

## 1 Social Choice Theory

Social choice theory deals with the aggregation of individual preferences over some number of alternatives. The arising aggregated preferences then reflect some sort of social preferences, and like individual preferences, we want these to satisfy certain conditions. For example, we may want unanimous agreement on the best option to be maintained in the social ranking by placing this option on top.

### 1.1 Notation

We consider the following  $\mathcal{E} = \{X, \succeq_i\}_{i=1}^n$  as representing a social group or economy, where  $n$  is the number of agents in the economy. We consider identical sets of alternatives  $X$  for each individual, and heterogeneous, rational (complete and transitive) preference relations. We denote with  $R$  the set of all possible rational (weak) preferences relations on  $X$  and  $P$  as the set of all strict rational preferences. A typical element of  $R \times R \times \dots \times R = R^n$  is  $(\succeq_1, \dots, \succeq_n)$ , which will be called a preference profile. If  $X$  is finite, then each profile is a collection of  $n$  rankings of the alternatives.

## 2 Social Welfare Functions

**Definition 2.1:** A **social welfare function** is a rule  $F : A \rightarrow R$  (with  $A = R^n$  or  $P^n$ ) that assigns a rational preferences relation  $F(\succeq_1, \dots, \succeq_n) \in R$ , interpreted as the social preference relation, to any profile of individual rational preferences in  $A$ .

**Definition 2.2:** A social welfare function  $F : A \rightarrow R$  is **Pareto-efficient** or **Paretian** if, for all alternatives  $x, y \in X$  and all preference profiles  $\succeq = (\succeq_1, \dots, \succeq_n) \in A$ , we have

$$\forall i \in n, x \succ_i y \implies x F(\succeq) y$$

**Definition 2.3:** A social welfare function  $F : A \rightarrow R$  satisfies **independence of irrelevant alternatives (IIA)** if for all  $x, y \in X$  and for all pairs of profiles  $\succeq, \succeq' \in A$  with the property

$$x \succeq_i y \iff x \succeq'_i y \quad \text{and} \quad y \succeq_i x \iff y \succeq'_i x$$

we have

$$x F(\succeq) y \iff x F(\succeq') y \quad \text{and} \quad y F(\succeq) x \iff y F(\succeq') x$$

**Definition 2.4:** A social welfare function  $F : A \rightarrow R$  is **dictatorial** if there is an agent  $h \in n$  such that for all  $x, y \in X$  and for all  $\succeq \in A$ , we have  $x \succeq_h y \implies x F(\succeq) y$ .

## 3 Social Choice Functions

Social preferences aren't that useful unless they are used to make some sort of social choice. Given that, we can represent the process of making a choice from a set of options based on a preference profile as a single function.

**Definition 3.1:** A **social choice function** is a rule  $f : A \rightarrow X$  (with  $A = R^I$  or  $P^I$ ) that assigns a chosen element  $f(\succeq_1, \dots, \succeq_I) \in X$  to every profile of individual rational preference relations in  $A$ .

At the moment, we are assuming a single-valued function, but this can readily be extended to a correspondence following Maskin.

**Definition 3.2:** A social choice function  $f : A \rightarrow X$  is **weakly Pareto efficient** is, for all preference profiles  $\succeq \in A$ , the choice  $f(\succeq) \in X$  is a weak Pareto optimum, i.e. for all  $x, y \in X$  such that  $x \succ_i y$  for all  $i \in n$ , then  $y \neq f(\succeq)$ .

**Definition 3.3:** The alternative  $x \in X$  **maintains its position** from  $\succeq \in R^n$  to  $\succeq' \in R^n$  if  $x \succeq_i y \implies x \succeq'_i y$  for all  $i \in n$  and  $y \in X$ . Equivalently, the lower contour sets of  $x$  in the preferences of the first profile are subsets of the lower contour sets of the corresponding preferences in the second profile.

**Definition 3.4:** A social choice function  $f : A \rightarrow X$  is **monotonic** if for all profiles  $\succeq, \succeq' \in A$  where  $x = f(\succeq)$  maintains its position from  $\succeq$  to  $\succeq'$ , we have  $f(\succeq') = x$  again.

**Definition 3.5:** A social choice function is **dictatorial** if there exists an agent  $h \in n$  such that for all profiles  $\succeq \in A$ , we have  $f(\succeq) = \{x \mid \forall y \in X, x \succeq_h y\}$ , i.e.  $f(\succeq)$  is a maximal element for  $h$  over  $X$ .

**Definition 3.6:** A social choice function is **strategy-proof** if for all  $i \in n$ , preferences profile  $\succeq \in A$ , and alternate preference  $\succeq'_i \in R$ , we have  $f(\succeq_i, \succeq_{-i}) \succ_i f(\succeq'_i, \succeq_{-i})$ .

## 4 Impossibility Results

**Theorem 4.1:** (Arrow 1950) Suppose  $|X| \geq 3$  and  $A = R^n$  or  $P^n$ . Then every social welfare function  $F : A \rightarrow R$  that is Pareto-efficient and satisfies IIA is dictatorial.

**Theorem 4.2:** (Muller and Satterthwaite 1977) Suppose  $|X| \geq 3$  and  $A = R^n$  or  $P^n$ . Then every social choice function  $f : A \rightarrow X$  that is weakly Pareto-efficient and monotonic is dictatorial.

**Theorem 4.3:** (Gibbard 1973, Satterthwaite 1977) If  $|X| \geq 3$  and the social choice function  $f : A \rightarrow X$  is onto (surjective) and strategy-proof, then  $f$  is dictatorial.

## 5 Mechanism Design

With social choice theory, we looked at how to aggregate preferences assuming these were observable. Gibbard-Satterthwaite points us toward the direction of imperfectly observable preferences that must be revealed by the individuals before they can be aggregated.

We again consider a set of agents  $n$  which will make a collective choice from a set of alternatives  $X$ . Utility functions  $u_i(x, \theta_i)$  are commonly known, but are determined by a privately known random variable  $\theta_i \in \Theta_i$  which is the agent's type. Typically there is a commonly known prior  $\phi(\cdot)$  representing a density over  $\theta \in \Theta_1 \times \dots \times \Theta_n$ . Since the utility functions themselves are common knowledge in this context, all uncertainty handled is through the type. Types can also encode the beliefs of an agent. Our new formulation of a social choice function is a straightforward extension based on this.

**Definition 5.1:** A **social choice function** is a rule  $f : \Theta_1 \times \dots \times \Theta_n \rightarrow X$  that assigns a chosen element  $f(\theta) \in X$  for each profile of individual preferences derived from the vector of types  $\theta$ .

Some properties are also simple extensions:

**Definition 5.2:** The social choice function  $f : \Theta \rightarrow X$  is **ex-post efficient** if for all  $\theta \in \Theta$  there does not exist an  $x \in X$  such that  $u_i(x, \theta_i) \geq u_i(f(\theta), \theta_i)$  for all  $i$  and  $u_i(x, \theta_i) > u_i(f(\theta), \theta_i)$  for some  $i$ .

## 6 Mechanisms and Implementation

Since agents' types are private information, they must communicate this information to the rest of the group or a central body. We could have any theory about how agents choose the messages they will send, but we typically assume messages will be chosen strategically. Once the messages are collected, an outcome will be chosen by the group. Together, the possible messages and the mapping from messages to outcomes define a mechanism.

**Definition 6.1:** A **mechanism**  $\xi = \{(M_1, \dots, M_n), g(\cdot)\}$  is a collection of  $n$  message sets (or strategy sets) and an outcome function  $g : M_1 \times \dots \times M_n \rightarrow X$ .

Mechanisms can be considered as procedures for how a collective decision is actually made, including a specification of strategy sets. We can now compare what outcomes are realized in equilibrium through a mechanism with the outcomes we "want" to see happen, as represented by a social choice function. If an outcome function  $g$  and a

social choice function  $f$  coincide on some restriction of their domain, then we say  $g$  implements  $f$ . The particular restriction depends on the equilibrium concept used.

**Definition 6.2:** The mechanism  $\xi = \{(M_1, \dots, M_n), g(\cdot)\}$  **implements** a social choice function  $f$  if there is an equilibrium strategy profile  $(m_1^*(\cdot), \dots, m_n^*(\cdot))$  of the game induced by  $\xi$  such that  $\forall \theta \in \Theta$ ,

$$g(m_1^*(\theta_1), \dots, m_n^*(\theta_n)) = f(\theta_1, \dots, \theta_n)$$

The term equilibrium here is still vague and we will extend this definition to talk about implementation in various contexts. The mechanism itself might have multiple equilibria. If the outcome of one equilibrium coincides with the scf, then the scf is *partially* implemented. If the outcomes of all equilibria coincide with the scf, then the scf is *fully* implemented. Figure 1 depicts the generic implementation problem.

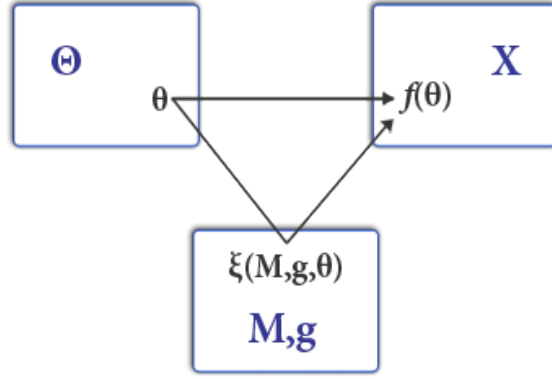


Figure 1: Reiter diagram of a mechanism  $\xi$  implementing  $f$

Due to the revelation principle, we can usually think of the message space as the type space, and most mechanisms we'll look at are direct revelation mechanisms. This also leads us to a notion of incentive compatibility.

**Definition 6.3:** A **direct revelation mechanism** for  $f$  is a mechanism in which  $M_i = \Theta_i$  for all  $i$  and  $g(\theta) = f(\theta)$  for all  $\theta \in \Theta_1 \times \dots \times \Theta_I$ .

**Definition 6.4:** The social choice function  $f$  is **truthfully implementable** if the direct revelation mechanism has an equilibrium  $(m_1^*(\cdot), \dots, m_n^*(\cdot))$  in which  $m_i^*(\theta_i) = \theta_i$  for all  $i$  and  $\theta$ , i.e. each agent telling the truth is an equilibrium of the induced game. Alternatively, such a mechanism is **incentive compatible**.

## 7 Dominant Strategy Implementation

**Definition 7.1:** A mechanism  $\xi = \{(M_1, \dots, M_n), g(\cdot)\}$  is a **dominant-strategy implementation** of a social choice function  $f$  if there is a dominant strategy equilibrium profile  $(m_1^*(\cdot), \dots, m_n^*(\cdot))$  of the game induced by  $\xi$  such that  $\forall \theta \in \Theta$ ,  $g(m_1^*(\theta_1), \dots, m_n^*(\theta_n)) = f(\theta_1, \dots, \theta_n)$ .

Here we start with some scf  $f$  and ask whether a mechanism exists such that the outcome from agents playing a dominant strategy matches the scf for all realization of types. If this was the case, then we would have a very robust way of implementing the scf, since we can be confident about what agents will want to do, regardless of what other agents do or what their types are.

Incentive compatibility in this context is also a simple extension.

**Definition 7.2:** The scf  $f(\cdot)$  is **truthfully implementable in dominant strategies** or **strategy-proof** if for all  $i$  and  $\theta_i$ ,

$$u_i(f(\theta), \theta_i, \theta_{-i}) \geq u_i(f(\hat{\theta}_i, \theta_{-i}), \hat{\theta}_i, \theta_{-i})$$

for all  $\hat{\theta}_i$  and  $\theta_{-i}$ .

Now we have the powerful, if nearly trivial, revelation principle for dominant strategy implementation:

**Theorem 7.1:** Suppose there exists a mechanism  $\xi$  that implements  $f$  in dominant strategies. Then  $f$  is truthfully implementable in dominant strategies.

Checking whether an arbitrary scf can be implemented in dominant strategies seems like a monumental task; imagine trying to search over all possible procedures for making a choice. By this theorem though, we know if there is such a mechanism, there is also a direct revelation mechanism where agents report their types honestly. Hence, checking implementability reduces to checking the inequalities in Definition 7.2. This holds for many other equilibrium concepts, which justifies looking only at revelation mechanisms. Of course, this assumes it's practical for agents to fully communicate their types, which could potentially include a full set of beliefs, beliefs about beliefs, etc, without cost, so we might be trading one intractability for another.

## 8 Nash Implementation

Maskin focuses on general preferences and social choice correspondences, which leads to the following definition:

**Definition 8.1:** A mechanism  $\{(S_1, \dots, S_I), g(\cdot)\}$  with  $g : \prod S_i \rightarrow X$  implements the social choice rule  $f : \prod R_i \rightarrow X$  in Nash equilibrium if

1.  $\forall x \in f(\succeq), \exists (s_1, \dots, s_I)$  such that  $x = g(s_1, \dots, s_I)$  and  $x \succeq_i g(s'_i, s_{-i})$  for all  $i$  and all  $s'_i$ .
2. If  $(s_1, \dots, s_I)$  is a Nash equilibrium, then  $g(s) \in f(\succeq)$ .

The first part extends our implementation to correspondences, while the second requires that all Nash equilibria result in outcomes in the social choice rule. Then, Maskin gives us the following result:

**Theorem 8.1:** If  $f : R^n \rightarrow X$  is implementable in Nash equilibrium, then  $f$  is monotonic.

This provides an easy first check to see whether a scr is Nash-implementable. However, the converse is not true. We'll need another condition:

**Definition 8.2:** A scr  $f : R^n \rightarrow X$  satisfies **no veto power (NVP)** if for all preference profiles  $\succeq$  and all  $x \in X$ , then for all agents  $i$  we have

$$\forall j \neq i, \forall y \in X, x \succeq_j y \implies x \in f(\succeq)$$

**Theorem 8.2:** With at least three players, if  $f$  is monotonic and satisfies NVP, then  $f$  is Nash-implementable.

## 9 Exercises

**Exercise 9.1** (p. 57 of *Toolbox for Economic Design*): Assume a social choice function  $f : \Theta \rightarrow X$  is dominant-strategy incentive compatible. Let  $\bar{f} : \bar{\Theta} \rightarrow X$ , where  $\bar{\Theta} \subset \Theta$ , be the restriction of  $f$  to  $\bar{\Theta}$ , i.e.  $\theta \in \bar{\Theta} \implies \bar{f}(\theta) = f(\theta)$ . Is  $\bar{f}$  also dominant-strategy incentive compatible?

**Exercise 9.2** (p. 56 of *Toolbox for Economic Design*): Imagine a society that consists of a single couple, Link and Zelda. The couple is trying to decide whether to paint their castle black, white, or silver. Zelda has one type  $\theta_Z$ . Her ranking is white > silver > black. Link, on the other hand, has two possible types:  $\theta_L$  and  $\bar{\theta}_L$ . When his type is  $\theta_L$ , his ranking is black > silver > white. When he is type  $\bar{\theta}_L$ , his ranking is silver > white > black.

1. Suppose we wish to implement the following social choice function:  $f(\theta_Z, \theta_L) = \text{silver}$  and  $f(\theta_Z, \bar{\theta}_L) = \text{white}$ . Will Link choose to truthfully reveal his type?
2. What if the desired social choice function is  $f(\theta_Z, \theta_L) = \text{black}$  and  $f(\theta_Z, \bar{\theta}_L) = \text{white}$ ?
3. What about  $f(\theta_Z, \theta_L) = \text{silver}$  and  $f(\theta_Z, \bar{\theta}_L) = \text{silver}$ ?

**Exercise 9.3** (Sprumont 1991): One useful restriction of preferences is single-peakedness. If the set of options  $X$  is single-dimensional and ordered, then a preference ordering  $\succ$  on  $X$  is single-peaked iff there is a maximum  $x^*$  of  $\succ$  and  $x^* < x < y$  or  $y < x < x^* \implies x \succ y$ .

Consider a partnership of  $n$  individuals who will invest in a project, with the benefits shared in proportion to each partner's investment. The project has a fixed cost of 1. The partners have single-peaked preferences over the amount they want to invest, with a peaks  $x_i^* \in [0, 1]$ . Because the sum of the peak amount of all partners may not be equal to the cost of the project, some partners may be forced to invest more or less than their ideal amounts. In this context, efficiency means that if the sum of ideal investments is less than the cost, everyone must invest at least their ideal amount, and vice versa.

In addition to strategy-proofness and efficiency, let's consider two other desirable properties of social choice functions. First, an scf is anonymous if  $f(\theta) = f(\pi(\theta))$  for all permutations  $\pi$  of the vector. Second, an scf is envy-free if for all  $i, j$ ,  $f_i(\theta) \succeq_i f_j(\theta)$ , i.e. the proportion allocated to  $i$  is preferred by  $i$  to all other agents' allocations.

Check whether the following are strategy-proof, efficient, anonymous, and/or envy-free:

1. The egalitarian rule  $f_i^e(\theta) = 1/n$ .
2. The proportional rule  $f_i^p(\theta) = x_i^* / \sum x_j^*$ .
3. The priority rule  $f_1^q(\theta) = x_1^*$  and  $f_i^q(\theta) = \min\{x_i^*, 1 - \sum_{j < i} x_j^*\}$  when  $\sum x_i^* \geq 1$ , and alternately  $f_j^q(\theta) = x_j^*$  for  $j < n$  and  $f_n^q(\theta) = 1 - \sum_{j=1}^{n-1} x_j^*$  when  $\sum x_i^* < 1$ .
4. The priority rule as defined above, but with a random order.
5. The uniform rule defined by the iterative process (for the case of  $\sum x_i^* > 1$ ):
  - (a) Start with all partners active and the full cost outstanding.
  - (b) Divide the outstanding cost equally among the active partners.
  - (c) If any active partner has an ideal below the equal share, set their share equal to their ideal and subtract this from the outstanding cost. These partners are now inactive.
  - (d) Repeat the previous two steps among the active partners until each has an ideal amount no less than the equal share of the outstanding cost.