

Análise de Caso: Vieses em IA na Triagem de Emergência do Hospital de Bordeaux

Sumário

- 1. Introdução
- 2. Contextualização do Problema: Os Nossos Vieses Repetidos no Código
 - 2.1. A Dinâmica Técnica: A IA é um Reflexo e um Amplificador dos seus Dados
 - 2.2. A Aplicação Prática: Riscos em Cenários Reais
 - 2.3. O Impacto Social: A Amplificação de Desigualdades
- 3. Análise sob a Ótica dos Fundamentos Éticos em Computação
 - 3.1. Viés e Justiça: A Quantificação da Desigualdade
 - 3.2. Transparência e Explicabilidade: Iluminando a "Caixa-Preta"
 - 3.3. Impacto Social e Direitos: O Dever de Proteção
 - 3.4. Responsabilidade e Governança: Uma Cadeia de Corresponsabilidade
- 4. Posicionamento e Recomendações
- 5. Conclusão
- Referências

1. Introdução

Este artigo, elaborado por nosso grupo, analisa os riscos sistêmicos de se treinar inteligências artificiais (IA) com dados enviesados. Adotamos os fundamentos da ética em computação como guia para nossa análise, aplicando-os a um estudo de caso robusto e com um volume de dados significativo para exame. Para isso, nos baseamos em dois materiais centrais: uma reportagem do portal Healthcare in Europe e o artigo científico original de Guerra-Adames et al., que serviu de fonte para a matéria. A pesquisa investiga como a IA pode ser empregada não apenas para detectar, mas também para mitigar vieses na triagem de emergência, utilizando um dataset coletado no pronto atendimento do Hospital Universitário de Bordeaux entre 2013 e 2021. A partir deste estudo, nosso grupo elaborou uma análise para compreender os desafios e, o que é mais importante, as oportunidades que esta

tecnologia apresenta. Adicionalmente, propomos sugestões práticas para assegurar que a implementação da IA como ferramenta de apoio seja mais justa e eficiente.

2. Contextualização do Problema: Os Nossos Vieses Repetidos no Código

Uma inteligência artificial não é neutra. Ela é um espelho dos dados que a alimentam e, se esses dados estiverem "contaminados" com preconceitos, o algoritmo não só os aprende, como os amplifica com uma eficiência assustadora.

É exatamente isso que o caso do Hospital de Bordeaux nos ensina. Ao analisar quase meio milhão de atendimentos, a pesquisa provou que a IA, treinada com decisões humanas, herdou e automatizou um viés de gênero. O resultado? Um sistema que, em vez de otimizar, passou a repetir erros da triagem do setor de emergência, com consequências potencialmente fatais.

O que está em jogo aqui vai além da tecnologia. Se trata de como a automação e multiplicação de um preconceito pode desgastar a confiança em sistemas críticos e violar o direito fundamental a um tratamento justo. O estudo não revela uma falha da IA, mas sim uma falha humana — um alerta de que, sem ética no design, estamos construindo um futuro que multiplica nossos piores erros em código.

2.1. A Dinâmica Técnica: A IA é um Reflexo e um Amplificador dos seus Dados

A inteligência artificial aprende a partir dos dados que a treinam. Isso é um princípio fundamental e, ao mesmo tempo, um alerta que não pode ser ignorado. Se essa base de dados, o dataset, já vem com vieses, o resultado é praticamente inevitável: o algoritmo vai acabar assimilando esses padrões e adotando eles como parte da sua lógica operacional. E isso não é só teoria — o caso do Hospital de Bordeaux mostra isso de forma bem concreta. A pesquisa de Guerra-Adames e seu grupo não ficou só na parte técnica; eles se basearam num volume impressionante de dados reais, com quase meio milhão de registros de atendimento sendo analisados. A conclusão foi direta, e preocupante: ao treinar uma IA com as decisões do grupo de enfermagem, o modelo virou um espelho fiel dos vieses de gênero que já influenciavam aquelas escolhas.

2.2. A Aplicação Prática: Riscos em Cenários Reais

O risco de usar uma IA treinada com dados tendenciosos é especialmente alto em domínios críticos como o pronto atendimento. A pesquisa focou na triagem de emergência, um processo que funciona como um filtro inicial de alta responsabilidade. É nesse momento que o grupo de enfermagem, usando sua experiência, decide com rapidez quem necessita de atenção imediata, garantindo que os pacientes em maior risco sejam atendidos prioritariamente. O estudo revelou um fator fundamental: a experiência profissional atua como um moderador do viés. Ao analisar a influência do tempo de carreira, os pesquisadores notaram que o viés de gênero nas decisões diminuía conforme os anos de experiência do grupo aumentavam. Isso sugere que, embora vieses possam existir, a maior experiência profissional (medidas em tempo de atuação na área) ajuda a reduzir a subestimação de risco em pacientes do sexo feminino. A chave para um sistema mais justo, portanto, parece ser a combinação da expertise humana com ferramentas de IA capazes de identificar e corrigir esses desvios. As consequências de ignorar isso são graves. Na saúde, um erro de viés pode levar à deterioração do quadro clínico e, em casos extremos, a desfechos fatais por falha na priorização do atendimento.

2.3. O Impacto Social: A Amplificação de Desigualdades

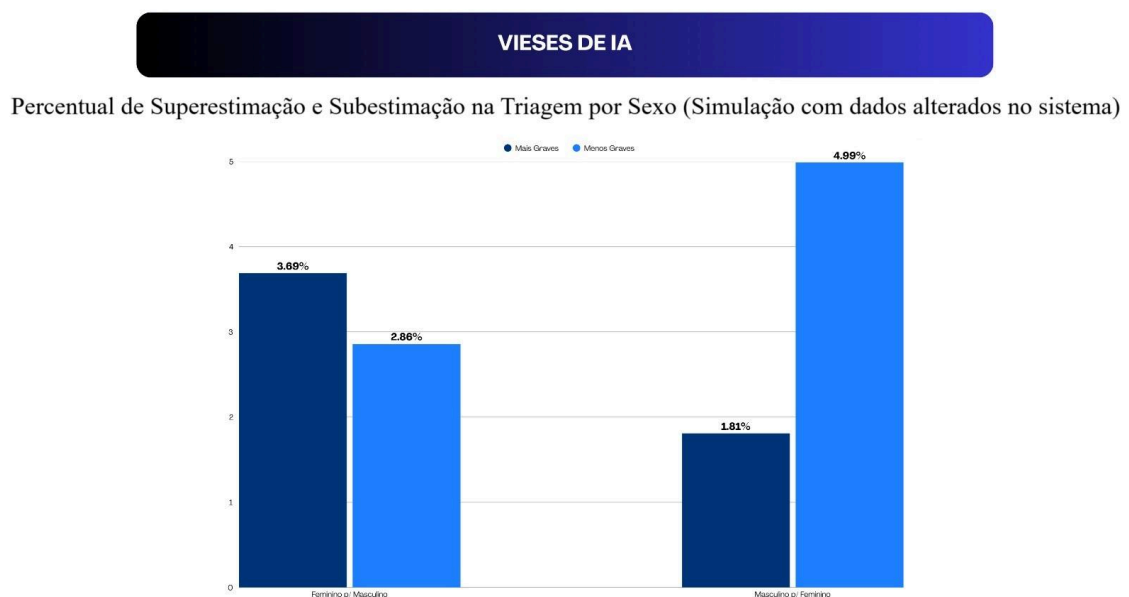
Quando tecnologias de uso diário reproduzem e amplificam preconceitos de gênero, raça ou idade, elas acabam por reforçar estereótipos e aprofundar desigualdades sociais já existentes. Um cenário como esse pode minar a confiança do público tanto na tecnologia quanto nas instituições que a utilizam. Além disso, a perpetuação de vieses em sistemas de saúde viola diretamente os princípios de equidade e o direito fundamental a um tratamento justo e igualitário. A questão transcende o técnico: a IA, ao refletir dados enviesados, torna-se um agente de risco em aplicações críticas, com potencial para causar danos significativos e reforçar a exclusão social.

3. Análise sob a Ótica dos Fundamentos Éticos em Computação

3.1. Viés e Justiça: A Quantificação da Desigualdade

O estudo de Guerra-Adames et al. comprovou com dados a existência de viés nas decisões de triagem. Utilizando um modelo de IA treinado com o histórico do hospital, a pesquisa revelou que a automação da classificação de prioridades gerou os seguintes resultados:

- * Mulheres apresentaram um risco aproximadamente **5% maior de serem subestimadas na triagem** (classificadas com um nível de gravidade inferior ao real) em comparação com homens;
- * Homens, em contrapartida, **foram superestimados** (classificados com um nível de gravidade superior) **em 3,7% dos casos**;
- * O viés diminuía de forma significativa com o aumento da experiência do grupo.



Essa análise demonstrou que o modelo de IA replicou os padrões humanos onde o sexo do paciente, e não apenas sua condição clínica, estava influenciando o resultado da triagem. Isso nos mostrou como um viés de dados, presente no histórico de atendimentos e agravado pela inexperiência dos atendentes do setor de triagem, se transforma em um viés algorítmico que falha em sua promessa de justiça, colocando o grupo feminino em uma posição de maior vulnerabilidade e

distribuindo os riscos de maneira inaceitável, expondo as mulheres a ainda mais riscos antes de terem um atendimento médico em situações de emergência.

3.2. Transparência e Explicabilidade: Iluminando a "Caixa-Preta"

Apesar do modelo de IA utilizado por aquele hospital ser do tipo "caixa-preta" (black box), cujos processos internos não são totalmente conhecidos ou facilmente explicados, os pesquisadores usaram a própria tecnologia para quantificar esse viés. Por meio de simulações, eles alteraram o sexo dos pacientes no programa e observaram as mudanças nos resultados. Essa abordagem de "engenharia reversa" contorna o problema de explicar como cada decisão foi tomada. Mesmo que ainda não seja possível explicar o raciocínio por trás de cada decisão tomada por aquele algoritmo, essas simulações tornaram o comportamento do algoritmo transparente.

3.3. Impacto Social e Direitos: O Dever de Proteção

Num contexto de saúde, a desigualdade pode ser fatal. Os resultados do estudo reforçam a necessidade de garantir que sistemas decisórios, especialmente os automatizados, não perpetuem ou agravem a discriminação. No cenário brasileiro, a Lei Geral de Proteção de Dados (LGPD) adiciona uma camada de complexidade, exigindo que a proteção de dados pessoais e sensíveis de saúde seja um pilar na concepção de sistemas éticos.

3.4. Responsabilidade e Governança: Uma Cadeia de Corresponsabilidade

A responsabilidade pela equidade de um sistema de IA é compartilhada. Ela recai sobre os desenvolvedores, os profissionais de saúde que geram os dados de treinamento e as instituições que implementam e validam esses sistemas. A pesquisa serve, portanto, como um guia para a criação de protocolos de governança mais justos. Um desafio apontado é a ausência de um "padrão-ouro" para a triagem, o que dificulta a avaliação objetiva dos modelos. O estudo sugere a criação de um "padrão-prata", baseado no consenso de especialistas, que poderia ser usado tanto para treinar a IA quanto para aprimorar a formação do grupo de enfermagem.

Este caso nos mostra uma falha no processo de desenvolvimento de um sistema para uma área crítica. Fica evidente que a responsabilidade começa no design, e que a aplicação dos princípios de "Ethical AI by Design" é uma etapa fundamental para garantir que esses sistemas sejam mais justos e seguros. Podemos aprender que a verdadeira falha foi não tratar a ética e a justiça como requisitos técnicos obrigatórios desde o início.

4. Posicionamento e Recomendações

Com base em nossa análise, concluímos que a IA não deve ser descartada da área da saúde, mas sim aprimorada, redesenhada e rigorosamente monitorada. O próprio estudo demonstra seu potencial como ferramenta para diagnosticar e corrigir o problema do viés.

Nossas três recomendações práticas são:

1. Auditoria Contínua: Implementar sistemas de IA para auditar em tempo real as decisões de triagem, identificando padrões de viés e permitindo intervenções rápidas;
2. Treinamento Humano Aprimorado: Utilizar simulações baseadas em IA para criar programas de treinamento que exponham os profissionais a cenários de viés, aprimorando sua capacidade de reconhecê-los e mitigá-los;
3. Ferramentas de Triagem Assistida: Desenvolver a IA não como um decisor autônomo, mas como uma ferramenta de apoio que oferece uma segunda opinião ou emite alertas quando um potencial viés é detectado na decisão humana.

5. Conclusão

A pesquisa utilizada no estudo aponta para um futuro onde a IA pode ser uma poderosa aliada na construção de um sistema de saúde mais justo. A tecnologia demonstrou ter o potencial não apenas de replicar, mas também de revelar e combater vieses sistêmicos, tornando a tomada de decisão mais equitativa.

Contudo, os riscos de modelos treinados com dados enviesados são significativos. A supervisão humana, o desenvolvimento ético e a melhoria contínua são, portanto, essenciais. Acreditamos que a sinergia entre a tecnologia e a ética é o caminho para um futuro onde a inteligência artificial na saúde seja sinônimo de justiça e equidade para todos.

Referências

1. Healthcare in Europe. *AI can help detect, reduce bias in emergency medicine*. Disponível em: <https://healthcare-in-europe.com/en/news/ai-detect-reduce-bias-emergency-medicine.html>. Acesso em: 30 ago. 2025, 15h 37min.
2. Guerra-Adames, M., et al. (2023). *Quantifying and Mitigating Gender Bias in Emergency Triage using Causal-based Counterfactuals*. In: Proceedings of the 2nd Conference on Health, Inference, and Learning (CHIL), Vol. 259. PMLR. Disponível em: <https://proceedings.mlr.press/v259/guerra-adames25a.html>. Acesso em: 30 ago. 2025, 15h 45min.