

1010

S1 File

1011 **6.1 Preliminary notation**

1012 Let N denote the total number of individuals (experts or participants), indexed by $i =$
1013 $1, \dots, N$ and C denote the total number of events or claims (hereafter, “claim”) to be
1014 assessed, usually indexed by $c = 1, \dots, C$. Each claim has outcome 1 if the claim is
1015 true, and 0 otherwise. For each claim c , each individual i provides assessments that

1016 the claim in question is true or false, by estimating three probabilities: $L_{i,c}$, which is a
 1017 lower bound ; $U_{i,c}$, an upper bound and $B_{i,c}$ which corresponds to the best estimate for
 1018 the probability given by individual i for claim c . These estimates satisfy the following
 1019 inequalities: $0 \leq L_{i,c} \leq B_{i,c} \leq U_{i,c} \leq 1$.

1020 Each claim is assessed by more than one individual and we aggregate their probabilities
 1021 to obtain a group probability, denoted \hat{p}_c . We will further denote $\hat{p}_c(\text{Method ID})$ as the
 1022 aggregated probability calculated using the aggregation method with a given ID . For
 1023 example, the first simple average (arithmetic mean) aggregation for claim c is:

$$\hat{p}_c(\text{ArMean}) = \frac{1}{N} \sum_{i=1}^N B_{i,c} \quad (5)$$

1024 6.2 Weights

1025 Given that many of the aggregation methods proposed involve weighted linear combina-
 1026 tions of individual assessments, we can define some standard notation to enhance clarity.

1027 We denote the unnormalized weights by w_{method} (with subscripts denoting corre-
 1028 sponding individuals or claims) and the normalised versions by $\tilde{w}_{\text{method}}$. All weights
 1029 need to be normalised (i.e., to sum to one), but as the process is the same for all of
 1030 them, we will give the formulae for the unnormalized weights. All differentially weighted
 1031 combinations will take the form:

$$\hat{p}_c(\text{Method ID}) = \sum_{i=1}^N \tilde{w}_{\text{method},i,c} B_{i,c} \quad (6)$$

1032 We note that while most weights will be calculated on a per claim, per individual basis
 1033 (i.e., judgements from the same individual may be weighted differently on any given
 1034 claim), in three cases, the weights will be calculated across all claims on a per individual
 1035 basis only. In these cases, weights for a given individual will not vary across claims and
 1036 the weights' subscript c from the right hand side of Equation (6) will be dropped.

1037 **6.3 Aggregation Methods**

1038 Method IDs will simply be abbreviations of the mathematical operations used to calculate
1039 the weights.

1040 **6.3.1 ArMean: Arithmetic mean of the best estimates**

1041 The simplest way to aggregate group estimates is to take the unweighted linear average
1042 (i.e., simply takes the average of the best estimates $B_{i,c}$ for each claim). As defined above,
1043 the aggregate estimate for claim c is therefore calculated using Equation (5).

1044 **6.3.2 Median: Median of the best estimates**

1045 Another approach that is often used due to its simplicity is to take the median of the
1046 individuals' best estimates.

$$\hat{p}_c(\text{Median}) = \text{Median} \{B_{i,c}\}_{i=1,\dots,N} \quad (7)$$

1047 **6.3.3 LOArMean: Arithmetic mean of the log odds transformed best esti-** 1048 **mates**

1049 Log odds are often used to model probabilities in generalised linear models and state
1050 estimation algorithms, typically due to the advantages of mapping probabilities onto a
1051 scale where very small values are still differentiable.

$$\text{LogOdds}_c = \frac{1}{N} \sum_{i=1}^N \log \left(\frac{B_{i,c}}{1 - B_{i,c}} \right)$$

1052 The average log odds estimate is then back transformed to give a final group estimate:

$$\hat{p}_c(\text{LOArMean}) = \frac{e^{\text{LogOdds}_c}}{1 + e^{\text{LogOdds}_c}} \quad (8)$$

1053 **6.3.4 BetaArMean: A beta-transformed arithmetic mean**

1054 This method takes the average of best estimates and transforms it using the cumulative
1055 distribution function of a beta distribution. This transformation makes the average more

1056 extreme, i.e. increases values larger than 0.5 and decreases values less than 0.5. The Beta
 1057 distribution is parameterised by two parameters α and β , and in this analysis, we chose
 1058 $\alpha = \beta$ and larger than one.

$$\hat{p}_c(\text{BetaArMean}) = H_\alpha \left(\frac{1}{N} \sum_{i=1}^N B_{i,c} \right) \quad (9)$$

1059 where H_α is the cumulative distribution function of the Beta distribution with two equal
 1060 parameters.

1061 6.3.5 DistribArMean: Arithmetic mean of the non-parametric distributions

1062 This method assumes that the elicited probabilities and bounds can be considered to
 1063 represent participants' subjective distributions associated with relative frequencies (rather
 1064 than unique events). That is to say that we considered that the lower bound of the
 1065 individual per claim corresponds to the 5% percentile of their subjective distribution on
 1066 the probability of replication, denoted $q_{5,i}$, the best estimate corresponds to the median
 1067 $q_{50,i}$, and the upper bound corresponds to the 95% percentile, $q_{95,i}$. With these three
 1068 percentiles, we can build the minimally informative non-parametric distribution that
 1069 spreads the mass uniformly between the three percentiles, such that the constructed
 1070 distribution agrees with participant's assessments and makes no extra assumptions. This
 1071 approach is inspired by methods for eliciting, constructing and aggregating quantities,
 1072 rather than probabilities [1].

$$F_i(x) = \begin{cases} 0, & \text{for } x < 0 \\ \frac{0.05}{q_{5,i}} \cdot x, & \text{for } 0 \leq x < q_{5,i} \\ \frac{0.45}{q_{50,i} - q_{5,i}} \cdot (x - q_{5,i}) + 0.05, & \text{for } q_{5,i} \leq x < q_{50,i} \\ \frac{0.45}{q_{95,i} - q_{50,i}} \cdot (x - q_{50,i}) + 0.5, & \text{for } q_{50,i} \leq x < q_{95,i} \\ \frac{0.05}{1 - q_{95,i}} \cdot (x - q_{95,i}) + 0.95, & \text{for } q_{95,i} \leq x < 1 \\ 1, & \text{for } x \geq 1. \end{cases}$$

We then average all such constructed distributions of participants for each claim

$$AvDistribution = \frac{1}{N} \sum_{i=1}^N F_i(x);$$

1073 and the aggregation is the median of the average distribution

$$\hat{p}_c(DistribArMean) = AvDistribution^{-1}(0.5) \quad (10)$$

1074 **6.3.6 IntWAgg: Weighted by interval width**

1075 The width of the interval provided by individuals may be an indicator of certainty, and
 1076 arguably of accuracy of the best estimate contained between the bounds of the interval. We
 1077 weight according to the interval width across individuals for that claim, defined as follows:

$$w_Interval_{i,c} = \frac{1}{U_{i,c} - L_{i,c}}$$

$$\hat{p}_c(IntWAgg) = \sum_{i=1}^N \tilde{w}_Interval_{i,c} B_{i,c} \quad (11)$$

1078 **6.3.7 IndIntWAgg: Weighted by the rescaled interval width (interval width** 1079 **relative to largest interval width provided by that individual)**

1080 Because of the variability in the widths of intervals participants give for different claims,
 1081 we may need to re-scale interval widths across all claims per individual. This results in a
 1082 re-scaled interval width weight, for individual i for claim c , relative to the widest interval
 1083 provided by that individual across all claims C :

$$w_nIndivInterval_{i,c} = \frac{1}{\frac{U_{i,c} - L_{i,c}}{\max(\{(U_{i,d} - L_{i,d}) : d=1, \dots, C\})}}$$

1084 where $U_{i,d} - L_{i,d}$ are individual i 's judgements for claim d . Then

$$\hat{p}_c(IndIntWAgg) = \sum_{i=1}^N \tilde{w}_nIndivInterval_{i,c} B_{i,c} \quad (12)$$

1085 **6.3.8 VarIndIntWAgg: Weighted by variation in individuals' interval widths**

1086 A related issue is that participants differ in how much they vary in their interval widths. A
 1087 higher variance may indicate a higher responsiveness to the existing supporting evidence
 1088 to different claims. Such responsiveness might be predictive of more accurate assessors.
 1089 We define:

$$w_varIndivInterval_i = var \{(U_{i,c} - L_{i,c}) : c = 1, \dots, C\}$$

1090 where the variance (*var*) is calculated across all claims for individual *i*. Then

$$\hat{p}_c (VarIndIntWAgg) = \sum_{i=1}^N \tilde{w}_varIndivInterval_i B_{i,c} \quad (13)$$

1091 **6.3.9 AsymWAgg: Weighted by asymmetry of intervals**

1092 We use the asymmetry of an interval relative to the corresponding best estimate to define
 1093 the following weights:

$$w_asym_{i,c} = \begin{cases} 1 - 2 \cdot \frac{U_{i,c} - B_{i,c}}{U_{i,c} - L_{i,c}}, & \text{for } B_{i,c} \geq \frac{U_{i,c} - L_{i,c}}{2} + L_{i,c} \\ 1 - 2 \cdot \frac{B_{i,c} - L_{i,c}}{U_{i,c} - L_{i,c}}, & \text{otherwise} \end{cases}$$

1094 Then

$$\hat{p}_c (AsymAg) = \sum_{i=1}^N \tilde{w}_asym_{i,c} B_{i,c} \quad (14)$$

1095 **6.3.10 IndIntAsymWAgg: Weighted by individuals' interval widths and asym-**
 1096 **metry**

1097 Assuming that we want to reward both asymmetric and narrow intervals, we would
 1098 need to formulate a weight that combines the weights calculated in the *AsymWAgg* and
 1099 *IndIntWAgg* methods. One simple way of achieving this is to multiply the previously
 1100 defined and normalised weights.

$$w_nIndivInterval_asym_{i,c} = \tilde{w}_nIndivInterval_{i,c} \cdot \tilde{w}_asym_{i,c}$$

$$\hat{p}_c(IndIntAsymWAg) = \sum_{i=1}^N \tilde{w}_nIndivInterval_asym_{i,c} B_{i,c} \quad (15)$$

1101 **6.3.11 KitchSinkWAgg: Weighted by everything but the kitchen sink**

1102 KitchSinkWAgg is an ad-hoc method developed and refined using a single dataset (later
 1103 used in the analysis as well). This method is informed by the intuition that we want
 1104 to reward narrow and asymmetric intervals, as well as variability between individuals'
 1105 interval widths (across their estimates).

$$w_kitchSink_{i,c} = \tilde{w}_nIndivInterval_{i,c} \cdot \tilde{w}_asym_{i,c} \cdot \tilde{w}_varIndivInterval_i$$

$$\hat{p}_c(KitchSinkWAg) = \sum_{i=1}^N \tilde{w}_kitchSink_{i,c} B_{i,c} \quad (16)$$

1106 **6.3.12 DistLimitWAgg: Weighted by the distance of the best estimate from**
 1107 **the closest certainty limit**

1108 We give greater weight to best estimates that are closer to certainty limits, as follows

$$w_distLimit_{i,c} = \max\{B_{i,c}, 1 - B_{i,c}\}$$

$$\hat{p}_c(DistLimitWAgg) = \sum_{i=1}^N \tilde{w}_distLimit_{i,c} B_{i,c} \quad (17)$$

1109 **6.3.13 ShiftWAgg: Weighted by judgments that shifted the most after dis-**
 1110 **ussion**

1111 When judgements are elicited using the IDEA protocol (or any other protocol which
 1112 allows experts to revisit their original estimates), the second round of estimates may
 1113 differ from the original first set of estimates an expert provides. Greater changes between

1114 rounds will be given greater weight, with more emphasis on changes in the best estimate
 1115 such that

$$w_shift_{i,c} = |B1_{i,c} - B_{i,c}| + \frac{|L1_{i,c} - L_{i,c}| + |U1_{i,c} - U_{i,c}|}{2}$$

$$\hat{p}_c(ShiftWAgg) = \sum_{i=1}^N \tilde{w}_shift_{i,c} B_{i,c} \quad (18)$$

1116 where $L1_{i,c}, B1_{i,c}, U1_{i,c}$ are the first round lower, best and upper estimates (prior to
 1117 discussion) and $L_{i,c}, B_{i,c}, U_{i,c}$ are the individual's revised second round estimates (after
 1118 discussion).

1119 **6.3.14 GranWAgg: Weighted by the granularity of best estimates**

1120 More skilled forecasters might be expected to have a finer grained appreciation of uncer-
 1121 tainty and thus make more granular forecasts.

1122 In our weighting scheme, individuals' received a score of one for each claim that their
 1123 best estimate was specified at a more granular level than 0.05 (i.e., not a multiple of
 1124 0.05), and a zero otherwise. The mean of scores per claim forms a weight per individual,
 1125 such that

$$w_gran_i = \frac{1}{C} \sum_{d=1}^C \left[\frac{B_{i,d}}{0.05} - \left\lfloor \frac{B_{i,d}}{0.05} \right\rfloor \right],$$

1126 where $\lfloor \cdot \rfloor$ and $\lceil \cdot \rceil$ are the mathematical floor and ceiling functions respectively.

$$\hat{p}_c(GranWAgg) = \sum_{i=1}^N \tilde{w}_gran_i B_{i,c} \quad (19)$$

1127 **6.3.15 EngWAgg: Weighted by the level of engagement as measured by the** 1128 **individuals' verbosity**

1129 When assessing claims, individuals have the chance to comment and engage in discussion
 1130 with other participants. We consider giving greater weight to best estimates that are
 1131 accompanied by a greater number of comments/justifications. We will consider $w_eng_{i,c}$

1132 to be the total number of words used by individual i in comments about their estimates
1133 for claim c .

$$\hat{p}_c(EngWAgg) = \sum_{i=1}^N \tilde{w}_{eng_{i,c}} B_{i,c} \quad (20)$$

1134 **6.3.16 ReasonWAgg: Weighted by the breadth of reasoning provided to sup-**
1135 **port the individuals' estimate**

1136 When individuals provide multiple unique reasons in support of their judgment, this may
1137 indicate a breadth of thinking, understanding and knowledge about the claim and its
1138 context, and may also reflect engagement and conscientiousness. Giving greater weight
1139 to best estimates that are accompanied by a greater number of supporting reasons may
1140 be beneficial. We will consider $w_{reason_{i,c}}$ to be the number of unique reasons provided
1141 by that individual i in support of their estimate for claim c .

$$\hat{p}_c(ReasonWAgg) = \sum_{i=1}^N \tilde{w}_{reason_{i,c}} B_{i,c} \quad (21)$$

1142 Qualitative statements made by individuals as they evaluate claims/studies were coded
1143 by the [repliCATS Reasoning team](#), according to a detailed coding manual developed to
1144 ensure analysts were each coding for common units of meaning in the same sets of textual
1145 data. This manual emerged through an iterative process, that included calculating the
1146 inter-coder-reliability (ICR), in the form of Krippendorff's alpha [2]. Roughly, ICR
1147 measures the extent to which different judges assign similar ratings to the evaluated
1148 characteristics, here in the form of reasoning categories. For context, an ICR (here
1149 Krippendorff's alpha) of 1 indicates perfect reliability, while 0 indicates the absence of
1150 reliability. Values less than 0 indicate systematic disagreements. From this manual, a
1151 subset of 25 codes were selected as reasoning categories, each of which were included
1152 in ReasonWAgg if the ICR was calculated at a minimum of 0.66 across two or more
1153 analysts, or an ICR between two analysts of at least 0.75 and a minimum overall ICR
1154 of 0.50. A quarter of the dataset was manually coded into these categories by multiple
1155 analysts (using the [NVivo Qualitative Data Analysis Software, Version 12, 2018](#)), and

1156 these datasets provided the training data for the remaining text to be auto-coded in
 1157 NVivo. Reasoning scores were calculated for individuals who received one point for each
 1158 of these 25 reasoning categories that they drew on over the course of their evaluation (i.e.,
 1159 statements from both IDEA rounds). The reasoning categories include: the plausibility
 1160 of claim, effect size, sample size, presence of a power analysis, transparency of reporting,
 1161 journal reputation.

1162 ReasonWAgg can be modified to incorporate not only the number of reasons, but also
 1163 their diversity across claims. This modified aggregation will be called ReasonWAgg2.
 1164 The latter component of this score will be calculated per individual from all the claims
 1165 they assessed, so it will be the same for each of the claims assessed by that individual.
 1166 We assume each individual answers at least two claims. If a participant has assessed
 1167 only one claim, for that claim we will default to the original ReasonWAgg.

1168 Table 3 shows a hypothetical example of the reasons used by one participant when
 1169 assessing four claims.

Table 3: The distribution of the reasons one participant mentioned in the comments they made when assessing four claims. A $(Claim, R)$ cell is 1 if the R was used to justify answers for $Claim$, and empty if R was not mentioned.

<i>Claims/Reasons</i>	R_1	R_2	R_3	\dots	R_{25}	Weighted “No. of Reasons”
$Claim_1$	1			1		$0.75 \cdot 1 + 1 \cdot 0.5 = 1.25$
$Claim_2$		1				0.75
$Claim_3$			1			0.5
$Claim_4$			1	1		1
Av. use of R_r	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{2}{4}$	$\frac{2}{4}$	$\frac{0}{4}$	
1 - Av. use of R_r	0.75	0.75	0.5	0.5	0	

1170 The penultimate row of the matrix showed in Table 3 gives the average use of reasons,
 1171 and the last row shows weights assigned to these reasons. These weights are then used
 1172 to calculate a final reasoning score per claim (per participant). This score is showed in
 1173 the last column of Table 3 and it is calculated as a weighted sum of the elements of the
 1174 vector of zeros and ones indicating use of reasons per claim.

1175 We will consider $w_varReason_{i,c}$ to be the weighted “number of unique reasons”
 1176 provided by participant i in support of their estimate for claim c . Assume there are

1177 25 unique reasons any participant can use to justify their numerical answers. Then, for
 1178 each participant i we can construct a matrix \mathbf{CR}_i with 25 columns, each corresponding
 1179 to a unique reason, r , and C rows, where C is the number of claims assessed by that
 1180 participant. Each element of this matrix $\mathbf{CR}_i(r, c)$ can be either 1 or 0. $\mathbf{CR}_i(r, c) = 1$
 1181 if reason R_r was used to justify the estimates assessed for claim c , and $\mathbf{CR}_i(r, c) = 0$ if
 1182 reason R_r was not mentioned when assessing claim c .

$$w_varReason_{i,c} = \sum_{r=1}^{25} \mathbf{CR}_i(c, r) \cdot \left(1 - \frac{\sum_{c=1}^C \mathbf{CR}_i(c, r)}{C} \right)$$

1183 ReasonWAgg2 will use the weights calculated as above in the prediction given by Equation
 1184 21.

1185 6.3.17 QuizWAgg: Weighted by performance on the quiz

1186 As part of the repliCATS project, individuals were asked to take a quiz before commencing
 1187 the main task of evaluating research claims. Hence, this aggregation method will
 1188 only apply to a dataset where such an exercise is undertaken prior to the elicitation. The
 1189 quiz aimed to gauge subject matter expertise, and in the case of the repliCATS project,
 1190 consisted of questions testing familiarity with previous research and concepts related to
 1191 assessments of replicability, e.g., understanding of statistical concepts, (false) positive
 1192 rates and replication rates in domain-relevant literature, and self-reported rates of ques-
 1193 tionable research practices. If answered, quiz responses would provide similar information
 1194 to that of seed questions, enabling differential performance-based weighted combinations
 1195 (giving greater weight to individuals with higher quiz scores). The quiz was encouraged,
 1196 but not compulsory, so choosing to take the quiz at all may also reflect engagement and
 1197 conscientiousness.

1198 The quiz contains $n_{quiz} = 22$ questions, 12 of which cover knowledge and understand-
 1199 ing of statistical concepts, and 10 are about meta-research. Questions that required less
 1200 effort to answer (i.e., 10 true/false questions and one two-part question) were assigned
 1201 half points.

Individuals provide answers for each question, resulting in a $N \times n_{quiz}$ matrix \mathbf{Q} , where

each element $\mathbf{Q}(i, h)$ is 1 if individual i answered question h correctly, and 0 otherwise. For each question answered correctly, the individual receives points, with the number of points received for a correct answer for each of the 22 questions specified in the points vector

$$\mathbf{v} = \begin{cases} 0.5, & \text{for questions 1 to 10, 16, 17} \\ 1, & \text{for questions 11 to 15, and 18 to 22} \end{cases}$$

1202 This results in a quiz score that ranges from 0 to 16, with higher scores indicating better
1203 performance. Then the un-normalised weight based on the quiz score is

$$w_quiz_i = \mathbf{Q} \cdot \mathbf{v}$$

1204 and the aggregated estimate is

$$\hat{p}_c(QuizWAgg) = \sum_{i=1}^N \tilde{w}_quiz_i B_{i,c} \quad (22)$$

1205 where w_quiz_i is the weight corresponding to the score of individual i on the quiz, as
1206 defined above.

1207 In the case of unanswered questions (missing data), individuals are assigned zero
1208 points for that question. Individuals who did not take the quiz at all will receive zero
1209 weight (and non-zero weight for those who responded to at least one item in the quiz).
1210 If only one person assessing a given claim took the quiz, the *QuizWAgg* aggregated
1211 estimate for that claim will be based solely on their judgment. If, however, nobody took
1212 the quiz, this aggregation method is impossible to construct.

1213 **6.3.18 CompWAgg: Weighted by the level of self-rated comprehension of** 1214 **the claim the individuals' report**

1215 In the repliCATS project, before assessing a claim, individuals were asked to assess how
1216 well they understood it. A 7-point scale, where 1 corresponds to “I have no idea what it
1217 means” and 7 corresponds to “It is perfectly clear to me” is used for this comprehensibility
1218 question. Intuitively, the numerical estimates of the individuals who are confident they

1219 understood the claim may be weighted more. We will consider $w_comp_{i,c}$ to be the number
 1220 assigned to the comprehensibility, as provided by individual i in support of their estimate
 1221 for claim c .

$$\hat{p}_c(CompWAgg) = \sum_{i=1}^N \tilde{w}_comp_{i,c} B_{i,c} \quad (23)$$

1222 **6.3.19 BayTriVar: Bayesian Triple-Variability Method**

1223 The last two aggregation methods proposed are Bayesian methods, and hence they use
 1224 the elicited probabilities differently, namely as data with which prior distributions are
 1225 updated.

Three kinds of variability around best estimates are considered: generic claim variability, generic participant variability, and claim - participant specific uncertainty (operationalised by bounds). The model takes the log odds transformed individual best estimates as input (data), uses a normal likelihood function and derives a posterior distribution for the probability of replication. That is to say, the log odds transformed best estimate data are assumed to follow a Normal distribution $\log\left(\frac{B_{i,c}}{1-B_{i,c}}\right) \sim N(\mu_c, \sigma_{i,c})$, where μ_c denotes the mean estimated probability of replication for claim c , and $\sigma_{i,c}$ denotes the standard deviation of the estimated probability of replication for claim c and individual i . Parameter $\sigma_{i,c}$ is calculated as:

$$\sigma_{i,c} = (U_{i,c} - L_{i,c} + 0.01) \times \sqrt{\sigma_i^2 + \sigma_c^2}$$

1226 with σ_i denoting the standard deviation of estimated probabilities of replication for indi-
 1227 vidual i and σ_c denoting the standard deviation of the estimated probability of replication
 1228 for claim c . The above formula for the standard deviation is derived using the statistical
 1229 rules for calculating the variances of a sum of two independent random variables. The
 1230 distribution of the best estimates is considered to be the convolution of the claim and
 1231 participant distributions (thought of as independent). The sum of these two variables is
 1232 then scaled by a constant (the width of an interval for a particular claim) which represents
 1233 the claim - participant specific uncertainty. The variance then is the scaled addition of

1234 the two variances.

1235 To complete the specification of the Bayesian model, priors need to be given for μ_c ,
1236 σ_i , and σ_c . These are defined as $\mu_c \sim N(0, 3)$, $\sigma_i \sim U(0, 10)$ and $\sigma_c \sim U(0, 10)$, with
1237 $U(0, 10)$ denoting the Uniform distribution on the interval from 0 to 10. The quantity of
1238 interest is the median of the posterior distribution of μ_c , the mean estimated probability
1239 of replication. In Bayesian statistics the posterior distribution is proportional to the
1240 product of the likelihood and the prior and in this instance a Monte Carlo Markov Chain
1241 algorithm [3] is used to sample from this posterior distribution. After obtaining the
1242 median of the posterior distribution of μ_c , we can back transform to obtain \hat{p}_c :

$$\hat{p}_c(\text{BayTriVar}) = \frac{e^{\mu_c}}{1 + e^{\mu_c}} \quad (24)$$

1243 **6.3.20 BayPRIORsAgg: Prior derived from predictive models, updated with** 1244 **best estimates**

1245 This BayPRIORsAgg method uses Bayesian updating to update a prior probability of
1246 replication estimated from a predictive model with an aggregate of the experts' best
1247 estimates for any given claim. The main difference between this method and the one
1248 presented in Section 6.3.19 is that the parameters of the prior distribution of μ_c are
1249 informed by the PRIORS model [4] which is a multilevel logistic regression model that
1250 predicts the probability of replication using attributes of the original study.

References

- [1] Cooke RM. Experts in uncertainty: Opinion and subjective probability in science. Environmental Ethics and Science Policy Series. Oxford University Press; 1991.
- [2] Krippendorff K. Content Analysis: An Introduction to Its Methodology, 4th Edition. Los Angeles: SAGE Publications; 2019.
- [3] Ravenzwaaij vD, Cassey P, Brown SD. A simple introduction to Markov Chain Monte-Carlo sampling. Psychon Bull Rev. 2018;25:143–154.
- [4] Gould E, Wilkinson DP, Willcox A, Groenewegen R, Vesik P, Fraser H, et al. Using model-based predictions to inform the mathematical aggregation of human-based predictions of replicability. MetaArXiv Preprints. 2021; Available from: <https://doi.org/10.31222/osf.io/f675q>.