

Manipulability in Matching Markets: Conflict and Coincidence of Interests*

Itai Ashlagi[†] and Flip Klijn[‡]

June 14, 2010

Abstract

We study comparative statics of manipulations by women in the men-proposing deferred acceptance mechanism in the two-sided one-to-one marriage market. We prove that if a group of women employs truncation strategies or weakly successfully manipulates, then all other women weakly benefit and all men are weakly harmed. We show that our results do not appropriately generalize to the many-to-one college admissions model.

1 Introduction

We study the effect of strategic agents on non-strategic agents in two-sided matching markets. Consider the marriage market introduced by Gale and Shapley (1962) where the two (finite) sides of the market are “men” and “women,” each agent having preferences over the other side of the market and the prospect of being alone. An outcome for a marriage market is a matching in which each agent either marries an agent from the other side of the market or remains single. A key property for a matching is stability. A matching is stable if each agent has an acceptable match and there is no pair of a man and a woman who like each other

*We thank Bettina Klaus and Alexandru Nichifor for detailed comments on a preliminary draft.

[†]Harvard Business School

[‡]*Corresponding author.* Institute for Economic Analysis (CSIC), Campus UAB, 08193 Bellaterra (Barcelona), Spain; During academic year 2009–2010: Harvard Business School, Baker Library | Bloomberg Center 437, Soldiers Field, Boston, MA 02163, USA; e-mail: flip.klijn@hbs.edu. He gratefully acknowledges a research fellowship from Harvard Business School and support from Plan Nacional I+D+i (ECO2008–04784), Generalitat de Catalunya (SGR2009–01142), the Barcelona GSE Research Network, and the Consolider-Ingenio 2010 (CSD2006–00016) program.

better than their current matches. Using their deferred acceptance algorithm, Gale and Shapley (1962) constructively proved that there exists a stable matching for each profile of preferences. Moreover, Knuth (1976) showed that the set of stable matchings is a distributive lattice with respect to the preferences of the agents. An important consequence is that on the set of stable matchings each side of the market has common interests that are in conflict with those of the other side.¹

In this note, we show that the conflict and coincidence of interests extends to the effects of manipulations in the direct-revelation games based on the deferred acceptance algorithm.² Consider the direct-revelation mechanism induced by the men-proposing deferred acceptance algorithm. It is in the best interest of each man to report his true preferences (Dubins and Freedman, 1981, and Roth, 1982), but women typically have incentives to misreport their true preferences. Concerning her strategic options, a woman needs to consider only truncation strategies, which are the strategies obtained by removing a tail of men (i.e., some least preferred men) from her (true) ordered list of acceptable men. More precisely, for any (general) manipulation by a woman, there is a truncation strategy which is at least as good. We show that under the men-proposing deferred acceptance mechanism,

- truncating preferences by some women is weakly beneficial to all other women and weakly harmful to all men (Proposition 3.2), and
- any weakly successful group manipulation³ by women is weakly beneficial to all other women and weakly harmful to all men (Proposition 3.3).

Finally, we consider extending our results to the many-to-one college admissions model where students have to be assigned to colleges (with possibly multiple seats). A minor adaptation of the proof of Proposition 3.2 shows that under the student-proposing deferred acceptance mechanism, any truncation of preferences by some colleges is weakly beneficial to the other colleges and weakly harmful to all students. However, Kojima and Pathak (2009) showed that under the student-proposing deferred acceptance mechanism, truncation strategies typically do not exhaust the strategic options of the colleges. They proved that so-called dropping strategies constitute a class of exhaustive strategies. A dropping strategy of a college is obtained by removing some students from its (true) ordered lists of acceptable students (i.e., not necessarily a tail of least preferred students). We show that neither of our results

¹See also Roth (1984) and Roth (1985b) for further results on polarization of interests in two-sided markets.

²For the important role of the deferred acceptance algorithm in both matching theory and many real-life applications we refer to Roth (2008).

³That is, none of the manipulating agents is strictly worse off.

extends to the college admissions model in an appropriate way: there are dropping strategies and successful manipulations that strictly harm some other college and strictly benefit some student.

Our results complement work by Crawford (1991) who studied general many-to-one matching markets and investigated the effect of the entrance of an agent on the welfare of the other agents. When restricted to the marriage market, his result is the particular case of our first result in which a woman submits an empty truncation strategy.

2 Model

In Gale and Shapley's (1962) marriage market there are two non-empty, finite, and disjoint sets of agents M (men) and W (women). A generic man, woman, and agent are denoted by m , w , and i , respectively. Each agent i has a complete, transitive, and strict preference relation P_i over the agents on the other side of the market and the prospect of being alone. Let $P = (P_i)_{i \in M \cup W}$ denote the profile of all agents' preferences.

For $w, w' \in W \cup \{m\}$, we write $w P_m w'$ if man m strictly prefers w to w' ($w \neq w'$), and $w R_m w'$ if m likes w at least as well as w' ($w P_m w'$ or $w = w'$). Similarly, we write $m P_w m'$ and $m R_w m'$. A woman w is acceptable to a man m if $w P_m m$. Analogously, m is acceptable to w if $m P_w w$.

With some abuse of notation we also represent a man m 's preferences P_m as an ordered list of the elements in $W \cup \{m\}$. For instance, $P_m = w_3 w_2 m w_1 \dots w_4$ indicates that m prefers w_3 to w_2 and he prefers remaining single to any other woman. Similarly, woman w 's preferences can be represented as an ordered list P_w of the elements in $M \cup \{w\}$. We often omit the unacceptable agents from agent i 's ordered list P_i .

A marriage market is a triple (M, W, P) , or P for short. A matching is a one-to-one function μ from $M \cup W$ to itself, such that for each $m \in M$ and for each $w \in W$ we have $\mu(m) = w$ if and only if $\mu(w) = m$, $\mu(m) \notin W$ implies $\mu(m) = m$, and similarly $\mu(w) \notin M$ implies $\mu(w) = w$. If $\mu(m) = w$, then man m and woman w are matched to one another. If $\mu(i) = i$, then agent i is unmatched or single. Agent $\mu(i)$ is called i 's match at μ . We sometimes use a vector of men (or women) to denote a matching, e.g., $\mu = (m_3, m_1, m_2)$ denotes the matching where w_1 is matched to m_3 , w_2 to m_1 , and w_3 to m_2 .

A matching μ is individually rational if $\mu(i) R_i i$ for all $i \in M \cup W$. A pair (m, w) is a blocking pair for a matching μ if $w P_m \mu(m)$ and $m P_w \mu(w)$. A matching is stable if it is individually rational and if there are no blocking pairs. Gale and Shapley (1962) proved constructively that each marriage market has at least one stable matching. For this they

introduced the deferred acceptance (DA) algorithm. Let Q be a profile of ordered lists of acceptable agents. The men-proposing DA algorithm applied to Q , denoted by $DA(Q)$ for short, finds a matching through the following steps.

STEP 1: Each man m proposes to the woman that is ranked first in Q_m (if there is no such woman then m remains single). Each woman w tentatively accepts the best man among her proposers (using the list Q_w). All other proposers are rejected.

STEP k , $k \geq 2$: Each man m that is rejected in Step $k - 1$ proposes to the next woman in his list Q_m (if there is no such woman then m remains single). She tentatively accepts the best man among the new proposers and the tentatively matched man from the previous step, if any (using the list Q_w). All other proposers are rejected.

The algorithm stops when no man is rejected. Then, all tentative matches become final. With some abuse of notation, let $\mu(Q)$ denote the matching. For $i \in M \cup W$, let $\mu(Q, i)$ denote the match of agent i at $\mu(Q)$. Gale and Shapley (1962) proved that for preference profile Q matching $\mu(Q)$ is the best (worst) stable matching for the men (women). Dubins and Freedman (1981) and Roth (1982) proved that under the direct-revelation mechanism induced by μ it is a weakly dominant strategy for the men to reveal their true preferences. Therefore, we will assume that men are truthful and that women are the only strategic agents. Whenever there are at least two stable matchings some woman have incentives to misreport their true preferences (see for instance Roth and Sotomayor, 1990, Corollary 4.12).

3 Results

Before we present our results on the direct-revelation mechanism induced by the men-proposing deferred acceptance algorithm, we first provide the formal definitions of two classes of manipulations.

Let P be a marriage market. A truncation strategy (Roth and Rothblum, 1999) of a woman w is a strategy (or equivalently, an ordered list) P'_w obtained from P_w by making a tail of acceptable men unacceptable. Formally, P'_w is a **truncation strategy** if for all $m, m' \in M$, (a) [if $m R'_w m' R'_w w$ then $m R_w m' R_w w$], and (b) [if $m P'_w w$ and $m' P_w m$ then $m' P'_w w$].

A **(group) manipulation** by a group of women W' is a strategy-profile $P_{W'} = (P_w)_{w \in W'}$. If $|W'| = 1$, then $P_{W'}$ is an **individual manipulation**. A manipulation is **weakly successful** if for all $w \in W'$, $\mu(P', w) R_w \mu(P, w)$ where $P' = (P'_{W'}, P_{-W'})$. A manipulation is **successful** if for all $w \in W'$, $\mu(P', w) R_w \mu(P, w)$ and for some $w' \in W'$, $\mu(P', w') P_w \mu(P, w')$.

Note that not every truncation strategy is a weakly successful manipulation. For instance, an empty truncation strategy leaves the woman unmatched. Likewise, not every weakly successful, individual manipulation is a truncation strategy (see, for instance, Example 1). However, truncation strategies are exhaustive in the sense that any weakly successful, individual manipulation can be replicated or improved upon by some truncation strategy.⁴

The following well-known result states that men and women have opposite interests whenever a manipulation leads to a stable matching.

Lemma 3.1. *Under the men-proposing DA mechanism, a group manipulation by some women W' is weakly beneficial to all women and weakly harmful to all men if the induced matching is stable. If the matching is not stable then each blocking pair contains a woman from W' .*

Proof. Let $P'_{W'}$ be a group manipulation and let $P' = (P'_{W'}, P_{-W'})$. By assumption, $\mu(P')$ is stable for the market P . Hence, by men-optimality of $\mu(P)$, all women weakly prefer $\mu(P')$ to $\mu(P)$ and all men weakly prefer $\mu(P)$ to $\mu(P')$. The second statement follows from the observation that $\mu(P')$ is stable for P' and that for each pair (m, w) with $w \notin W'$, $P_m = P'_m$ and $P_w = P'_w$. \square

The following example illustrates that a manipulation may lead to an unstable matching, even if the manipulating women are strictly better off at the new matching.

Example 1. (A successful manipulation that yields an unstable matching.)

Consider the matching market with 3 men, 3 women, and preferences P given by the columns in the table below. For instance, $w_3 P_{m_1} w_1 P_{m_1} w_2 P_{m_1} m_1$. One easily verifies that $\mu(P) =$

Men			Women		
m_1	m_2	m_3	w_1	w_2	w_3
w_3	w_2	w_1	m_1	m_1	m_3
w_1	w_1	w_3	m_2	m_2	m_1
w_2	w_3	w_2	m_3	m_3	m_2

(m_3, m_2, m_1) — the boxed matching in the table. Suppose that woman w_1 submits the list $P'_{w_1} = m_2$. Then, $\mu(P') = (m_2, m_1, m_3)$ — the boldfaced matching in the table. Note P'_{w_1} is

⁴To see this, let P'_w be an individual manipulation. Let $m = \mu((P'_w, P_{-w}), w) \in M \cup \{w\}$. Consider the truncation strategy P''_w obtained from P_w by making all men that are strictly less preferred than m unacceptable. One easily verifies that $\mu((P''_w, P_{-w}), w) R_w \mu((P'_w, P_{-w}), w)$.

a successful manipulation since $\mu(P', w_1) = m_2 P_{w_1} m_3 = \mu(P, w_1)$. But $\mu(P')$ is not stable with respect to the true preferences P (the unique blocking pair is (m_1, w_1)). \diamond

In Example 1, all women that do not manipulate weakly benefit and all men are weakly harmed. Since the resulting matching is not stable this observation does not follow from Lemma 3.1. Nevertheless, we will prove that the observed opposed interests are a feature of two interesting classes of group manipulations: group truncation strategies and weakly successful group manipulations.

For marriage markets, the next proposition generalizes the results of Crawford (1991) from an *individual empty* truncation strategy to *arbitrary group* truncation strategies. We include the proof, which is similar to that of Crawford (1991), for two reasons. First, for marriage markets the arguments are shorter and more transparent. Second, it will be useful in pointing out why the same arguments do not immediately carry over to other manipulations.

To prove our results we introduce the following additional notation. For every integer $k \geq 1$, let $X(Q, w, k)$ be the set of men that will have proposed to woman w by step k under DA(Q), i.e., in some step $l \in \{1, \dots, k\}$ of DA(Q). Let $X(Q, w)$ be the set of men that will have proposed to w by the last step of DA(Q), i.e., $X(Q, w) = \cup_k X(Q, w, k)$.

Proposition 3.2. *Under the men-proposing DA mechanism, any group manipulation by women that consists of truncation strategies is weakly beneficial to the other women and weakly harmful to all men.*

Proof. Let $P'_{W'}$ be a group manipulation of some women W' such that for each $w' \in W'$, $P'_{w'}$ is a truncation strategy. Let $P' = (P'_{W'}, P_{-W'})$. It is sufficient to show that for each woman w and each step k , $X(P, w, k) \subseteq X(P', w, k)$. For $k = 1$ the inclusion is in fact an equality since at step 1 of DA(P) and DA(P') each man proposes to exactly the same woman.

Assume that the inclusion holds for k . We will show that the inclusion also holds for $k + 1$. Let $m \in X(P, w, k + 1)$. If $m \in X(P, w, k)$, then by induction, $m \in X(P', w, k)$, and hence $m \in X(P', w, k) \subseteq X(P', w, k + 1)$. So, assume $m \in X(P, w, k + 1) \setminus X(P, w, k)$. Then, in DA(P), man m proposed to w at step $k + 1$ but not at step k . So, m was rejected by some woman $\bar{w} \neq w$ at step k of DA(P). By the induction hypothesis, $m \in X(P, \bar{w}, k) \subseteq X(P', \bar{w}, k)$. If $\bar{w} \notin W'$ then \bar{w} will also have rejected m by step k of DA(P') since $P'_{\bar{w}} = P_{\bar{w}}$. If $\bar{w} \in W'$ then \bar{w} will also have rejected m by step k of DA(P') since $P'_{\bar{w}}$ is a truncation strategy obtained from $P_{\bar{w}}$. Since m makes his proposals in the same order in DA(P) and DA(P'), he will have proposed to w by step $k + 1$ of DA(P'). Hence, $m \in X(P', w, k + 1)$. \square

The following example shows that if we replace (possibly unsuccessful) truncation strategies in the statement of Proposition 3.2 by weakly successful manipulations then the key

argument in the proof does no longer work.

Example 2. (A successful manipulation with a “rejection lag.”)

Consider the matching market with 3 men, 3 women, and preferences P given by the columns in the table below. One easily verifies that $\mu(P) = (m_1, m_3, m_2)$ — the boxed

Men			Women		
m_1	m_2	m_3	w_1	w_2	w_3
w_1	w_1	w_2	m_3	m_2	m_1
w_2	w_3	w_1	m_1	m_1	m_2
w_3	w_2	w_3	m_2	m_3	m_3

matching in the table. Suppose that woman w_1 submits the list $P'_{w_1} = m_3, m_2, m_1$. Then, $\mu(P') = (m_3, m_1, m_2)$ — the boldfaced matching in the table. Note that P'_{w_1} is a successful manipulation since $\mu(P', w_1) = m_3 P_{w_1} m_1 = \mu(P, w_1)$.

Note that all other women weakly benefit and all men are weakly hurt by the manipulation. However, the arguments in the proof of Proposition 3.2 cannot be directly applied here. This can be seen as follows. In $DA(P)$, woman w_1 rejects m_2 in the first step (after which he proposes to w_3 , gets accepted, and the algorithm halts). In $DA(P')$, the manipulating woman w_1 will reject m_2 in a later step (i.e., not in the first step). Therefore, $X(P, w_1, 1) \not\subseteq X(P', w_1, 1)$. Hence, the arguments of Proposition 3.2 cannot be applied to tackle successful manipulations that are not truncation strategies. \diamond

Our second result shows that the conflict and coincidence of interests as observed in Example 2 holds in fact for any weakly successful group manipulation. In other words, we can replace the (possibly unsuccessful) truncation strategies in Proposition 3.2 by weakly successful manipulations.

Proposition 3.3. *Under the men-proposing DA mechanism, any weakly successful group manipulation by women is weakly beneficial to the other women and weakly harmful to all men.*

Proof. Let $P'_{W'}$ be a weakly successful manipulation of a group of women W' and let $P' = (P'_{W'}, P_{-W'})$. It is sufficient to show that for each woman w and each step k , $X(P, w, k) \subseteq X(P', w)$. For $k = 1$ the inclusion is obvious since at step 1 of $DA(P)$ and $DA(P')$ each man proposes to exactly the same woman.

Assume that the inclusion holds for k . We will show that the inclusion also holds for $k + 1$. Let $m \in X(P, w, k + 1)$. If $m \in X(P, w, k)$, then by induction, $m \in X(P', w)$. So,

assume $m \in X(P, w, k + 1) \setminus X(P, w, k)$. Then, in $\text{DA}(P)$, man m proposed to w at step $k + 1$ but not at step k . So, m was rejected by some woman $\bar{w} \neq w$ at step k of $\text{DA}(P)$. By the induction hypothesis, $m \in X(P, \bar{w}, k) \subseteq X(P', \bar{w})$. If $\bar{w} \notin W'$ then \bar{w} will also reject m in $\text{DA}(P')$ since $P'_{\bar{w}} = P_{\bar{w}}$. If $\bar{w} \in W'$ then $\mu(P', \bar{w}) R_{\bar{w}} \mu(P, \bar{w}) P_{\bar{w}} m$, which implies that in the last step of $\text{DA}(P')$ woman \bar{w} is matched to a man she strictly prefers to m (according to her true preferences). Therefore, in either case \bar{w} will also eventually reject m in $\text{DA}(P')$. Since m makes his proposals in the same order in $\text{DA}(P)$ and $\text{DA}(P')$, he will have proposed to w by the last step of $\text{DA}(P')$. Hence, $m \in X(P', w)$. \square

Finally, we consider extending our results to the many-to-one college admissions model where students have to be assigned to colleges with possibly multiple seats, strict preferences over individual students, and responsive preferences over groups of students.⁵ Note that the men-proposing DA algorithm and some of its properties can be straightforwardly generalized to college admissions (such that the men “become” students, and the women “become” colleges with possibly multiple seats). In particular, it is a weakly dominant strategy for the students to submit their true preferences (Roth, 1985a, Theorem 5*) under the mechanism induced by the student-proposing DA algorithm, which will be denoted by μ . Note that a college can manipulate not only its ordered list of students but also the number of available seats, i.e., the strategy space is much richer than in one-to-one markets.

A minor adaptation of the proof of Proposition 3.2 shows that under the student-proposing DA mechanism, any group manipulation by colleges that consists of truncation strategies is weakly beneficial to the other colleges and weakly harmful to all students. However, Kojima and Pathak (2009) showed that under the student-proposing deferred acceptance mechanism, truncation strategies typically do *not* exhaust the strategic options of the colleges. More precisely, they presented a many-to-one market in which for some college there is a strategy such that any truncation strategy yields a strictly worse match. They also proved that so-called dropping strategies constitute a class of exhaustive strategies. A dropping strategy of a college is obtained by removing some students from its (true) ordered lists of acceptable students (i.e., not necessarily a tail of least preferred students).⁶ Formally, for a college c with preferences P_c over individual students, P'_c is a **dropping strategy** if for all students s, s' , $[s R'_c s' R'_c \emptyset$ implies $s R_c s' R_c \emptyset]$.

Therefore, a possible appropriate extension of Proposition 3.2 to college admissions would involve dropping strategies rather than truncation strategies. The next example, however,

⁵For a formal definition of the college admissions model and responsiveness in particular, see Roth (1985a).

⁶The fact that dropping strategies are exhaustive implies that it suffices to focus on each college’s submittable ordered lists of students.

shows that neither of our results extends to the college admissions model in an appropriate way: there are dropping strategies and successful manipulations that strictly harm some other college and strictly benefit some student.

Example 3. (Propositions 3.2 and 3.3 cannot be appropriately generalized to college admissions.)

Consider the following matching market with students s_1, s_2, s_3 , and s_4 , and colleges c_1 and c_2 . Each college has two seats. The preferences P over *individual* agents are given by the columns in the table below. We assume that the colleges' preferences over sets of students are responsive to the preferences over individual students and that both colleges prefer $\{s_1, s_4\}$ to $\{s_2, s_3\}$.⁷ One easily verifies that $\mu(P) = (c_2, c_1, c_1, c_2)$ —the boxed matching

Students				Colleges	
s_1	s_2	s_3	s_4	c_1	c_2
c_2	c_1	c_1	c_1	s_1	s_4
c_1	c_2	c_2	c_2	s_2	s_2
				s_3	s_3
				s_4	s_1

in the table. Suppose that college c_1 submits the dropping strategy $P'_{c_1} = s_1, s_4$. Then, $\mu(P') = (c_1, c_2, c_2, c_1)$ —the boldfaced matching in the table. Note that P'_{c_1} is a successful dropping strategy since college c_1 prefers $\{s_1, s_4\}$ to $\{s_2, s_3\}$. Since college c_2 is strictly worse off and student s_4 is strictly better off under $\mu(P')$ it follows that Propositions 3.2 and 3.3 cannot be appropriately extended to college admissions. \diamond

Remark 1. In fact, using the many-to-one market in Example 3 one can construct a marriage market in which an individual (unsuccessful) dropping strategy of a woman makes another woman strictly worse off and some man strictly better off (cf. Proposition 3.3).⁸ For two reasons we do not provide further details and present Example 3 instead. First, the class of dropping strategies contains the strictly smaller class of truncation strategies, which is already exhaustive for one-to-one markets. Second, the market in Example 3 shows not only the impossibility of appropriately generalizing Proposition 3.2 but also the impossibility of generalizing Proposition 3.3. \diamond

Finally, we note that Example 3 uncovers another difference between marriage markets and college admissions and adds to those already identified in Roth (1985a).

⁷Note that preferring $\{1, 4\}$ to $\{2, 3\}$ is compatible with responsiveness.

⁸We thank Bettina Klaus for pointing this out.

References

- [1] V.P. Crawford (1991). Comparative Statics in Matching Markets. *Journal of Economic Theory* 54(1): 389–400.
- [2] L.E. Dubins and D.A. Freedman (1981). Machiavelli and the Gale-Shapley Algorithm. *American Mathematical Monthly* 88(7): 485–494.
- [3] D. Gale and L.S. Shapley (1962). College Admissions and the Stability of Marriage. *American Mathematical Monthly* 69(1): 9–15.
- [4] D.E. Knuth (1976). *Mariages Stables*. Montréal: Les Presses de l’Université de Montréal.
- [5] F. Kojima and P.A. Pathak (2009). Incentives and Stability in Large Two-Sided Matching Markets. *American Economic Review*, 99(3): 608–627.
- [6] A.E. Roth (1982). The Economics of Matching: Stability and Incentives. *Mathematics of Operations Research* 7(4): 617–628.
- [7] A.E. Roth (1984). Stability and Polarization of Interests in Job Matching. *Econometrica* 52(1): 47–58.
- [8] A.E. Roth (1985a). The College Admission Problem is not Equivalent to the Marriage Problem. *Journal of Economic Theory* 36(2): 277–288.
- [9] A.E. Roth (1985b). Conflict and Coincidence of Interest in Job Matching: Some New Results and Open Questions. *Mathematics of Operations Research* 10(3): 379–389.
- [10] A.E. Roth (2008). Deferred Acceptance Algorithms: History, Theory, Practice, and Open Questions. *International Journal of Game Theory* 36(3): 537–569.
- [11] A.E. Roth and U.G. Rothblum (1999). Truncation Strategies in Matching Markets – In Search of Advice for Participants. *Econometrica* 67(1): 21–43.
- [12] A.E. Roth and M.A.O. Sotomayor (1990). *Two-Sided Matching: A Study in Game-Theoretic Modeling and Analysis*. Econometric Society Monograph Series. New York: Cambridge University Press.