEEE/ACM 1st Workshop on AI Engineering-Software Engineering for AI (WAIN)*. IEEE, 109–112.

[17] Vincenzo Morabito. 2015. Big data governance. *Big data and analytics* (2015), 83–104.

[18] Andrei Paleyes, Raoul-Gabriel Urma, and Neil D Lawrence. 2020. Challenges in deploying machine learning: a survey of case studies. *ACM Computing Surveys (CSUR)* (2020).

[19] Johan Perols. 2011. Financial statement fraud detection: An analysis of statistical and machine learning algorithms. *Auditing: A Journal of Practice & Theory* 30, 2 (2011), 19–50.

[20] Ashish Sarraju, Andrew Ward, Sukyung Chung, Jiang Li, David Scheinker, and Fátima Rodríguez. 2021. Machine learning approaches improve risk stratification for secondary cardiovascular disease prevention in multiethnic patients. *Open heart* 8, 2 (2021), e001802.

[21] D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips, Dietmar Ebner, Vinay Chaudhary, Michael Young, Jean-Francois Crespo, and Dan Dennison. 2015. Hidden Technical Debt in Machine Learning Systems. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2 (Montreal, Canada) (NIPS'15)*. MIT Press, Cambridge, MA, USA, 2503–2511.

[22] João Vitor Ramos de Souza. 2021. Adoção de MLOps: desafios de gerenciar código, modelo e dados automaticamente. (2021).

[23] Kyle Wiggers. 2019. *IDC: For 1 in 4 companies, half of all AI projects fail, 2019*. Retrieved Oct 14, 2022 from <https://venturebeat.com/ai/idc-for-1-in-4-companies-half-of-all-ai-projects-fail/>