

# Towards a Logic of Social Welfare

Thomas Ågotnes<sup>1</sup>  
joint work with Wiebe van der Hoek<sup>2</sup> and Michael  
Wooldridge<sup>2</sup>

<sup>1</sup>Department of Computer Engineering  
Bergen University College, Norway

<sup>2</sup>Department of Computer Science  
University of Liverpool, UK

1<sup>st</sup> International Workshop on Computational Social Choice  
(COMSOC 06)

# Motivation: Formal Reasoning about Social Choice

Social Choice Theory		Formal SCT
Concept	Example	Concept
Social welfare function (SWF)		Model $M$
Possible property of SWFs	Pareto optimality	Formula $\phi$
Fundamental property	Transitivity	Axiom $\phi$
Theorem	Arrow's theorem	Derivable formula $\vdash \phi$
Proof		Formal derivation from axioms

# Preferences and Social Welfare

- $A$ : set of **alternatives**
- **Preference relations**  $L(A)$ : total orders  $R \subseteq A \times A$  (antisymm., trans., refl.).  $R^s$  denotes the irreflexive version.
- **Preference profiles** for  $n$  agents:  $L(A)^n$
- **Social Welfare Function (SWF)**:

$$F : L(A)^n \rightarrow L(A)$$

# Expressing IIA

## Independence of Irrelevant Alternatives (IIA)

$$\forall (R_1, \dots, R_n) \in L(A)^n \forall (S_1, \dots, S_n) \in L(A)^n \forall a \in A \forall b \in A \\ (\forall i \in \Sigma (aR_i b \Leftrightarrow aS_i b)) \Rightarrow (aF(R_1, \dots, R_n)b \Leftrightarrow aF(S_1, \dots, S_n)b)$$

Which constructs would we need in a logical language, in order to be able to express, e.g., IIA? It seems that we need to be able to express (in a single formula):

- Quantification over **alternatives**
- Quantification over preference relations, i.e., over **sets of alternatives**
- Properties of preference relations for **different agents**
- Properties of different preference relations for **the same agent**
- Comparisons of different preference relations
- The preference relation resulting from applying a SWF to other preference relations

# Expressing IIA

## Independence of Irrelevant Alternatives (IIA)

$$\forall (R_1, \dots, R_n) \in L(A)^n \forall (S_1, \dots, S_n) \in L(A)^n \forall a \in A \forall b \in A \\ (\forall i \in \Sigma (aR_i b \Leftrightarrow aS_i b)) \Rightarrow (aF(R_1, \dots, R_n)b \Leftrightarrow aF(S_1, \dots, S_n)b)$$

Which constructs would we need in a logical language, in order to be able to express, e.g., IIA? It seems that we need to be able to express (in a single formula):

- Quantification over **alternatives**
- Quantification over preference relations, i.e., over **sets of alternatives**
- Properties of preference relations for **different agents**
- Properties of different preference relations for **the same agent**
- Comparisons of different preference relations
- The preference relation resulting from applying a SWF to other preference relations

# A Logic of SWFs

$$\phi ::= r \mid r_i \mid \neg\phi \mid \phi \wedge \phi \mid \Box\phi \mid \Box\Box\phi$$

where  $r \in \Pi$  (propositions) and  $i \in \Sigma$  (agents). Define

$$\Diamond\phi \equiv \neg\Box\neg\phi, \quad \Diamond\Box\phi \equiv \neg\Box\neg\Box\phi.$$

Satisfaction: let  $F$  be a SWF,  $\delta : \Pi \rightarrow L(A)^n$  and  $a, b \in A$ :

$$(A, F, \delta, (a, b)) \models r_i \quad \Leftrightarrow \quad (a, b) \in \delta_i(r)$$

$$(A, F, \delta, (a, b)) \models r \quad \Leftrightarrow \quad (a, b) \in F(\delta(r))$$

$$(A, F, \delta, (a, b)) \models \Box\phi \quad \Leftrightarrow \quad \forall_{\delta'} (A, F, \delta', (a, b)) \models \phi$$

$$(A, F, \delta, (a, b)) \models \Box\Box\phi \quad \Leftrightarrow \quad (\forall_{(a' \neq b') \in A \times A} (A, F, \delta, (a', b')) \models \phi)$$

$(A, F) \models \phi$  iff  $(A, F, \delta, (a, b)) \models \phi$  for all  $\delta, (a, b)$ , etc.

# A Logic of SWFs

$$\phi ::= r \mid r_i \mid \neg\phi \mid \phi \wedge \phi \mid \Box\phi \mid \Box\Box\phi$$

where  $r \in \Pi$  (propositions) and  $i \in \Sigma$  (agents). Define

$$\Diamond\phi \equiv \neg\Box\neg\phi, \quad \Diamond\Box\phi \equiv \neg\Box\neg\Box\phi.$$

Satisfaction: let  $F$  be a SWF,  $\delta : \Pi \rightarrow L(A)^n$  and  $a, b \in A$ :

$$(A, F, \delta, (a, b)) \models r_i \quad \Leftrightarrow \quad (a, b) \in \delta_i(r)$$

$$(A, F, \delta, (a, b)) \models r \quad \Leftrightarrow \quad (a, b) \in F(\delta(r))$$

$$(A, F, \delta, (a, b)) \models \Box\phi \quad \Leftrightarrow \quad \forall_{\delta'} (A, F, \delta', (a, b)) \models \phi$$

$$(A, F, \delta, (a, b)) \models \Box\Box\phi \quad \Leftrightarrow \quad (\forall_{(a' \neq b') \in A \times A} (A, F, \delta, (a', b')) \models \phi)$$

$(A, F) \models \phi$  iff  $(A, F, \delta, (a, b)) \models \phi$  for all  $\delta, (a, b)$ , etc.

# Pareto Optimality

## Pareto Optimality (PO)

$$\forall (R_1, \dots, R_n) \in L(A)^n \forall a \in A \forall b \in A ((\forall i \in \Sigma a R_i^s b) \Rightarrow a F(R_1, \dots, R_n)^s b)$$

$$PO = \Box \Box ((r_1 \wedge \dots \wedge r_n) \rightarrow r)$$

## Proposition

$(A, F) \models PO$  iff  $F$  is pareto optimal



# Pareto Optimality

## Pareto Optimality (PO)

$$\forall (R_1, \dots, R_n) \in L(A)^n \forall a \in A \forall b \in A ((\forall i \in \Sigma a R_i^s b) \Rightarrow a F(R_1, \dots, R_n)^s b)$$

$$PO = \Box \Box ((r_1 \wedge \dots \wedge r_n) \rightarrow r)$$

## Proposition

$(A, F) \models PO$  iff  $F$  is pareto optimal

# Non-Dictatorship

## Non-Dictatorship (ND)

$$\neg \exists i \in \Sigma \forall (R_1, \dots, R_n) \in L(A)^n F(R_1, \dots, R_n) = R_i$$

$$ND = \bigwedge_{i \in \Sigma} \Diamond \Diamond \neg (r \leftrightarrow r_i)$$

## Proposition

$(A, F) \models ND$  iff  $F$  does not have a dictator

# Non-Dictatorship

## Non-Dictatorship (ND)

$$\neg \exists i \in \Sigma \forall (R_1, \dots, R_n) \in L(A)^n F(R_1, \dots, R_n) = R_i$$

$$ND = \bigwedge_{i \in \Sigma} \Diamond \Diamond \neg (r \leftrightarrow r_i)$$

## Proposition

$(A, F) \models ND$  iff  $F$  does not have a dictator

# Independence of Irrelevant Alternatives

## Independence of Irrelevant Alternatives (IIA)

$$\forall (R_1, \dots, R_n) \in L(A)^n \forall (S_1, \dots, S_n) \in L(A)^n \forall a \in A \forall b \in A \\ (\forall i \in \Sigma (a R_i b \Leftrightarrow a S_i b)) \Rightarrow (a F(R_1, \dots, R_n) b \Leftrightarrow a F(S_1, \dots, S_n) b)$$

$$IIA = \Box \Box ((\bigwedge_{i \in \Sigma} (r_i \leftrightarrow s_i)) \rightarrow (r \leftrightarrow s))$$

## Proposition

$(A, F) \models IIA$  iff  $F$  has the IIA property

# Independence of Irrelevant Alternatives

## Independence of Irrelevant Alternatives (IIA)

$$\forall (R_1, \dots, R_n) \in L(A)^n \forall (S_1, \dots, S_n) \in L(A)^n \forall a \in A \forall b \in A \\ (\forall i \in \Sigma (a R_i b \Leftrightarrow a S_i b)) \Rightarrow (a F(R_1, \dots, R_n) b \Leftrightarrow a F(S_1, \dots, S_n) b)$$

$$IIA = \Box \Box ((\bigwedge_{i \in \Sigma} (r_i \leftrightarrow s_i)) \rightarrow (r \leftrightarrow s))$$

## Proposition

$(A, F) \models IIA$  iff  $F$  has the IIA property

# Arrow's Theorem

$$MT2 = \Diamond(\Diamond(r_1 \wedge s_1) \wedge \Diamond(r_1 \wedge \neg s_1))$$

## Proposition

$(A, F) \models MT2$  iff  $|A| > 2$

## Theorem (Arrow)

$$\models MT2 \rightarrow \neg(PO \wedge ND \wedge IIA)$$

# Arrow's Theorem

$$MT2 = \Diamond (\Diamond (r_1 \wedge s_1) \wedge \Diamond (r_1 \wedge \neg s_1))$$

## Proposition

$(A, F) \models MT2$  iff  $|A| > 2$

## Theorem (Arrow)

$$\models MT2 \rightarrow \neg (PO \wedge ND \wedge IIA)$$

# Judgment Aggregation

- **Underlying logic**  $\mathbf{L}$  with language  $\mathcal{L}$
- **Agenda**  $\mathcal{A} \subseteq \mathcal{L}$  (closed under single negation)
- **Judgment sets**  $J(\mathcal{A}, \mathbf{L})$ : consistent and complete  $A_i \subseteq \mathcal{A}$
- **Judgment Aggregation Rule (JAR)**  $f$ :  
 $f(A_1, \dots, A_n) \in J(\mathcal{A}, \mathbf{L})$

Interpretation of our language in JARs: let  $\mathcal{A}$  be an agenda,  $f$  be a JAR,  $\delta : \Pi \rightarrow J(\mathcal{A}, \mathbf{L})^n$  and  $p \in \mathcal{A}$ :

$$(\mathcal{A}, f, \delta, p) \models_{\mathbf{L}} r_i \quad \Leftrightarrow \quad p \in \delta_i(r)$$

$$(\mathcal{A}, f, \delta, p) \models_{\mathbf{L}} r \quad \Leftrightarrow \quad p \in f(\delta(r))$$

$$(\mathcal{A}, f, \delta, p) \models_{\mathbf{L}} \Box \phi \quad \Leftrightarrow \quad \forall_{\delta'} (\mathcal{A}, f, \delta', p) \models_{\mathbf{L}} \phi$$

$$(\mathcal{A}, f, \delta, p) \models_{\mathbf{L}} \Box \phi \quad \Leftrightarrow \quad (\forall_{p \in \mathcal{A}} (\mathcal{A}, f, \delta, p) \models_{\mathbf{L}} \phi)$$



# Judgment Aggregation

- **Underlying logic**  $\mathbf{L}$  with language  $\mathcal{L}$
- **Agenda**  $\mathcal{A} \subseteq \mathcal{L}$  (closed under single negation)
- **Judgment sets**  $J(\mathcal{A}, \mathbf{L})$ : consistent and complete  $A_i \subseteq \mathcal{A}$
- **Judgment Aggregation Rule (JAR)**  $f$ :  

$$f(A_1, \dots, A_n) \in J(\mathcal{A}, \mathbf{L})$$

Interpretation of our language in JARs: let  $\mathcal{A}$  be an agenda,  $f$  be a JAR,  $\delta : \Pi \rightarrow J(\mathcal{A}, \mathbf{L})^n$  and  $p \in \mathcal{A}$ :

$$\begin{aligned}
 (\mathcal{A}, f, \delta, p) \models_{\mathbf{L}} r_i &\Leftrightarrow p \in \delta_i(r) \\
 (\mathcal{A}, f, \delta, p) \models_{\mathbf{L}} r &\Leftrightarrow p \in f(\delta(r)) \\
 (\mathcal{A}, f, \delta, p) \models_{\mathbf{L}} \Box \phi &\Leftrightarrow \forall_{\delta'} (\mathcal{A}, f, \delta', p) \models_{\mathbf{L}} \phi \\
 (\mathcal{A}, f, \delta, p) \models_{\mathbf{L}} \Box \phi &\Leftrightarrow (\forall_{p \in \mathcal{A}} (\mathcal{A}, f, \delta, p) \models_{\mathbf{L}} \phi)
 \end{aligned}$$

# Example

Majority voting on a proposition:

$$MV = r \leftrightarrow \bigvee_{G \subseteq \Sigma, |G| > \frac{n}{2}} \bigwedge_{i \in G} r_i$$

The Discursive Dilemma

$\models_{\mathcal{L}} \neg \Box \Box MV$

# Example

Majority voting on a proposition:

$$MV = r \leftrightarrow \bigvee_{G \subseteq \Sigma, |G| > \frac{n}{2}} \bigwedge_{i \in G} r_i$$

## The Discursive Dilemma

$$\models_{\mathbf{L}} \neg \Box \Box MV$$

# In order to achieve completeness, we extend the language

Extend the language with an atom

$\mathbf{h}_p$

for each  $p \in \mathcal{A}$

$$(\mathcal{A}, f, \delta, p) \models_{\mathbf{L}} \mathbf{h}_q \Leftrightarrow p = q$$

# Axiomatisation

Given underlying logic  $\mathbf{L}$ , the logic  $JAL(\mathbf{L})$  is:

$\neg(\mathbf{h}_p \wedge \mathbf{h}_q)$ if $p \neq q$	<i>Atmost</i>		
$\bigvee_{p \in \mathcal{A}} \mathbf{h}_p$	<i>Atleast</i>	all inst. of prop. taut.	Taut
$\diamond \mathbf{h}_p$	<i>Agenda</i>	$\blacksquare(\psi_1 \rightarrow \psi_2) \rightarrow (\blacksquare\psi_1 \rightarrow \blacksquare\psi_2)$	<i>K</i>
$\diamond(\mathbf{h}_p \wedge \varphi) \rightarrow \Box(\mathbf{h}_p \rightarrow \varphi)$	<i>Once</i>	$\blacksquare\psi \rightarrow \psi$	<i>T</i>
$\diamond(\mathbf{h}_p \wedge x) \vee \diamond(\mathbf{h}'_p \wedge x)$	<i>CpJS</i>	$\blacksquare\psi \rightarrow \blacksquare\blacksquare\psi$	4
$(\diamond i \wedge \diamond \neg j) \rightarrow \bigwedge_{o \in O} \diamond o$	<i>C</i>	$\neg \blacksquare\psi \rightarrow \blacksquare \neg \blacksquare\psi$	5
$\Box \Box \psi \leftrightarrow \Box \Box \psi$	<i>(COMM)</i>		

From  $p_1, \dots, p_n \vdash_{\mathbf{L}} q$  infer

$\diamond(\mathbf{h}_{p_1} \wedge x) \wedge \dots \wedge \diamond(\mathbf{h}_{p_n} \wedge x) \rightarrow \Box(\mathbf{h}_q \rightarrow x) \wedge \Box(\mathbf{h}'_q \rightarrow \neg x)$  *Closure*

From  $\varphi \rightarrow \psi$  and  $\varphi$  infer  $\psi$

*MP*

From  $\psi$  infer  $\blacksquare\psi$

*Nec*

where  $\blacksquare \in \{\Box, \Box\}$ ,  $x \in \{r, r_j\}$ ,  $O = \{x_1, \dots, x_k : x_j = (\neg)r_j\}$

## Theorem

$JAL(\mathbf{L})$  is sound and complete wrt. JARs over finite agendas.

# Axiomatisation

Given underlying logic  $\mathbf{L}$ , the logic  $JAL(\mathbf{L})$  is:

$\neg(\mathbf{h}_p \wedge \mathbf{h}_q)$ if $p \neq q$	<i>Atmost</i>		
$\bigvee_{p \in \mathcal{A}} \mathbf{h}_p$	<i>Atleast</i>	all inst. of prop. taut.	Taut
$\diamond \mathbf{h}_p$	<i>Agenda</i>	$\blacksquare(\psi_1 \rightarrow \psi_2) \rightarrow (\blacksquare\psi_1 \rightarrow \blacksquare\psi_2)$	<i>K</i>
$\diamond(\mathbf{h}_p \wedge \varphi) \rightarrow \Box(\mathbf{h}_p \rightarrow \varphi)$	<i>Once</i>	$\blacksquare\psi \rightarrow \psi$	<i>T</i>
$\diamond(\mathbf{h}_p \wedge x) \vee \diamond(\mathbf{h}'_p \wedge x)$	<i>CpJS</i>	$\blacksquare\psi \rightarrow \blacksquare\blacksquare\psi$	4
$(\diamond i \wedge \diamond \neg j) \rightarrow \bigwedge_{o \in O} \diamond o$	<i>C</i>	$\neg \blacksquare\psi \rightarrow \blacksquare \neg \blacksquare\psi$	5
$\Box \Box \psi \leftrightarrow \Box \Box \psi$	<i>(COMM)</i>		

From  $p_1, \dots, p_n \vdash_{\mathbf{L}} q$  infer

$\diamond(\mathbf{h}_{p_1} \wedge x) \wedge \dots \wedge \diamond(\mathbf{h}_{p_n} \wedge x) \rightarrow \Box(\mathbf{h}_q \rightarrow x) \wedge \Box(\mathbf{h}'_q \rightarrow \neg x)$  *Closure*

From  $\varphi \rightarrow \psi$  and  $\varphi$  infer  $\psi$

*MP*

From  $\psi$  infer  $\blacksquare\psi$

*Nec*

where  $\blacksquare \in \{\Box, \Box\}$ ,  $x \in \{r, r_j\}$ ,  $O = \{x_1, \dots, x_k : x_j = (\neg)r_j\}$

## Theorem

$JAL(\mathbf{L})$  is sound and complete wrt. JARs over finite agendas.

# Preference vs. Judgment aggregation

Dietrich and List (2006):

- PA can be embedded in JA
- Given a set of alternatives  $A$ , we can define the underlying logic  $\mathbf{L}^A$  such that **preference relations correspond to judgment sets**

## Corollary

$JAL(\mathbf{L}^A)$  is a sound and complete axiomatisation of SWFs over finite sets of alternatives  $A$ .

# Preference vs. Judgment aggregation

Dietrich and List (2006):

- PA can be embedded in JA
- Given a set of alternatives  $A$ , we can define the underlying logic  $\mathbf{L}^A$  such that **preference relations correspond to judgment sets**

## Corollary

$JAL(\mathbf{L}^A)$  is a sound and complete axiomatisation of SWFs over finite sets of alternatives  $A$ .



# Summary

- Language interpreted in SWFs or JARs
- Syntactically simple, yet expressive - can, e.g., express
  - Rules such as majority voting
  - Properties such as Pareto Optimality
  - Results such as Arrow's theorem, the discursive paradox, Condorcet's paradox
- Sound and complete axiomatisation (finite alternatives/agenda)
- Sheds light on the logical principles of judgment- and preference aggregation
- Sheds light on the differences between the logical principles behind judgement- and preference aggregation