

Disentangling and Operationalizing AI Fairness at LinkedIn

Joaquin Quiñonero-Candela*
joaquin@linkedin.com
LinkedIn
USA

Yuwen Wu
yuwwu@linkedin.com
LinkedIn
USA

Brian Hsu
bhsu@linkedin.com
LinkedIn
USA

Sakshi Jain
sjain2@linkedin.com
LinkedIn
USA

Jen Ramos
jramos@linkedin.com
LinkedIn
USA

Jon Adams
jnadams@linkedin.com
LinkedIn
USA

Robert Hallman
rhallman@linkedin.com
LinkedIn
USA

Kinjal Basu†
basukinjal@gmail.com
LinkedIn
USA

ABSTRACT

Operationalizing AI fairness at LinkedIn’s scale is challenging not only because there are multiple mutually incompatible definitions of fairness but also because determining what is fair depends on the specifics and context of the product where AI is deployed. Moreover, AI practitioners need clarity on what fairness expectations need to be addressed at the AI level. In this paper, we present the evolving AI fairness framework used at LinkedIn to address these three challenges. The framework disentangles AI fairness by separating out equal treatment and equitable product expectations. Rather than imposing a trade-off between these two commonly opposing interpretations of fairness, the framework provides clear guidelines for operationalizing equal AI treatment complemented with a product equity strategy. This paper focuses on the equal AI treatment component of LinkedIn’s AI fairness framework, shares the principles that support it, and illustrates their application through a case study. We hope this paper will encourage other big tech companies to join us in sharing their approach to operationalizing AI fairness at scale, so that together we can keep advancing this constantly evolving field.

KEYWORDS

AI Fairness strategy, operationalization, equity, large-organizational process

*Corresponding author.

†Work done while the author was with LinkedIn.

This is the author’s version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in ACM Fairness, Accountability, and Transparency (FAccT ’23), <http://dx.doi.org/10.1145/3593013.3594075>. 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 the author(s) 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.

FAccT ’23, June 12–15, 2023, Chicago, IL, USA

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0192-4/23/06...\$15.00

<https://doi.org/10.1145/3593013.3594075>

ACM Reference Format:

Joaquin Quiñonero-Candela, Yuwen Wu, Brian Hsu, Sakshi Jain, Jen Ramos, Jon Adams, Robert Hallman, and Kinjal Basu. 2023. Disentangling and Operationalizing AI Fairness at LinkedIn. In *2023 ACM Conference on Fairness, Accountability, and Transparency (FAccT ’23)*, June 12–15, 2023, Chicago, IL, USA. ACM, New York, NY, USA, 16 pages. <https://doi.org/10.1145/3593013.3594075>

1 INTRODUCTION

LinkedIn’s products aim to live up to the company’s vision to “*create economic opportunity for every member of the global workforce*”¹. This would be impossible without leveraging AI at scale. We use AI to power many product offerings, whether to recommend new job openings to job seekers [49], show qualified candidates to recruiters [34, 65], share relevant news and articles to our members [5, 7, 8], or recommend new connections to help members grow their network [37]. Leveraging AI for a plethora of products and creating opportunities for every member means it is critical to work towards building AI that treats every member fairly.

But what does it mean for AI to be fair, or not biased? Who should decide what definition of AI fairness applies to a given product? Given a definition of AI fairness, how do we operationalize measurement and mitigation at scale? These questions aren’t rhetorical, they are motivated by substantial challenges we have found when trying to operationalize AI fairness at LinkedIn’s scale. These challenges include:

- The multiple reasonable definitions of AI fairness that are mutually incompatible, and choosing one vs. another is consequential and dependent on the product [21, 60, 64].
- Choosing a particular definition of AI fairness for a specific product requires a deep understanding of the context in which that product is deployed, an assessment of the associated benefits, risks and potential unintended consequences, and thoughtful oversight and governance [22].

Imagine a recruiter searches for qualified candidates for a job opening where gender is irrelevant to being qualified for this job. The recruiter gets a list of 100 candidates with 80 males and 20 females (we’re using binary gender in this example). Is this unfair because we didn’t get 50 males and 50 females, or something closer to a 50/50 distribution? But would 50/50 be fair if in the current

¹<https://about.linkedin.com/>

social context, of all people who are qualified for this job (e.g. have the relevant skills) only 20% are female? Should an AI powering recruiter search aim to accurately predict whether a candidate is qualified, irrespective of their gender, reflecting potential societal imbalances?² Or should the AI instead take gender into consideration and compensate for societal imbalances by increasing female representation?

One way to reason about this tension is to consider the difference between equality and equity when we think about what is fair. In our example, a principle of equality would mean treating every candidate the same, and ensuring the AI predictions correctly predict real-world qualifications as accurately for female as for male candidates. In contrast, a principle of equity would investigate whether there have been historical barriers to females acquiring the qualifications required for the job, or whether the qualifications required are too narrow and would exclude females capable of performing that job as well as males. An equity principle may require having a female representation in the search results that is higher than the 20% baseline. It is important to note that a female representation higher than 20% would require treating female candidates differently than male ones explicitly based on their gender, and therefore would be incompatible with a principle of equality (or equal treatment).

The tension between equality and equity is mirrored in the inherent tradeoffs between AI fairness metrics like predictive parity (or calibration) which aligns with equal treatment, and demographic parity (and to some extent equalized odds) which align with equitable outcomes, see [51]. AI practitioners working on AI fairness are faced with a very difficult choice: what definition of AI fairness to use for measurement and mitigation? This decision boils down to making a consequential tradeoff on the spectrum between equal treatment and equitable outcomes. We take the position that:

- (1) AI practitioners cannot be expected to have all necessary understanding of the context in which a given product is deployed, and therefore should not be put in a position to make consequential AI fairness tradeoff decisions.
- (2) Addressing AI fairness comprehensively requires acting both on the AI components and also on the broader product powered by the AI, including product goals, policies and design as part of a comprehensive product equity strategy.

We therefore propose to disentangle AI fairness by clearly separating out equal treatment and equitable outcome expectations, so that both can be addressed explicitly, and so AI practitioners have clear guidelines for operationalizing AI fairness at scale. Our framework can be summarized with this simplified equation:

$$\text{AI Fairness} = \text{Equal AI Treatment} + \text{Product Equity}.$$

This approach demands of the AI component that it treats everyone equally (we will detail what we mean by equal AI treatment in the remainder of the paper), but it also acknowledges that equal treatment is not sufficient to ensure equitable product experiences and outcomes. Hence, our AI Fairness approach requires a complementary product equity strategy that considers additional efforts

focused on understanding potential harms or barriers that disproportionately impact certain groups, as well as potential investments that can make the product accessible and beneficial to all.

In the example above, an AI that satisfies equal AI treatment will not address the social problem of female under-representation in the pool of qualified candidates, and as a result, females will still be underrepresented in the search results. A complementary product equity strategy will for example invest in awareness and actionable suggestions for recruiters to increase the diversity of their candidate pool, such as “diversity nudges” that let recruiters know females are underrepresented in their search, and suggestions to consider additional relevant skills, and broader or different search criteria [3]. Another complementary product equity strategy would be to re-rank the candidates so every slate of candidates (say the 20 candidates that fit in each of the 5 pages that together contain all candidates) matches the proportions of the underlying qualified candidate distribution [34].

In this paper, we focus on the equal AI treatment component of our AI Fairness framework and share the principles that support it. While we emphasize where a product equity strategy needs to complement equal AI treatment, we aim to present details of our product equity strategy in future work. Meta and Google have shared aspects of their approaches to AI fairness in Bakalar et al. [13] and Beutel et al. [16] respectively. Like our work, Meta’s approach offers a comprehensive end to end framework for thinking about AI fairness. Although it also builds on a notion of predictive parity, it focuses on a binary classification fairness criterion around the decision boundary. Google’s approach clearly articulates a measurement and mitigation approach, but it does not capture the full process of evaluating and selecting the appropriate mitigation strategy, or potential complementary product changes. Like Meta, it focuses on the binary classification case, but with a different definition of fairness, “conditional equality,” that extends equality of opportunity and therefore aims to close false positive rate (FPR) gaps. What distinguishes our work is that unlike in Meta or Google’s case, we offer a justifiability framework for AI unfairness mitigation and demonstrate its application on a real-world LinkedIn product in Section 4.

We recognize that AI fairness is an evolving and complex area, with no consensus on definitions, goals, or mitigation strategies. We expect that our principles and our implementation strategies will develop over time, and we hope that an increasing number of large public companies will also share their AI fairness strategy. The rest of the paper is organized as follows. We first present LinkedIn’s principles for equal AI treatment in Section 2. Section 3 then details the operationalization guidelines that allow the applications of these principles in practice. We illustrate the principles and their application through a real-world case study and share learnings and results in Section 4. We end with a discussion in Section 5.

2 EQUAL AI TREATMENT PRINCIPLES AT LINKEDIN

When developing principles for equal AI treatment at LinkedIn, we started by disentangling AI fairness by separating equal treatment and equitable outcome expectations. We set equal treatment expectations as the default bar for the AI component of our products so AI practitioners have clear guidelines that can be operationalized

²Note that in this thought example, the AI is an assistive tool for the recruiter, it doesn’t make hiring decisions, the recruiter does.

at scale. We then addressed the important challenge that equal AI treatment is not sufficient for achieving equitable outcomes and made sure the principles reflect the need for a complementary product equity strategy. Finally, recognizing that AI fairness is a fast evolving field, we committed to evolving our AI fairness strategy transparently and with the input of external stakeholders. These are LinkedIn’s Equal AI Treatment principles:

- (1) **We will measure and work to mitigate algorithmic bias so that our AI systems treat everyone equally.** This means measuring and, where appropriate, mitigating systematically unequal predictions or errors affecting member demographics. In other words, similarly qualified members should receive similar opportunities irrespective of what group they belong to. This principle does not extend to interventions intended to create equitable product experiences or outcomes.
- (2) **We will not consider equal AI treatment the end of our work but will treat it as the foundation of a broader fairness and equity strategy.** Additional measures – such as product features and design changes – and coordination between product and AI teams are key to achieving LinkedIn’s vision of “creating economic opportunity for every member of the global workforce”, and to addressing unintended consequences of our equal AI treatment efforts.
- (3) **We will validate our approach externally and lead with transparency in this developing field.** Across our equal AI treatment and our broader equity and fairness initiatives, we will share learnings and case studies, and we will leverage collaboration and feedback from members, customers, advocacy groups, and social scientists. We aspire to learn from and guide other organizations in delivering equal AI treatment and exploring broader equity initiatives.

3 FROM PRINCIPLES TO PRACTICE

The three principles above are the foundation on which we develop practical strategies to make sure the AI in our products treats every member fairly. In this section, we share the operationalization guidelines for the principles as a template that’s applicable across all AI-powered products at LinkedIn.

3.1 Principle 1: We will measure and work to mitigate algorithmic bias so that our AI systems treat everyone equally

3.1.1 Definitions and Measurement: Equal AI treatment means that AI predictions treat people the same way, from a statistical perspective, irrespective of their demographic group membership. Conversely, algorithmic or AI bias means that people are not being treated the same way across groups. This could mean that the accuracy of predictions is systematically worse for a particular demographic group, or that systematic errors benefit some groups and disadvantage others. Formally, we define and measure equal AI treatment in different AI systems slightly differently depending on the context, but it relies on appropriate and accurate demographic group data in all cases.

- When AI allocates scarce opportunities, equal AI treatment by default means “equal opportunity for equally qualified

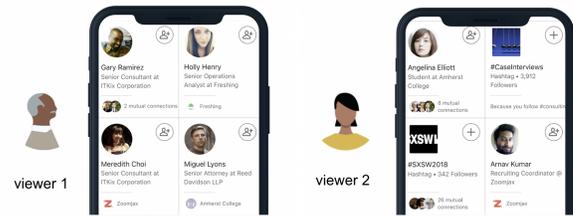


Figure 1: Viewers are members who get to see the recommendations. Here viewer 1 (on left) and viewer 2 (on right) see different recommendations. Equal AI treatment for them is through the standard quality of service. Recommended members are those that are being shown to viewers. For these recommended members equal AI treatment means equal opportunity for equally qualified and is measured by predictive parity.

candidates”. Mathematically, we define it as predictive parity, where equally qualified candidates get similar predicted scores irrespective of what group they belong to. Formally, this is expressed as, $Y \perp\!\!\!\perp A | \hat{Y}$ where, \hat{Y} is the prediction score, A is the random variable signifying the demographic group membership and Y is the outcome of interest. For details see [14, 23]. Product examples here include members who are ranked in People You May Know (PYMK) [9], Recruiter Search [34], as well as content being ranked in the Newsfeed [8], where we are considering fairness with respect to content creators. We also differentiate how predictive parity is computed depending on whether the AI model does pointwise [23] or listwise inference [66].

REMARK 1. *Other approaches to algorithmic fairness include demographic parity and equalized odds, amongst many others [14]. It is also very well known that the above three crucial measures are at odds with one another [19, 39, 44, 51]. While there is not yet an industry consensus, our approach to choosing predictive parity is consistent with the general standard for AI fairness outlined in Microsoft’s Responsible AI Standard [4].*

- When AI does not allocate scarce opportunities, “equal AI treatment” by default means “standard quality of service” and is mathematically defined as a comparison of relevant model performance and product metrics computed across the identified demographic groups to an absolute minimum standard. Product examples include members viewing job recommendations or searching for jobs, members viewing content on their feed, viewing other members to connect to grow their network, etc. See Figure 1 for more details.

In adopting this definition, we recognize that “equal AI treatment” does not focus on, or guarantee, equitable outcomes as detailed in Section 1. That being said, we are still choosing this as our starting point to be able to clearly disentangle the equal treatment approach from equitable outcomes. We understand that choosing this measure can potentially only maintain the status quo in society, and can even be unintentionally detrimental to underrepresented populations, i.e. members who were not getting equitable outcomes. But that is a *conscious choice*. Maintaining equal AI treatment as

defined above ensures that we are not compounding an inherent bias in society with our algorithms. This separation is crucial for solving the much larger societal challenge. If our algorithms enforce equal AI treatment, we can bring forth product changes toward achieving equitable outcomes, while ensuring the AI component isn't introducing any unfair bias. We discuss strategies around this in detail in the next section.

We also recognize that for a global service like LinkedIn, there are an unlimited number of demographic groupings (such as race/ethnicity, gender, disability, religion, ...) and sub-groupings (e.g. specific racial or ethnic identities) that could be leveraged for delivering "equal AI treatment," and significant challenges to address, including a lack of accurate data on membership in relevant groupings, regional differences in group definitions or prioritization, evolving norms of group identity, and the difficulty of harmonizing efforts to apply equal treatment across a wide range of interconnected groupings. We don't view these challenges as a reason not to act based on the data available, but we expect to expand and develop our approach to demographic groupings over time. As a first step, we have launched the Self-ID initiative [1, 25] to be able to collect such information (see Figure 9 in the appendix) in an ongoing manner in order to make LinkedIn a more equitable platform.

REMARK 2. *Note that AI can also be used without the context of a member for example in generative models or text-based content moderation. In such situations, the mathematical definition of equal AI treatment can differ and be highly context specific. We do not go into such details here due to brevity. Please see [47, 69, 71] and [33] for details.*

3.1.2 Mitigation via a Justifiability Framework: For non-demographic groups, if a model is found to be performing poorly across particular groups, the standard approach is to simply add the particular group membership as a feature. However, there are multiple reasons why this approach may not be appropriate or even feasible when considering sensitive demographic data like binary gender, including privacy concerns about data access and usage, and possible unintended consequences of including demographic information in models. Thus even if it is mathematically possible to define mitigation strategies for enforcing predictive parity [23], we must take a cautious approach since from a legal and policy perspective [73, 75], there aren't clear processes to follow that rely on the use of demographic data.

We present a justifiability framework that we propose to follow to determine whether there are means to mitigate algorithmic bias, first, without using demographic information. If not, then we deem that the use of demographic information is justifiable (albeit subject to appropriate guardrails, as a last resort, and only if the benefits outweigh the risks). The framework is not meant to be exhaustive or necessarily include theoretical guarantees; instead, the goal is to create a standard set of questions to ask and investigations to be performed. The major steps of the framework are highlighted below:

- (1) **Root-Cause Analysis:** Preparing a comprehensive root cause analysis checklist to determine, to the extent feasible, why we observed a prediction parity gap or a difference in the quality of service. We do not want to blindly mitigate the

situation without understanding why it happened in the first place. For some initial work in this area please see [11, 35].

- (2) **Mitigation strategies without demographic information:** We consider alternative equal treatment implementations without using demographic information that results in a smaller predictive parity gap or difference in the quality of service. Examples include collecting more and higher quality training data, for example, to create a more balanced representation of diverse groups [38, 45], training larger capacity models, or finding relevant features that are confounding factors of demographic group membership, etc. Please see [40, 53] and the references therein for more state-of-the-art methods in reducing bias without demographic information.
- (3) **Mitigation using demographic information:** In many instances it may not be possible to close the gap without using demographic information [12, 27]. In such cases, we adopt a least granular intervention approach where mitigation requires using demographic data at prediction time (e.g. to calibrate a model [15, 41, 63]). Accordingly, we will generally prefer post-processing [50, 55, 62] to in-processing [74] approaches due to their interpretability and accountability of what exactly is happening under the hood.
- (4) **Unintended consequences:** We will assess potential product risks of the proposed AI fairness mitigation (especially when using demographic information), and quantify those risks through experimentation [31, 52, 67]. Only when the benefits clearly outweigh the risks will we implement the mitigation for all members of the platform.
- (5) **Mitigation Guardrails:** Mitigations to achieve equal AI treatment will not involve affirmative action or interventions seeking equitable outcomes for particular demographic groups. We consider such interventions as better suited to product equity initiatives for several reasons.
 - (a) Pursuing equitable outcomes requires qualitative user research, a deep understanding of product features and how they affect different demographic groups, assessment of benefits and risks, and close governance and oversight.
 - (b) AI models can change with time and can suffer from a lack of monitoring or maintenance, leading to potential unintended consequences of affirmative action interventions within a model.
 - (c) The legal and social context in which a product exists—as well as the equity challenges and priorities of a product—can evolve over time, demanding that any affirmative actions at an AI level be kept fully synchronized with those changes, which represents a significant operational challenge.

We highlight each of these mitigation steps in detail in Section 4 where we deep-dive into a particular case study.

3.1.3 Privacy and Security: One of the major steps in being able to accurately measure and potentially mitigate algorithmic bias relies on access to accurate group information. Although we at LinkedIn began asking for gender and disability demographic data globally many years ago, we realized that it did not encompass the broad spectrum of identities of our members. In 2021, we introduced the "Self-ID" initiative to expand on those efforts in the U.S. [1, 25].

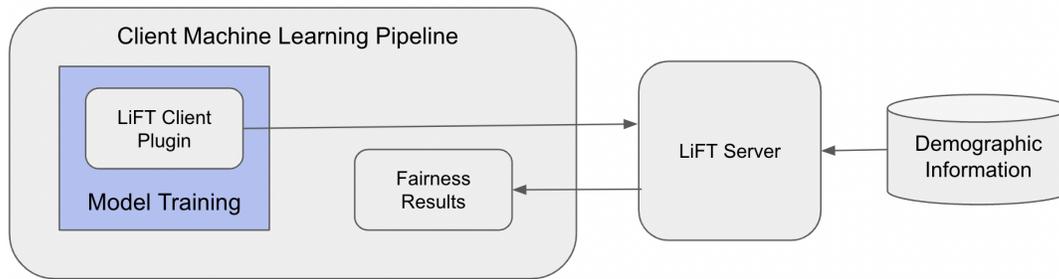


Figure 2: Client-Server architecture which allows a model owner to evaluate the fairness metrics for their model without needing access to demographic information data. LiFT here refers to the LinkedIn Fairness Toolkit [72] which houses the measurement algorithms.

When members choose to provide this data, we handle this data under strict privacy and security measures. By default and unless they have a legitimate business need, LinkedIn engineers do not have access to this data in an unencrypted format. However, the same protections used to safeguard this sensitive data from unauthorized use or disclosure, also create obstacles in using it for measurement of algorithmic bias across our relevant products. To address this challenge, new and innovative work is underway both from an engineering and an algorithmic aspect to be able to leverage this data without undue risk. Initiatives such as privacy-preserving machine learning [10, 56, 77], homomorphic encryption [30, 70, 78], A/B testing under differential privacy [32], etc., are all being studied across various teams.

Given one of our focus areas is group-based measurement, it may be possible to use differential privacy [26, 28, 29] to obfuscate the true group membership information, but still be able to estimate predictive parity differences with some degree of accuracy. We can mathematically derive an unbiased estimate of the fairness metric by knowing the noise that is being added to the group information data [17, 29]. Moreover, for an extra layer of protection, none of the product-focused AI teams would need access to this data for evaluating their models. We have developed the pipelines through a server-client architecture (see Figure 2) such that each model owner would be able to evaluate their model for fairness without ever needing access to this information [58].

Although it might be possible to use demographic data for measurement purposes, the problem is much more complex and nuanced for mitigation. The safest mitigation technique is when we are able to mitigate a predictive parity gap without needing access to demographic attributes. However, in instances where that is not possible, we do plan to use demographic information in order to mitigate our models, albeit under limited use cases and with protections designed to ensure limited access for this use case.

An obvious question arises, why not use obfuscated information (similar to the differentially private approach for measurement) in doing mitigation? The answer to that lies in being very conscious about our choice. Our post-processing mitigation solution [23] for enforcing predictive parity changes the score based on group membership. In the case where we are using noisy group information data, we would know for certain that for those fraction of noisy cases, we have changed their score in a wrongful manner. At the end of the day, these score changes happen at an individual

member level and we strongly believe in the fact that **we should not knowingly give a wrong score leading to undeserving opportunities.**

Overall, operationalizing this framework requires various engineering aspects to ensure that we adhere to the relevant laws and our privacy commitments, including but not limited to internal controls regarding employees who have access, reducing access to only when necessary, safeguards to prevent data misuse and system infrastructure such that equity can only be achieved without giving internal stakeholders direct access to such private information.

REMARK 3. Data availability: We may not have the demographic data for many of our members, which makes it difficult to apply the post-processing techniques. We are currently researching new methods that may work without the strong dependency on demographic data. For details please see [43].

3.2 Principle 2: We will not consider equal AI treatment the end of our work, but will treat it as the foundation of a broader fairness and equity strategy.

We consider equal AI treatment, as a necessary foundation, but not itself sufficient for achieving equitable product experiences or outcomes. For example, equal AI treatment would not address a group’s under-representation in recruiter search results or connection recommendation if that under-representation stems from a real-world under-representation in a given sector. Achieving the broader equity goals requires that equal AI treatment be complemented with a product equity strategy and in partnership with customers who choose to pursue equity goals. Other measures – such as product features and design changes – are outside the scope of AI initiatives and require different inputs and assessments.

We have already started our journey toward achieving equitable product experiences. The first step of this journey undoubtedly starts with understanding how our members experience our platform, especially those from historically and systemically marginalized communities. Through the Self-ID initiative, more than 10 million members [24] have shared some aspect of their identity on LinkedIn. As more members join us on this equity journey and Self-ID, we’ll be able to launch new products and experiences to help drive more equitable outcomes for those members facing barriers. One such product feature that we have already shipped is “diversity nudges” in LinkedIn Recruiter and LinkedIn Learning [3].

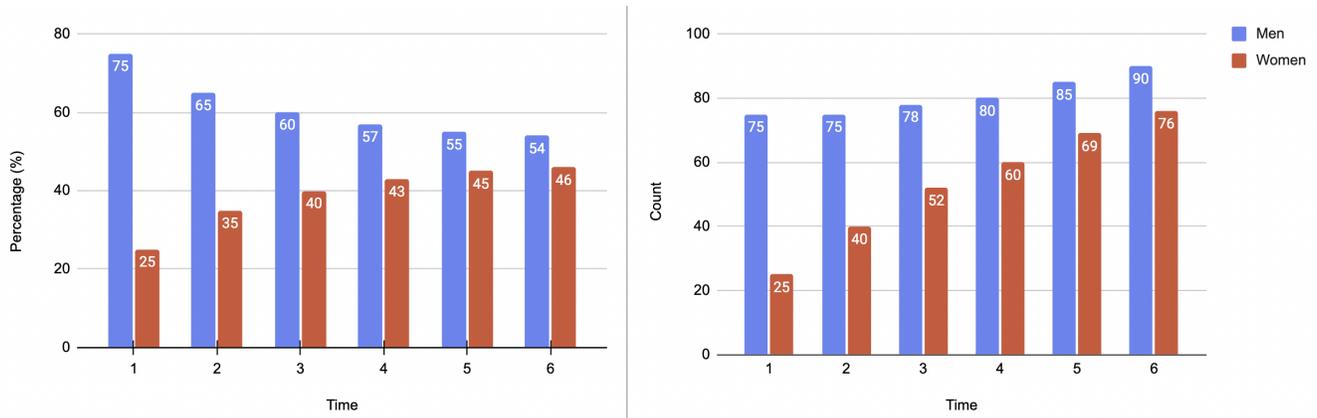


Figure 3: A hypothetical situation over time where we start with an initial distribution of 75% and 25% and through exploration ultimately converge to 54% and 46% (Figure on left). Note that in many situations it may not be possible to reach an absolute equal distribution. At the same time, we are never sacrificing on the total number. The figure on the right shows how their outcomes are growing over time thereby bringing us to a more equitable future. Moreover, throughout this time frame, we are ensuring that equal AI treatment holds for the members in situations where we do not have active exploration.

These features allow recruiters and hiring managers to increase gender representation in their qualified candidate pools. We also have a feature to hide names and photos of candidates to reduce unconscious bias [2].

Beyond having a product equity strategy, it is also crucial to understand the potential unintended consequences that interventions aimed at ensuring equal AI treatment can have. Given how tightly interconnected AI and product outcomes are, we take an end-to-end view. Any AI change, although only targeted at ensuring equal AI treatment, can have a strong impact on product metrics. As a result, it is necessary for product teams to commit to supporting equal AI treatment mitigation efforts and monitoring outcomes against product equity strategy and goals over time. For example, addressing a prediction gap may result in increased visibility or participation for a given group (e.g. female profiles) with attendant risks of harassment which must be assessed.

3.2.1 Towards Equitable Outcomes via AI. Note that although we stated that AI systems are expected by default to meet equal AI treatment, it is also possible to devise AI changes with the explicit aim of achieving equitable outcomes. Such interventions are thought of as product design changes that leverage AI rather than direct AI interventions aimed at equitable outcomes. An example would be a product that reserves a portion of the screen real estate to show underrepresented people or content, using AI to retrieve the highest ranked members or items of that group. This is AI serving a product design decision driven by an equity strategy, rather than a modification of the AI algorithm to satisfy an equity strategy. This conceptual separation is crucial to be able to structure and adhere to the principles.

As we mentioned earlier, the notion of equal AI treatment by enforcing predictive parity can reflect and doesn't disrupt the status quo in society. It is a valid criticism that in a world where predictive parity holds, we might not be giving the deserving opportunity to some members only because their group in the past, did not engage with the platform. More formally, if there are groups of individuals

who rarely visited or engaged with the platform (thereby mostly having $Y = NA$), enforcing predictive parity would give them a very low score ($\hat{Y} \approx 0$) and hence they would either be very low in the ranked list or not shown at all. As a result, if there is a member who is highly qualified (with a potential of getting $Y = 1$)³, but belonged to this group, enforcing predictive parity will not give the deserving opportunity to this member.

We believe such situations should be solved through an explore/exploit strategy [6, 54] as part of the product equity strategy and goes beyond the equal AI treatment notion. Consider a situation in time, where our models have enforced predictive parity. We now do not know if there are members who might be qualified but are not being shown due to our enforcement of predictive parity. At this stage, we consider a budget for exploration. Let's assume that now for 95% of all member sessions, we are maintaining predictive parity, while for 5% of sessions, we are actively trying to explore members who we think might be qualified, but we lack the data to accurately predict it. There are several ways of doing this, for details please see [20, 46, 48]. Based on this exploration data, if we can truly find some of the missing qualified members, then their data is automatically passed into the training model and over time, the model understands that these are qualified members and automatically shows them without the need for artificial exploration. Overall through these iterations, the main goal is to maintain equal AI treatment for the majority of the sessions (say 95%) and slowly improve the actual metrics toward achieving equitable outcomes.

As a concrete example, consider a case where an AI researcher sees 75% men and 25% women in their connection recommendations and the ranking algorithm satisfies predictive parity. Through the explore/exploit paradigm, our aim would be to improve this and potentially make it 65% men and 35% female or even better in the long run, while maintaining predictive parity. Figure 3 (on the left) shows the cycle of improvement over time. It is important to note

³Unless specifically mentioned we are usually working in the regime of binary classification.

that in the above example, recommended members being viewed by this AI researcher meet the feature criteria to be recommended. Thus, we are never replacing a qualified member with someone who is unqualified. Moreover, we are only working with percentages above and not raw counts. In an ideal situation, we would expect the total count to grow over time. For example, if we initially started with 75 men and 25 women, we can call our strategy a success if we now have 75 men and 40 women, raising the total count but reducing the percentage gap. See Figure 3 (right).

We believe that tracking these proportional and count statistics can give insight into the long-term effects of the explore/exploit strategy and indicate if we are reinforcing the status quo or moving towards a more equitable environment. Still, there are potential pitfalls in this strategy that our random sampling-based exploration does not explicitly consider. For example, it is possible that the underrepresented group is simply harder to learn or lacks samples (even with randomization). Hence, if we reach a male-to-female 70/30 steady-state proportion, it is not clear if this is the true long-term behavior or if it is because a model lacks the capacity to adapt to the newly gathered samples. One method for addressing this could involve ranking based on confidence intervals of the expected outcomes $E[Y|S]$ rather than the point estimates as suggested by Salem and Gupta [68], which could better account for model uncertainty and the potential for a candidate to rank higher, rather than the point estimate. In a similar vein, bandit-based approaches that utilize the confidence intervals can pose a more targeted strategy than random sampling. Lastly, if there are a-priori contextual reasons that a group should have higher outcomes than they currently do, it is also possible to empirically test this hypothesis by incrementally boosting the groups' scores and running a post-hoc analysis to determine if they have been under-scored. Learning under-selection bias remains an open research problem.

Overall, such AI strategies, along with specific product launches are also part of our overall equity strategy toward achieving equitable outcomes. The prior work on generating a representative list in talent search [34] also falls into this framework as an example of an AI intervention targeted towards equitable outcomes, where the authors tried to match the underlying gender distribution of a search query to the ranked results.

3.3 Principle 3: We will validate our approach externally and lead with transparency.

AI fairness is a developing field where LinkedIn is well-positioned to make significant contributions. We aspire to help and inspire other organizations as they seek to deliver equal AI treatment and explore broader equity strategies. We will share our learnings through case studies, articles, and other presentations, and we will leverage input and collaboration from our members, our customers, academia, industry partners, and advocacy groups, as well as Microsoft's RAI and social science resources. We will explore approaches to algorithmic auditing by vetted external expert groups with the goal of delivering an external quantitative assessment of our measurement and mitigation approach.

4 CASE-STUDY

This section illustrates how the methods described in Section 3 are applied in a real-world example. We focus on a people recommendation product at LinkedIn called "People Also Viewed" (PAV, see figure 4a and the appendix). This product recommends additional profiles you might be interested to connect with or learn from based on the current profile you are viewing and is an example of the scarce resource allocation paradigm (limited number of recommended profile slots). The product is the second largest source of traffic for profile views, therefore any negative bias (failing to surface qualified recommendations) could adversely affect members' abilities to grow their network. On the other hand, 'over-recommending' members can lead to harm such as unwanted connection requests and spam. In this case study we focus on measuring and mitigating bias for binary gender. Equal AI treatment in this case means that the model that recommends related profiles for PAV produces ranking scores that satisfy predictive parity across binary gender.

4.1 Measurement

The motivation for this use case is an observed gap in predictive parity between male and female members⁴. Figure 4b illustrates the real-world impact of lack of predictive parity. Consider a group of female and male members that share a similar AI relevance score: the females in that group will see higher real-world outcomes of actual network formation. This means that the score should have been higher for the females in that group. The model is therefore under-predicting for females. Meeting equal AI treatment would require the two curves in 4b to overlap.

The remainder of this Section details our approach to mitigating this predictive parity gap. Although we focus on binary gender, the framework can be extended to other demographic groups, assuming there is sufficient data.

4.2 Mitigation via Justifiability Framework

We follow the steps in Section 3 to mitigate the observed gap in predictive parity.

4.2.1 Root-Cause Analysis. Our framework provides a standard set of investigations that can provide insights into the mechanisms leading to a (gender) gap⁵.

- (1) **Distribution of binary gender:** The obvious first step is breaking down the data to see the ratio of binary gender; in the most trivial case, if one gender is missing from the data, it would explain why the model is miscalibrated for this group. In the less extreme case, binary gender skew could still contribute to bias if the relevant signal for one gender is not large enough for the model to learn predictions for each gender. In the case of PAV, we did not see a large gender skew for the source or destination members (see Table 1).⁶
- (2) **Feature justifiability:** Building AI models responsibly requires that we only include features if they are justifiable.

⁴We did not observe a viewer-side gap in this product since the model does not use viewer-side features.

⁵Note that these methodologies can be applied to general model improvement, and are not limited to just closing a fairness gap.

⁶There was a 29pp difference between men and women for the viewer gender, but given the model did not use viewer features at all, we did not see mitigate on that.

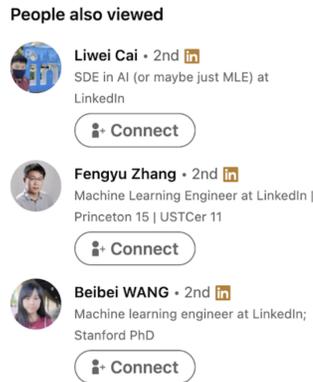


Figure 4(a): The People Also Viewed product that shows a ranked list of recommended members.

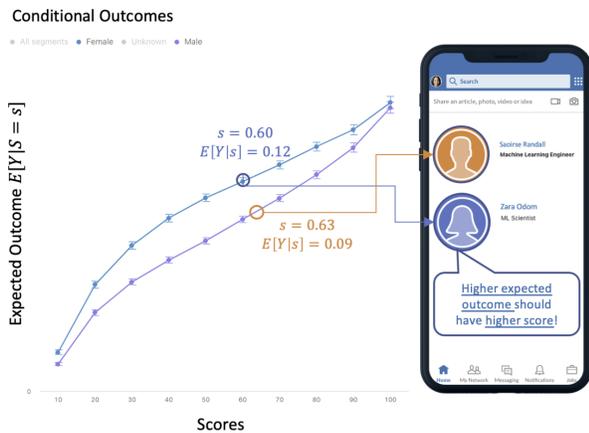


Figure 4(b): How predictive parity gap changes the ranking order in products such as People Also Viewed.

	Men	Women
Source	54.50%	45.50%
Destination	54.70%	45.30%

Table 1: The distribution of gender across different member groups. The results are normalized to binary gender. Unknown, non-disclosure and non-binary gender are not shown.

Although we do not give a formal definition, the use of a feature is justified, for example, if the feature is relevant to the modeling task. Domain knowledge and common sense are required to judge that relevance in context.

The justifiability bar is highest for demographic information or other highly confidential data. In addition to privacy and security concerns, including these data in a model could exacerbate biases or lead to other unintended consequences. This is why we ask that the unjustifiable features be removed even if model performance drops. Removing a feature that isn't justified doesn't mean all other features correlated with it also need to be removed. On the contrary, as we discuss in Section 4.2.2, one way to close the predictive parity gap is to

add relevant correlated features on a causal chain. An audit of the PAV model features did not reveal any unjustifiable features.

- (3) **Cohort-level Error Analysis:** We have observed that PAV is miscalibrated when segmenting on binary gender. However, the metric is a population-level average, meaning we do not have insight into the level of miscalibration for different subgroups. To get a more granular understanding of what is causing this binary gender-based gap in predictive performance, we perform a cohort-level analysis.

We build a tree model to split the data using a user-defined error metric. The splits are determined such that the partitions of the data maximize the differences in the error metric. The resulting cohorts and the features used in the splits can give more granular insight into model performance compared to standard feature importance methods [42], as the data is automatically segmented into high and low-error cohorts. Furthermore, we can include non-model features as candidates in the cohort splits, and we can also use a non-model error metric as the splitting criterion [76].

For the fairness use case, we set the metric as the residual between the label and the predicted score. This is a proxy for miscalibration error, which is essentially what predictive parity tries to measure. As the goal is human-interpretable cohorts, we limit the depth of the tree to be three, so that only eight cohorts are generated. For candidate features, we include the gender of any member in the model (for PAV, these would be the viewer, current profile member, and recommended member) in addition to the actual model features. Figure 7 (see appendix) shows the result of running the error cohort analyzer on PAV data. The error cohort model used the binary gender of both the recommended member and the current profile member as top feature splits, with MALE cohorts generally having lower residuals compared to FEMALE cohorts. This suggests binary gender could be a driving factor for the measured bias, increasing the likelihood that gender will need to be directly used for mitigation.

4.2.2 Mitigation experiments without demographic information.

Mitigation of fairness violations is typically studied in contexts that allow access to demographic information during model training and inference. However, a key component of our justifiability framework is our proposition that demographic data should only be used when other methods are demonstrably inadequate, and when we were unable to find negative unintended consequences from mitigating using demographic data. We therefore evaluated a wide range of methods for mitigation without demographic data at inference time, from feature-selection methods to fairness-constrained in-processing methods. This section provides an overview of these techniques, and presents our rationale for ultimately choosing to mitigate with demographic data.

The problem of closing the gender calibration gap without using gender as a feature can be cast as one of identifying missing features that are relevant and correlate with gender. We motivate this approach by visualizing gender and model features in a causal graph [61], where the terminal node is model prediction. Our goal is then to determine if there are additional, non-gender features

that could ‘block’ the effects of gender on the gap such that the gap is reduced.

Let’s consider a toy model with only two features; we show the causal relationships between the features, model predictions, and binary gender (which is not included as a feature) in a directed acyclic graph (DAG) (see Figure 5a). The original model is biased because binary gender is causally related to the predictions through a direct path. Imagine we identify a new feature, click probability, that is a descendant of binary gender and fully captures its effect. When click probability is included in an updated model version, binary gender no longer has a direct effect on the model predictions (see Figure 5b). We can therefore theoretically close any gender-based bias by including missing features in the direct causal path of binary gender. A caveat of this approach is that our toy example

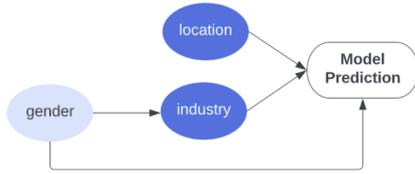


Figure 5(a): Original model: Gender affects model prediction both through a direct path as well as an indirect path through industry

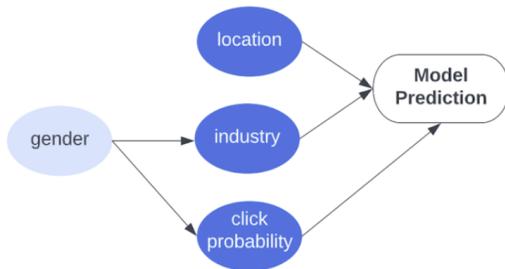


Figure 5(b): New model: Click probability identified as a new feature and Gender affects model prediction only through indirect paths.

only shows the ideal case where the new feature closes the gap and does not have additional interaction effects. In reality, any new node in a large graph can drastically change the DAG by adding new edges and therefore changing causal relationships [36]. New features in the direct path of binary gender may even widen the gap, therefore we cannot indiscriminately add features we think are correlated with binary gender and outcome.

One approach for a more intentional mitigation strategy is to leverage domain expertise to identify relevant missing features; however, it is difficult to scale these ad-hoc solutions. We therefore first created a superset of common LinkedIn member features and tested automated feature selection algorithms to pick the subset that could close the gap. The two main strategies we tested were Quantile Prediction Drift (QPD) and a method based on causal effect decomposition [18]. We also tested more ‘standard’ techniques such as brute force search of the feature space and imputation of missing feature values (identified as another possible issue from root cause analysis). Finally, we ran experiments including binary gender in the model to get a direct comparison between gender-aware and gender-blind techniques. The results are outlined in Table 2 and we

observe that adding generic LinkedIn features showed no reduction in the fairness gap, regardless of the feature selection strategy. Gender-aware approaches consistently outperformed gender-blind options to close the gap.

While we have focused on data-based mitigation methods, there are also other avenues of mitigation. Some examples include tuning model parameters to optimize for calibration, training with calibration-motivated loss functions, or post-processing with group-agnostic calibrators. As part of our efforts to identify effective demographic data-restricted methods for mitigation, we conducted a wide survey of methods that varied in both the point of intervention (pre-/in-/post-processing) as well as the volume of demographic data required (e.g., only in training vs. in training and inference). We present this survey in a separate paper, see Hsu et al. [43], but note here that the overwhelming conclusion is that using gender in post-processes was by far the most effective strategy for achieving predictive parity fairness.

Given that our thorough experimentation with gender-blind and gender-limited techniques did not lead to successful mitigation, we felt **justified in using the post-processing gender-based calibration** to mitigate the algorithmic bias in PAV.

4.2.3 *Mitigation experiments with Demographic Information.* We use the post-processing technique detailed in DiCiccio et al. [23] as our bias mitigation training (BMT) algorithm which requires the use of demographic information. Intuitively, BMT fulfills predictive parity by setting $E[Y|\hat{Y} = s, A = a] = s$ for all groups A (here s is the predicted score from the model). As shown in Table 2, BMT (row 1) successfully reduced the gap in offline experiments.

Based on our root cause analysis and offline mitigation experiments, we concluded that BMT was effective (offline) and should be experimented online. We launched an online A/B test comparing the baseline PAV model with a BMT-mitigated version (Figure 6a and 6b) that also showed a reduction in the predictive parity gap. Furthermore, the BMT model showed significant 2% lift in total profile actions (e.g., clicking the recommended profile), which is the main product metric, confirming our hypothesis that equal AI treatment need not sacrifice general model performance.

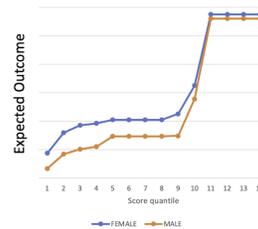


Figure 6(a): Online results of base model

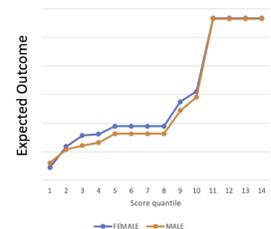


Figure 6(b): Online results of Bias Mitigated Model.

4.2.4 *Measuring unintended consequences.* Ensuring equal AI treatment in our case study requires that the gap between AI scores attributed to males and females be closed. Relative to male scores, female scores will now be higher. The product consequence of this is that a higher proportion of female profiles will be shown on the PAV product. Before we can ship this equal AI treatment intervention, aimed at eliminating predictive bias, we need to make sure there are no negative real-world consequences. One potential

Mitigation strategy	AUROC (%)	Predictive Parity Gap (%)	Reduction from baseline (pp)
Post-Processing gender-based calibration (offline) [23]	66.3	0	13.3
Add gender*	65	2.8	10.5
New feature superset + gender*	66.8	4.6	8.7
Imputation with gender-specific means*	65	8.4	4.8
QPD selection**	65.5	10.7	2.6
Causal effect decomposition**	65.2	11.8	1.5
New feature superset	66.8	13.1	0.18
Imputation with population means	64.8	13.2	0.05
Baseline	65.1	13.3	NA

Table 2: A condensed summary of mitigation experiments. *Not a justifiable mitigation strategy, but useful for offline comparison. **Best combination of features

risk is that females might now receive more unwanted messages, invites, or even harassment. To quantify this risk, we built edge-level (member → female) metrics that measure the rate at which females accept connection request, respond to messages received and report unwanted connection requests. We compared the equal AI treatment mitigation described above ("treatment") to the model before mitigation ("control").⁷

Metric	ControlRate	TreatmentRate	RelativeDiff
accept	0.20275	0.234024	0.154251
reply	0.312765	0.306352	-0.020507
report	4.70e-5	4.60e-5	-0.017529

Table 3: Edge-level response rates across metrics comparing the base model with the BMT model.

Table 3 summarizes the experimental results. The treatment group accepts connection invitations at a higher rate and reports problematic invitations at a lower rate. This means higher network formation and fewer unwanted connection requests. We observe a slight decrease in messaging response rates, tolerable given the other two positive outcomes. We therefore **did not find unintended consequences that arise from ensuring equal AI treatment** for PAV.

4.3 Key takeaways from the Case Study

We were able to measure and mitigate the predictive parity gap that was observed in our PAV model. Through initial root-cause analysis, we were confident that binary gender was a driving factor in the parity gap that was observed. We tried extensive experimentation to try to close that gap without using gender information, however, none of the gender-blind techniques were successful. This led us to

⁷Technical note: Our mitigation is applied on the viewer side and thus it is possible for a recommended member to have their score treated by BMT in some PAV tabs but not others. This can lead to two-way interference effects in the measurement of edge-level metrics such as response rate. Although methods exist for AB testing in two-sided marketplaces such as [57, 59], they are an open area of research and require highly customized online experimentation implementations. Hence, while the edge-level analysis we proposed and conducted does not rigorously account for two-way interactions, we feel that it provides a reasonable heuristic for gauging these treatment effects.

justifiably use gender to close the gap as a last resort through post-processing. Before shipping the mitigation, we measured to ensure the mitigation resulted in the product working better and without unintended consequences for the under-served female group. Fairness is a process, and mitigations need to be reevaluated and revisited periodically going through the steps of the justifiability framework.

5 DISCUSSION

We have presented LinkedIn’s framework and principles for equal AI treatment. The framework disentangles the definition of AI fairness and brings operational clarity on what AI practitioners are accountable for. Recognizing that AI fairness isn’t only an AI problem, a robust product equity strategy needs to complement equal AI treatment.

Mitigating AI fairness deficiencies suffers from an additional level of complexity: whether or not it is justifiable to use demographic information to close AI bias gaps, and how precisely to use this data. We have also presented a “justifiability framework” that addresses this question. The framework requires a thorough root-cause analysis to understand predictive parity gaps and demands exploring alternative mitigations that don’t require using demographic data. If the only way to close the predictive parity gap is to directly use demographic data, the framework asks for the minimal possible intervention and for experimentation to detect possible unintended consequences and to confirm the resulting product outcomes are beneficial, and that these benefits outweigh potential risks.

An apparent paradox emerges: using demographic information for mitigation may seem like no longer treating every group the same. The opposite is true: by using demographic information we ensure the scores an AI gives to members of different groups map equally well to real-world outcomes. Equal AI treatment means equally good scores for every group independent of group membership.

Finally, fairness is never “done”. Alternative mitigations that don’t require using demographic information should continue to be explored and should be preferred. However, unaddressed bias causes real-world harm, and when considering how to mitigate, we need to take into account the cost of inaction.

We hope this framework can be leveraged broadly and we encourage other organizations to share and collaborate.

REFERENCES

- [1] 2021. Together we can improve equal access to opportunity. (2021). <https://members.linkedin.com/equal-access>
- [2] 2022. Hide candidate names and photos in LinkedIn Recruiter. (2022). <https://www.linkedin.com/help/recruiter/answer/a481559>
- [3] 2022. Improve Gender Representation in Your Candidate Pools with These Diversity Features. (2022). <https://www.linkedin.com/pulse/improve-gender-representation-your-candidate-/>
- [4] 2022. Microsoft Responsible AI Standard. (2022). <https://query.prod.cms.rt.microsoft.com/cms/api/am/binary/RE5cmFl>
- [5] Deepak Agarwal, Kinjal Basu, Souvik Ghosh, Ying Xuan, Yang Yang, and Liang Zhang. 2018. Online parameter selection for web-based ranking problems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 23–32.
- [6] Deepak Agarwal, Bee-Chung Chen, and Pradheep Elango. 2009. Explore/exploit schemes for web content optimization. In *2009 Ninth IEEE International Conference on Data Mining*. IEEE, 1–10.
- [7] Deepak Agarwal, Bee-Chung Chen, Rupesh Gupta, Joshua Hartman, Qi He, Anand Iyer, Sumanth Kolar, Yiming Ma, Pannagadatta Shivaswamy, Ajit Singh, et al. 2014. Activity ranking in LinkedIn feed. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. 1603–1612.
- [8] Deepak Agarwal, Bee-Chung Chen, Qi He, Zhenhao Hua, Guy Lebanon, Yiming Ma, Pannagadatta Shivaswamy, Hsiao-Ping Tseng, Jaewon Yang, and Liang Zhang. 2015. Personalizing linkedin feed. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 1651–1660.
- [9] Parag Agrawal, Ankan Saha, Yafei Wang, Aastha Nigam, and Eric Lawrence. 2020. Building a heterogeneous social network recommendation system. (2020). <https://engineering.linkedin.com/blog/2020/building-a-heterogeneous-social-network-recommendation-system>
- [10] Mohammad Al-Rubaie and J Morris Chang. 2019. Privacy-preserving machine learning: Threats and solutions. *IEEE Security & Privacy* 17, 2 (2019), 49–58.
- [11] Kiana Alikhademi, Brianna Richardson, Emma Drobina, and Juan E Gilbert. 2021. Can explainable AI explain unfairness? A framework for evaluating explainable AI. *arXiv preprint arXiv:2106.07483* (2021).
- [12] McKane Andrus, Elena Spitzer, Jeffrey Brown, and Alice Xiang. 2021. What We Can't Measure, We Can't Understand: Challenges to Demographic Data Procurement in the Pursuit of Fairness. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. 249–260.
- [13] Chloé Bakalar, Renata Barreto, Stevie Bergman, Miranda Bogen, Bobbie Chern, Sam Corbett-Davies, Melissa Hall, Isabel Kloumann, Michelle Lam, Joaquin Quiñero Candela, et al. 2021. Fairness on the ground: Applying algorithmic fairness approaches to production systems. *arXiv preprint arXiv:2103.06172* (2021).
- [14] Solon Barocas, Moritz Hardt, and Arvind Narayanan. 2017. Fairness in machine learning. *Nips tutorial* 1 (2017), 2.
- [15] Antonio Bella, César Ferri, José Hernández-Orallo, and María José Ramírez-Quintana. 2010. Calibration of machine learning models. In *Handbook of Research on Machine Learning Applications and Trends: Algorithms, Methods, and Techniques*. IGI Global, 128–146.
- [16] Alex Beutel, Jilin Chen, Tulsee Doshi, Hai Qian, Allison Woodruff, Christine Luu, Pierre Kreitmann, Jonathan Bischof, and Ed H. Chi. 2019. Putting Fairness Principles into Practice: Challenges, Metrics, and Improvements. <https://arxiv.org/pdf/1901.04562.pdf>
- [17] Miranda Bogen, Pushkar Tripathi, Aditya Srinivas Timmaraju, Mehdi Mashayekhi, Qi Zeng, Rabyd (Rob) Roudani, Sean Gahagan, Andrew Howard, and Isabella Leone. 2023. Toward fairness in personalized ads. (2023). <https://about.fb.com/news/2023/01/an-update-on-our-ads-fairness-efforts/>
- [18] Abhishek Chakraborty, Preetam Nandy, and Hongzhe Li. 2018. Inference for individual mediation effects and interventional effects in sparse high-dimensional causal graphical models. *arXiv preprint arXiv:1809.10652* (2018).
- [19] Alexandra Chouldechova. 2017. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data* 5, 2 (2017), 153–163.
- [20] Alexandra Chouldechova and Aaron Roth. 2020. A snapshot of the frontiers of fairness in machine learning. *Commun. ACM* 63, 5 (2020), 82–89.
- [21] Sam Corbett-Davies and Sharad Goel. 2018. The measure and mismeasure of fairness: A critical review of fair machine learning. *arXiv preprint arXiv:1808.00023* (2018).
- [22] Kate Crawford. 2021. *The atlas of AI: Power, politics, and the planetary costs of artificial intelligence*. Yale University Press.
- [23] Cyrus DiCiccio, Brian Hsu, YinYin Yu, Preetam Nandy, and Kinjal Basu. 2022. Predictive Rate Parity Testing and Mitigation. *arXiv preprint arXiv:2204.05947* (2022).
- [24] Imani Dunbar. 2022. Mythbusting the Feed: How We Work to Address Bias. (2022). <https://blog.linkedin.com/2022/november/1/mythbusting-the-feed-how-we-work-to-address-bias>

- [25] Imani Dunbar. 2022. Scaling Self-ID to Better Connect Members to Opportunities. (2022). <https://www.linkedin.com/pulse/scaling-self-id-better-connect-members-opportunities-imani-dunbar/?trk=pulse-article>
- [26] Cynthia Dwork. 2008. Differential privacy: A survey of results. In *International conference on theory and applications of models of computation*. Springer, 1–19.
- [27] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. 2012. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*. 214–226.
- [28] Cynthia Dwork, Aaron Roth, et al. 2014. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science* 9, 3–4 (2014), 211–407.
- [29] Ferdinando Fioretto, Cuong Tran, Pascal Van Hentenryck, and Keyu Zhu. 2022. Differential Privacy and Fairness in Decisions and Learning Tasks: A Survey. *arXiv preprint arXiv:2202.08187* (2022).
- [30] Caroline Fontaine and Fabien Galand. 2007. A survey of homomorphic encryption for nonspecialists. *EURASIP Journal on Information Security* 2007 (2007), 1–10.
- [31] Rina Friedberg, Stuart Ambler, and Guillaume Saint-Jacques. 2022. Representation-Aware Experimentation: Group Inequality Analysis for A/B Testing and Alerting. *arXiv preprint arXiv:2204.12011* (2022).
- [32] Rina Friedberg and Ryan Rogers. 2022. Privacy Aware Experimentation over Sensitive Groups: A General Chi Square Approach. *arXiv preprint arXiv:2208.08564* (2022).
- [33] Sahaj Garg, Vincent Perot, Nicole Limtiaco, Ankur Taly, Ed H Chi, and Alex Beutel. 2019. Counterfactual fairness in text classification through robustness. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*. 219–226.
- [34] Sahin Cem Geyik, Qi Guo, Bo Hu, Cagri Ozcaglar, Ketan Thakkar, Xianren Wu, and Krishnaram Kenthapadi. 2018. Talent search and recommendation systems at LinkedIn: Practical challenges and lessons learned. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*. 1353–1354.
- [35] Avijit Ghosh, Aalok Shanbhag, and Christo Wilson. 2022. Faircanary: Rapid continuous explainable fairness. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*. 307–316.
- [36] Ruocheng Guo, Lu Cheng, Jundong Li, P Richard Hahn, and Huan Liu. 2020. A survey of learning causality with data: Problems and methods. *ACM Computing Surveys (CSUR)* 53, 4 (2020), 1–37.
- [37] Ido Guy and Luiz Pizzato. 2016. People recommendation tutorial. In *Proceedings of the 10th ACM Conference on Recommender Systems*. 431–432.
- [38] Xudong Han, Timothy Baldwin, and Trevor Cohn. 2021. Balancing out bias: Achieving fairness through training reweighting. *arXiv preprint arXiv:2109.08253* (2021).
- [39] Moritz Hardt, Eric Price, and Nati Srebro. 2016. Equality of opportunity in supervised learning. *Advances in neural information processing systems* 29 (2016).
- [40] Tatsunori Hashimoto, Megha Srivastava, Hongseok Namkoong, and Percy Liang. 2018. Fairness without demographics in repeated loss minimization. In *International Conference on Machine Learning*. PMLR, 1929–1938.
- [41] Ursula Hébert-Johnson, Michael Kim, Omer Reingold, and Guy Rothblum. 2018. Multicalibration: Calibration for the (computationally-identifiable) masses. In *International Conference on Machine Learning*. PMLR, 1939–1948.
- [42] Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. 2018. Evaluating feature importance estimates. (2018).
- [43] Brian Hsu, Xiaotong Chen, Ying Han, Hongseok Namkoong, and Kinjal Basu. 2023. An Operational Perspective to Fairness Interventions: Where and How to Intervene. *arXiv preprint arXiv:2302.01574* (2023).
- [44] Brian Hsu, Rahul Mazumder, Preetam Nandy, and Kinjal Basu. 2022. Pushing the limits of fairness impossibility: Who’s the fairest of them all? *arXiv preprint arXiv:2208.12606* (2022).
- [45] Vasileios Iosifidis and Eirini Ntoutsi. 2019. Adafair: Cumulative fairness adaptive boosting. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. 781–790.
- [46] Shahin Jabbari, Matthew Joseph, Michael Kearns, Jamie Morgenstern, and Aaron Roth. 2017. Fairness in reinforcement learning. In *International conference on machine learning*. PMLR, 1617–1626.
- [47] Ajil Jalal, Sushrut Karmalkar, Jessica Hoffmann, Alex Dimakis, and Eric Price. 2021. Fairness for image generation with uncertain sensitive attributes. In *International Conference on Machine Learning*. PMLR, 4721–4732.
- [48] Matthew Joseph, Michael Kearns, Jamie H Morgenstern, and Aaron Roth. 2016. Fairness in learning: Classic and contextual bandits. *Advances in neural information processing systems* 29 (2016).
- [49] Krishnaram Kenthapadi, Benjamin Le, and Ganesh Venkataraman. 2017. Personalized job recommendation system at linkedin: Practical challenges and lessons learned. In *Proceedings of the eleventh ACM conference on recommender systems*. 346–347.
- [50] Michael P Kim, Amirata Ghorbani, and James Zou. 2019. Multiaccuracy: Black-box post-processing for fairness in classification. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*. 247–254.
- [51] Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. 2017. Inherent Trade-Offs in the Fair Determination of Risk Scores. In *8th Innovations in Theoretical Computer Science Conference (ITCS 2017) (Leibniz International Proceedings in Informatics (LIPIcs), Vol. 67)*, Christos H. Papadimitriou (Ed.), Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, Dagstuhl, Germany, 43:1–43:23. <https://doi.org/10.4230/LIPIcs.ITCS.2017.43>
- [52] Ron Kohavi, Alex Deng, Brian Frasca, Toby Walker, Ya Xu, and Nils Pohlmann. 2013. Single controlled experiments at large scale. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*. 1168–1176.
- [53] Preeti Lahoti, Alex Beutel, Jilin Chen, Kang Lee, Flavien Prost, Nithum Thain, Xuezhi Wang, and Ed Chi. 2020. Fairness without demographics through adversarially reweighted learning. *Advances in neural information processing systems* 33 (2020), 728–740.
- [54] Honglei Liu, Anuj Kumar, Wenhai Yang, and Benoit Dumoulin. 2018. Explore-Exploit: A Framework for Interactive and Online Learning. *arXiv preprint arXiv:1812.00116* (2018).
- [55] Pranay K Lohia, Karthikeyan Natesan Ramamurthy, Manish Bhide, Diptikalyan Saha, Kush R Varshney, and Ruchir Puri. 2019. Bias mitigation post-processing for individual and group fairness. In *Icassp 2019-2019 IEEE international conference on acoustics, speech and signal processing (icassp)*. IEEE, 2847–2851.
- [56] Payman Mohassel and Yupeng Zhang. 2017. Secureml: A system for scalable privacy-preserving machine learning. In *2017 IEEE symposium on security and privacy (SP)*. IEEE, 19–38.
- [57] Preetam Nandy, Kinjal Basu, Shaunak Chatterjee, and Ye Tu. 2020. A/B testing in dense large-scale networks: design and inference. *Advances in Neural Information Processing Systems* 33 (2020), 2870–2880.
- [58] Preetam Nandy, Yunsong Meng, Cyrus DiCiccio, Heloise Logan, Amir Sepelri, Divya Venugopalan, Kinjal Basu, and Nouredine El Karoui. 2021. Using the LinkedIn Fairness Toolkit in large-scale AI systems. (2021). <https://engineering.linkedin.com/blog/2021/using-the-linkedin-fairness-toolkit-large-scale-ai>
- [59] Preetam Nandy, Divya Venugopalan, Chun Lo, and Shaunak Chatterjee. 2021. A/B testing for recommender systems in a two-sided marketplace. *Advances in Neural Information Processing Systems* 34 (2021), 6466–6477.
- [60] Arvind Narayanan. 21. Fairness definitions and their politics. In *Tutorial presented at the Conf. on Fairness, Accountability, and Transparency*.
- [61] Judea Pearl. 2009. *Causality*. Cambridge university press.
- [62] Felix Petersen, Debarghya Mukherjee, Yuekai Sun, and Mikhail Yurochkin. 2021. Post-processing for individual fairness. *Advances in Neural Information Processing Systems* 34 (2021), 25944–25955.
- [63] Geoff Pleiss, Manish Raghavan, Felix Wu, Jon Kleinberg, and Kilian Q Weinberger. 2017. On fairness and calibration. *Advances in neural information processing systems* 30 (2017).
- [64] Manish Raghavan, Solon Barocas, Jon Kleinberg, and Karen Levy. 2020. Mitigating bias in algorithmic hiring: Evaluating claims and practices. In *Proceedings of the 2020 conference on fairness, accountability, and transparency*. 469–481.
- [65] Rohan Ramanath, Hakan Inan, Gungor Polatkan, Bo Hu, Qi Guo, Cagri Ozcaglar, Xianren Wu, Krishnaram Kenthapadi, and Sahin Cem Geyik. 2018. Towards deep and representation learning for talent search at LinkedIn. In *Proceedings of the 27th ACM international conference on information and knowledge management*. 2253–2261.
- [66] Jonathan Roth, Guillaume Saint-Jacques, and YinYin Yu. 2022. An Outcome Test of Discrimination for Ranked Lists. In *2022 ACM Conference on Fairness, Accountability, and Transparency*. 350–356.
- [67] Guillaume Saint-Jacques, Amir Sepelri, Nicole Li, and Igor Perisic. 2020. Fairness through Experimentation: Inequality in A/B testing as an approach to responsible design. *arXiv preprint arXiv:2002.05819* (2020).
- [68] Jad Salem and Swati Gupta. 2019. Secretary Problems with Biased Evaluations using Partial Ordinal Information. *Available at SSRN 3444283* (2019).
- [69] Chenglei Si, Zhe Gan, Zhengyuan Yang, Shuohang Wang, Jianfeng Wang, Jordan Boyd-Graber, and Lijuan Wang. 2022. Prompting gpt-3 to be reliable. *arXiv preprint arXiv:2210.09150* (2022).
- [70] Xiaoqiang Sun, Peng Zhang, Joseph K Liu, Jianping Yu, and Weixin Xie. 2018. Private machine learning classification based on fully homomorphic encryption. *IEEE Transactions on Emerging Topics in Computing* 8, 2 (2018), 352–364.
- [71] Shuhan Tan, Yujun Shen, and Bolei Zhou. 2020. Improving the fairness of deep generative models without retraining. *arXiv preprint arXiv:2012.04842* (2020).
- [72] Sriram Vasudevan and Krishnaram Kenthapadi. 2020. Lift: A scalable framework for measuring fairness in ml applications. In *Proceedings of the 29th ACM international conference on information & knowledge management*. 2773–2780.
- [73] Sandra Wachter, Brent Mittelstadt, and Chris Russell. 2021. Why fairness cannot be automated: Bridging the gap between EU non-discrimination law and AI. *Computer Law & Security Review* 41 (2021), 105567.
- [74] Mingyang Wan, Daochen Zha, Ninghao Liu, and Na Zou. 2022. In-Processing Modeling Techniques for Machine Learning Fairness: A Survey. *ACM Transactions on Knowledge Discovery from Data (TKDD)* (2022).
- [75] Alice Xiang and Inioluwa Deborah Raji. 2019. On the legal compatibility of fairness definitions. *arXiv preprint arXiv:1912.00761* (2019).
- [76] Feiyu Xu, Hans Uszkoreit, Yangzhou Du, Wei Fan, Dongyan Zhao, and Jun Zhu. 2019. Explainable AI: A brief survey on history, research areas, approaches and challenges. In *Natural Language Processing and Chinese Computing: 8th CCF*

International Conference, NLPCC 2019, Dunhuang, China, October 9–14, 2019, Proceedings, Part II 8. Springer, 563–574.

- [77] Kaihe Xu, Hao Yue, Linke Guo, Yuanxiong Guo, and Yuguang Fang. 2015. Privacy-preserving machine learning algorithms for big data systems. In *2015 IEEE 35th international conference on distributed computing systems*. IEEE, 318–327.
- [78] Xun Yi, Russell Paulet, Elisa Bertino, Xun Yi, Russell Paulet, and Elisa Bertino. 2014. *Homomorphic encryption*. Springer.

A ERROR COHORT ANALYSIS

The Sankey diagram (Figure 7) shows the flow of members from the entire data (left) to individual cohorts (right). Cohorts are numbered in ascending order based on residual error. For example, cohort 1 has the most negative error (corresponding to more overpredicted samples) and is the brightest blue. On the other hand, cohort 8 has the most positive error (corresponding to more ‘underpredicted’ samples) and is the brightest red. Text describes the features used in the splits.

B EDGE-LEVEL ANALYSIS FOR MEASURING UNINTENDED CONSEQUENCES

Figure 8 shows a toy example to highlight the analysis.

We compute the edge-level response rates (RR) for the female members “A” and “C.” Green edges indicate positive responses (e.g. a connection invite was sent in the arrow direction and the invite was accepted) while red edges indicate negative responses (e.g. ignored/rejected invite). Our aim is to detect if the positive response rate to treatment viewers is lower, which may indicate that females are receiving lower-quality interactions on the platform.

The full edge-based analysis results are highlighted in the Table 4.

C SELF-ID

Figure 9 shows one of the entry points into the form as well as some of the initial questions in the form. For more details on self-id please see [1, 24, 25].

D PEOPLE ALSO VIEWED

This People Also Viewed product recommends additional profiles you might be interested in based on the current profile you are viewing and is an example of the scarce resource allocation paradigm when considering the recommended members. The product is the second largest source of traffic for profile views, therefore any negative bias (failing to surface qualified recommendations) could adversely affect members’ abilities to grow their network.

The list of recommendations appear when a viewer is looking at a particular member’s (source) profile. See the Figure 10 for the details.

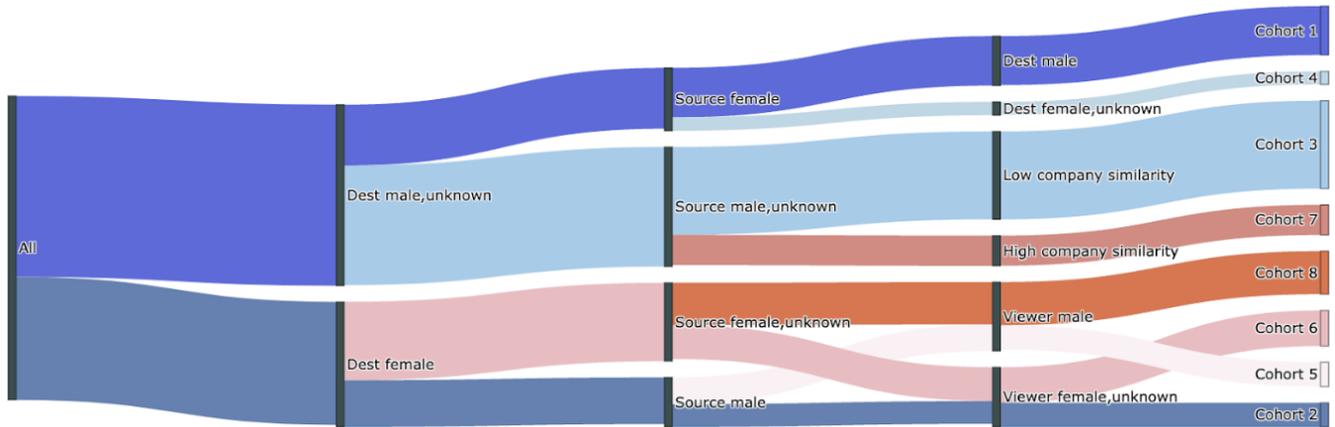


Figure 7: The Sankey diagram for People Also Viewed

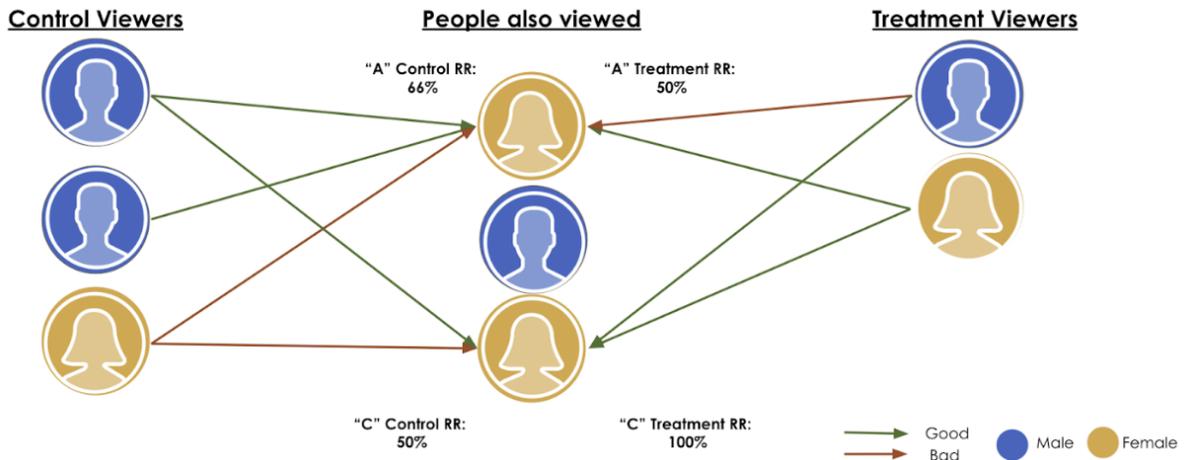


Figure 8: We look at all edges from control to females and treatment to female and check if those edges are “good” at a similar rate. “Good” edges are formed when the recommended member accepts the invite, replies to a message and doesn’t report the viewer. “Bad” edges are formed on the contrary. RR stands for the response rate.

Edge Type	Metric	ControlRate	TreatmentRate	RelativeDiff
All Edges	accept	0.211385	0.244152	0.155011
Edges to Female	accept	0.20275	0.234024	0.154251
Edges to Male	accept	0.218781	0.250646	0.145647
All Edges	reply	0.313753	0.305175	-0.027337
Edges to Female	reply	0.312765	0.306352	-0.020507
Edges to Male	reply	0.310661	0.300983	-0.031154
All Edges	report	7.00e-5	7.60e-5	0.077257
Edges to Female	report	4.70e-5	4.60e-5	-0.017529
Edges to Male	report	8.30e-5	9.20e-5	0.115913

Table 4: Edge-level response rates across metrics comparing the base model with the BMT model.

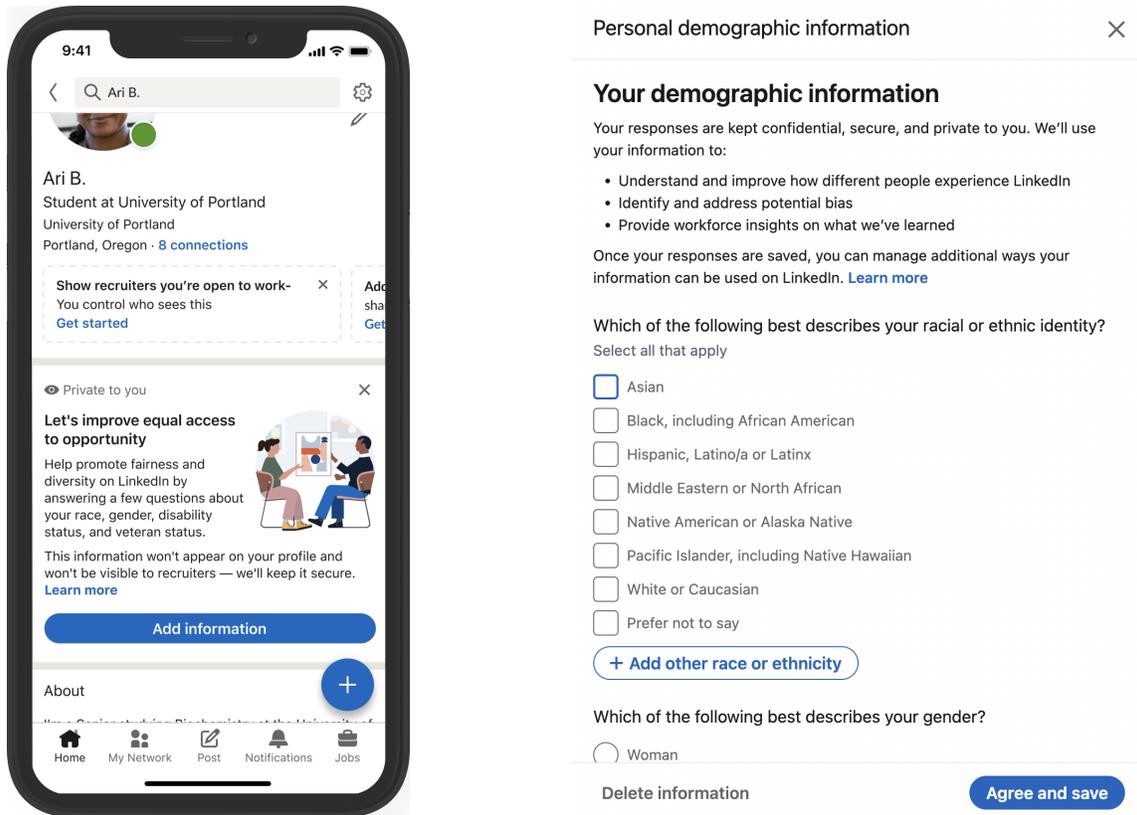


Figure 9: Self-ID form (partially visible) which is accessible through the profile update feature on LinkedIn

The image shows a LinkedIn profile page for Yifan (Eva) Xue. At the top, there is a navigation bar with icons for Home, My Network, Jobs, Messaging, Notifications, Me, Work, and Learning. The profile header features a yellow background with a network diagram. The profile information includes a circular profile picture, the name 'Yifan (Eva) Xue (She/Her)', the title 'Machine Learning Engineer at LinkedIn Retention and Ecosystem AI', and the affiliation 'University of Pittsburgh School of Medicine'. Below this, there are statistics for followers and connections, and buttons for 'Open to', 'Add profile section', and 'More'. Two interactive cards are visible: one for showing recruiters you're open to work and another for sharing hiring status. The 'Suggested for you' section is partially visible. The 'People also viewed' section, highlighted with a red border, lists three recommended members: Liwei Cai (SDE in AI at LinkedIn), Fengyu Zhang (Machine Learning Engineer at LinkedIn | Princeton 15 | USTC'er 11), and Beibei WANG (Machine learning engineer at LinkedIn; Stanford PhD). Each member has a 'Connect' button.

Figure 10: In this example, Yifan Xue is the source member, the recommended members in the highlighted “People also viewed” card are dest members, and whoever is viewing the profile is the viewer.