

# On Decentralized Affirmative Action Policies and Their Duration\*

Philippe Jehiel<sup>†</sup>  
Paris School of Economics  
& University College London

Matthew V. Leduc<sup>‡</sup>  
Paris School of Economics  
& Université Paris 1

June 30, 2022

## Abstract

Successive decentralized policy makers must decide whether to implement an affirmative action policy aimed at improving the performance distribution of future generations of a targeted group. Workers receive wages corresponding to their expected performance, suffer a feeling of injustice when getting less than their actual performance, and employers do not observe directly whether workers benefited from affirmative action. We find that welfare-maximizing policy makers choose to implement affirmative action *perpetually*, despite the resulting feeling of injustice eventually dominating the anticipated benefits to the targeted group's performance. This contrasts with the first-best that requires affirmative action to be temporary.

**Keywords:** Affirmative Action, General Equilibrium, Loss Aversion, Prospect Theory, Moral Hazard, Game Theory

**JEL codes:** D40; I28; I30; J15

---

\*We thank seminar participants at Université Paris 1, the Paris School of Economics, the Institut Henri Poincaré, Cambridge University, the Stockholm School of Economics and meetings of the Econometric Society and the Royal Economic Society, namely Agnieszka Rusinowska, Emily Tanimura, Gabrielle Demange, Jean-Marc Tallon, Leonardo Pejsachowicz, Francis Bloch, Frédéric Koessler, Olivier Tercieux, Tristan Tomala, Nicholas Vieille, Marie Laclau, Mikhail Safronov, Xin Gao. We also thank Jörgen Weibull, Andrea Galeotti, Matthew L. Elliott and Larry Samuelson for useful comments. Jehiel thanks the ERC (grant no 742816) for funding.

<sup>†</sup>Email: philippe.jehiel@gmail.com

<sup>‡</sup>(Corresponding author) Email: mattvleduc@gmail.com. Mailing address: Paris Jourdan Sciences Economiques, 48 boulevard Jourdan, 75014 Paris, France. Tel.: (+33) 01 80 52 16 60

# 1 Introduction

The original rationale for affirmative action was to help underrepresented groups close achievement gaps and such policies were often anticipated to be temporary. Decades after their inception, affirmative action policies however often remain in place. Such policies have generated a deep interest and the literature spawned over the past decades is large. Some work focuses on how such policies can be used to attempt to close achievement gaps, while other work analyzes certain inefficiencies mostly related to market distortions, such as mismatches between workers and jobs. See, for example, Fang and Moro (2011) or Sowell (2005) for broad studies.

A form of inefficiency that has gathered less attention is the devaluing effect an affirmative action policy can have on the perception of a worker's curriculum vitae. Indeed, when an affirmative action policy is in place, the mere possibility that a worker may have benefited from it can have a devaluing effect on his diplomas. The policy can also be perceived as decreasing the quality of diplomas through other channels, such as lowering academic standards by relaxing entrance requirements. If the quality of a worker's curriculum vitae is perceived by employers to be lower than his actual skills level, the worker can then experience a stigma or a feeling of injustice.

The current article will thus attempt to provide a novel explanation for the apparent stickiness of affirmative action policies in a model that takes explicitly into account their effect on devaluing the perception of a worker's curriculum vitae.

In our approach, we will consider a decentralized setting in which successive policy makers in many different districts have to decide whether or not to implement affirmative action policies. Each district's population is composed of a group  $A$  (the main group) and a group  $B$  (the group targeted by the affirmative action policy). The policy makers can be thought of as local government representatives or as school or university managers in charge of choosing which pupils or students to admit. Moreover, the policy maker of a given district anticipates that an affirmative action policy improves the talent distribution of group  $B$  in future periods in that district. This is in line with popular role model theories (see, for example, Chung (2000)), according to which witnessing certain members of an underrepresented group achieving success would lead other group members to achieve higher success in the future. After pupils/students have completed their schooling/university period, they enter the labor market, which is assumed to mix pupils/students from all districts.

The successive policy makers are assumed to be benevolent and we study their incentives to implement affirmative action policies during their tenure. We do so using a repeated game setting, with each successive localized policy maker seeking to maximize its local welfare.

In the main part of the paper, we will suppose for simplicity that employers cannot condition wages on group identity. Although it is not necessary for our results to hold, it is in line with many anti-discrimination policies. In a perfectly competitive labor market, each employer pays a worker a wage equal to his expected performance given the district he comes from. The employer does not observe whether the worker benefited from affirmative action or not and can only estimate this performance based on a curriculum vitae (which may be artificially improved by affirmative action), as well as some aggregate statistics describing the average level of affirmative action policy implemented over a range of districts. Paying workers a wage equal to their expected performance thus means that non-beneficiaries of affirmative action will get a wage below their true performance level. These non-beneficiaries can include members of both groups  $A$  and  $B$  since the affirmative action policy typically does not reach all members of the target group  $B$ . We postulate that in such

a case, the worker suffers from a *feeling of injustice* that is proportional to the difference between his true performance (which the worker knows) and his wage.

Although our model is stylized, recall that this depressed wage can be understood, more broadly, as being associated with the devaluation of a worker's diplomas (or even career promotions), which results from the mere possibility that he may have benefited from affirmative action. We believe such a feeling of injustice is very common. In the case of group  $B$ , this feeling can often be associated with the stigmatization felt by workers who did not benefit from affirmative action (or in more practical situations, even by those who did not need the policy in order to be accepted in a school or university), but are yet underrated due to the mere possibility that some members of their group may have benefited from the policy. In the case of group  $A$ , this feeling of injustice is also in line with not being favored by the policy.

In a first-best scenario, this depressed wage given to non-beneficiaries of affirmative action (and the associated feeling of injustice) means that affirmative action should not last permanently. The optimal duration would be determined by a number of parameters, namely by the weights in welfare assigned to members of the main and the targeted groups, the propensity of non-beneficiaries to experience the feeling of injustice, etc. However, affirmative action would necessarily be ended at some point, as long as non-beneficiaries suffer some (even very small) feeling of injustice in the long-term.

To the contrary, we show that decentralized policy makers *always* choosing to implement affirmative action policies is the unique equilibrium. The intuition is that, in our setting, affirmative action policies are not observed district by district by employers. Thus, if a policy maker were to not implement an affirmative action policy in some period and in some district, this policy maker could deviate without being observed, implement the policy, and this would have no effect on depressing wages. Since such a deviation would be anticipated to improve the future performance distribution of the targeted group (through a role model argument), the policy maker would do it, thereby showing that affirmative action policies are perpetually implemented in all districts. In other words, the non-transparency of the affirmative action policies creates a moral hazard environment, by which each policy maker necessarily chooses to implement an affirmative action policy and fails to internalize the effect that it has on devaluing diplomas (and thus on depressing wages).

We believe that our non-transparency assumption is justified when affirmative action decisions are implemented at a decentralized level as considered in our model, as it is often very difficult in practice to determine whether a specific policy maker actually implemented an affirmative action policy or not. For example, in the United States, these policies are complex, they vary from state to state, even from school to school, and when they are not officially implemented, they may actually take place through private channels (e.g. non-governmental diversity enhancement programs, etc.).

The paper is organized as follows. In Section 2, we introduce the basic setting and define the workers' utilities and welfare. In Section 3, we study how employers set the wages they pay to workers and show that it leads to a feeling of injustice felt by non-beneficiaries of affirmative action (of both groups  $A$  and  $B$ ). In Section 4, we analyze each policy maker's welfare maximization problem and present the two central results: (i) perpetual affirmative action as an equilibrium policy and (ii) the first-best policy where affirmative action is ultimately ended. In Section 5, we discuss how our assumptions can be relaxed, as well as model extensions. We also compare our model with the existing literature. Proofs are relegated to Section 6. A supplementary appendix in Section 7

extends our model to an even more general setting allowing for strategic behavior by workers.

## 2 Setting

There is a continuum  $J$  of districts (or jurisdictions), indexed by  $j$ , where  $j$  is uniformly distributed on  $(0, \bar{J})$ . At each time  $t \in \mathbb{N}$ , district policy makers must each decide whether to implement an affirmative action policy in their district for the duration of their tenure (one period). That is, the policy maker of district  $j \in J$  chooses an action  $\sigma_t^j \in \{0, 1\}$ , where  $\sigma_t^j = 0$  corresponds to no affirmative action and  $\sigma_t^j = 1$  corresponds to affirmative action. One can think of policy makers as local government representatives or as private authorities such as school principals. In the following, we will be assuming that policy makers' interests are aligned with total welfare so that the inefficiencies we highlight cannot be attributed to conflicts of interests.

In each district, a population of workers consists of two groups: group  $A^j$  (the main group) and group  $B^j$  (the targeted group). A worker has a performance level  $c \in [0, 1]$ . This can be understood, for instance, as his result in a standardized university admission test.

At any time  $t$ , group  $A^j$ 's performance density is  $f_{A^j}(c)$  while group  $B^j$ 's performance<sup>1</sup> density is  $f_{B^j, n_t^j}(c)$ , where  $n_t^j = \sum_{s < t} \sigma_s^j$  is the number of times previous policy makers have implemented affirmative action policies in district  $j$ .  $f_{A^j}(c)$  and  $f_{B^j, n_t^j}(c)$  have support  $[0, 1]$  and are non-degenerate. We will describe later how  $f_{B^j, n_t^j}(c)$  varies with  $n_t^j$  but intuitively as  $n_t^j$  increases,  $f_{B^j, n_t^j}(c)$  shifts lower values of  $c$  to higher values, resulting in first-order stochastic dominance. Each agent lives for only one period<sup>2</sup>. At each time  $t$ , a mass  $|A^j|$  and a mass  $|B^j|$  of new agents from groups  $A^j$  and  $B^j$  respectively are born in district  $j$  to replace the ones that have expired, with performance levels drawn according to  $f_{A^j}(c)$  and  $f_{B^j, n_t^j}(c)$ .

### 2.1 Effect of affirmative action policy

An affirmative action policy has two effects. First, it gives an immediate artificial boost to the curriculum vitae of a worker benefiting from it. This models the fact that a beneficiary of affirmative action has expanded opportunities in terms of education (university admissions or other professional formations) compared to a non-beneficiary, thereby artificially enhancing the quality of his curriculum vitae. Second, it is also anticipated by policy makers to have long-term, positive effects on the performance distribution of group  $B^j$ . This anticipated long-term effect is in line with popular role model theories (e.g. Chung (2000)). This second effect will be captured by the dependence of  $f_{B^j, n_t^j}(c)$  on  $n_t^j$ .

It is important to note that an affirmative action policy can be interpreted<sup>3</sup> as anything that artificially increases the quality of a curriculum vitae (immediate effect) and improves the performance distribution of future generations (anticipated role model effect).

In a given period  $t$  where it is implemented, we will allow the affirmative action policy to only reach a fraction  $\xi \in (0, 1]$  of the members of the targeted group  $B^j$ . Indeed, in practice, not all

<sup>1</sup>In the applications we will have in mind, it is reasonable to think that group  $A^j$ 's performance distribution initially differs from that of group  $B^j$ , although this plays no role in our analysis.

<sup>2</sup>The model can easily be extended to allow agents to live for more than one period and to have overlapping generations. Since such elaborations would play no role in our analysis, we have chosen the simpler setting in which agents just live for one period.

<sup>3</sup>See Section 5.1.3 for a discussion of how our model can accommodate even more general interpretations of affirmative action.

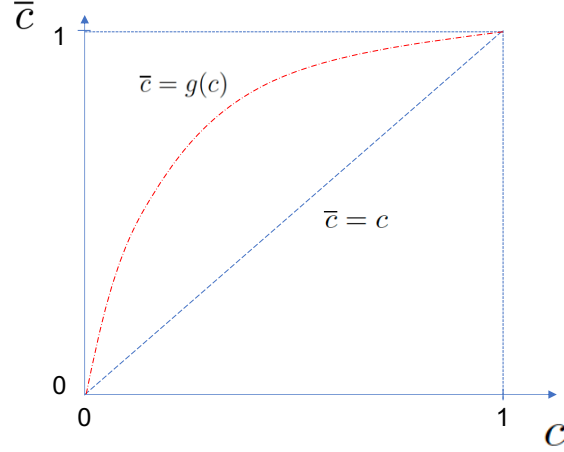


Figure 1: Effect of affirmative action on curriculum vitae quality  $\bar{c}$ . A beneficiary of affirmative action has curriculum vitae quality higher than his actual performance level:  $\bar{c} = g(c) > c$  (red curve). A non-beneficiary has curriculum vitae quality corresponding to his actual performance level:  $\bar{c} = c$  (blue line).

members of a targeted group may benefit from the policy<sup>4</sup>.

### 2.1.1 Effect of affirmative action policy on curriculum vitae quality

When  $\sigma_t^j = 1$ , with probability  $\xi \in (0, 1]$ , a member of group  $B^j$  with performance level  $c \in [0, 1]$  will have a curriculum vitae quality  $\bar{c} = g(c)$ , where  $g$  is an increasing function such that  $g(c) > c$ ,  $\forall c \in (0, 1)$ , and  $g(0) = 0$ ,  $g(1) = 1$ . The support of  $\bar{c}$  is thus also  $[0, 1]$ . With probability  $1 - \xi$ , a member of group  $B^j$  with performance level  $c$  will have a curriculum vitae quality corresponding to his actual performance level:  $\bar{c} = c$ .

When  $\sigma_t^j = 0$ , a member of group  $B^j$  with performance level  $c$  will have a curriculum vitae quality corresponding to his actual performance level:  $\bar{c} = c$ .

Whether  $\sigma_t^j = 0$  or 1, a member of group  $A^j$  with performance level  $c$  always has a curriculum vitae quality corresponding to his actual performance level:  $\bar{c} = c$ .

An affirmative action policy therefore increases the curriculum vitae quality of a beneficiary above his actual performance level, while it has no effect on the curriculum vitae quality of members of group  $A^j$  nor on those of members of group  $B^j$  who did not benefit from the affirmative action policy. That is, their curriculum vitae quality corresponds to their actual performance level. This is illustrated in Figure 1.

### 2.1.2 Effect of affirmative action policy on actual performance

We suppose that if  $\sigma_t^j = 1$ , then the next period's performance distribution of group  $B^j$  is shifted so that  $f_{B^j, n_{t+1}^j}(c) \succ f_{B^j, n_t^j}(c)$ , where  $\succ$  indicates strong first-order stochastic dominance. Note that the effect of the shift is permanent, i.e. the improvement remains in all future periods. This purported improvement in the performance of future cohorts of workers is consistent with the role model argument.

<sup>4</sup>Note that our results require that not all members of the targeted group may be reached, i.e.  $\xi < 1$ , only when considering the case in which wages are allowed to depend on the group identity. More on this in Section 5.1.1.

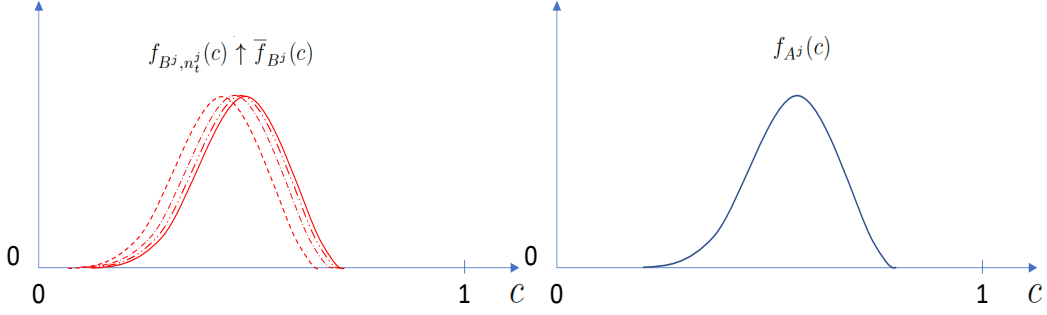


Figure 2: Effect of affirmative action policy on actual performance. If  $\sigma_t^j = 1$ , then the next period's performance distribution of group  $B^j$  is shifted so that  $f_{B^j, n_{t+1}^j}(c) \succ f_{B^j, n_t^j}(c)$ . If  $\sigma_t^j = 1$  for all  $t$ , then  $f_{B^j, n_t^j}(c)$  converges to a limiting distribution  $f_{B^j, n_t^j}(c) \uparrow \bar{f}_{B^j}(c)$ . Group  $A^j$ 's performance distribution  $f_{A^j}(c)$  is not affected.

If  $\sigma_t^j = 1$  for all  $t$ , then  $f_{B^j, n_t^j}(c) \uparrow \bar{f}_{B^j}(c)$ . Since  $f_{B^j, n_t^j}(c)$  converges from below to a limiting distribution  $\bar{f}_{B^j}(c)$ , this implies that the distributional improvements become smaller and smaller as policy makers keep implementing affirmative action policies. Group  $A^j$ 's performance distribution  $f_{A^j}(c)$  does not vary with  $t$ . This is illustrated in Figure 2. Observe that the densities  $f_{B^j, n_t^j}(c)$  could depend on  $\xi$ , as the larger  $\xi$  the more individuals in group  $B^j$  are likely to be exposed to the effect of the affirmative action policy in district  $j$ . Since our results do not rely on varying  $\xi$ , we omit an explicit reference to this dependence.

## 2.2 Utilities and welfare

A worker is of type  $\theta = (c, \bar{c}, G^j)$ , where  $c$  is his true performance level,  $\bar{c}$  is his curriculum vitae quality and  $G^j \in \{A^j, B^j\}$  is the group  $G$  this worker belongs to and the district  $j$  he comes from. A time  $t$  worker knows his type and the wage function  $\omega_t^j(\bar{c})$  set by employers, which is the wage the worker earns based on the information on his curriculum vitae (i.e. the curriculum vitae quality  $\bar{c}$  and the district  $j$  the worker comes from).<sup>5</sup> This is formalized in the following definition.

**Definition 1** A wage function  $\omega_t^j : [0, 1] \rightarrow [0, 1]$  determines, at time  $t$ , the wage a worker coming from district  $j$  earns when presenting a curriculum vitae of quality  $\bar{c}$  to an employer.

Note here that we chose not to allow employers to condition wages on the group  $A$  or  $B$  to which a worker belongs. This is motivated on grounds that such group-based discrimination is in general forbidden. Our results are however robust to conditioning wages on group identity, i.e. giving a wage  $\omega_t^{G^j}$  instead of  $\omega_t^j$ . This is further discussed in section 5.1.1.

### 2.2.1 Utility

The utility of a type  $(c, \bar{c}, G^j)$  worker at time  $t$  is

$$u_{G^j, t}(\bar{c}, c) = \omega_t^j(\bar{c}) - \gamma_{G^j} \max\{c - \omega_t^j(\bar{c}), 0\} \quad (1)$$

<sup>5</sup>If workers were to live several periods, we could envision a more elaborate model in which the wage earned in later periods would also depend on the true performance assumed to be partly observed then. Our qualitative insights would be unaffected.

where  $\gamma_{G^j} \max\{c - \omega_t^j(\bar{c}), 0\}$ , for some  $\gamma_{G^j} > 0$ , captures the fact that a feeling of “injustice” is suffered when a worker gets a salary that is below his true performance level. Note that we allow  $\gamma_{A^j} \neq \gamma_{B^j}$  so as to capture that the feeling of injustice may differently affect groups  $A$  and  $B$  in district  $j$ .

In particular, the utility of a type  $(c, \bar{c}, G^j)$  worker who benefits from affirmative action has the form

$$u_{G^j,t}(\bar{c}, c) = \omega_t^j(g(c)) - \gamma_{G^j} \max\{c - \omega_t^j(g(c)), 0\}$$

since  $\bar{c} = g(c)$ , while the utility of a type  $(c, \bar{c}, G^j)$  worker who does not benefit from affirmative action has the form

$$u_{G^j,t}(\bar{c}, c) = \omega_t^j(c) - \gamma_{G^j} \max\{c - \omega_t^j(c), 0\}$$

since  $\bar{c} = c$ . We will often denote by  $u_{B^j,t}(g(c), c)$  (respectively, by  $u_{B^j,t}(c, c)$ ) the utility of a group  $B^j$  worker benefiting (respectively, not benefiting) from affirmative action, while we will denote by  $u_{A^j,t}(c, c)$  the utility of a group  $A^j$  worker.

In the above, we assume that workers have the correct perception of their performance level  $c$ . We also note that there are no extra positive effects on utility of receiving a wage greater than the performance level. Such an asymmetry in the utility assessment of wages above or below the performance level is in line with well documented psychological studies (see in particular the prospect theory of Kahneman and Tversky (1979)), which suggest a different assessment for payoff realizations above or below the reference point (here naturally identified with the performance level).

### 2.2.2 Welfare

The welfare of each group in district  $j$  at time  $t$  is defined by taking the aggregate utility of that group. We thus have,

$$W_{A^j,t} = |A^j| \int_0^1 u_{A^j,t}(c, c) f_{A^j}(c) dc$$

$$W_{B^j,t} = |B^j| \int_0^1 \left( \xi \sigma_t^j u_{B^j,t}(g(c), c) + (1 - \xi \sigma_t^j) u_{B^j,t}(c, c) \right) f_{B^j, n_t^j}(c) dc$$

where  $\sigma_t^j$  is the actual policy decision made by the time  $t$  policy maker of district  $j$ .

Total welfare in district  $j$  at time  $t$  is defined by

$$W_t^j = W_{A^j,t} + \lambda_{B^j} W_{B^j,t}$$

where the weight on  $B^j$ ,  $\lambda_{B^j}$ , is non-negative and typically no greater than 1.  $\lambda_{B^j} = 1$  corresponds to the standard total welfare criterion and  $\lambda_{B^j} < 1$  reflects a preference for the main group  $A^j$  in the policy maker’s objective. We will also comment on the case when  $\lambda_{B^j} > 1$ , which reflects a preference for the targeted group  $B^j$ .

Letting  $\delta$  denote the common discount factor, total welfare in district  $j$  over all periods is then

defined by

$$W^j = \sum_{t=1}^{\infty} \delta^t W_t^j$$

and total welfare in the economy is defined by

$$W = \int_{j \in J} W^j dj.$$

### 3 Effect of affirmative action policy on wage levels

We model a fully mixing labor market, where workers educated in all districts match freely with employers and are paid wages by the latter.

#### 3.1 Informational environment

While it is plausible to assume that employers observe some aggregate statistics about the decentralized affirmative action policy decisions, we believe that in many applications it is natural to assume that employers do not observe each  $\sigma_t^j$  separately. Indeed, in the face of a large number of districts, it would be very difficult to keep track of all decentralized policy decisions.

To formalize this idea most simply, we assume that employers at time  $t$  can observe an aggregate statistic  $\bar{\sigma}_t$  of all policy decisions made by the different districts. Here we take  $\bar{\sigma}_t$  to be simply the average policy across all districts, that is

$$\bar{\sigma}_t = \frac{1}{|J|} \int_{j \in J} \sigma_t^j dj. \quad (2)$$

Thus, they know the sequence  $\{\bar{\sigma}_s\}_{s=1}^t$  of all average policy decisions made over time (up to time  $t$ ).

Not knowing for sure whether affirmative action took place in a particular district, they may not be able to tell for certain whether a worker benefited from affirmative action or not. They can however compute the probability that a worker benefited from affirmative action, conditional upon observing his curriculum vitae quality, the district this worker comes from (e.g., where he graduated school or university), the aggregate policy statistics and considering a putative strategy played by policy makers.

Note that the form of the observed aggregate statistic in Eq. (2) can be generalized. For example, an employer could observe a more localized average of the policies practiced around district  $j$ , such as  $\bar{\sigma}_t^j = \frac{1}{2\epsilon} \int_{i=j-\epsilon}^{j+\epsilon} \sigma_t^i di$ . What is key is that no inference can be made from the observed statistics about the value of a given  $\sigma_t^j$ , which sounds plausible when the number of districts is very large. We discuss this further in Section 5.1.2.

#### 3.2 Setting wages

We consider a perfectly competitive labor market, where an employer pays a worker a wage equal to his expected performance level. In Section 5.1.5, this reduced-form approach is micro-founded based on a Bertrand-type model of competition between employers.

As mentioned earlier, we assume that employers are not allowed to take group information ( $A$  or  $B$ ) into account when giving a wage to a particular worker. This is consistent with anti-



discrimination laws enacted in many countries and occupational areas (although it is not necessary for our results to hold, as previously mentioned and as discussed in Section 5.1.1). Thus, they set a wage conditioned only on the curriculum vitae quality  $\bar{c}$ , the district  $j$  a worker is from, the observed sequence of aggregate policy statistics  $\{\bar{\sigma}_s\}_{s=1}^t$  and a given putative policy sequence assumed by them  $\sigma = \{\{\sigma_s^j\}_{j \in J}\}_{s=1}^\infty$ . The wage  $\omega_t^j(\bar{c})$  paid to a worker of type  $(c, \bar{c}, A^j)$  or to a worker of type  $(c, \bar{c}, B^j)$  is thus the conditional expectation  $\mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$  of the worker's true performance level  $c$ , expressed in the following lemma.

**Lemma 1 (Wage function)** *Given some putative policy sequence  $\sigma = \{\{\sigma_s^j\}_{j \in J}\}_{s=1}^\infty$  and an observation of aggregate policy statistics  $\{\bar{\sigma}_s\}_{s=1}^t$  consistent with  $\sigma$ , the wage paid at time  $t$  to a worker with curriculum vitae quality  $\bar{c}$  coming from district  $j$  has the form*

$$\omega_t^j(\bar{c}) = \mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) \cdot g^{-1}(\bar{c}) + (1 - \mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma)) \cdot \bar{c}$$

where

$$\mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) = \frac{|B^j| \xi \sigma_t^j f_{B, n_t^j}(g^{-1}(\bar{c}))/g'^{-1}(\bar{c})}{|A^j| f_A(\bar{c}) + |B^j| (1 - \xi \sigma_t^j) f_{B, n_t^j}(\bar{c}) + |B^j| \xi \sigma_t^j f_{B, n_t^j}(g^{-1}(\bar{c}))/g'^{-1}(\bar{c})}.$$

and  $\{aa\}$  is the event that a worker benefited from affirmative action.

In words,  $\omega_t^j(\bar{c})$  is a convex combination between  $g^{-1}(\bar{c})$  and  $\bar{c}$ , where the weight assigned to  $g^{-1}(\bar{c})$  is the probability that a worker with curriculum vitae  $\bar{c}$  and coming from district  $j$  benefited from affirmative action at time  $t$  (taking into account the policy sequence  $\sigma$  believed to be followed by policy makers, hence the expression for  $\mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma)$ ).

We make the following assumption on  $\mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$  for simplicity of exposition. Our results do not depend on it, but it will allow us to present them in a simpler manner, since we can rule out strategic behavior by which an agent could present a curriculum vitae of lower<sup>6</sup> quality than  $\bar{c}$ . An extension where a worker is allowed to present a curriculum vitae of a different quality than  $\bar{c}$  is presented in a supplementary appendix (Section 7.1), where the robustness of our results to such strategic behavior is established in a more general context, and which thus removes the need for Assumption 1.

**Assumption 1** *The conditional expectation  $\mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$ , and thus the wage function  $\omega_t^j(\bar{c})$ , is non-decreasing in  $\bar{c}$ .*

Although this assumption may appear, at first sight, to depend on  $\sigma$  (more specifically  $\sigma_t^j$ ), which is endogenous, it is actually easily satisfied for *any*  $\sigma$  under some conditions, e.g. when the likelihood ratio  $\frac{f_{A^j}(c)}{f_{B^j}(g(c))}$  is increasing or when the mass  $|A^j|$  is sufficiently larger than the mass  $|B^j|$ .

Relying on the expression of equilibrium wage derived in Lemma 1, we note that whether the earned wage lies above or below the performance level solely depends on whether or not the worker benefited from affirmative action:

**Lemma 2 (Wage versus performance level)**

<sup>6</sup>Indeed, if the wage function  $\omega_t^j(\bar{c}) = \mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$  is decreasing on some parts of the support  $[0, 1]$ , a worker could earn a higher wage by presenting a curriculum vitae of lower quality than  $\bar{c}$ .

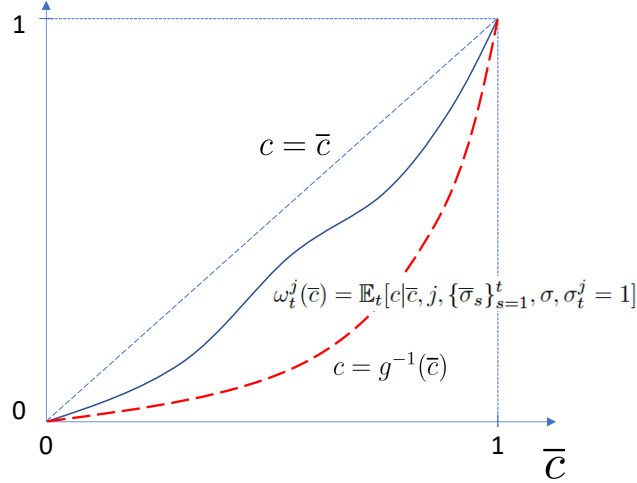


Figure 3: Illustration of Lemma 2. The curriculum vitae quality  $\bar{c}$  is on the horizontal axis. (i) When affirmative action is implemented ( $\sigma_t^j = 1$ ), we see on the vertical axis that the wage  $\omega_t^j(\bar{c}) = \mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma, \sigma_t^j = 1] < \bar{c}$  (full blue curve) is lower than the performance level of a non-beneficiary (thinly dotted blue line) and higher than the performance level of a beneficiary (thickly dotted red curve). (ii) When no affirmative action is implemented ( $\sigma_t^j = 0$ ), then the wage  $\omega_t^j(\bar{c}) = \mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma, \sigma_t^j = 0] = \bar{c}$  corresponds to the performance level of any worker (i.e. thinly dotted blue line).

- (i) Suppose  $\sigma_t^j = 1$ . Then any district- $j$  worker gets a wage lower than his curriculum vitae quality (i.e.  $\bar{c} > \omega_t^j(\bar{c})$ ). Moreover, a worker benefiting from affirmative action gets a wage higher than his performance level (i.e.  $c = g^{-1}(\bar{c}) < \omega_t^j(\bar{c})$ ), while a worker not benefiting from affirmative action gets a wage lower than his performance level (i.e.  $c = \bar{c} > \omega_t^j(\bar{c})$ ).
- (ii) Suppose  $\sigma_t^j = 0$ . Then any district- $j$  worker gets a wage equal to his curriculum vitae quality and his performance level (i.e.  $c = \bar{c} = \omega_t^j(\bar{c})$ ).

This lemma is illustrated in Figure 3.

### 3.3 Feeling of injustice and broader interpretation of the depressed wage

In our model, the wage is depressed due to the possibility that a worker benefited from affirmative action. This represents the fact that a certain curriculum vitae quality is, in expectation, no longer associated with the same performance level as if there were no affirmative action policy. Indeed, an affirmative action policy has the effect of devaluing the diplomas or promotions that figure on a worker's curriculum vitae, if there is only some chance that the worker may have benefited from such a policy.

Using Lemma 2, we now make the following observation.

#### Observation 1 (Feeling of injustice)

- (i) Suppose  $\sigma_t^j = 1$ .

The utility of a type  $(c, \bar{c}, G^j)$  worker not benefiting from affirmative action can be written as

$$\begin{aligned} u_{G^j,t}(c, c) &= \omega_t^j(c) - \gamma_{G^j} \max\{c - \omega_t^j(c), 0\} \\ &= \omega_t^j(c) - \gamma_{G^j}(c - \omega_t^j(c)) \end{aligned}$$

since  $\bar{c} = c$  and  $\omega_t^j(c) < c$ . Such a worker gets a wage lower than his performance level and suffers a feeling of injustice.

By contrast, the utility of a type  $(c, \bar{c}, G^j)$  worker benefiting from affirmative action can be written as

$$\begin{aligned} u_{G^j,t}(g(c), c) &= \omega_t^j(g(c)) - \gamma_{G^j} \max\{c - \omega_t^j(g(c)), 0\} \\ &= \omega_t^j(g(c)) \end{aligned}$$

since  $\bar{c} = g(c) > c$  and  $\omega_t^j(g(c)) > c$ . Such a worker gets a wage higher than his performance level and does not suffer a feeling of injustice.

(ii) Suppose  $\sigma_t^j = 0$ .

The utility of a type  $(c, \bar{c}, G^j)$  worker can be written as

$$\begin{aligned} u_{G^j,t}(c, c) &= \omega_t^j(c) - \gamma_{G^j} \max\{c - \omega_t^j(c), 0\} \\ &= c \end{aligned}$$

since  $\bar{c} = c$  and  $\omega_t^j(c) = c$ . Such a worker gets a wage equal to his performance level and does not suffer a feeling of injustice.

It is important to emphasize that workers from both groups ( $A$  and  $B$ ) can experience a feeling of injustice. In the case of group  $B$ , this feeling can often be associated with the stigmatization felt by workers who did not benefit from affirmative action (or in more practical situations, even those who did not need it in order to be accepted in a school or university), but are yet underrated due to the mere possibility that some members of their group may have benefited from the policy. This is suggested by a good deal of empirical evidence (see, for example, Leslie et al. (2014), Heilman et al. (1997) or Heilman et al. (1992)).

## 4 The policy maker's decision problem

### 4.1 Informational environment

At any time  $t$ , the policy maker of district  $j$ , knows the purported effect of  $n_t^j$  on  $f_{B^j, n_t^j}(c)$  and thus anticipates that choosing an affirmative action policy  $\sigma_t^j = 1$  will improve the future performance distribution of group  $B^j$ :  $f_{B^j, n_{t+1}^j}(c) \succ f_{B^j, n_t^j}(c)$ .

### 4.2 Policy decisions

For all  $t \geq 1$ , a district- $j$  policy maker wants to choose a policy  $\sigma_t^j$  in order to maximize the following objective function:

$$\max_{\sigma_t^j \in \{0,1\}} \sum_{s=t}^{\infty} \delta^{s-t} (W_{A^j,s} + \lambda_{B^j} W_{B^j,s}) \quad (3)$$

That is, we assume that the objective of the district- $j$  policy maker coincides with the local welfare as aggregated over the remaining time periods.

Given some putative policy sequence  $\sigma = \{\{\sigma_s^j\}_{j \in J}\}_{s=1}^{\infty}$  followed by policy makers across time, a district- $j$  policy maker<sup>7</sup> is able to compute  $W_{A^j,s}$  and  $W_{B^j,s}$  for  $s > t$ , where  $\omega_s^j(\bar{c})$  and  $f_{B^j,n_s^j}(c)$  are taken to be consistent with  $\sigma$ .

Our first main result, Proposition 1, states that all district policy makers choosing to *perpetually* implement an affirmative action policy is the unique equilibrium.

**Proposition 1 (Permanent affirmative action in equilibrium)** *Let  $\lambda_{B^j} > 0$ . Then there exists  $\bar{\gamma}_{B^j}$  such that for any  $\gamma_{B^j} < \bar{\gamma}_{B^j}$  the unique equilibrium is  $\sigma_t^{j*} = 1$  for all  $t$  and  $j$ .*

The intuition behind Proposition 1 is that any district- $j$  policy maker anticipates that implementing an affirmative action policy improves the performance distribution of future cohorts of  $B^j$  workers. Thus, the only reason it would choose not to implement such a policy would be to have an uplifting effect on the wage function  $\omega_t^j(\bar{c})$  chosen by employers (which, as we know, is depressed by the possibility that a worker has benefited from affirmative action). However, due to its small size (measure zero), a district- $j$  policy maker cannot have any impact on the aggregate (average) statistic  $\bar{\sigma}_t$ , which is the only policy information observed by employers when setting wages. Therefore, there is no reason why a particular policy maker would deviate, by choosing  $\sigma_t^j = 0$ , from an equilibrium policy  $\sigma_t^{j*} = 1$  in which it implements an affirmative action policy. Conversely, a deviation from a putative equilibrium in which  $\sigma_t^{j*} = 0$  to  $\sigma_t^j = 1$  would increase the average performance of future cohorts of  $B^j$  workers, without having a worsening impact on the wage, since that deviation will not be reflected in the aggregate statistic  $\bar{\sigma}_t$  observed by employers and thus on the wage function  $\omega_t^j(\bar{c})$ . This establishes  $\sigma_t^{j*} = 1$  as the unique equilibrium.

The intuition behind the sufficient (and not necessary) condition  $\gamma_{B^j} < \bar{\gamma}_{B^j}$  is that although a deviation from a putative equilibrium in which  $\sigma_t^{j*} = 0$  to  $\sigma_t^j = 1$  would increase the average performance of future cohorts of  $B^j$  workers, it could also potentially increase the average feeling of injustice felt by  $B^j$  workers not benefiting from affirmative action in future periods. Indeed, the feeling of injustice could worsen following an increase in the performance level, if the latter increases faster than the wage received at a higher performance level (recall that the feeling of injustice is  $\gamma_{B^j} \max\{c - \omega_s^{j*}(c), 0\}$ ). A sufficient condition for the positive effect to dominate the negative one is that the parameter  $\gamma_{B^j}$  be small enough<sup>8</sup>.

Our second main result, Proposition 2, states that in the first-best scenario, affirmative action policies *always* end after a finite number of periods.

<sup>7</sup>In the welfare, we have not included the firms' profits. Note however that these profits are null on the equilibrium path, due to our assumption of perfect competition. Including firms' profits in the policy maker's objective would affect the assessment of deviations, but the qualitative insights presented below would be unaffected.

<sup>8</sup>It is interesting to note that it is enough that such a parameter  $\gamma_{B^j}$ , capturing the feeling of injustice felt by members of group  $B^j$  not benefiting from affirmative action, corresponds to one chosen by the policy maker and it need not be the actual one felt in population  $B^j$ . Indeed, recall that the equilibrium wage  $\omega_s^{j*}$  actually does not depend on  $\gamma_{B^j}$ . Only the welfare  $W_{B^j,s}$  of group  $B^j$  does.

**Proposition 2 (Temporary affirmative action as first-best policy)** *Suppose that at time  $t = 0$ , a single centralized policy maker announces (and commits to) the policy plan  $\hat{\sigma} = \{\{\hat{\sigma}_t^j\}_{j \in J}\}_{t=1}^{\infty}$  that maximizes the welfare function  $\sum_{t=1}^{\infty} \int_{j \in J} \delta^t (W_{A^j,t} + \lambda_{B^j} W_{B^j,t}) dj$ , and assume  $\gamma_{A^j} \neq 0$  (or likewise  $\gamma_{B^j} \neq 0$  when  $\xi < 1$ ). Then for any  $\lambda_{B^j} \in [0, 1]$ , there exists  $\bar{\delta} \in (0, 1)$  such that for all  $\delta \in (\bar{\delta}, 1)$ ,  $\hat{\sigma}$  has a threshold form:  $\hat{\sigma}_t^j = 1$  for  $t < \bar{T}^j$  and  $\hat{\sigma}_t^j = 0$  for  $t \geq \bar{T}^j$ , for some (finite)  $\bar{T}^j \in \mathbb{N}$ .*

Proposition 2 essentially means that if different policy makers were able to coordinate their actions over time periods so as to maximize global welfare, they would never choose to make affirmative action permanent. The intuition is quite simple: After a certain number of periods the improvement in the performance distribution becomes marginal, while the depressing effect on wages (corresponding to curricula vitae being devalued) is not. As a matter of fact,  $f_{B^j, n_t^j}(c)$  converges from below to a limiting distribution  $\bar{f}_{B^j}(c)$ , implying that the distributional improvements become smaller and smaller as affirmative action policies are implemented over time.

The optimal threshold  $\bar{T}^j$ , while always finite, depends on the relative weight placed by policy makers on the welfare of the targeted group  $B^j$  relative to the main group  $A^j$ , i.e. on  $\lambda_{B^j}$ , as well as on the intensity of the feeling of injustice felt by non-beneficiaries of both groups, i.e. on  $\gamma_{A^j}$  and  $\gamma_{B^j}$ , and on the fraction  $\xi$  of group  $B^j$  reached by the affirmative action policy.

Note that when  $\lambda_{B^j} < 1$  (i.e. when the policy maker cares relatively more about group  $A^j$  than group  $B^j$ ), the parameter  $\gamma_{A^j}$  governing the feeling of injustice of group  $A^j$  can be 0 and the first-best policy will still prescribe stopping affirmative action after a finite number of periods, because the depressed wage penalizes group  $A^j$  sufficiently while the performance distribution of group  $B^j$  is only marginally improved.

When  $\lambda_{B^j} = 1$  (i.e. when the policy maker cares equally about group  $A^j$  and group  $B^j$ ), then since the average wage is equal to the average performance level across the district (i.e.  $\mathbb{E}_t[\omega_t^j(\bar{c})] = \mathbb{E}_t[c|j]$ ), an affirmative action policy effectively represents just a transfer of welfare from the non-beneficiaries to the beneficiaries. Indeed, this transfer of welfare takes place through non-beneficiaries of both groups  $A^j$  and  $B^j$  receiving wages lower than their performance levels while beneficiaries receive wages higher than their performance levels. In this case, as long as the parameter  $\gamma_{A^j}$  (or likewise  $\gamma_{B^j} \neq 0$  when  $\xi < 1$ ) is *strictly greater* than 0 (no matter how small it is), a first-best policy will prescribe stopping affirmative action after a finite number of periods because otherwise the feeling of injustice felt by non-beneficiaries would become worse than the improvement in the performance distribution of group  $B^j$  after sufficiently many implementations of the affirmative action policy.

Finally, if  $\lambda_{B^j}$  were to be strictly greater than 1 (i.e. when the policy maker cares relatively more about group  $B^j$  than group  $A^j$ ), then we might need  $\gamma_{A^j}$  to be sufficiently positive in order to justify stopping affirmative action in the case when  $\xi = 1$  (i.e. when all group  $B^j$  members benefit from affirmative action). Else, when  $\xi < 1$ , any feeling of injustice felt by non-beneficiaries of group  $B^j$  (i.e.  $\gamma_{B^j} \neq 0$ ) will justify stopping affirmative action at some point.

## 5 Some extensions

### 5.1 Discussion of assumptions

In the next subsections, we show that our main results are quite robust and most often hold, even if we relax the assumptions made in the main part of the paper. We explain that all we really need for our main results to hold is that an employer cannot be certain that a worker from group  $B^j$  has not benefited from affirmative action.

#### 5.1.1 Allowing employers to condition wages on group identity

In our main model, we have not allowed employers to condition wages on the group  $A$  or  $B$  to which a worker belongs. This was motivated on grounds that such discrimination is in general forbidden. If conditioning wages on group identity were allowed, our results would actually hold as long as some members of the targeted group  $B$  do not benefit from affirmative action (a fairly weak assumption). In such a case, the feeling of injustice is suffered entirely by them (and not also by members of group  $A$ ) and this is enough for our results to hold.

More specifically, a group  $G^j$  worker would receive a wage  $\omega_t^{G^j}(\bar{c}) = \mathbb{E}_t[c|\bar{c}, G^j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$ . In this case, a group  $A^j$  worker would receive a wage  $\omega_t^{A^j}(\bar{c}) = \bar{c} = c$  equal to his performance level, since group  $A$  workers do not benefit from affirmative action. Group  $A$  workers would then suffer no feeling of injustice and be unaffected by the affirmative action policy. When affirmative action only reaches a fraction  $\xi \in (0, 1)$  of group  $B^j$ , then the share  $1 - \xi$  of group  $B^j$  workers not benefiting from affirmative action would still get a wage  $\omega_t^{B^j}(\bar{c}) < \bar{c} = c$  lower than their performance level and thus suffer a feeling of injustice. The share  $\xi$  of group  $B^j$  workers benefiting from the policy would still get a wage  $\omega_t^{B^j}(\bar{c}) > g^{-1}(\bar{c}) = c$  higher than their performance level and thus would not suffer a feeling of injustice.

We see that allowing employers to condition wages on the group  $A$  or  $B$  to which a worker belongs leads to similar insights as in the main model, although the key tensions now take place entirely within group  $B$ .

The particular case of  $\xi = 1$ , that is when affirmative action reaches *all* members of group  $B^j$ , is worth commenting on. In this particular case,  $\omega_t^{B^j}(\bar{c}) = g^{-1}(\bar{c}) = c$ , and workers from both groups  $A$  and  $B$  would receive wages equal to their performance levels, as there would be no uncertainty as to whether they benefited from the policy or not. As a result, the equilibrium would coincide with first-best.

#### 5.1.2 Variations on the informational assumptions

We could consider several variations in what we assume employers can condition their wages on.

As already said, in our environment with a continuum of districts, assuming that employers observe the average policy over various ranges of districts would not affect the analysis, as long as the averages are taken over positive Lebesgue measures of districts. This is so because a deviation by a single district policy maker would not affect such statistics, even if finer than the one considered in the main model.

From another perspective, one might request that employers do not condition their wage on the district  $j$  the worker comes from. Such constraints may be the result of anti-discrimination consid-

erations, this time based on location rather than group membership, as it may reasonably be argued that employers would naturally have access to where the worker completed his studies. Our equilibrium analysis would be unchanged in this setting. One might have thought that having the same wage across districts would create an additional externality between the policy makers of the various districts, to the extent that a choice of policy in some districts could now adversely affect the wages received by workers in other districts. However, in our setting where only aggregates are observed, the equilibrium already involves perpetual affirmative action even without this externality. It then follows that the same property of perpetual affirmative actions would a fortiori hold when such an externality is present<sup>9</sup>.

The nature of the first-best policy (a threshold form) would be unaffected by such modifications. This establishes the robustness of our main insights with respect to a broad class of observational environments.

### 5.1.3 Affirmative action as a biased promotion process

The leading interpretation of our model so far is that affirmative action takes the form of favoring in their school/university studies (some share of the) members of group  $B$  as opposed to members of group  $A$ , and our model has emphasized the decentralized nature of affirmative action decisions so as to motivate our key informational assumption that affirmative action policies are not observable, district by district, by employers.

Sometimes affirmative action is instead thought of in terms of biased promotion (here in favor of group  $B$ ) rather than in terms of biased school admission, and one may wonder whether a logic similar to that developed in our main model would be at work in such contexts. We now suggest that a similar phenomenon may be at work. Think now of the life of a worker as having two phases, the early phase and the mature phase. In the early phase, we assume employment takes place in the worker's own district, while in the mature phase employment takes place in a fully-mixing labor market. That is, we have in mind that workers in their early phase go to the local labor market and then get rematched to new firms in the mature phase, and that this rematching is not localized. It is not difficult to see how our main model would transpose in such a variant. Companies, when considering workers in the early phase, may decide to promote more easily the workers from group  $B$  in an attempt to increase the ability distribution of group  $B$  workers (through increased motivation, say) in the firm. How a given company favors the promotion of  $B$  workers would hardly be known to outsiders, which is in line with our view that, at the rematching stage, it would be difficult to determine whether a  $B$  worker benefitted from an artificial boost in his early career. On the other hand, such biased promotions in the early careers of (some share of)  $B$  workers will depreciate the assessment of early-career promotions in the mature phase, leading to a feeling of injustice among the non-beneficiaries of such biased early-career promotions. Given that a particular firm will not have control on how wages are set in the mature phase, they will unambiguously go for the biased promotions of  $B$  workers in the early phase, similarly as in the main model (and this would overall be inefficient, at least when the gains in the ability distribution of group  $B$  stabilize).

---

<sup>9</sup>By contrast, such an externality would have an effect (resulting in shifting towards more affirmative action policies being implemented) in contexts where the affirmative action policies would be observed, district by district, by employers.

### 5.1.4 Accounting for labor market congestion

Here we introduce a labor market congestion externality caused by affirmative action, which allows us to capture at least to some extent the fact that jobs obtained by beneficiaries of affirmative action are no longer available to non-beneficiaries. We show that our main insights go through even in the presence of such elaborations.

In our model, non-beneficiaries of affirmative action suffer from receiving a wage that is lower than their actual performance level, while beneficiaries of affirmative action receive a wage that is higher than their actual performance level. Therefore, a transfer of utility between beneficiaries and non-beneficiaries arises through the wage channel.

From another perspective, affirmative action is often thought of as an allocation problem, e.g. allocating a finite number of jobs between two groups, which would result in extra transfers between beneficiaries and non-beneficiaries of affirmative action in addition to the wage effect considered in our main model. While modeling a full-scale matching process is beyond the scope of this paper, our model can be extended in such a direction by adding a labor market congestion externality. This will be represented by a positive term in the utility function of a beneficiary and a negative term in the utility function of a non-beneficiary.

In this section, for clarity of exposition, we will suppose that for all  $j$ ,  $|A^j| = \bar{A}$  and  $|B^j| = \bar{B}$ . That is, all districts have the same mass of  $A$  workers and the same mass of  $B$  workers.

The utility of a beneficiary will take the form

$$\begin{aligned}\tilde{u}_{B^j,t}(g(c), c) &= u_{B^j,t}(g(c), c) + \eta \\ &= \omega_t^j(g(c)) + \eta\end{aligned}$$

where  $\eta$  is a parameter measuring the magnitude of the allocation advantage in the labor market (e.g. the advantage of having a reserved slot in the labor market).

It is easy to see that the aggregate transfer of utility from non-beneficiaries to beneficiaries, due to labor market congestion, is simply

$$\int_{j \in J} \bar{B} \xi \sigma_t^j \eta dj = \bar{J} \bar{B} \xi \bar{\sigma}_t \eta,$$

recalling that  $|J| = \bar{J}$  and that  $\bar{\sigma}_t = \frac{1}{|\bar{J}|} \int_{j \in J} \sigma_t^j dj$ .

The utility of a non-beneficiary will then take the form

$$\begin{aligned}\tilde{u}_{G^j,t}(c, c) &= u_{G^j,t}(c, c) - K(\bar{\sigma}_t) \eta \\ &= \omega_t^j(c) - \gamma_{G^j}(c - \omega_t^j(c)) - \frac{\bar{J} \bar{B} \xi \bar{\sigma}_t \eta}{\bar{J}(\bar{A} + \bar{B}(1 - \xi \bar{\sigma}_t))}\end{aligned}$$

where  $K(\bar{\sigma}_t) = \frac{\bar{J} \bar{B} \xi \bar{\sigma}_t}{\bar{J}(\bar{A} + \bar{B}(1 - \xi \bar{\sigma}_t))}$  is a term reflecting the congestion externality faced by a non-beneficiary of affirmative action in the labor market, due to certain slots being reserved for beneficiaries (the total mass of non-beneficiaries being given by  $\bar{J}(\bar{A} + \bar{B}(1 - \xi \bar{\sigma}_t))$ ).



The welfare of non-beneficiaries of group  $A^j$  at time  $t$  is

$$\begin{aligned}\bar{A} \int_0^1 \tilde{u}_{A^j,t}(c, c) f_{A^j} dc &= \bar{A} \int_0^1 (u_{A^j,t}(c, c) - K(\bar{\sigma}_t)\eta) f_{A^j}(c) dc \\ &= W_{A^j,t} - \bar{A}K(\bar{\sigma}_t)\eta.\end{aligned}\quad (4)$$

The welfare of non-beneficiaries of group  $B^j$  at time  $t$  is

$$\begin{aligned}\bar{B}(1 - \xi\sigma_t^j) \int_0^1 \tilde{u}_{B^j,t}(c, c) f_{B^j, n_t^j} dc &= \bar{B}(1 - \xi\sigma_t^j) \int_0^1 (u_{B^j,t}(c, c) - K(\bar{\sigma}_t)\eta) f_{B^j, n_t^j}(c) dc \\ &= \bar{B}(1 - \xi\sigma_t^j) \int_0^1 u_{B^j,t}(c, c) f_{B^j, n_t^j}(c) dc - \bar{B}(1 - \xi\sigma_t^j)K(\bar{\sigma}_t)\eta\end{aligned}\quad (5)$$

while the welfare of beneficiaries of group  $B^j$  at time  $t$  is

$$\begin{aligned}\bar{B}\xi\sigma_t^j \int_0^1 \tilde{u}_{B^j,t}(g(c), c) f_{B^j, n_t^j} dc &= \bar{B}\xi\sigma_t^j \int_0^1 (u_{B^j,t}(g(c), c) + \eta) f_{B^j, n_t^j}(c) dc \\ &= \bar{B}\xi\sigma_t^j \int_0^1 u_{B^j,t}(g(c), c) f_{B^j, n_t^j}(c) dc + \bar{B}\xi\sigma_t^j\eta.\end{aligned}\quad (6)$$

Using Eqs. (4) to (6), we obtain that the welfare in district  $j$  at time  $t$  is

$$\tilde{W}_{A^j,t} + \lambda_{B^j}\tilde{W}_{B^j,t} = W_{A^j,t} + \lambda_{B^j}W_{B^j,t} - (\bar{A} + \lambda_{B^j}\bar{B})K(\bar{\sigma}_t)\eta + \lambda_{B^j}\bar{B}\xi\sigma_t^j\eta(K(\bar{\sigma}_t) + 1)\quad (7)$$

where  $W_{A^j,t}$  and  $W_{B^j,t}$  are the welfare of groups  $A^j$  and  $B^j$  at time  $t$ , *absent* the congestion externality, while the additional terms represent the welfare associated to the transfer of utility from non-beneficiaries to beneficiaries due to labor market congestion. We see from Eq. (7) that choosing  $\sigma_t^j = 1$  results in an additional benefit. Indeed, the additional labor market allocation benefit to the beneficiaries of affirmative action in district  $j$  is positive, whereas there is no additional labor market congestion felt by non-beneficiaries since district  $j$  has measure zero and the decision  $\sigma_t^j = 1$  therefore cannot influence  $\bar{\sigma}_t$  (and  $K(\bar{\sigma}_t)$ ).

A time- $t$  policy maker's objective function (evaluated at some putative policy sequence  $\sigma$ ) can now be written as

$$\sum_{s=t}^{\infty} \delta^{s-t} \left( W_{A^j,s} + \lambda_{B^j}W_{B^j,s} - (\bar{A} + \lambda_{B^j}\bar{B})K(\bar{\sigma}_s)\eta + \lambda_{B^j}\bar{B}\xi\sigma_s^j\eta(K(\bar{\sigma}_s) + 1) \right).$$

We therefore have the following analogue of Proposition 1.

**Observation 2 (Equilibrium policy with congestion)** *Proposition 1 (permanent affirmative action in equilibrium) holds in the presence of labor market congestion.*

We will now show that, in the presence of labor market congestion, the first-best policy also involves temporary affirmative action. For simplicity of exposition, we will suppose here that  $\lambda_{B^j} = \lambda_B$  for all  $j \in J$ . That is, all district policy makers place the same weight on the welfare of group  $B$  relative to that of group  $A$ .

A centralized policy maker's objective function is now

$$\sum_{t=1}^{\infty} \delta^t \int_{j \in J} \left( W_{A^j,t} + \lambda_B W_{B^j,t} - (\bar{A} + \lambda_B \bar{B}) K(\bar{\sigma}_s) \eta + \lambda_B \bar{B} \xi \sigma_t^j \eta (K(\bar{\sigma}_t) + 1) \right) dj,$$

which can be simplified as

$$\sum_{t=1}^{\infty} \delta^t \left( \int_{j \in J} (W_{A^j,t} + \lambda_B W_{B^j,t}) dj - \bar{J} (\bar{A} + \lambda_B \bar{B} (1 - \xi \bar{\sigma}_t)) K(\bar{\sigma}_t) \eta + \lambda_B \bar{J} \bar{B} \xi \bar{\sigma}_t \eta \right). \quad (8)$$

The first term in Eq. (8) is the welfare in the *absence* of congestion. It is easy to show that the second and third terms sum to 0 when  $\lambda_B = 1$ , since the congestion externality amounts to a transfer of utility between beneficiaries and non-beneficiaries. It is also easy to verify that, when  $\lambda_B < 1$ , the sum of the second and third terms is less than 0, as the welfare transfer to beneficiaries is weighted relatively less than what is taken from non-beneficiaries.

It then follows that when  $\lambda_B \in [0, 1]$ , labor market congestion represents an additional cost of implementing affirmative action in each period. The argument of Proposition 2 thus still holds and affirmative action will be stopped after a finite number of periods.

**Observation 3 (First-best policy with congestion)** *Proposition 2 (temporary affirmative action as first-best policy) holds in the presence of labor market congestion.*

### 5.1.5 Micro-foundations: Wage setting with Bertrand competition

We suppose that each firm produces a numeraire good of price equal to 1 with a constant return to scale technology and using labor as the input. The quantity of the numeraire good produced by a unit mass of workers of performance level  $c$  is thus simply  $c$ . The profit generated by a unit mass of workers of performance level  $c$ , when they are paid a wage  $\omega_t^j(\bar{c})$ , is thus

$$\pi = c - \omega_t^j(\bar{c}).$$

Since a firm only observes the curriculum vitae quality  $\bar{c}$  of a district- $j$  worker it hires, the expected profit generated by a unit mass of district- $j$  workers with such curriculum vitae is then

$$\mathbb{E}_t[\pi|\bar{c}, j] = \mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] - \omega_t^j(\bar{c})$$

where, as we know,  $\mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$  is the expected performance level of a district- $j$  worker presenting a curriculum vitae  $\bar{c}$ , given some putative affirmative action policy sequence  $\sigma$ .

If the firm hires a mass  $q$  of district- $j$  workers with curriculum vitae qualities having a density function  $f_t^j(\bar{c})$ , then its expected profit is

$$\begin{aligned} \Pi &= q \mathbb{E}_t[\pi|j] \\ &= q \int_{\bar{c}} \mathbb{E}_t[\pi|\bar{c}, j] f_t^j(\bar{c}) d\bar{c} \\ &= q \int_{\bar{c}} \left( \mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] - \omega_t^j(\bar{c}) \right) f_t^j(\bar{c}) d\bar{c} \end{aligned} \quad (9)$$

where  $\Pi$  is also the realized profit, since each worker has zero measure.

A firm will thus maximize this profit by choosing an optimal wage function  $\omega_t^j$ . Note that the profit in Eq. (9) is additively separable across  $\bar{c}$ . A firm thus chooses, for each curriculum vitae quality  $\bar{c}$ , the wage  $\omega_t^j(\bar{c})$  that maximizes

$$\mathbb{E}_t[\pi|\bar{c}, j] = \mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] - \omega_t^j(\bar{c}).$$

Since we consider a perfectly competitive Bertrand setting, it follows that the optimal wage will be equal to a worker's expected performance level, i.e.  $\omega_t^j(\bar{c}) = \mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$ , which is the worker's marginal productivity. Indeed, giving a wage higher than  $\mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$  would result in a negative profit from hiring workers of that curriculum vitae quality, while giving a wage lower than  $\mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$  would result in another employer hiring the workers away with a slightly higher wage. It also follows that a firm's profit is zero, i.e.  $\Pi = 0$ , on the equilibrium path. Thus, even if the policy makers care about the firms' welfare, the latter will not appear in their objective function (i.e. in Eq. (3)) along the equilibrium path. Including profits in the policy maker's objective function would affect the assessment of deviations, but the qualitative insights presented throughout the paper would be unaffected.

### 5.1.6 Other elaborations

In an Appendix (Section 7.1), we formulate a generalized model where we allow for strategic behavior by workers, by which they can present a curriculum vitae of any chosen quality. We show that the wage chosen by employers is then a non-decreasing function of the curriculum vitae quality. This generalization formally removes the need for Assumption 1.

## 5.2 Comparisons with existing literature

We mainly depart from the existing literature on affirmative action by studying the incentives of decentralized policy makers to implement affirmative action policies. Indeed, most of the literature focuses on other incentives: those linked to hiring decisions made by employers or to investments in human capital made by workers, which may be reduced by an affirmative action policy (e.g. Lundberg and Startz (1983), Coate and Loury (1993a) or Coate and Loury (1993b); see also Fang and Moro (2011) for a survey on discrimination and affirmative action).

The existing literature on affirmative action is vast and often tries to describe or explain inequalities between groups. Early developments include taste-based theories of discrimination (e.g. Becker (1957)), which suppose that exogenous preferences generate wage differences between groups, although the latter are unlikely to persist in competitive markets. Statistical discrimination theories, on the other hand, mainly attempt to explain outcome differences using imperfect information about the workers' performance levels, which leads to different wages being rationally paid to workers of different groups (e.g. Phelps (1972), Arrow (1973), Lundberg and Startz (1983), Coate and Loury (1993a) or Coate and Loury (1993b)). Such models often also link these different wages to the workers' incentives to invest in human capital, thus sustaining a performance gap between groups.

Our argument is based on a novel moral hazard consideration on the policy makers' part and it complements other (more direct) critiques such as those voiced by Sowell (2005).

## 6 Proofs

### Proof of Lemma 1 (Wage function).

Note that  $\omega_t^j(\bar{c})$  is the conditional expectation of a district- $j$  worker's actual performance level at time  $t$  when declaring a curriculum vitae of quality  $\bar{c}$ , given a putative policy sequence  $\sigma = \{\{\sigma_t^j\}_{j \in J}\}_{t=1}^\infty$  assumed by employers and given observed aggregate (average) policy statistics  $\{\bar{\sigma}_s\}_{s=1}^t$  consistent with  $\sigma$  (which is the case on the equilibrium path). Thus,

$$\begin{aligned}\omega_t^j(\bar{c}) &= \mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] \\ &= \mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) \cdot g^{-1}(\bar{c}) + (1 - \mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma)) \cdot \bar{c}\end{aligned}$$

Now to express  $\mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma)$ , we first express  $\mathbb{P}_t(\{aa\}|\bar{c} \in N(\bar{c}, \epsilon), j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma)$ , where  $N(\bar{c}, \epsilon)$  is an  $\epsilon$ -neighborhood of  $\bar{c}$ :

$$\begin{aligned}\mathbb{P}_t(\{aa\}|\bar{c} \in N(\bar{c}, \epsilon), j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) &= \frac{\mathbb{P}_t(\{\bar{c} \in N(\bar{c}, \epsilon)\} \cap \{aa\}|j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma)}{\mathbb{P}_t(\bar{c} \in N(\bar{c}, \epsilon)|j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma)} \\ &= \frac{\mathbb{P}_t(\bar{c} \in N(\bar{c}, \epsilon) \cap B^j|j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) \cdot \xi \sigma_t^j}{\mathbb{P}_t(\bar{c} \in N(\bar{c}, \epsilon)|j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma)} \\ &= \frac{\mathbb{P}_t(\bar{c} \in N(\bar{c}, \epsilon)|B^j, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) \cdot \mathbb{P}(B^j) \cdot \xi \sigma_t^j}{\mathbb{P}_t(\bar{c} \in N(\bar{c}, \epsilon)|j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma)} \\ &= \frac{\int_{\bar{c} \in N(g^{-1}(\bar{c}), \epsilon/g'^{-1}(\bar{c}))} f_{B^j, n_t^j}(g^{-1}(\bar{c})) d\bar{c} \frac{|B^j|}{|A^j|+|B^j|} \xi \sigma_t^j}{\int_{\bar{c} \in N(\bar{c}, \epsilon)} f_{A^j}(\bar{c}) d\bar{c} \frac{|A^j|}{|A^j|+|B^j|} + \int_{\bar{c} \in N(\bar{c}, \epsilon)} f_{B^j, n_t^j}(\bar{c}) d\bar{c} \frac{|B^j|}{|A^j|+|B^j|} (1 - \xi \sigma_t^j) + \int_{\bar{c} \in N(g^{-1}(\bar{c}), \epsilon/g'^{-1}(\bar{c}))} f_{B^j, n_t^j}(g^{-1}(\bar{c})) d\bar{c} \frac{|B^j|}{|A^j|+|B^j|} \xi \sigma_t^j} \\ &= \frac{|B^j| \xi \sigma_t^j \int_{\bar{c} \in N(g^{-1}(\bar{c}), \epsilon/g'^{-1}(\bar{c}))} f_{B^j, n_t^j}(g^{-1}(\bar{c})) d\bar{c}}{|A^j| \int_{\bar{c} \in N(\bar{c}, \epsilon)} f_{A^j}(\bar{c}) d\bar{c} + |B^j| (1 - \xi \sigma_t^j) \int_{\bar{c} \in N(\bar{c}, \epsilon)} f_{B^j, n_t^j}(\bar{c}) d\bar{c} + |B^j| \xi \sigma_t^j \int_{\bar{c} \in N(g^{-1}(\bar{c}), \epsilon/g'^{-1}(\bar{c}))} f_{B^j, n_t^j}(g^{-1}(\bar{c})) d\bar{c}}\end{aligned}$$

Then, we take the limit as  $\epsilon \rightarrow 0$ :

$$\begin{aligned}\mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) &= \lim_{\epsilon \rightarrow 0} \mathbb{P}_t(\{aa\}|\bar{c} \in N(\bar{c}, \epsilon), j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) \\ &= \lim_{\epsilon \rightarrow 0} \frac{|B^j| \xi \sigma_t^j f_{B^j, n_t^j}(g^{-1}(\bar{c})) 2\epsilon/g'^{-1}(\bar{c})}{|A^j| f_{A^j}(\bar{c}) 2\epsilon + |B^j| (1 - \xi \sigma_t^j) f_{B^j, n_t^j}(\bar{c}) 2\epsilon + |B^j| \xi \sigma_t^j f_{B^j, n_t^j}(g^{-1}(\bar{c})) 2\epsilon/g'^{-1}(\bar{c})} \\ &= \frac{|B^j| \xi \sigma_t^j f_{B^j, n_t^j}(g^{-1}(\bar{c})) / g'^{-1}(\bar{c})}{|A^j| f_{A^j}(\bar{c}) + |B^j| (1 - \xi \sigma_t^j) f_{B^j, n_t^j}(\bar{c}) + |B^j| \xi \sigma_t^j f_{B^j, n_t^j}(g^{-1}(\bar{c})) / g'^{-1}(\bar{c})}\end{aligned}$$

■

### Proof of Lemma 2 (Wage versus performance level).

From Lemma 1 we know that

$$\omega_t^j(\bar{c}) = \mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) \cdot g^{-1}(\bar{c}) + (1 - \mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma)) \cdot \bar{c}$$

Part (i): When  $\sigma_t^j = 1$ , then  $\mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) > 0$ . Since  $g^{-1}(\bar{c}) < \bar{c}$ , it follows immediately that  $g^{-1}(\bar{c}) < \omega_t^j(\bar{c}) < \bar{c}$ . Thus, if the worker does not benefit from affirmative action (i.e.  $c = \bar{c}$ ), then  $\omega_t^j(\bar{c}) < c$  and he gets a wage lower than his performance level. On the other hand,

if the worker benefits from affirmative action (i.e.  $c = g^{-1}(\bar{c})$ ), then  $c < \omega_t^j(\bar{c})$  and he gets a wage higher than his performance level.

Part (ii): When  $\sigma_t^j = 0$ , then  $\mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) = 0$ . Thus,  $\omega_t^j(\bar{c}) = \bar{c}$  and  $\bar{c} = c$  since no one benefits from affirmative action. ■

### Proof of Proposition 1 (Equilibrium policy).

We first show that  $\{\{\sigma_s^{j*}\}_{j \in J}\}_{s=1}^\infty = \{\{1\}_{j \in J}\}_{s=1}^\infty$  is an equilibrium.

Given some equilibrium decision profile  $\sigma^* = \{\{\sigma_s^{j*}\}_{j \in J}\}_{s=1}^\infty = \{\{1\}_{j \in J}\}_{s=1}^\infty$ , any deviation  $\sigma_t^{j'}$  at some time  $t$  has no impact on the wage function since this deviation is unobserved by employers. Indeed, employers form a wage  $\omega_t^{j*}(\bar{c}) = \mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma^*]$  after observing the aggregate (average) policy statistics  $\{\bar{\sigma}_s\}_{s=1}^t$ . A district  $j$ 's actual policy  $\sigma_t^{j'}$  cannot be traced back to it from observing this aggregate statistic. Moreover, by virtue of having measure zero, an individual district's policy  $\sigma_t^{j'}$  does not influence the aggregate statistic  $\bar{\sigma}_t$ .

Therefore,  $\sum_{s=t}^\infty W_{A^j, s} \delta^{s-t}$ , the discounted future welfare of group  $A^j$ , is completely unaffected by an unobserved deviation to  $\sigma_t^{j'}$ . Indeed,  $W_{A^j, s} = |A^j| \int_0^1 u_{A^j, t}(c, c) f_{A^j}(c) dc$ , where the density function  $f_{A^j}(c)$  is constant through time and thus not impacted by  $\sigma_t^{j'}$ , while  $u_{A^j, t}(c, c) = \omega_t^{j*}(c) - \gamma_{A^j}(c - \omega_t^{j*}(c))$  and the wage  $\omega_t^{j*}(c)$  is unaffected by an unobserved deviation to  $\sigma_t^{j'}$ .

On the other hand,  $\sum_{s=t}^\infty \lambda_{B^j} W_{B^j, s} \delta^{s-t}$ , the discounted future welfare of group  $B^j$ , is strictly lower following an unobserved deviation from  $\sigma_t^{j*} = 1$  to  $\sigma_t^{j'} = 0$ . Indeed, at time  $t$ ,

$$\begin{aligned} W_{B^j, t | \sigma_t^{j'} = 0} &= |B^j| \int_0^1 u_{B^j, t}(c, c) f_{B^j, n_t^j}(c) dc \\ &= |B^j| \int_0^1 (\omega_t^{j*}(c) - \gamma_{B^j}(c - \omega_t^{j*}(c))) f_{B^j, n_t^j}(c) dc \\ &< |B^j| \int_0^1 \left[ \xi \omega_t^{j*}(g(c)) + (1 - \xi)(\omega_t^{j*}(c) - \gamma_{B^j}(c - \omega_t^{j*}(c))) \right] f_{B^j, n_t^j}(c) dc \\ &= |B^j| \int_0^1 \left[ \xi u_{B^j, t}(g(c), c) + (1 - \xi) u_{B^j, t}(c, c) \right] f_{B^j, n_t^j}(c) dc \\ &= W_{B^j, t | \sigma_t^{j*} = 1} \end{aligned}$$

where we have used Observation 1(i), the fact that the wage function is not affected by an unobserved deviation and the facts that  $\omega_t^{j*}(c) < \omega_t^{j*}(g(c))$ .

Moreover, at times  $s > t$ ,  $f_{B^j, n_s^j | \sigma_t^{j'} = 0}(c) < f_{B^j, n_s^j | \sigma_t^{j*} = 1}(c)$  since a deviation to  $\sigma_t^{j'} = 0$  has the effect of not changing the distribution of performance at time  $t+1$  compared to the previous period  $t$ . Thus, using Observation 1(i), and the fact that the wage is not affected by an unobserved deviation, then for all  $s > t$ ,

$$\begin{aligned} W_{B^j, s | \sigma_t^{j'} = 0} &= |B^j| \int_0^1 \left[ \xi u_{B^j, s}(g(c), c) + (1 - \xi) u_{B^j, s}(c, c) \right] f_{B^j, n_s^j | \sigma_t^{j'} = 0}(c) dc \\ &= |B^j| \int_0^1 \left[ \xi \omega_s^{j*}(g(c)) + (1 - \xi)(\omega_s^{j*}(c) - \gamma_{B^j}(c - \omega_s^{j*}(c))) \right] f_{B^j, n_s^j | \sigma_t^{j'} = 0}(c) dc \\ &< |B^j| \int_0^1 \left[ \xi \omega_s^{j*}(g(c)) + (1 - \xi)(\omega_s^{j*}(c) - \gamma_{B^j}(c - \omega_s^{j*}(c))) \right] f_{B^j, n_s^j | \sigma_t^{j*} = 1}(c) dc \\ &= |B^j| \int_0^1 \left[ \xi u_{B^j, s}(g(c), c) + (1 - \xi) u_{B^j, s}(c, c) \right] f_{B^j, n_s^j | \sigma_t^{j*} = 1}(c) dc \\ &= W_{B^j, s | \sigma_t^{j*} = 1} \end{aligned}$$

where the inequality follows from  $f_{B^j, n_s^j | \sigma_t^{j'} = 0}(c) \prec f_{B^j, n_s^j | \sigma_t^{j'} = 1}(c)$  and the fact that by Assumption 1 (wage function  $\omega_s^{j*}$  is non-decreasing),  $\xi \omega_s^{j*}(g(c)) + (1 - \xi)(\omega_s^{j*}(c) - \gamma_{B^j}(c - \omega_s^{j*}(c)))$  is non-decreasing when  $\xi$  is high enough.

It follows that as long as  $\lambda_{B^j} > 0$ , then  $\sigma_t^{j*} = 1$  for all  $t$  and  $j$  will be an equilibrium.

To show that this is the unique equilibrium, we now have to show that a deviation to  $\sigma_t^{j'} = 1$ , from a putative equilibrium in which  $\sigma_t^{j*} = 0$ , is always desirable for a district- $j$  policy maker at time  $t$ . For that purpose, suppose that  $\sigma_t^{j*} = 0$  for some  $j, t$ . Then, we must show that  $\sum_{s=t}^{\infty} \lambda_{B^j} W_{B^j, s} \delta^{s-t}$  is strictly higher following a deviation from  $\sigma_t^{j*} = 0$  to  $\sigma_t^{j'} = 1$ .

Consider first the effect of this deviation on the welfare at time  $t$  of members of group  $B^j$ . The same argument as before can be used to show that  $W_{B^j, t | \sigma_t^{j'} = 1} > W_{B^j, t | \sigma_t^{j*} = 0}$ .

Consider now the effect of this deviation on the welfare, at any future time  $s > t$ , of members of group  $B^j$ . We know that  $f_{B^j, n_s^j | \sigma_t^{j*} = 0}(c) \prec f_{B^j, n_s^j | \sigma_t^{j'} = 1}(c)$  for all  $s > t$  since a deviation to  $\sigma_t^{j'} = 1$  has the effect of shifting (in a strict first-order stochastic dominance sense) the future performance distributions of group  $B^j$ .

Then for all  $s > t$ ,

$$\begin{aligned} W_{B^j, s | \sigma_t^{j*} = 0} &= |B^j| \int_0^1 [\xi \sigma_s^{j*} \omega_s^{j*}(g(c)) + (1 - \xi \sigma_s^{j*}) (\omega_s^{j*}(c) - \gamma_{B^j}(c - \omega_s^{j*}(c)))] f_{B^j, n_s^j | \sigma_t^{j*} = 0}(c) dc \\ &< |B^j| \int_0^1 [\xi \sigma_s^{j*} \omega_s^{j*}(g(c)) + (1 - \xi \sigma_s^{j*}) (\omega_s^{j*}(c) - \gamma_{B^j}(c - \omega_s^{j*}(c)))] f_{B^j, n_s^j | \sigma_t^{j'} = 1}(c) dc \\ &= W_{B^j, s | \sigma_t^{j'} = 1} \end{aligned}$$

where we made use of  $\xi \sigma_s^{j*} \omega_s^{j*}(g(c)) + (1 - \xi \sigma_s^{j*}) (\omega_s^{j*}(c) - \gamma_{B^j}(c - \omega_s^{j*}(c)))$  being non-decreasing in  $c$  (follows from Assumption 1 and  $\xi$  being high enough) and  $f_{B^j, n_s^j | \sigma_t^{j*} = 0}(c) \prec f_{B^j, n_s^j | \sigma_t^{j'} = 1}(c)$  for all  $s > t$ . ■

### Proof of Proposition 2 (First-best policy).

We start with the following lemma.

**Lemma 3** *Consider some particular district  $j$ . Let  $\sigma^{j'} = \{\sigma_t^{j'}\}_{t=1}^{\infty}$  be a policy plan with  $\sigma_t^{j'} = 0$  and  $\sigma_{\tau+1}^{j'} = 1$  for some  $\tau$ . Let  $\sigma^j = \{\sigma_t^j\}_{t=1}^{\infty}$  be another policy plan with  $\sigma_\tau^j = 1$ ,  $\sigma_{\tau+1}^j = 0$  and  $\sigma_t^j = \sigma_t^{j'}$  for all other  $t$ . Then there exists  $\bar{\delta} \geq 0$  such that for all  $\delta \in (\bar{\delta}, 1)$ ,  $\sigma^j$  yields a strictly higher welfare for district  $j$  than  $\sigma^{j'}$ .*

### Proof of Lemma 3.

First note that for any group  $G^j \in \{A^j, B^j\}$ ,

$$\sum_{t=1}^{\infty} \delta^t W_{G^j, t} = \sum_{t=1}^{\tau-1} \delta^t W_{G^j, t} + \delta^\tau W_{G^j, \tau} + \delta^{\tau+1} W_{G^j, \tau+1} + \sum_{t=\tau+2}^{\infty} \delta^t W_{G^j, t},$$

where only the terms  $\delta^\tau W_{G^j, \tau}$  and  $\delta^{\tau+1} W_{G^j, \tau+1}$  are different under policies  $\sigma^j$  versus  $\sigma^{j'}$ . We thus only need to compare these two terms under the two policies.

Suppose for now that  $\delta = 1$ .

For group  $A^j$ , the sum  $\delta^\tau W_{A^j, \tau} + \delta^{\tau+1} W_{A^j, \tau+1}$  is the same under policies  $\sigma^j$  and  $\sigma^{j'}$ .

For group  $B^j$ , on the other hand,  $\delta^\tau W_{B^j, \tau} + \delta^{\tau+1} W_{B^j, \tau+1}$  is strictly greater under plan  $\sigma^j$  than

under  $\sigma^{j'}$ . To see this, note that under  $\sigma^j$

$$\delta^\tau W_{B^j, \tau} + \delta^{\tau+1} W_{B^j, \tau+1} = \delta^\tau |B^j| \int_0^1 [\xi \omega_\tau^j(g(c)) + (1-\xi)(\omega_\tau^j(c) - \gamma_{B^j}(c - \omega_\tau^j(c)))] f_{B^j, n_\tau^j}(c) dc + \delta^{\tau+1} |B^j| \int_0^1 \omega_{\tau+1}^j(c) f_{B^j, n_{\tau+1}^j}(c) dc$$

while under  $\sigma^{j'}$

$$\delta^\tau W'_{B^j, \tau} + \delta^{\tau+1} W'_{B^j, \tau+1} = \delta^\tau |B^j| \int_0^1 \omega_\tau^{j'}(c) f_{B^j, n_\tau^{j'}}(c) dc + \delta^{\tau+1} |B^j| \int_0^1 [\xi \omega_{\tau+1}^{j'}(g(c)) + (1-\xi)(\omega_{\tau+1}^{j'}(c) - \gamma_{B^j}(c - \omega_{\tau+1}^{j'}(c)))] f_{B^j, n_{\tau+1}^{j'}}(c) dc.$$

The fact that  $\delta^\tau W_{B^j, \tau} + \delta^{\tau+1} W_{B^j, \tau+1} > \delta^\tau W'_{B^j, \tau} + \delta^{\tau+1} W'_{B^j, \tau+1}$ , when  $\delta = 1$ , follows from the facts that  $\omega_{\tau+1}^j(c) = \omega_{\tau+1}^{j'}(c) = c$ , that  $\xi \omega_\tau^j(g(c)) + (1-\xi)(\omega_\tau^j(c) - \gamma_{B^j}(c - \omega_\tau^j(c))) = \xi \omega_{\tau+1}^{j'}(g(c)) + (1-\xi)(\omega_{\tau+1}^{j'}(c) - \gamma_{B^j}(c - \omega_{\tau+1}^{j'}(c)))$ , that  $f_{B^j, n_\tau^j}(c) = f_{B^j, n_\tau^{j'}}(c)$  and that  $f_{B^j, n_{\tau+1}^j}(c) \succ f_{B^j, n_{\tau+1}^{j'}}(c)$ .

By continuity, it then follows that there exists  $\bar{\delta} \in (0, 1)$  such that for all  $\delta \in (\bar{\delta}, 1)$ , the total welfare is also higher under plan  $\sigma^j$  than under  $\sigma^{j'}$ . ■

Therefore, when  $\delta$  is high enough, it follows by iterative application of Lemma 3 that the optimal policy in district  $j$  has a threshold form  $\hat{\sigma}_t^j = 1$  for  $t < \bar{T}^j$  and  $\hat{\sigma}_t^j = 0$  for  $t \geq \bar{T}^j$  for some  $\bar{T}^j \in \mathbb{N} \cup \infty$ .

We will now rule out the case where  $\bar{T}^j$  could be infinite and thus show that  $\bar{T}^j \in \mathbb{N}$ .

Let us thus compare the welfare of some (large)  $\bar{T}^j < \infty$  to that of the case  $\bar{T}^{j'} = \infty$ . In what follows, the quantities with a prime ( ' ) will be the ones associated to  $\bar{T}^{j'} = \infty$ .

We need to show that

$$\sum_{t=1}^{\infty} \delta^t (W_{A^j, t} + \lambda_{B^j} W_{B^j, t}) > \sum_{t=1}^{\infty} \delta^t (W'_{A^j, t} + \lambda_{B^j} W'_{B^j, t}). \quad (10)$$

Equivalently, it will be convenient to multiply the welfare by the constant  $\frac{1}{|A^j| + |B^j|}$  and verify that

$$\frac{1}{|A^j| + |B^j|} \left( \sum_{t=1}^{\infty} \delta^t (W_{A^j, t} + \lambda_{B^j} W_{B^j, t}) - \sum_{t=1}^{\infty} \delta^t (W'_{A^j, t} + \lambda_{B^j} W'_{B^j, t}) \right) > 0$$

$$\begin{aligned} \sum_{t=1}^{\infty} \frac{\delta^t}{|A^j| + |B^j|} ((W_{A^j, t} + \lambda_{B^j} W_{B^j, t}) - (W'_{A^j, t} + \lambda_{B^j} W'_{B^j, t})) &= \sum_{t=1}^{\infty} \delta^t \left( \frac{|A^j|}{|A^j| + |B^j|} \int \omega_t^j(c) f_{A^j}(c) dc \right. \\ &+ \frac{\lambda_{B^j} |B^j|}{|A^j| + |B^j|} \int [\xi \sigma_t^j \omega_t^j(g(c)) + (1 - \xi \sigma_t^j) \omega_t^j(c)] f_{B^j, n_t^j}(c) dc \\ &- \frac{|A^j|}{|A^j| + |B^j|} \int \omega_t^{j'}(c) f_{A^j}(c) dc \\ &- \frac{\lambda_{B^j} |B^j|}{|A^j| + |B^j|} \int [\xi \omega_t^{j'}(g(c)) + (1 - \xi) \omega_t^{j'}(c)] f_{B^j, n_t^{j'}}(c) dc \\ &+ \sum_{t=1}^{\infty} \delta^t \left( \frac{|A^j|}{|A^j| + |B^j|} \gamma_{A^j} \int (\omega_t^j(c) - \omega_t^{j'}(c)) f_{A^j}(c) dc \right. \\ &+ \frac{\lambda_{B^j} |B^j|}{|A^j| + |B^j|} \int (1 - \xi \sigma_t^j) \gamma_{B^j} [\omega_t^j(c) - c] f_{B^j, n_t^j}(c) dc \\ &\left. - \frac{\lambda_{B^j} |B^j|}{|A^j| + |B^j|} \int (1 - \xi) \gamma_{B^j} [\omega_t^{j'}(c) - c] f_{B^j, n_t^{j'}}(c) dc \right) \quad (11) \end{aligned}$$

The case  $\lambda_{B^j} = 1$  is interesting and worth examining first. In that case, note that the first two

terms of the right-hand side of Eq. (11) rewrite as

$$\left( \frac{|A^j|}{|A^j| + |B^j|} \int \omega_t^j(c) f_{A^j}(c) dc + \frac{|B^j|}{|A^j| + |B^j|} \int [\xi \sigma_t^j \omega_t^j(g(c)) + (1 - \xi \sigma_t^j) \omega_t^j(c)] f_{B^j, n_t^j}(c) dc \right) = \mathbb{E}_t[c|j],$$

since the time  $t$  average wage in district  $j$  under policy  $\bar{T}^j < \infty$  is equal to the time  $t$  average performance level in district  $j$  under policy  $\bar{T}^j < \infty$  (here denoted by  $\mathbb{E}_t[c|j]$ ).

Likewise, the third and fourth terms rewrite as

$$-\left( \frac{|A^j|}{|A^j| + |B^j|} \int \omega_t^{j'}(c) f_{A^j}(c) dc + \frac{|B^j|}{|A^j| + |B^j|} \int [\xi \omega_t^{j'}(g(c)) + (1 - \xi) \omega_t^{j'}(c)] f_{B^j, n_t^{j'}}(c) dc \right) = -\mathbb{E}'_t[c|j],$$

since the time  $t$  average wage in district  $j$  under policy  $\bar{T}^{j'} = \infty$  is equal to the time  $t$  average performance level in district  $j$  under policy  $\bar{T}^{j'} = \infty$  (here denoted by  $\mathbb{E}'_t[c|j]$ ).

We then have that the right-hand side of Eq. (11) can be written as

$$\begin{aligned} \sum_{t=1}^{\infty} \delta^t (\mathbb{E}_t[c|j] - \mathbb{E}'_t[c|j]) + & \sum_{t=1}^{\infty} \delta^t \left( \frac{|A^j|}{|A^j| + |B^j|} \gamma_{A^j} \int (\omega_t^j(c) - \omega_t^{j'}(c)) f_{A^j}(c) dc + \right. \\ & \frac{\lambda_{B^j} |B^j|}{|A^j| + |B^j|} \int (1 - \xi \sigma_t^j) \gamma_{B^j} [\omega_t^j(c) - c] f_{B^j, n_t^j}(c) dc \\ & \left. - \frac{\lambda_{B^j} |B^j|}{|A^j| + |B^j|} \int (1 - \xi) \gamma_{B^j} [\omega_t^{j'}(c) - c] f_{B^j, n_t^{j'}}(c) dc \right) \end{aligned}$$

We must now verify if this is greater than 0. We first make the following observations:

- The first summation term is negative and converges to 0 as  $\bar{T}^j \rightarrow \infty$ . Indeed,  $\mathbb{E}_t[c|j] < \mathbb{E}'_t[c|j]$  for  $t \geq \bar{T}^j$ , since the time  $t$  average performance level keeps increasing as affirmative action gets implemented for more periods. This term converges to 0 as  $\bar{T}^j \rightarrow \infty$  since  $\mathbb{E}_t[c|j] = \mathbb{E}'_t[c|j]$  for  $t < \bar{T}^j$  and  $\sup_{t \geq \bar{T}^j} |\mathbb{E}_t[c|j] - \mathbb{E}'_t[c|j]| \xrightarrow{\bar{T}^j \rightarrow \infty} 0$ , reflecting the fact that the improvements in the performance distribution of group  $B^j$  become marginal after a while.
- The second summation term is positive and bounded away from 0 as  $\bar{T}^j \rightarrow \infty$ . This captures the gain to the non-beneficiaries (of both groups  $A^j$  and  $B^j$ ) of stopping affirmative action after a finite number of periods. Indeed, under a policy of permanent affirmative action  $\bar{T}^{j'} = \infty$ ,

$$\begin{aligned} \omega_t^{j'}(c) &= \mathbb{P}_t(\{aa\}|c, j, \{\bar{\sigma}_s\}_{s=1}^t, \hat{\sigma}^{j'}) g^{-1}(c) + (1 - \mathbb{P}_t(\{aa\}|c, j, \{\bar{\sigma}_s\}_{s=1}^t, \hat{\sigma}^{j'})) c \\ &< c \end{aligned}$$

since  $\mathbb{P}_t(\{aa\}|c, j, \{\bar{\sigma}_s\}_{s=1}^t, \hat{\sigma}^{j'}) > 0$  for all  $t$ . Thus, for  $t \geq \bar{T}^j$ ,  $\omega_t^j(c) - \omega_t^{j'}(c) = c - \omega_t^{j'}(c) > \Delta$  for some  $\Delta > 0$ , while  $\omega_t^j(c) - c = c - c = 0$ .

From the above observations, we can formally state that  $\forall \epsilon > 0$ , there exists  $\bar{T}^j < \infty$  large enough and  $\bar{\delta}(\bar{T}^j) \in (0, 1)$  such that  $\forall \delta \in (\bar{\delta}(\bar{T}^j), 1)$

$$\sum_{t=1}^{\infty} \delta^t |\mathbb{E}_t[c|j] - \mathbb{E}'_t[c|j]| < \epsilon,$$



and

$$\begin{aligned} & \sum_{t=1}^{\infty} \delta^t \left( \frac{|A^j|}{|A^j| + |B^j|} \gamma_{A^j} \int (\omega_t^j(c) - \omega_t^{j'}(c)) f_{A^j}(c) dc + \right. \\ & \quad \left. \frac{\lambda_{B^j} |B^j|}{|A^j| + |B^j|} \int (1 - \xi \sigma_t^j) \gamma_{B^j} [\omega_t^j(c) - c] f_{B^j, n_t^j}(c) dc \right. \\ & \quad \left. - \frac{\lambda_{B^j} |B^j|}{|A^j| + |B^j|} \int (1 - \xi) \gamma_{B^j} [\omega_t^{j'}(c) - c] f_{B^j, n_t^{j'}}(c) dc \right) > 2\epsilon \end{aligned}$$

from which it follows that the right-hand side of Eq. (11) is positive and thus that Eq. (10) is verified.

To complete the proof, we now turn to the case when  $\lambda_{B^j} < 1$ .

First note that when  $\delta$  is high enough, unsurprisingly, group  $A^j$  gains from stopping affirmative action whereas at least a fraction of group  $B^j$  loses. Thus, rearranging the left-hand side of Eq. (11) as follows

$$\sum_{t=1}^{\infty} \frac{\delta^t}{(|A^j| + |B^j|)} ((W_{A^j, t} - W'_{A^j, t}) + \lambda_{B^j} (W_{B^j, t} - W'_{B^j, t})),$$

we notice that decreasing the weight  $\lambda_{B^j}$  placed on the welfare of group  $B^j$  to values strictly smaller than 1 keeps this quantity positive. We can thus conclude that it will still be worth stopping affirmative action after  $\bar{T}^j < \infty$  periods as opposed to continuing it forever. The first-best optimal policy  $\bar{T}_{\lambda_{B^j}}^j$  for some  $\lambda_{B^j} < 1$  will thus be such that  $\bar{T}_{\lambda_{B^j}}^j \leq \bar{T}_{\lambda_{B^j}=1}^j < \infty$ .

Since optimal policies are separable across districts, it follows that the above is true for any  $j \in J$ . This completes the proof. ■

## References

- ARROW, K. (1973): “The theory of discrimination. In: Discrimination in labor markets,” *Orley Achenfelter and Albert Rees (eds.), Princeton–New Jersey*.
- BECKER, G. S. (1957): *The economics of discrimination*, University of Chicago press.
- CHUNG, K.-S. (2000): “Role models and arguments for affirmative action,” *American Economic Review*, 90, 640–648.
- COATE, S. AND G. LOURY (1993a): “Antidiscrimination enforcement and the problem of patronization,” *American Economic Review*, 83, 92–98.
- COATE, S. AND G. C. LOURY (1993b): “Will affirmative-action policies eliminate negative stereotypes?” *American Economic Review*, 1220–1240.
- FANG, H. AND A. MORO (2011): “Theories of statistical discrimination and affirmative action: A survey,” *Handbook of social economics*, 1, 133–200.
- HEILMAN, M. E., C. J. BLOCK, AND J. A. LUCAS (1992): “Presumed incompetent? Stigmatization and affirmative action efforts.” *Journal of Applied Psychology*, 77, 536.

- HEILMAN, M. E., C. J. BLOCK, AND P. STATHATOS (1997): “The affirmative action stigma of incompetence: Effects of performance information ambiguity,” *Academy of Management Journal*, 40, 603–625.
- KAHNEMAN, D. AND A. TVERSKY (1979): “Prospect Theory: An Analysis of Decision under Risk,” *Econometrica*, 47, 263–292.
- LESLIE, L. M., D. M. MAYER, AND D. A. KRAVITZ (2014): “The stigma of affirmative action: A stereotyping-based theory and meta-analytic test of the consequences for performance,” *Academy of Management Journal*, 57, 964–989.
- LUNDBERG, S. J. AND R. STARTZ (1983): “Private discrimination and social intervention in competitive labor market,” *American Economic Review*, 73, 340–347.
- PHELPS, E. S. (1972): “The statistical theory of racism and sexism,” *American Economic Review*, 62, 659–661.
- SOWELL, T. (2005): *Affirmative Action Around the World: An Empirical Study*, Yale University Press.

## 7 Supplementary appendix

### 7.1 A more general model allowing for strategic behavior by workers

We present here a more general model, of which the model presented in the main part of the paper is a particular case. We show that in equilibrium, this more general model endogenously generates a wage function that is non-decreasing in the curriculum vitae quality, thus formally removing the need for Assumption 1.

Here, we allow agents to choose the curriculum vitae quality that they present to employers. This allows us to treat the more general case where the conditional expectation  $\mathbb{E}_t[c|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$  may not be monotone. We illustrate that the results presented in the main part of the paper still hold, since they are just a particular case of this more general setting (i.e. the case when agents truthfully declare their curriculum vitae quality).

In this general model, a wage function  $\omega_t^j(\hat{c})$  set by employers is the wage the worker earns when declaring a curriculum vitae of quality  $\hat{c} \in [0, 1]$  to the employer. Here, we see that a worker can declare a curriculum vitae of quality not necessarily equal to his actual quality  $\bar{c}$ . This is formalized in the following definition.

**Definition 2** *A wage function  $\omega_t^j : [0, 1] \rightarrow [0, 1]$  determines the wage a worker earns when declaring a curriculum vitae of quality  $\hat{c}$  to the employer.*

The utility of a type  $(c, \bar{c}, G^j)$  worker, when presenting a curriculum vitae of quality  $\hat{c} \in [0, 1]$ , is thus

$$u_{G^j,t}(\hat{c}, c) = \omega_t^j(\hat{c}) - \gamma_G \max\{c - \omega_t^j(\hat{c}), 0\} - \kappa \max\{\hat{c} - \bar{c}, 0\} \quad (12)$$

where  $\kappa \max\{\hat{c}-\bar{c}, 0\}$ , with  $\kappa > 0$ , is a penalty suffered for cheating (i.e. presenting a curriculum vitae quality higher than the actual one  $\bar{c}$ ). Note that no penalty is suffered for presenting a curriculum vitae of lower quality than  $\bar{c}$ .

A worker thus chooses to present a curriculum vitae of quality  $\hat{c}$  such that

$$\hat{c} \in \operatorname{argmax}_{\tilde{c} \in [0,1]} u_{G^j,t}(\tilde{c}, c)$$

**Definition 3** Given a wage function  $\omega_t^j : [0, 1] \rightarrow [0, 1]$ , a curriculum vitae declaration function  $\mu_t^j : [0, 1] \rightarrow [0, 1]$  assigns a declared curriculum vitae quality  $\hat{c}$  to an actual curriculum vitae quality  $\bar{c}$ , that is  $\hat{c} = \mu_t^j(\bar{c})$ .

**Definition 4** Given a putative policy sequence  $\sigma$ , a labor market equilibrium  $(\omega_t^{j*}, \mu_t^{j*})$  is a continuous wage function and a curriculum vitae declaration function such that

$$\omega_t^{j*}(\hat{c}) = \mathbb{E}_t[c|\hat{c}, j, \mu_t^{j*}, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$$

and

$$\mu_t^{j*}(\bar{c}) \in \operatorname{argmax}_{\tilde{c} \in [0, \bar{c}]} u_{G^j,t}(\tilde{c}, c).$$

Recall from Eq.(12) that the utility  $u_{G^j,t}(\hat{c}, c)$  depends on the wage  $\omega_t^{j*}(\hat{c})$ .

If  $\kappa$  is high enough, a continuous wage function  $\omega_t^j(\hat{c})$  will prevent cheating since the marginal penalty of presenting a curriculum vitae quality greater than  $\bar{c}$  will exceed the marginal benefit in terms of increased wage. A sufficient condition for this to hold is that  $\kappa > \frac{\omega_t^j(\hat{c}) - \omega_t^j(\bar{c})}{\hat{c} - \bar{c}}$  for any  $\hat{c} > \bar{c}$ .

We thus have the following lemma.

**Lemma 4** Suppose  $\kappa$  is high enough. Given a putative policy sequence  $\sigma$ , there exist intervals  $\{(c_l^L, c_l^H)\}_{l=1}^{\bar{l}}$  with  $\bar{l} \geq 0$ , so that the (weakly) increasing wage function

$$\omega_t^{j*}(\hat{c}) = \begin{cases} \mathbb{E}_t[c|\bar{c} = \hat{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] & \text{if } \hat{c} \notin \bigcup_l (c_l^L, c_l^H) \\ \mathbb{E}_t[c|\bar{c} \in (c_l^L, c_l^H), j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] & \text{if } \hat{c} \in (c_l^L, c_l^H) \end{cases} \quad (13)$$

and the curriculum vitae declaration strategy

$$\mu_t^{j*}(\bar{c}) = \begin{cases} \bar{c} & \text{if } \bar{c} \notin \bigcup_l (c_l^L, c_l^H) \\ c_l^L & \text{if } \bar{c} \in (c_l^L, c_l^H) \end{cases} \quad (14)$$

constitute a labor market equilibrium.

In the above,

$$\mathbb{E}_t[c|\bar{c} = \hat{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] = \mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) \cdot g^{-1}(\bar{c}) + (1 - \mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma)) \cdot \bar{c},$$

with

$$\mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) = \frac{|B^j| \xi \sigma_t^j f_{B^j, n_t^j}(g^{-1}(\bar{c})) / g'^{-1}(\bar{c})}{|A^j| f_{A^j}(\bar{c}) + |B^j| (1 - \xi \sigma_t^j) f_{B^j, n_t^j}(\bar{c}) + |B^j| \xi \sigma_t^j f_{B^j, n_t^j}(g^{-1}(\bar{c})) / g'^{-1}(\bar{c})},$$

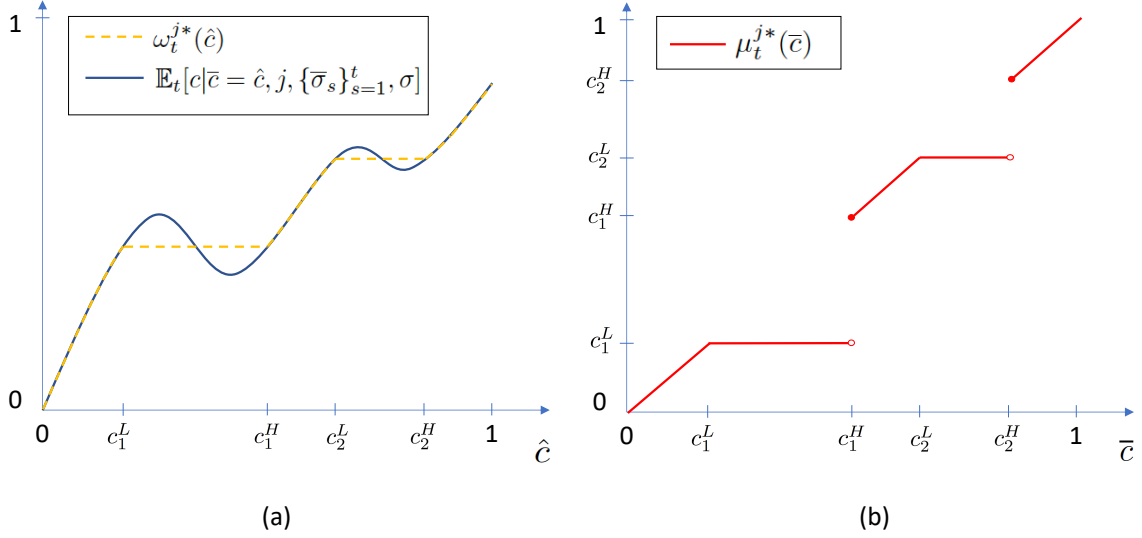


Figure 4: Equilibrium wage function  $\omega_t^{j*}$  (panel (a)) and curriculum vitae declaration function  $\mu_t^{j*}$  (panel (b)).

$\{aa\}$  being the event that a worker benefited from affirmative action, while

$$\mathbb{E}_t[c|\bar{c} \in (c_1^L, c_1^H), j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] = \int_{\bar{c}=c_1^L}^{c_1^H} \mathbb{E}_t[c|\bar{c} = \hat{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] f_t^j(\bar{c}) d\bar{c}$$

where

$$f_t^j(\bar{c}) = \frac{1}{|A^j| + |B^j|} \left( |A^j| f_{A^j}(\bar{c}) + |B^j| \xi \sigma_t^j f_{B^j, n_t^j}(g^{-1}(\bar{c})) / g'^{-1}(\bar{c}) + |B^j| (1 - \xi \sigma_t^j) f_{B^j, n_t^j}(\bar{c}) \right)$$

is the overall population density for the curriculum vitae quality at time  $t$  in district  $j$ .

The equilibrium wage function stated in Lemma 4 has the form described in Figure 4(a). We see that it is weakly increasing, but strictly increasing in certain sections. In the particular case when  $\bar{l} = 0$ , then it can be strictly increasing over the whole domain, as in the case presented earlier in the main part of the paper. The equilibrium curriculum vitae declaration function has the form described in Figure 4(b). It is such that a worker truthfully declares his curriculum vitae quality, i.e.  $\hat{c} = \bar{c}$ , when  $\bar{c}$  is in an interval where the wage function is strictly increasing, since declaring anything lower would yield a lower salary. On the other hand, when  $\bar{c}$  is in an interval where the wage function is flat, the worker declares the lowest curriculum vitae quality  $\hat{c}$  on that flat interval, i.e.  $\hat{c} = c_1^L$ . Indeed, declaring such a curriculum vitae quality  $\hat{c} \leq \bar{c}$  provides the worker with the same salary as he would get when declaring the actual one:  $\omega_t^{j*}(\hat{c}) = \omega_t^{j*}(\bar{c})$ . In the particular case where  $\bar{l} = 0$  and the wage function is strictly increasing, then all workers would always declare their true curriculum vitae quality ( $\mu_t^{j*}(\bar{c}) = \bar{c}$ , as in the case presented earlier in the main part of the paper). Note that, as required by the equilibrium definition, an employer correctly sets the wage equal to the conditional expectation of a worker's performance (i.e.  $\omega_t^{j*}(\hat{c}) = \mathbb{E}_t[c|\hat{c}, j, \mu_t^{j*}, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$ ).

**Lemma 5** *The equilibrium wage function  $\omega_t^{j*}(\hat{c})$  is weakly increasing, but strictly increasing at least*

on some regions of the support<sup>10</sup>  $[0, 1]$ .

**Lemma 6** *Let  $h(c)$  be any weakly increasing function that is strictly increasing at least on some opened subinterval of its support  $[0, 1]$  and is differentiable almost everywhere. If  $f \succ \tilde{f}$ , where  $f$  and  $\tilde{f}$  are probability density functions on  $[0, 1]$  and  $\succ$  indicates strict first-order stochastic dominance, then  $\int_0^1 h(c)f(c)dc > \int_0^1 h(c)\tilde{f}(c)dc$ .*

The next lemma is simply a more general version of Lemma 2(i) of the main part of the paper, adapted to the labor market equilibrium concept defined in Definition 4.

**Lemma 7** *If a worker benefits from affirmative action (i.e.  $c = g^{-1}(\bar{c})$ ), then he gets a wage higher than his performance level (i.e.  $c < \omega_t^{j*}(\mu_t^{j*}(\bar{c}))$ ). If a worker does not benefit from affirmative action (i.e.  $c = \bar{c}$ ), then he gets a wage lower than his performance level (i.e.  $c > \omega_t^{j*}(\mu_t^{j*}(\bar{c}))$ ).*

Using Lemma 7, we can make the same observations as in the main part of the paper, namely that non-beneficiaries of affirmative action (of either group  $A$  or  $B$ ) suffer a feeling of injustice, while beneficiaries do not.

Lemmas 5, 6 and 7 are all the ingredients needed to confirm that Proposition 1 (permanent affirmative action in equilibrium) and Proposition 2 (temporary affirmative action in the first-best case) of the main part of the paper hold in this more general model. The proofs are otherwise identical.

## 7.2 Proofs of results in Section 7.1

**Proof of Lemma 4.** Throughout this proof, we suppose  $\kappa$  is high enough to prevent cheating. A sufficient condition for this to hold is that  $\kappa > \frac{\omega_t^j(\hat{c}) - \omega_t^j(\bar{c})}{\hat{c} - \bar{c}}$  for any  $\hat{c} > \bar{c}$ . In such a case, the marginal penalty of presenting a curriculum vitae quality greater than  $\bar{c}$  will exceed the marginal benefit in terms of increased wage.

*Step I: Compute the wage  $\tilde{\omega}_t^j$  assuming truthful declaration of  $\bar{c}$ .*

Suppose first that workers truthfully declare their curriculum vitae quality, i.e.  $\hat{c} = \mu_t^j(\bar{c}) = \bar{c}$ . Under such a declaration function  $\mu$ , call  $\tilde{\omega}_t^j(\hat{c}) = \mathbb{E}_t[c|\hat{c}, j, \mu_t^j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$  the conditional expectation of the actual performance level when declaring a curriculum vitae of quality  $\hat{c}$ . Then,

$$\begin{aligned} \tilde{\omega}_t^j(\hat{c}) &= \mathbb{E}_t[c|\hat{c}, j, \mu_t^j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] \\ &= \mathbb{E}_t[c|\hat{c} = \bar{c}, j, \mu_t^j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] \\ &= \mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) \cdot g^{-1}(\bar{c}) + (1 - \mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma)) \cdot \bar{c} \end{aligned}$$

and the same argument as in the proof of Lemma 1 of the main part of the paper allows us to state that

$$\mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) = \frac{|B^j|\xi\sigma_t^j f_{B^j, n_t^j}(g^{-1}(\bar{c}))/g'^{-1}(\bar{c})}{|A^j|f_{A^j}(\bar{c}) + |B^j|(1 - \xi\sigma_t^j)f_{B^j, n_t^j}(\bar{c}) + |B^j|\xi\sigma_t^j f_{B^j, n_t^j}(g^{-1}(\bar{c}))/g'^{-1}(\bar{c})},$$

<sup>10</sup>This is actually stronger than needed. For Propositions 1 and 2 to hold in this more general model,  $\omega_t^{j*}(\hat{c})$  only needs to have these properties for the  $\hat{c}$ 's being played in equilibrium (i.e.  $\hat{c} = \mu_t^{j*}(\bar{c})$ ).

$\{aa\}$  being the event that a worker benefited from affirmative action.

*Step II: Such a wage function  $\tilde{\omega}_t^j$  cannot in general be part of an equilibrium.*

Suppose that  $\tilde{\omega}_t^j$  is increasing for  $c \in [0, c_1]$  and decreasing over some interval  $[c_1, c_1']$ . If the wage function is  $\tilde{\omega}_t^j$ , then a worker with an actual curriculum vitae quality  $\bar{c} \in (c_1, c_1']$  will choose to declare a curriculum vitae quality  $\hat{c} < \bar{c}$  since he can obtain a higher wage  $\tilde{\omega}_t^j(\hat{c}) > \tilde{\omega}_t^j(\bar{c})$  by doing so. It follows that  $\mu_t^j(\bar{c}) = \bar{c}$  cannot be part of an equilibrium since  $\mu_t^j(\bar{c}) \notin \operatorname{argmax}_{\tilde{c} \in [0, \bar{c}]} u_{G^j, t}(\tilde{c}, c)$  for such  $\bar{c}$ .

Since  $\mu_t^j(\bar{c}) = \bar{c}$  is not part of an equilibrium, it follows that  $\tilde{\omega}_t^j(\hat{c}) = \mathbb{E}_t[c | \bar{c} = \hat{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$  is not equal to the correct conditional expectation  $\mathbb{E}_t[c | \hat{c}, j, \mu_t^{j*}, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$  where  $\mu_t^{j*}$  is an equilibrium declaration function. Thus,  $\tilde{\omega}_t^j(\hat{c})$  cannot in general be the equilibrium wage function.

*Step III: Building a weakly increasing wage function  $\omega_t^{j*}(\hat{c})$  using  $\tilde{\omega}_t^j(\hat{c})$ .*

On the other hand, there exist  $c_1^L < c_1$  and  $c_1^H \geq c_1'$  such that a wage

$$\omega_t^{j*}(\hat{c}) = \begin{cases} \tilde{\omega}_t^j(\hat{c}), & \text{if } \hat{c} \in [0, c_1^L] \\ \tilde{\omega}_t^j(c_1^L) & \text{when } \hat{c} \in (c_1^L, c_1^H] \end{cases} \quad (15)$$

corresponds to  $\mathbb{E}_t[c | \hat{c}, j, \mu_t^{j*}, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$ , where  $\mu_t^{j*}$  is as in the statement of the lemma. Such a pair  $\{c_1^L, c_1^H\}$  satisfies

$$\tilde{\omega}_t^j(c_1^L) = \int_{\bar{c}=c_1^L}^{c_1^H} \tilde{\omega}_t^j(\bar{c}) f_t^j(\bar{c}) d\bar{c} \quad (16)$$

$$\tilde{\omega}_t^j(c_1^H) = \int_{\bar{c}=c_1^L}^{c_1^H} \tilde{\omega}_t^j(\bar{c}) f_t^j(\bar{c}) d\bar{c} \quad (17)$$

and

$$\int_{\bar{c}=c_1^H}^1 \tilde{\omega}_t^j(\bar{c}) f_t^j(\bar{c}) d\bar{c} > \tilde{\omega}_t^j(c_1^H). \quad (18)$$

where

$$f_t^j(\bar{c}) = \frac{1}{|A^j| + |B^j|} \left( |A^j| f_{A^j}(\bar{c}) + |B^j| \xi \sigma_t^j f_{B^j, n_t^j}(g^{-1}(\bar{c})) / g'^{-1}(\bar{c}) + |B^j| (1 - \xi \sigma_t^j) f_{B^j, n_t^j}(\bar{c}) \right)$$

is simply the overall population density for the curriculum vitae quality  $\bar{c}$  at time  $t$  in district  $j$ .

By construction,  $\omega_t^{j*}(\hat{c})$  is strictly increasing for  $\hat{c} \in [0, c_1^L]$  and flat for  $\hat{c} \in (c_1^L, c_1^H]$ . This is pictured in Figure 4(a). We will generalize this in Step V below.

*Step IV: Verifying that  $(\omega_t^{j*}, \mu_t^{j*})$  is a labor market equilibrium for  $\bar{c} \in [0, c_1^H]$ .*

For any worker with an actual curriculum vitae quality  $\bar{c} \in [0, c_1^L]$ , the best response to such a wage function is  $\mu_t^{j*}(\bar{c}) = \bar{c} = \operatorname{argmax}_{\tilde{c} \in [0, \bar{c}]} u_{G^j, t}(\tilde{c}, c)$  since  $\omega_t^{j*}(\hat{c})$  is strictly increasing over that range and thus the worker chooses to declare  $\hat{c} = \bar{c}$  to maximize his wage. Therefore,  $\omega_t^{j*}(\hat{c}) = \mathbb{E}_t[c | \hat{c}, j, \mu_t^{j*}, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] = \mathbb{E}[c | \hat{c} = \bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] = \tilde{\omega}_t^j(\hat{c})$  for  $\bar{c} \in [0, c_1^L]$ . It follows that  $\omega_t^{j*}$  and  $\mu_t^{j*}$  satisfy the labor market equilibrium condition for  $\bar{c} \in [0, c_1^L]$ .

Moreover, for any worker with an actual curriculum vitae quality  $\bar{c} \in (c_1^L, c_1^H]$ , the best response set to a such a wage function is  $[c_1^L, \bar{c}] = \operatorname{argmax}_{\tilde{c} \in [0, \bar{c}]} u_{G^j, t}(\tilde{c}, c)$ . A worker is indeed indifferent about declaring any  $\hat{c} \in [c_1^L, \bar{c}]$ , since it yields a salary  $\omega_t^{j*}(\hat{c}) = \tilde{\omega}_t^j(c_1^L)$ , which is the maxi-

mum the worker can obtain. It follows that  $\mu_t^{j*}(\bar{c}) = c_1^L \in \operatorname{argmax}_{\bar{c} \in [0, \bar{c}]} u_{G^j, t}(\bar{c}, c)$ . Since,  $\omega_t^{j*}(c_1^L) = \mathbb{E}_t[c|\hat{c}, j, \mu_t^{j*}, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] = \mathbb{E}_t[c|\hat{c} = c_1^L, j, \mu_t^{j*}, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] = \mathbb{E}_t[c|\bar{c} \in [c_1^L, c_1^H], j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] = \tilde{\omega}_t^j(c_1^L)$ , it follows that  $\omega_t^{j*}$  and  $\mu_t^{j*}$  satisfies the labor market equilibrium condition for  $\bar{c} \in (c_1^L, c_1^H]$ .

*Step V: Generalizing to  $\bar{c} \in [0, 1]$ .*

If  $c_1^H < 1$  and  $\tilde{\omega}_t^j(\bar{c})$  is decreasing over some range(s) in  $[c_1^H, 1]$ , then an iterative application of conditions (16), (17) and (18) allows to find other pairs  $\{c_l^L, c_l^H\}$  such that

$$\omega_t^{j*}(\hat{c}) = \begin{cases} \mathbb{E}_t[c|\bar{c} = \hat{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] & \text{if } \hat{c} \notin \bigcup_l (c_l^L, c_l^H) \\ \mathbb{E}_t[c|\bar{c} \in (c_l^L, c_l^H), j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma] & \text{if } \hat{c} \in (c_l^L, c_l^H) \end{cases}$$

and the analysis of Steps II, III and IV generalizes to the rest of the support. ■

**Proof of Lemma 5.** This is a corollary of Lemma 4.

Lemma 4 states that  $\omega_t^{j*}(\hat{c}) = \mathbb{E}_t[c|\bar{c} \in (c_l^L, c_l^H), j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$  for any  $\hat{c} \in (c_l^L, c_l^H)$ , implying that  $\omega_t^{j*}(\hat{c})$  is flat for such  $\hat{c}$  (since  $\mathbb{E}_t[c|\bar{c} \in (c_l^L, c_l^H), j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$  is a constant).

On the other hand, Lemma 4 states that  $\omega_t^{j*}(\hat{c}) = \mathbb{E}_t[c|\bar{c} = \hat{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma]$  when  $\hat{c} \notin \bigcup_l (c_l^L, c_l^H)$  and Steps III and V of the proof of Lemma 4 show that  $\omega_t^{j*}(\hat{c})$  is constructed so as to be strictly increasing over such intervals. ■

**Proof of Lemma 6.** The inequality rewrites

$$\int_0^1 h(c) [f(c) - \tilde{f}(c)] dc > 0.$$

After integrating by parts, this can be written as

$$\left[ h(c) [F(c) - \tilde{F}(c)] \right] \Big|_0^1 - \int_0^1 h'(c) [F(c) - \tilde{F}(c)] dc$$

where  $F$  and  $\tilde{F}$  are the CDFs associated with the PDFs  $f$  and  $\tilde{f}$ . The first term is equal to 0 since  $F(0) = \tilde{F}(0) = 0$  and  $F(1) = \tilde{F}(1) = 1$ . Moreover, since  $h'(c) \geq 0$  almost everywhere with  $h'(c) > 0$  on non-trivial parts of the support, the last term is strictly greater than 0 if  $F(c) < \tilde{F}(c)$  for all  $c \in (0, 1)$ , i.e. if  $f \succ \tilde{f}$ . ■

**Proof of Lemma 7.** When  $\bar{c} \notin \bigcup_l (c_l^L, c_l^H)$ , then from Lemma 4 we know that a worker truthfully declares a curriculum vitae quality  $\hat{c} = \bar{c}$  and gets a wage

$$\omega_t^{j*}(\bar{c}) = \mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma) \cdot g^{-1}(\bar{c}) + (1 - \mathbb{P}_t(\{aa\}|\bar{c}, j, \{\bar{\sigma}_s\}_{s=1}^t, \sigma)) \cdot \bar{c}$$

Since  $g^{-1}(\bar{c}) < \bar{c}$ , it follows immediately that  $g^{-1}(\bar{c}) < \omega_t^{j*}(\bar{c}) < \bar{c}$ .

Thus, if the worker does not benefit from affirmative action (i.e.  $c = \bar{c}$ ), then  $\omega_t^{j*}(\bar{c}) < c$  and he gets a wage lower than his performance level. On the other hand, if the worker benefits from affirmative action (i.e.  $c = g^{-1}(\bar{c})$ ), then  $c < \omega_t^{j*}(\bar{c})$  and he gets a wage higher than his performance level.

We now show that this is also true when  $\bar{c} \in \bigcup_l (c_l^L, c_l^H)$ .

Recall from Lemma 4 that the wage function is flat over  $[c_l^L, c_l^H]$  and equal to  $\omega_t^{j*}(c_l^L)$ . Thus, a worker of performance level  $c_l^L$  who does not benefit from affirmative action gets a wage  $\omega_t^{j*}(c_l^L)$  with  $\omega_t^{j*}(c_l^L) < c'$  and a worker of performance level  $c_l^H$  who does not benefit from affirmative action also gets a wage  $\omega_t^{j*}(c_l^L)$  and  $\omega_t^{j*}(c_l^L) < c''$ . Consider now a worker who does not benefit from

affirmative action and  $\bar{c} \in (c_l^L, c_l^H)$ . Then,  $c = \bar{c}$  with  $c' < c < c''$  and the worker gets a wage  $\omega^*(c_l^L)$ . It follows that  $\omega_t^{j*}(c_l^L) < c$  and he gets a wage lower than his performance level. This applies to any  $\bar{c} \in \bigcup_l(c_l^L, c_l^H)$ .

Now again, recall from Lemma 4 that the wage function is flat over  $[c_l^L, c_l^H]$  and equal to  $\omega_t^{j*}(c_l^L)$ . Thus, a worker of performance level  $g^{-1}(c_l^L)$  who benefits from affirmative action gets a wage  $\omega_t^{j*}(c_l^L)$  with  $g^{-1}(c_l^L) < \omega_t^{j*}(c_l^L)$  and a worker of performance level  $g^{-1}(c_l^H)$  who benefits from affirmative action also gets a wage  $\omega_t^{j*}(c_l^L)$  and  $g^{-1}(c_l^H) < \omega_t^{j*}(c_l^L)$ . Consider now a worker who benefits from affirmative action and  $\bar{c} \in (c_l^L, c_l^H)$ . Then,  $c = g^{-1}(\bar{c})$  with  $g^{-1}(c_l^L) < g^{-1}(\bar{c}) < g^{-1}(c_l^H)$  and the worker gets a wage  $\omega_t^{j*}(c_l^L)$ . It follows that  $c = g^{-1}(\bar{c}) < \omega_t^{j*}(c_l^L)$  and he gets a wage higher than his performance level. This applies to any  $\bar{c} \in \bigcup_l(c_l^L, c_l^H)$ . ■