

Marcelo Veloso Maciel

**Um estudo de caso do uso de mineração de
dados e aprendizado de máquina no
aprimoramento de inspeções de estações radio
base**

Brasil

Marcelo Veloso Maciel

**Um estudo de caso do uso de mineração de dados e
aprendizado de máquina no aprimoramento de inspeções
de estações radio base**

Trabalho de conclusão

Universidade de Pernambuco – UPE

Residência Tecnológica em Inteligência Artificial

Brasil

List of Figures

Figure 1 – Número de Abonados vis-à-vis Avaliados pré e pós balanceamento . . .	8
Figure 2 – Distribuição de acúracias. Acurácia média anotada em cada caixa. . . .	9

List of Tables

Contents

	Introdução	5
1	Descrição do Caso	6
2	Solução proposta	7
	Conclusão	10
	Bibliography	11

Introdução

Nas últimas décadas a temática do impacto social da inteligência artificial vem tomando centralidade no imaginário prospectivo do cidadão médio, da comunidade científica e dos agentes estatais (CAMERON; WISHER, 1991; COCKBURN; HENDERSON; STERN, 2018; MAKRIDAKIS, 2017). A ascensão do assunto na opinião pública não é desconexa de mudanças no contexto econômico e político (KOGUT, 2003). A difusão da internet na sociedade, culminando nas tecnologias IoT (GUBBI et al., 2013), faz com que dados passem a ser consideradas pela *The Economist*¹ o novo petróleo.

Esse papel dos dados pressupõe a capacidade dos agentes econômicos de extrair valor deles. É essa a seara de inserção dos algoritmos de inteligência computacional, particularmente os de aprendizado de máquina. Algoritmos de aprendizado de máquina são aqueles que aprendem com uma experiência com relação a alguma tarefa e uma medida de performance se a performance na tarefa melhora com a experiência (CARBONELL; MITCHELL; MICHALSKI, 1984). Se os dados são o novo petróleo então os algoritmos utilizados para extrair informação e aprender com esses dados podem ser considerados os novos motores da economia.

Embora grandes empresas de tecnologia como Google, Facebook e Amazon façam uso de grandes arquiteturas de redes neurais artificiais as quais necessitam de dezenas de horas de treinamento em unidades de processamento gráfico, a realidade da maior parte das empresas que buscam se inserir nessa nova era algorítmica difere em escopo (CANZIANI; PASZKE; CULURCIO, 2016). Se por um lado a inteligência artificial traz a possibilidade de uma riqueza de aplicações e otimizações no processo produtivo das empresas, por outro lado se faz necessária uma infraestrutura de dados que permita a aplicação dessas técnicas e uma “pipeline” de mineração e recuperação de informação (SCHÜTZE; MANNING; RAGHAVAN, 2007). Ademais a restrição orçamentária e computacional e o imperativo da interpretabilidade² do funcionamento dos algoritmos nos direciona, nesses casos medianos, à algoritmos mais bem estabelecidos e simples em comparação aos de alta publicização (DREISEITL; OHNO-MACHADO, 2002).

O presente estudo apresenta um caso de sucesso da aplicação de sistemas inteligentes de recuperação e análise de informação de relativa simplicidade no aprimoramento de um processo rotineiro na indústria de telecomunicações: a inspeção de estações rádio base.

¹ Fonte: <<https://tinyurl.com/y39u52kk>>. Acessado em 1 de Novembro de 2019.

² No contexto de aprendizado de máquina a interpretabilidade é definida por Doshi-Velez e Kim (2017, p.2) "como a habilidade de explicar ou apresentar em termos compreensíveis para humanos". Uma definição equivalente de interpretabilidade é: o grau no qual um humano pode compreender a causa de uma decisão (MILLER, 2018).

1 Descrição do Caso

Como referenciado anteriormente o sistema alvo de interesse do nosso estudo está inserido no âmbito da indústria de telecomunicações. Na rede de celulares a mediação entre o celular dos usuários e as companhias telefônicas é feita pelas Estações Rádio Base (doravante ERB ou sítio celular). São nesses sítios que estão instalados os equipamentos necessários para a comunicação entre aparelhos celulares e as centrais de comunicação das agências telefônicas. Nesses ambientes são realizadas vistorias frequentes tendo em vista sua relevância para a qualidade do serviço de telefonia. Nessas vistorias são checados itens referentes às chaves do sítio, à rua de acesso, alarmes externos, aterramento, baterias, cabos, fontes de energia, antenas, dentre centenas outros. Essa vistoria é um trabalho conjunto entre técnicos que visitam os sítios e engenheiros de telecomunicação que analisam as informações. Atualmente essa troca de informação é feita da seguinte maneira: o técnico visita a ERB e para cada item de um *checklist*, que tem até 600 itens a depender da empresa de telefonia detentora do sítio, tiram fotos que são enviadas a um sistema, onde são aceitas ou rejeitadas pelos engenheiros na central. Contudo, nem todo item precisa ser checado a depender de condições particulares da ERB. Estes itens são, portanto, abonados.

Em conversas com técnicos e engenheiros responsáveis pelas inspeções foram identificadas ao menos duas possibilidades de aplicação de inteligência computacional no aperfeiçoamento do processo: a definição de quais itens são abonados e quais são aprovadas ou rejeitadas. O problema da dispensa do item, enfoque do presente trabalho, é que os técnicos não sabem de antemão quais itens devem ser abonados em um determinado sítio. Ao chegarem a ERB, desta forma, primeiro devem checar dentre centenas de itens em uma lista quais são dispensáveis e só então iniciam o trabalho da vistoria propriamente dita. Isso contribui drasticamente para a lentidão da atividade. Nossa contribuição para a redução do tempo despendido nessa checagem é descrita em seguida.

2 Solução proposta

Temos por problema a determinação de quais itens de um checklist são passíveis de abono. Isso pode ser modelado como um problema de classificação binária : dado um conjunto de características de um sítio e qual o item desejamos prever se ele é da classe “abonado” ou não (JAMES et al., 2013). Especialistas apontaram a seguinte lista de características de um sítio que os próprios técnicos usam para abonar manualmente os itens:

- Tipo de site: Gabinete ou Container¹;
- Tipo de tecnologia: WCDMA, LTE, GSM;
- Frequência: 450Mhz, 700Mhz, 850Mhz, 1800Mhz, 2100Mhz, 2600Mhz;
- Equipamentos de radiofrequência (RF): Diplex, Triplex, Quadriplex, EHCU, Filtro, TMA, DTMA

Essas informações, contudo, não estão prontamente disponíveis. Uma fonte possível de informação são os Projetos Preliminares de Instalação (PPIs). Eles estão disponíveis em um sistema interno das empresas de telefonia, ao qual nos foi dado acesso, em formato pdf. Tivemos acesso também à base de checklists dos sítios. Identificamos 602 ERBs cadastrada nesse sistema das quais baixamos cerca de 150 PPIs e os checklists de fevereiro a setembro. Dentre os PPIs foram identificados 3 padrões de documento. Como um esforço inicial trabalhamos na extração de informação de um único padrão. Dado esse recorte de um único tipo de documento, a intersecção entre o grupo de sítios que tínhamos tanto o checklist quanto o PPI tem uma cardinalidade de 44.

As características das ERBs estavam de presentes de forma não estruturada em tabelas e textos nos ppis. A informação não contida nas tabelas, extraídas por meio de pacotes especializados, foi obtida por meio da tokenização dos textos. Desta forma geramos automaticamente uma base de características dos sítios. A partir da intersecção entre a base de características e a base de itens geramos um banco de dados de 19000 observações. Na base temos 322 itens únicos, com uma mediana de 243 itens por ERB, e 19 atributos ('Items', 'Status', 'Operadora', 'Tipo de Site', 'WCDMA', 'LTE', 'GSM', '450Mhz', '700Mhz', '850Mhz', '1800Mhz', '2100Mhz', '2600Mhz', 'Diplex', 'Triplex', 'Quadriplexer', 'EHCU', 'Filtro', 'TMA', 'DTMA'), onde todos menos “Item” e “Tipo de Site” são variáveis binárias.

¹ Na verdade existem 4 variantes de site: 'RT', 'GF', 'RF', 'IN' **PRECISO QUE RESPONDAM O EMAIL**

Uma inspeção inicial na base nos permitiu identificar um desbalanceamento no número de itens avaliados x os abonados, no “Status” do item. O desbalanceamento das classes impacta na performance preditiva de modelos, na medida em que o modelo ganha um viés para a classe majoritária simplesmente pelo maior número de observações dessa classe, aumentando, portanto, o número de falso negativos (FACELI et al., 2011). Como demonstrado na Figura 1 o número de itens abonados era mais do dobro dos itens avaliados, de forma que optamos pela sobreamostragem da classe minoritária por meio de um método de interpolação padrão: o SMOTE (Synthetic Minority Over-sampling Technique) (CHAWLA et al., 2002).

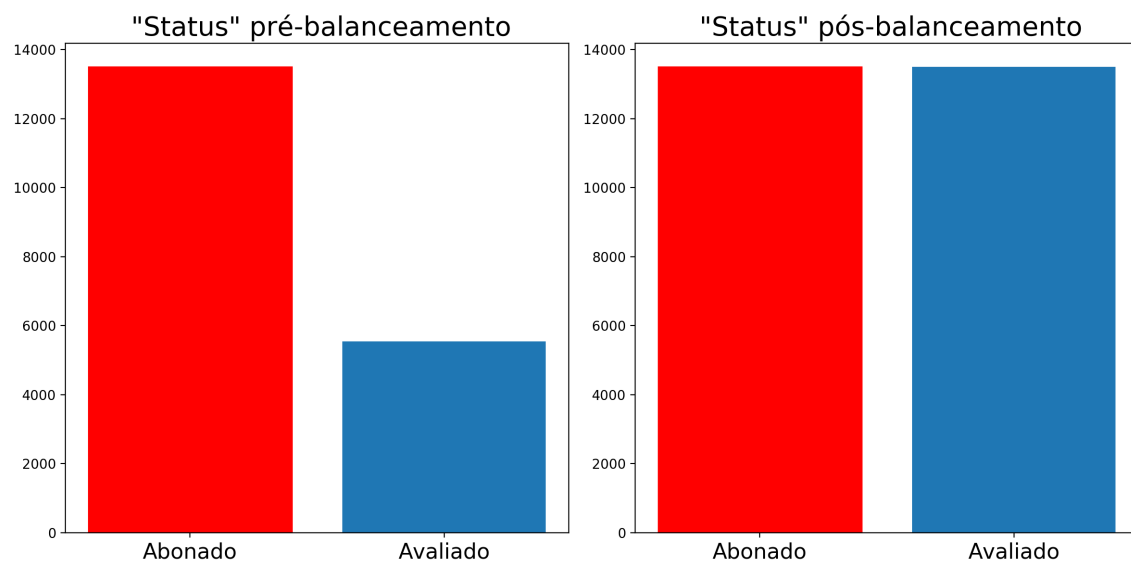


Figure 1 – Número de Abonados vis-à-vis Avaliados pré e pós balanceamento

Após o rebalanceamento codificamos os atributos “Item” e “Tipo de Site” por meio de one-hot-encoding. Uma outra opção de codificação seria atribuir um número inteiro a cada item, mas essa estratégia confundiria o modelo ao implicitamente atribuir ordem a uma variável nominal (FACELI et al., 2011). Uma vez concluído o préprocessamento partimos para o uso de modelos preditivos de aprendizado de máquina. Começamos a análise por meio da validação cruzada, k-fold com 10 folds, dos seguintes modelos: Decision Tree, Multilayer Perceptron, Logistic Regression, Random Forest, Xgboost. **apresentar cada um e quais parâmetros usados** A Figura 2 demonstra a distribuição de acurácias dos classificadores.

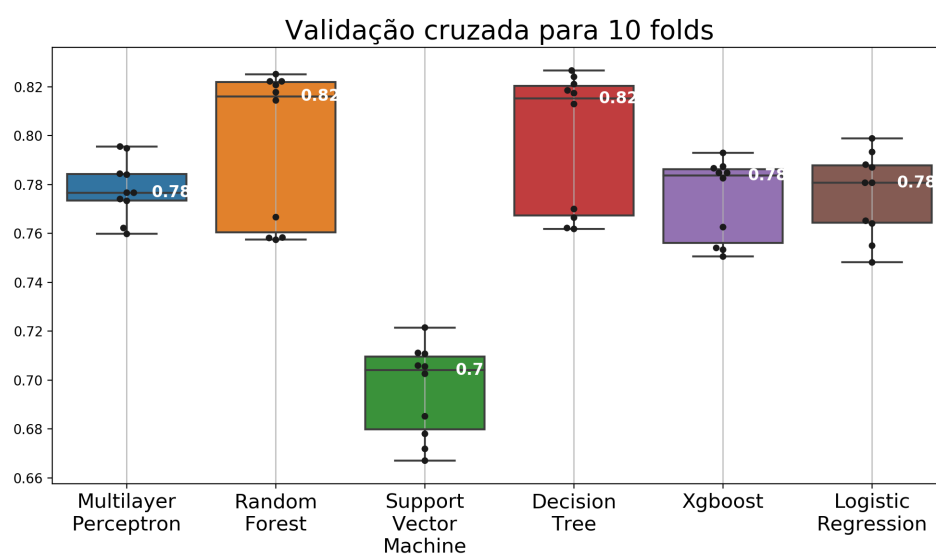


Figure 2 – Distribuição de acúracias. Acurácia média anotada em cada caixa.

Conclusão

Bibliography

- CAMERON, J.; WISHER, W. *Terminator 2: Judgment Day*. [S.l.]: USA, 1991. Citado na página 5.
- CANZIANI, A.; PASZKE, A.; CULURCIELLO, E. An analysis of deep neural network models for practical applications. *arXiv preprint arXiv:1605.07678*, 2016. Citado na página 5.
- CARBONELL, J. G.; MITCHELL, T. M.; MICHALSKI, R. S. *Machine learning: An artificial intelligence approach*. [S.l.]: Springer-Verlag, 1984. Citado na página 5.
- CHAWLA, N. V. et al. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, v. 16, p. 321–357, 2002. Citado na página 8.
- COCKBURN, I. M.; HENDERSON, R.; STERN, S. *The impact of artificial intelligence on innovation*. [S.l.], 2018. Citado na página 5.
- DOSHI-VELEZ, F.; KIM, B. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*, 2017. Citado na página 5.
- DREISEITL, S.; OHNO-MACHADO, L. Logistic regression and artificial neural network classification models: a methodology review. *Journal of biomedical informatics*, Elsevier, v. 35, n. 5-6, p. 352–359, 2002. Citado na página 5.
- FACELI, K. et al. Inteligência artificial: Uma abordagem de aprendizado de máquina. 2011. Citado na página 8.
- GUBBI, J. et al. Internet of things (iot): A vision, architectural elements, and future directions. *Future generation computer systems*, Elsevier, v. 29, n. 7, p. 1645–1660, 2013. Citado na página 5.
- JAMES, G. et al. *An introduction to statistical learning*. [S.l.]: Springer, 2013. Citado na página 7.
- KOGUT, B. M. *The global internet economy*. [S.l.]: MIT Press, 2003. Citado na página 5.
- MAKRIDAKIS, S. The forthcoming artificial intelligence (ai) revolution: Its impact on society and firms. *Futures*, Elsevier, v. 90, p. 46–60, 2017. Citado na página 5.
- MILLER, T. Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence*, Elsevier, 2018. Citado na página 5.
- SCHÜTZE, H.; MANNING, C. D.; RAGHAVAN, P. *An introduction to information retrieval*. [S.l.]: Cambridge University Press,, 2007. Citado na página 5.