



Group Manipulation in Judgment Aggregation



Sirin Botan Arianna Novaro Ulle Endriss

ILLC – UNIVERSITY OF AMSTERDAM

Judgment aggregation is a formal framework for integrating the views of several agents into a single collective view. This is the first study of strategic behaviour by groups of agents in judgment aggregation. We introduce the concept of **group manipulation** – where a coalition of agents can cooperate to manipulate together – and characterise the family of **aggregation rules** for which group manipulation can be avoided.

Judgment Aggregation

Agenda $\Phi := \Phi^+ \cup \{\neg\varphi \mid \varphi \in \Phi^+\}$

- finite set of formulas of propositional logic
- only non-negated formulas in the *pre-agenda* Φ^+
- **atomic** if Φ^+ only contains atomic propositions

Judgment Set for Φ $J \subseteq \Phi$

- *complete* if $\varphi \in J$ or $\neg\varphi \in J$ for all $\varphi \in \Phi^+$
- *consistent* if it is logically consistent
- $\mathcal{J}(\Phi)$ is the set of complete & consistent judgment sets over Φ

Agents and Profiles

- $\mathcal{N} = \{1, \dots, n\}$ is a finite set of **agents**
- $\mathbf{J} = (J_1, \dots, J_n)$ is a **profile**, vector of *individual* judgment sets
- $N_\varphi^{\mathbf{J}} = \{i \in \mathcal{N} \mid \varphi \in J_i\}$ is the *coalition of supporters* of φ in \mathbf{J}
- (\mathbf{J}_{-i}, J'_i) is a profile like \mathbf{J} , except that J'_i replaced J_i
- \mathbf{J} and \mathbf{J}' are **C-variants**, for $C \subseteq \mathcal{N}$, if $J_i = J'_i$ for all $i \in \mathcal{N} \setminus C$

Flipping

- $J^{\neg\varphi}$ means replacing φ by $\neg\varphi$ or $\neg\varphi$ by φ
- $\mathbf{J}^{\neg S}$ means flipping formulas in S in all judgment sets in \mathbf{J}

Aggregation Rules $F : \mathcal{J}(\Phi)^n \rightarrow 2^\Phi$

- **uniform quota rules** $F_q(\mathbf{J}) = \{\varphi \in \Phi \mid \#N_\varphi^{\mathbf{J}} \geq q\}$ for quota q
 - **nomination rule** if $q = 1$
 - **weak majority rule** if $q = \lceil \frac{n}{2} \rceil$
 - **unanimity rule** if $q = n$

Axioms for Aggregation Rule F

- **independence** $N_\varphi^{\mathbf{J}} = N_\varphi^{\mathbf{J}'}$ implies $\varphi \in F(\mathbf{J}) \Leftrightarrow \varphi \in F(\mathbf{J}')$
- **monotonicity** $\varphi \in J'_i \setminus J_i$ implies $\varphi \in F(\mathbf{J}) \Rightarrow \varphi \in F(\mathbf{J}_{-i}, J'_i)$
- **neutrality** $N_\varphi^{\mathbf{J}} = N_\psi^{\mathbf{J}}$ implies $\varphi \in F(\mathbf{J}) \Leftrightarrow \psi \in F(\mathbf{J})$
- **unbiasedness** $F(\mathbf{J}^{\neg S}) = F(\mathbf{J})^{\neg S}$ for any $\mathbf{J} \in \mathcal{J}(\Phi)^n$ and $S \subseteq \Phi^+$ where $\mathbf{J}^{\neg S} \in \mathcal{J}(\Phi)^n$

Preferences

- J_i is the most preferred judgment set of agent i
- preference ranking in terms of distance to J_i
- **Hamming Distance** $H(J, J') = |J \setminus J'| + |J' \setminus J|$
- *weak order* on judgment sets $J \succsim_i^J J' \Leftrightarrow H(J, J_i) \leq H(J', J_i)$

Example. If agent 3 only cares about the conclusion $(p \wedge q)$ she can manipulate the outcome in her favour by rejecting q .

	p	q	$p \wedge q$
Agent 1	✓	✓	✓
Agent 2	✓	×	×
Agent 3	×	✓	×
PB-Rule	✓	✓	✓

Single-Agent Strategyproofness

A rule is strategyproof if no agent has an incentive to manipulate by reporting an untruthful opinion.

Definition 1. A rule F is **strategyproof**, if for all profiles $\mathbf{J} \in \mathcal{J}(\Phi)^n$, agents $i \in \mathcal{N}$, and judgment sets $J'_i \in \mathcal{J}(\Phi)$ it is the case that $F(\mathbf{J}) \succsim_i^J F(\mathbf{J}_{-i}, J'_i)$.

Some rules, e.g. uniform quota rules, are strategyproof.

Theorem 1. A neutral and unbiased aggregation rule F is single-agent strategyproof **iff** it is both independent and monotonic.

Group Strategyproofness

A rule is group-strategyproof if no coalition of manipulators has an incentive to report untruthful judgments.

Definition 2. A rule F is **group-strategyproof** against coalitions of up to k manipulators, if for all profiles $\mathbf{J} \in \mathcal{J}(\Phi)^n$, coalitions $C \subseteq \mathcal{N}$ with $|C| \leq k$, and C -variants $\mathbf{J}' \in \mathcal{J}(\Phi)^n$ of \mathbf{J} it is the case that $F(\mathbf{J}) \succsim_i^J F(\mathbf{J}')$ for all agents $i \in C$.

Example. If the first three agents form a coalition, they will benefit from flipping their judgments on the indicated formulas.

	φ_1	φ_2	φ_3	$\neg\varphi_1$	$\neg\varphi_2$	$\neg\varphi_3$
Agent 1	⊗	✓	✓	✓	×	×
Agent 2	✓	⊗	✓	×	✓	×
Agent 3	✓	✓	⊗	×	×	✓
Agent 4	×	×	×	✓	✓	✓
Agent 5	×	×	×	✓	✓	✓
Majority	⊗	⊗	⊗	✓	✓	✓

Almost no rule is group-strategyproof.

Theorem 2. Suppose the agenda Φ is atomic. Then a neutral and unbiased aggregation rule F is group-strategyproof against coalitions of up to 3 manipulators **iff** F is independent and monotonic, and if none of the restrictions of F to 3 agents and 3 pre-agenda formulas is either the nomination rule or the unanimity rule.

Uniform quota rules are not group-strategyproof.

Corollary 3. No uniform quota rule F_q with a quota q satisfying $3 \leq q \leq n$ or $1 \leq q \leq n - 2$ that is defined on an atomic agenda Φ is group-strategyproof.

Strategyproofness for Fragile Coalitions

A manipulator may decide to unilaterally opt-out of a manipulation.

Definition 3. A rule F is **group-strategyproof against fragile coalitions** of up to k manipulators, if for all profiles $\mathbf{J} \in \mathcal{J}(\Phi)^n$, coalitions $C \subseteq \mathcal{N}$ with $|C| \leq k$, and C -variants $\mathbf{J}' \in \mathcal{J}(\Phi)^n$ of \mathbf{J} with $F(\mathbf{J}') \succsim_i^J F(\mathbf{J})$ and $F(\mathbf{J}'_{-i}, J_i) \neq F(\mathbf{J}')$ for all $i \in C$ it is the case that $F(\mathbf{J}'_{-i}, J_i) \succsim_i^J F(\mathbf{J}')$ for some $i \in C$.

If agents can opt-out, strategyproof rules are group-strategyproof.

Theorem 4. A neutral and unbiased aggregation rule F is group-strategyproof against fragile coalitions of manipulators **iff** it is independent and monotonic.