

တို့ရ

Comparing Top Forecasters to Domain Experts

Reviewing the evidence

An Arb Research Report March 2022

The hypotheses

The superforecasting phenomenon — that certain teams of forecasters are better than other prediction mechanisms (large crowds, simple statistical prediction rules, etc) seems sound. But serious interest in superforecasting stems from the <u>reported triumph</u> of forecaster *generalists* over non-forecaster *experts*. (Another version says that they also outperform analysts with classified information.) So distinguish some claims:

- 1) Forecasters > public
- 2) Forecasters > <u>simple models</u>
- 3) Forecasters > experts
 - a) Forecasters > experts with classified info
 - b) Averaged forecasters > experts
 - c) <u>Aggregated</u> forecasters > experts

Is (3) true? We review studies comparing the performance of trained forecasters and experts in some domain. We also tried to cover prediction markets vs experts.

Summary

We are fairly pessimistic (but uncertain):

- We think that claim (1) is true with 99% confidence¹ and that claim (2) is true with 95% confidence. But surprisingly few studies compare experts to generalists (3). The analysis quality and transparency leave much to be desired. <u>The best study</u> found that forecasters and health professionals performed similarly. In other studies, experts had goals besides accuracy, or there were too few to aggregate well.
- 3a) There's a common misconception that superforecasters outperformed intelligence analysts by "30%". Instead: <u>Goldstein et al</u> showed that [EDIT: the Good Judgment Project's best-performing aggregation method]² outperformed the intelligence community, but this was partly due to the different aggregation technique used (the GJP weighting algorithm seems to perform better than prediction markets, given the apparently low volumes of the ICPM market). The forecaster *prediction market* performed about as well as the intelligence analyst prediction market; in general, prediction pools outperform prediction markets under current conditions (e.g. subsidies, volume, incentives, demographics). [85%]

¹ This is almost a trivial claim, since forecasters are by definition more interested in current affairs than average, and much more interested in epistemics than average. So we'd select for the subset of "the public" who should outperform simply through increased effort, even if all the practice and formal calibration training did nothing, which it probably doesn't.

² Previously this section said "superforecasters"; after discussion, it seems more prudent to say "the Good Judgment Project's best-performing aggregation method". See this <u>comment</u> for details.

- 3b) In the same study, the forecaster average was notably *worse* than the intelligence community.
- 3c) Ideally, we would pit a crowd of forecasters against a crowd of experts. Only an unpublished extension of Sell et al. manages this; it found a small (~3%) forecaster advantage.
- The bar may be low. Expertise, plus basic forecasting training and active willingness to forecast regularly, were enough to be on par with the best forecasters. [33%]³
- In more complex domains, like ML, there could be significant returns to expertise. So it might be better to broaden focus from *generalist forecasters* to *competent ML pros* who are excited about forecasting. [40%]

³ Our exact probability hinges on what's considered low and on how good trained Hypermind forecasters are. This is less obvious than it seems: in a <u>similar tournament</u>, the CSET Foretell Top Forecasters were not uniformly good; a team including Misha finished with a 4x better relative Brier score than the "top forecasters team." Further, our priors are mixed: (a) common sense makes us favor experts, (b) but common sense also somewhat favors expert forecasters, (c) Tetlock's work on expert political judgment pushes us away from politics experts, and finally (d) we have first-hand experience about superforecasting being real.

Table of studies

	Comparison	Result			Notes	
Geopolitics						
<u>Goldstein</u> <u>et al</u> (2015)	US Intelligence Community Prediction Market (ICPM) Good Judgement Project (GJP): an average, vs prediction market (PM), vs best method (selected post hoc) ⁴ . Participants are rewarded for accuracy. ICPM was low stakes: play-money ⁵ , while GJP participants "were paid a small honorarium for their active participation".	Equal performant best aggregation method was sele methods perform <i>Mean of means of</i> Goldstein et al (2015) ICPM GJP (avg) GJP (PM)	ce for expe method wi octed post h ned within 2 of daily Brie MMBD .23 .32 .21	ert and forecaste as notably bette noc among 20, b 2% of the best. er scores ⁶ (MMB 95% CI (.19, .27) (.29, .35)*** (.17, .26)	er prediction markets. The r than ICPM. The best but several of the other	Unpublished document used to justify the famous "Supers are 30% better than the CIA" claim. The most direct comparison between forecasters (GJP PM) and experts (ICPM) finds similar performance (insignificant diff). Prediction markets seem worse than super-aggregating opinion pools (see <u>Appendix</u> <u>A</u>); this study itself shows a large gap between GJP (PM) and GJP (best).

⁴ All Surveys Logit "takes the most recent forecasts from a selection of individuals in GJP's survey elicitation condition, weights them based on a forecaster's historical accuracy, expertise, and psychometric profile, and then extremizes the aggregate forecast (towards 1 or 0) using an optimized extremization coefficient." Note that this method was selected *post hoc*, which raises the question of multiple comparisons; the authors respond that "several other GJP methods were of similar accuracy (<2% difference in accuracy)."

⁵ There is some inconclusive research comparing real- and play-money: <u>Servan-Schreiber et al. (2004)</u> find no significant difference for predicting NFL (American football); <u>Rosenbloom & Notz (2006)</u> find that in non-sports events, real-money markets are more accurate and that they are comparably accurate for sports markets; and <u>Slamka et al. (2008)</u> finds real- and play-money prediction markets comparable for UEFA (soccer).

⁶ MMBD is not a proper scoring rule (one incentivizing truthful reporting). If a question has a chance of resolving early (e.g., all questions of the form "will X occur by date?"), the rule incentivizes forecasters to report higher probabilities for such outcomes. This could have affected GJP (avg and best) predictors, who were rewarded for it; but should have not affected ICPM and GJP (PM), as these used the Logarithmic Market Scoring Rule. See <u>Sempere & Lawsen</u> (2021) for details.

	N=139 geopolitical questions	GJP (be	st) .	15	(.10, .	21)***						
		<u>Mean Per</u>	<u>centage</u>	of Days L	Directio	nally Accura	ate (MPDDA)					
			Ν	IPDDA	95% 0	CI						
		ICPM	8	1.58	(76.63	8, 86.54)						
		GJP (av	g) 7	9.74	(73.82	2, 85.66)						
		GJP (PM	1) 8	3.45	(78.83	8, 88.06)						
		GJP (be	st) 8	8.20	(83.87	7, 92.5)***						
		** p < .00	1 vs ICPI	M.	-							
<u>Stastny &</u>	Qualitative forecasts from intelligence reports. Seasoned professional analysts produced ⁷ :	Mean abs Moreover better (p<	Mean absolute error of ICPM was better (p <.001) than in the reports. Moreover, the initial forecasts by seasoned intelligence analysts were better (p <.05) than the forecasts imputed by them from the reports. Note that Initial forecasts were almost as good as ICPM forecasts .						Mandel (2019) critiques the study. Table 1 is illuminating:			
<u>(2018)</u>	 initial personal probabilities, 	Mean abs	Mean absolute error						all	non-f.		
	 probabilities imputed in the reports 		All q.	Non-fuz	zy q.	Fuzzy q.			Initial personal	.194	.200	
	 imputed⁸ probabilities 	Initial	0.317	_	-		—		Imputed	.252	.254	
	events,	Imputed	0.416	0.412		0.427			Imputed upd.	.238	.243	
	 updated personal 	ICPM	0.302	0.305		0.3						

⁷ Note that our understanding is that these were not averaged. On average there have been ~2.5 imputation predictions per report.

⁸ Unclear if imputers did a reasonable job of separating their personal views from their imputations. <u>Mandel (2019)</u> notes that Person correlation between mean Brier scores for personal and imputed forecasts is very high, r(3)=.98, p=.005. Imputers average Brier scores ranged from .145 to .362 suggesting that traditional analysis' apparent accuracy depends on whether interpreters are better or worse forecasters. <u>Lehner and Stastny (2019)</u> responded. We don't take a stance on their dispute.

						-			
	probabilities Aforementioned ICPM N=99 geopolitical questions, 28 of which had a "fuzzy" resolution criteria	Initial and imputed selected on the da imputed probabilit Due to the posting report authors. <i>Ha</i> advantage. Both ICPM probal with Calibration In than .025, .014, a	d probabilities were ays on which the r ties. g delay, ICPM had <i>owever</i> , longer pos bilities and impute idexes of .047 and nd .016 from othe	e compared to ICF eaders submitted information not av sting delays would d estimates were 1 .097 respectively r studies).	Personal upd. ICPM Updated person ICPM (p=.087). I analysts performer market. Note that scores ranged fro for selection. (See fn 8 for whe about the quality	.150 .188 al fored Data sug ed comp t their in om .145 other we of intell	.158 .195 casts did ggests th parably to itial aver to .362 s can con igence re	d better than at seasoned o the prediction age Brier so there is room	
<u>Kajdasz</u> <u>et al.</u> (<u>2014)</u>	ICPM v. InTrade v. "10 best IC experts we could identify on each topic" N=10 geopolitical questions ⁹ N="152 individual forecasts from the ICPM, InTrade, and individual IC experts over approximately matching topics and time horizons."	Note that the thre ("approximately m The market prices experts. No statis InTrade. Brier score summ ICPM	e groups answere natching topics") s provided significa tical difference ¹¹ in <i>ary statistics</i> n _{forecasts} 48	d different questio antly ¹⁰ more accur accuracy betwee mean .0746	Different number suggests that pre have been well b better if every for with forecasts fro same day and or ICPM's .075 Brie across many que al. (2015). And In suggests that trac predicting confide	of n _{forecc} ediction alancec ecast of m ICPM the sau the sau r is 3x lo estions r Trade is ders we ently, an	asts is con from gro l; it would f IC SME I and InT me time I ower that eported is at .036 ire rarely id so wei	fusing (as it ups might not d have been was matched frade on the horizon). ¹² n its average in Goldstein et 6, which (if ever) re rarely on the	

⁹ "We replicated some of these markets in the ICPM, or identified closely analogous predictions if they existed, so that direct comparisons between the two prediction markets could be made over time."

¹⁰ They report $F_{A \text{ comp}}$ (1, 149) = 19.85, p<.01, $\hat{W}^2 \psi$ = 0.1095.

¹¹ They report F(1, 149) = 1.33, n.s.

¹² Authors write: "We repeatedly collected forecasts from our markets and our experts to sample various time horizons, ranging from very near-term forecasts to as long as 4 months before a subject was resolved. All told, we collected 152 individual forecasts from the ICPM, InTrade, and individual IC experts over approximately matching topics and time horizons."

		InTrade	50	.0366	.0634	wrong side of maybe. ¹³	
		IC SME	54	.1895	.2529		
Beadle (2022); summary here	 465,673 predictions over 3 years 1,375 participants 150 resolved questions, 240 in total "The average time perspective in the FFI tournament was 521 days, i.e. four times as long as the questions in GJP." 	Unfortunately FFI the 150 same qu much more likely. "The standardised almost perfectly s superforecasters their scores after update on their pr "Based on the firs participants in FF 0.11)." No better t "Moreover, even t worse at predictio updating, training superforecasters defence research was strikingly sim	-superforecaster estions, which ma d Brier scores of F imilar to that of the in GJP (0.37)". "N updating. However redictions". It 150 questions, the han predicting 50° hough regular fore on than GJP foreca or grouping), the s compared to re hers with access ilar."	s were selected a akes regression to FI superforecaster e initial forecasts o ote that GJP forec er, the FFI forecast he average Brier s 0.52 (SD: % on all questions ecasters in the FFI asters overall (prot relative accuracy gular forecasters to classified info	and evaluated on ward the mean rs (0.36) were of asters improved ers could not core of the tournament were pably due to not of FFI's (-0.06), and to prmation (-0.1)	"In 2017 the Norwegian Research Defence Establishment (FFI) ran a forecasting tournament intended to investigate if GJP's findings in the ACE tournament would replicate on questions with a longer time horizon, and in a Nordic context." "An important difference from GJP is that FFI's tournament was open to anyone who wanted to participate."	

¹³ A perfectly calibrated forecaster expects on average p - p² Brier points from their prediction. So this average Brier suggests that a "typical" InTrade prediction was either <4% or >96%. From experience, this feels too confident and suggests that questions were either biased towards low noise or that luck is partly responsible for such good performance.

¹⁴ Per <u>comment</u>: "the 60 FFI supers were selected and evaluated on the same 150 questions (Beadle, 2022, 169-170). Beadle also identified the top 100 forecasters based on the first 25 questions, and evaluated their performance on the basis of the remaining 125 questions to see if their accuracy was stable over time, or due to luck. Similarly to the GJP studies, he found that they were consistent over time (Beadle, 2022, 128-131)."

<u>Tetlock et</u> <u>al. (2023)</u>	 For the 25 year timeframe, there were 42 questions about nuclear proliferation, and 40 about boundaries. The study drew on EPJ studies of experts' probabilistic forecasts on slow-motion variables: Base rate of change lower than 25% over 25 years. Base rate of at least 5%. Exercises that had at least 25 forecasts per-participant. 	Experts were more accurate than non-experts in nuclear proliferations questions. The Brier of non-experts was 60 % (= 0.08/0.05 - 1) higher for the 25 year timeframe. Table 2 Accuracy Scores for Nuclear Proliferation (NP) Forecasts $\frac{1}{1000} \frac{1000}{1000} 1000$	 "Cruder operational definition that treated forecasters as experts if they had been educated at the post-graduate level in relevant disciplines and if they saw the topic as central to their professional identity" "The study has many methodological shortcomings [since it "was never a priority in the larger EPJ project"]: small sample sizes, inadequate measures of expertise, a flawed probability scale¹⁵, and a rushed schedule that gave forecasters little time to deliberate."
		<i>Note.</i> For Log Scores, we set probabilities of 0 at .01 and probabilities of 1 at .99.	

Pandemics

¹⁵ "Forecasters used the same 11-point, zero-to-one, subjective probability scale as in other EPJ exercises, with equal 0.1 spacing between levels (0, .1, .2, ..., .9, 1)."



¹⁶ Personal communication with Servan-Schreiber.

		And the 1st place went to one of the very <i>Key problem</i> : experts got busy with the p <i>From the talk:</i> • Individually forecasters are 3% • Crowds performed similarly. • Mixed group +7% over experts	y few public-health professionals who was bandemic, so forecasters updated their for b better (.454 v .467, p=.01). alone.	also a skilled Hypermind forecaster.				
McAndre	Experienced life science pros (n=10) Top-1% Metaculus forecasters (n=11)	Trained forecasters had the highest log scores on average, followed by consensus models, and then subject-matter experts (nonsignificantly, the study is underpowered).						
<u>W,</u>			25th and 75th percentiles for log score	Mean scaled rank ¹⁷				
Cambeir		individuals all	[0.42, 2.98]	—				
<u>Besiroglu</u>	the 2 groups	individual forecasters	_	.56 80CI: [.18, .94]				
(2021)	Only 6 out 23 questions have	individual experts	—	.48 80Cl: [.08, .98]				
	a COVID-19 vaccine.	consensus all	[0.98, 2.96]	.58 80Cl: [.49, .63]				
		consensus forecasters	[1.24, 2.90]	.56 80CI: [.43, .72]				
		consensus experts	[0.65, 3.07]	.53 80Cl: [.35, .73]				

¹⁷ Given N log scores, scaled rank assigns a value of 1/N to the smallest log score, a value of 2/N to the second smallest log score, and so on, assigning a value of 1 to the highest log score. (As with log scores, here computed from probability density functions, — the higher rank the better.)

<u>Bosse et</u> al. (2021)	Two semi-mechanistic models Ensemble of <u>all models</u> submitted to the Forecast Hub	Crowd consistently ou the Hub ensemble wh deaths. <i>Weighted Interval Sco ensemble</i> Two weeks ahead:	utperformed en nen forecastin ore (WIS, the Cases	pidemiolog og cases bu <i>lower the b</i> Deaths	Not clear how good Forecast Hub models were but their affiliations are impressive. Irrespective of the above, suggests that crowd forecasting might be useful in practice.	
	n=32 forecasters (17 are self-identified experts in	Hub ensemble	1	1		
	forecasting or epidemiology)	Renewal model	1.40	1.79		
		Convolution model	_	1.22		
		Crowd	0.89	1.26		
		For cases, our contributions) conformation on the forecasting horizons (For deaths, contributions) conformation of the formation of the formation of the forecast in the forecast were largely for the forecast w	outions (comp isistently impl ce.g., rel. WIS ons from the rformance on asingly negati 1 two weeks is from both t negative.	pared to the roved perfo 0.9, two we renewal mo ly for one w ive impact o ahead, 1.0 he renewal	Hub ensemble without rmance across all eeks ahead). del and crowd forecast veek ahead predictions on performance for longer 5 four weeks ahead). model and the crowd	
	A single superforecaster	Forecaster did 10% b	etter than exp	perts as jud	ged by Brier score:	As usual, it's unclear if the panel faced other
<u>Liptay</u> (2021)	CDC-funded panel of experts	Superforecaster		.246		
	n=28 pandemic-related questions from UMass	Experts		.268		

Movies										
	Movie critics: n=40	Prediction market <u>RMSE</u> was 10%+ better than pundits.								
Pathak et	Betfair, a prediction market:	RMSE for 2013 Oscar								
<u>al (2015)</u>	variable n, including "low liquidity markets"		Days before	Categories	Categories Experts		r			
	Predicting Oscar winners	Average, n=40	3	24	.20	.18				
	5	Nate Silver	3	6	.26	.18				
		Ben Zauzmer	8-9	21	.25	.20				
		(Hollywood Stock Exchange seems to be doing 10%50% worse than Betfair, Intrade, and PredicWise.)								
Spann &	Hollywood Stock Exchange, a virtual-points prediction market	HSX is much bett recalibrated HSX (n=140).	ter than BOR in prediction is no	terms of <u>MAPE</u> nsignificantly di	(n=24). And fferent from B0	МС				
<u>SKIELA</u> (2003)	Two expert predictions: Box	<u>MAPE, n=24</u>								
(2003)	Office Mojo, Box Office Report.	HSX		40.62	40.62					
		HSX, recalibrate	ed	36.48	36.48					
		ВОМ		35.30	35.30					
		BOR	BOR 53.40							
		MAPE, n=140								
		HSX		31.11	31.11					
		HSX, recalibrate	ed	28.40						

		BOM 28.05		
SCOTUS				
<u>Katz et</u> <u>al. (2017)</u>	7,000 participants 600,000 predictions 450 cases	Built an impressively accurate mode predictions, and from Ruger et al. (2 outperform experts.	el on top of FantasySCOTUS 2004) we know that simple models	FantasySCOTUS
<u>Blackma</u> <u>n et al.</u> (2012)	The Forecasting Project's decision tree vs FantasySCOTUS vs The Forecasting Project's experts FantasySCOTUS most active users vs other users	75% v 64.7% v 59.1% — the comparation of the power predictor average, 7.93 average, 7.25 points. And "The result the power predictors' forecasts we Although the power predictors generate make rather strong predictions to br	arison is between different terms. ¹⁸ points, was higher than the crowd ults do not conclusively prove that ere superior to those of the crowd. erally do better, the crowd is able to ridge the gap."	Most (seems like at least 75%) active FantasySCOTUS betters sometimes have specialized backgrounds. See a <u>blogpost</u> and ¶35-6 of <u>the paper</u> .
<u>Ruger et</u> al. (2004)	Fairly simple decision tree vs subject matter experts	The model predicted 75% of the cas accurate than their experts with 59.	ses correctly, which was more 1%.	The Forecasting Project, SCOTUS
Elections				

¹⁸ It's unclear to me how well they did compare to a prior based on how often SCOTUS reverses the decisions. The historical average is ~70% with ~80% reversals in 2008, the relevant term.

	Hypermind and 7 statistical	Moor	Daily Brian	Scoro			Low n and orrors are somewhat correlated as
Servan-S	models	.34	Hypermind				not particularly informative.
chreiber	6 questions on U.S. 2014	.41	Daily Kos				
<u>&</u> Atapasov	majority-control of the	.43	Huffington I	Post			
(2015)	most-undecided states.	.43	PredictWise	e			
· · · · · ·		.45	Models Me	an			
		.46 Washington Post					
		.46 FiveThirtyEight					
		.48 New York Times					
		.68	.68 Princeton Election Consortium				
Miscellaneous							
Corporate setting: demand			prediction m	arket / MSE expe	rts at firms		
Cowgill &	forecasting, project completion, project quality,	Ford		Google	¹⁹	_	
Zitzewitz	external events	0.742 0.727		0.727	0.924 0.908		
(2015)							

¹⁹ Anonymous basic materials conglomerate.

Search criteria

We were given a set of initial studies to branch out from.

- Good Judgement Project
- <u>Tom McAndrew studies</u>
- <u>Hypermind + Johns Hopkins</u>

And some general suggestions for scholarship:

- look for review articles
- look for textbooks and handbooks or companions
- find key terms
- go through researchers' homepages/google scholar

Superforecasting began with IARPA's ACE tournament.²⁰ We think the evidence in Tetlock's *Expert Political Judgment* doesn't fit: there were no known skilled-amateur forecasters at that point. See <u>Appendix C: Tetlock's Expert Political Judgment</u>.

A Google Scholar <u>search</u> for studies funded by IARPA ACE yielded no studies. We looked at other IARPA projects (ForeST, HCT, and OSI), which sounded remotely relevant to our goals.

We searched Google Scholar for (non-exhaustive list): "good judgment project", "superforecasters", "collective intelligence", "wisdom of crowds", "crowd prediction", "judgemental forecasting", ..., and various combinations of these, and "comparison", "experts", ...

We got niche prediction markets from the Database of Prediction Markets and searched for studies mentioning them. Hollywood SX and FantasySCOTUS paid off as a result. We also searched for things people commonly predict: sports, elections, Oscars, and macroeconomics.

In the process, we read the papers for additional keywords and references. We also looked for other papers from the authors we encountered.

Regarding AI forecasting

• In more complex domains, like ML, there could be significant returns to knowledge and expertise.

It seems to us that moving from *generalist forecasters* to *competent ML practitioners/researchers* might be better because:

• To predict e.g. scaling laws and emerging capabilities, people need to understand them, which requires some expertise and understanding of ML

²⁰ For completeness we could mention <u>Galton (1907)</u>, the first demonstration of the wisdom of crowds.

- It's unclear whether general forecasters actually outperform experts in a legible domain, even though we believe in the phenomenon of superforecasting, (that some people are much better forecasters than most). We also liked David Manheim's <u>take</u> on Superforecasting.
- We think that this will plausibly reduce ML researchers' aversion to forecasting proposals and if we were to execute it, we would be selecting good forecasters based on their performance anyway. It seems <u>potentially feasible</u>.

Finally, we note that the above reasoning is heavily limited by a lack of data (lack of it collected and a lack of it made available). We hope that the experimental data gets reanalyzed.

Thanks to Emile Servan-Schreiber, Luke Muehlhauser, and Javier Prieto for comments. These commenters don't necessarily endorse anything in this post, and mistakes are our own. Research funded by Open Philanthropy.

Appendix A: prediction markets vs. opinion pools

Appendix B: Table of less relevant studies

Appendix C: Tetlock's Expert Political Judgment

See also: Database of Prediction Markets

<u>Changelog</u>