

# Collectivizing Justice: A Novel Argument for Quota-Based Affirmative Action

György Barabás and András Szigeti

**Abstract.** We offer a new argument in favor of quota-based affirmative action and a mathematical demonstration of that argument. Consider a minority group  $A$  forming part of a larger group  $G$ . We show that even an equitable hiring policy can fail to improve the underrepresentation of  $A$ . Under the plausible assumption that members of  $A$  are more likely to quit  $G$  due to being marginalized, intragroup distribution will be dominated by “points of recalcitrance” preventing the move towards a more proportionate distribution. Consequently, if affirmative action is to succeed, the use of quota-based recruitment procedures is in practice often unavoidable because only strongly preferential hiring of the minority group will make it possible to push the intragroup distribution past the point of recalcitrance. A general implication for moral theory is that some duties owed to individuals can only be discharged by targeting the group to which these individuals belong.

“In order to get beyond racism, we must first take account of race. There is no other way.” (Justice Harry Blackmun in *Regents of the University of California v. Bakke*)

## I INTRODUCTION

This paper offers a new argument in favor of quota-based affirmative action and a mathematical demonstration of that argument. The argument purports to show that if affirmative action is to succeed, then the use of quota-based recruitment procedures is in practice often *unavoidable*—whatever the aims and justification of affirmative action may be. A general implication of the argument is that some duties generated by individual rights or entitlements can only be discharged effectively by addressing the group to which these individuals belong. We may sometimes have to pay to the group what we owe to individuals.

By affirmative action we mean measures undertaken to increase the proportion of some smaller group  $A$  (e.g., female employees) in a larger group  $G$  (e.g., employees at a workplace) (see Fullinwider 2018). Affirmative action is typically thought of as involving the preferential selection of group  $A$  at the expense of another group  $B$  (e.g., male employees) or groups ( $B, C, D, \dots$ )<sup>1</sup> constituting  $G$ . Justifications of affirmative action have appealed to compensatory and/or distributive justice, and/or considerations of social utility. As will be seen, our argument is of wide applicability insofar as it does not require commitment to any of these justifications. The claim is simply that whatever the motivation for affirmative action policies may be, these policies can fail to achieve an equitable representation due to the dynamic characteristics of the relative distribution of  $A$  and  $B$  within  $G$ . This is because under a small set of plausible and realistic assumptions—in particular, the assumption that members of  $A$  are more likely to quit  $G$  due to being marginalized—the intragroup distribution will be dominated by certain “points of recalcitrance” which prevent a shift towards the desired, more proportionate distribution.

---

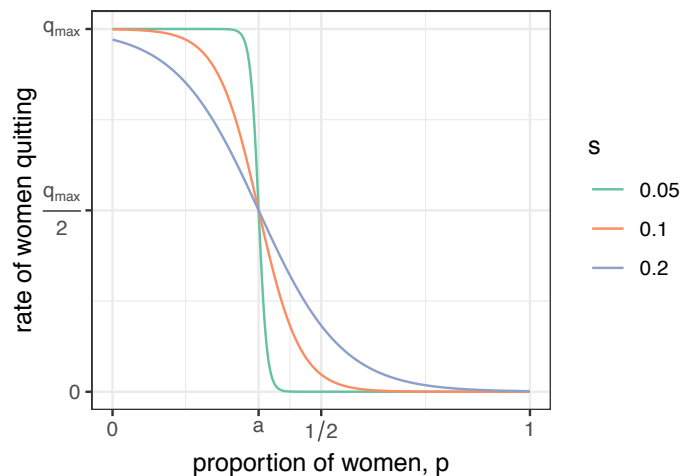
<sup>1</sup>In what follows, we ignore the possibility of more than two subgroups constituting  $G$ .

The surprising upshot is that *a hiring policy which is in principle fair can conserve the minority status of the smaller group*, and this will remain the case even if all relevant stakeholders are well-intentioned and committed to the aims of affirmative action. We see this finding as strongly supporting quota-based affirmative action because it shows that only preferential hiring of relatively large numbers of the minority group will make it possible to push the intragroup distribution past the points of recalcitrance.

While we are not committed to any substantive theory in the ethics of affirmative action, our finding does have some bearing on the debate concerning what is perhaps the most common *objection* to affirmative action, namely that it amounts to reverse discrimination (Cowan 1972, Nunn 1974, Goldman 1975, Simon 1978; see also Fullinwider 2018, Lippert-Rasmussen 2018). If the argument presented here is correct, then that debate is largely redundant. As long as one recognizes as valid *any* normative reason (based on considerations of justice and/or social utility) for increasing the representation of the minority group, reliance on quota-based affirmative action will be practically inescapable. However, all advocates of the reverse discrimination objection do recognize the validity of at least some such reasons, as is clear from the very name given to their objection which implies that discrimination is a morally objectionable phenomenon.

## II THE MODEL

Consider an institution hiring employees at a relative rate  $h$  (i.e.,  $h$  is the rate of new hires, divided by the number of employees). A fraction  $f$  of all new hires are women. Employees retire at a per capita rate  $r$ , assumed to be equal across genders. However, crucially, employees may also quit for reasons of marginalization caused by being a minority at work. These rates depend on the current gender proportion at the workplace: the higher the proportion of one gender, the less likely it is for them to leave before retirement. For women, this rate is given by a function  $q(p)$ , where  $p$  is their fraction at the workplace. We also assume that premature quitting is symmetric with respect to gender, i.e., the probability of a woman quitting a workplace with only 10% women is the same as the probability of a man quitting a workplace



**Figure 1:** The rate of women prematurely quitting their job, as a function of their fraction  $p$  at the workplace. The function takes on its maximum  $q_{\max}$  when  $p = 0$ , then decreases to near zero for  $p$  sufficiently high. The function drops to half its maximum value,  $q_{\max}/2$ , at  $p = a$ ;  $a$  is a measure of how high the fraction of women must be to provide a work environment which does not compel them to quit. The speed of transition from high to low quitting rates is governed by the parameter  $s$  (color legend): the smaller its value, the more abrupt the transition will be.

with 10% men. With these assumptions, one can write down a model to trace the rate of change of women’s proportion in the group,  $dp/dt$ , as a function of the actual proportion  $p$  of women:<sup>2</sup>

$$\frac{dp}{dt} = (f - p)h - p(1 - p)[q(p) - q(1 - p)]$$

The function  $q(p)$  is shown in Figure 1. It transitions from its maximum  $q_{\max}$  to zero around the point  $a$ , with an abruptness measured by the parameter  $s$ . In words,  $q_{\max}$  is the rate of women quitting due to being in minority when their fraction is very low;  $a$  is the threshold fraction below which women are truly a minority (in the sense that their incentive to quit prematurely is strong); and  $s$  measures how fast the quitting rates transition from high to low values, with  $s = 0$  meaning an infinitely fast (“flip”) transition.

The model is fully specified after designating the starting proportion of women and the five parameters  $f$ ,  $h$ ,  $q_{\max}$ ,  $a$ , and  $s$ . We now turn to the discussion of four instructive scenarios which can be realized within the model depending on parameter values.

<sup>2</sup>See the Appendix for a detailed derivation.

### III FOUR SCENARIOS OF INTRAGROUP GENDER DYNAMICS

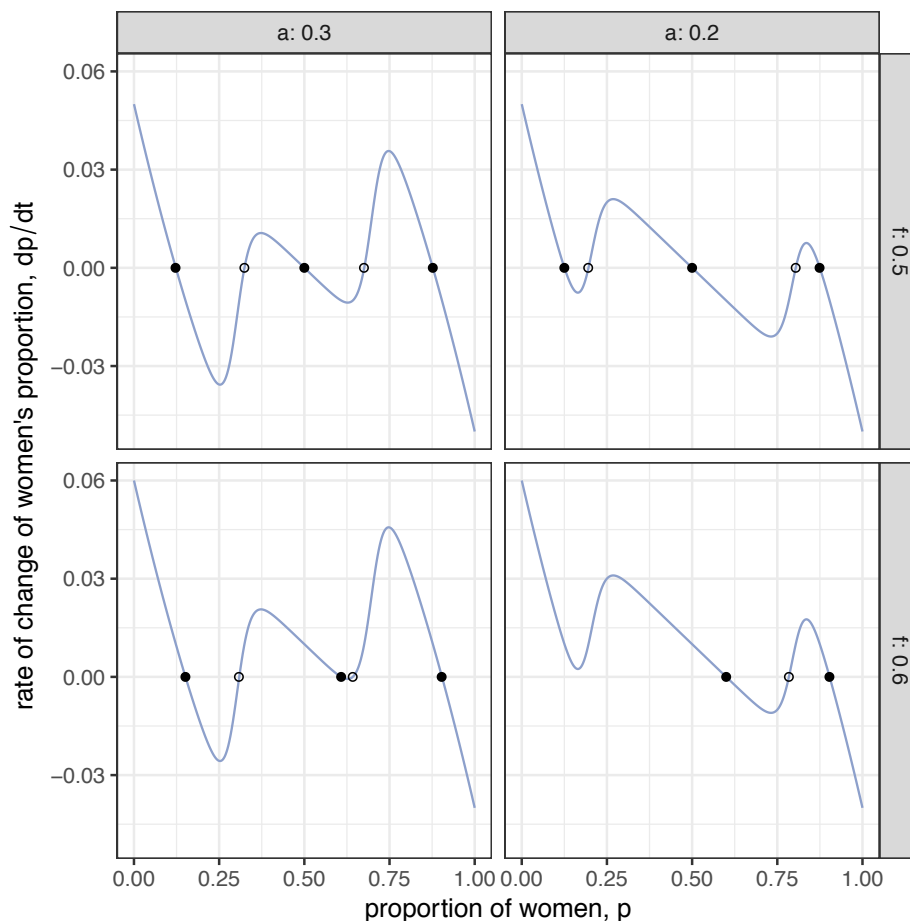
In the first scenario (Figure 2, top left), the hiring rates of men and women are equal. If there are no female employees initially, the rate of change will be positive for a while. Consequently, the proportion of women will increase until a stable equilibrium is reached at 12%. By contrast, the model reveals that the rate of change will be negative between 12% and 32%. That is, crucially, the proportion of women will decrease, and the fraction of women will eventually slide back to 12%. If the initial proportion is higher than 32% but lower than 68%, then the proportion of women will converge to a stable equilibrium at 50%, the point of gender equality. Finally, if the initial proportion exceeds 68%, then (since the rate of change is positive) the proportion will increase to 88%. This would lead to a reversal of gender roles: now men are the minority with a high propensity to quit prematurely.

These observations illustrate our central claim about points of recalcitrance in the intragroup distribution. If the initial fraction of female employees is low (in this scenario, lower than 32%), then their fraction will eventually stabilize at 12%—despite the fact that hiring rates of women and men are equal. This is because when women are a minority, they are more likely to quit their jobs prematurely. As a result, their proportion will not reach 50%. In short, although the hiring regime is fair it nevertheless conserves the disadvantage of the minority group.

The second scenario (Figure 2, top right) has the same equal hiring rate, but quitting rates from being in minority are reduced by placing the point of transition from high to low quitting rates earlier. In terms of our model, this means that the parameter  $a$  is lowered.<sup>3</sup> That is to say, in this scenario women are assumed to be somewhat less likely to leave due to marginalization. What the model shows for this parameterization is that the proportion of women needs to exceed 20% to converge to equal representation at 50%. Significantly, however, for initial proportions less than 20%, female representation will still stabilize at a low 12.5% despite the fact that women are less affected by marginalization in this scenario.

---

<sup>3</sup>See Figure 1 and the Appendix.



**Figure 2:** The rate of change of the proportion of women,  $dp/dt$ , as a function of the actual proportion  $p$ , for four different model parameterizations. Solid/open dots show stable/unstable equilibria (see Appendix). Rows show scenarios with the same relative hiring rate ( $f = 0.5$  for the top and  $0.6$  for the bottom row); columns show scenarios with the same inflection point  $a$  of the quitting function ( $a = 0.3$  in the left and  $0.2$  in the right column; see Figure 1 and the Appendix for details). All other parameters are fixed and equal across scenarios:  $h = 0.1$ ,  $q_{\max} = 0.35$ ,  $s = 0.04$ .

In the third scenario (Figure 2, bottom left), quitting rates are as high as in the first one, but the hiring rate is skewed to favor women ( $f = 0.6$ ). It is interesting to observe that, by itself, a *moderate* affirmative action regime such as this one will not improve matters significantly. The point of recaltrance will still only move to around 15%, and the proportion of women needs to exceed 31% to increase above 15% in the long run.

In fact, one can calculate that in this example, the hiring rate of women would have to be pushed over 87% to eliminate the lower point of recaltrance—that is, to reach a high representation of women even if their initial proportion is low. Since enforcing such a quota,

in which approximately 9 out of every 10 hires are women, would constitute a rather radical form of affirmative action, it is worth taking a look at the potential impact of a combination of measures. For example, one could skew hiring rates to favor women as in the third scenario *and* assume a reduced effect of marginalization as in the second. This possibility is examined in our fourth and final scenario (Figure 2, bottom right). Thanks to this combination of measures, it is suddenly possible for women to achieve high representation even if their initial proportion is low. As seen, the lower point of recalcitrance is eliminated, and instead the system approaches the stable equilibrium at 60%. Observe that this scenario does not yield a gender-equal distribution. However, balance can be achieved by resetting the hiring rate to  $f = 0.5$  and stopping affirmative action after the proportion of women has reached 50%. Once this is done, the dynamics of the second scenario takes over, and so gender parity becomes robust and stable in the long run.

These four scenarios are instructive because they allow us to quantify the impact of quota-based affirmative action. They illustrate how to eliminate the stable equilibrium—the point of recalcitrance—where the proportion of women is low. Sometimes this can only be done by means of a strongly preferential hiring regime. If other parameters can also be adjusted, then affirmative action may be combined with measures aimed at reducing the probability that women will leave owing to their marginalization in the group.

## IV DISCUSSION

We are assuming that women and men are equally good meritocratically speaking and that the two pools of talent from which women and men are respectively recruited are equally large.<sup>4</sup> With these assumptions in place, all we are operating with is the highly intuitive idea that every person of equal merit should have an equal chance of being part of the group (an idea also embraced by opponents of affirmative action or any form of compensatory justice).

---

<sup>4</sup>The latter assumption may be questioned among others because the lack of fair representation may be partly attributable to a “pipeline problem”. However, note that if that assumption was found to be wrong (and the first assumption is accepted), then that would only strengthen the justification for quota-based affirmative action. Below, we will argue that the implications of our model support those who argue that even if there is a pipeline problem, it is not the only major reason for the lack of fair representation.

What we show is how difficult it is to level the playing field. In fact, even the aim of equal representation has only been chosen for simplicity's sake in our model. Points of recalcitrance will emerge for most other desired distributions as well, e.g., if the desirable share of the minority group is higher or lower than 50%. Similarly, both the assumption that retirement rates are gender-independent, and that the effects of marginalization are symmetrical can be relaxed without eliminating points of recalcitrance.

A crucial assumption is that members of the minority group are more likely to quit due to the fact that they are in the minority. Naturally, this need not be true. In some situations, it is even possible that being in the minority will decrease the probability of quitting—for example, minority members may subjectively or objectively benefit from their minority position. Furthermore, the S-shape of the quitting function in Figure 1 is essential for generating points of recalcitrance, which would not occur with e.g. a linear quitting function.

Nevertheless, it seems to us eminently plausible that in many cases being marginalized significantly increases the likelihood that individuals constituting the minority will leave the group. For instance, the expression of important aspects of one's identity (e.g., language, religion, gender) and the representation of certain concerns and needs within a larger group are dependent on there being a sufficient number of other members sharing the same identity, concerns, and needs. Other things being equal, marginalization is also more likely to give rise to segregation, stigmatization, and discrimination (Lippert-Rasmussen 2018). Finally, there is an increased risk of an unfair distribution of resources due to the latent or explicit tyranny of the majority. Note that since our model looks at the likelihood that members of the minority group will *feel* like or entertain thoughts of leaving, it is not even necessary here to consider whether or not being in the minority does in fact entail such adverse consequences in a given group. What matters is that awareness of such risks can be expected to engender a feeling of precarity and thereby make it more likely that members of the minority will quit.

There is also ample empirical evidence supporting our assumption. To quote from three areas of particular interest: it has been argued, first, that the underrepresentation of women in philosophy is not just a “pipeline problem” (Dodds and Goddard 2013; see also Allen and



Castleman 2001, Schiebinger 2000). That is, the problem of underrepresentation cannot be resolved just by making sure that there are enough female PhDs and junior female philosophers in the profession. On the contrary, it seems that several other factors contribute significantly to entrenching underrepresentation, and many of these have to do with marginalization such as a hostile or “chilly” institutional culture or an insufficient recognition of and respect for the minority group (Dodds and Goddard 2013). Second, studies of women in legislatures show that a “critical mass” of women may be necessary to effectively promote a feminist political agenda broadly understood (Beckwith and Cowell-Meyers 2007). It can be safely predicted that, frustrated in their political aims due to the lack of a critical mass, marginalized women will be discouraged to stay in politics. And third, empirical data shows that African-American students suffer high attrition rates at law schools despite<sup>5</sup> the fact that many of these schools use affirmative action to recruit students (Sander 2004; see also Fullinwider 2018).

The evidence and the theoretical considerations provided above are of course piecemeal and contextual. Our sole purpose here is to show that the key assumption of our model—that marginalization increases the propensity to leave—is not just intuitive but also consistent with readily available data. Future work can determine to what extent the model’s predictions are borne out in various areas where affirmative action is practiced. We will now go on to consider some of the implications of the model for moral theory.

## V GROUP JUSTICE

While in itself uncommitted to any specific view about the justification of affirmative action, the above finding about the dynamics of intragroup distribution has significant implications for the debate about affirmative action in moral philosophy. In particular, it bears on the controversy around what is perhaps the best rehearsed objection to affirmative action, namely that affirmative action is unjustifiable because it itself amounts to discrimination.

---

<sup>5</sup>Or as some would argue precisely *because* (or at least partly because) these schools use affirmative action (see Thernstrom and Thernstrom 1997) as the students benefiting from affirmative action end up in academic environments where they underperform.

The most theoretically challenging way to understand this objection is that affirmative action involves discrimination of some group (e.g., white males) because members of that group suffer a disadvantage solely based on their membership in a racial, gender, etc. group. And conversely, members of the group targeted by affirmative action do enjoy an advantage solely based on their membership in that group. This is morally objectionable because this kind of group membership is extraneous to whether someone deserves or has a legitimate claim to the kind of position or resource allocated by means of affirmative action. Further, this is morally objectionable as well because even those members of the beneficiary group who are not at all disadvantaged ("the African-American brain surgeon's son") will benefit from affirmative action, and conversely, even disadvantaged members of the "losing" group (the equally proverbial "white miner's son") will be further disadvantaged by affirmative action (Nunn 1974, Goldman 1975, 2015).

In short, the problem with affirmative action is supposed to be that it amounts to "substituting group considerations for evaluation of individual cases" (Simon 1978, 41) or "individuals [being] regarded merely as members of that group rather than in their individuality" (Cowan 1972, 11). The standard reply to this by proponents of affirmative action is that while the property  $P_1$  of being a group member in a racial, gender, etc. group in itself is indeed morally irrelevant in theory, in practice it is not. This is because  $P_1$  overlaps with the property  $P_2$  of being someone who should<sup>6</sup> receive preferential treatment through affirmative action to a sufficient degree to justify "administering justice" at the group level (Lippert-Rasmussen 2018, Nagel 1973, Nickel 1972, 1974, Thomson 1973). It is true that administering justice at the level of each individual would be better morally speaking, but determining whether or not an individual  $X$  who is  $P_1$  is also  $P_2$  may just be prohibitively costly (see esp. Nickel 1974).

---

<sup>6</sup>As we have noted, proponents of affirmative action differ as to what entitles one to be the beneficiary of affirmative action. Being the victim of past or present discrimination is the most frequent reason given (Thomson 1973), but other reasons to do with social utility are also frequently mentioned (Nagel 1973). Proponents of affirmative action may also have different views about the extent to which  $P_1$  and  $P_2$  overlap. For the claim that  $P_1 = P_2$  in the case of some minority groups, notably that of African-Americans, see for example Thomson 1973, 381). The idea is, however, that even in the absence of *complete* overlap there will be sufficient overlap in the relevant cases to render it "likely" (Lippert-Rasmussen 2018, 273) that if someone instantiates  $P_1$ , then she will also instantiate  $P_2$  (e.g., in the case of women in academia).

What our model shows, however, is that the reason for “administering justice” at the group level is not the cost or difficulty, *but rather the practical impossibility of doing so at the individual level*. In order to push past the points of recalcitrance identified by our model, it is necessary to employ some more or less radical form of quota-based affirmative action policy. This means that justice rendered through affirmative action must be collectivized if it is rendered at all.<sup>7</sup> Note also that justice here is understood in the broadest possible sense. As long as one recognizes *any* reason for increasing the representation of the minority group (based on considerations of justice, social utility, or anything else), reliance on quota-based affirmative action may be inevitable from a practical point of view.

## VI CONCLUSION

Naturally, even if we are right that quota-based affirmative action is inescapable for the reasons discussed above, this does not mean that affirmative action is necessarily justified. Some may feel that despite the peculiar dynamics of intragroup distribution, quota-based affirmative action is too high a price to pay for achieving the desired level of representation because the individuals who stand to lose from affirmative action have an unconscionably large burden to bear. As regards the justifiability of quota-based affirmative action, it is worth keeping in mind, however, two further points.

First, recall that in two scenarios we used to illustrate our model in Section III, we could eliminate (4th scenario) or at least push in the right direction (2nd scenario) the point of recalcitrance by setting the probability of quitting due to marginalization lower. It is thus possible to combine affirmative action with other measures, such as those making members of the minority group less inclined to leave. In fact, some measures recommended in the empirical literature we cited (Dodds and Goddard 2013; see also Van Norden 2017) explicitly target changing the relevant group’s culture by increasing the recognition of minorities’ needs and

---

<sup>7</sup>We thus agree with Taylor (1973) and Bayles (1973) that the target of affirmative action must be the collective. However, we reach this conclusion without relying on the contentious claim, which is central to their approach, that a group *qua* group can be harmed.

concerns. A combination of such available measures could reduce the resistance to affirmative action.

Second, while in this article we have focused on affirmative action in the usual setting (education, employment), it is worth noting that the model is applicable in other "nontraditional sites" of affirmative action such as health and housing as well (see Lippert-Rasmussen 2018, 276-7). In fact, while at this day and age the model does lend support to a certain progressive political agenda, it is in principle applicable in any situation or context where achieving a certain intragroup distribution of two or more sub-groups within the larger group seems desirable for whatever reason.

## REFERENCES

- Allen, M., Castleman, T., 2001. Fighting the pipeline fallacy. In: Brooks, A., Mackinnon, A. (Eds.), *Gender and the Restructured University*. Open University Press, pp. 151–165.
- Bayles, M. D., 1973. Reparations to wronged groups. *Analysis* 33, 182–184.
- Beckwith, K., Cowell-Meyers, K., 2007. Sheer numbers: Critical representation thresholds and women's political representation. *Perspectives on Politics* 5, 553–565.
- Cowan, J. L., 1972. Inverse discrimination. *Analysis* 33, 10–12.
- Dodds, S., Goddard, E., 2013. Not just a pipeline problem. In: Hutchison, K., Jenkins, F. (Eds.), *Women in Philosophy*. Oxford University Press, pp. 143–164.
- Fullinwider, R., 2018. Affirmative action. *Stanford Encyclopedia of Philosophy*.
- Goldman, A. H., 1975. Reparations to individuals or groups? *Analysis* 35, 168–170.
- Goldman, A. H., 2015. *Justice and Reverse Discrimination*. Princeton University Press.
- Lippert-Rasmussen, K., 2018. The ethics of anti-discrimination policies. In: Lever, A., Poama, A. (Eds.), *The Routledge Handbook of Ethics and Public Policy*. Routledge, pp. 267–280.

- Nagel, T., 1973. Equal treatment and compensatory discrimination. *Philosophy & Public Affairs* 2, 348–363.
- Nickel, J. W., 1972. Discrimination and morally relevant characteristics. *Analysis* 32, 113–114.
- Nickel, J. W., 1974. Should reparations be to individuals or to groups? *Analysis* 34, 154–160.
- Nunn, W. A., 1974. Reverse discrimination. *Analysis* 34, 151–154.
- Sander, R. H., 2004. A systemic analysis of affirmative action in American law schools. *Stanford Law Review* 57, 367.
- Schiebinger, L., 2000. Has feminism changed science? *Signs: Journal of Women in Culture and Society* 25, 1171–1175.
- Simon, R. L., 1978. Statistical justification of discrimination. *Analysis* 38, 37–42.
- Taylor, P. W., 1973. Reverse discrimination and compensatory justice. *Analysis* 33, 177–182.
- Thernstrom, S., Thernstrom, A., 1997. *America in Black and White: One Nation, Indivisible*. Simon and Schuster.
- Thomson, J. J., 1973. Preferential hiring. *Philosophy & Public Affairs* 2, 364–384.
- Van Norden, B. W., 2017. *Taking Back Philosophy: A Multicultural Manifesto*. Columbia University Press.

## APPENDIX

Here we present the derivation of our mathematical model of intragroup dynamics. Consider an institution hiring employees, at gender-specific rates  $h_w$  (women) and  $h_m$  (men). Employees retire at a per capita rate  $r$ ; this is assumed to be equal across genders. Finally, employees may also quit for reasons of marginalization caused by being a minority at work. These rates, denoted  $q_w$  for women and  $q_m$  for men, depend on their current proportions at the workplace: the higher the proportion of one gender, the less likely it is for them to leave before retirement.

Assuming that the above three processes (hires, retirements, and premature quitting) are the only factors changing the number of employees, the model may be cast in differential equation form:

$$\frac{dw}{dt} = h_w - wr - wq_w \quad (1)$$

$$\frac{dm}{dt} = h_m - mr - mq_m \quad (2)$$

Here  $w$  and  $m$  denote the number of women and men employees, and  $t$  is time. We are interested in the fraction of women,  $p = w/(w + m)$ . Applying Leibniz's rule, the chain rule of differentiation, and Eqs. 1-2, we write

$$\begin{aligned} \frac{dp}{dt} &= \frac{d}{dt} \left( \frac{w}{w+m} \right) \\ &= \frac{1}{w+m} \frac{dw}{dt} - \frac{w}{(w+m)^2} \frac{d(w+m)}{dt} \\ &= \frac{1}{w+m} \frac{dw}{dt} - \frac{w}{(w+m)^2} \left( \frac{dw}{dt} + \frac{dm}{dt} \right) \\ &= \frac{1}{w+m} (h_w - wr - wq_w) - p \frac{1}{w+m} (h_w - wr - wq_w + h_m - mr - mq_m) \\ &= \frac{h_w}{w+m} - pr - pq_w - p \left( \frac{h_w}{w+m} - pr - pq_w + \frac{h_m}{w+m} - (1-p)r - (1-p)q_m \right) \\ &= \frac{h_w}{w+m} - p \frac{h_w + h_m}{w+m} - p(1-p)(q_w - q_m) \end{aligned} \quad (3)$$

We introduce a slight reparameterization. First, we define  $h = (h_w + h_m)/(w + m)$ , the total rate of hires relative to the total number of employees  $w + m$ . Second, we define the fraction of female hires  $f = h_w/(h_w + h_m)$ . With their use, Eq. (3) reads

$$\frac{dp}{dt} = (f - p)h - p(1-p)(q_w - q_m) \quad (4)$$

To move forward, we assume that premature quitting is symmetric with respect to gender: the probability of a woman quitting a workplace with only 10% women is the same as the probability of a man quitting a workplace with 10% men. Mathematically, since the fraction of women is  $p$  (and therefore the fraction of men  $1 - p$ ), this means that we can express

both  $q_w$  and  $q_m$  through the same function:  $q_w = q(p)$  and  $q_m = q(1 - p)$ , where  $q(p)$  is the rate of quitting by women given that they make up a fraction  $p$  of the workplace. With this assumption, Eq. (4) can be written

$$\frac{dp}{dt} = (f - p)h - p(1 - p)[q(p) - q(1 - p)] \quad (5)$$

which is the form of the model used in the main text.

A reasonable form of the function  $q(p)$  and the one we adopt here is given by

$$q(p) = \frac{q_{\max}}{2} \left[ 1 - \tanh\left(\frac{p - a}{s}\right) \right] \quad (6)$$

where  $\tanh(\cdot)$  is the hyperbolic tangent function. It is shown and explained in Figure 1 of the main text.

Importantly, one can gain a broad overview of the model's behavior without having to actually solve the differential equation. All one needs to do is plot the rate of change  $dp/dt$ , given by the right hand side of Eq. (5), against  $p$ . This was done in Figure 2 for four different parameterizations. The derivative  $dp/dt$  measures the instantaneous rate of change of the proportion  $p$  of women. If  $dp/dt$  is positive,  $p$  increases; if negative,  $p$  decreases. If it is exactly zero,  $p$  does not change: the gender ratio is at equilibrium. An equilibrium point is stable if small deviations away from it decay and the system eventually returns to the same equilibrium, and unstable otherwise (indicated by shaded/open dots in Figure 2, respectively). Keeping these facts in mind, one can trace the projected fate of the fraction  $p$  in time by following along the  $x$ -axis, moving in the direction indicated by the value of  $dp/dt$ .

We have also created a web application where every parameter can be individually manipulated to examine their effect on the model. (<https://dysordys.shinyapps.io/shinyapp/>).