

Um Estudo sobre Aceitabilidade de Resultados de Sistemas Aproximados em Redes Neurais Artificiais

Guilherme Saides Serbai, Rogério Aparecido Gonçalves, João Fabrício Filho

¹Departamento Acadêmico de Computação
Universidade Tecnológica Federal do Paraná (UTFPR) – Campus Campo Mourão
Caixa Postal 271 – 87301-899 – Campo Mourão – PR – Brazil

guilhermesaidesserbai@alunos.utfpr.edu.br, {rogerioag, joaof}@utfpr.edu.br

Abstract. *Approximate computing is a Computer Science field that improves performance and energy efficiency at the expense of controlled precision reduction. In the context of training and validating classification neural networks, investigating the impact of approximation on image quality is essential to determine how much data degradation can be tolerated without compromising result validity. This work examines the acceptability and quality of results under different levels of image approximation. We employ a residual neural network (ResNet-50) [He et al. 2016] and evaluate it across various training and validation scenarios using both approximated and non-approximated images from the Imagenette2 dataset [FastAI 2019]. Our objective is to investigate data acceptability thresholds and their relationship with the network's prediction quality. The results demonstrate the correlation between acceptability and accuracy in neural network validation and training with approximated images. The ResNet-50 [He et al. 2016] achieved accuracy ranging from ($\geq 13.7\%$ in scenarios with high training divergence) to ($\geq 68.6\%$ in conditions similar or identical to training), proving that approximate computing can be viable when data similarity is maintained - a crucial factor for energy-efficient and high-performance systems.*

Resumo. *A computação aproximada é uma área da Ciência da Computação que busca aumentar o desempenho e a eficiência energética às custas de uma redução controlada na precisão. No contexto do treinamento e validação de redes neurais de classificação, investigar o impacto da aproximação na qualidade das imagens é crucial para determinar até que ponto podemos degradar os dados sem comprometer a validade dos resultados. Este trabalho investiga a aceitabilidade e qualidade dos resultados frente a diferentes níveis de aproximação em imagens. Para isso, foi utilizado uma rede neural residual (ResNet-50) [He et al. 2016] e submetida a diferentes cenários de treinamento e validação com imagens aproximadas e não aproximadas do dataset Imagenette2 [FastAI 2019]. O objetivo é investigar os níveis de aceitabilidade dos dados e sua relação com a qualidade das previsões da rede. Os resultados obtidos evidenciam a relação entre aceitabilidade e qualidade na validação e no treinamento de redes neurais com imagens aproximadas. A ResNet-50 [He et al. 2016] apresentou acurácia de ($\geq 13.7\%$ cenários com grande divergência do treinamento) a ($\geq 68.6\%$ condições próximas ou iguais ao treinamento), demonstrando que a computação aproximada pode ser viável quando mantida a similaridade dos dados - crucial para sistemas eficientes em energia e desempenho.*

1. Introdução

Sistemas aproximados oferecem maior eficiência na Computação, ao custo de precisão no resultado final [Felzmann et al. 2021]. Nos dias atuais, a demanda por energia e poder de processamento cresce exponencialmente [LoganKugler 2015], principalmente no cenário da Inteligência Artificial (IA). Redes Neurais complexas exigem grande quantidade de tempo e energia para serem treinadas, o que inviabiliza sua aplicação dependendo das especificações do *hardware* disponível.

O cenário atual da computação aproximada, em redes neurais, foca em tornar as aproximações mais seguras e eficientes. Soluções como o *AXNet* [Peng et al. 2018] integram aproximador e preditor em uma única rede treinável, otimizando desempenho. Ferramentas como o *ApproxANN* [Zhang et al. 2015] melhoram a tolerância a erros e a eficiência energética no aprendizado profundo. Além disso, métodos de treinamento especializado para *hardware* aproximado como [Li et al. 2023] permitem acelerar o processo em até 18 vezes. O trabalho de [Felzmann et al. 2021] ressalta também que a qualidade dos dados aproximados deve ser avaliada com base na utilidade contextual, e não apenas em métricas quantitativas.

A relação entre aceitabilidade e qualidade dos resultados obtidos pela validação de uma rede neural treinada com aproximação em nível de dados ainda é pouco explorada. Compreender como a qualidade dos resultados se comporta diante desse tipo de abordagem é essencial para avançarmos nos estudos sobre computação aproximada e suas aplicações. Nesse sentido, é fundamental construir um conjunto de dados adequado, que possibilite uma análise aprofundada sobre o impacto da computação aproximada em redes neurais.

Esta pesquisa tem como objetivo investigar a relação entre a aceitabilidade dos resultados e os níveis de aproximação dos dados de entrada de uma Rede Neural residual (*ResNet-50*). Busca-se demonstrar que a computação aproximada é uma estratégia consciente e segura, desde que adaptada ao contexto da aplicação, contribuindo para uma maior eficiência computacional sem comprometer significativamente a confiabilidade das previsões.

A metodologia utilizada nesta pesquisa se fundamentou no uso de uma rede neural residual (*ResNet-50*) [He et al. 2016] e no *dataset Imagenette2* [FastAI 2019]. Foram criados diferentes cenários de aproximação, tendo como principal técnica o algoritmo *Bitflip*.

Os resultados mostram que a acurácia da rede neural depende diretamente do tipo de dado usado no treinamento. Quando os dados de teste são semelhantes aos dados aproximados usados no treino, a acurácia é maior. Por outro lado, ao testar com dados precisos, a acurácia caiu de forma significativa. Também foi observado que quanto maior o nível de aproximação, mais difícil foi para a rede aprender.

Este trabalho demonstra que a computação aproximada pode ser uma estratégia eficiente e segura para redes neurais quando adequada ao contexto de aplicação, comprovando que é possível manter resultados aceitáveis mesmo com dados degradados.

2. Materiais e Métodos

Os principais materiais utilizados na pesquisa foram o *dataset Imagenette2* e uma rede neural residual (*ResNet-50*), previamente treinada no IMAGENET1K, disponível por meio do módulo `torchvision.models` da biblioteca `PyTorch`. Sobre o *Imagenette2*, foi aplicada uma técnica de computação aproximada, sendo ela o algoritmo *Bitflip*. O algoritmo *Bitflip* consiste em alterar um *bit* aleatório nos três octetos de cada canal RGB dos *pixels* da imagem. Nesse contexto, utilizando o *Bitflip*, foram construídos diferentes cenários nos quais foram afetados 25%, 50% e 100% dos *pixels* da imagem. Para cada percentual, foram criados subcenários em que a função de modificação foi aplicada 1x, 2x, 3x e 4x vezes, com o objetivo de simular diferentes níveis de perturbação nos dados. Vale ressaltar que em técnicas de aproximação configuráveis, quanto maior a quantidade de erros, como *bitflips*, maiores os possíveis ganhos energéticos e de desempenho [Felzmann et al. 2020]. Por fim, a *ResNet-50* foi treinada em todos os cenários, incluindo o original (sem aproximação), com o conjunto de treino (*train*) do *Imagenette2*. Os modelos foram comparados usando o conjunto de validação (*val*).

3. Resultados e Discussão

Ficou evidente, a partir dos resultados obtidos, que a computação aproximada no contexto de redes neurais é fortemente dependente do cenário de treinamento. A Figura 1 mostra que, quanto mais os dados de teste se assemelham ao contexto aproximado utilizado durante o treinamento, maior é a acurácia obtida. Nesse sentido, um modelo treinado com dados aproximados apresentou acurácia significativamente menor ao ser testado com dados precisos, enquanto obteve melhor desempenho com dados também aproximados. Isso evidencia que a rede neural tende a apresentar maior acurácia quanto mais próximo o cenário de inferência estiver do cenário de treinamento, como ilustrado na Figura 1.

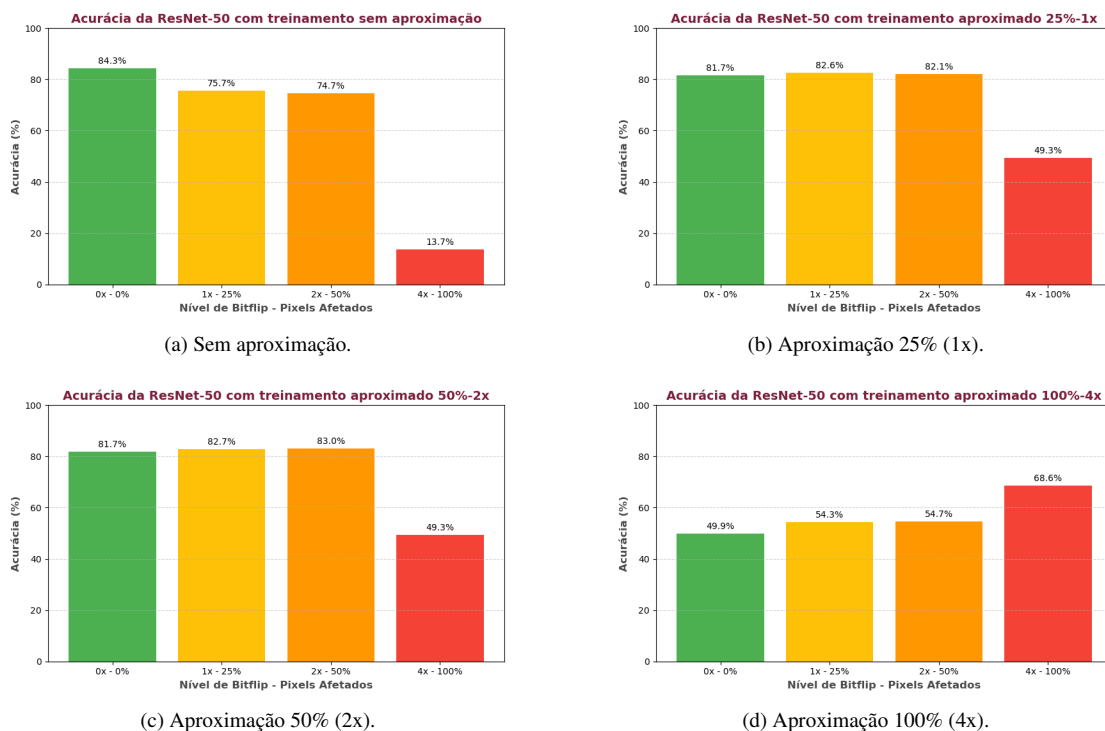


Figura 1. Acurácia dos diferentes modelos da (*ResNet-50*) em relação aos *datasets* aproximados e não aproximados.

Nesse contexto, vale destacar que, à medida que o nível de aproximação se tornava mais crítico, maior era a dificuldade encontrada pela rede para alcançar uma acurácia de treinamento comparável à obtida com dados submetidos a aproximações mais leves ou até mesmo sem aproximação, como demonstram as Figuras 1d e 1a. O cenário sem aproximação apresentou uma acurácia de 84,3% (dentro do contexto de treinamento) 1a, enquanto os cenários aproximados (fora do contexto de treinamento) obtiveram acurácias menores, variando de 13,7% a 75,7%. Vale destacar que, mesmo com 100% de ruído aplicado 4 vezes, resultando em uma acurácia de 13,7% (1a), ainda foi possível obter resultados válidos, pois as alterações aleatórias preservaram traços visuais nas imagens. Esse cenário se inverteu à medida que o contexto de aproximação se aproximava do contexto de treinamento: a acurácia aumentava, enquanto o modelo sem aproximação apresentava queda, chegando, no pior caso, a 49,9% (1d). Esses resultados reforçam o que foi discutido anteriormente sobre a dependência do desempenho da rede em relação ao contexto de treinamento.

4. Considerações Finais

O trabalho destacou que a relação entre aceitabilidade dos resultados está intimamente ligada ao contexto e aos testes adequados, reforçando a hipótese de que a computação aproximada pode ser segura e eficiente, desde que adaptada ao contexto e à utilidade dos dados dentro desse mesmo cenário. Vale ressaltar que em técnicas de aproximação configuráveis, quanto maior a quantidade de erros, como *bitflips*, maiores os possíveis ganhos energéticos e de desempenho [Felzmann et al. 2020].

References

- FastAI (2019). Imagenette2 dataset.
- Felzmann, I., Fabrício Filho, J., de Oliveira, J. R., and Wanner, L. (2021). Special Session: How much quality is enough quality? A case for acceptability in approximate designs. In *2021 IEEE 39th International Conference on Computer Design (ICCD)*, pages 5–8, Los Alamitos, CA, USA. IEEE Computer Society.
- Felzmann, I., Fabrício Filho, J., and Wanner, L. (2020). Risk-5: Controlled approximations for risc-v. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 39(11):4052–4063.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778.
- Li, T., Li, S., and Gupta, P. (2023). Training neural networks for execution on approximate hardware.
- LoganKugler (2015). Is “good enough” computing good enough? the energy-accuracy trade-off in approximate computing. *Communications of the ACM*, 58(5):12.
- Peng, Z., Chen, X., Xu, C., Jing, N., Liang, X., Lu, C., and Jiang, L. (2018). Axnet: Approximate computing using an end-to-end trainable neural network.
- Zhang, Q., Wang, T., Tian, Y., Yuan, F., and Xu, Q. (2015). Approxann: An approximate computing framework for artificial neural network. In *2015 Design, Automation Test in Europe Conference Exhibition (DATE)*, pages 701–706.