

# AI Tooling for Bias-Free Product Ranking and Recommendation in Retail Platforms

Udit Agarwal

udit15@gmail.com

## Abstract:

The digital retail economy is increasingly governed by recommendation engines and product ranking algorithms that curate consumer experiences and allocate market visibility. While these systems are designed to optimize for accuracy and engagement, they frequently manifest systemic biases that distort competition, marginalize niche suppliers, and reinforce societal inequalities. This research report provides a comprehensive examination of the AI tooling and architectural frameworks required to achieve bias-free product ranking in modern retail platforms. By synthesizing contemporary scholarly research and industrial technical disclosures, the analysis categorizes the multi-dimensional nature of algorithmic bias—including popularity, position, exposure, and demographic biases—and evaluates the effectiveness of pre-processing, in-processing, and post-processing mitigation strategies. The report delves into specific case studies from leading platforms such as Amazon, eBay, and Etsy, detailing their transition from binary feedback models to multi-relevance architectures and the implementation of large language models (LLMs) as evaluative judges. Furthermore, the study explores the role of microservices-oriented search architectures and specialized open-source toolkits in facilitating algorithmic transparency and governance. The findings suggest that achieving neutrality in product ranking requires a socio-technical approach that integrates rigorous statistical de-biasing, domain-driven architecture, and continuous lifecycle monitoring to preserve market integrity and consumer trust.

**Keywords:** Recommender Systems, Algorithmic Fairness, Product Ranking, Popularity Bias, Position Bias, Retail AI Governance, Machine Learning Fairness Tooling, Socio-technical Systems.

## Introduction

In the current era of hyper-personalized e-commerce, recommendation engines and search ranking algorithms have moved beyond their original role as convenience features to become the primary gatekeepers of the global retail marketplace. The economic gravity of these systems is immense, with industry-leading platforms such as Amazon attributing approximately thirty-five percent of their total revenue directly to the outputs of their recommendation engines. These algorithmic curators shape the choices of billions of users, determining which products are discovered and which remain in the "data shadows" of the long tail. However, as these models grow in complexity—leveraging deep learning and multi-tower architectures—they have increasingly been found to exhibit and amplify systematic biases that were either inherent in the training data or introduced through specific architectural and optimization choices.

The challenge of ensuring bias-free ranking is not merely a technical hurdle but a socio-technical imperative. Algorithmic bias in retail can be defined as systematic favoritism or discrimination that results in unfair outcomes, arbitrarily disadvantaging certain products, sellers, or user demographics. These biases manifest in various forms, from popularity bias—which favors already-successful products regardless of their actual relevance to a specific user—to more insidious forms of demographic discrimination that reflect historical human prejudices. When a platform's ranking logic is compromised by these factors, it creates a "blockbuster effect" where a small number of items dominate visibility, stifling innovation and disadvantaging small and medium-sized enterprises (SMEs) that lack the historical data or promotional resources to overcome algorithmic inertia.

Furthermore, the lack of transparency in many modern recommendation systems—often characterized as "black boxes"—poses significant risks to brand integrity and legal compliance. With the emergence of regulations like the European Union's AI Act and the NIST AI Risk Management Framework, platforms are

now required to demonstrate algorithmic accountability and fairness. This necessitates the deployment of specialized AI tooling designed to detect, measure, and mitigate bias throughout the entire software development lifecycle. This report provides an exhaustive analysis of the contemporary landscape of these tools, the methodologies they employ, and the architectural strategies utilized by major retail platforms to maintain a neutral and equitable digital marketplace. By examining the interplay between data science, system design, and economic policy, this research offers a blueprint for the next generation of responsible retail AI.

### **The Multi-Dimensional Taxonomy of Algorithmic Bias in Retail**

To construct a framework for de-biasing, it is essential to first understand the mechanisms and categories of bias that permeate retail platforms. Research indicates that bias arises from a combination of data collection flaws, algorithmic design choices, and user interaction patterns. These are often categorized as data-driven, model-driven, and behavioral biases, each requiring distinct mitigation strategies.

Popularity bias is perhaps the most documented form of model-based bias in recommender systems. It occurs when algorithms prioritize items with a high volume of historical interactions, creating a self-reinforcing loop where popular items gain more exposure and thus even more interactions. This phenomenon is deeply rooted in Zipf's Law, which observes that in many cultural and retail markets, a tiny fraction of items receives the vast majority of consumer attention, while the "long tail" remains largely invisible. While recommending popular items can be a strong baseline for general accuracy, it often fails users who have niche interests and discourages the discovery of novel products. In retail environments, this results in the blockbuster effect, where the diversity of the product catalog is artificially constricted, harming both the consumer experience and the economic viability of smaller sellers.

Position bias represents a critical behavioral bias inherent in the user interface design. Consumers are conditioned to focus their attention on the top-ranked results or the most prominent positions in a grid-based display. This leads to higher click-through rates for top-positioned items regardless of their actual relevance or quality. When interaction logs are used to train subsequent models without correcting for this position bias, the algorithm mistakenly learns that top items are inherently superior, further entrenching the existing ranking and creating a feedback loop that is difficult to break. This is particularly challenging in e-commerce websites where results are often displayed in two-dimensional grids rather than linear lists, as user attention decay follows more complex patterns involving row-skipping and horizontal browsing behaviors.

Exposure bias and selection bias occur because users are only exposed to a small fraction of the total product catalog. In most retail systems, user feedback is "missing not at random" (MNAR); that is, we only observe interactions for products the system chose to show the user. This creates a truncated view of user preferences. If a user only ever sees high-end luxury goods due to a platform's assumptions, the lack of interaction with budget-friendly items is not an indicator of disinterest but a result of non-exposure. Furthermore, historical and societal biases can be encoded into datasets if the training data reflects past discriminatory practices. For instance, if historical lending or hiring data favored a specific demographic, an AI system trained on that data will likely replicate those patterns, treating historical imbalances as rules for success.

### **AI Tooling and Frameworks for Bias Mitigation**

The operationalization of fairness in retail requires a robust suite of tools capable of auditing and correcting these biases. Modern AI practitioners utilize specialized libraries and frameworks that integrate into the machine learning pipeline, offering techniques categorized by their stage of intervention: pre-processing, in-processing, and post-processing.

IBM's AI Fairness 360 (AIF360) is one of the most comprehensive open-source toolkits in this domain. It provides data scientists with over seventy fairness metrics and eleven bias mitigation algorithms that can be applied across the entire model lifecycle. AIF360 supports fairness checks during data ingestion, model training, and post-prediction, allowing for a comparative analysis of different protected attributes such as race, gender, or age. For retail platforms, AIF360's reweighing algorithm is particularly useful in the pre-processing stage, where it adjusts the weights of different training examples to ensure that underrepresented groups or niche product categories are given equitable influence before the model is even trained.

Microsoft's Fairlearn is another dominant library, focusing on identifying and mitigating harms such as under-representation or disparate impact. It integrates seamlessly with the popular scikit-learn environment, providing a MetricFrame dashboard that allows developers to evaluate model performance across different

sensitive subgroups simultaneously. Fairlearn is often used in the post-processing phase, where algorithms like the Threshold Optimizer can adjust the decision boundaries of a recommendation model to equalize outcomes for different groups without requiring a full retraining of the system. This is especially valuable in high-velocity retail environments where model retraining is computationally expensive.

For measuring and mitigating bias specifically in recommender algorithms, the Holistic AI Library provides a specialized Python-based toolkit. It offers risk mitigation roadmaps that guide users through five key areas: bias, efficacy, robustness, privacy, and explainability. The library's mitigation module helps practitioners apply fairness constraints during the actual learning process (in-processing), which can result in models that are inherently more fair and less likely to pick up on spurious correlations in the data. Additionally, tools like SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-agnostic Explanations) are critical for algorithmic transparency. While not direct mitigation tools, they allow researchers to understand which features are driving a specific product ranking. If a model is found to be relying heavily on a problematic proxy feature that correlates with a protected attribute, these interpretability tools can expose the underlying mechanism of bias, prompting further refinement of the feature engineering process.

### **Advanced Methodologies for Neutral Ranking in Retail**

Beyond general toolkits, retail platforms have developed specialized methodologies to address the unique challenges of e-commerce search and recommendation. These involve complex statistical models designed to "de-convolve" the feedback loops and interaction biases that are endemic to digital marketplaces.

Inverse Propensity Scoring (IPS) has emerged as a foundational technique for de-biasing click logs in search ranking. As research from Etsy indicates, the propensity of a user to interact with a product is heavily influenced by its display position. To obtain an unbiased estimate of a product's relevance, IPS assigns a weight to each interaction that is inversely proportional to the likelihood of the item being seen in that specific position. This requires sophisticated click models that can account for row-skipping and non-linear browsing patterns in 2D grid layouts. By normalizing interaction signals in this way, platforms can train rankers that reflect the true attractiveness of a product rather than its positional luck.

Another significant advancement is the shift from binary feedback to multi-relevance ranking models. Historically, many retail systems used binary labels like "purchased" or "not purchased" to train their rankers. However, research at eBay suggests this approach is fundamentally flawed because it ignores the wealth of intent signals that precede a purchase. Users may click on several high-quality items before selecting only one to buy, leading to "false negatives" for the other relevant products. To solve this, eBay implemented a multi-relevance scale that weights actions like "Add to Watchlist," "Add to Cart," and "Make Offer" based on their historical correlation with final conversion. This provides a much denser signal for the model, reducing the impact of data sparsity and allowing for more diverse and novel recommendations that are not solely reliant on rare purchase events.

Furthermore, the use of large language models (LLMs) as evaluative judges is becoming a standard in governance. Amazon's Nova LLM-as-a-Judge framework is used to provide impartial assessments of AI outputs across various real-world scenarios. This framework is built on a diverse dataset of human preferences across ninety languages and is designed to reflect broad consensus rather than individual viewpoints. To ensure fairness, the framework benchmarking includes rigorous validation steps like reporting results for "positionally swapped" responses to eliminate the potential for the judge model itself to exhibit positional bias. By integrating these LLM judges into their MLOps pipelines, platforms can monitor for fairness regressions and ensure that their systems continue to provide equitable guidance as they evolve.

### **Socio-Technical Architecture and System Scalability**

The implementation of de-biasing tools is deeply tied to the underlying technical architecture of the retail platform. As these systems scale to handle millions of queries per second, the architecture must support both high-throughput inference and rigorous fairness governance.

A microservices-oriented approach is increasingly favored for its ability to decouple search and recommendation from other domain boundaries. In a well-designed architecture, each domain—such as products, reviews, or promotions—owns its own metadata and behavior. By treating search as a standalone domain, platforms can ingest and normalize data from multiple sources through a unified pipeline. This decoupling allows for centralized tuning of ranking models and the integration of specialized fairness checks

without disrupting the operational efficiency of the product or inventory services. Event-driven communication, often mediated by platforms like Kafka, ensures that changes in one part of the system—such as a price update or a new review—are reflected in the ranking models through eventual consistency. To handle the computational demands of modern deep learning rankers, platforms like eBay utilize specialized GPU clusters, such as the Krylov platform, which supports Airflow-based scheduling and high-volume traffic production environments. These systems allow for the deployment of complex architectures like "MicroBERT," a distilled version of the BERT language model optimized for low-latency inference on CPUs. By using knowledge distillation, retailers can maintain the sophisticated semantic understanding required for fair recommendations while ensuring the system remains responsive to user interactions.

### **Economic Implications and Market Competition**

The transition to bias-free ranking has profound economic consequences, particularly regarding market competition and the visibility of small and medium-sized enterprises (SMEs). Research into the impact of neutral versus non-neutral (biased) algorithms suggests that biased systems can lead to a decrease in the intensity of price competition among merchants. When algorithms favor "prominent" firms or those with the resources to buy their way into top positions, smaller sellers who rely on competitive pricing to gain traction are effectively shut out of the market.

Biased ranking also correlates with increased price dispersion, as prominent sellers can maintain higher prices without losing market share to more affordable but less visible competitors. While neutral algorithms generally enhance platform profits and social welfare by matching consumers with the best products, biased systems often prioritize short-term ad revenue at the expense of long-term ecosystem health. Therefore, maintaining algorithmic neutrality is not just an ethical goal but a strategy for ensuring a healthy, competitive marketplace where resources are allocated based on product quality and service improvement rather than just "paid prominence".

### **Conclusion**

Achieving bias-free product ranking and recommendation in retail platforms is a complex undertaking that requires the integration of sophisticated AI tooling, robust statistical methodologies, and domain-driven architectural strategies. As this research has elucidated, the types of bias prevalent in e-commerce are multifaceted, stemming from unrepresentative data, algorithmic design flaws, and behavioral interaction patterns. The emergence of specialized toolkits like IBM's AI Fairness 360 and Microsoft's Fairlearn provides practitioners with the foundational means to measure and mitigate these biases throughout the model lifecycle. Furthermore, advanced techniques such as Inverse Propensity Scoring for grid layouts and multi-relevance feedback loops have demonstrated significant success in de-biasing interaction signals and improving the diversity of product recommendations.

The case studies from Amazon, eBay, and Etsy highlight a broader industry shift toward algorithmic accountability and governance. By leveraging LLM-as-a-judge frameworks and Automated Reasoning checks, these platforms are moving toward a more transparent and equitable curator role. However, the long-term sustainability of these efforts depends on a socio-technical approach that considers the entire AI lifecycle and involves multi-stakeholder collaboration. As regulatory requirements become more stringent, the ability to demonstrate bias-free outcomes will become a key competitive advantage, fostering consumer trust and ensuring that the digital retail marketplace remains an open and fair arena for all participants. The future of retail AI lies in systems that not only optimize for the next click but also preserve the fundamental principles of fairness and market integrity.

### **REFERENCES:**

1. PMC (2022). Algorithmic bias in retail recommendation systems types popularity exposure position bias research paper. National Center for Biotechnology Information.
2. ArXiv (2023). Popularity Bias in Recommender Systems: A Review of Origins and Mitigation. Cornell University.
3. CEUR Workshop Proceedings (2021). Vol-3078, Paper-69: Empirical evaluation of popularity bias in recommendation algorithms.

4. Greenlining Institute (2021). Algorithmic Bias Explained: Report on socioeconomic opportunity and bias in machine learning.
5. ResearchGate (2024). Model Bias in Recommendation Systems: Understanding Impact and Mitigation Techniques.
6. Frontiers in Big Data (2025). Building equitable AI systems: Integrated framework for bias mitigation across the lifecycle.
7. ResearchGate (2025). Algorithmic Bias in Recommendation Systems and Its Social Impact on User Behavior.
8. AI Multiple (2024). Responsible AI Platform: Enterprise and Open-source Tooling for Bias and Governance.
9. NIST (2022). SP 1270: Towards a Standard for Managing AI Bias. National Institute of Standards and Technology.
10. GitHub: RecDebiasing (2023). Survey on bias and de-biasing strategies in Recommender Systems (SIGIR, KDD, RecSys).
11. ArXiv (2024). Position bias and item popularity: De-convolving the feedback loop in ranking systems.
12. ResearchGate (2020). Debiasing Grid-based Product Search in E-commerce. Presented at KDD 2020.
13. Etsy Data Science (2020). Publication Summary: Debiasing Grid-based Product Search in E-commerce.
14. IJFMR (2025). Cognitive biases and their impact on consumer decision-making in online shopping.
15. ResearchGate (2024). Empirical examination of biased personalized product recommendations on consumers.
16. ResearchGate (2024). Ethical Implications and Bias Mitigation in Generative AI Models.
17. Holistic AI (2024). Technical resources for bias mitigation: Python Bias Toolkit and Risk Mitigation.
18. Keymakr (2024). Data bias in AI e-commerce: Identifying and mitigating risks for fairness.
19. PMC (2023). Framework for AI bias detection and mitigation as a defect management process.
20. Springer Nature (2015). How algorithmic popularity bias hinders or promotes quality content in cultural markets.
21. RePEc (2025). The impact of different recommendation algorithms on consumer search and merchant competition.
22. Medium (2026). Breaking AI Bias: Tools and Practical Guide to AIF360 and Fairlearn.
23. AWS Builders (2024). Designing Search in a Microservices Architecture: Owning the right problems.
24. IJCEM (2021). Implementing Microservices Architecture in Retail Application Development.
25. Amazon Science (2026). Evaluating generative AI models with Amazon Nova LLM-as-a-Judge.
26. Amazon Science (2025). FiSCo: Making fairness in LLMs observable, quantifiable, and governable.
27. AWS Machine Learning Blog (2025). Minimize generative AI hallucinations with Amazon Bedrock automated reasoning checks.
28. eBay Innovation (2024). Evolving recommendations: A personalized user-based ranking model and the Krylov GPU cluster.
29. eBay Innovation (2024). How eBay created a language model with three billion item titles using MicroBERT and Knowledge Distillation.
30. Medium: eBay Tech Blog (2022). Multi-relevance ranking model for similar item recommendation.