

Robust recalibration of aggregate probability forecasts using meta-beliefs*

Cem Peker^{†1} and Tom Wilkening²

¹Divison of Social Science, New York University Abu Dhabi

²Department of Economics, University of Melbourne

August, 2024

Abstract

Previous work suggests that aggregate probabilistic forecasts on a binary event are often conservative. Extremizing transformations that adjust the aggregate forecast away from the uninformed prior of 0.5 can improve calibration in many settings. However, such transformations may be problematic in decision problems where forecasters share a biased prior. In these problems, extremizing transformations can introduce further miscalibration. We develop a two-step algorithm where we first estimate the prior using each forecasters' belief about the average forecast of others. We then transform away from this estimated prior in each forecasting problem. Our algorithm works in single-question forecasting problems and does not require past data. Evidence from experimental prediction tasks suggest that the resulting average probability forecast is robust to biased priors and improves calibration.

Keywords— judgment aggregation, wisdom of crowds, forecasting, extremization, recalibration, meta-beliefs

*We thank the audiences at 2022 INFORMS Annual Meeting and European Decision Sciences Seminar for helpful comments. Cem Peker gratefully acknowledges financial support from the NYUAD Center for Behavioral Institutional Design (C-BID) under the NYUAD Research Institute Award CG005. Tom Wilkening gratefully acknowledges financial support from the Australian Research Council (Future Fellowship Research Grant, FT190100630).

[†]**E-mail addresses:** cem.peker@nyu.edu (C. Peker), tom.wilkening@unimelb.edu.au (T. Wilkening).

1 Introduction

Problems of practical decision-making often require probabilistic forecasts of uncertain events. Knowledge regarding the true likelihood of the event is often scattered across multiple individuals leading to an information aggregation problem where individual forecasts must be combined into a single forecast. Constructing the best aggregation method is difficult because forecasters may make errors when updating their information, may differ in expertise, and may vary in the overlap of the information they have available.

In data-rich environments, it is often possible to use information from training data or other forecasts to better understand the structure of information that exists amongst forecasters. In ideal settings, training data from past forecasts of known outcomes can be used to empirically estimate the diversity of information across individuals and aggregate unknown events accordingly (Atanasov et al., 2017; Breiman, 1996; Dana et al., 2019; Raftery et al., 1997; Satopää, Baron, et al., 2014; Satopää, Jensen, et al., 2014). Alternatively, in cases where forecasters are answering many questions, it may be possible to use answers from many questions to estimate features of the data-generating process that are useful to improving aggregation (Lichtendahl Jr et al., 2022; Satopää et al., 2017).

Unfortunately, decision-makers may not always have access to data that is relevant to the questions of importance. For example, the performance of forecasters on problems with known outcomes may not be relevant to the unknown problem of interest, and collecting relevant data on similar problems may be impractical (Clemen, 1989; Genre et al., 2013). The challenge in these “single-question” forecasting problems is to make the best forecast possible with data related only to the question being asked. We concentrate on the single-question problems for the rest of the paper.

The simple average is a common method to aggregate probability forecasts in the single-question domain (Winkler et al., 2019). Combining independent judgments from many forecasters can lead many individual-specific errors to cancel out leading to improved forecasts via the “wisdom of crowds” effect (Larrick & Soll, 2006; Surowiecki, 2005). However,

28 previous work suggests that the average probability forecast has a major shortcoming: ag-
 29 gregated forecasts tend to be too conservative with the probability of unlikely events being
 30 over-predicted and the probability of near-certain events being under-predicted (Ariely et al.,
 31 2000; Turner et al., 2014). This aggregate conservatism naturally arises in settings where
 32 information is scattered and forecasters have access to different sets of information (Baron
 33 et al., 2014). It also arises even when individual forecasts are well-calibrated since the linear
 34 combination of probability forecasts is always theoretically miscalibrated and lacks sharpness
 35 (Ranjan & Gneiting, 2010).

One way to address the conservative bias is to recalibrate aggregate forecasts using an
 extremization function. Consider the linear log odds (LLO) transformation

$$t(p) = \frac{\delta p^\gamma}{\delta p^\gamma + (1-p)^\gamma}, \quad (1)$$

36 where p and $t(p)$ are the original and transformed probabilities, and $\{\delta, \gamma\}$ are parameters.¹
 37 Extremizing transformations of the LLO form typically improve the accuracy of aggregate
 38 probabilistic forecasts (Atanasov et al., 2017; Budescu et al., 1997; Han & Budescu, 2022).
 39 However, a second potential issue arises in cases where the prior is biased. In many “wicked”
 40 forecasting problems, majority is wrong (Prelec et al., 2017; Wilkening et al., 2022) and/or
 41 inaccurate forecasters express higher confidence (Hertwig, 2012; Koriat, 2008, 2012; Lee &
 42 Lee, 2017). In these cases, average forecast often falls on the wrong side of 0.5. Extremizing
 43 wrong-sided average forecasts using the LLO transformation has the potential of pushing
 44 the forecast away from the true probability and can increase miscalibration rather than
 45 improving accuracy.

¹The LLO transformation follows from a linear log-odds model

$$\log\left(\frac{t(p)}{1-t(p)}\right) = \gamma \log\left(\frac{p}{1-p}\right) + \tau, \quad (2)$$

where γ is the slope and $\tau = \log(\delta)$ gives the intercept (Turner et al., 2014). A simplified implementation sets $\delta = 1$ (Erev et al., 1994; Karmarkar, 1978; Shlomi & Wallsten, 2010), which is shown to improve calibration of the aggregate probability in forecasting geopolitical events (Mellers et al., 2014).

46 In this paper, we ask whether it is possible to estimate the prior in a single-question
47 framework and to use this as the starting point for recalibration. Our main contribution is
48 to show that the common prior can be estimated in the single-question domain by eliciting
49 forecasts and meta-predictions about the forecasts of others. We demonstrate how this
50 information can be used to improve recalibration over standard single-question recalibration
51 methods, and discuss its performance relative to other single-question algorithms that have
52 recently been developed.

53 We consider an environment in which individuals share a common prior that an event
54 may occur, which may be biased.² Forecasters receive independent signals conditional on the
55 actual state, leading to an average probability forecast that puts a higher probability on the
56 actual state than the prior. When the prior that the event occurs is 0.5, the average forecast
57 in these problems always falls on the correct side of 0.5 as the overall crowd size grows large,
58 but the resulting forecast is always conservative. Thus, in these cases, extremization away
59 from 0.5 can improve calibration. However, in a biased decision problem, wrong-sidedness
60 can occur. For example, if the prior is 0.7, there exists cases where the posterior is below 0.7
61 but above 0.5. In these cases, the LLO transformation would extremize the average forecast
62 towards 1, even though the information contained in forecaster’s private signals suggest a
63 lower probability than the prior.

64 To address this issue, we elicit each forecaster’s estimate on the average forecast of others
65 (referred to as their meta-prediction) as well as their probabilistic forecast. We show that
66 these two measures can be combined to estimate the prior in our setting, and then implement
67 an LLO transformation that recalibrates away from the estimated prior rather than using a
68 neutral prior of 0.5.

69 To evaluate how well our robust recalibration algorithm calibrates, we estimate calibration
70 curves across a variety of decision problems related to general knowledge, sports, and

²We are agnostic as to where this bias might come from, but the setup is consistent with one where all forecasters initially observe the same common-signal and then receive a private idiosyncratic one. The common signal leads to the initial prior that differs from 0.5.

71 the price of art works. For recalibration parameters in the range of those suggested in Baron
72 et al. (2014), we find that our algorithm generally improves calibration relative to a variety
73 of alternative algorithms that have been explored in the literature. These include the min-
74 imal pivoting algorithm (Palley & Soll, 2019), the knowledge weighting mechanism (Palley
75 & Satopää, 2023), the meta probability weighting algorithm (Martinie et al., 2020), and
76 the surprising overshoot (SO) algorithm (Peker, 2023). Robust recalibration also generates
77 very low brier scores across decision problems, suggesting that it has very good accuracy
78 characteristics overall.

79 The rest of this paper is organized as follows: Section 2 reviews the recalibration literature
80 and summarizes the other single-question algorithms that we compare our algorithm with.
81 Section 3 introduces the Bayesian framework. Section 4 discusses the existence of wrong-side
82 average forecasts in biased decision problems and develops the robust recalibration method
83 that utilizes meta-predictions. Section 5 provides empirical evidence from experimental
84 prediction tasks. Section 6 provides an overview of our contribution and concludes.

85 **2 Related Literature**

86 Recalibration approaches that seek to account for the partial overlap in shared informa-
87 tion amongst forecasters have been shown in a variety of settings to improve outcomes over
88 techniques that allow only for a weighted average of individual predictions (Baron et al.,
89 2014; Turner et al., 2014). Recalibration typically involves the use of an extremization func-
90 tion, which adjusts forecasts toward extreme outcomes. The most popular choices are logit
91 and probit transformations (Baron et al., 2014; Satopää, Baron, et al., 2014; Satopää et al.,
92 2016; Turner et al., 2014).

93 Recalibration functions are typically symmetric around 0.5. However, as noted in Turner
94 et al. (2014), it is possible and often beneficial to allow for more flexible calibration ap-
95 proaches by extremizing from a different initial prior. A challenge in improving calibration

96 is therefore to incorporate information about the prior into the aggregation algorithm (Di-
97 etrich, 2010; Satopää, 2022). This has been accomplished in multiple-question forecasting
98 environments by using a Bayesian framework and multiple predictions within the same sur-
99 vey to estimate a non-uniform prior across a range of prediction tasks (Lichtendahl Jr et al.,
100 2022; Satopää et al., 2017).³

101 Our approach within the recalibration literature is similar to Lichtendahl Jr et al. (2022),
102 which also stress the importance of using a value other than 0.5 as the basis for extremiza-
103 tion. In their paper, the authors explore data-generating processes in which experts observe
104 multiple independent and identically distributed signals from a joint distribution along with
105 multiple commonly observed private signals. The authors show that with multiple forecasts
106 and historical data, it is possible to develop estimation procedures that are well calibrated
107 and which “antiextremizes” the average in a large number of cases.

108 We see the empirical approach taken in Lichtendahl Jr et al. (2022) as being highly
109 useful in environments where there is substantial historical data to estimate base rates and
110 some confidence in the error structures generated from the data generating process. Our
111 approach, which estimates the prior from meta-predictions and predictions alone, is likely to
112 be more valuable in environments where forecasters have limited historical data and where
113 there is significant uncertainty about the underlying data generating process. We note the
114 two approaches are not mutually exclusive: it is an open and interesting question of how to
115 best combine the two approaches when historical data, training data, and meta-prediction
116 data are available.

117 Our paper also contributes to the emerging literature on forecast aggregation methods
118 that rely on higher order beliefs (Chen et al., 2021; Martinie et al., 2020; Palley & Satopää,
119 2023; Palley & Soll, 2019; Peker, 2023; Prelec et al., 2017; Wilkening et al., 2022). The
120 elicitation of higher-order beliefs allows the researcher additional information about the

³In settings where forecasters have heterogeneous preferences over the extent to which their forecast conforms or contrasts to the reports of others, it may also be possible to estimate the prior using only choice data. See Jia et al. (2024) for an approach to improving forecasts in this alternative setting.

121 signals that individuals receive. Such information can be useful in cases where signals are
122 either correlated or noisy, and where forecasters themselves have more information about
123 the data-generating process than the aggregator.

124 Meta-prediction algorithms have been developed both for binary classification problems
125 (e.g., Chen et al., 2021; Prelec et al., 2017; Wilkening et al., 2022) and in settings like
126 ours where the aggregator wishes to make a probabilistic forecast. In this second class of
127 problems, four main alternative approaches have been proposed: meta-probability weighting,
128 minimal pivoting, knowledge weighting, and the surprising overshoot (SO) algorithm. Meta-
129 probability weighting aims to use forecasters’ meta-prediction as well as their prediction
130 to deal with biased priors or shared information. Forecasters whose prediction and meta-
131 prediction diverge receive higher weights in the subsequent weighted average of predictions
132 (Martinie et al., 2020). Minimal pivoting adjusts the average predictions based on how much
133 it differs from the average meta-prediction (Palley & Soll, 2019). The adjustment corrects for
134 the shared-information bias in the aggregate resulting from forecasters’ common information.
135 Knowledge-weighting proposes a weighted aggregation that seeks to overweight forecasters
136 who are better at predicting the forecasters of their peers (Palley & Satopää, 2023). Finally,
137 the surprising overshoot algorithm corrects for shared information using the observation that
138 the prediction and meta-prediction of an individual should both fall on the same side of a
139 well-calibrated average (Peker, 2023).

140 Our formal framework is similar to Wilkening et al. (2022) and Martinie et al. (2020) in
141 that individuals receive private noisy signals but share a common biased prior. This frame-
142 work naturally introduces conservative forecasts since all individuals have only imperfect
143 information about the true state. Palley and Soll (2019), Palley and Satopää (2023) and
144 Peker (2023) use an alternative framework that allows for intermediate types of shared infor-
145 mation, but places stronger restrictions on the types of signals received. The framework used
146 in knowledge weighting lies between the two approaches and considers an environment where
147 forecasters make noisy predictions and meta-predictions based on their true information.

148 Although it is not emphasized in the previous literature, the framework used in Palley
149 and Soll (2019) is one in which the meta-prediction and prediction correspondences are linear
150 and where the intersection of these lines corresponds to the common prior that exists after
151 accounting for publicly observable information. As a result, the ordering of the prediction
152 and meta-prediction correspondences switch at the uninformative prior. An implication of
153 this is that the minimum pivoting mechanism—which uses the difference in the average pre-
154 diction and meta-prediction to adjust forecasts—is fundamentally an extremizing procedure
155 that adjusts forecasts away from the common prior. As seen in the results section, our algo-
156 rithm with the suggested extremizing parameters of Baron et al. (2014) is more aggressive
157 than the adjustment made in the pivot mechanism and performs better. Thus, at least in
158 the data sets considered, our results suggest that the minimum pivot mechanism is too con-
159 servative. This finding is similar to the contemporaneous work presented in Rilling (2024)
160 that explores a neutral pivoting mechanism that is more aggressive than the original minimal
161 pivot mechanism.

162 Our recalibration procedure relies on a regression approach that is similar to the fit-
163 ting technique used in Palley and Satopää (2023) that seeks to estimate a meta-prediction
164 function using reported predictions and meta-predictions. Regression approaches have also
165 been proposed by Libgober (2023) to identify information regarding the underlying data-
166 generating process.

167 **3 Framework**

168 Our framework is similar to Wilkening et al. (2022) but adapted to the forecasting do-
169 main. We are interested in predicting the probability that a binary event E will occur. The
170 probability that this event occurs varies with a state that is unobservable to the decision
171 maker. However, forecasters receive signals regarding the underlying state and have common
172 knowledge regarding the likelihood of each potential signal in each potential state.

173 We consider a setting where there are two potential underlying states. Let $\omega \in \{\omega_G, \omega_B\}$
174 be the state of the world where G and B represent “Good” and “Bad” states respectively.
175 Event E occurs with probability $Pr(E|\omega_G) = g$ in the good state and with probability
176 $Pr(E|\omega_B) = b$ in the bad state, satisfying $g > b$. Nature determines the state with unknown
177 probability $Pr(\omega = \omega_G)$. Thus, a probability forecast g of E when the state is good and b
178 when the state is bad would be a perfectly well-calibrated forecast.

179 An aggregator elicits and aggregates judgments from a crowd of N forecasters. Forecast-
180 ers share a common prior that the state is good, q , resulting in a common prior belief that
181 the event E will occur with probability $Pr(E|q) = qg + (1 - q)b$.⁴ Each forecaster k receives
182 a signal σ_k from $S \equiv \{s_1, \dots, s_m\} \cup \{s_\emptyset\}$ regarding the underlying state. Without loss of
183 generality, signals are normalized so that $s_i := p(\omega_G|s_i)$, where $p(\omega_G|s_i)$ is forecaster k 's pos-
184 terior belief on the probability of the true state being ω_G when $\sigma_k = s_i$. The uninformative
185 signal satisfies $s_\emptyset := q$ and the signal space is bounded in $[0, 1]$.

186 Let $p(s_i|\omega)$ denote the probability of a signal s_i in state ω , satisfying $\sum_{s_i \in S} p(s_i|\omega) = 1$ for
187 each $\omega \in \{\omega_G, \omega_B\}$. The conditional distribution of signals is represented by a likelihood
188 matrix $[Q_{\omega j}]_{2 \times (m+1)}$. The first and second rows give the likelihoods of each signal in states ω_G
189 and ω_B respectively. Thus, $Q_{\omega_G i} = Q_{1i} \equiv p(s_i|\omega_G)$. We will assume there exists at least one
190 signal $s_l \in \{s_1, \dots, s_m\}$, where $Q_{\omega l} \in (0, 1)$, which implies that at least one signal provides
191 noisy information about the underlying state.⁵ Consistent with our naming convention of
192 states, we also assume $E[\sigma_k|\omega_G] > s_\emptyset > E[\sigma_k|\omega_B]$, which implies that signals are informative
193 and the expected posterior belief is higher in the good state than the bad state.

194 It is useful at this point to note a distinction that we are making regarding events and
195 states. In our framework, the values b and g represent the best prediction that could be made
196 by an aggregator in the corresponding state if he knew the structure of the information service

⁴As can be seen here, there is a one-to-one correspondence between the prior q on ω_G and the prior $qg + (1 - q)b$ on the event E . A similar one-to-one correspondence exists between posteriors on ω_G and E . We will use the words prior and posterior to refer to beliefs over both states and events and will differentiate between them if there is potential ambiguity.

⁵This assumption implies that the signal distribution is non-degenerate in either state since $\sum_j Q_{\omega j} = 1$.

197 and observed an infinite number of draws from it. In some settings, such as asking about
 198 the answer to an objective true/false knowledge question, signals may be fully revealing and
 199 we could set g and b to 1 and 0 respectively. However, in other settings, such as predicting
 200 whether someone will be convicted of a crime, some aspects of the problem (e.g., the detailed
 201 knowledge of the jurists) may be unobservable. In these cases g and b represent the best
 202 possible predictions that could be made about the event based on all possible information
 203 available.

Given a signal s_i such that $p(s_i|\omega_G) + p(s_i|\omega_B) > 0$, the posterior belief that the state is ω_G is given by

$$p(\omega_G|s_i) = \frac{p(\omega_G)p(s_i|\omega_G)}{p(\omega_G)p(s_i|\omega_G) + p(\omega_B)p(s_i|\omega_B)} = s_i.$$

204 Given $p(\omega_G|\sigma_k) = \sigma_k$ for a forecaster with signal σ_k , posterior belief on the occurrence of
 205 event E is given by $Pr(E|\sigma_k) = \sigma_k g + (1 - \sigma_k)b$.

206 The signal densities $\{Q_{G_i}, Q_{B_i}\}$, prior q , and state-conditional event probabilities $\{g, b\}$
 207 are common knowledge to the forecasters but unknown to the aggregator. Each forecaster k
 208 is asked to report i) a *prediction* P_k on the probability of event E and ii) a *meta-prediction*
 209 M_k on the average of others' predictions. Since the likelihood of E depends on the state, a
 210 forecaster's probability prediction is dependent on the forecaster's signal. We will assume
 211 that all forecasters report their best estimate for prediction and meta-prediction, and it is
 212 common knowledge that they do so. Let $P(\sigma_k)$ denote the prediction function of event E ,
 213 where

$$P(\sigma_k) = \sigma_k g + (1 - \sigma_k)b. \tag{3}$$

214 Further, let P_i be the prediction of forecaster i and let $\bar{P}_{-k} = \frac{1}{N-1} \sum_{i \neq k} P_i$ be the average
 215 prediction made by the other $N - 1$ forecasters. Forecaster k 's meta-prediction is given by
 216 $M_k = \mathbb{E}[\bar{P}_{-k}|\sigma_k]$.

For a given outcome state ω , the expected prediction made by a randomly selected other

forecaster is given by

$$\mathbb{E}[P|\omega] \equiv \sum_{s_i \in \mathcal{S}} p(s_i|\omega)[gs_i + b(1 - s_i)].$$

217 Noting that we have assumed that signals are independent once we have conditioned on the
 218 state, $\mathbb{E}[\bar{P}_{-k}|\omega] = \mathbb{E}[P|\omega]$ for all k . Thus, the meta-prediction function, denoted by $M(\sigma_k)$,
 219 can be written as

$$M(\sigma_k) = \sigma_k \mathbb{E}[P|\omega_G] + (1 - \sigma_k) \mathbb{E}[P|\omega_B]. \quad (4)$$

220 Figure 1 plots $P(\sigma_k)$ and $M(\sigma_k)$ in the space of predictions and signals. We note three
 221 general properties that are the basis for our recalibration algorithm. First, both functions
 222 increase linearly in σ_k with the prediction line being more steep than the meta-prediction
 223 line. Note that $P(\sigma_k) \in [b, g]$ and $M(\sigma_k) \in [\mathbb{E}[P|\omega_B], \mathbb{E}[P|\omega_G]]$. We also have $\mathbb{E}[P|\omega_B] > b$
 224 and $\mathbb{E}[P|\omega_G] < g$, i.e. the average prediction will be too conservative in our setting in both
 225 states. To illustrate, consider the case $\omega = \omega_G$ where the true probability of the event is
 226 g . Then, a forecaster k has a perfectly calibrated prediction $P(\sigma_k) = g$ only if $\sigma_k = 1$
 227 and predictions are conservative for all $\sigma_k < 1$. Recall that at least one noisy signal about
 228 the state occurs with strictly positive probability by assumption. Thus, in a large enough
 229 sample, there will always exist forecasters with $\sigma_k < 1$, leading to an average prediction lower
 230 than g . Furthermore, it is common knowledge that forecasters with $\sigma_k < 1$ exist. Forecasters
 231 with $\sigma_k = 1$ expect average prediction to be more conservative than their own prediction,
 232 implying $M(\sigma_k) < P(\sigma_k) = g$ for $\sigma_k = 1$. A similar reasoning holds for $\omega = \omega_B$, resulting in
 233 conservatism in average prediction and a relatively more steep prediction line.

234 Second, the prediction and meta-prediction lines cross exactly once. Figure 1 illustrates
 235 this result. Both functions are monotonically increasing, linear in σ_k , and the range of
 236 meta-predictions is a subset of predictions, resulting in a unique crossing point. Lemma 1
 237 shows that this crossing point occurs at the uninformative prior. All proofs are included in
 238 Appendix A.

239 **Lemma 1.** $M(s_\emptyset) = P(s_\emptyset)$, i.e. a forecaster k 's meta-prediction is equal to her prediction

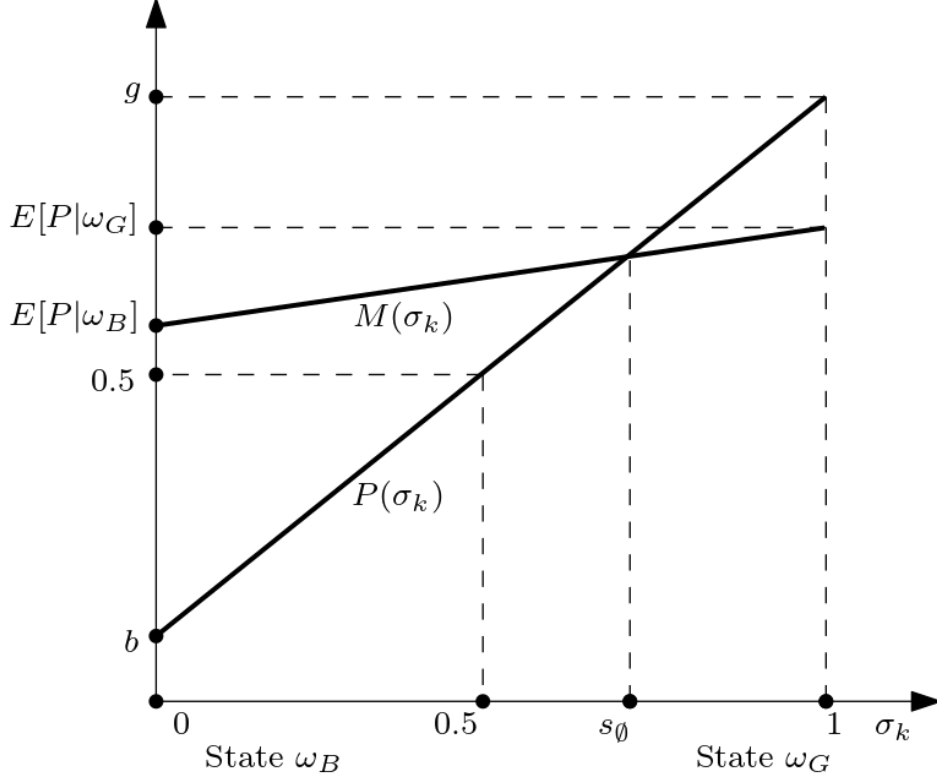


Figure 1: Prediction and meta-prediction functions for a case of $s_\theta > 0.5$. Note that, in this example, the average forecast is higher than 0.5 in both the good and the bad state. Section 4 will explore a potential pitfall in recalibrating such forecasts.

240 *at the prior.*

Finally, since both lines are linear, it is possible to identify $P(s_\theta)$ when there are at least two forecasters with different signals using the crossing point property and a projection. To see this, note that it is possible to rewrite the prediction function as:

$$\sigma_k = \frac{P(\sigma_k) - b}{g - b}.$$

241 Substituting this in Equation 4 yields

$$M(\sigma_k) = \alpha(Q, q, g, b) + \beta(Q, q, g, b)P(\sigma_k), \quad (5)$$

242 where $\alpha(Q, q, g, b) := \frac{g\mathbb{E}[P|\omega_B] - b\mathbb{E}[P|\omega_G]}{g - b}$ and $\beta(Q, q, g, b) := \frac{\mathbb{E}[P|\omega_G] - \mathbb{E}[P|\omega_B]}{g - b}$ are con-

243 starts that do not vary with σ_k . Using any two forecasts and meta-predictions that differ,
244 the terms $\alpha(Q, q, g, b)$ and $\beta(Q, q, g, b)$ can be solved. Prior prediction $P(s_\theta)$ can then be
245 identified by finding the point where $M(s_\theta) = P(s_\theta)$.

246 Before turning to our recalibration strategy, we note that our model presents an ideal
247 environment in which all forecasters perfectly map their signals to predictions and meta-
248 predictions and there are exactly two states. Previous work suggests that the crossing point
249 property between the meta-prediction and prediction correspondence is reasonably robust to
250 systematic individual-level miscalibrations. Wilkening et al. (2022) show that the crossing
251 property holds in decision problems where probability forecasts are miscalibrated as long
252 as miscalibrated forecasts are common knowledge. Chen et al. (2021) show that the same
253 property continues to hold in decision problems where signals are correlated.⁶ Nonetheless,
254 it is likely that there is idiosyncratic noise, particularly in the report of meta-predictions.
255 As seen below, we use regression approaches to estimate the prediction and meta-prediction
256 correspondences in order to help reduce the impact of such noise.

257 In Appendix B, we extend the theoretical discussion and provide two examples that show
258 that the properties of the algorithm are not guaranteed when there are more than two states.
259 The first example shows that the prediction and meta-prediction lines may cross multiple
260 times when we expand the state space and that the estimated prior may not be correct.
261 Nonetheless, the example demonstrates that the algorithm may still function well as long as
262 the estimated prior still identifies the correct direction for extremization.

263 The second example identifies a situation where our algorithm fails to extremize in the
264 correct direction for one of the states. The counter-example highlights a case where signals
265 are very informative about the signals of others but only weakly informative about the
266 underlying likelihood of the event. We see such situations as being quite rare: it requires

⁶Both of these papers explore prediction algorithms that try to correctly predict the correct state rather than make a probabilistic forecast. Wilkening et al. (2022) use the ordering of the average prediction and average meta-prediction to the left and the right of the prior to make predictions. Chen et al. (2021) predict $\mathbb{E}[\bar{P}|\omega]$ in each state using the relationship between predictions and meta predictions and selects the state where the average prediction is closest to the predicted average.

267 a very specific signal structure where the event of interest is only weakly connected to the
268 signals. Nonetheless, the possibility of such cases warrants a careful empirical exploration of
269 the algorithm to assess its applicability in real-world settings.

270 **3.1 A special case with $b = 0$ and $g = 1$**

271 In the empirical section below, we study true-false questions where there is an objectively
272 correct answer. In these questions, it is possible that a very well-informed forecaster could
273 know the state with certainty. Thus, these types of questions might be seen as a special case
274 of our model where $b = 0$ and $g = 1$. In this special case, the prediction correspondence
275 is $P(\sigma_k) = \sigma_k$, and the meta-prediction correspondence is as given by Equation 4 where
276 $\mathbb{E}[P|\omega] \equiv \sum_{s_i \in S} p(s_i|\omega)s_i$ for $\omega \in \{\omega_G, \omega_B\}$. The prediction line is predicted to travel along
277 the 45 degree line in the space of signals and predictions. Thus, the prior corresponds to the
278 point where the meta-prediction correspondence crosses the 45-degree line.

279 In our empirical analyses, we do not directly impose that the prediction line is equal
280 to the 45-degree line since testing this relationship would require information related to
281 signals that are unobservable in empirical data. Instead, we estimate the two parameters in
282 Equation 5 using linear regressions. We then use these estimates to predict the point where
283 $M(s_\theta) = P(s_\theta)$. This approach is valid for any $0 \leq b < g \leq 1$ and therefore nests the special
284 case where $b = 0$ and $g = 1$.

285 The second step of our algorithm involves extremizing the data away from this estimated
286 prior. As discussed below, our algorithm can overshoot the true state when $0 < b < g < 1$
287 but not when $b = 0$ and $g = 1$. We therefore discuss the theoretical properties of the
288 algorithm both for the general case and the special case below.

4 Robust recalibration

As discussed in Section 1, the traditional approach to extremizing compares the average probability prediction $\bar{P} = \frac{1}{N} \sum_{i=1}^N P_i$ to the threshold of 0.5 for determining whether forecasts are extremized towards 0 or 1. This approach can improve forecasts that are too conservative, but problems can arise in some settings where the prior is not 0.5. Figure 1 illustrates the potential problem. The prior is biased towards 1 and the average prediction in the bad state is above 0.5. As seen in Equation 1, the LLO transformation leads to either $t(\bar{P}) > \bar{P} > 0.5$ or $t(\bar{P}) < \bar{P} < 0.5$ for $\bar{P} \neq 0.5$. Figure 1 depicts an example where $E[P|\omega_B] > 0.5$ while $b < 0.5$. Thus, in state ω_B , $t(\bar{P})$ is expected to be even more inaccurate than the original average probability. We refer to such problems as being wrong sided:

Definition 1 (Wrong-sided average prediction). *Average prediction \bar{P} is wrong-sided if i) $\omega = \omega_G$ and $\bar{P} < 0.5 < g$ or, ii) $\omega = \omega_B$ and $\bar{P} > 0.5 > b$.*

Extremization away from 0.5 increases the miscalibration in a wrong-sided average prediction. When can the average prediction be wrong-sided? First, it must be the case that $P(s_\emptyset) \neq 0.5$ for forecasts to be wrong-sided as the sample size grows infinitely large. To see this, note that in a two-state environment, $E[P|\omega_B] < P(s_\emptyset) < E[P|\omega_G]$ and the average prediction will be the expected prediction in each state as the sample grows large. Second, wrong-sidedness can only occur in one of the two states. This follows from the fact that the prior is always between 0 and 1 and the expected posterior is equal to the prior. This implies that on average extremization away from 0.5 can still be beneficial (as found in the literature) but suggests that an algorithm that better identifies cases where wrong-sidedness may occur can improve accuracy.

To account for situations where the average prediction can be wrong-sided, we propose the following **Robust Recalibration** procedure. We first use the data to estimate the prior. Following a similar approach to Palley and Satopää (2023), we allow for random noise ϵ in

314 reported meta-predictions and assume:

$$M_k = \beta_0 + \beta_1 P_k + \epsilon. \quad (6)$$

315 Denoting the estimates $\{\hat{\beta}_0, \hat{\beta}_1\}$, the predicted probability at the prior is found by finding
 316 the probability where the prediction and meta-prediction are equal. This will be given by
 317 $\hat{P}(s_\emptyset) = \hat{\beta}_0 / (1 - \hat{\beta}_1)$ for $\hat{\beta}_1 \neq 1$.

318 Next, using the estimated uninformed prediction $\hat{P}(s_\emptyset)$, we propose a transformation
 319 function $t_{RR}(\bar{P})$ that satisfies the following expression:

$$\log \left(\frac{t_{RR}(\bar{P})}{1 - t_{RR}(\bar{P})} \right) = \log \left(\frac{\bar{P}}{1 - \bar{P}} \right) + \gamma \left[\log \left(\frac{\bar{P}}{1 - \bar{P}} \right) - \log \left(\frac{\hat{P}(s_\emptyset)}{1 - \hat{P}(s_\emptyset)} \right) \right]. \quad (7)$$

320 Equation 7 suggests a linear transformation in log odds where (i) $\bar{P} \geq \hat{P}(s_\emptyset)$ is adjusted
 321 towards 1 and (ii) $\bar{P} < \hat{P}(s_\emptyset)$ is adjusted towards zero 0 when $\gamma \geq 0$. Note that for
 322 $\hat{P}(s_\emptyset) = 0.5$, Equation 7 is the same as Equation 2 with a reparametrization of the slope—
 323 $1 + \gamma$ instead of γ —and an intercept of zero. Thus, in the special case of the estimated prior
 324 being unbiased ($\hat{P}(s_\emptyset) = 0.5$), t_{RR} reduces to the LLO transformation away from 0.5 with
 325 $\delta = 1$, also known as the Karmarkar equation (Karmarkar, 1978).

Solving Equation 7 for $t_{RR}(\bar{P})$, we get

$$t_{RR}(\bar{P}) = \frac{\delta \bar{P}^{1+\gamma}}{\delta \bar{P}^{1+\gamma} + (1 - \bar{P})^{1+\gamma}} \quad (8)$$

326 where $\delta = [(1 - \hat{P}(s_\emptyset)) / \hat{P}(s_\emptyset)]^\gamma$. Unlike simple extremization away from 0.5, $t_{RR}(\bar{P})$ is robust
 327 to wrong-side average predictions. The average is transformed away from $\hat{P}(s_\emptyset)$ instead
 328 of 0.5. If $\hat{P}(s_\emptyset)$ estimates the unknown $P(s_\emptyset)$ accurately, we should expect t_{RR} to adjust
 329 wrong-sided average predictions in the correct direction.

330 Our algorithm essentially uses two pieces of information to transform the average pre-
 331 diction. The first is the estimated common prior which reflects all the commonly shared

332 information in the system. We treat this information as being important to prediction, but
 333 do not recalibrate it as it reflects information that is common across all forecasters. The sec-
 334 ond is the difference between the actual prediction and the common prior. This value reflects
 335 the average change in prediction based on the private signals available to the forecasters. As
 336 these signals are likely to have less overlap, using the average is likely to be conservative.
 337 Thus, by extremizing the difference, we hope to improve the outcome of the estimate.

338 We note an important factor in the relative performance of robust recalibration. The
 339 extent of transformation in robust recalibration depends on both γ and $\hat{P}(s_\emptyset)$, while simple
 340 extremization always uses $\hat{P}(s_\emptyset) = 0.5$. Thus, the step size for adjustment in the two
 341 methods may differ for the same value of γ . Note that for $0 < b < g < 1$, either method
 342 may over-adjust and produce more extreme probabilities than b and g in the corresponding
 343 state. The remainder of this section provides a comparative discussion on the properties of
 344 robust recalibration.

345 In problems that are wrong-sided, simple extremization will adjust predictions away from
 346 the true probability of the event while robust recalibration will adjust predictions in the
 347 direction of the true probability. As mentioned above, the extent of transformation is also
 348 a factor in accuracy. Proposition 1 compares simple extremization and robust recalibration
 349 in wrong-sided problems.

350 **Proposition 1.** *Suppose that the decision problem is wrong-sided. Then, there exists a*
 351 *threshold $g'(b')$ in state $\omega_G(\omega_B)$ such that, for all $g > g'$ ($b < b'$), robust recalibration leads*
 352 *to a lower average Brier score than extremization away from 0.5 for any identical tuning*
 353 *parameter γ in the limit as the sample size goes to infinity. The threshold becomes more*
 354 *extreme (g' to 1, b' to 0) as $|\bar{P} - P(s_\emptyset)|$ increases.*

355 Proposition 1 establishes that robust recalibration achieves higher accuracy in wrong-
 356 sided problems where the true probabilities in the good and bad state are sufficiently extreme.
 357 In other problems, however, there is the potential that robust recalibration “overshoots”
 358 the true probability. To illustrate with a numerical example, suppose $\omega = \omega_G$, $g = 0.55$,

359 $P(s_\theta) = 0.3$, $\bar{P} = 0.49$ and let $\gamma = 1$. Robust-recalibrated probability is 0.68, while simple
 360 extremization leads to 0.48. Robust recalibration transforms in the correct direction, but
 361 the overadjustment produces a less accurate prediction. Such overadjustment becomes more
 362 likely when the adjustment of robust recalibration is larger, which occurs when \bar{P} is further
 363 away from $P(s_\theta)$.

364 The benefits of a proper probability transformation are highest when the true probabil-
 365 ity is close to 0 or 1. These problems represent situations where extremizing in the wrong
 366 direction is very costly in terms of accuracy and where there is little chance that robust re-
 367 calibration overshoots the true probability. A special case where robust recalibration always
 368 improves the Brier score is one where $b = 0$ and $g = 1$. In these problems, it is not possible
 369 to overshoot the true probability through recalibration and a more extreme forecast is better
 370 on average as the sample grows large.

371 **Proposition 2.** *Suppose that the decision problem is wrong-sided, $b = 0$ and $g = 1$. Then*
 372 *robust recalibration leads to a strictly lower average Brier score than extremization away from*
 373 *0.5 for any identical tuning parameter γ in the limit as the sample size goes to infinity.*

374 Proposition 2 follows from the observation that, unlike simple extremization, robust
 375 recalibration transforms wrong-sided average forecasts towards the correct extreme. Since
 376 over-adjustment is not a concern for $b = 0$ and $g = 1$, robust recalibration achieves strictly
 377 higher accuracy.

378 In decision problems where average forecast is not wrong-sided, both robust recalibration
 379 and simple extremization will adjust forecasts in the direction of the true state and therefore
 380 will lead to relatively similar forecasts. However, the intensity of adjustment could differ due
 381 to prior. This may affect the relative accuracy of the two algorithms depending on the extent
 382 to which the average forecast needs to be extremized. To illustrate, consider a simple example
 383 where $\omega = \omega_G$, $g = 0.75$, $\bar{P} = 0.6$, $P(s_\theta) = 0.4$ and $\gamma = 1$. As the sample of forecasters
 384 grow to infinity, robust recalibration recovers $P(s_\theta)$ and transforms according to Equation 8
 385 with $\delta = 1.5$, which leads to a robust-recalibrated probability of 0.77. Simple extremization

386 applies the same transformation with $\delta = 1$ and produces an extremized probability of 0.69.
 387 Since $g = 0.75$, robust recalibration achieves higher accuracy. Now suppose $P(s_\emptyset) = 0.55$
 388 instead. Then, robust-recalibrated probability becomes 0.65 and simple extremization is
 389 more accurate. The opposite result would be true if g is closer to 0.5 and thus, requires a
 390 smaller extremizing adjustment. As a result, we can establish a general result only for the
 391 special case of $b = 0$ and $g = 1$.

392 **Proposition 3.** *Suppose that the decision problem is not wrong-sided, $b = 0$ and $g = 1$.
 393 Then, robust recalibration achieves a lower average Brier score than extremizing away from
 394 0.5 for any identical tuning parameter γ if $|\bar{P} - P(s_\emptyset)| > |\bar{P} - 0.5|$ in the limit as the sample
 395 size goes to infinity.*

396 In Proposition 3, $|\bar{P} - P(s_\emptyset)| > |\bar{P} - 0.5|$ is simply a condition for a larger extrem-
 397 izing adjustment in robust recalibration than simple extremization. Since, extremizing is
 398 always beneficial and over-adjustment is not a concern, the algorithm with a more intensive
 399 extremization achieves higher accuracy.

400 Taking these propositions together, robust recalibration is likely to improve accuracy in
 401 most wrong-sided decision problems. Robust recalibration is strictly preferable in particular
 402 for questions where a binary truth (conditional on the state) exists and extremizing adjust-
 403 ments cannot overshoot the true probability. In problems where the average forecast is not
 404 wrong-sided, relative performance depends on the size of the extremizing adjustment, which
 405 is determined by how the prior prediction compares to 0.5. We may expect similar perfor-
 406 mance to simple extremization when estimated priors are in the vicinity of the uninformative
 407 prior.

408 Before continuing to the empirical section of the paper, it is useful to discuss how we have
 409 set the tuning parameter γ in our empirical analysis. Recall that γ controls the intensity of
 410 extremization away from the estimated prior. As shown in Figure 1, the expected prediction
 411 in states $\{\omega_B, \omega_G\}$ satisfies $b < E[P|\omega_B] < P(s_\emptyset) < E[P|\omega_G] < g$. Perfect calibration is
 412 achieved when extremization away from $P(s_\emptyset)$ is such that the transformed probability is b

413 in state ω_B and g in state ω_G . The optimal value of γ depends on the level of conservatism
414 in the average prediction and informativeness of the prior prediction. To illustrate, suppose
415 the actual state is ω_G . Given $P(s_\emptyset) < E[P|\omega_G] < g$, optimal γ is lower if $P(s_\emptyset)$ is closer to
416 g . In contrast, optimal γ would be higher if the prior is biased towards b .

417 Robust recalibration does not know the optimal value of γ as b and g are unknown, and
418 additional information (such as past data) that may allow estimation of γ is assumed to be
419 unavailable within a single-question aggregation problem. In what follows, we present a wide
420 range of values of γ to investigate how sensitive our approach is to the tuning parameter.
421 Further, when making performance comparisons to other single-question algorithms, we have
422 restricted attention to the tuning parameter range suggested in Baron et al. (2014) and show
423 that our algorithm outperforms the others for both the largest and smallest parameter in
424 this range.

425 Section 5 tests the robust recalibration method $t_{RR}(\bar{P})$ using a variety of experimental
426 data sets. Note that the case of $\hat{P}(s_\emptyset) = 0.5$ (Karmarkar equation) corresponds to the ex-
427 tremizing transformation proposed by Baron et al. (2014). Their LLO extremization can
428 be considered as an implementation of t_{RR} where all decision problems are considered unbi-
429 ased. Thus, we will consider $t_{RR}(\bar{P})$ with $\hat{P}(s_\emptyset) = 0.5$ in all problems as a benchmark that
430 represents “always extremize away from 0.5”. This benchmark allows us to evaluate if the
431 use of meta-predictions to estimate $P(s_\emptyset)$ improves the calibration. The analysis will then
432 compare t_{RR} with various single-question aggregation mechanisms that generate probability
433 forecasts.

434 5 Empirical evidence

435 This section presents empirical evidence for the effectiveness of robust recalibration. We
436 use data from experimental prediction tasks where subjects are asked to report a meta-
437 prediction as well as their prediction. Section 5.1 introduces the data sets. Section 5.2

438 presents preliminary evidence on the existence of wrong-sided average predictions and dis-
439 cusses estimated priors. Section 5.3 offers a comparative analysis on the calibration of
440 transformed probabilities.⁷

441 5.1 Data Sets

442 We investigate the empirical performance of robust recalibration using four distinct types
443 of experimental tasks taken from Wilkening et al. (2022) and Howe et al. (2024). Appendix
444 C provides example questions from each data set.

445 The first set of data consists of simple true/false scientific statements. For each statement,
446 participants report a probabilistic prediction on the statement being true as well as a meta-
447 prediction on the average of other participants’ predictions. Wilkening et al. (2022) collected
448 data from 500 such statements while Howe et al. (2024) replicated the experiment using a
449 subset of these statements. Each implementation recruited a new sample of participants.
450 Thus, we treat each statement-forecasting crowd combination as a distinct forecasting task.
451 The resulting “Science” data set includes 680 tasks in total and the number of participants
452 in a task varied between 79 and 98.

453 The second data set, referred to as “States” data, was also collected by Wilkening et al.
454 (2022). Each task presented a statement on the largest city of a U.S. state being the capital
455 city of the corresponding state. As seen in Prelec et al. (2017), many people erroneously
456 predict that the largest city is highly likely to be the state capital when they do not know
457 the true answer. As such, the dataset is naturally biased towards true. The States data set
458 includes 50 tasks. In each task, a total of 89 subjects reported probabilistic predictions and
459 meta-predictions on the truth of each statement.

460 Howe et al. (2024) collected predictions and meta-predictions on various other domains
461 and we use their questions related to art and NFL trivia. In the “Artwork” data set, subjects
462 saw a picture of a drawing and were asked to predict how likely it is that the market value

⁷Supplemental material includes the datasets and R scripts to reproduce all results (Neuwirth, 2022; R Core Team, 2023; RStudio Team, 2020; Wickham, 2007; Wickham et al., 2019).

463 was more than \$10000. Our data includes 40 decision problems that were repeated in two
464 separate experiments to produce 80 total tasks. The sample size for each task varied between
465 79 and 87 forecasters. The “NFL” domain tasks presented 50 trivia statements about the
466 NFL draft to a US-based subject pool. Similar to the Artwork data, two runs produced 100
467 tasks in total. The sample size per task was either 98 or 99.

468 We note that in two tasks of the Science data, the estimated priors used in the robust
469 recalibration algorithm were outside $(0, 1)$. This can be considered as a failure to estimate
470 $P(s_\emptyset)$ accurately. Appendix D provides the estimated meta-prediction functions and reveals
471 that these were questions where almost all forecasters perfectly predicted the correct answer.
472 Thus, it is likely that these are problems where there is very limited amounts of private
473 information regarding the true state and where idiosyncratic noise in meta-predictions played
474 a large role. We exclude these two science tasks from the results in Section 5.3 and discuss
475 the issue as a potential limitation of our approach in Section 6.⁸

476 Excluding the two science questions, we had a total of 908 tasks in our data.

477 **5.2 Preliminary evidence on priors and wrong-sided average pre-** 478 **dictions**

479 Robust recalibration is expected to improve over simple extremization in transforming
480 wrong-sided average probabilities. Thus, a first step in the analysis is to evaluate the extent
481 to which wrong-sidedness is a problem in the data.

482 As with most practical forecasting problems, we cannot directly observe the correctly
483 calibrated values of g and b in each of our decision problems. Thus, to classify problems as
484 being wrong-sided, we have to make an assumption regarding these values. In this section,
485 we will assume that $b = 0$ and $g = 1$ so that the state corresponds to the true answer. This
486 assumption is based on the fact that the majority of decision problems are questions that

⁸Alternative approaches to dealing with these two observations such as ignoring the bounds on the prior and running the algorithm or using the original prediction do not change the significance of any test in the paper.

487 have an objectively correct answer that could be known by a very well-informed forecaster.
 488 Thus, the true state could potentially be predicted by a forecaster who receives an infinite
 489 number of draws from the potential information system. For $b = 0$ and $g = 1$, Propositions
 490 2 and 3 predict that the robust recalibration algorithm achieves higher accuracy than simple
 491 extremization in wrong-sided problems, while performance could be comparable in others.
 492 Thus, we expect robust recalibration to improve accuracy on average.

493 Figure 2 shows the number of tasks in each data set where the average prediction is
 494 wrong-sided under the above assumption that $b = 0$ and $g = 1$. As seen, the average
 495 prediction is wrong-sided in a considerable number of tasks in each of the data sets. Further,
 496 wrong-sided averages are more common in false statements in all task types, suggesting that
 497 there is a bias towards true in all datasets.

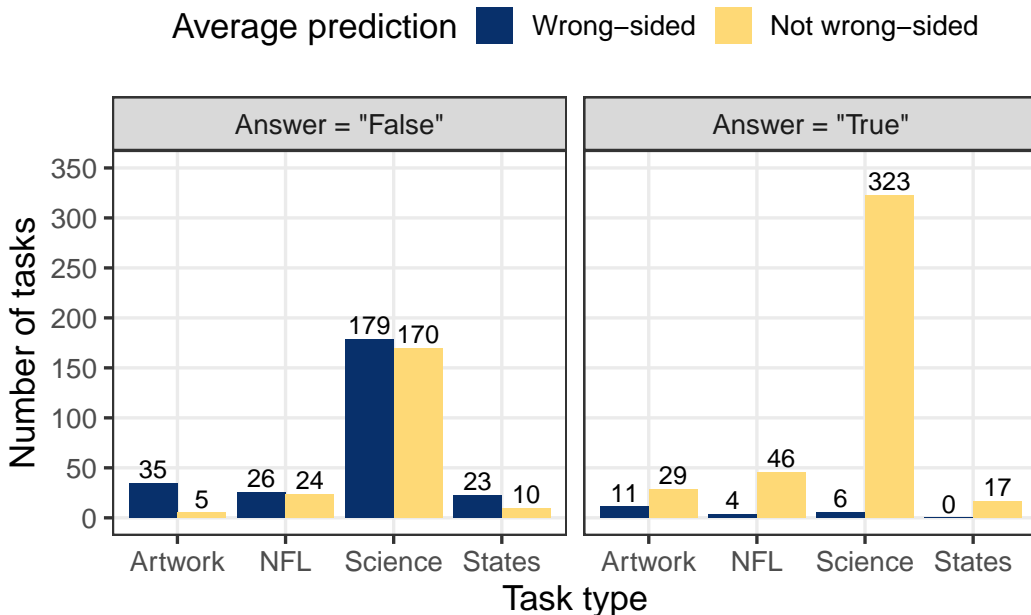


Figure 2: The number of wrong-sided averages in each data set.

498 Figure 3 estimates the prior using the first stage of our robust recalibration procedure and
 499 also supports the conjecture that there is a bias towards true in the data. Estimated priors are
 500 typically higher than 0.5. As such, there are likely to be cases where the robust recalibration
 501 algorithm transforms an average prediction above 0.5 towards 0 while extremization pushes

502 the same average further towards 1.

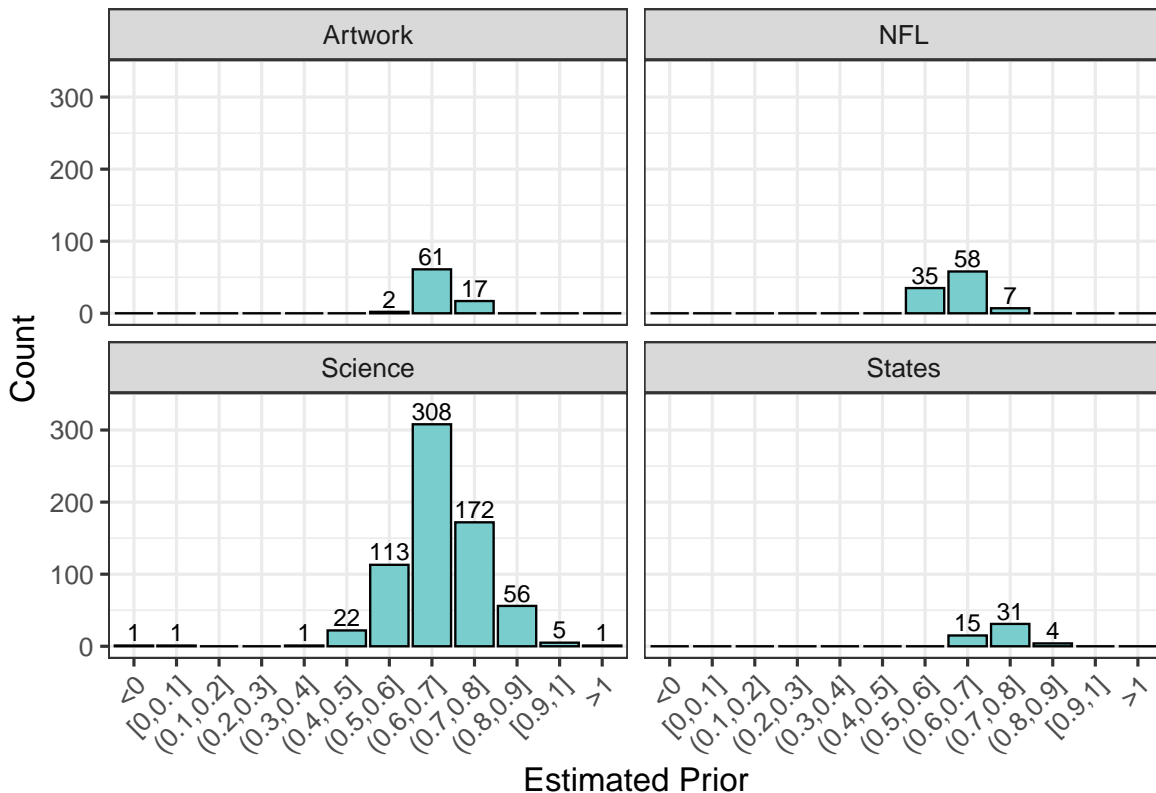


Figure 3: The distribution of estimated priors in each data set.

503 To understand how the estimated priors influence extremization, we also report the num-
 504 ber of decision problems where standard recalibration and robust recalibration procedure
 505 recalibrate forecasts towards and away from the true outcome. Tables 1a and 1b show how
 506 average predictions compare to 0.5 and the estimated priors respectively. Observations along
 507 the diagonal are extremized in the correct direction while observations in the off-diagonal
 508 are adjusted in the wrong direction. As can be seen, there are 263 observations in which
 509 the average prediction is above 0.5 but the correct answer is false. Of these, the robust
 510 recalibration algorithm correctly anti-extremizes 223 observations, while the remaining 40
 511 are still transformed towards 1 as the average prediction is above the estimated prior as well.
 512 There are also 415 observations in which the average prediction is above 0.5 and the correct
 513 answer is true. Of these, the robust recalibration algorithm incorrectly anti-extremizes 146

514 observations and the remaining 269 are correctly transformed towards 1. We evaluate how
 515 these differences in prediction affect accuracy and calibration in the next section.

(a)				(b)			
Correct answer				Correct answer			
	True	False	Total		True	False	Total
$\bar{P} > 0.5$	415	263	678	$\bar{P} > \hat{P}(s_\theta)$	269	40	309
$\bar{P} < 0.5$	21	209	230	$\bar{P} < \hat{P}(s_\theta)$	167	432	599
Total	436	472	908	Total	436	472	908

Table 1: Average prediction vs. 0.5 or estimated prior for “True” and “False” statements

516 5.3 Results

517 This section investigates the accuracy and calibration of the robust-recalibrated proba-
 518 bility forecasts. We run comparative analyses where alternative methods are implemented
 519 as benchmarks. The first analysis compares robust recalibration to the average prediction
 520 and the average extremized away from 0.5. The former is the untransformed simple aver-
 521 age of predictions while the latter transforms the average prediction using Equation 8 with
 522 $\hat{P}(s_\theta) = 0.5$, which corresponds to $\delta = 1$. We consider $\gamma \in \{0.5, 1, 1.5, 2, 2.5, 3\}$ in our
 523 implementations of Equation 8 for both extremization and robust recalibration.

524 Our second analysis compares robust recalibration to various alternative single-question
 525 aggregation algorithms that use meta-predictions to improve accuracy. To make comparisons
 526 here meaningful, we restrict attention to the range of parameters suggested in Baron et al.
 527 (2014) and report results using $\gamma \in \{1.5, 2\}$, which correspond to the suggested lowest and
 528 highest values in our reparametrization. We will consider our algorithm as outperforming
 529 an alternative if it achieves higher accuracy for both values of γ considered.

530 The main text reports the analysis when all 908 tasks are used as the basis of the analysis.
 531 We provide summary statistic tables for the figures provided in the main text in Appendix E.
 532 We also provide an alternative analysis where we compare performance for each of the four

533 prediction tasks separately in Appendix F.

534 **5.3.1 A comparison of robust recalibration to the average prediction and the** 535 **average extremized away from 0.5**

536 Figure 4 shows the distribution of Brier scores of the average prediction, extremized
537 average and robust-recalibrated prediction across all tasks.⁹ Lower scores indicate more
538 accurate forecasts. Each row in the 3×6 grid shows the implementation of extremization away
539 from 0.5 and robust recalibration for various values of γ . We also classify the tasks in terms
540 of how extreme the untransformed average prediction is. Average probability predictions
541 above 0.5 correspond to the confidence for “True”, while for an average probability below
542 0.5, one minus the probability gives the confidence for “False”. The coloring in Figure 4
543 breaks down the distribution of score for five different confidence levels of the corresponding
544 average prediction.

545 Figure 4 demonstrates that extremizing the average prediction away from 0.5 increases
546 the expected accuracy. This result agrees with previous findings on extremization (Han &
547 Budescu, 2022). The robust recalibration procedure offers additional improvements in Brier
548 score over both the average and standard extremization approach for all potential γ param-
549 eters that we explored. As seen in Table 2, the performance difference between extremization
550 and robust recalibration is significant for all values of γ in a paired Wilcoxon sign rank
551 test that treats each decision problem as an observation. Table F1 in Appendix F performs
552 pairwise tests separately for each data set and compares standard extremization to simple
553 average of predictions as well. Robust recalibration achieves substantial and significant im-
554 provement in the Science and States tasks, while the level of accuracy is similar to standard
555 extremization in the Artwork and NFL trivia tasks.

⁹Summary statistics for this analysis is provided in Appendix E. Additional task-level analysis is available in Appendix F.

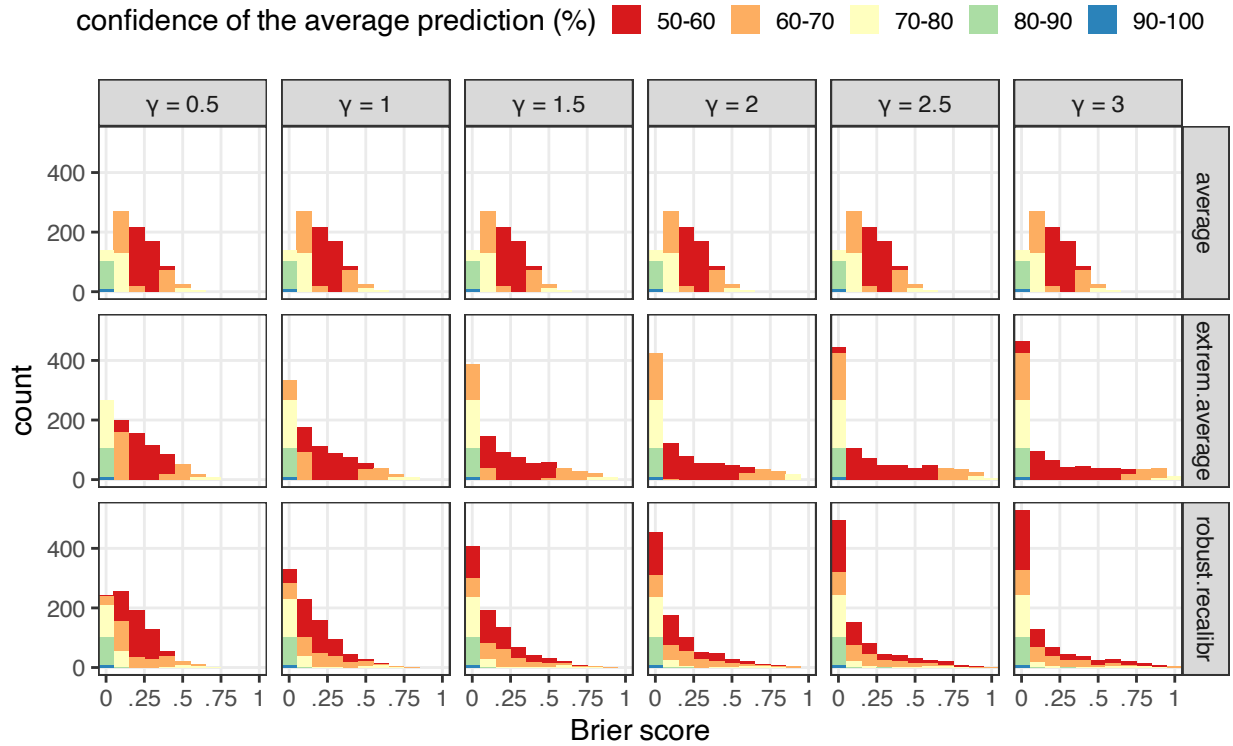


Figure 4: Brier scores of simple average, extremized average and robust-recalibrated probabilities, 908 observations in each panel

γ	Methods		Avg.diff	Med.diff	Test stat.	p-value
0.5	robust.recalibr	extrem.average	-0.0249	-0.0072	V=137,029	<0.0001
1	robust.recalibr	extrem.average	-0.0431	-0.0052	V=143,280	<0.0001
1.5	robust.recalibr	extrem.average	-0.0563	-0.0022	V=148,088	<0.0001
2	robust.recalibr	extrem.average	-0.0658	-0.0008	V=151,761	<0.0001
2.5	robust.recalibr	extrem.average	-0.0728	-0.0003	V=154,699	<0.0001
3	robust.recalibr	extrem.average	-0.0778	-0.0001	V=157,007	<0.0001

Table 2: Two-sided paired Wilcoxon signed rank test of Brier scores, Robust recalibration vs Extremizing away from 0.5. Negative differences indicate higher accuracy for robust recalibration.

556 Figure 4 also suggests that robust recalibration is particularly effective in transforming
 557 low-confidence average predictions. Robust recalibration achieves lower Brier scores when
 558 the corresponding average prediction is 50-60% confident, while extremization away from

559 0.5 leads to higher Brier scores for many such average predictions. Gains in accuracy are
 560 especially strong for larger γ . Figure 5 graphs pairwise difference in Brier scores between
 561 extremization and robust recalibration. In most tasks where robust recalibration achieves
 562 lower Brier scores than simple extremization, the corresponding average prediction is 50-60%
 563 confident.

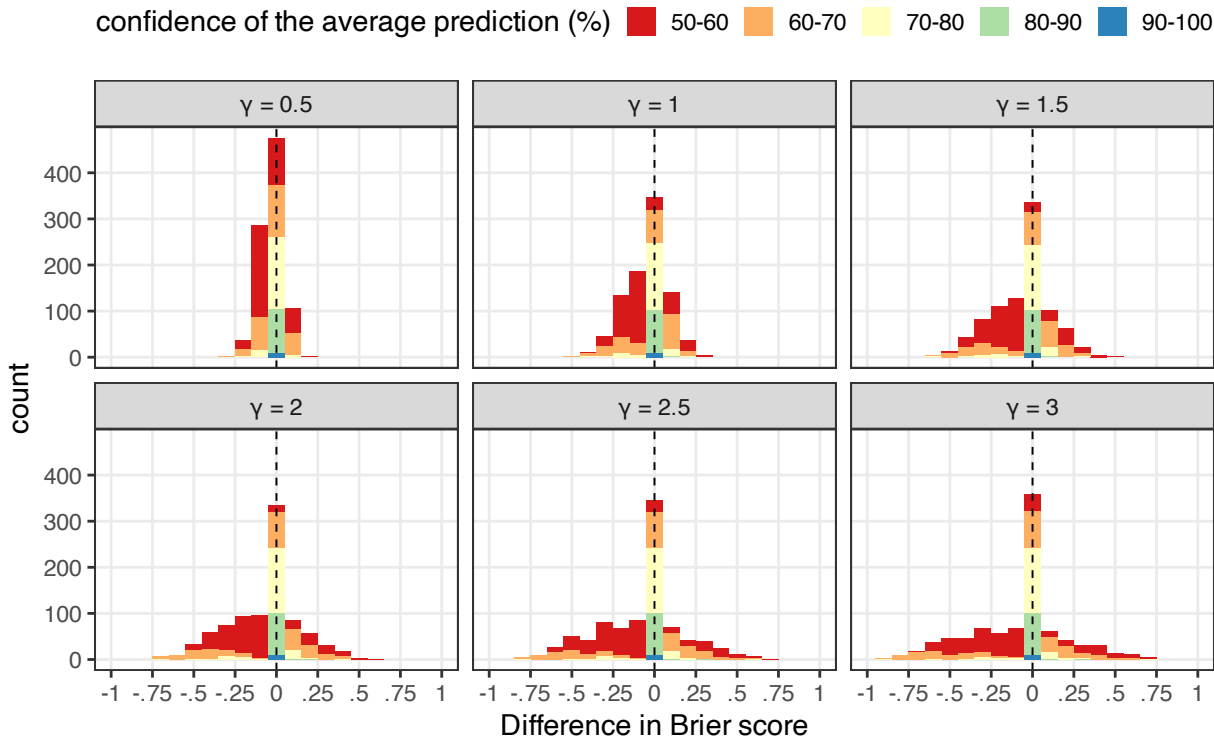


Figure 5: Pairwise differences in Brier score, robust recalibration vs extremized average for $\gamma \in \{0.5, 1, 1.5, 2, 2.5, 3\}$. Negative differences indicate higher accuracy for robust recalibration.

564 Why does robust recalibration make the most difference in low-confidence average pre-
 565 dictions? Table 3 shows the number of wrong-sided average predictions by confidence across
 566 all tasks and reveals that most wrong-sided averages are within the 50-60% confidence cate-
 567 gory. Recall that wrong-sided averages occur mostly in false statements in our experimental
 568 prediction tasks (Table 1) and that estimated priors tend to be above 0.5. As such, simple
 569 extremization wrongly transforms these average prediction into high-confidence true pre-
 570 dictions. Robust recalibration, by contrast, pushes the average prediction away from the

571 estimated prior instead. This anti-extremization produces better Brier scores on average.

	Confidence of the average prediction (%)					Total
	50-60	60-70	70-80	80-90	90-100	
Wrong-sided	182	85	17	0	0	284
Not wrong-sided	198	160	163	94	9	624
Total	380	245	180	94	9	908

Table 3: Number of wrong-sided average predictions by confidence level.

572 As we noted in the previous section, robust recalibration also incorrectly anti-extremizes
 573 some observations that were true and that had an average prediction above 0.5. Such incor-
 574 rect recalibrations hurt accuracy relative to the theoretical optimal, but may or may not affect
 575 the overall calibration of the algorithm depending on the resulting predicted probabilities.
 576 To better understand how well the algorithm calibrates forecasts, we constructed calibration
 577 curves for each method by first separating the data into bins of $\{[0, 0.1], (0.1, 0.2], \dots, (0.9, 1]\}$
 578 based on the predictions of each method. We then plotted the predicted probability of true
 579 in each bin against the actual proportion of problems where true was the correct answer.

580 Figure 6 shows the calibration curves with a separate panel for each γ in the analysis
 581 set. The shaded regions represent the range of proportion true at which the probability
 582 predictions in the corresponding bin are considered well-calibrated. Intuitively, the shaded
 583 regions are analogous to the 45-degree line of perfect calibration.

584 Figure 6 suggests that the transformed probabilities from robust recalibration achieve
 585 better calibration than standard extremization and the average. In particular for $\gamma \geq 1.5$,
 586 robust-recalibrated probabilities on true closely reflect the actual frequency of true in most
 587 bins. In contrast, for extremized averages, the actual proportion of true is typically lower
 588 than the predicted probability in the corresponding bin. In other words, extremized averages
 589 typically overestimate the probability of true. Figures 4 and 6 together imply that the robust
 590 recalibration presents a probability transformation that manages to improve both accuracy
 591 and calibration.

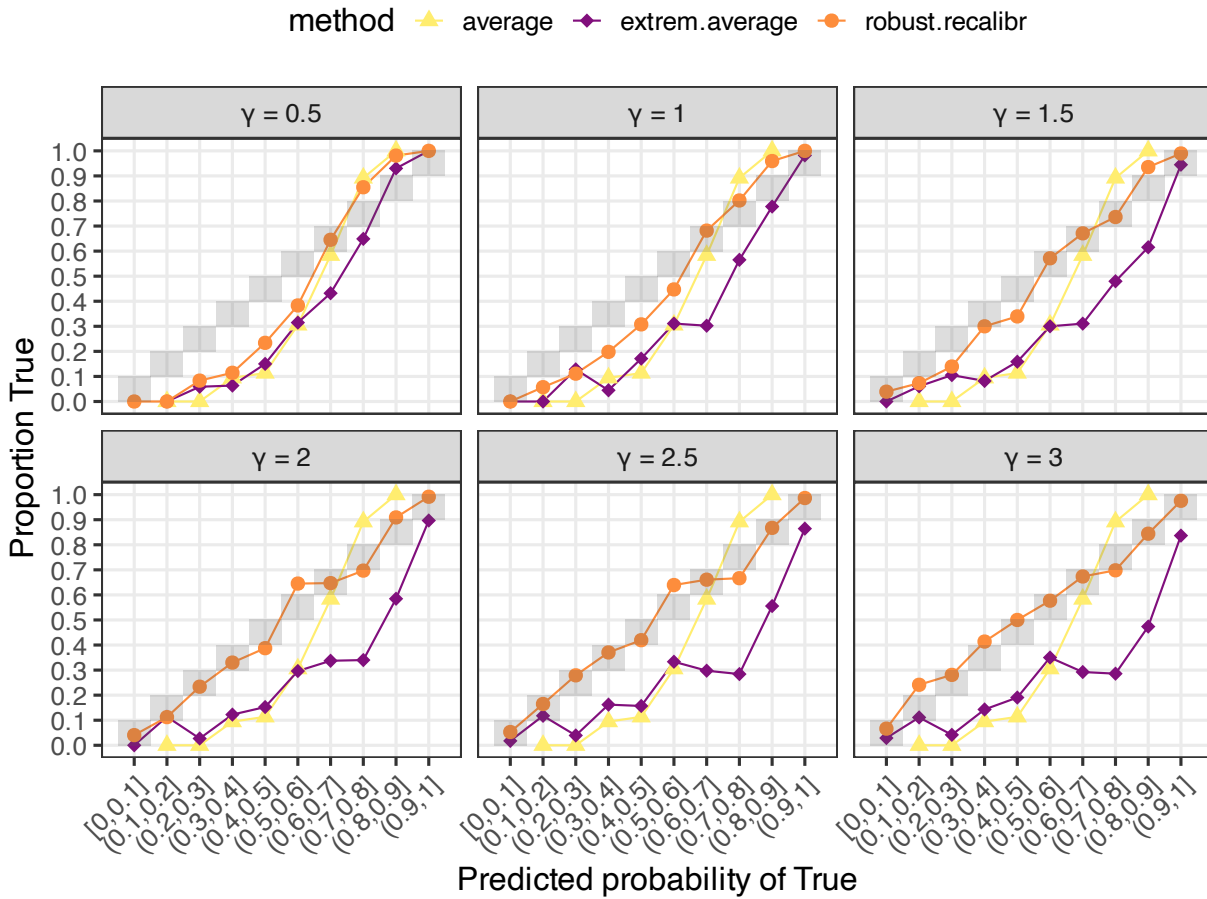


Figure 6: Calibration curves for simple average, extremized average and robust-recalibrated probabilities.

592 **5.3.2 A comparison of robust recalibration to other forecasting algorithms that**
 593 **use meta-predictions**

594 Our analysis thus far compared robust recalibration to methods that do not use meta-
 595 prediction data. One might wonder how it performs against alternative existing methods
 596 that seek to use meta-predictions to produce forecasts. To answer this question, we formed
 597 predictions using a number of alternative algorithms that exist in the literature. We elaborate
 598 on how these algorithms were constructed before continuing on to our second comparative
 599 analysis.

600 We consider four alternative algorithms that seek to exploit meta-predictions to improve
 601 forecasts:

- 602 1. **Meta-probability weighting:** This algorithm constructs a weighted average of prob-
603 abilistic forecasts, where a forecaster’s weight is proportional to the absolute difference
604 between her prediction and meta-prediction (Martinie et al., 2020). Consider the sce-
605 nario where the average forecast is wrong-sided because only a minority of forecasters
606 endorse the correct state. If accurate forecasters anticipate that they are in the mi-
607 nority, we may observe a larger absolute difference between their own forecast and
608 meta-prediction on the average forecast of others. In that case, such forecasters would
609 be weighted more heavily, potentially transforming a wrong-sided forecast correctly in
610 the opposite direction of extremization.
- 611 2. **Knowledge-weighting:** This algorithm, developed in (Palley & Satopää, 2023), seeks
612 to construct optimal weights that minimize the “peer-prediction gap”. This gap mea-
613 sures the difference between a weighted average of forecasters meta-predictions and
614 the actual realization of the average forecast. If forecasters use their information opti-
615 mally in forming meta-predictions, the weights that minimize the peer-prediction gap
616 minimize the error in aggregate forecast as well. Intuitively, if the accurate minority
617 of forecasters are also more accurate in their meta-predictions, knowledge-weighting
618 is expected to put a higher weight on their forecasts, which may transform a wrong-
619 sided average forecast in the correct direction. Knowledge-weighting is applicable in all
620 forms of continuous variables, including non-probabilistic predictions. The knowledge-
621 weighted prediction was outside of $[0, 1]$ in some of our tasks. We winsorize these
622 predictions such that aggregates below 0 (above 1) are set at 0 (1).
- 623 3. **Minimal pivoting:** This algorithm uses meta-prediction data to correct for a poten-
624 tial shared-information bias in the average forecast (Palley & Soll, 2019). Information
625 commonly available to forecasters may bias probabilistic forecasts in a particular direc-
626 tion, which could lead to a wrong-side average forecast. Minimal pivoting adjusts the
627 average forecast according to the difference between average forecast and the average

628 meta-prediction. Meta-predictions are expected to be influenced more heavily by the
629 shared information because forecasters anticipate that their peers will also incorporate
630 it in their forecasts. The pivoting procedure estimates the shared and private informa-
631 tion in the crowd wisdom, and moves the average away from the shared component.
632 Since shared information contains the prior, correction for the shared-information bias
633 is analogous to an extremization away from the prior and it may improve the calibra-
634 tion as well. Similar to the knowledge-weighting algorithm, transformed probabilities
635 that are outside of $[0, 1]$ are winsorized.

636 **4. Surprising Overshoot (SO) algorithm:** This algorithm is another aggregation
637 method that addresses the shared-information problem (Peker, 2023). Information
638 available to a forecaster determines the meta-prediction as well as the prediction, result-
639 ing in a positive correlation between the two. Then, prediction and meta-prediction of
640 an individual should typically fall on the same side of a well-calibrated average predic-
641 tion. As mentioned above, shared information biases meta-predictions more strongly.
642 A significant difference between the percentage of predictions and meta-predictions
643 that overshoot the average prediction would constitute an “overshoot surprise”, which
644 suggests a miscalibration in the average prediction itself. The SO algorithm produces
645 an aggregate forecast that corrects for the shared-information bias using the informa-
646 tion in the size and direction of an overshoot surprise.

647 As can be seen from the description above, the alternative meta-prediction methods do
648 not have a tuning parameter and thus comparing these algorithms to the robust recalibration
649 method with an extremization parameter that is optimized using a subset of the data is not
650 a fair comparison. To avoid this issue, we instead compare methods using the upper and
651 lower bounds of the parameters that are recommended in the literature. Baron et al. (2014)
652 estimated that the optimal parameter value in the standard LLO transformation (Equation 2)
653 for the average forecast is between 2.5 and 3, depending on the expertise of forecasters. In
654 our transformation (Equation 7), this would correspond to $\gamma \in [1.5, 2]$, as we define the

655 tuning parameter as $1 + \gamma$. When making direct comparisons, we report comparisons using
 656 both the lower and upper value in this set and consider the robust recalibration algorithm
 657 as an improvement only if it generates an improvement for both of these bounds.¹⁰

658 Figure 7 presents the frequency distribution of Brier scores for each of the benchmark
 659 algorithms and our robust recalibration method. Panels in the second and third rows show
 660 the results for robust recalibration for each $\gamma \in \{0.5, 1, 1.5, 2, 2.5, 3\}$. Similar to Figure 4, we
 661 color-coded the confidence levels of the average prediction in the corresponding prediction
 662 task to identify potential patterns over types of decision problems.

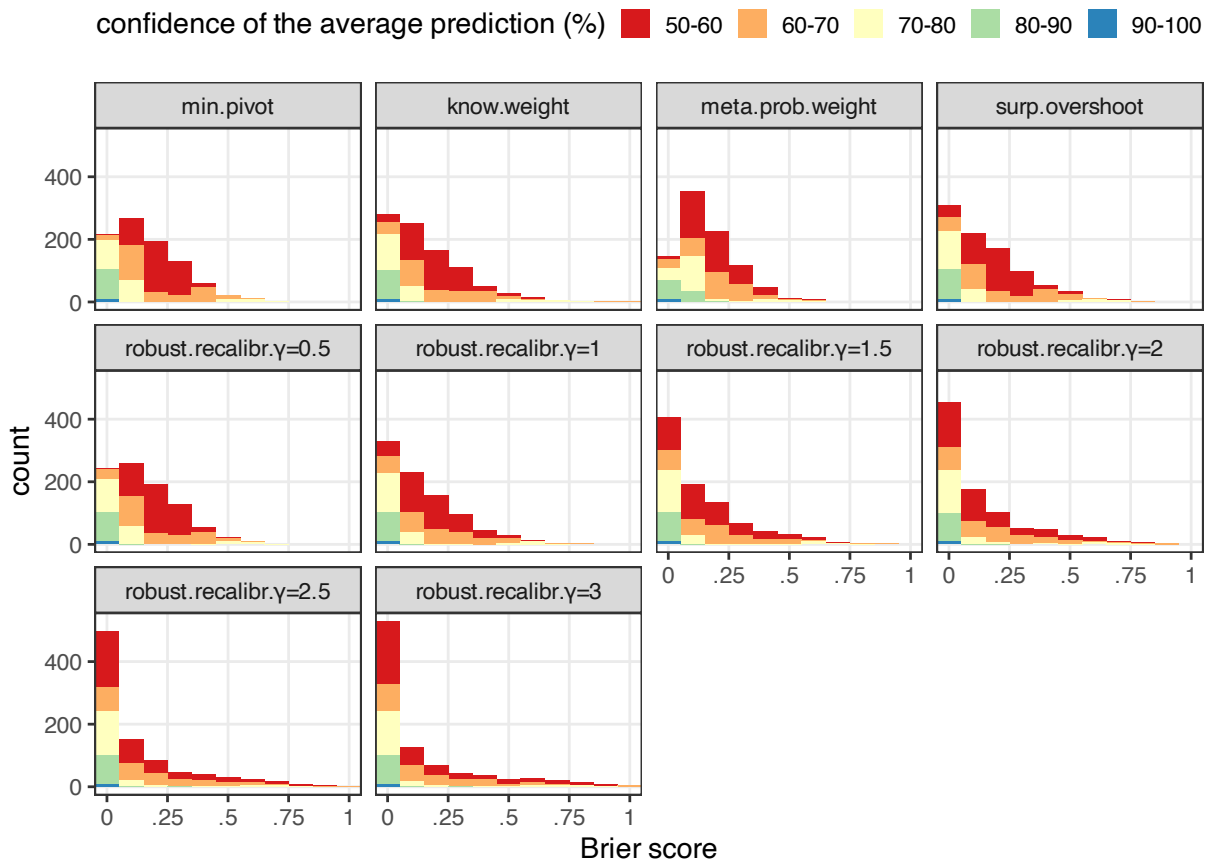


Figure 7: Brier scores of simple average, extremized average and robust-recalibrated probabilities.

663 Figure 7 demonstrates that robust recalibration achieves very small Brier scores more
 664 often than the benchmarks, in particular for $\gamma \geq 1$. The difference between the Brier scores

¹⁰Table F3 in Appendix F provides comparisons for all $\gamma \in \{0.5, 1, 1.5, 2, 2.5, 3\}$ for completeness.

665 of algorithms is significant (ANOVA test, F-value = 5.371, $p < 0.0001$).

666 We next look at pairwise comparisons of the robust recalibration method with $\gamma \in \{1.5, 2\}$
 667 to the other methods. Table 4 shows that the robust recalibration method achieves higher
 668 accuracy against all benchmarks for both values of γ . Table F4 in Appendix F reports the
 669 same pairwise tests for each dataset separately. We observe significantly higher accuracy
 670 for robust recalibration in the Science and States tasks but find that performance is similar
 671 between algorithms in the Arts and NFL trivia tasks. Thus the performance differences
 672 between algorithms are likely to relate to characteristics of the underlying data generating
 673 process.

Method	Benchmark	Avg.diff	Med.diff	Test stat.	p-value	Signif. better?
robust.recalibr. $\gamma=1.5$	know.weight	-0.0230	-0.0150	V=96,184	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1.5$	meta.prob.weight	-0.0212	-0.0363	V=103,043	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1.5$	min.pivot	-0.0296	-0.0257	V=103,024	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1.5$	surp.overshoot	-0.0197	-0.0118	V=123,548	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2$	know.weight	-0.0257	-0.0216	V=102,362	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2$	meta.prob.weight	-0.0239	-0.0467	V=107,335	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2$	min.pivot	-0.0323	-0.0328	V=110,455	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2$	surp.overshoot	-0.0224	-0.0188	V=122,617	<0.0001	robust.recalibr

Table 4: Comparison of Brier scores, two-sided paired Wilcoxon signed rank tests, robust recalibration with $\gamma \in \{1.5, 2\}$ vs benchmarks.

674 In addition to the Brier score, we also constructed the calibration curve for each algorithm
 675 to understand how each algorithm is reshaping the predictions. These calibration curves are
 676 presented in Figure 8 and were constructed using the same methodology as Figure 6. As
 677 seen in the diagram, robust recalibration achieves better calibration than the alternatives in
 678 most bins for $\gamma \in \{1.5, 2, 2.5, 3\}$. Predicted probabilities of robust-recalibrated aggregates
 679 are very close to the actual frequencies. Similar to the results in accuracy above, robust
 680 recalibration with sufficiently high γ appears to improve calibration over the alternatives.

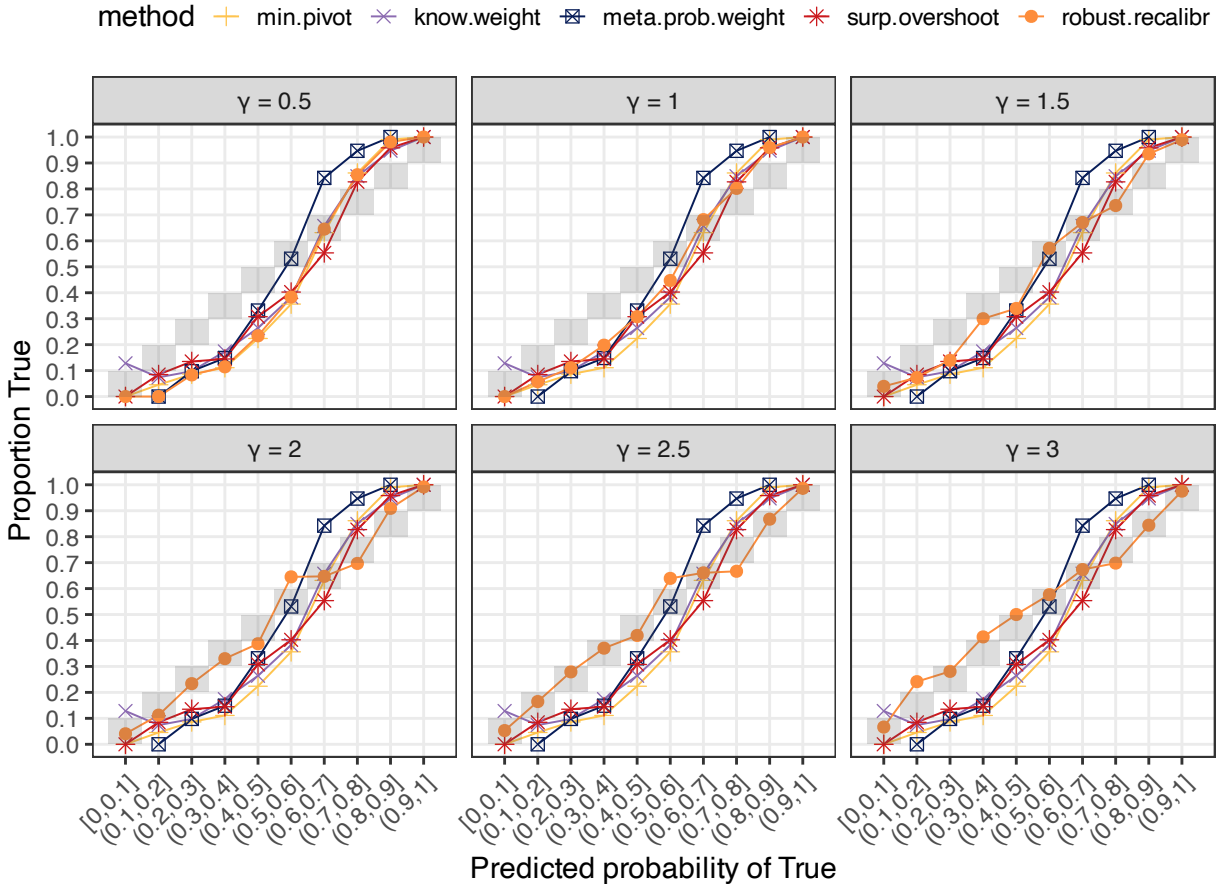


Figure 8: Calibration curves for simple average, extremized average and robust-recalibrated probabilities.

681 **6 Conclusion**

682 Probabilistic forecasts are often too conservative, which leads to average probability fore-
 683 casts not being sufficiently extreme. Previous work documented that extremizing transfor-
 684 mations that adjust the average away from 0.5 improve calibration. However, such transfor-
 685 mations may have shortcomings. In some forecasting problems, the crowd may have a biased
 686 prior that favors a certain outcome. Then, the average forecast may put a higher probabil-
 687 ity on the wrong outcome even when individuals receive informative signals conditional on
 688 the correct outcome. Extremizing a wrong-sided average forecast would introduce further
 689 miscalibration.

690 We show that forecasters' meta-beliefs on others' predictions can be used to estimate

691 the prior in single-question forecasting problems. We then propose a recalibration function
692 that transforms the average away from the estimated prior instead of 0.5. A bias in crowd’s
693 prior probability is reflected in the estimated prior. Thus, unlike simple extremization away
694 from 0.5, robust recalibration is capable of correctly transforming wrong-side averages in the
695 opposite direction of extremization, which should produce aggregate probability forecasts
696 with better calibration.

697 We test the performance of robust recalibration using prediction and meta-prediction
698 data from four distinct experimental tasks. We implement robust recalibration with var-
699 ious values of γ , which is a tuning parameter that controls the intensity of extremization
700 away from the estimated prior. Our findings suggest that robust recalibration is effective in
701 improving the accuracy and calibration of probability forecasts. We first demonstrate that
702 robust recalibration outperforms simple extremization away from 0.5 for all values of γ we
703 explored. Robust-recalibrated probabilities achieve lower Brier scores in most tasks and pre-
704 dict the actual frequency of occurrence more accurately than extremized averages. Robust
705 recalibration is particularly effective in transforming wrong-sided averages which are close
706 to 50%, which characterize most wrong-sided averages in our data set. We show that, unlike
707 simple extremization, prior estimation using meta-predictions can detect and transform such
708 wrong-sided averages towards the correct extreme.

709 We also compared robust recalibration to four single-question aggregation algorithms
710 developed by recent work (Martinie et al., 2020; Palley & Satopää, 2023; Palley & Soll,
711 2019; Peker, 2023). These algorithms also rely on meta-predictions as well as predictions,
712 but unlike robust recalibration, they do not require a tuning parameter. Thus, they present
713 natural alternatives to our algorithm when meta-prediction data are available. We find that
714 robust recalibration achieves significantly higher accuracy in most tasks when using tuning
715 parameters suggested in the literature. The method also improves calibration provided that
716 γ is sufficiently high. Intuitively, the aggregation algorithms we considered are expected
717 to achieve some improvement in accuracy over simple averaging. Robust recalibration real-

718 izes further gains when transformation away from the estimated prior is sufficiently strong,
719 implying that prior estimation is effective in finding the correct direction to transform the
720 average prediction.

721 Similar to the benchmark algorithms, robust recalibration considers a single forecasting
722 problem where no data other than predictions and meta-predictions are available. Optimal
723 value of γ in a given problem is unknown. Our results suggest that the aggregator may
724 prefer to be aggressive rather than cautious in extremizing away from the estimated prior.
725 Subsequent work may test if this result generalizes to a larger set of forecast aggregation
726 problems. Furthermore, task-level analysis suggests that there is heterogeneity in the relative
727 effectiveness of our algorithm across the tasks studied. Robust recalibration achieved higher
728 accuracy in Science and States tasks, while we see a similar performance to other benchmarks
729 in Artwork and NFL tasks. Future work may investigate if the gains in accuracy differ in
730 various other domains of forecasting as well.

731 Robust recalibration procedure may have practical limitations due to the prior estima-
732 tion stage. In two tasks out of 910 in our original data set, the estimated prior probability
733 is not within $(0, 1)$. Appendix D shows that the estimated meta-prediction functions in
734 these two tasks imply meta-predictions outside $(0, 1)$, leading to invalid prior estimates. We
735 observe that in both tasks, predictions are clustered at the correct extreme (0 or 1 depend-
736 ing on the correct answer). In other words, a strong majority of the forecasters were very
737 accurate in their predictions. Robust recalibration uses a linear regression model to esti-
738 mate the parameters. The actual meta-prediction function may not be estimated accurately
739 when predictions are heavily clustered or the sample of forecasters is small. As discussed in
740 Section 5.2, prior estimation is inaccurate if the estimated meta-prediction function implies
741 meta-predictions outside of the probability scale. Thus, in practical applications, the aggre-
742 gator can use the information from the estimation procedure to decide on the applicability
743 of robust recalibration.

References

- Ariely, D., Tung Au, W., Bender, R. H., Budescu, D. V., Dietz, C. B., Gu, H., Wallsten, T. S., & Zauberman, G. (2000). The effects of averaging subjective probability estimates between and within judges. *Journal of Experimental Psychology: Applied*, *6*(2), 130–147. <https://doi.org/10.1037/1076-898X.6.2.130>
- Atanasov, P., Rescober, P., Stone, E., Swift, S. A., Servan-Schreiber, E., Tetlock, P., Ungar, L., & Mellers, B. (2017). Distilling the wisdom of crowds: Prediction markets vs. prediction polls. *Management Science*, *63*(3), 691–706. <https://doi.org/10.1287/mnsc.2015.2374>
- Baron, J., Mellers, B. A., Tetlock, P. E., Stone, E., & Ungar, L. H. (2014). Two reasons to make aggregated probability forecasts more extreme. *Decision Analysis*, *11*(2), 133–145. <https://doi.org/10.1287/deca.2014.0293>
- Breiman, L. (1996). Stacked regressions. *Machine learning*, *24*, 49–64. <https://doi.org/10.1007/BF00117832>
- Budescu, D. V., Wallsten, T. S., & Au, W. T. (1997). On the importance of random error in the study of probability judgment. part ii: Applying the stochastic judgment model to detect systematic trends. *Journal of Behavioral Decision Making*, *10*(3), 173–188. [https://doi.org/10.1002/\(SICI\)1099-0771\(199709\)10:3<173::AID-BDM261>3.0.CO;2-6](https://doi.org/10.1002/(SICI)1099-0771(199709)10:3<173::AID-BDM261>3.0.CO;2-6)
- Chen, Y.-C., Mueller-Frank, M., & Pai, M. M. (2021). The wisdom of the crowd and higher-order beliefs. <https://arxiv.org/abs/2102.02666>
- Clemen, R. T. (1989). Combining forecasts: A review and annotated bibliography. *International Journal of Forecasting*, *5*(4), 559–583. [https://doi.org/10.1016/0169-2070\(89\)90012-5](https://doi.org/10.1016/0169-2070(89)90012-5)
- Dana, J., Atanasov, P., Tetlock, P., & Mellers, B. (2019). Are markets more accurate than polls? the surprising informational value of “just asking”. *Judgment and Decision Making*, *14*(2), 135–147. <https://doi.org/10.1017/S1930297500003375>

771 Dietrich, F. (2010). Bayesian group belief. *Social choice and welfare*, *35*, 595–626. <https://doi.org/10.1007/s00355-010-0453-x>

772

773 Erev, I., Wallsten, T. S., & Budescu, D. V. (1994). Simultaneous over-and underconfidence:
774 The role of error in judgment processes. *Psychological review*, *101*(3), 519. <https://doi.org/10.1037/0033-295X.101.3.519>

775

776 Genre, V., Kenny, G., Meyler, A., & Timmermann, A. (2013). Combining expert forecasts:
777 Can anything beat the simple average? *International Journal of Forecasting*, *29*(1),
778 108–121. <https://doi.org/10.1016/j.ijforecast.2012.06.004>

779 Han, Y., & Budescu, D. V. (2022). Recalibrating probabilistic forecasts to improve their
780 accuracy. *Judgment and Decision Making*, *17*(1), 91–123. [https://doi.org/10.1017/](https://doi.org/10.1017/S1930297500009049)
781 [S1930297500009049](https://doi.org/10.1017/S1930297500009049)

782 Hertwig, R. (2012). Tapping into the wisdom of the crowd—with confidence. *Science*, *336*(6079),
783 303–304. <https://doi.org/10.1126/science.1221403>

784 Howe, P. D., Martinie, M., & Wilkening, T. (2024). Using cross-domain expertise to aggregate
785 forecasts when within-domain expertise is unknown. *Decision*, *11*(1), 35–59. <https://doi.org/10.1037/dec0000212>

786

787 Jia, Y., Keppo, J., & Satopää, V. (2024). The wisdom of strategically diverse crowds. *Avail-*
788 *able at SSRN 4855714*. <https://ssrn.com/abstract=4855714>

789 Karmarkar, U. S. (1978). Subjectively weighted utility: A descriptive extension of the ex-
790 pected utility model. *Organizational behavior and human performance*, *21*(1), 61–72.
791 [https://doi.org/10.1016/0030-5073\(78\)90039-9](https://doi.org/10.1016/0030-5073(78)90039-9)

792 Koriat, A. (2008). Subjective confidence in one’s answers: The consensuality principle. *Jour-*
793 *nal of Experimental Psychology: Learning, Memory, and Cognition*, *34*(4), 945. <https://doi.org/10.1037/0278-7393.34.4.945>

794

795 Koriat, A. (2012). When are two heads better than one and why? *Science*, *336*(6079), 360–
796 362. <https://doi.org/10.1126/science.1216549>

- 797 Larrick, R. P., & Soll, J. B. (2006). Intuitions about combining opinions: Misappreciation
798 of the averaging principle. *Management Science*, *52*(1), 111–127. [https://doi.org/10.](https://doi.org/10.1287/mnsc.1050.0459)
799 [1287/mnsc.1050.0459](https://doi.org/10.1287/mnsc.1050.0459)
- 800 Lee, M. D., & Lee, M. N. (2017). The relationship between crowd majority and accuracy for
801 binary decisions. *Judgment and Decision Making*, *12*(4), 328–343. [https://doi.org/](https://doi.org/10.1017/S1930297500006227)
802 [10.1017/S1930297500006227](https://doi.org/10.1017/S1930297500006227)
- 803 Libgober, J. (2023). Identifying wisdom (of the crowd): A regression approach. [https://arxiv.](https://arxiv.org/abs/2105.07097)
804 [org/abs/2105.07097](https://arxiv.org/abs/2105.07097)
- 805 Lichtendahl Jr, K. C., Grushka-Cockayne, Y., Jose, V. R., & Winkler, R. L. (2022). Extrem-
806 izing and antiextremizing in bayesian ensembles of binary-event forecasts. *Operations*
807 *Research*, *70*(5), 2998–3014. <https://doi.org/10.1287/opre.2021.2176>
- 808 Martinie, M., Wilkening, T., & Howe, P. D. (2020). Using meta-predictions to identify experts
809 in the crowd when past performance is unknown. *Plos one*, *15*(4), e0232058. [https:](https://doi.org/10.1371/journal.pone.0232058)
810 [//doi.org/10.1371/journal.pone.0232058](https://doi.org/10.1371/journal.pone.0232058)
- 811 Mellers, B., Ungar, L., Baron, J., Ramos, J., Gurcay, B., Fincher, K., Scott, S. E., Moore,
812 D., Atanasov, P., Swift, S. A., et al. (2014). Psychological strategies for winning a
813 geopolitical forecasting tournament. *Psychological science*, *25*(5), 1106–1115. [https:](https://doi.org/10.1177/0956797614524255)
814 [//doi.org/10.1177/0956797614524255](https://doi.org/10.1177/0956797614524255)
- 815 Neuwirth, E. (2022). *Rcolorbrewer: Colorbrewer palettes* [R package version 1.1-3]. [https:](https://doi.org/10.32614/CRAN.package.RColorBrewer)
816 [//doi.org/10.32614/CRAN.package.RColorBrewer](https://doi.org/10.32614/CRAN.package.RColorBrewer)
- 817 Palley, A., & Satopää, V. A. (2023). Boosting the wisdom of crowds within a single judgment
818 problem: Weighted averaging based on peer predictions. *Management Science*, *69*(9),
819 5128–5146. <https://doi.org/10.1287/mnsc.2022.4648>
- 820 Palley, A., & Soll, J. (2019). Extracting the wisdom of crowds when information is shared.
821 *Management Science*, *65*(5), 2291–2309. <https://doi.org/10.1287/mnsc.2018.3047>
- 822 Peker, C. (2023). Extracting the collective wisdom in probabilistic judgments. *Theory and*
823 *Decision*, *94*(3), 467–501. <https://doi.org/10.1007/s11238-022-09899-4>

824 Prelec, D., Seung, H. S., & McCoy, J. (2017). A solution to the single-question crowd wisdom
825 problem. *Nature*, *541*(7638), 532–535. <https://doi.org/10.1038/nature21054>

826 R Core Team. (2023). *R: A language and environment for statistical computing*. R Foundation
827 for Statistical Computing. Vienna, Austria. <https://www.R-project.org/>

828 Raftery, A. E., Madigan, D., & Hoeting, J. A. (1997). Bayesian model averaging for linear
829 regression models. *Journal of the American Statistical Association*, *92*(437), 179–191.
830 <https://doi.org/10.1080/01621459.1997.10473615>

831 Ranjan, R., & Gneiting, T. (2010). Combining probability forecasts. *Journal of the Royal*
832 *Statistical Society Series B: Statistical Methodology*, *72*(1), 71–91. [https://doi.org/](https://doi.org/10.1111/j.1467-9868.2009.00726.x)
833 [10.1111/j.1467-9868.2009.00726.x](https://doi.org/10.1111/j.1467-9868.2009.00726.x)

834 Rilling, J. (2024). Neutral pivoting: Strong bias correction for shared information. [https:](https://arxiv.org/abs/2404.17737)
835 [//arxiv.org/abs/2404.17737](https://arxiv.org/abs/2404.17737)

836 RStudio Team. (2020). *Rstudio: Integrated development environment for r*. RStudio, PBC.
837 Boston, MA. <http://www.rstudio.com/>

838 Satopää, V. A. (2022). Regularized aggregation of one-off probability predictions. *Operations*
839 *Research*, *70*(6), 3558–3580. <https://doi.org/10.1287/opre.2021.2224>

840 Satopää, V. A., Baron, J., Foster, D. P., Mellers, B. A., Tetlock, P. E., & Ungar, L. H. (2014).
841 Combining multiple probability predictions using a simple logit model. *International*
842 *Journal of Forecasting*, *30*(2), 344–356. [https://doi.org/10.1016/j.ijforecast.2013.09.](https://doi.org/10.1016/j.ijforecast.2013.09.009)
843 [009](https://doi.org/10.1016/j.ijforecast.2013.09.009)

844 Satopää, V. A., Jensen, S. T., Mellers, B. A., Tetlock, P. E., & Ungar, L. H. (2014). Probabil-
845 ity aggregation in time-series: Dynamic hierarchical modeling of sparse expert beliefs.
846 *The Annals of Applied Statistics*, *8*(2), 1256–1280. [https://doi.org/10.1214/14-](https://doi.org/10.1214/14-AOAS739)
847 [AOAS739](https://doi.org/10.1214/14-AOAS739)

848 Satopää, V. A., Jensen, S. T., Pemantle, R., & Ungar, L. H. (2017). Partial information
849 framework: Model-based aggregation of estimates from diverse information sources.

850 *Electronic Journal of Statistics*, 11(2), 3781–3814. <https://doi.org/10.1214/17->
851 EJS1346

852 Satopää, V. A., Pemantle, R., & Ungar, L. H. (2016). Modeling probability forecasts via in-
853 formation diversity. *Journal of the American Statistical Association*, 111(516), 1623–
854 1633. <https://doi.org/10.1080/01621459.2015.1100621>

855 Shlomi, Y., & Wallsten, T. S. (2010). Subjective recalibration of advisors’ probability es-
856 timates. *Psychonomic bulletin & review*, 17(4), 492–498. [https://doi.org/10.3758/](https://doi.org/10.3758/PBR.17.4.492)
857 PBR.17.4.492

858 Surowiecki, J. (2005). *The wisdom of crowds*. Anchor.

859 Turner, B. M., Steyvers, M., Merkle, E. C., Budescu, D. V., & Wallsten, T. S. (2014). Forecast
860 aggregation via recalibration. *Machine learning*, 95(3), 261–289. [https://doi.org/10.](https://doi.org/10.1007/s10994-013-5401-4)
861 1007/s10994-013-5401-4

862 Wickham, H. (2007). Reshaping data with the reshape package. *Journal of Statistical Soft-*
863 *ware*, 21(12), 1–20. <https://doi.org/10.18637/jss.v021.i12>

864 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., Grolemund,
865 G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T. L., Miller, E., Bache,
866 S. M., Müller, K., Ooms, J., Robinson, D., Seidel, D. P., Spinu, V., . . . Yutani, H.
867 (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43), 1686.
868 <https://doi.org/10.21105/joss.01686>

869 Wilkening, T., Martinie, M., & Howe, P. D. (2022). Hidden experts in the crowd: Using meta-
870 predictions to leverage expertise in single-question prediction problems. *Management*
871 *Science*, 68(1), 487–508. <https://doi.org/10.1287/mnsc.2020.3919>

872 Winkler, R. L., Grushka-Cockayne, Y., Lichtendahl, K. C., & Jose, V. R. R. (2019). Prob-
873 ability forecasts and their combination: A research perspective. *Decision Analysis*,
874 16(4), 239–260. <https://doi.org/10.1287/deca.2019.0391>

875 Appendices

876 A Proofs

877 Lemma 1

878 This result is due to the fact that the expected posterior prediction generated from an
879 information service is equal to the prediction that would be made at the prior. At the prior:

$$\begin{aligned} P(s_\emptyset) = P(E|\sigma_k = s_\emptyset) &= \sum_i [P(E|s_i)P(s_i|s_\emptyset)] \\ &= \sum_i [qP(E|s_i)P(s_i|\omega_G) + (1-q)P(E|s_i)P(s_i|\omega_B)] \\ &= q \sum_i [P(E|s_i)P(s_i|\omega_G)] + (1-q) \sum_i [P(E|s_i)P(s_i|\omega_B)] \\ &= q\mathbb{E}[P|\omega_G] + (1-q)\mathbb{E}[P|\omega_B]. \end{aligned}$$

In the main text, we showed that

$$M(\sigma_k) = \sigma_k\mathbb{E}[P|\omega_G] + (1-\sigma_k)\mathbb{E}[P|\omega_B].$$

and thus

$$M(s_\emptyset) = q\mathbb{E}[P|\omega_G] + (1-q)\mathbb{E}[P|\omega_B].$$

880 It follows immediately that $P(s_\emptyset) = M(s_\emptyset)$.

881 Proposition 1

Consider the case $\omega = \omega_G$. Following the notation in Equation 8, let $t_{RR}(\bar{P})$ denote the robust-recalibrated probability. Also let $t_E(\bar{P})$ be simple-extremized probability ($\delta = 1$ in Equation 8) with the same tuning parameter. Robust recalibration would achieve lower average Brier score if $(t_{RR}(\bar{P}) - g)^2 < (t_E(\bar{P}) - g)^2$, i.e. when the robust-recalibrated

probability is more accurate. This expression reduces to $\frac{1}{2}(t_{RR}(\bar{P}) + t_E(\bar{P})) < g$, which gives

$$\frac{1}{2} \left(\frac{\delta \bar{P}^{1+\gamma}}{\delta \bar{P}^{1+\gamma} + (1 - \bar{P})^{1+\gamma}} + \frac{\bar{P}^{1+\gamma}}{\bar{P}^{1+\gamma} + (1 - \bar{P})^{1+\gamma}} \right) < g$$

882 Note that $\lim_{N \rightarrow \infty} \hat{P}(s_\theta) = P(s_\theta)$, i.e. estimated prior converges to the actual prior prediction
 883 at the limit. Thus, $\delta = [(1 - P(s_\theta))/P(s_\theta)]^\gamma$. Also note that $\lim_{N \rightarrow \infty} \bar{P} = E[P|\omega_G]$ in state ω_G .
 884 Since we consider wrong-sided problems, we have $P(s_\theta) < \bar{P} < 0.5$. Then, we have $\delta > 1$ for
 885 any γ .

886 We can define $g'(\delta) = \frac{1}{2}(t_{RR}(\bar{P}) + t_E(\bar{P}))$ as the threshold such that, for $g > g'(\delta)$, $t_{RR}(\bar{P})$
 887 is strictly more accurate than $t_E(\bar{P})$ for any \bar{P} . Furthermore, $g'(\delta)$ increases as δ increases
 888 for all $\delta > 1$. Since δ increases as $P(s_\theta)$ decreases, $g'(\delta)$ increases with $|\bar{P} - P(s_\theta)|$.

889 A similar result can be obtained for $\omega = \omega_B$. Robust recalibration is more accurate if
 890 $(t_{RR}(\bar{P}) - b)^2 < (t_E(\bar{P}) - b)^2$ is satisfied, which reduces to $\frac{1}{2}(t_{RR}(\bar{P}) + t_E(\bar{P})) > b$. We
 891 now have $\lim_{N \rightarrow \infty} \bar{P} = E[P|\omega_B]$, $0.5 < \bar{P} < P(s_\theta)$, and $\delta < 1$ for any γ . We can define
 892 $b'(\delta) = \frac{1}{2}(t_{RR}(\bar{P}) + t_E(\bar{P}))$, which decreases as $P(s_\theta)$ increases, implying that the threshold
 893 $b'(\delta)$ decreases with $|\bar{P} - P(s_\theta)|$.

894 **Proposition 2**

895 From the proof of Proposition 1, we know that $\delta > 1$ for $\omega = \omega_G$ and $\delta < 1$ for $\omega = \omega_B$.
 896 Then, we simply have $t_E(\bar{P}) < t_{RR}(\bar{P}) < g = 1$ for $\omega = \omega_G$ and $b = 0 < t_{RR}(\bar{P}) < t_E(\bar{P})$. In
 897 both states, robust-recalibrated probability is strictly more accurate.

898 **Proposition 3**

899 Consider the case $\omega = \omega_G$. Since average forecast is not wrong sided, we have $0.5 <$
 900 $\bar{P} < 1$. As in the proof of Proposition 1, let $t_{RR}(\bar{P})$ and $t_E(\bar{P})$ denote robust-recalibrated
 901 and extremized probabilities. We have $t_E(\bar{P}) < t_{RR}(\bar{P}) < 1 \iff \delta > 1$, which requires
 902 $P(s_\theta) < 0.5$. Thus, $t_{RR}(\bar{P})$ is strictly more accurate when $|\bar{P} - P(s_\theta)| > |\bar{P} - 0.5|$. Similarly

903 for $\omega = \omega_B$, we have $0 < t_{RR}(\bar{P}) < t_E(\bar{P}) \iff \delta < 1$, which holds for $P(s_\emptyset) > 0.5$, and
 904 $0 < \bar{P} < 0.5$. So, $t_{RR}(\bar{P})$ outperforms $t_E(\bar{P})$ when $|\bar{P} - P(s_\emptyset)| > |\bar{P} - 0.5|$.

905 **B Robust Recalibration with more than two states**

906 In the main text, we showed that it is always possible to correctly estimate the prior using
 907 prediction and meta-predictions in an environment where there are exactly two states. This
 908 ensured that the algorithm would always identify the correct direction for extremization
 909 in large sample. In this section, we use two examples to show that the properties of the
 910 algorithm are not guaranteed when there are more than two states. The first example shows
 911 that the prediction and meta-prediction lines may cross multiple times when we increase the
 912 state space and that the estimated prior may not be correct. Nonetheless, the algorithm
 913 may still function well as long as the estimated prior still identifies the correct direction for
 914 extremization.

915 The second example identifies a situation where our algorithm fails to extremize in the
 916 correct direction for one of the states. The counter-example highlights a case where the
 917 monotone likelihood ratio principal is violated and where signals are very informative about
 918 the signals of others but only weakly informative about the underlying likelihood of an event.
 919 In such cases, it is possible to construct situations where the meta-prediction line is non-
 920 linear and create perverse cases where the algorithm fails. We see such situations as being
 921 quite rare, but the possibility of such cases warrant an empirical exploration of the algorithm.

922 In both examples, we use a general likelihood matrix \mathbf{Q} where the rows correspond to
 923 states and the columns relate to signals. Predictions and meta-predictions can be written
 924 using the posterior beliefs for each state just as in Section 3.

925 **Example 1: Multiple Cross Points where the estimated posterior is incorrect**
 926 **but the direction of extremization is correct.** Suppose there are four states with
 927 probabilities of E given by $\{.8, .6, .4, .2\}$. For simplicity, we will refer to the states by using

928 the corresponding probability. Forecasters have a prior of $\{1/4, 1/4, 1/4, 1/4\}$ over the states.
 929 Each forecaster receives a signal from $\{s_1, s_2, s_\emptyset, s_3, s_4\}$. The likelihood matrix is given by

$$\mathbf{Q} = \begin{bmatrix} 0 & 0 & 0 & \frac{1}{3} & \frac{2}{3} \\ 0 & 0 & 0 & \frac{2}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{2}{3} & 0 & 0 & 0 \\ \frac{2}{3} & \frac{1}{3} & 0 & 0 & 0 \end{bmatrix}.$$

930 Rows 1 to 4 (top to bottom) give the likelihoods for states 0.8, 0.6, 0.4 and 0.2 respectively
 931 while columns 1 to 5 (left to right) represents the signals $s_1, s_2, s_\emptyset, s_3$ and s_4 . Unlike the binary
 932 framework, the signals do not represent the posterior beliefs on one of the states. However,
 933 signals with a higher index indicate a weakly higher posterior probability on the “best” state
 934 (i.e. state 0.8). In this example, $\{s_3, s_4\}$ are generated when we are in state .8 or .6, while
 935 $\{s_1, s_2\}$ occur in states .4 and .2. Posterior belief on state 0.8 is highest for s_4 , followed by
 936 s_3 and s_1, s_2 where the last two imply zero probability. Figure B1 depicts the corresponding
 937 prediction and meta-prediction functions.

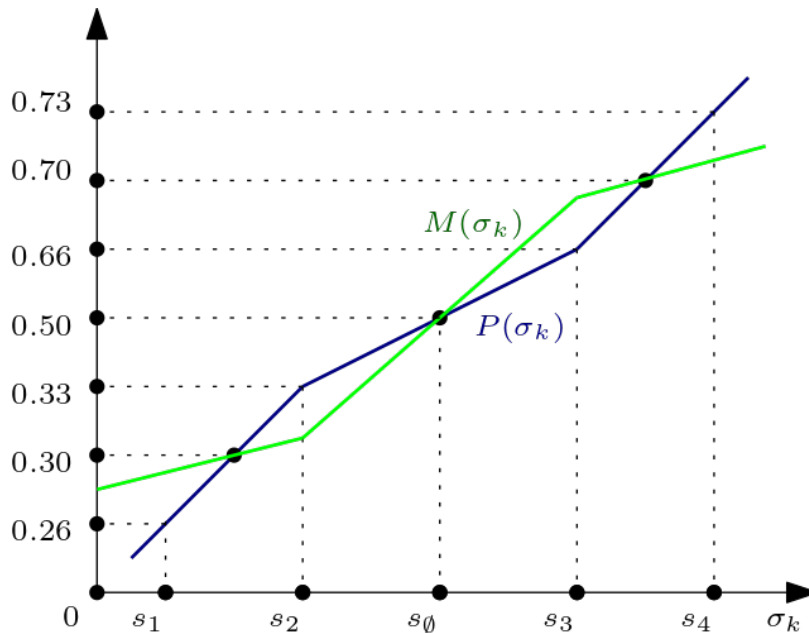


Figure B1: Example 1 prediction and meta-prediction functions (linear extrapolations from the predictions and meta-predictions at $\sigma_k \in \{s_1, s_2, s_\emptyset, s_3, s_4\}$).

938 The prediction and meta-prediction functions intersect at two distinct values other than
939 s_\emptyset . Thus, solving for $M(x) = P(x)$ does not uniquely recover the prior. Nevertheless, this
940 example demonstrates that robust recalibration could transform the average in the correct
941 direction despite the inaccuracy in estimating $P(s_\emptyset)$. To see this, we first calculate the
942 average prediction, which are $\{0.71, 0.69, 0.31, 0.29\}$ in states $\{0.8, 0.6, 0.4, 0.2\}$ respectively.
943 If the true state is 0.2 or 0.4, we get $\sigma_k \in \{s_1, s_2\}$. Then, the estimated prior will be
944 0.3, as it would be the unique intersection of the prediction and meta-prediction functions
945 in the corresponding range. Robust recalibration transforms 0.29 and 0.31 away from 0.3,
946 which could lead to transformed probabilities closer to the true probability (0.2 and 0.4
947 respectively). In contrast, extremizing away from 0.5 adjusts 0.31 in the wrong direction in
948 state 0.4. A similar result holds in states 0.6 and 0.8. Then, the estimated prior will be 0.7.
949 Average predictions of 0.69 and 0.71 are robust-recalibrated in the correct direction while
950 extremizing away from 0.5 pushes 0.69 further away from the true probability of the event
951 in state 0.6.

952 Note that the robust recalibration procedure is effective even though it does not produce
953 an accurate estimate of the actual prior ($P(s_\emptyset)$) in any state. The likelihood matrix suggests
954 that the forecasters have a non-zero posterior probability for two states only. The prediction
955 and meta-prediction functions are locally linear and estimated prior gives the intersection.

Example 2: Violation of MLRP. Consider an example with three states with prob-
abilities $\{0.7, 0.4, 0\}$. Forecasters have a uniform prior $\{1/3, 1/3, 1/3\}$ over the states. Prior
prediction is given by $P(s_\emptyset) = \frac{1}{3} 0.7 + \frac{1}{3} 0.4 + \frac{1}{3} 0 = 0.367$. Each forecaster receives a signal
from $\{s_1, s_\emptyset, s_2, s_3\}$ according to the following likelihood matrix:

$$\mathbf{Q} = \begin{bmatrix} .3 & 0 & \frac{1}{3} & .367 \\ 0 & 0 & \frac{2}{3} & \frac{1}{3} \\ .7 & 0 & 0 & .3 \end{bmatrix}$$

956 Rows 1 to 3 give the likelihoods of each signal in states 0.7, 0.4 and 0 respectively. Signals

957 are ordered in the implied posterior belief on the best state (i.e. state 0.7) as $s_3 > s_2 > s_1$.
 958 The prediction function satisfies $P(s_1) = 0.21$, $P(s_2) = 0.5$ and $P(s_3) = 0.39$.

959 For meta-predictions, we first calculate the average prediction in each state, which leads
 960 to $E[\bar{P}|state = 0] = 0.264$, $E[\bar{P}|state = 0.4] = 0.463$ and $E[\bar{P}|state = 0.7] = 0.373$. For any
 961 agent with signal $\sigma_k \in \{s_1, s_\emptyset, s_2, s_3\}$, $M(\sigma_k)$ will be a convex combination of $E[\bar{P}|state]$ with
 962 weights being the posterior probabilities over the states. The resulting meta-prediction func-
 963 tion satisfies $M(s_1) = 0.296$, $M(s_\emptyset) = 0.367$, $M(s_2) = 0.433$ and $M(s_3) = 0.37$. Figure B2
 964 depicts the prediction and meta-prediction functions.

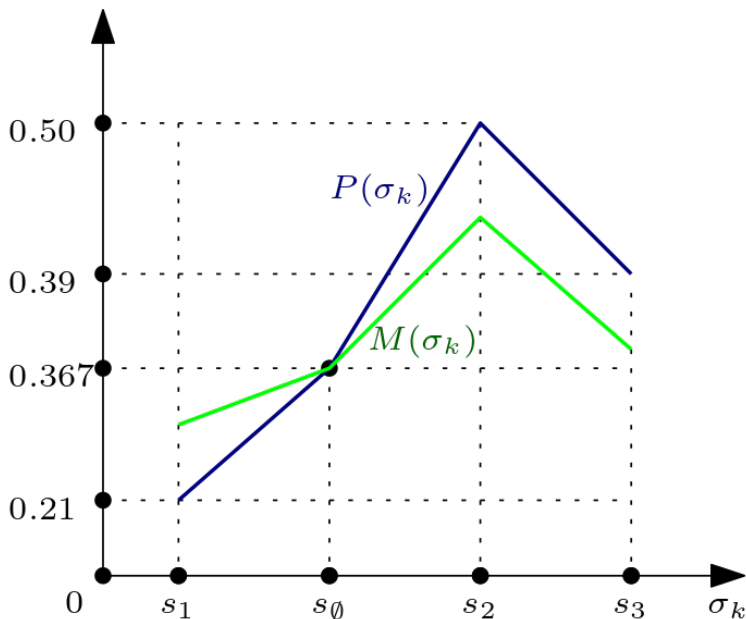


Figure B2: Example 2 prediction and meta-prediction functions

965 To see how robust recalibration performs, we randomly draw a sample of 10000 pre-
 966 dictions and meta-predictions according to the functions in Figure B2. Then, we intro-
 967 duce random noise in meta-predictions and estimate the prior as described in Section 4.
 968 This procedure is repeated 100 times. Average estimated priors in each state is given by
 969 $\{0.366, 0.344, 0.357\}$ with standard errors strictly smaller than 0.001. Recall that the average
 970 predictions are 0.264, 0.463 and 0.373 in states 0, 0.4 and 0.7 respectively. Thus, the average
 971 should be recalibrated down in states 0 and 0.4 and up in state 0.7. Robust recalibration

972 transforms the average predictions in states 0 and 0.7 in the correct direction. However, in
973 state 0.4, the robust recalibration procedure transforms the average in the wrong direction
974 while extremization away from 0.5 would push the average towards 0.4.

975 The miscalibration in state 0.4 is a result of s_2 being very informative about the predic-
976 tions of others and the likelihood that the state is not 0. Recall that the posterior beliefs
977 for states $\{0.7, 0.4, 0\}$ following s_3 and s_2 are $\{0.367, 1/3, 0.3\}$ and $\{1/3, 2/3, 0\}$ respectively.
978 Signal s_3 leads to the highest posterior belief on state 0.7 (followed by s_2 and s_1). However,
979 s_2 rules out the worst state and leads to a higher probability prediction and meta-prediction
980 overall. Since s_2 is more frequent in state 0.4, the resulting average prediction on the occur-
981 rence of the event is higher in state 0.4 than state 0.7, even though the event is more likely
982 in the latter.

983 In the binary framework, signals can be normalized to represent the posterior beliefs on
984 the good state (ω_G). When the true state is ω_G , signals favor a higher probability for the
985 occurrence of E . Therefore, $E[\bar{P}|\omega_G] > E[\bar{P}|\omega_B]$ always holds. The same is not necessarily
986 true for the “best state” in a multiple state framework where a signal is informative for
987 beliefs on more than one state. Likelihoods in state 0.4 (second row of \mathbf{Q}) suggest that all
988 forecasters observe s_2 or s_3 , and the corresponding predictions are 0.5 and 0.39. However,
989 in state 0.7 (first row of \mathbf{Q}), 30% of forecasters will observe s_1 and predict 0.21. As a result,
990 $E[\bar{P}|state = 0.4] > E[\bar{P}|state = 0.7]$. In other words, the information conveyed by signals
991 in state 0.4 favors high states (and hence, a higher probability for the event) more than
992 the information in state 0.7 on average. Unlike the binary framework, average prediction is
993 higher at a lower state. Such information structures are likely to be rare in practice, because
994 it would imply that the evidence itself is expected to incorrectly suggest a higher probability
995 for the occurrence of the event in a lower state. Thus, we expect robust recalibration to
996 perform well in most applications with more than two states.

997 **C Prediction tasks**

Table C1: Sample statements from Science and States data. See the supplemental material of Wilkening et al. (2022) for full list of statements

Data set	Statement
Science	Scurvy and anemia are diseases not caused by bacteria or viruses
Science	Secondary industries dominate the market in emerging economies
Science	Earthquakes and volcanoes typically occur at the boundaries of tectonic plates
Science	A substance with a pH of 8 is a strong acid
Science	Hamsters hate to run
Science	Plant cells are easier to clone than animal cells
Science	Convex lenses are used to correct for short-sightedness
Science	Darwin’s theory was not widely accepted when it was first published in the late 19th century
Science	Increasing the number of impermeable rocks in rivers help decrease the flood risk
States	Jacksonville is the capital city of Florida
States	Los Angeles is the capital city of California
States	Denver is the capital city of Colorado

Table C2: Sample NFL statements

Statement
In the 2018 NFL draft, Mark Andrews was drafted by the Minnesota Vikings
In the 2018 NFL draft, the New York Giants were the only team to draft a player out of FCS champion North Dakota State University
In the 2017 NFL draft, the Big Ten was one of the athletic conferences where no players were drafted that year
In the 2016 NFL draft, Rico Gathers was drafted by the Oakland Raiders
In the 2016 NFL draft, David Onyemata was drafted by the New Orleans Saints
In NFL rules, a player who wears illegal equipment is to be suspended for the next two games
In NFL rules, a delay of game penalty at the start of either half is a 5-yard penalty
In NFL rules, the penalty for attempting to use more than 3 timeouts in a half is 5 yards
In NFL, a “Hail Mary” is a play in which the receivers are all sent downfield towards the end zone
In NFL, a “two-point conversion” is a play a team attempts instead of kicking a one-point conversion immediately after it scores a touchdown

Figure C1: Sample items from the Artwork data set



998 **D Two tasks where robust recalibration failed to esti-**
 999 **mate the prior**

1000 Figure D1 shows the estimated meta-prediction function for the two Science tasks where
 1001 estimated prior lies outside $(0, 1)$. The statements are “Centimetres are a measure of length”
 1002 and “Fish have fur to keep them warm” with correct answers being true and false respectively.

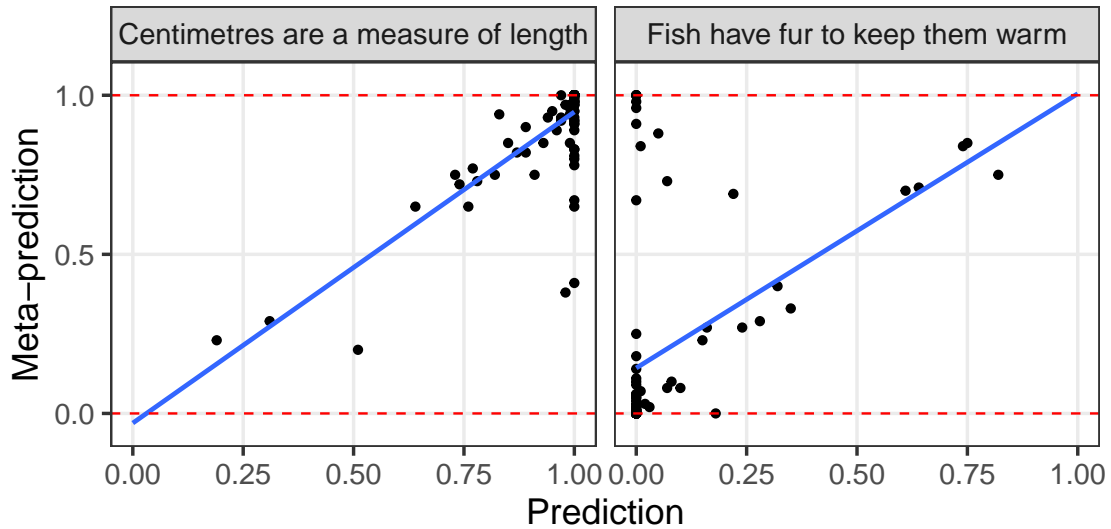


Figure D1: Estimated meta-prediction functions (blue line) in two tasks where estimated prior is not within $(0, 1)$

1003 Estimated meta-prediction functions (as in Equation 6) are $M_k = -0.0302 + 0.9778P_k$
 1004 (left panel) and $M_k = 0.1428 + 0.8622P_k$ (right panel). Note that $\hat{\beta}_0 < 0$ for “Centimetres are
 1005 a measure of length”, which leads to a negative estimated prior of -1.3602 from $\hat{\beta}_0/(1 - \hat{\beta}_1)$.
 1006 In “Fish have fur to keep them warm”, we have $\hat{\beta}_0 + \hat{\beta}_1 = 1.0049 > 1$, which leads to an
 1007 estimated prior of 1.0359 . Estimated prior probabilities are not within $(0, 1)$.

E Summary statistics and additional figures

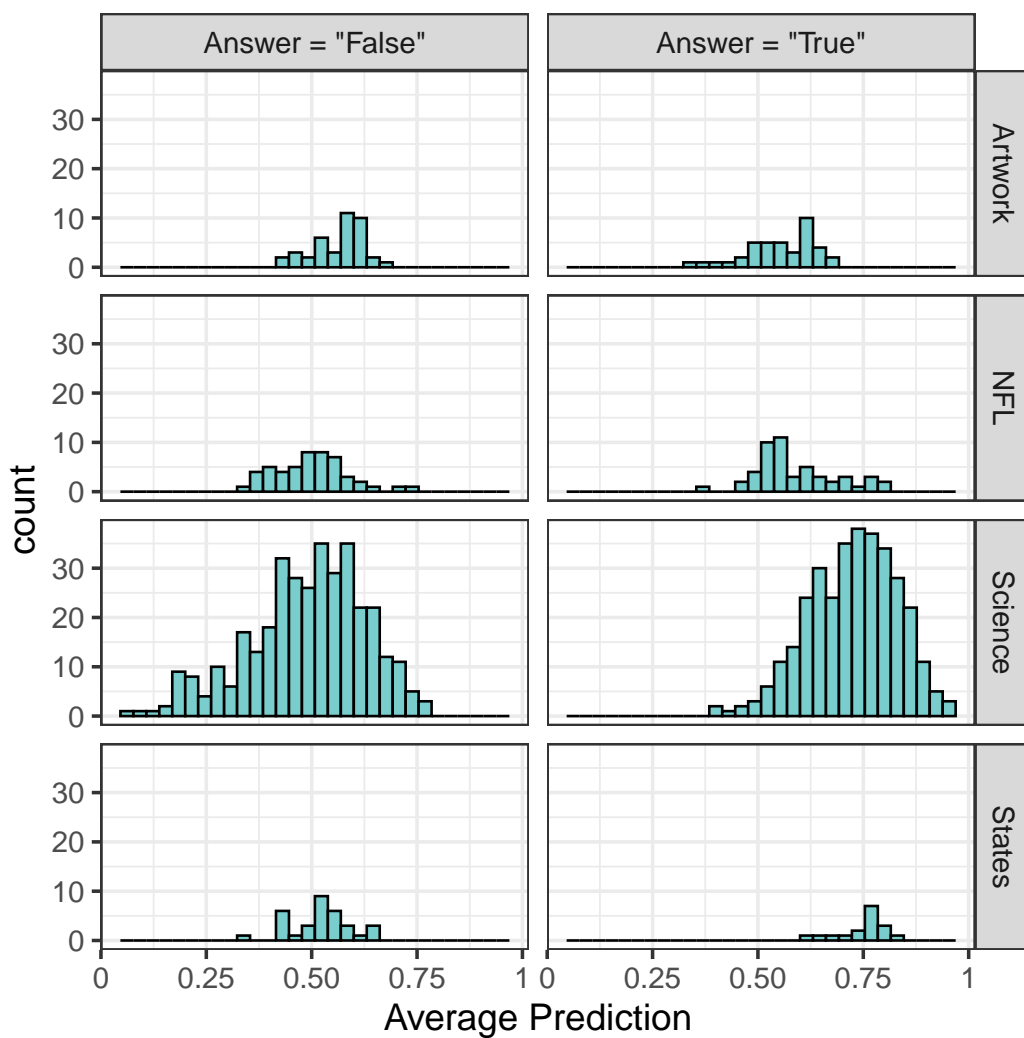


Figure E1: The distribution of average predictions for “True” and “False” statements in each data set.

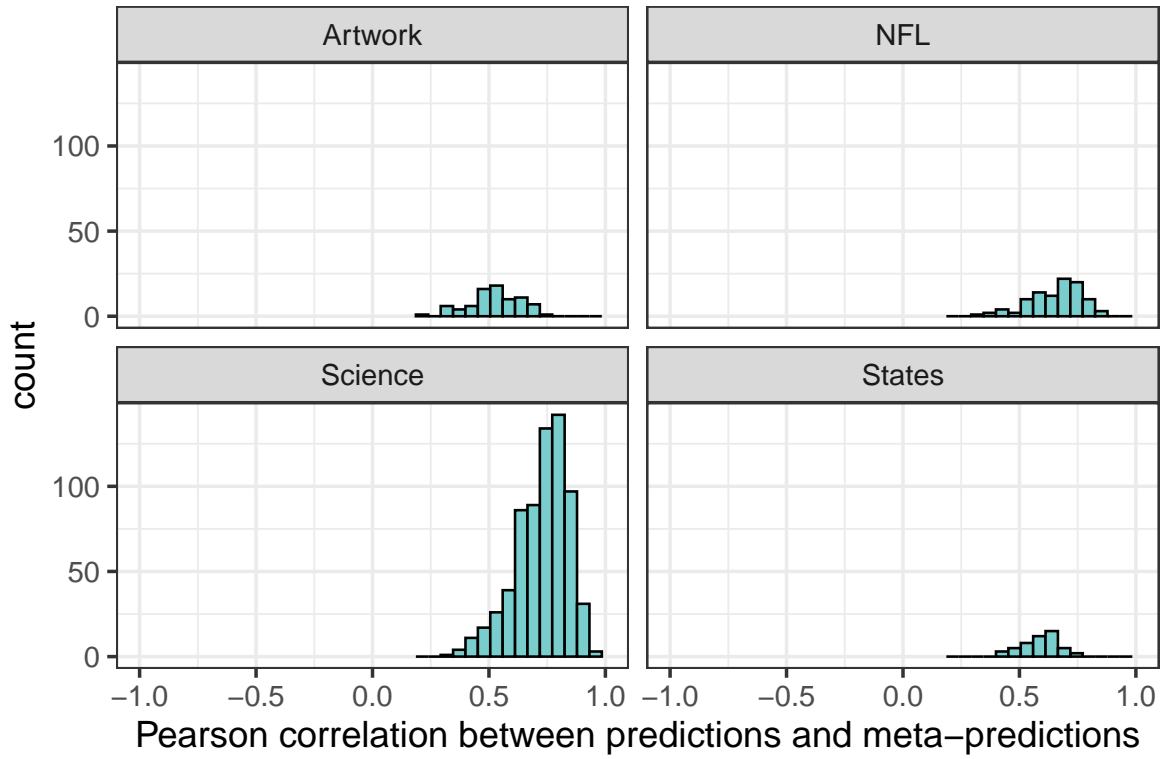


Figure E2: Correlation between predictions and meta-predictions. Each data point represents a task, 910 in total.

method	γ	min	max	mean	lower quartile	median	upper quartile
average		0.0018	0.5878	0.1901	0.0769	0.1737	0.2821
extrem.average	0.5	0.0001	0.7331	0.1859	0.0369	0.1418	0.2987
extrem.average	1	0.0000	0.8376	0.1886	0.0165	0.1143	0.3158
extrem.average	1.5	0.0000	0.9051	0.1944	0.0070	0.0909	0.3332
extrem.average	2	0.0000	0.9459	0.2012	0.0029	0.0715	0.3509
extrem.average	2.5	0.0000	0.9696	0.2083	0.0011	0.0556	0.3688
extrem.average	3	0.0000	0.9831	0.2150	0.0004	0.0428	0.3869
robust.recalibr	0.5	0.0001	0.6529	0.1610	0.0478	0.1314	0.2405
robust.recalibr	1	0.0000	0.7755	0.1455	0.0269	0.0968	0.2224
robust.recalibr	1.5	0.0000	0.8793	0.1381	0.0141	0.0689	0.2037
robust.recalibr	2	0.0000	0.9380	0.1354	0.0068	0.0494	0.1918
robust.recalibr	2.5	0.0000	0.9689	0.1355	0.0031	0.0370	0.1809
robust.recalibr	3	0.0000	0.9846	0.1372	0.0014	0.0259	0.1715

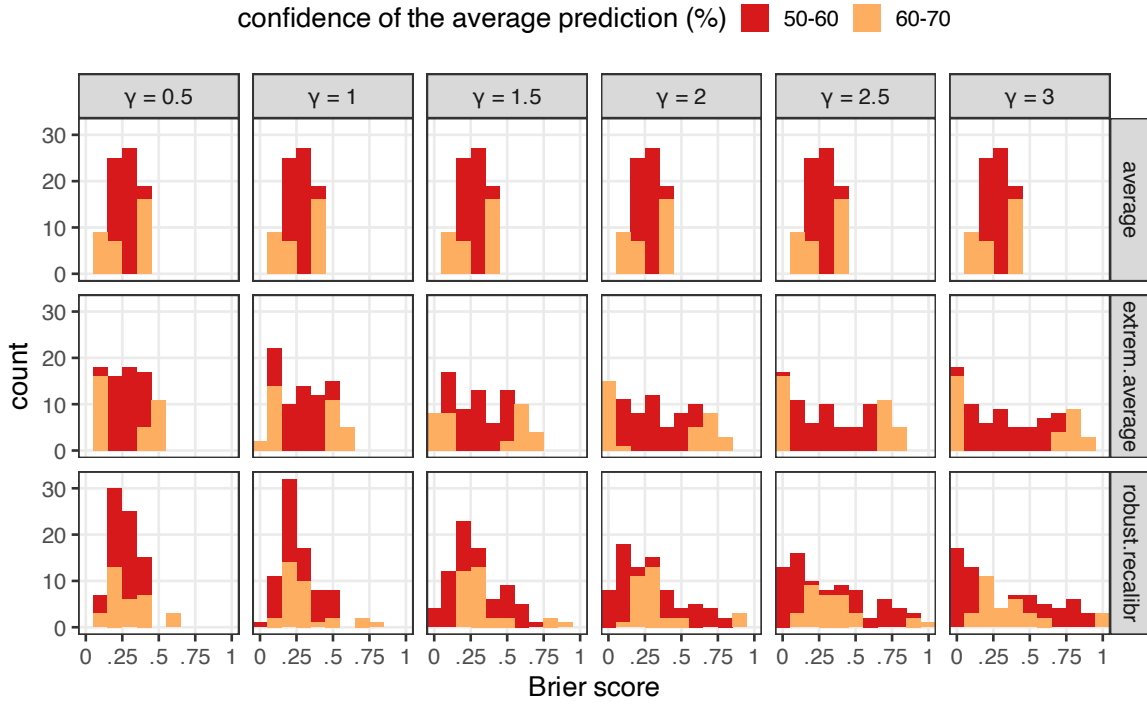
Table E1: Summary statistics, Brier scores in Figure 4.

method	γ	min	max	mean	lower quartile	median	upper quartile
min.pivot		0.0000	0.7031	0.1677	0.0527	0.1399	0.2512
know.weight		0.0000	1.0000	0.1611	0.0366	0.1136	0.2377
meta.prob.weight		0.0014	0.6384	0.1593	0.0723	0.1315	0.2207
surp.overshoot		0.0000	0.7569	0.1578	0.0324	0.1024	0.2500
robust.recalibr	0.5	0.0001	0.6529	0.1610	0.0478	0.1314	0.2405
robust.recalibr	1	0.0000	0.7755	0.1455	0.0269	0.0968	0.2224
robust.recalibr	1.5	0.0000	0.8793	0.1381	0.0141	0.0689	0.2037
robust.recalibr	2	0.0000	0.9380	0.1354	0.0068	0.0494	0.1918
robust.recalibr	2.5	0.0000	0.9689	0.1355	0.0031	0.0370	0.1809
robust.recalibr	3	0.0000	0.9846	0.1372	0.0014	0.0259	0.1715

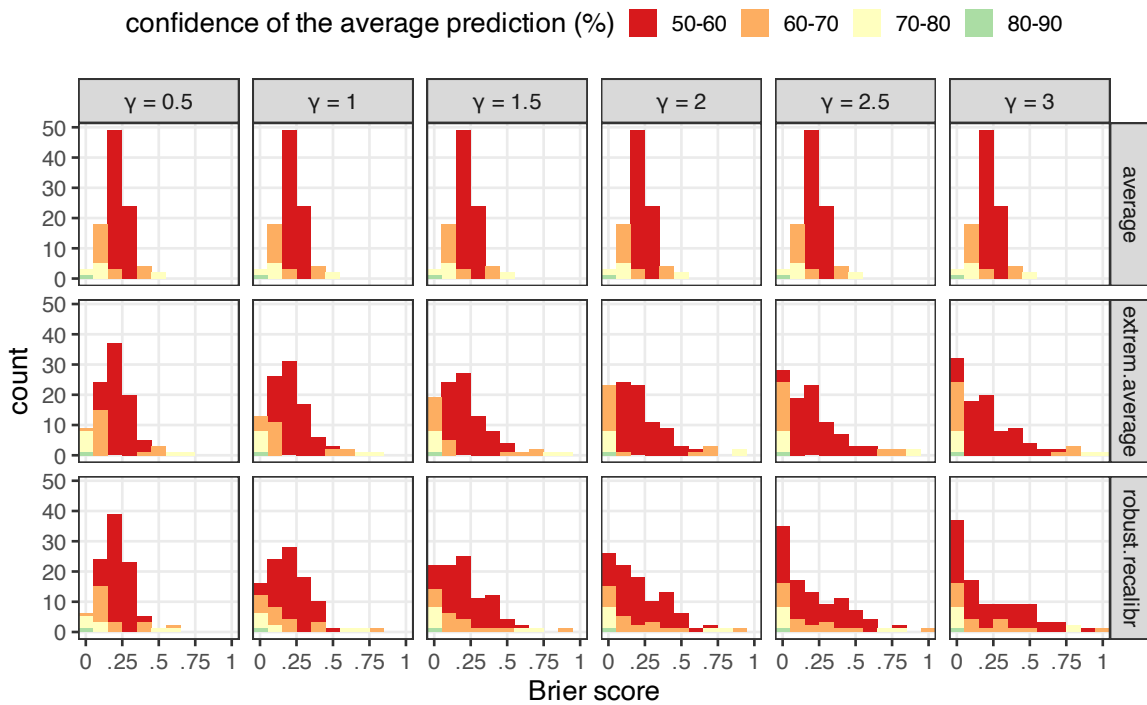
Table E2: Summary statistics, Brier scores in Figure 7.

1009 **F** Results by data set

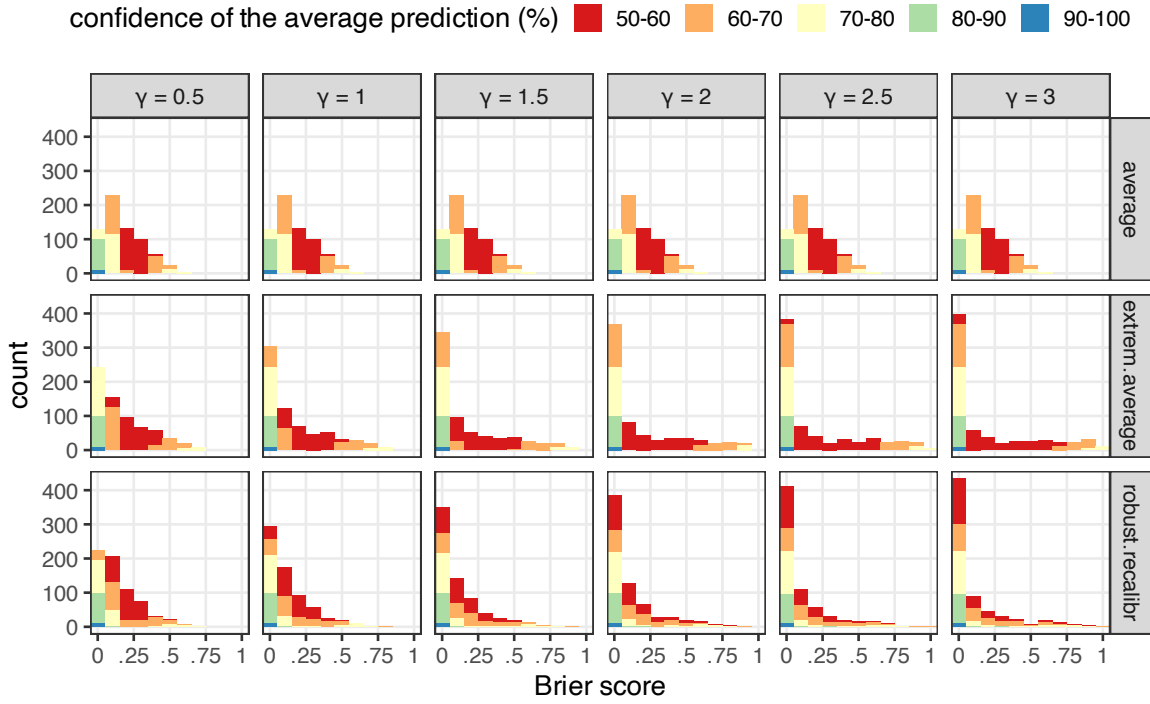
(a) Brier scores, Artwork data only.



(b) Brier scores, NFL data only.



(c) Brier scores, Science data only.



(d) Brier scores, States data only.

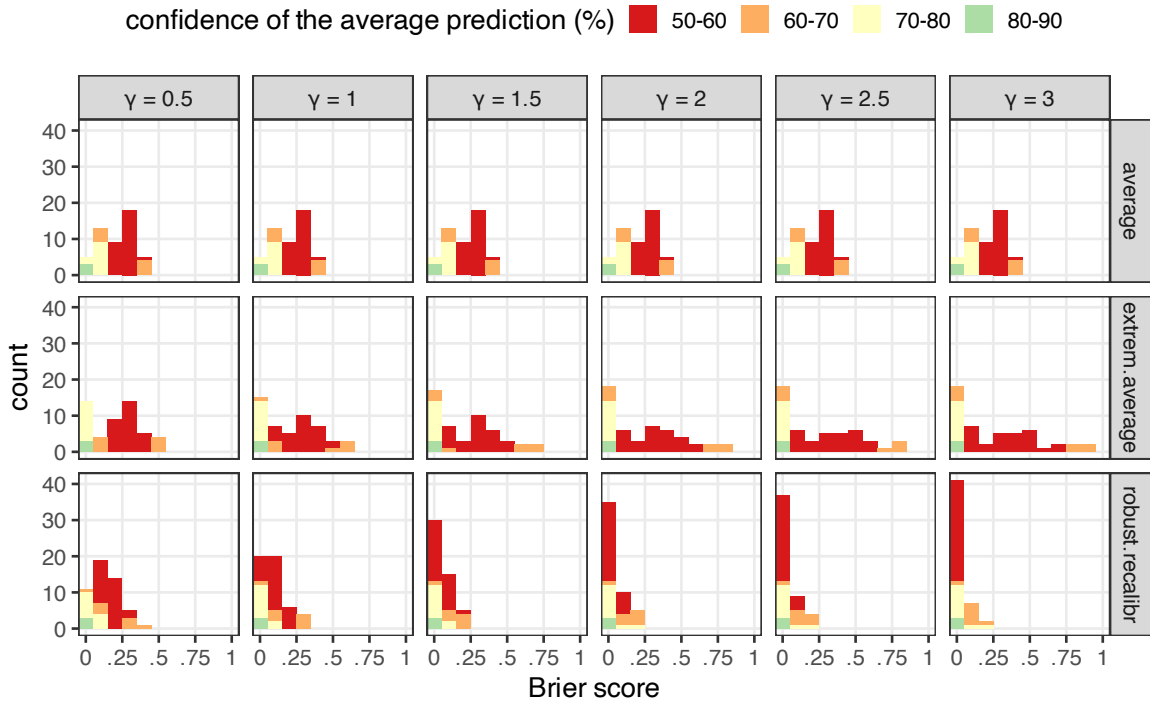
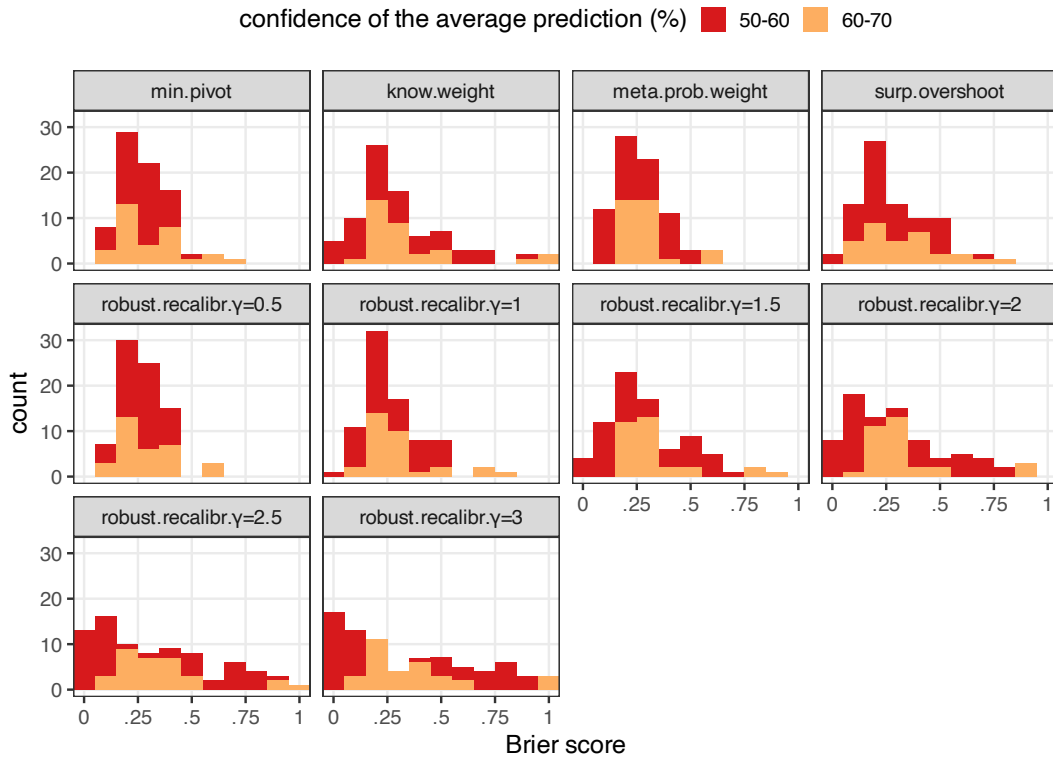
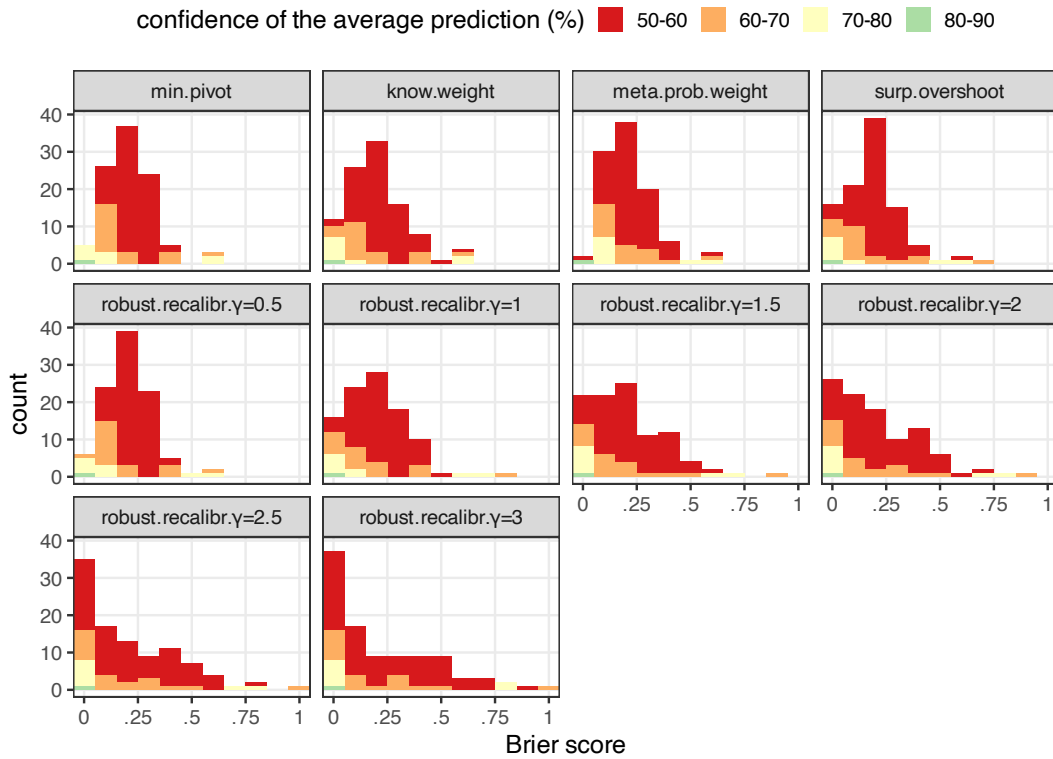


Figure F1: Brier scores of simple average, extremized average and robust-recalibrated probabilities.

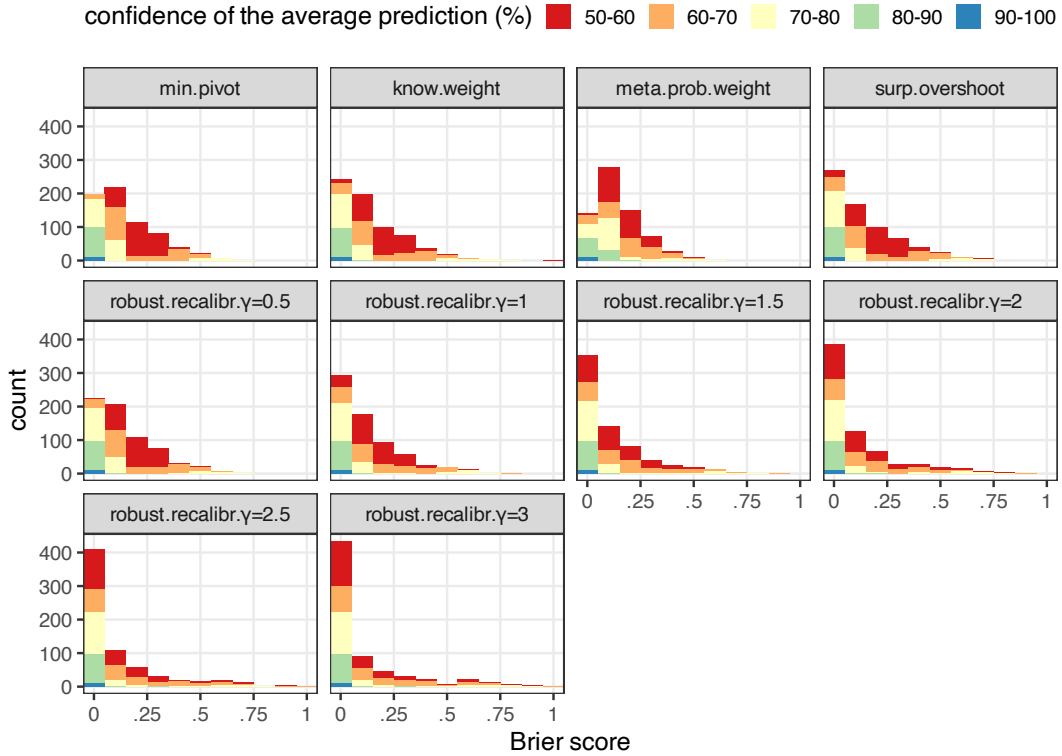
(a) Brier scores, Artwork data only.



(b) Brier scores, NFL data only.



(c) Brier scores, Science data only.



(d) Brier scores, States data only.

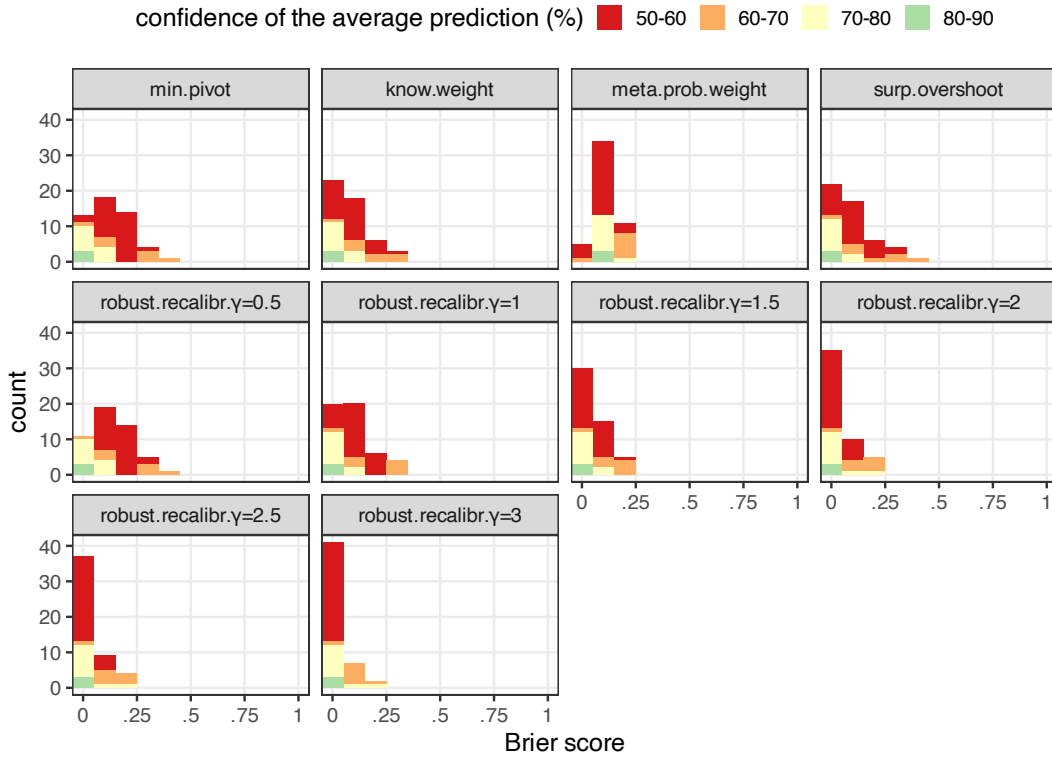


Figure F2: Brier scores of robust recalibration and other benchmarks.

(a) Artwork data only

γ	Method.1	Method.2	Avg.diff	Med.diff	Test stat.	p-value	Signif. better?
0.5	extrem.average	average	0.0135	0.0096	V=2,121	0.0164	Method.2
0.5	robust.recalibr	extrem.average	-0.0105	-0.0032	V=1,215	0.0524	No diff.
1	extrem.average	average	0.0292	0.0193	V=2,149	0.0112	Method.2
1	robust.recalibr	extrem.average	-0.0169	0.0021	V=1,261	0.0855	No diff.
1.5	extrem.average	average	0.0460	0.0291	V=2,174	0.0079	Method.2
1.5	robust.recalibr	extrem.average	-0.0206	0.0130	V=1,334	0.1709	No diff.
2	extrem.average	average	0.0630	0.0391	V=2,213	0.0045	Method.2
2	robust.recalibr	extrem.average	-0.0224	0.0265	V=1,379	0.2487	No diff.
2.5	extrem.average	average	0.0795	0.0492	V=2,234	0.0033	Method.2
2.5	robust.recalibr	extrem.average	-0.0232	0.0281	V=1,414	0.3243	No diff.
3	extrem.average	average	0.0951	0.0594	V=2,249	0.0026	Method.2
3	robust.recalibr	extrem.average	-0.0230	0.0212	V=1,446	0.4053	No diff.

(b) NFL data only

γ	Method.1	Method.2	Avg.diff	Med.diff	Test stat.	p-value	Signif. better?
0.5	extrem.average	average	-0.0067	-0.0129	V=1,557	0.0009	Method.1
0.5	robust.recalibr	extrem.average	-0.0051	-0.0079	V=2,130	0.1750	No diff.
1	extrem.average	average	-0.0098	-0.0254	V=1,627	0.0020	Method.1
1	robust.recalibr	extrem.average	-0.0062	-0.0097	V=2,303	0.4463	No diff.
1.5	extrem.average	average	-0.0106	-0.0373	V=1,699	0.0045	Method.1
1.5	robust.recalibr	extrem.average	-0.0044	-0.0080	V=2,440	0.7714	No diff.
2	extrem.average	average	-0.0102	-0.0452	V=1,772	0.0097	Method.1
2	robust.recalibr	extrem.average	-0.0007	-0.0055	V=2,508	0.9548	No diff.
2.5	extrem.average	average	-0.0089	-0.0531	V=1,849	0.0202	Method.1
2.5	robust.recalibr	extrem.average	0.0042	-0.0034	V=2,571	0.8757	No diff.
3	extrem.average	average	-0.0072	-0.0622	V=1,900	0.0318	Method.1
3	robust.recalibr	extrem.average	0.0098	-0.0020	V=2,604	0.7872	No diff.

(c) Science data only

γ	Method.1	Method.2	Avg.diff	Med.diff	Test stat.	p-value	Signif. better?
0.5	extrem.average	average	-0.0063	-0.0254	V=81,582	<0.0001	Method.1
0.5	robust.recalibr	extrem.average	-0.0264	-0.0050	V=74,929	<0.0001	Method.1
1	extrem.average	average	-0.0045	-0.0377	V=87,242	<0.0001	Method.1
1	robust.recalibr	extrem.average	-0.0461	-0.0024	V=78,104	<0.0001	Method.1
1.5	extrem.average	average	0.0006	-0.0431	V=91,266	<0.0001	Method.1
1.5	robust.recalibr	extrem.average	-0.0608	-0.0007	V=80,416	<0.0001	Method.1
2	extrem.average	average	0.0069	-0.0471	V=94,089	<0.0001	Method.1
2	robust.recalibr	extrem.average	-0.0718	-0.0002	V=82,239	<0.0001	Method.1
2.5	extrem.average	average	0.0134	-0.0489	V=96,155	0.0002	Method.1
2.5	robust.recalibr	extrem.average	-0.0801	-0.0001	V=83,672	<0.0001	Method.1
3	extrem.average	average	0.0195	-0.0510	V=97,698	0.0007	Method.1
3	robust.recalibr	extrem.average	-0.0864	-0.0000	V=84,804	<0.0001	Method.1

(d) States data only

γ	Method.1	Method.2	Avg.diff	Med.diff	Test stat.	p-value	Signif. better?
0.5	extrem.average	average	0.0002	-0.0116	V=584	0.6089	No diff.
0.5	robust.recalibr	extrem.average	-0.0667	-0.0808	V=155	<0.0001	Method.1
1	extrem.average	average	0.0071	-0.0224	V=640	0.9846	No diff.
1	robust.recalibr	extrem.average	-0.1183	-0.1256	V=161	<0.0001	Method.1
1.5	extrem.average	average	0.0170	-0.0276	V=688	0.6293	No diff.
1.5	robust.recalibr	extrem.average	-0.1566	-0.1465	V=171	<0.0001	Method.1
2	extrem.average	average	0.0279	-0.0316	V=708	0.4992	No diff.
2	robust.recalibr	extrem.average	-0.1850	-0.1593	V=187	<0.0001	Method.1
2.5	extrem.average	average	0.0388	-0.0350	V=725	0.401	No diff.
2.5	robust.recalibr	extrem.average	-0.2069	-0.1604	V=192	<0.0001	Method.1
3	extrem.average	average	0.0494	-0.0357	V=741	0.3201	No diff.
3	robust.recalibr	extrem.average	-0.2244	-0.1563	V=196	<0.0001	Method.1

Table F1: Two-sided paired Wilcoxon signed rank tests of Brier scores in each data set. Compares robust recalibration, extremizing away from 0.5 and simple average.

Data set	Degrees of Freedom	Mean Sq. Error	F-stat	p-value
Artwork	9	0.0438	1.097	0.362
NFL	9	0.00388	0.142	0.998
Science	9	0.1919	8.125	< 0.0001
States	9	0.07304	13.99	< 0.0001

Table F2: One-way ANOVA test of Brier scores across 10 methods (four benchmark algorithms and robust recalibration with $\gamma \in \{0.5, 1, 1.5, 2, 2.5, 3\}$) in each data set. Results suggest significant differences in Science and States data.

Method	Benchmark	Avg.diff	Med.diff	Test stat.	p-value	Signif. better?
robust.recalibr. $\gamma=0.5$	know.weight	-0.0001	0.0021	V=247,540	<0.0001	know.weight
robust.recalibr. $\gamma=0.5$	meta.prob.weight	0.0017	-0.0075	V=200,532	0.4623	No difference
robust.recalibr. $\gamma=0.5$	min.pivot	-0.0067	-0.0017	V=121,239	<0.0001	robust.recalibr
robust.recalibr. $\gamma=0.5$	surp.overshoot	0.0032	0.0053	V=246,687	<0.0001	surp.overshoot
robust.recalibr. $\gamma=1$	know.weight	-0.0156	-0.0056	V=123,231	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1$	meta.prob.weight	-0.0138	-0.0238	V=121,218	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1$	min.pivot	-0.0222	-0.0164	V=93,364	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1$	surp.overshoot	-0.0123	-0.0047	V=153,070	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1.5$	know.weight	-0.0230	-0.0150	V=96,184	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1.5$	meta.prob.weight	-0.0212	-0.0363	V=103,043	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1.5$	min.pivot	-0.0296	-0.0257	V=103,024	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1.5$	surp.overshoot	-0.0197	-0.0118	V=123,548	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2$	know.weight	-0.0257	-0.0216	V=102,362	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2$	meta.prob.weight	-0.0239	-0.0467	V=107,335	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2$	min.pivot	-0.0323	-0.0328	V=110,455	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2$	surp.overshoot	-0.0224	-0.0188	V=122,617	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2.5$	know.weight	-0.0256	-0.0240	V=110,829	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2.5$	meta.prob.weight	-0.0238	-0.0550	V=114,400	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2.5$	min.pivot	-0.0322	-0.0383	V=116,401	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2.5$	surp.overshoot	-0.0223	-0.0220	V=125,542	<0.0001	robust.recalibr
robust.recalibr. $\gamma=3$	know.weight	-0.0239	-0.0274	V=118,513	<0.0001	robust.recalibr
robust.recalibr. $\gamma=3$	meta.prob.weight	-0.0221	-0.0588	V=120,723	<0.0001	robust.recalibr
robust.recalibr. $\gamma=3$	min.pivot	-0.0305	-0.0421	V=121,302	<0.0001	robust.recalibr
robust.recalibr. $\gamma=3$	surp.overshoot	-0.0206	-0.0244	V=129,139	<0.0001	robust.recalibr

Table F3: Comparison of Brier scores, two-sided paired Wilcoxon signed rank tests, robust recalibration with $\gamma \in \{0.5, 1, 1.5, 2, 2.5, 3\}$ vs benchmarks.

(a) Artwork data only

Method	Benchmark	Avg.diff	Med.diff	Test stat.	p-value	Signif. better?
robust.recalibr. $\gamma=0.5$	know.weight	-0.0395	-0.0050	V=1,368	0.2277	No difference
robust.recalibr. $\gamma=0.5$	meta.prob.weight	-0.0038	-0.0070	V=1,535	0.6853	No difference
robust.recalibr. $\gamma=0.5$	min.pivot	-0.0046	-0.0011	V=1,281	0.1045	No difference
robust.recalibr. $\gamma=0.5$	surp.overshoot	-0.0162	-0.0010	V=1,413	0.3220	No difference
robust.recalibr. $\gamma=1$	know.weight	-0.0302	-0.0039	V=1,275	0.0985	No difference
robust.recalibr. $\gamma=1$	meta.prob.weight	0.0054	-0.0005	V=1,710	0.6677	No difference
robust.recalibr. $\gamma=1$	min.pivot	0.0047	0.0070	V=1,645	0.9065	No difference
robust.recalibr. $\gamma=1$	surp.overshoot	-0.0069	0.0036	V=1,480	0.5034	No difference
robust.recalibr. $\gamma=1.5$	know.weight	-0.0170	-0.0119	V=1,203	0.0458	robust.recalibr
robust.recalibr. $\gamma=1.5$	meta.prob.weight	0.0186	-0.0124	V=1,731	0.5961	No difference
robust.recalibr. $\gamma=1.5$	min.pivot	0.0178	0.0133	V=1,799	0.3919	No difference
robust.recalibr. $\gamma=1.5$	surp.overshoot	0.0062	-0.0010	V=1,718	0.6400	No difference
robust.recalibr. $\gamma=2$	know.weight	-0.0019	-0.0289	V=1,387	0.2648	No difference
robust.recalibr. $\gamma=2$	meta.prob.weight	0.0337	-0.0051	V=1,845	0.2816	No difference
robust.recalibr. $\gamma=2$	min.pivot	0.0329	0.0198	V=1,928	0.1403	No difference
robust.recalibr. $\gamma=2$	surp.overshoot	0.0214	-0.0070	V=1,926	0.1428	No difference
robust.recalibr. $\gamma=2.5$	know.weight	0.0139	-0.0027	V=1,642	0.9179	No difference
robust.recalibr. $\gamma=2.5$	meta.prob.weight	0.0495	-0.0029	V=1,977	0.0873	No difference
robust.recalibr. $\gamma=2.5$	min.pivot	0.0487	0.0264	V=2,047	0.0408	min.pivot
robust.recalibr. $\gamma=2.5$	surp.overshoot	0.0372	-0.0096	V=2,048	0.0403	robust.recalibr
robust.recalibr. $\gamma=3$	know.weight	0.0296	0.0099	V=1,840	0.2924	No difference
robust.recalibr. $\gamma=3$	meta.prob.weight	0.0652	-0.0104	V=2,106	0.0199	robust.recalibr
robust.recalibr. $\gamma=3$	min.pivot	0.0645	0.0332	V=2,118	0.0170	min.pivot
robust.recalibr. $\gamma=3$	surp.overshoot	0.0529	0.0176	V=2,115	0.0177	surp.overshoot

(b) NFL data only

Method	Benchmark	Avg.diff	Med.diff	Test stat.	p-value	Signif. better?
robust.recalibr. $\gamma=0.5$	know.weight	-0.0005	0.0030	V=3,060	0.0661	No difference
robust.recalibr. $\gamma=0.5$	meta.prob.weight	-0.0014	0.0000	V=2,550	0.9329	No difference
robust.recalibr. $\gamma=0.5$	min.pivot	-0.0011	-0.0004	V=2,222	0.2983	No difference
robust.recalibr. $\gamma=0.5$	surp.overshoot	0.0083	0.0077	V=3,441	0.0016	No difference
robust.recalibr. $\gamma=1$	know.weight	-0.0047	-0.0016	V=2,198	0.2616	No difference
robust.recalibr. $\gamma=1$	meta.prob.weight	-0.0056	-0.0132	V=1,933	0.0420	robust.recalibr
robust.recalibr. $\gamma=1$	min.pivot	-0.0053	-0.0110	V=1,970	0.0566	No difference
robust.recalibr. $\gamma=1$	surp.overshoot	0.0041	0.0003	V=2,673	0.6120	No difference
robust.recalibr. $\gamma=1.5$	know.weight	-0.0037	-0.0105	V=1,981	0.0617	No difference
robust.recalibr. $\gamma=1.5$	meta.prob.weight	-0.0046	-0.0253	V=2,015	0.0798	No difference
robust.recalibr. $\gamma=1.5$	min.pivot	-0.0044	-0.0204	V=2,148	0.1955	No difference
robust.recalibr. $\gamma=1.5$	surp.overshoot	0.0050	-0.0062	V=2,445	0.7846	No difference
robust.recalibr. $\gamma=2$	know.weight	0.0004	-0.0168	V=2,173	0.2268	No difference
robust.recalibr. $\gamma=2$	meta.prob.weight	-0.0004	-0.0402	V=2,210	0.2795	No difference
robust.recalibr. $\gamma=2$	min.pivot	-0.0002	-0.0268	V=2,307	0.4546	No difference
robust.recalibr. $\gamma=2$	surp.overshoot	0.0092	-0.0119	V=2,472	0.8568	No difference
robust.recalibr. $\gamma=2.5$	know.weight	0.0066	-0.0218	V=2,319	0.4798	No difference
robust.recalibr. $\gamma=2.5$	meta.prob.weight	0.0057	-0.0511	V=2,332	0.5080	No difference
robust.recalibr. $\gamma=2.5$	min.pivot	0.0060	-0.0291	V=2,415	0.7065	No difference
robust.recalibr. $\gamma=2.5$	surp.overshoot	0.0153	-0.0158	V=2,518	0.9822	No difference
robust.recalibr. $\gamma=3$	know.weight	0.0139	-0.0250	V=2,454	0.8085	No difference
robust.recalibr. $\gamma=3$	meta.prob.weight	0.0130	-0.0558	V=2,454	0.8085	No difference
robust.recalibr. $\gamma=3$	min.pivot	0.0133	-0.0313	V=2,517	0.9794	No difference
robust.recalibr. $\gamma=3$	surp.overshoot	0.0227	-0.0191	V=2,586	0.8352	No difference

(c) Science data only

Method	Benchmark	Avg.diff	Med.diff	Test stat.	p-value	Signif. better?
robust.recalibr. $\gamma=0.5$	know.weight	0.0005	0.0014	V=135,238	<0.0001	know.weight
robust.recalibr. $\gamma=0.5$	meta.prob.weight	0.0005	-0.0087	V=105,406	0.0577	No difference
robust.recalibr. $\gamma=0.5$	min.pivot	-0.0084	-0.0024	V=55,092	<0.0001	robust.recalibr
robust.recalibr. $\gamma=0.5$	surp.overshoot	0.0017	0.0045	V=133,503	0.0003	surp.overshoot
robust.recalibr. $\gamma=1$	know.weight	-0.0174	-0.0068	V=53,859	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1$	meta.prob.weight	-0.0175	-0.0272	V=57,205	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1$	min.pivot	-0.0264	-0.0166	V=39,850	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1$	surp.overshoot	-0.0163	-0.0058	V=73,182	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1.5$	know.weight	-0.0269	-0.0162	V=43,809	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1.5$	meta.prob.weight	-0.0270	-0.0389	V=47,981	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1.5$	min.pivot	-0.0359	-0.0253	V=43,628	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1.5$	surp.overshoot	-0.0258	-0.0123	V=55,148	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2$	know.weight	-0.0316	-0.0216	V=46,463	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2$	meta.prob.weight	-0.0317	-0.0481	V=48,503	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2$	min.pivot	-0.0406	-0.0327	V=46,822	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2$	surp.overshoot	-0.0305	-0.0192	V=54,264	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2.5$	know.weight	-0.0334	-0.0244	V=49,251	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2.5$	meta.prob.weight	-0.0335	-0.0557	V=50472	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2.5$	min.pivot	-0.0424	-0.0378	V=49,365	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2.5$	surp.overshoot	-0.0323	-0.0225	V=55,183	<0.0001	robust.recalibr
robust.recalibr. $\gamma=3$	know.weight	-0.0336	-0.0278	V=51,837	<0.0001	robust.recalibr
robust.recalibr. $\gamma=3$	meta.prob.weight	-0.0337	-0.0576	V=52,322	<0.0001	robust.recalibr
robust.recalibr. $\gamma=3$	min.pivot	-0.0426	-0.0416	V=51,598	<0.0001	robust.recalibr
robust.recalibr. $\gamma=3$	surp.overshoot	-0.0325	-0.0254	V=56,356	<0.0001	robust.recalibr

(d) States data only

Method	Benchmark	Avg.diff	Med.diff	Test stat.	p-value	Signif. better?
robust.recalibr. $\gamma=0.5$	know.weight	0.0551	0.0463	V=1,246	<0.0001	know.weight
robust.recalibr. $\gamma=0.5$	meta.prob.weight	0.0337	0.0322	V=932	0.0045	meta.prob.weight
robust.recalibr. $\gamma=0.5$	min.pivot	0.0019	0.0008	V=798	0.1225	No difference
robust.recalibr. $\gamma=0.5$	surp.overshoot	0.0448	0.0210	V=1,167	<0.0001	surp.overshoot
robust.recalibr. $\gamma=1$	know.weight	0.0104	0.0039	V=911	0.0084	know.weight
robust.recalibr. $\gamma=1$	meta.prob.weight	-0.0110	-0.0182	V=417	0.0337	robust.recalibr
robust.recalibr. $\gamma=1$	min.pivot	-0.0429	-0.0537	V=44	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1$	surp.overshoot	0.0001	0.0071	V=696	0.5756	No difference
robust.recalibr. $\gamma=1.5$	know.weight	-0.0180	-0.0124	V=273	0.0004	robust.recalibr
robust.recalibr. $\gamma=1.5$	meta.prob.weight	-0.0394	-0.0419	V=84	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1.5$	min.pivot	-0.0712	-0.0868	V=46	<0.0001	robust.recalibr
robust.recalibr. $\gamma=1.5$	surp.overshoot	-0.0283	-0.0132	V=318	0.0021	robust.recalibr
robust.recalibr. $\gamma=2$	know.weight	-0.0356	-0.0272	V=138	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2$	meta.prob.weight	-0.0570	-0.0590	V=4	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2$	min.pivot	-0.0889	-0.1092	V=51	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2$	surp.overshoot	-0.0459	-0.0220	V=178	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2.5$	know.weight	-0.0465	-0.0327	V=106	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2.5$	meta.prob.weight	-0.0679	-0.0675	V=1	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2.5$	min.pivot	-0.0998	-0.1152	V=52	<0.0001	robust.recalibr
robust.recalibr. $\gamma=2.5$	surp.overshoot	-0.0569	-0.0295	V=146	<0.0001	robust.recalibr
robust.recalibr. $\gamma=3$	know.weight	-0.0533	-0.0361	V=99	<0.0001	robust.recalibr
robust.recalibr. $\gamma=3$	meta.prob.weight	-0.0748	-0.0740	V=7	<0.0001	robust.recalibr
robust.recalibr. $\gamma=3$	min.pivot	-0.1066	-0.1174	V=58	<0.0001	robust.recalibr
robust.recalibr. $\gamma=3$	surp.overshoot	-0.0637	-0.0351	V=138	<0.0001	robust.recalibr

Table F4: Comparison of Brier scores, two-sided paired Wilcoxon signed rank tests, robust recalibration vs benchmarks in each data set.