

Logic and Social Choice Theory

Ulle Endriss

Institute for Logic, Language and Computation
University of Amsterdam

[<http://www.illc.uva.nl/~ulle/teaching/esslli-2013/>]

Plan for Today

References to “logic” in classical social choice theory are mostly about the axiomatic method, which is logic-like in spirit but doesn’t make use of a formal language with an associated semantics and proof theory.

Today’s lecture is about *logics for social choice*: embedding parts of the theory of social choice into a logical system.

We will first review various arguments for *why this is useful* and then see three concrete approaches that use different logics to model the Arrovian framework of *preference aggregation*:

- an approach based on a specifically designed *modal logic*;
- an approach using *classical first-order logic*; and
- an approach using *classical propositional logic*.

Logics for Social Choice

Our goal today will be to embed part of SCT into a formal logic.

Roughly, *models* of the logic should encode *aggregators* and *formulas* should encode their *properties*.

Why would we want to do this? Standard answers for any such an exercise in formalisation include:

- Because the act of formalisation has the potential to help us gain a *deeper understanding* of the domain we are formalising.
- Because we are *interested in a particular logical system* and want to explore its expressive power.

These are valid arguments, but there is more.

Verification

Logic has long been used to formally specify computer systems, enabling formal and automatic verification. Maybe we can apply a similar methodology to social choice mechanisms?

Parikh has coined the term “*social software*” for this research agenda.

Besides checking whether a given mechanism satisfies a given property (\rightsquigarrow *model checking*), we may also try to formally verify theorems from social choice theory (\rightsquigarrow *automated theorem proving*).

Example: Arrow’s original proof was not entirely correct. Nowadays this is not an issue anymore, but it could be for new results.

R. Parikh. Social Software. *Synthese*, 132(3):187–211, 2002.

Formal Minimalism

Pauly (2008) argues that when judging the appropriateness of an axiom in social choice theory, besides its *normative appeal* and its *mathematical strength*, we should also consider the *expressivity* of the language used to define it: the less the better (*formal minimalism*).

A related point:

- IIA, making reference to both the profile under consideration and another counterfactual profile, is less appealing than the
- Pareto condition, which just says what to do in the profile at hand.

To make such issues precise, we need a formal language for axioms.

M. Pauly. On the Role of Language in Social Choice Theory. *Synthese*, 163(2):227–243, 2008.

Modelling the Arrovian Framework

Recall the Arrovian framework of *social welfare functions*, for a finite set \mathcal{N} of individuals and an arbitrary set \mathcal{X} of alternatives:

A SWF is a function $F : \mathcal{L}(\mathcal{X})^n \rightarrow \mathcal{L}(\mathcal{X})$ mapping any given profile of preference orders (i.e., linear orders) to a collective preference order.

F is defined on all profiles in $\mathcal{L}(\mathcal{X})^n$ (*universal domain* assumption).

Arrow suggested the following axioms (desirable properties of F):

- *Pareto*: if all individuals rank $x \succ y$, then so does society
- *IIA*: whether society ranks $x \succ y$ depends only on who ranks $x \succ y$
- *Nondictatorship*: F does not just copy the \succ of a fixed individual

Arrow's Theorem establishes that no SWF F satisfies all three axioms, if there are ≥ 3 alternatives. This holds for any finite set of individuals.

► Can we express these things in a suitable logic?

Approach 1: Modal Logic

One approach to take is to develop a *new logic* specifically aimed at modelling the aspect of social choice theory we are interested in.

Modal logic looks like a useful technical framework for doing this.

It is intuitively clear that we can (somehow) devise a modal logic that can capture the Arrovian framework of SWF's, but how to do it exactly is less clear and finding a good way of doing this is a real challenge.

Adopting a semantics-guided approach, we first have to decide:

- what do we take to be our possible worlds?, and
- what accessibility relation(s) should we define?

Next, we shall review a specific proposal due to Ågotnes et al. (2011).

T. Ågotnes, W. van der Hoek, and M. Wooldridge. On the Logic of Preference and Judgment Aggregation. *J. Auton. Agents Multiagent Sys.*, 22(1):4–30, 2011.

Frames

Given: fixed (and finite) \mathcal{N} (n individuals) and \mathcal{X} (m alternatives)

Each *possible world* consists of

- a profile R and
- a pair (x, y) of alternatives.

There are two *accessibility relations* defined on the possible worlds:

- Two worlds are related via relation PROF if their associated pairs are identical (i.e., only their profiles differ, if anything).
- Two worlds are related via relation PAIR if their associated profiles are identical (i.e., only their pairs differ, if anything).

A *frame* $\langle \mathcal{L}(\mathcal{X})^n \times \mathcal{X}^2, \text{PROF}, \text{PAIR} \rangle$ consists of the set of worlds and the two accessibility relations (all induced by \mathcal{N} and \mathcal{X}).

Language

The language of the logic has the following *atomic* propositions:

- p_i for every individual $i \in \mathcal{N}$
Intuition: p_i is true at world $\langle \mathbf{R}, (x, y) \rangle$ if $x \succ y$ according to R_i
- $q_{(x,y)}$ for every pair of alternatives $(x, y) \in \mathcal{X}^2$
Intuition: $q_{(x',y')}$ is true at world $\langle \mathbf{R}, (x, y) \rangle$ if $(x, y) = (x', y')$
- a special proposition σ
Intuition: σ is true at world $\langle \mathbf{R}, (x, y) \rangle$ if society ranks $x \succ y$

The set of *formulas* φ is defined as follows:

$$\varphi ::= p_i \mid q_{(x,y)} \mid \sigma \mid \neg\varphi \mid \varphi \wedge \varphi \mid [\text{PROF}]\varphi \mid [\text{PAIR}]\varphi$$

Disjunction, implication, and diamond-modalities are defined in the usual manner (e.g., $\langle \text{PROF} \rangle \varphi \equiv \neg[\text{PROF}]\neg\varphi$).

Semantics

In modal logic, a *valuation* determines which atomic propositions are true in which world, and a frame and a valuation together define a *model*. For this logic, the valuation of p_i and $q_{(x,y)}$ is fixed and the valuation of σ will be defined in terms of a SWF F .

So, for given and fixed \mathcal{N} and \mathcal{X} (and thus for a fixed frame), we now define *truth* of a formula φ at a world $\langle \mathbf{R}, (x, y) \rangle$ wrt. a SWF F :

- $F, \langle \mathbf{R}, (x, y) \rangle \models p_i$ iff $(x, y) \in R_i$
- $F, \langle \mathbf{R}, (x, y) \rangle \models q_{(x', y')}$ iff $(x, y) = (x', y')$
- $F, \langle \mathbf{R}, (x, y) \rangle \models \sigma$ iff $(x, y) \in F(\mathbf{R})$
- $F, \langle \mathbf{R}, (x, y) \rangle \models \neg\varphi$ iff $F, \langle \mathbf{R}, (x, y) \rangle \not\models \varphi$
- $F, \langle \mathbf{R}, (x, y) \rangle \models \varphi \wedge \psi$ iff $F, \langle \mathbf{R}, (x, y) \rangle \models \varphi$ and $F, \langle \mathbf{R}, (x, y) \rangle \models \psi$
- $F, \langle \mathbf{R}, (x, y) \rangle \models [\text{PROF}]\varphi$ iff $F, \langle \mathbf{R}', (x, y) \rangle \models \varphi$ for all profiles \mathbf{R}'
- $F, \langle \mathbf{R}, (x, y) \rangle \models [\text{PAIR}]\varphi$ iff $F, \langle \mathbf{R}, (x', y') \rangle \models \varphi$ for all pairs (x', y')

That is, the operator $[\text{PROF}]$ is a standard box-modality wrt. the relation PROF and $[\text{PAIR}]$ is a standard box-modality wrt. the relation PAIR.

Decidability

Formula φ is *satisfiable* if there are an F and a world w s.t. $F, w \models \varphi$.

The logic discussed here is *decidable*, i.e., there exists an effective algorithm that will decide whether a given formula is satisfiable:

- First, recall that *the frame is fixed*: to even write down a formula, we need to fix the language, which means fixing \mathcal{N} and \mathcal{X} .
- Second, observe that the number of possible SWF's is (huge but) *bounded*: there are exactly $m!(m!^n)$ possibilities.
- Third, observe that *model checking is decidable*: there is an effective algorithm for deciding $F, w \models \varphi$ for given F, w, φ .
- Thus, for a given φ we can “just” try model checking for every possible SWF F and every possible world w .

Of course, this is not a practical algorithm. Ågotnes et al. consider complexity questions in more depth and also provide an axiomatisation.

Modelling: The Pareto Condition

We can model the *Pareto condition* as follows:

$$\text{PARETO} := [\text{PROF}][\text{PAIR}](p_1 \wedge \cdots \wedge p_n \rightarrow \sigma)$$

That is, in every world $\langle \mathbf{R}, (x, y) \rangle$ it must be the case that, whenever all individuals rank $x \succ y$ (i.e., all p_i are true), then also society will rank $x \succ y$ (i.e., σ is true).

Write $F \models \varphi$ if $F, w \models \varphi$ for all worlds w .

We have: $F \models \text{PARETO}$ iff F satisfies the Pareto condition.

Remark: The nesting $[\text{PROF}][\text{PAIR}]$ amounts to a *universal modality* (you can reach every possible world).

Modelling: Independence of Irrelevant Alternatives

Notation: For any coalition $C \subseteq \mathcal{N}$, define p_C as

$$p_C := \bigwedge_{i \in C} p_i \wedge \bigwedge_{i \in \mathcal{N} \setminus C} \neg p_i.$$

We can now express **IIA**:

$$\text{IIA} := [\text{PROF}][\text{PAIR}] \bigwedge_{C \subseteq \mathcal{N}} (p_C \wedge \sigma \rightarrow [\text{PROF}](p_C \rightarrow \sigma))$$

That is, in every world $\langle \mathbf{R}, (x, y) \rangle$ it must be the case that, if exactly the individuals in the group C rank $x \succ y$ (i.e., p_C is true) and society also ranks $x \succ y$ (i.e., σ is true), then for any other profile \mathbf{R}' under which still exactly those in C rank $x \succ y$ society also must rank $x \succ y$.

We have $F \models \text{IIA}$ iff F satisfies IIA.

Modelling: Dictatorships

Finally, we can model what it means for F to be *dictatorial*:

$$\text{DICTATORIAL} := \bigvee_{i \in \mathcal{N}} [\text{PROF}][\text{PAIR}](p_i \leftrightarrow \sigma)$$

That is, there exists an individual i (the dictator) such that it is the case that, to whichever world $\langle \mathbf{R}, (x, y) \rangle$ we move, society will rank $x \succ y$ (i.e., σ will be true) if and only if i ranks $x \succ y$ (i.e., p_i is true).

We have $F \models \neg \text{DICTATORIAL}$ iff F is nondictatorial.

Modelling Arrow's Theorem

Write $\models \varphi$ if $F \models \varphi$ for all SWF's F (for the fixed sets \mathcal{N} and \mathcal{X}).

We are now ready to state *Arrow's Theorem*:

If $|\mathcal{X}| \geq 3$, then $\models \neg(\text{PARETO} \wedge \text{IIA} \wedge \neg\text{DICTATORIAL})$.

Note that this does *not* mean that we have a proof within this logic, although the completeness result of Ågotnes et al. regarding their axiomatisation means that such a proof is feasible in principle.

Remark: To be precise, the above is only a statement of Arrow's Theorem for a fixed (but arbitrary) choice of \mathcal{N} and \mathcal{X} . As none of the formulas involved refer to any $q_{(x,y)}$, this is not a major limitation as far as alternatives are concerned. But wrt. individuals it is a limitation. We cannot state that the theorem holds *for all* finite sets of individuals and we cannot make the restriction to finite electorates explicit.

Approach 2: First-Order Logic

Instead of designing a new logic specifically for our needs, we may ask whether what we want can be expressed in a given standard logic.

Next, we will explore to what extent classical *first-order logic* can be used to model the Arrovian framework of social welfare functions.

Initial considerations:

- FOL is a natural logic to speak about *binary relations*, such as those used to model preference orders.
- Some aspects of the Arrovian framework (e.g., IIA speaking about *all* profiles with particular properties) seem to have a certain *higher-order feel* to them, which *could* be a problem.
- FOL cannot express *finiteness*, which *will* be a problem.

For details on the approach presented next, see the paper cited below.

U. Grandi and U. Endriss. First-Order Logic Formalisation of Impossibility Theorems in Preference Aggregation. *J. Phil. Log.*, 42(4):595-618, 2013.

Language

A key idea is to not talk about profiles (with their internal structure) directly, but to instead introduce the notion of *situation*.

Introduce these *predicate symbols* (with their intuitive meaning):

- $N(z)$: z is an individual
- $X(x)$: x is an alternative
- $S(u)$: u is a situation (referring to a profile)
- $p(z, x, y, u)$: individual z ranks x above y in situation/profile u
- $w(x, y, u)$: society ranks x above y in situation/profile u

Modelling: Social Welfare Functions

We can now write axioms forcing the intended interpretations, e.g.:

- Individual and collective preferences need to be *linear orders*. For instance, p must be interpreted as a *transitive* relation:

$$\forall z. \forall x_1. \forall x_2. \forall x_3. \forall u. [N(z) \wedge X(x_1) \wedge X(x_2) \wedge X(x_3) \wedge S(u) \rightarrow (p(z, x_1, x_2, u) \wedge p(z, x_2, x_3, u) \rightarrow p(z, x_1, x_3, u))]$$

- The predicates N , X and S must *partition* the domain. That is, any object must belong to exactly one of them:

$$\forall x. [N(x) \vee X(x) \vee S(x)] \wedge \forall x. [N(x) \rightarrow \neg X(x) \wedge \neg S(x)] \wedge \dots$$

Together with a few other simple axioms like this, we can ensure that any model satisfying them must correspond to a SWF (see paper).

The only critical issue is to ensure that models are not too small: we need to ensure that the *universal domain* assumption gets respected.

Modelling: Universal Domain Assumption

The universal domain assumption can be modelled, but it's not pretty:

$$\begin{aligned}
& \forall z. \forall x. \forall y. \forall u. [p(z, x, y, u) \rightarrow \exists v. [S(v) \wedge p(z, y, x, v) \wedge \\
& \quad \forall x_1. [p(z, x, x_1, u) \wedge p(z, x_1, y, u) \rightarrow p(z, x_1, x, v) \wedge p(z, y, x_1, v)] \wedge \\
& \quad \forall x_1. [(p(z, x_1, x, u) \rightarrow p(z, x_1, y, v)) \wedge (p(z, y, x_1, u) \rightarrow p(z, x, x_1, v))] \wedge \\
& \quad \forall x_1. \forall y_1. [x_1 \neq x \wedge x_1 \neq y \wedge y_1 \neq y \wedge y_1 \neq x \rightarrow \\
& \quad \quad \quad (p(z, x_1, y_1, u) \leftrightarrow p(z, x_1, y_1, v))] \wedge \\
& \quad \forall z_1. \forall x_1. \forall y_1. [z_1 \neq z \rightarrow (p(z_1, x_1, y_1, u) \leftrightarrow p(z_1, x_1, y_1, v))] \quad]]
\end{aligned}$$

That is, if there exists a situation u in which individual z ranks x above y , then there must exist a situation v where z ranks y above x and everything else remains the same. Once we ensure the existence of at least one situation, this generates a universal domain.

Modelling: Arrow's Axioms

Modelling Arrow's axioms is fairly easy.

The Pareto condition:

$$S(u) \wedge X(x) \wedge X(y) \rightarrow [\forall z.(N(z) \rightarrow p(z, x, y, u)) \rightarrow w(x, y, u)]$$

Independence of irrelevant alternatives (IIA):

$$S(u_1) \wedge S(u_2) \wedge A(x) \wedge A(y) \rightarrow \\ [\forall z.(N(z) \rightarrow (p(z, x, y, u_1) \leftrightarrow p(z, x, y, u_2))) \rightarrow \\ (w(x, y, u_1) \leftrightarrow w(x, y, u_2))]$$

Being nondictatorial:

$$\neg \exists z. N(z) \wedge \forall u. \forall x. \forall y. [S(u) \wedge X(x) \wedge X(y) \wedge p(z, x, y, u) \rightarrow w(x, y, u)]$$

Note: All free variables are understood to be universally quantified.

Modelling: Arrow's Theorem

Let T_{SWF} be the set of axioms defining the theory of SWF's (those shown here and those only given in the paper, including one that ensure that there are ≥ 3 alternatives). Let T_{ARROW} be the union of T_{SWF} and the three axioms on the previous slide.

We are now ready to state Arrow's Theorem:

T_{ARROW} does not have a finite model.

A shortcoming of this approach is that we cannot reduce this to a statement about some formula being a theorem of FOL. Only if we are willing to fix the number n of individuals, then we can do this (easily).

Thus, for fixed n this approach, in principle, permits a proof of Arrow's Theorem in FOL; and given the availability of complete theorem provers for FOL such a proof can, in principle, be found automatically. However, to date no such proof has been realised in practice.

Approach 3: Propositional Logic

For the special case of $n = 2$ and $m = 3$ (or indeed any fixed sizes) we can rewrite the FOL representation in propositional logic:

- predicates $p(z, x, y, u)$ becomes atomic propositions $p_{z,x,y,u}$
- predicates $w(x, y, u)$ become atomic propositions $w_{x,y,u}$
- universal quantifications become conjunctions and existential quantifications become disjunction

That is, we need $2 \cdot 3^2 \cdot (3!)^2 + 3^2 \cdot (3!)^2 = 972$ propositional variables.

Direct rewriting of all axioms into CNF leads to an exponential blowup, but clever rewriting using auxiliary variables leads to a formula with around 35,000 variables and 100,000 clauses (Tang and Lin, 2009).

P. Tang and F. Lin. Computer-aided Proofs of Arrows and other Impossibility Theorems. *Artificial Intelligence*, 173(11):1041–1053, 2009.

Computer-aided Proof of Arrow's Theorem

Tang and Lin (2009) prove two inductive lemmas:

- If there exists an Arrovian SWF for n individuals and $m+1$ alternatives, then there exists one for n and m (if $n \geq 2$, $m \geq 3$).
- If there exists an Arrovian SWF for $n+1$ individuals and m alternatives, then there exists one for n and m (if $n \geq 2$, $m \geq 3$).

That is, Arrow's Theorem holds iff its "*base case*" for 2 individuals and 3 alternatives is true—which we've modelled in propositional logic.

Despite being huge, a modern *SAT solver* can verify the inconsistency of the set of clauses corresponding to $\text{ARROW}(2, 3)$ in < 1 second!

Discussion: Opens up opportunities for quick sanity checks of hypotheses regarding new impossibility theorems.

P. Tang and F. Lin. Computer-aided Proofs of Arrows and other Impossibility Theorems. *Artificial Intelligence*, 173(11):1041–1053, 2009.

Related Work (1)

- A logic based on PDL, applied to the modelling of a standard algorithm from the cake-cutting literature (Parikh, 1985).
- Another PDL-like logic to model negotiation over indivisible goods, preferences, and Pareto efficiency (Endriss and Pacuit, 2006).
- A modal logic that can be used to characterise the majority rule in judgment aggregation (Pauly, 2007).
- A modal logic for social choice functions (Troquard et al., 2011).

R. Parikh. The Logic of Games and its Applications. *Annals of Discrete Mathematics*, 24:111–140, 1985.

U. Endriss and E. Pacuit. Modal Logics of Negotiation and Preference. Proc. JELIA-2006.

M. Pauly. Axiomatizing Collective Judgment Sets in a Minimal Logical Language. *Synthese*, 158(2):233–250, 2007

N. Troquard, W. van der Hoek, and M. Wooldridge. Reasoning about Social Choice Functions. *Journal of Philosophical Logic*, 40(4):473–498, 2011.

Related Work (2)

- Wiedijk (2007) and Nipkow (2009) formalise and verify known *proofs* of Arrow's Theorem using the higher-order logic interactive proof assistants MIZAR and ISABELLE, respectively.
- In the (relatively simple) domain of *ranking sets of objects*, the fully automated derivation of new theorems is possible, using a SAT solver (Geist and Endriss, 2011).

F. Wiedijk. Arrow's Impossibility Theorem. *Formalized Mathematics*, 15(4):171–174, 2007.

T. Nipkow. Social Choice Theory in HOL. *Journal of Automated Reasoning*, 43(3):289–304, 2009.

C. Geist and U. Endriss. Automated Search for Impossibility Theorems in Social Choice Theory: Ranking Sets of Objects. *J. Artif. Intell. Res.*, 40:143–174, 2011.

Summary

We have seen three approaches to *modelling* certain aspects of social choice (here, the classical Arrowian framework) *in logic*, providing different degrees of support for *automated reasoning*:

- modal logic (specifically designed for this job)
- first-order logic (for arbitrary numbers of individuals/alternatives)
- propositional logic (for fixed sets of individuals/alternatives)

We are left with (at least) these questions and challenges:

- don't fix the *set of individuals* (and alternatives) in the language
- model the *universal domain* assumption in an elegant manner
- support *automated reasoning*