Who Reviews The Reviewers? A Multi-Level Jury Problem

Ben Abramowitz¹, Omer Lev², Nicholas Mattei¹

¹Tulane University
²Ben-Gurion University of the Negev
babramow@tulane.edu, omerlev@bgu.ac.il, nsmattei@tulane.edu

Abstract

We consider the problem of determining a binary ground truth using advice from a group of independent reviewers (experts) who express their guess about a ground truth correctly with some independent probability (competence) p_i . In this setting, when all reviewers are competent with $p \geq 0.5$, the Condorcet Jury Theorem tells us that adding more reviewers increases the overall accuracy, and if all p_i 's are known, then there exists an optimal weighting of the reviewers.

However, in practical settings, reviewers may be noisy or incompetent, i.e., $p_i \leq 0.5$, and the number of experts may be small, so the asymptotic Condorcet Jury Theorem is not practically relevant. In such cases we explore appointing one or more chairs (judges) who determine the weight of each reviewer for aggregation, creating multiple levels. However, these chairs may be unable to correctly identify the competence of the reviewers they oversee, and therefore unable to compute the optimal weighting.

We give conditions when a set of chairs is able to weight the reviewers optimally, and depending on the competence distribution of the agents, give results about when it is better to have more chairs or more reviewers. Through numerical simulations we show that in some cases it is better to have more chairs, but in many cases it is better to have more reviewers.

1 Introduction

People have been struggling with finding the *correct* answer for millennia¹. In ancient times, when faced with a problem that required discovering a ground truth, two main approaches dominated. The first, less common today, was to approach deities and either ask them to intervene on the randomness of the world (as in the Book of Joshua, Chapter 7), which is a bit akin to sortition (Flanigan et al. 2021); or to ask the deity's wisdom directly (e.g., the Oracle at Delphi). The second approach, still in widespread use today, is to try to assess the known information and draw a conclusion. This can either be done by laymen, the basic premise of the jury system as established by Magna Carta, or by people with expertise. In both cases, groups of people are used (instead of single individuals) to increase the reliability and accuracy of the answers, building on the "wisdom of the crowds".

Copyright © 2024, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

Mathematical analysis of using a group of agents – a jury, or a set of experts – to assess information and make a decision has been done since at least Condorcet's time, in the late 18th century, when he established the Condorcet Jury Theorem (Condorcet 1785; Dietrich and Spiekermann 2023). In the standard jury setting, agents vote on a binary ground truth and the objective is to aggregate their votes, using a voting rule, to maximize the probability of the outcome being correct. In this setting it is typical to assume that agents guess the ground truth correctly with some independent probability (competence) p_i , we call agents competent when p > 0.5 and incompetent when $p \le 0.5^2$. According to the Condorcet Jury Theorem (CJT), when the agents are competent, the collective accuracy of their majority vote tends to correctness as the number of agents increases. Even with a relatively small number of highly competent agents, accuracy can be very high. However, this result, which basically tells us that groups are less prone to mistake than individuals, rests on a knife's edge. If the agents are even minimally incompetent then, as the population grows, their collective accuracy under majority voting tends to 0, and a small group of highly incompetent agents stands no chance.

In the world around us, this idea is used everywhere – in judicial settings (juries), in academic conferences (peer evaluation), in voting for political leaders or in referendums, and even in settings with inanimate agents, such as aggregating sensor outputs into a single reading or indicator.

The precariousness of the Condorcet Jury Theorem stems from the underlying aggregation procedure, majority voting, being anonymous, thus treating all agents equally, regardless of their competence. When agent competences can be different, majority rule is generally sub-optimal, and if one knows exactly the agents' level of expertise, the optimal aggregation method for maximizing accuracy with any number of independent experts and any competences is to use a weighted majority rule in which each agent's weight is the log-odds of their competence (Shapley and Grofman 1984; Nitzan and Paroush 1982). Somewhat surprisingly, the optimal weight of each agent does not depend on the competences of the other agents or even on the total number of agents. However, the assumption that the competence of each agent is exactly known by others is highly unrealistic.

¹Fans of *The Hitchhiker's Guide to the Galaxy* know it is **42**.

²The term competent is not meant to express a value judgment.

We consider a variant of the classic jury setting, inspired by the domain of academic peer review (Shah 2022), which attempts to address these issues. Since the quality of reviewers may not be known by the conference Program Chairs, many conferences (e.g., AAAI) appoint more senior researchers as SPCs (or chairs) to evaluate the reviewers and decide on how to aggregate their views. Such a multi-level process inspired our model: There are not only reviewers, who we will call *experts*, but also chairs, who we will call *judges*, that evaluate the experts and assign them weights.

Analyzing this setting is particularly interesting when the number of agents is relatively small, and therefore we cannot rely on the asymptotic guarantees of the CJT; as well as when there is a potential for significant deviation among agent competency, even when the particular competence values are unknown; or when agents can be incompetent, i.e., they will make the wrong decision most of the time³. We examine when such a two-level system works well, under what conditions it might be worthwhile to implement it, and when is it better to have an expert become a judge.

Contribution We propose and investigate a model of multi-level jury problems for use when we have a small number of possibly unreliable agents. We show that when we know the agent competences exactly (or even approximately), we can find an optimal aggregation procedure, as long as the judges are competent. When the agents' (experts and judges) competences are unknown, we provide a set of numerical experiments demonstrating that adding more than a single judge is rarely helpful, and, indeed, in some cases, the potential damage of a less competent judge is enough to prefer to avoid judges completely.

2 Related Work

There is a long history of studying the Condorcet jury model and its extensions (e.g., Berend and Paroush (1998); Ben-Yashar and Paroush (2000); Grofman (1978); Feld and Grofman (1984)), which appears in many areas, including computer science, philosophy, and economics. Our investigation is primarily based off the literature on weighting experts in both the offline (Shapley and Grofman 1984; Nitzan and Paroush 1982) and online settings (Cesa-Bianchi et al. 1997; Vovk 1990; Freeman et al. 2020; Berend and Kontorovich 2015), in this work we restrict our focus to a single decision. The overall CJT model can be seen in portfolio solver techniques where slower, more reliable algorithms evaluate ensembles of faster, less accurate algorithms (Thornton et al. 2013) and in boosting techniques from machine learning, where one aggregates weakly competent classifiers into a better overall classifier (Schapire and Freund 2013).

In settings with repeated decisions, the competences of the experts can be estimated based on their voting history. Their competence might be estimated by their similarity to other agents, of how often they agreed with the decision outcomes in the past (Grofman, Owen, and Feld 1983; Baharad et al. 2012; Romeijn and Atkinson 2011). In our setting, however, we do not have access to this history and cannot use it to estimate competences. Indeed, in peer review, one may have a notion of other reviewers' competency, but rarely does one co-review with another to form a precise estimate.

Our emphasis on imperfect judges is also inspired by work on "wisdom of the crowds" and crowdsourcing (Surowiecki 2005; Brabham 2013, 2015), proxy voting (Abramowitz and Mattei 2019; Pivato and Soh 2020) and truth-tracking in Liquid Democracy (Zhang and Grossi 2022; Becker et al. 2021). But, as noted above, a major inspiration has been the work on academic peer evaluation (Shah 2022), in which experts assess each other's competences. Some treat this matrix as a Markov chain, and use its eigenvector values as the experts' weights (Grofman and Feld 1983) in a manner reminiscent of some peer-evaluation models (Page et al. 1999; Walsh 2014; Lev et al. 2023). This contrasts with our setting in which the set of agents who vote and the set of agents who weight the voters are disjoint.

Finally, the problem of partitioning agents into judges and experts is also related to the problem of computing optimal committee sizes (Magdon-Ismail and Xia 2018; Revel, Lin, and Halpern 2021). There has also been attention paid to how group accuracy depends on the size of the group and their mean competence (Grofman 1978; Grofman, Feld, and Owen 1984), which is reflected, in part, in our simulations.

3 Model and Notation

Our model has two types of agents; judges and experts. Let E be a set of m experts and J be a set of n judges. The experts vote on a single binary issue where there is only one correct (ground truth) outcome. Without loss of generality, let the options be represented as $\{1,0\}$ where 1 is correct and 0 is incorrect. Each expert $e \in E$ has a competence, or probability p_e of voting correctly, independent of all other experts. We associate each expert's index with their vote, so expert $e \in E$ casts a vote $v_e \in \{1,0\}$ with competence $p_e = P(v_e = 1)$. If an agent's competence is above $^1\!/^2$ we will say that they are competent, and call them incompetent otherwise. We assume no one is always correct or always incorrect, and so $0 < p_e < 1$.

Weighted Majority Rules. For any competence vector and aggregation rule we refer to the probability of producing the correct outcome as the *accuracy*, and reserve the term *competence* to refer to individual agents' probabilities of voting correctly; i.e. accuracy means collective competence.

Definition 1 (Weighted Majority Rule). A weighted majority rule gives each expert $e \in E$ a weight $w_e \in \mathbb{R}$ and selects 1 as the winner if $\sum_{v_e=1} w_e > \sum_{v_e=0} w_e$, selects 0 as the winner if $\sum_{v_e=1} w_e < \sum_{v_e=0} w_e$, and uses a tie-breaking rule (e.g. coin flip) for the edge case where these sums are equal.

Definition 2 (Simple Majority Rule). Simple majority rule refers to the weighted majority rule where all weights are equal and positive and ties are broken randomly.

The Condorcet Jury Theorem tells us that if $p_e \ge 0.5 + \epsilon$ for some $\epsilon > 0$ for all experts, then with simple majority

³In academic peer evaluation it is uncommon for reviewers to be given negative weights, as is required for the log-odds rule. But it is common in other settings, e.g., sensors or proxy voting scenarios where one might want to always do the opposite of a political rival.

accuracy tends to 1 asymptotically as the number of experts tends to infinity. A weighting function maps vectors of values in (0,1) (i.e. competences) to equal length vectors of real values. For any set of experts, including incompetent ones, the optimal aggregation method of experts' votes, is to apply the log-odds weighting function to the experts' competences and use the corresponding weighted majority rule (Shapley and Grofman 1984; Nitzan and Paroush 1982).

Definition 3 (Log-Odds Weighting Function). Given a vector of values in the open unit interval $\vec{p}=(p_1,\ldots,p_m)$, the log-odds weighting function returns the vector $\vec{w}=(w_1,\ldots,w_m)$ where $w_e=\log(\frac{p_e}{1-p_e})$ for all $1\leq e\leq m$.

Any weighting of the agents implies a collection of winning coalitions - subsets of agents who, if they all vote the same way, determine the outcome regardless of the votes of the remaining agents (Taylor and Zwicker 1992). Different weightings may yield the same rule because they imply the exact same winning coalitions. For example, with 5 agents there are exactly 7 distinct weighted majority rules (Karotkin, Nitzal, and Paroush 1988; Karotkin and Paroush 1994). Multiplying the weights of all agents by a constant does not change the winning coalitions and therefore does not change the rule. Similarly, perturbing agent weights by small amounts may not change the winning coalitions. Therefore, while the weights may vary continuously, the accuracy under various weightings will change in discrete steps. In practice, weights may be finite precision rather than true real numbers, and this is also the case in our simulations that use floating point arithmetic, but as long as the rounding tends to be to small to change winning coalitions for most instances its effect will be negligible.

The log-odds weighting rule assigns a positive weight to competent experts when $p_e > 0.5$, weight of zero if $p_e = 0.5$, and negative weight to any incompetent expert with $p_e < 0.5$. In some settings it may be inappropriate to allow negative weights and better to assume any such weights are rounded up to zero. Bounding weights below by zero has the effect of ignoring the incompetent experts and is therefore qualitatively similar to assuming all experts are competent, though with a smaller number of experts. We therefore focus on the more informative setting where weights can be negative. Negative weights also have real-world motivation. A remote sensor may have drifted so far off its initial calibration to be reliably wrong, as has happened with many spacecraft (Bar-Itzhack and Harman 2003). However, we would like to believe that peer reviewers, jurors, and the like are not so reliably wrong that negative weights would be needed.

Multi-Level Jury Problems. Each judge $j \in J$ estimates the competence of each expert $e \in E$ as $p_{je} \in (0,1)$ and assigns them a weight w_{je} based on this estimate. When assigning weights to the experts, our judges always use the log-odds weighting function on their competence estimates. Intuitively, our model's judges are trying to implement the optimal rule using their estimates of the experts' competences. Formally, judge j assigns expert e a weight $w_{je} = \log(\frac{p_{je}}{1-p_{je}})$. When there are multiple judges, the weight of an expert will be the average weight assigned to

them by the judges $w_e = \frac{1}{n} \sum_j w_{je}$.

Perceived Competences. Our theoretical results depend on judges using the log-odds weighting but do not depend on how the judges form their competence estimates p_{ie} . To perform empirical analysis we must make assumptions about where these estimates come from. Rather than drawing the estimates p_{ie} from some named distribution with mean p_e , we take an approach inspired by peer review, and assume that the judges are fundamentally similar to the experts in the same way chairs are similar to reviewers. Each judge jhas competence p_i just like the experts and estimates competence of expert e as $p_{je} = (p_j \cdot p_e) + (1 - p_j)(1 - p_e)$, i.e., the probability that expert e agrees with them. When p_{je} is derived in this way, we will refer to it as the judge's perceived competence of the expert. As with peer review, a judge may estimate an expert's competency from knowing them professionally, but may not have observed many or any of their past reviews. A judge could also reach this estimate of competency if they observe enough votes from the expert but the ground truth is never revealed, as is the case in some peer prediction settings (Witkowski and Parkes 2012).

Example 1. Suppose we have 5 experts with competences $\vec{p}_E = (0.6, 0.6, 0.6, 0.7, 0.9)$. The optimal log-odds weighting is approximately $\vec{w}_E^* = (0.41, 0.41, 0.41, 0.85, 2.2)$. With these weights the most competent expert $(p_e = 0.9)$ receives a weight $(w_e = 2.2)$ that makes them a dictator in a weighted majority vote, since their weight is greater than all other experts combined. Hence, the accuracy under the log-odds weighting is exactly 0.9. If instead we use simple majority, the accuracy drops to 0.82.

A judge with competence 0.6 would assign the experts weights of approximately $\vec{w}_E^{0.6} = (0.08, 0.08, 0.08, 0.16, 0.323)$ using the log-odds weighting of perceived competences. Note that the fifth expert is no longer a dictator. How high of a competence would the judge need to have to assign weights that results in the optimal weighting? The judge's competence would have to be greater than 0.962; a far cry from 0.6 and higher than all the experts. And yet, the judge's sub-optimal weighting yields an accuracy of 0.898, which is a great improvement over simple majority, and extremely close to optimal at 0.9! Example 1 is illustrated in Figure 1 where we plot the accuracy, sweeping p_j from 0.0 to 1.0.

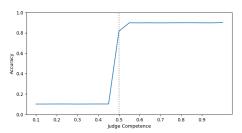


Figure 1: Accuracy of perceived optimal weightings from a single judge with the expert competences as in Example 1.

In Example 1 we rounded the values of the weights to two decimal places, which did not change the rule implemented. Similarly, whether judges are human, sensors, or algorithms, they do not (always) need to provide high precision of weights. More specifically, the smaller the number of agents, the less chance there is for small perturbations or rounding of the weights to change the rule. This is also why Figure 1 will be piecewise linear regardless of the step size we choose for the judge's competence.

4 Optimality and Robustness

With a single judge, if $p_{je}=p_e$ for all $e\in E$, then all experts receive their optimal weight. As noted, small perturbations to the weights do not change the rule because the winning coalitions determined by the weights do not change. Thus, if p_{je} is close enough to p_e for all experts, they will still produce the optimal weighting. We now establish sufficient conditions for an ensemble of judges to produce the optimal weighting, and provide a condition under which the difference between an expert's weight and their optimal weight tends to be small.

Proposition 1 shows that when the geometric mean of the judges' perceived competence of experts odds is their true competence odds, all experts are assigned their optimal weights, $w_e = w_e^*$. This does not require individual judges to know the experts' true competences, and does not depend on the number of experts nor the number of judges.

Proposition 1. If each judge uses the log-odds weighting function on their estimates of expert competences, and the geometric mean of the judges' estimates of each expert's competence odds is the expert's true odds, then the weighted majority rule using the judges' average weights to weight each expert is exactly the optimal weighted majority rule.

Proof. Since judge j gives each expert a weight of $w_{je} = \log(\frac{p_{je}}{1-p_{je}})$,

$$\begin{split} w_e &= \frac{1}{n} \sum_j w_{je} = \frac{1}{n} \sum_j \log(\frac{p_{je}}{1 - p_{je}}) = \\ &= \frac{1}{n} \log(\prod_j \frac{p_{je}}{1 - p_{je}}) = \log((\prod_j \frac{p_{je}}{1 - p_{je}})^{\frac{1}{n}}) \end{split}$$

We assume the geometric mean of judge' estimates of the experts' competence odds is correct, i.e., $(\frac{p_e}{1-p_e})=(\prod_j \frac{p_{je}}{1-p_{je}})^{\frac{1}{n}}.$ Thus, $w_e=\log(\frac{p_e}{1-p_e})=w_e^*.$

This result requires judges' competence estimates that must hold for *all* experts. However, suppose there are errors in these collective competence estimates. We want to know how sensitive the weight of a single expert is to such errors.

Corollary 1. If the geometric mean of judge estimates of competence is off by some multiplicative factor α for some expert, then the error of that expert's weight is only $\log(\alpha)$.

Proof. In the proof above, assume instead that $\alpha\left(\frac{p_e}{1-p_e}\right) = (\prod_j \frac{p_{je}}{1-p_{je}})^{\frac{1}{n}}$. Then $w_e = \log(\alpha \cdot \frac{p_e}{1-p_e}) = w_e^* + \log(\alpha)$. \square

Admittedly, it is not clear in what settings the conditions of Proposition 1 should be expected to hold. Neither can we make claims about what multiplicative factors are realistic in Corollary 1. But ultimately, what we care about most is the sensitive of the accuracy to errors in competence estimates, which has to do with the set of winning coalitions induced by the weights, not the sensitivity of the weights themselves, although the sensitivity of the weights gives some intuition.

Looking back at Example 1, we see that for all $p_j > 0.55$ the accuracy rivals that of the optimal rule, with a nearly imperceptible difference. The effect that dominates Figure 1 is when we move from $p_j < 0.5$ to $p_j = 0.5$ (when the rule becomes simple majority), with another slight bump with a move to $p_j > 0.55$.

Recent work (Baharad, Nitzan, and Segal-Halevi 2022) shows that when expert competences are drawn from certain distributions over the range (1/2,1), simple majority achieves an accuracy close to optimal. However, as one might expect, when experts can be incompetent the majority rule is no longer a good approximation to the optimal weighted majority rule. Thus, if judges can at least differentiate the competent from incompetent experts, the weighting they produce should be expected to outperform simple majority rule when there are incompetent agents. In our model, any minimally competent judge with $p_j > 0.5$ is able to achieve this. We will discuss this more in the next section with an array of experiments with a single judge, but for now we introduce some basic theoretical observations that help us understand the phenomenon.

Example 2 (Two Experts). Suppose there are two experts with competences (p_1, p_2) such that $p_1 > p_2$. If $p_2 > {}^1/2$, i.e., both experts are competent, the optimal aggregation rule is to make p_1 dictator. However, if $p_1 > {}^1/2 > p_2$ then the optimal rule is either to make p_1 a dictator if $p_1 \ge 1 - p_2$, or else make p_2 an anti-dictator using negative weight such that the outcome is the opposite of however p_2 votes. If ${}^1/2 > p_1 > p_2$, then the optimal rule makes p_2 the anti-dictator symmetrically with the first case.

From Example 2, we see that even with only two experts, if a judge can determine which experts are competent, and order their competences correctly, this is enough information to produce the optimal rule. With more experts, the situation is more complicated, but our experiments reveal that merely separating the competent experts from incompetent ones creates a large improvement in overall accuracy. Any chair with $p_j > 1/2$ using log-odds weightings of the perceived competences can achieve this improvement.

Proposition 2 (Correct Sign). If $sign(p_{je} - 0.5) = sign(p_e - 0.5)$, then $sign(w_{je}) = sign(w_{je}^*)$.

Proof. Proposition 2 follows directly from the fact that $\frac{p}{1-p}>1$ if and only if p>1/2, and therefore $\log(\frac{p}{1-p})>0$ if and only if p>1/2. Symmetrically for p<1/2. \qed

Proposition 3 (Correct Order). If p_{je} is a strictly monotonic increasing function of p_e , then the order of expert weights given by judge j is the order of the experts' competences.

Proposition 3 follows from the monotonicity of the logodds weighting. When a single judge applies the log-odds weighting to their perceived competences of a small set of experts then we can make the following observations.

Observation 1. If $p_j = 1/2$ then all experts are equally weighed. If $p_j = 1$ then experts are optimal weighed.

If $p_j=1/2$ then the judge will perceive all experts as having a competence of 1/2, and therefore assign them the weight of 0, which we treat as giving them equal weight. When $p_j=1$, the judge knows exact competences of the experts and therefore assign them their optimal weights.

Observation 2. If
$$p_j > 1/2$$
, then $p_{je} > p_{je'}$ iff $p_e > p_{e'}$.

This means that a judge's perceived competences of the experts preserves the order of their true competences if $p_j > 1/2$. When weights are based on perceived competences, this means whenever $p_j > 0.5$, the judge will assign all experts' weights with the correct sign and in the correct order. This is because when $p_j > 0.5$, p_{je} is monotonically increasing in p_e and $p_{je} > 0.5$ iff $p_e > 0.5$. Thus, even a single barely competent judge might give us an edge over simple majority.

Theorem 1 (Minimal Competent Single Judge). If $p_j > 0.5$ and the judge assigns experts weights according to their perceived competences, the weights given by the judge will have the correct sign and the correct order.

Proof. Let
$$p_j=1/2+\varepsilon_j$$
 and $p_e=1/2+\varepsilon_e$.
$$p_{je}=\left(1/2+\varepsilon_j\right)\left(1/2+\varepsilon_e\right)+\left(1/2-\varepsilon_j\right)\left(1/2-\varepsilon_e\right)$$
$$=1/2+2\varepsilon_j\varepsilon_e$$

If $\varepsilon_j>0$ and $\varepsilon_e>0$, then this value is greater than $^1\!/_2$, if $\varepsilon_j>0$ and $\varepsilon_e<0$, then this value is less than $^1\!/_2$, and if $\varepsilon_e=0$ then this value is exactly $^1\!/_2$. The proposition then follows from Proposition 3 and Proposition 2.

We can generalize Proposition 2 by replacing the requirement that $p_j > 1/2$ with the requirement that the geometric mean of judges' estimated competence odds is greater than 1. Notice that the requirements for this theorem are far weaker than for the optimality demanded by Proposition 1.

Proposition 4 (Correct Sign). If the geometric mean of every expert's estimated competence from the judges is greater than 1 whenever $p_i > 1/2$, less than 1 whenever $p_i < 1/2$, equal to 1 when $p_i = 1/2$, every expert will be assigned a weight with the same sign as their optimal weight.

Proof. Suppose $p_e>1/2$. Their optimal weight is positive, and so we need the following to hold:

$$\sum_{j \in J} \log \left(\frac{p_{je}}{1 - p_{je}} \right) > 0$$

$$\prod_{j \in J} \frac{p_{je}}{1 - p_{je}} > 1$$

$$\left(\prod_{j} \frac{p_{je}}{1 - p_{je}} \right)^{\gamma} > 1$$

for $\gamma>0$. When $\gamma=\frac{1}{n}$ this is the geometric mean. The case for $p_e<1/2$ is symmetric with flipped inequality, and for $p_e=1/2$ is the same but with strict equality. \square

It is straightforward to see that Proposition 3 and Theorem 1 can similarly be generalized to multiple judges, but we omit these here due to space constraints.

5 Single Judge

To better understand the behavior we see in Example 1 and Section 4, we undertake a set of numerical experiments to investigate how the accuracy varies with the competence of a single judge. In our simulations, the experts' competences are drawn i.i.d from various distributions. All experiments were run for 100,000 iterations for each parameterization of the problem instance so that the variances are negligible.

We consider three distributions of expert competences: uniform, truncated normal, and truncated exponential. The uniform distribution reflects settings where the experts can equally have any competence; the exponential distribution models settings where the expertise tends to be rare (Berend and Sapir 2003); and the normal distribution is appropriate for a common expertise, coalescing around a mean value (Taleb 2007).

The top row of Figure 2 shows accuracy as a function of the single judge's competence when expert competences are distributed over the interval [0.001, 0.999] according to the uniform, truncated normal $(\mathcal{N}(^1/2, \text{varying }\sigma))$, and truncated exponential distributions (using the density function $\frac{e^{-x}}{1-e^{-b}}$ for varying values of b) respectively. Only Figures 2a and 2b exhibit true symmetry because competences are drawn from a symmetric distribution with mean 0.5, and Figures 2d and 2e show behavior most similar to Figure 1.

In the top row of Figure 2 we can have highly incompetent experts, but even in this setting whenever the judge has competence $p_j > 1/2$, high overall accuracy is achieved. This is because the ability of the judges to differentiate competent experts from incompetent ones is of primary importance, and Proposition 4 shows that a judge using perceived competences is able to do this.

In Figures 2a-2c, once the judge passes a minimum threshold of competence, little is gained from increasing p_j . Interestingly, in Figure 2b we see that when expert competences are distributed normally with mean $^1/2$, higher variance leads to higher collective accuracy. This appears to be because a judge with sufficiently high competence can differentiate between highly competent and minimally competent experts, and then leverage the benefits of having highly competent experts when they are present.

In the bottom row of Figure 2 we show accuracy as a function of the single judge's competence when expert competences are distributed over the interval [0.501, 0.999] according to the uniform, truncated normal, and truncated exponential distributions respectively. This is closer to prior work in the literature where all experts are assumed to be competent. Unlike in Figure 2a and 2b, which were based on symmetrical distributions, we now see a distinctive asymmetry around $p_j=0.5$. When the judge's competence is 1/2, the judge gives all experts the same weight, so when experts competence is symmetrical around 1/2 (as in the upper row of Figure 2), this results in an accuracy of 1/2, but when they cannot be incompetent, the resulting accuracy is higher – al-

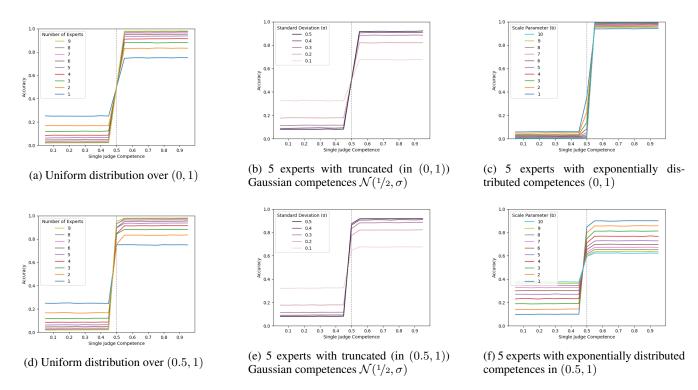


Figure 2: Accuracy with a single judge and expert competences drawn i.i.d. from a distribution with support [0.001, 0.999] (top row) or support [0.501, 0.999] (bottom row).

most optimal (Baharad, Nitzan, and Segal-Halevi 2022). In contrast to the top row of Figure 2, when all experts are competent, there is a large difference in accuracy for truncated exponential distributions with different scale parameters.

6 Should We Add a Judge or an Expert?

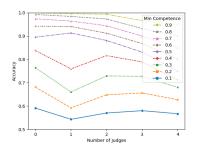
Empirically, with a single judge the accuracy improves as the judge's competence grows, and we know from the Condorcet Jury Theorems that as the number of experts increases, if the experts are competent, accuracy will increase. It follows immediately that if the conditions of Proposition 4 hold, then accuracy will increase as the number of experts increases whether they are competent or incompetent, as long as they have competences that are not equal to 1/2.

We examine the balance between the benefits of increasing the number of judges and increasing the number of experts. That is, with a fixed set of agents of unknown competences, how should they be partitioned between experts and judges? This problem is faced by any scientific conference with a hierarchical structure: how to divide its Program Committee between reviewers and SPCs, ACs, etc.

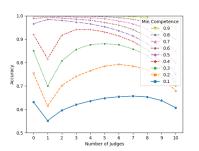
We first draw agents' competences from uniform distributions with varying lower bounds and examine the optimal number of agents to set aside as non-voting judges rather than experts when we have 5 and 11 agents, respectively (Figures 3a and 3d). For both number of agents, we see that setting aside a single agent as judge diminishes the accuracy compared the simple majority rule in almost all cases. This

is more pronounced when there is a possibility the judge will have competence below 1/2, i.e., when a lone judge is incompetent they give all competent experts negative weights and incompetent experts positive weights. The only case where a judge is helpful is when the minimum value of agents is 1/2, perhaps because there is high enough chance that the judge will be helpful, and the agents' competence is not guaranteed to be high enough that losing the judge as an expert is too big a hit. Even adding more judges, at best, returns the accuracy to the level of a simple majority rule, though most commonly it does not. In the 11 agent case, Figure 3d, this effect is even more pronounced than in 5 agent case. In the 11 agent case, adding enough judges can eventually bring peak accuracy to slightly above the simple majority, though it requires roughly an even split between judges and experts (or even slightly more judges).

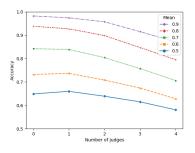
In contrast to the uniform distribution, when drawing competences from the normal and the exponential distributions, things are a bit different. They show that a single judge can be productive. With normal distributions, when the mean is high but not extremely so (Figures 3b and 3e), adding a judge helps. When the mean is very high (0.8 and above), aggregating all agents as experts seems to be better than having a judge, for whom there is still a probability of being bad. But when agents are with a lower mean, having a judge seems to help, and this is true even for a mean of 0.5, in which there is a probability of 0.5 that the judge will be incompetent. This pattern appears for 5 agents, but, as in the



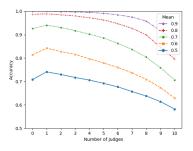
(a) 5 agents with competence from uniform distribution over (min, 1)



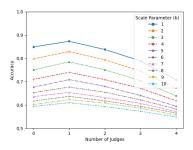
(d) 11 agents with competence from uniform distribution over $(\min,1)$



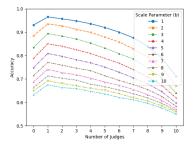
(b) 5 agents with competence from truncated (in (0,1)) Gaussian competences $\mathcal{N}(\text{mean},0.1)$



(e) 11 agents with competence from truncated (in (0,1)) Gaussian competences $\mathcal{N}(\text{mean},0.1)$



(c) 5 agents with competence from truncated (in (0.5,1)) exponential distribution



(f) 11 agents with competence from truncated (in (0.5,1)) exponential distribution

Figure 3: Accuracy from partitioning a set of agents randomly into judges and experts given that all agent competences come i.i.d. from the same distribution.

uniform case, it is more pronounced for 11 agents.

For the exponential distributions (Figures 3c and 3f), this property is stronger – it is *always* beneficial, for our parameters, for agents to have one judge, and that improves over a simple majority. This is likely due to the fact that the small loss of accuracy from having one fewer experts is made up for by the ability of even a minimally competent judge (and all agents are competent in this distribution) to distinguish highly competent experts from less competent ones. Unlike in the uniform case, adding more judges (after a single one) is never helpful compared to a single judge (though sometimes two judges are better than simple majority).

This exploration implies that the division of people in scientific conference is counter-intuitive: the existence of multiple layers above regular experts (e.g., SPC, AC, which are, de facto, multiple judges) does not seem to be helpful. It seems better to have a flat hierarchy (i.e., fewer judges) and use simple majority, despite it being often frowned upon as a strict measure of a paper's quality. We did not investigate the case where a judge is explicitly better than experts, but it is not clear that we are better off in using a judge in such a case, as losing a top expert incurs a cost. Indeed, it is not at all clear that the best judge is the one with the highest competence, and we leave this open question to future research.

7 Discussion and Future Work

We consider a multi-level jury problem in which experts are given weights according to estimates from judges of their competence. We focus on settings where there is a small number of agents, so the classic asymptotic results from the literature do not apply, as well as cases where it is possible for agents to be incompetent (i.e., their chance of being correct might be less than 1/2).

We prove several conditions guaranteeing good outcomes, as well as some which give some minimal guarantees on the quality of the result. Moreover, we show some cases where judges bring a meaningful benefit to the process. However, our results regarding how to divide a group of agents – a particularly relevant issue for scientific conferences – indicate that multiple judges may be unhelpful, and there are cases (e.g., uniform distributions) in which an additional expert is more valuable than a judge.

There are several interesting future directions. One is to reconsider the problem we have presented here when the weights given by experts must all be non-negative, or when it is required for each judge j that $\sum_{e \in E} w_{je} = 1$ (as required in Aziz et al. (2016, 2019) for the setting of peer evaluation). Another is to examine what happens when the hierarchy level is increased by adding an additional layers (as in large conferences, which have Area Chairs in charge of SPCs, in charge of PC members). At what point does it no

longer become helpful (or begin to be helpful)? Can a guarantee of minimal quality of judges change the value proposition of having them?

8 Acknowledgements

Nicholas Mattei was supported by NSF Awards IIS-RI-2007955, IIS-III-2107505, and IIS-RI-2134857, as well as an IBM Faculty Award and a Google Research Scholar Award. This work was supported by the Tulane University Jurist Center for Artificial Intelligence and the Tulane University Center for Community-Engaged Artificial Intelligence. Ben Abramowitz was supported by the NSF under Grant #2127309 to the Computing Research Association for the CIFellows Project. Omer Lev was supported, in part, by NSF-BSF grant #2021659, and by Israel Science Fund (ISF) grants #1965/20 and #3152/20.

References

Abramowitz, B.; and Mattei, N. 2019. Flexible representative democracy: an introduction with binary issues. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, 3–10.

Aziz, H.; Lev, O.; Mattei, N.; Rosenschein, J. S.; and Walsh, T. 2016. Strategyproof Peer Selection: Mechanisms, Analyses, and Experiments. In *Proceedings of the 30th Conference on Artificial Intelligence (AAAI)*, 397–403. Phoenix, Arizona.

Aziz, H.; Lev, O.; Mattei, N.; Rosenschein, J. S.; and Walsh, T. 2019. Strategyproof peer selection using randomization, partitioning, and apportionment. *Artificial Intelligence (AIJ)*, 275: 295–309.

Baharad, E.; Goldberger, J.; Koppel, M.; and Nitzan, S. 2012. Beyond Condorcet: Optimal aggregation rules using voting records. *Theory and decision*, 72(1): 113–130.

Baharad, R.; Nitzan, S.; and Segal-Halevi, E. 2022. One person, one weight: when is weighted voting democratic? *Social Choice and Welfare*, 1–27.

Bar-Itzhack, I.; and Harman, R. 2003. The effect of Sensor Failure on the Attitude and Rate Estimation of the MAP Spacecraft. In *AIAA Guidance, Navigation, and Control Conference and Exhibit*, 5485.

Becker, R.; D'angelo, G.; Delfaraz, E.; and Gilbert, H. 2021. Unveiling the Truth in Liquid Democracy with Misinformed Voters. In *International Conference on Algorithmic Decision Theory*, 132–146. Springer.

Ben-Yashar, R.; and Paroush, J. 2000. A nonasymptotic Condorcet jury theorem. *Social Choice and Welfare*, 17(2): 189–199.

Berend, D.; and Kontorovich, A. 2015. A finite sample analysis of the Naive Bayes classifier. *J. Mach. Learn. Res.*, 16(1): 1519–1545.

Berend, D.; and Paroush, J. 1998. When is Condorcet's jury theorem valid? *Social Choice and Welfare*, 15(4): 481–488.

Berend, D.; and Sapir, L. 2003. Between the expert and majority rules. *Advances in Applied Probability*, 35(4): 941–960.

Brabham, D. C. 2013. *Using crowdsourcing in government*. IBM Center for the Business of Government Washington, DC

Brabham, D. C. 2015. *Crowdsourcing in the public sector*. Georgetown University Press.

Cesa-Bianchi, N.; Freund, Y.; Haussler, D.; Helmbold, D. P.; Schapire, R. E.; and Warmuth, M. K. 1997. How to use expert advice. *Journal of the ACM (JACM)*, 44(3): 427–485.

Condorcet, M. 1785. Essai sur l'application de l'analyse à la probabilité des décisions rendues à la pluralité des voix. De l'Imprimerie Royale (Paris).

Dietrich, F.; and Spiekermann, K. 2023. Jury Theorems. In Zalta, E. N.; and Nodelman, U., eds., *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Spring 2023 edition.

Feld, S. L.; and Grofman, B. 1984. The accuracy of group majority decisions in groups with added members. *Public Choice*, 42(3): 273–285.

Flanigan, B.; Gölz, P.; Gupta, A.; Hennig, B.; and Procaccia, A. D. 2021. Fair algorithms for selecting citizens' assemblies. *Nature*, 596: 548–552.

Freeman, R.; Pennock, D.; Podimata, C.; and Vaughan, J. W. 2020. No-regret and incentive-compatible online learning. In *International Conference on Machine Learning*, 3270–3279. PMLR.

Grofman, B. 1978. Judgmental competence of individuals and groups in a dichotomous choice situation: Is a majority of heads better than one? *Journal of Mathematical Sociology*, 6(1): 47–60.

Grofman, B.; and Feld, S. L. 1983. Determining optimal weights for expert judgment. In *Information Pooling and Group Decision Making: Proceedings of the Second University of California, Irvine, Conference on Political Economy*, 167–72. JAI Press Greenwich, CT.

Grofman, B.; Feld, S. L.; and Owen, G. 1984. Group size and the performance of a composite group majority: Statistical truths and empirical results. *Organizational Behavior and Human Performance*, 33(3): 350–359.

Grofman, B.; Owen, G.; and Feld, S. L. 1983. Thirteen theorems in search of the truth. *Theory and decision*, 15(3): 261–278.

Karotkin, D.; Nitzal, S.; and Paroush, J. 1988. The essential ranking of decision rules in small panels of experts. *Theory and Decision*, 24: 253–268.

Karotkin, D.; and Paroush, J. 1994. Variability of decisional ability and the essential order of decision rules. *Journal of Economic Behavior & Organization*, 23(3): 343–354.

Lev, O.; Mattei, N.; Turrini, P.; and Zhydkov, S. 2023. Peer-Nomination: A novel peer selection algorithm to handle strategic and noisy assessments. *Artificial Intelligence (AIJ)*, 316: 103843.

Magdon-Ismail, M.; and Xia, L. 2018. A mathematical model for optimal decisions in a representative democracy. *Advances in Neural Information Processing Systems*, 31.

- Nitzan, S.; and Paroush, J. 1982. Optimal decision rules in uncertain dichotomous choice situations. *International Economic Review*, 289–297.
- Page, L.; Brin, S.; Motwani, R.; and Winograd, T. 1999. The PageRank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab.
- Pivato, M.; and Soh, A. 2020. Weighted representative democracy. *Journal of Mathematical Economics*, 88: 52–63.
- Revel, M.; Lin, T.; and Halpern, D. 2021. The Optimal Size of an Epistemic Congress. *arXiv preprint arXiv:2107.01042*.
- Romeijn, J.-W.; and Atkinson, D. 2011. Learning juror competence: A generalized Condorcet jury theorem. *Politics, Philosophy & Economics*, 10(3): 237–262.
- Schapire, R. E.; and Freund, Y. 2013. Boosting: Foundations and algorithms. *Kybernetes*, 42(1): 164–166.
- Shah, N. B. 2022. Challenges, experiments, and computational solutions in peer review. *Communications of the ACM*, 65(6): 76–87.
- Shapley, L.; and Grofman, B. 1984. Optimizing group judgmental accuracy in the presence of interdependencies. *Public Choice*, 43(3): 329–343.
- Surowiecki, J. 2005. The wisdom of crowds. Anchor.
- Taleb, N. N. 2007. The Black Swan: The Impact of the Highly Improbable. Random House.
- Taylor, A.; and Zwicker, W. 1992. A characterization of weighted voting. *Proceedings of the American mathematical society*, 115(4): 1089–1094.
- Thornton, C.; Hutter, F.; Hoos, H. H.; and Leyton-Brown, K. 2013. Auto-WEKA: Combined selection and hyperparameter optimization of classification algorithms. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, 847–855.
- Vovk, V. G. 1990. Aggregating strategies. *Proc. of Computational Learning Theory*, 1990.
- Walsh, T. 2014. The PeerRank Method for Peer Assessment. In *Proceedings of the 21st European Conference on Artificial Intelligence (ECAI)*, 909–914. Prague, Czech Republic.
- Witkowski, J.; and Parkes, D. C. 2012. Peer prediction without a common prior. In *Proceedings of the 13th ACM Conference on Electronic Commerce*, 964–981.
- Zhang, Y.; and Grossi, D. 2022. Tracking Truth by Weighting Proxies in Liquid Democracy. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems*, 1482–1490.