

Predictable Noise and patterns from millions of questions

Sida Wang
Research Scientist, FAIR at Meta

Data contribution: Sean OBrien, Alex Gu, Lovish Madaan, Dieuwke Hupkes, Jiawei Liu, Yuxiang Wei, Naman Jain, Yuhang Lai, Sten Sootla

Outline

- Intro
- Motivation: benchmarks are getting harder and smaller
- Exploratory data analysis: What does the evaluation say beyond single accuracy
- Implications on the noise levels and why?
 - Recommendations about noise intervals
- How can we make this better + examples

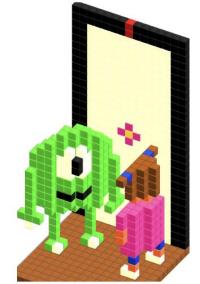
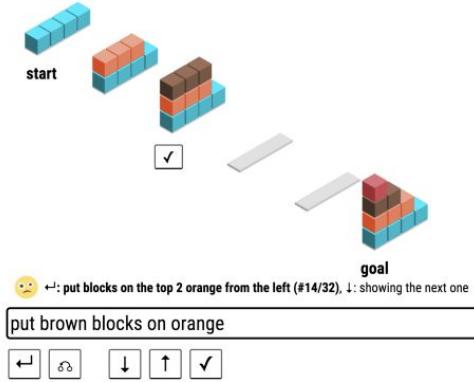
Quick self-intro

Learning. When the human provides feedback in the form of a particular y , the system forms the following loss function:

$$\ell(\theta, x, y) = -\log p_\theta(y | x, s) + \lambda \|\theta\|_1, \quad (4.2)$$

$$p_\theta(y | x, s) = \sum_{z: \|z\|_1=y} p_\theta(z | x). \quad (4.3)$$

2016: learn from human feedback using marginal likelihood, also originated here at UCB. an alternative to RL without emphasis on exploration



Monsters, Inc: initial – move forward – add green brown floor – add girl – go back and down – add d (some steps for moving are omitted)

CWM: An Open-Weights LLM for Research on Code Generation with World Models: RL for SWE-bench, stabilizing RL

Code LLM: [AST-FIM](#), [SWE-RL](#), [eval-arena](#), [LEVER](#), [Coder-reviewer](#), [InCoder](#), [MBR-exec](#).

Benchmarks: [SWE-bench M](#), [Spider 2.0](#), [LiveCodeBench](#), [SAFIM](#), [CRUXEval](#), [DS-1000](#)

A bit of reflection on evaluation

Do 10% better on ImageNet, and no more doubts about deep learning

- necessary condition: no doubt about statistical significance

Do we similarly trust our current LLM evaluations?

Many papers report 2-10% improvements on HumanEval, or a few % on SWEbench

How can we interpret this?

- old answer is statistical testing: 2.5% on HumanEval is not stats significant

Test set sizes are now much smaller than MNIST

Deep learning was measured and validated by MNIST and ImageNet

Testset sizes: MNIST: **10k**, ImageNet: **100k**

LLM multiple choice: MMLU ~10k

Generative test-based eval (generate a short program)

HumanEval: 164, MBPP+: 378, CruxEval: 800, DS-1000: 1000, DMC: 165,
LiveCodeBench ~1000, Alder: 225

Agent Evals (often need 100k+ tokens)

SweBench: 2.3k, -Verified: 500, -lite: 300, -multimodal: 517

T-Bench: 80

But each question is more informative!

- To answer each question, a lot of text/code is generated and evaluated by tests. Is this more informative than True/False and multiple choice?

```
def is_simple_power(x, n):
    """Your task is to write a function that returns true if a number x is a simple
    power of n and false in other cases.
    x is a simple power of n if n**int=x
    For example:
    is_simple_power(1, 4) => true ..
```

- For humans, you can ask a few good questions and get good information
 - [IOI](#) / [IMO](#): 6 questions, Putnam: 12 questions
 - has an interesting range of difficulties, but then maybe humans are still fairly similar

Solve an open problem?

- solve even one major open problem is significant!
 - no one will object that a sample size of 1 is not statistically significant

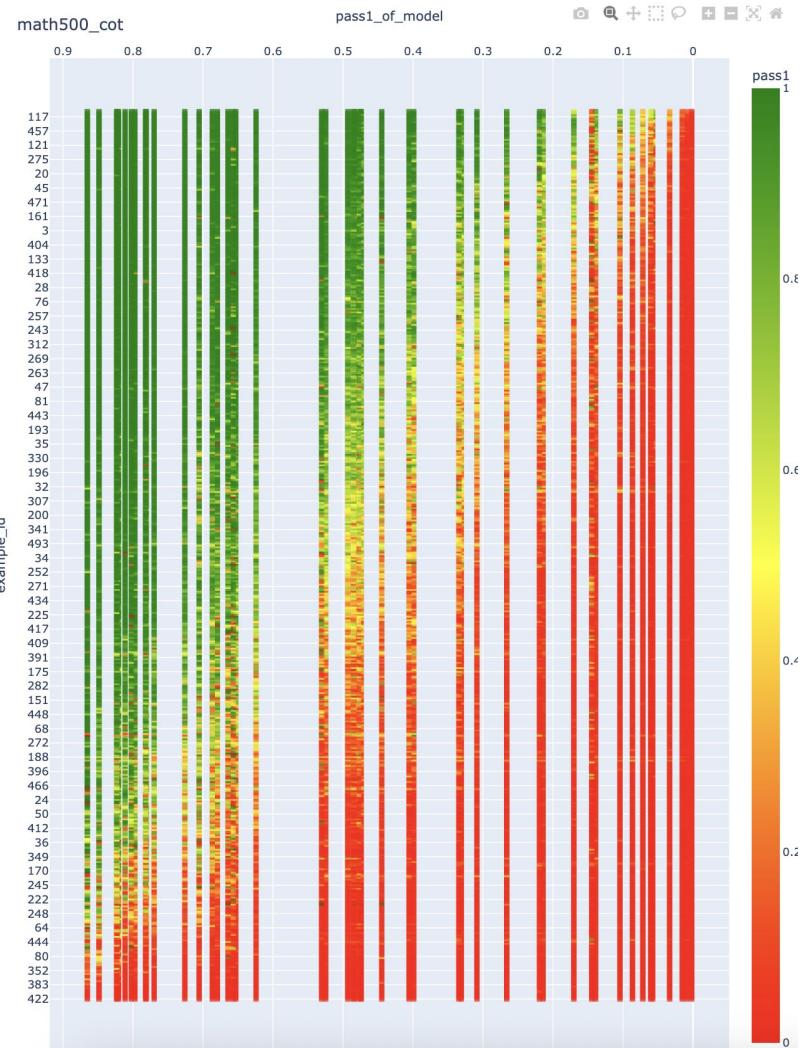
Many daily agentic tasks are more economically important

So what's the case?

- A: small benchmarks are not reliable even if they contains some hard generative/agentic problems.
- B: Solving a single hard problem could be significant already, so we can and should use small but very hard benchmarks

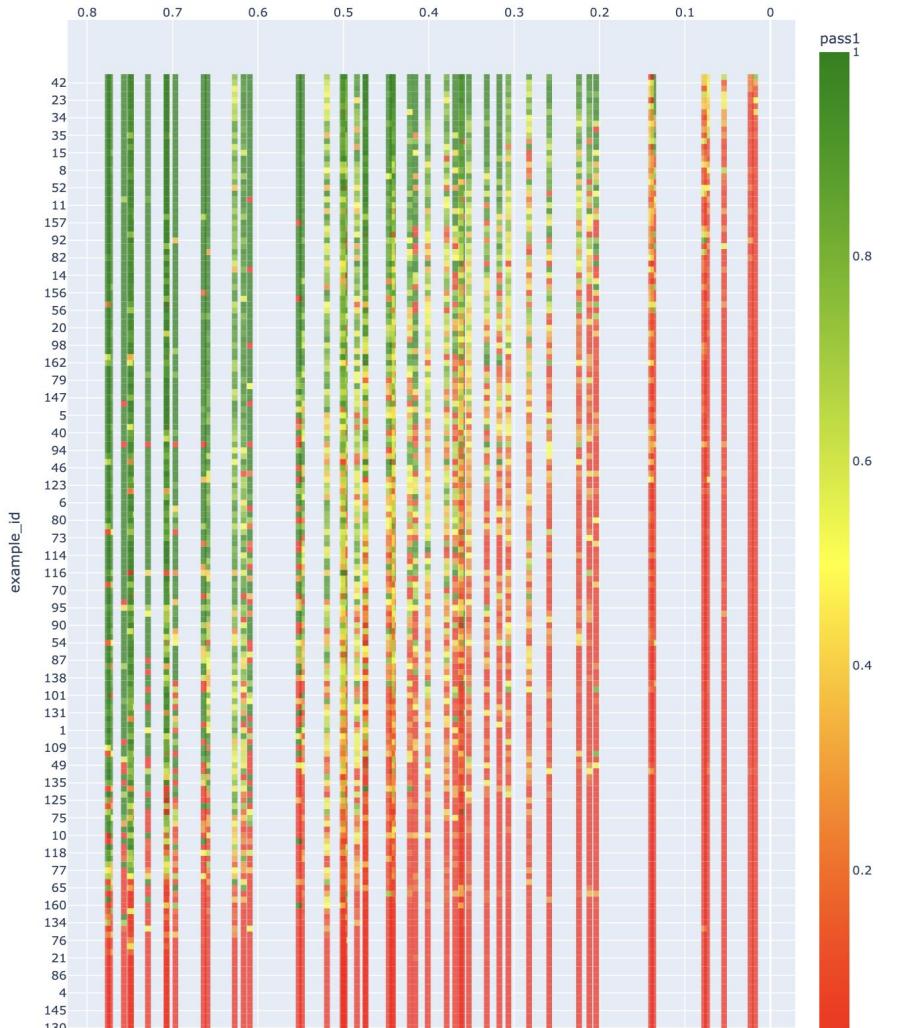
Let's do some exploratory data analysis

- Each column is a model, x axis is their pass@1
- Each row is a question, sorted by their difficulty. From easiest at top to hardest at bottom
- Color is pass@1 of model on a question, over 10+ samples whenever possible



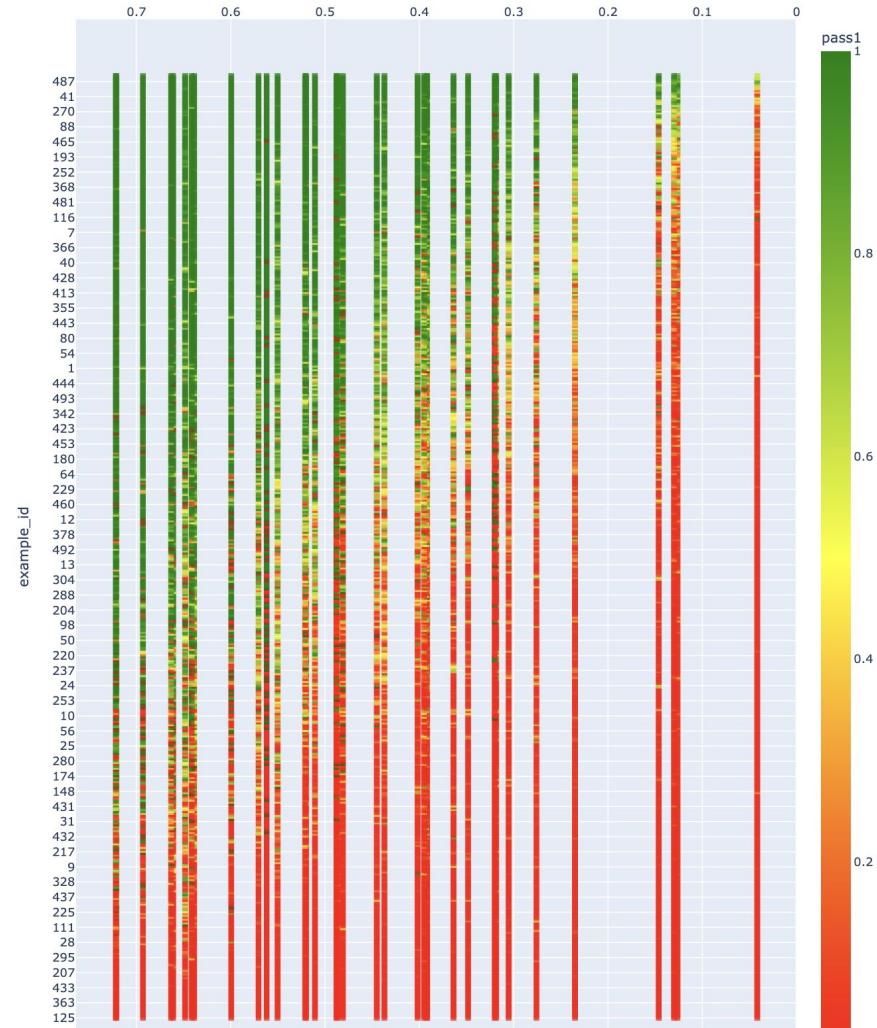
human_eval_plus

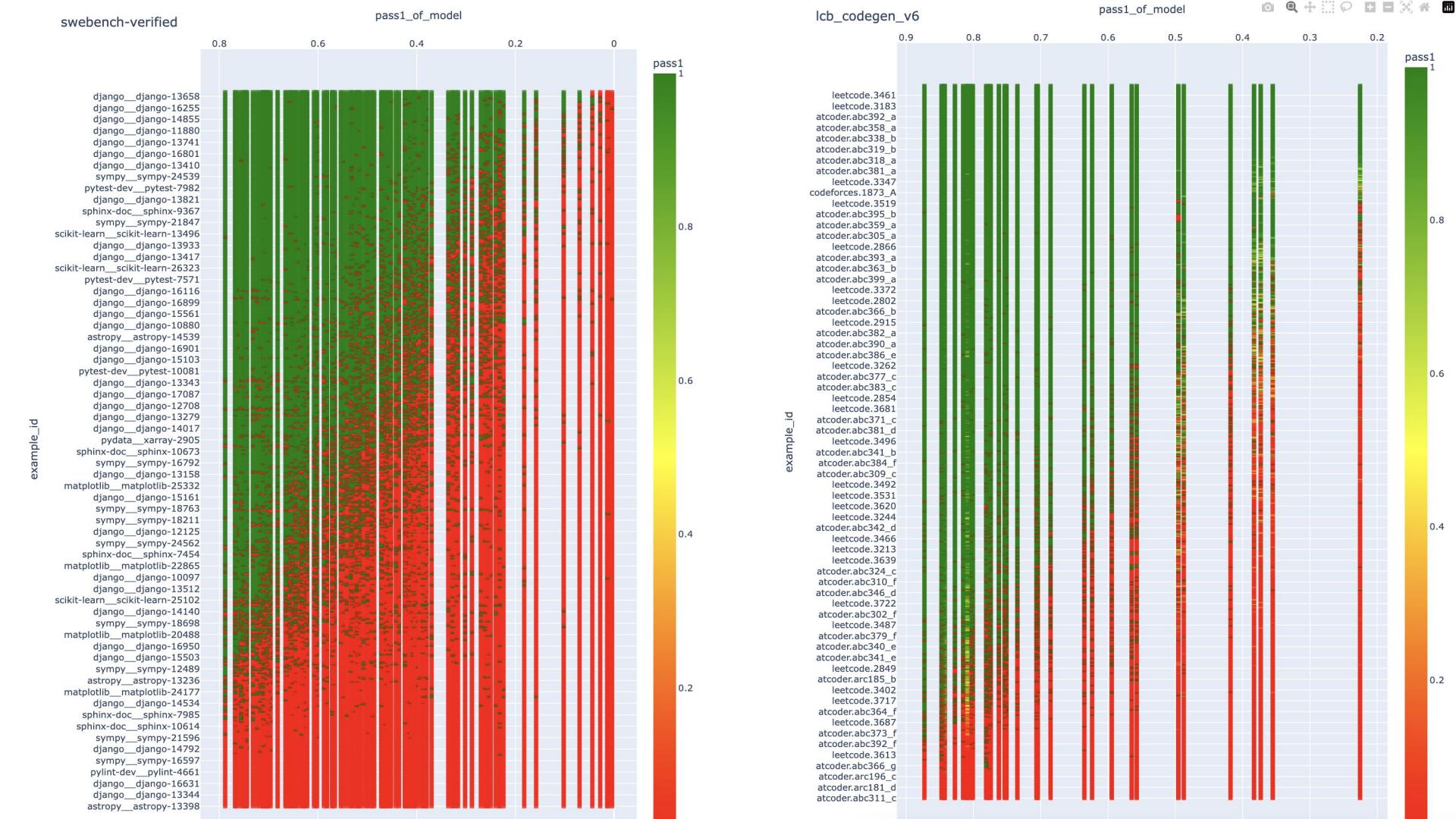
pass1_of_model



mbpp

pass1_of_model





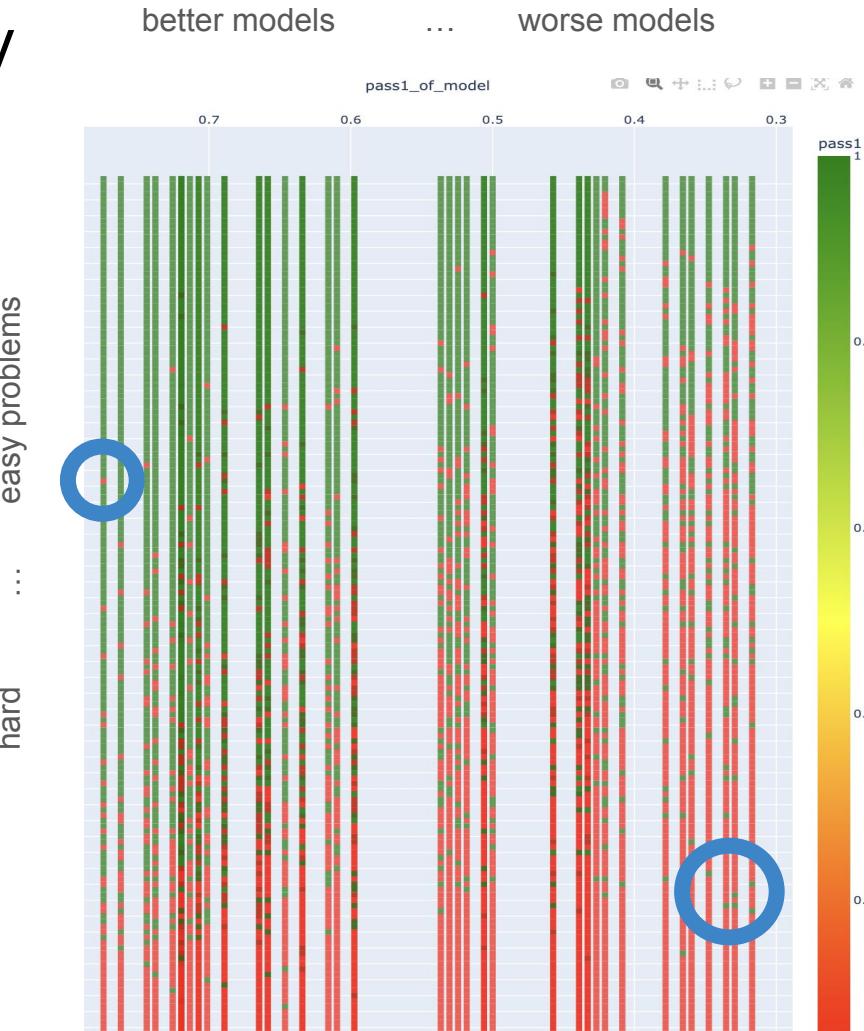
There is a lot of inconsistency

HumanEval

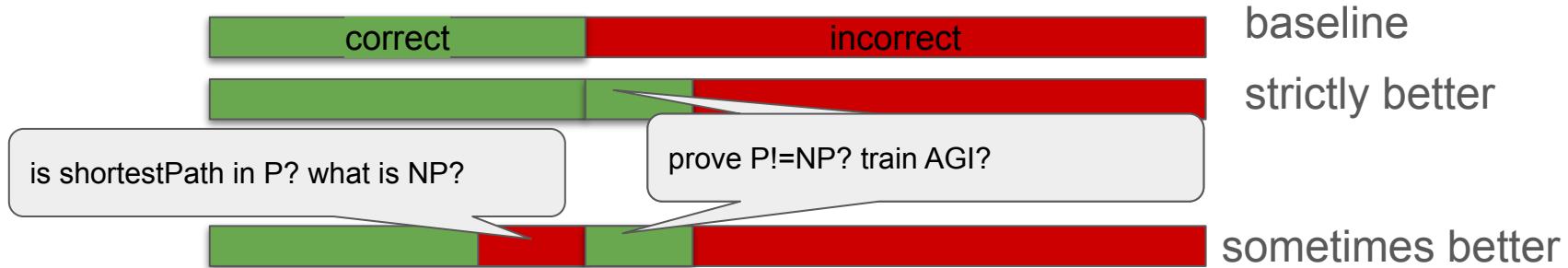
The worst models sometimes succeeds on hard problem

The best models still fails on easy problems

So perhaps the problems are not too special



Can we make up for the low statistical power with hard questions with long answers?



Not if the answers are inconsistent:

- wrote down a hard and long proof of $P \neq NP$ but cannot define NP ?
copied it somewhere?

red are not sharp 0s, but just lower probability

This question cannot be a very hard open problem if a much weaker model just have a smaller probability of solving the problem

Increasing from 1% correct to 99% may not be that significant (especially on verifiable problems). The first 1% is the most impressive.



So what's the case?

- A: small benchmarks are not reliable even if they contains some hard generative/agentic problems
- B: ~~Solving a single hard problem could be significant already, so we can and should use small and hard benchmarks~~

Help fix my kettle



label the parts for me

- kettle (boils water)
 - agent says it's a robotic vacuum
 - should I trust its later advice (could still be correct)
-
- Hard problems often have many steps
 - Probability of error increases exponentially as the number of steps increases

⌚ Thought process

Dissected robotic vacuum parts for comprehensive labeling

Looking at this disassembled device (which appears to be a robot vacuum wheel/drive assembly), here are the main components I can identify:

Key Parts:

1. **Drive Motor** - The central circular component with the brass/bronze colored ring - this provides power to rotate the wheel
2. **Motor Housing/Mount** - The black and

So probably statistical analysis is still important

How does it work? Let's go over the basics

I'll put down some math, and convey the intuition with some examples

The super-population (basic stats, Miller 2024)

- With sample questions from a super population, then the mean and variance of the population can be estimated from these sample

On a evaluation, we compute the average of all N questions to get the mean accuracy

$$\bar{A} := \frac{1}{N} \sum_{i=1}^N A(x_i) \approx \text{E}_x[A(x)].$$

The estimated variance of A and standard error (SE) of the mean \bar{A} are respectively

$$\begin{aligned}\text{Var}_x[A(x)] &\approx \frac{1}{N} \sum_{i=1}^N (A(x_i) - \bar{A})^2, \\ \text{SE}(A) &= N^{-1/2} \cdot \text{Var}[A]^{1/2}.\end{aligned}$$

Paired comparison. $A(x)$ and $B(x)$ are correlated, and in an actual benchmark, the same set of question x is used for both $A(x)$ and $B(x)$. To be more powerful for free, we can use the paired variance

$$\text{Var}_x[A(x) - B(x)] = \text{Var}_x[A] + \text{Var}_x[B] - 2\text{Cov}_x[A(x), B(x)]. \quad (4)$$

The Cov term may potentially reduce the paired variance to 0 when A and B are perfectly correlated.

— —

Ask if this result is likely just due to chance

Once we have the standard error, we obtain the z-score $z = \frac{\bar{A} - \bar{B}}{\text{SE Var}[A - B]}$ which follows a standard normal distribution. For example, $\Pr[|z| > 1] = 0.32$ is very weak, $\Pr[|z| > 5] < 10^{-6}$ is beyond doubt, and $\Pr[|z| > 1.96] = 0.05$ for a p -value of 5% is a conventional and reasonable standard.

Noise seems pretty large for small datasets

Example 1 (HumanEval). For the popular HumanEval dataset, $N = 164$, say a model has $\bar{A} = 0.5$ accuracy. The SE is 4%. you'd need to produce a difference of $4\% \times \sqrt{2} \times 1.96 = 10.8\%$ to get 0.05 p -value using the unpaired test. With a paired test, the difference SE is typically still 4%, so a difference of $4\% \times 1.96 = 7.84\%$ is needed to produce a 0.05 p -value. Both are a larger difference than the gains reported by many papers.

Paired vs. Unpaired: intuition

Example 2 (Paired vs unpaired). In the unpaired case, suppose A scored 50% on this year's exam, B scored 60% on last year's exam and the exams contain different exact questions drawn from the same distribution. So B might have done better because last year's exam was easier. In the paired case, suppose B scored 60% on the same exact exam as A with exactly the same questions, then the same 10% difference is more indicative of the better performance by B .

Paired: $x_i \sim \text{Data} \rightarrow A(x_i), B(x_i)$

Unpaired: $x_i, x_j \sim \text{Data} \rightarrow A(x_i), B(x_j)$

Use the ability to draw independent samples?

This is a remarkable ability of the models, humans can't seem to do this. A common confusion is to consider this the only type of noise.

To capture this setting, we use ϵ for the seed generating the noise, which is independent of the question x , thus allowing the concise notation $\text{Var}_x[\mathbb{E}_\epsilon[A]] = \text{Var}_x[\mathbb{E}_{\epsilon|x}[A(x, \epsilon) \mid x]]$. As before, but now also averaging over the K samples for each question x , we consider the average score

$$\bar{A} = \frac{1}{N} \sum_{i=1}^N \frac{1}{K} \sum_{j=1}^K A(x_i, \epsilon_{ij}) \approx \mathbb{E}_{x,\epsilon}[A(x, \epsilon)], \quad (5)$$

The estimated variance is then

$$\text{Var}_{x,\epsilon}[A(x, \epsilon)] \approx \frac{1}{N} \sum_{i=1}^N \frac{1}{K} \sum_{j=1}^K (A(x_i, \epsilon_{ij}) - \bar{A})^2. \quad (6)$$

Total variance to standard errors

We have the law of total variance satisfied by the components

$$\text{Var}_{x,\epsilon}[A] = \text{Var}_x[\text{E}_\epsilon[A]] + \text{E}_x[\text{Var}_\epsilon[A]]. \quad (7)$$

$\text{Var}_{x,\epsilon}[A]$ is the total variance of first drawing a question x from the pool of questions, and then drawing a sample answer for this question. $\text{Var}_x[\text{E}_\epsilon[A]]$ is variance of sampling from the data, which is also the variance of the expected score $\text{E}_\epsilon[A]$. $\text{E}_x[\text{Var}_\epsilon[A]]$ is the variance due to sampling from the model. All of which can be normalized to be the standard errors,

$$\text{SE}(A) = N^{-1/2}(\text{Var}_{x,\epsilon}[A])^{1/2}, \quad (8)$$

$$\text{SE}_x(A) = N^{-1/2}(\text{Var}_x[\text{E}_\epsilon[A]])^{1/2}, \quad (9)$$

$$\text{SE}_\epsilon(A) = N^{-1/2}(\text{E}_x[\text{Var}_\epsilon[A]])^{1/2}. \quad (10)$$

Paired comparison

Paired comparison. When we have two models A and B , we can consider the paired standard deviation $\text{Var}_{x,\epsilon}[A(x, \epsilon) - B(x, \epsilon)]$, which also satisfy the law of total variance on $A - B$,

$$\text{Var}_{x,\epsilon}[A - B] = \text{Var}_x[\mathbb{E}_\epsilon[A - B]] + \mathbb{E}_x[\text{Var}_\epsilon[A - B]]. \quad (11)$$

This formula allows us to compute from the samples directly.

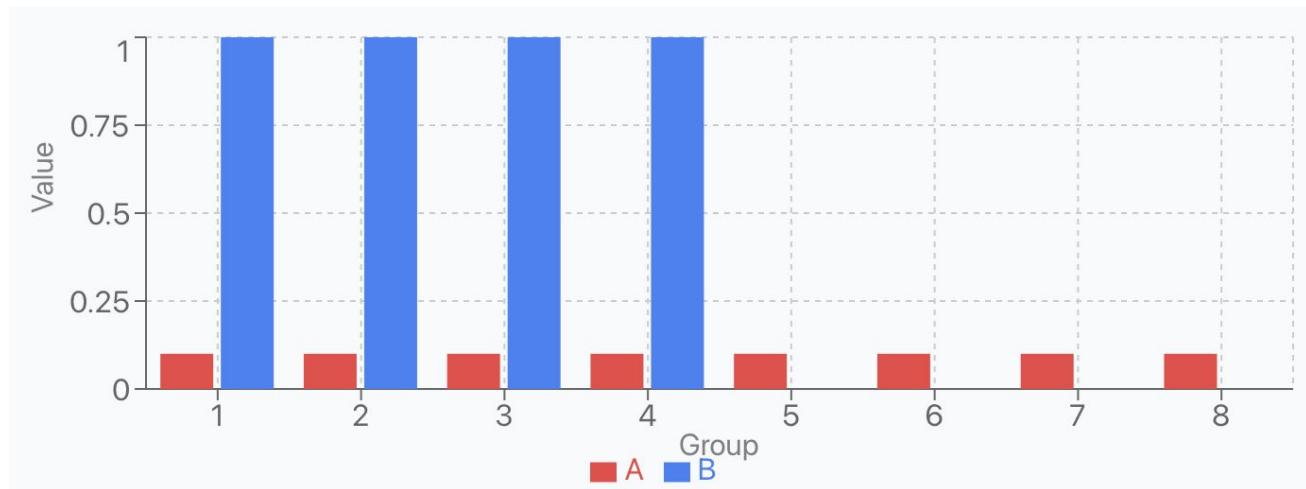
$$\text{Var}_{x,\epsilon}[A - B] = \text{Var}_{x,\epsilon}[A] + \text{Var}_{x,\epsilon}[B] - 2\text{Cov}_{x,\epsilon}[A, B] \quad (12)$$

$$= \text{Var}_{x,\epsilon}[A] + \text{Var}_{x,\epsilon}[B] - 2\text{Cov}_x[\mathbb{E}_\epsilon[A], \mathbb{E}_\epsilon[B]]. \quad (13)$$

By the law of total covariance, $\text{Cov}_{x,\epsilon}[A, B] = \text{Cov}_x[\mathbb{E}_\epsilon[A], \mathbb{E}_\epsilon[B]] + \mathbb{E}_x[\text{Cov}_\epsilon[A, B]]$ where the second term is 0 since different random samples are independent.

Correct with probability 0.8 leads to some weirdness

Example 4 (better mean vs. better max). Suppose B gains a consistent improvement over A on all questions i , $E_\epsilon[A(x_i, \epsilon)] = 0.6$, $E_\epsilon[B(x_i, \epsilon)] = 0.8$. While this can easily be statistically significant, it might be also be achievable by majority voting, verification or just sharpening the distribution. Now suppose that on half of the questions, $E_\epsilon[A] = 0.1$, $E_\epsilon[B] = 1$, but on the other half we have $E_\epsilon[A] = 0.1$, $E_\epsilon[B] = 0$. While B has a much higher average, A might actually be much better. For example, if the questions are problems that can be verified such as search problems or proofs, then A is successful 100% of the time with a verifier in the loop.



Compute with simple array operations

~~The answer is $p(1-p)$, let's compute with bootstrap...~~

Name	formula	code	Bernoulli
total variance	$\text{Var}_{x,\epsilon}[A]$	<code>var(A)</code>	$\bar{p}(1 - \bar{p})$
data variance	$\text{Var}_x[\mathbb{E}_\epsilon[A]]$	<code>var(mean(A, axis=1))</code>	$\frac{1}{N} \sum_i (p_i - \bar{p})^2$
model variance	$\mathbb{E}_x[\text{Var}_\epsilon[A]]$	<code>mean(var(A, axis=1))</code>	$\frac{1}{N} \sum_i p_i(1 - p_i)$

Table 1: translation of unpaired to code where $A \in \mathbb{R}^{N \times K}$ with K samples for each of N questions.

Formula	Code
$\text{Var}_{x,\epsilon}[A - B]$	<code>var(A)+var(B)-2*cov(mean(A, axis=1), mean(B, axis=1))</code>
$\text{Var}_x[\mathbb{E}_\epsilon[A - B]]$	<code>var(mean(A, axis=1)-mean(B, axis=1))</code>
$\mathbb{E}_x[\text{Var}_\epsilon[A - B]]$	<code>mean(var(A, axis=1)+var(B, axis=1))</code>

Table 2: code for paired estimators (with bias)

The other viewpoints gives about the same results

Bootstrap: Drawing samples (with replacement) from the empirical distribution

$\hat{=}$ Drawing more samples from unseen super population

Sign test: wherever predictions are different, how likely is it to be due to randomness (same for A and B that are close. if not close, then already statistically different)

Eval-Arena ([doc/code](#))

Results: <https://crux-eval.github.io/eval-arena/main/>

For each benchmark, do pairwise comparisons / statistical testing on all models pairs. Establish that no model is consistently better than another model (among statistically questionable models).

- Assuming no special problems and the model has some chance of error on each problem, we can measure noise meaningfully for each dataset
- **you can just read off the numbers in eval-arena instead of trusting the authors to do their testing!**

Predictable noise

$SE(A) \sim SE(A-B)$

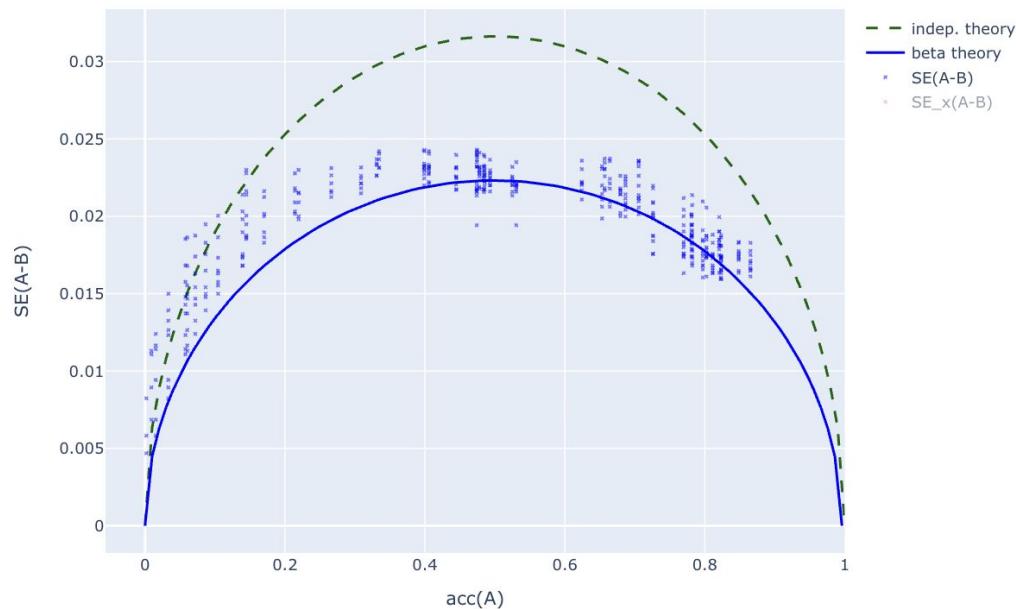
are roughly similar

-> correlation ~ 0.5

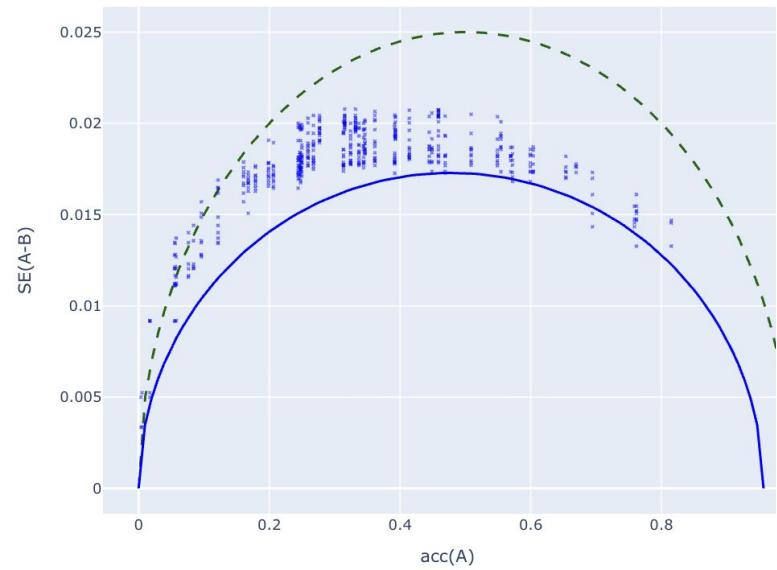
benchmark_id	size	models	SE(A)	SE(A-B)	corr(A,B)
CRUXEval-input-T0.2	800	37	<u>1.7</u>	<u>1.8</u>	<u>56</u>
CRUXEval-input-T0.8	800	31	<u>1.6</u>	<u>1.8</u>	<u>69</u>
CRUXEval-output-T0.2	800	37	<u>1.7</u>	<u>1.6</u>	<u>64</u>
CRUXEval-output-T0.8	800	31	<u>1.7</u>	<u>1.6</u>	<u>75</u>
DS1000	1000	105	<u>1.3</u>	<u>1.5</u>	<u>46</u>
agi_english	2546	35	<u>0.93</u>	<u>1.1</u>	<u>25</u>
arc_challenge	1165	38	<u>1.4</u>	<u>1.3</u>	<u>62</u>
gsm8k	1319	37	<u>1.1</u>	<u>1.3</u>	<u>38</u>
hellaswag	10042	36	<u>0.41</u>	<u>0.26</u>	<u>77</u>
humaneval	164	78	<u>3.5</u>	<u>3.7</u>	<u>47</u>

There is a dependence on overall acc(A)

math500_cot

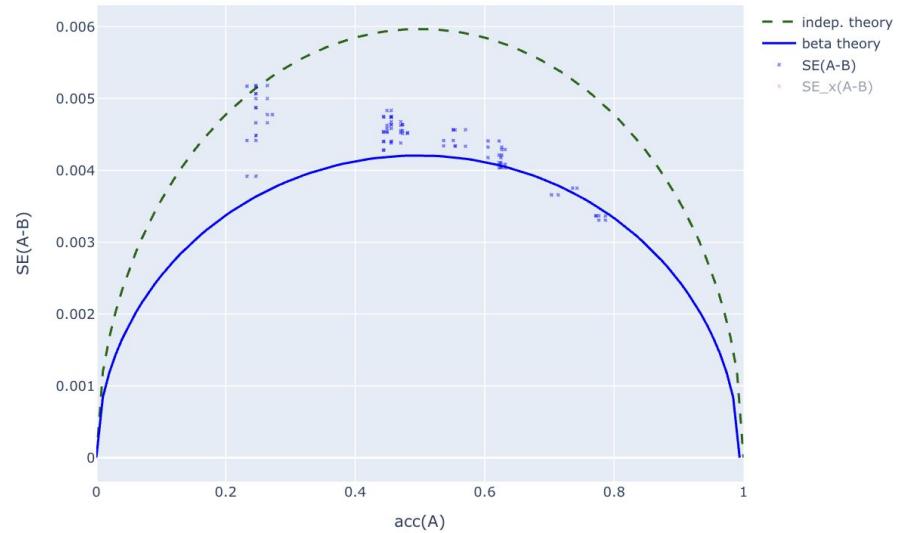


cruxeval_output_cot

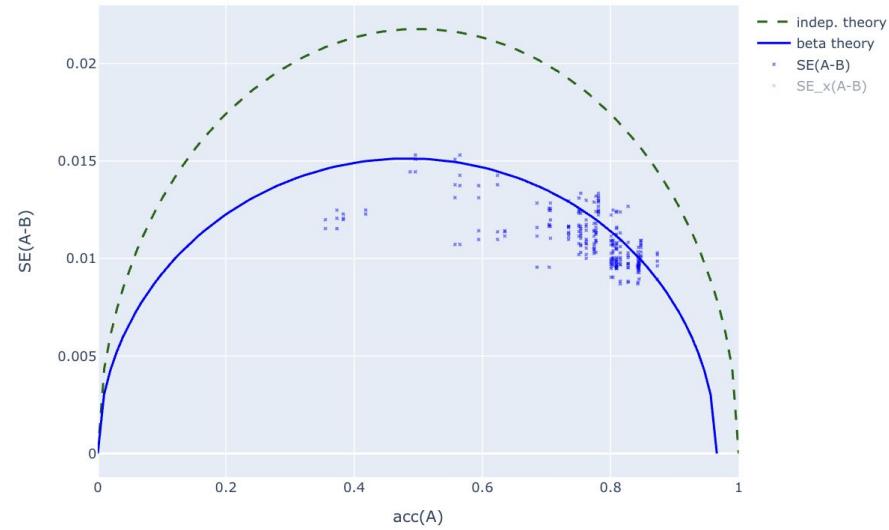


There is a dependence on overall acc(A)

mmlu

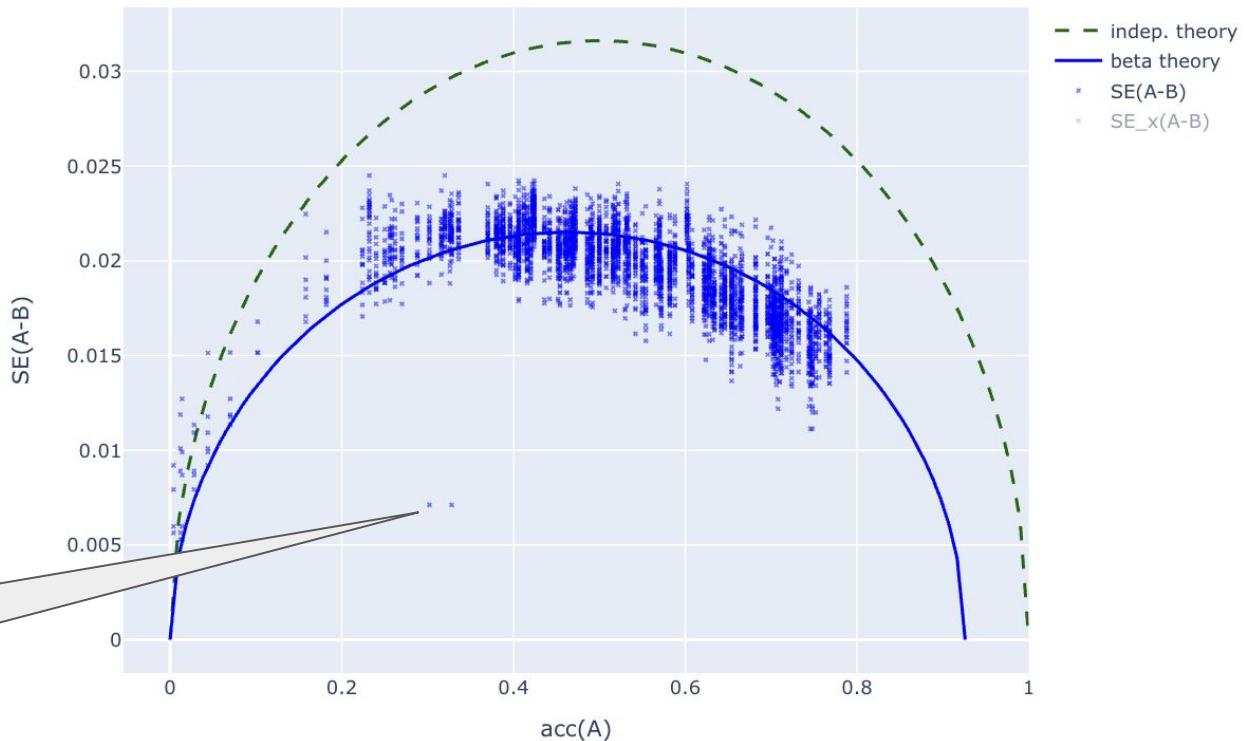


lcb_codegen_v6



Swebench verified

swebench-verified

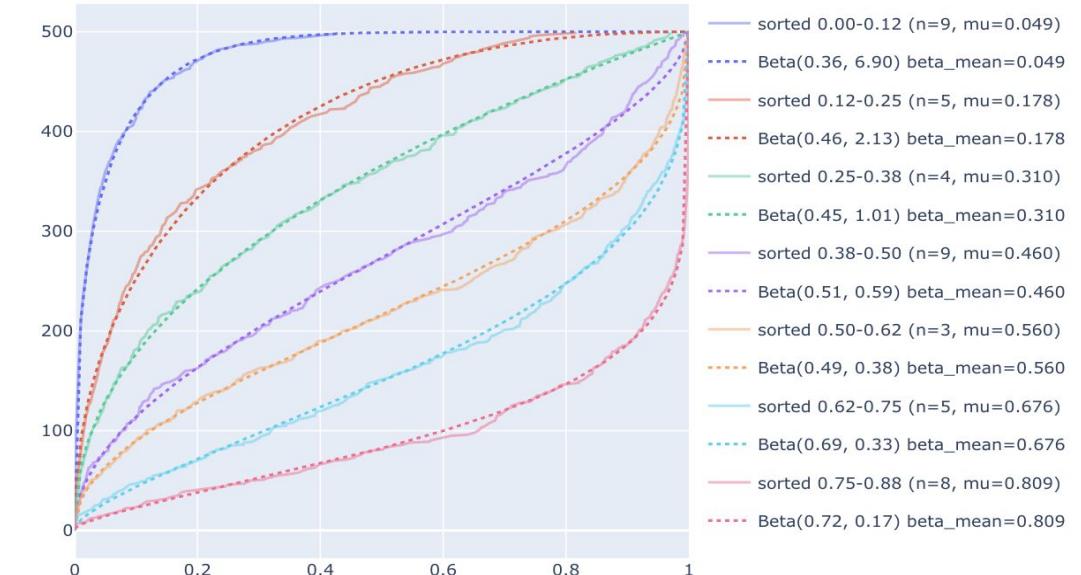
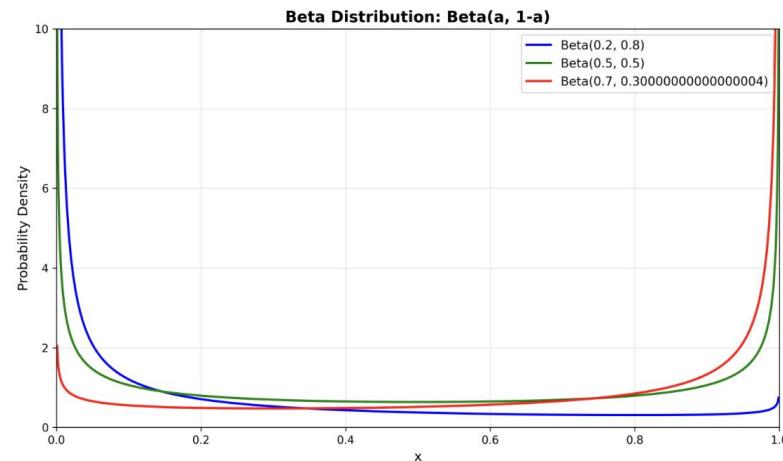


SWEFixer, vs.
SWEFixer+P2P filter

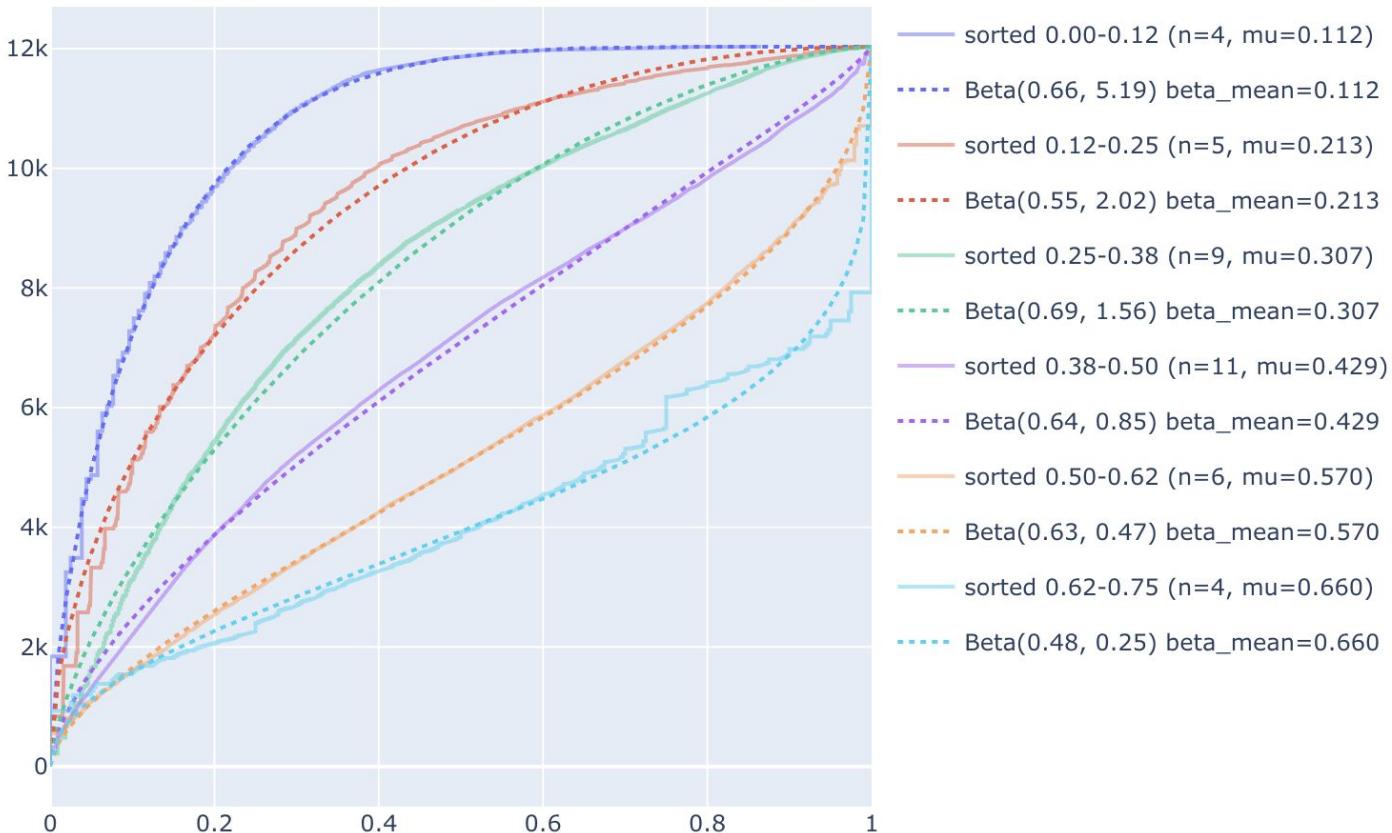
Empirical probabilities follow a Beta distribution

Once performance is good, then the bimodal $\sim \text{Beta}(a, 1-a)$

cdf on math500_cot



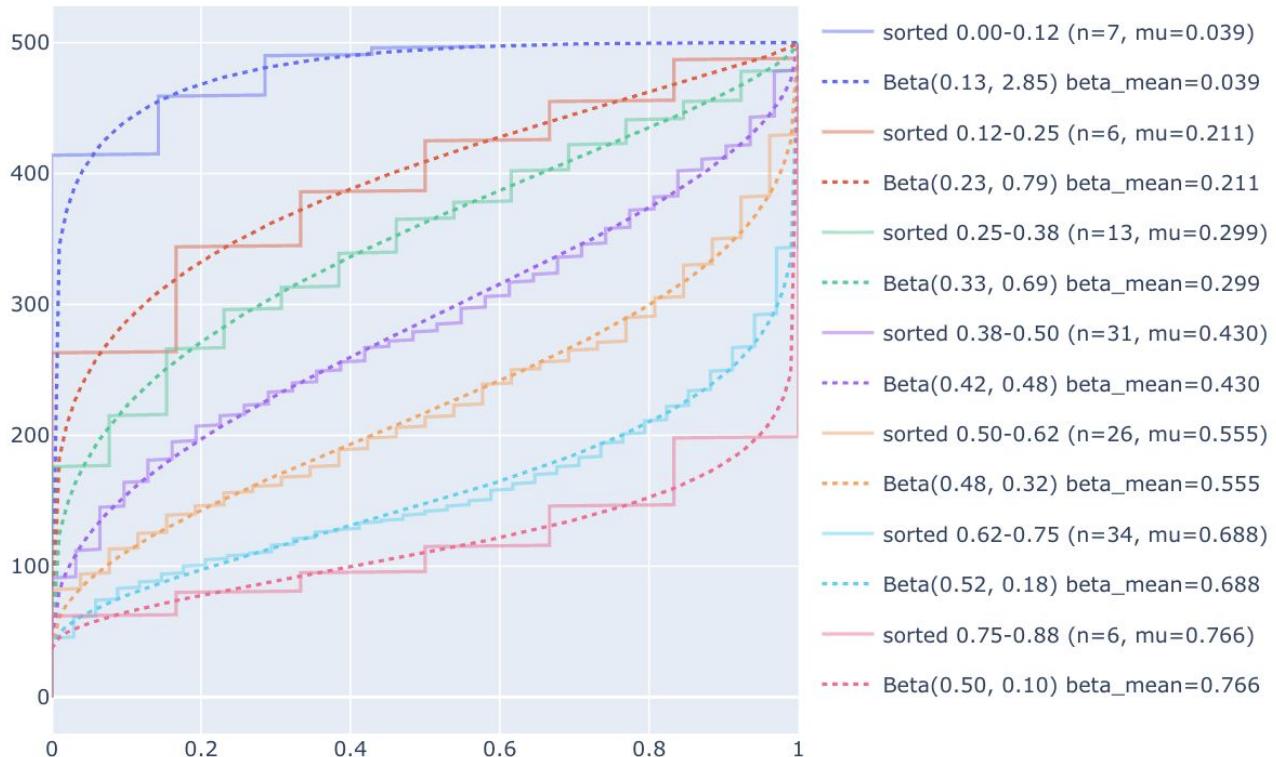
More beta



Betas

only binary preds

cdf on swebench-verified



The failure: why we cannot just boost the signal

My original goal is to see if I can make current benchmarks less noisy

- Still just getting 1 bit of information per hard example.
- Model performance is very noisy and unreliable
 - indirect evidence of memorization when a bad model solves a hard problem
 - unless it's all open problems, there could always be memorization. "we deduplicate all n-grams overlap"
- Modelling difficulty and error of each problem failed to get much signal. Item Response Model failed to do better than baseline.
- Other's attempt to correct errors did not lead to any measurable improvements
 - Subjectively poor quality dataset (MBPP) has better signal to noise than subjectively higher quality dataset (HumanEval+).
 - HumanEval+, MBPP+ is not any better than the original

Attempting to boost signal to noise by filtering, reweighting problems is not promising (because the models are the bigger sources of inconsistency than the benchmarks)

Recommend: Report on more benchmarks. Collect bigger datasets or at least attempt to collect **more information per example** (i.e. > 1 bit) on complex and expensive tasks i.e. SWE-Bench.

To be reconsidered when models are less random!

Conclusions for measuring noise

- run with multiple seeds: <https://arxiv.org/abs/2406.10229v1>
 - usually an underestimate, but can produce significant results. compatible to this work.
- variance of the training curve (underestimate)
- We recommend: in increasing accuracy, and decreasing convenience
 - read off our table
 - read off our figure in your accuracy range
 - Export your results and run our code, provides context
 - benchmark/leaderboard builders should share the full question level results
 - Do your own statistical tests and calculations?? fairly error prone

Other features of [Eval-Arena \(docs\)](#)

Noise level references and signal-to-noise ratio (this talk)

Model [performance ranked](#) by Elo, win-rate, accuracy (best)

Tool for finding suspect [examples](#) (we found more errors / inconsistencies on MBPP+) and some ways. Explore if time allows

Finding suspicious models

Signal to noise

The noise values allows us to interpret reported results. To get a sense of the quality of the benchmark, we need signal to noise ratio.

Signal: the gain from doubling the size of model in each series

Noise: implied standard deviation by statistical testing (same as bootstrapping)

$$\text{sig/noise} = |A-B| / \text{SE}(A-B)$$

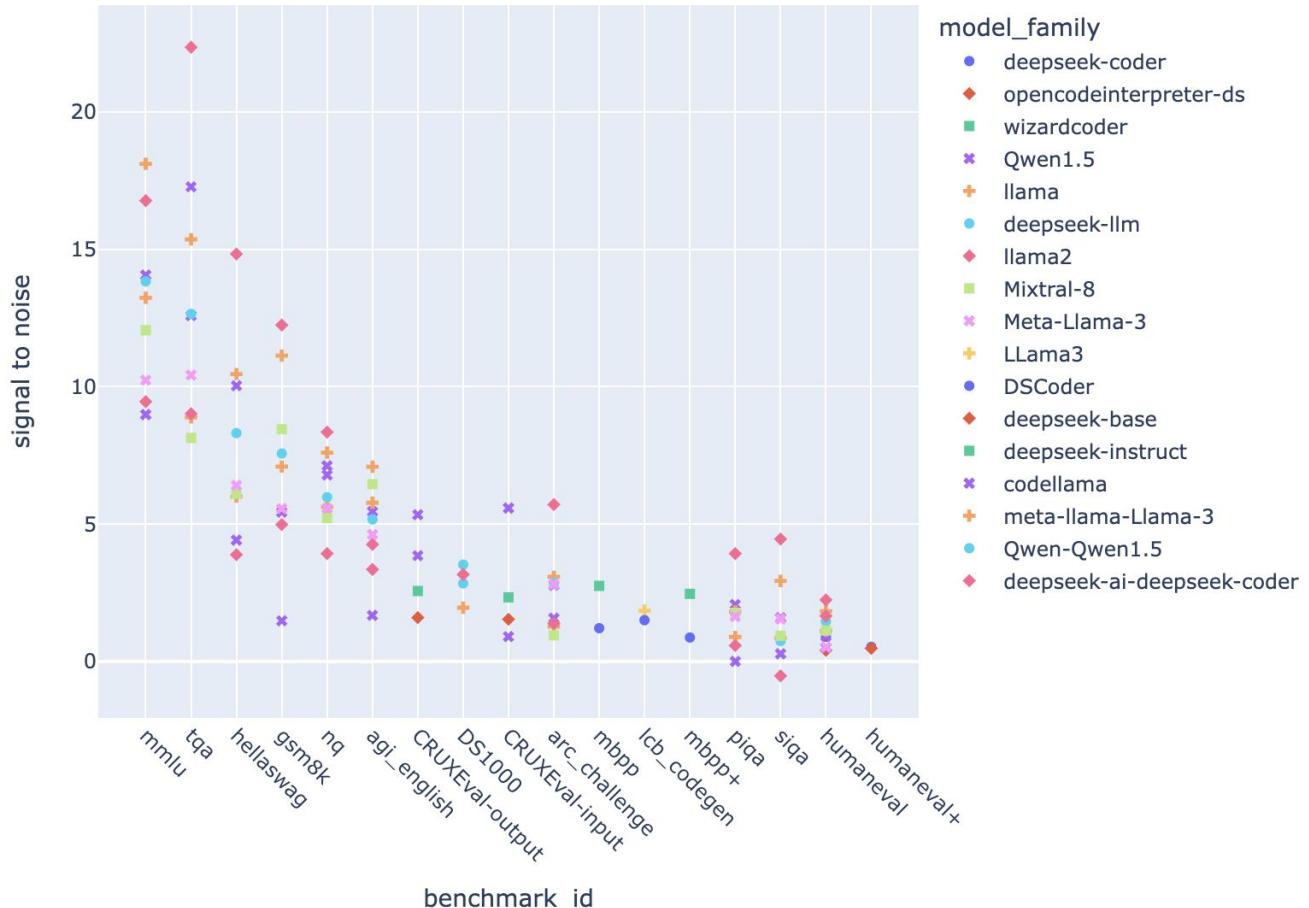
sig/noise < 2: cannot reliably measure the improvement from doubling model size. True on most **code benchmarks!** [examples of inversion](#)

Simple multiple choice benchmarks have much better signal to noise ratio than complex generation benchmarks

benchmark_id	size	sig/noise
mmlu	14042	13.2
tqa	11313	12.6
gsm8k	1319	7.1
hellaswag	10042	6.4
nq	3610	5.9
agi_english	2546	5.1
CRUXEval-output	800	3.2
DS1000	1000	3.0
arc_challenge	1165	2.7
mbpp	378	1.9
CRUXEval-input	800	1.9
pipa	1838	1.7
lcb_codegen	400	1.6
mbpp+	378	1.6
humaneval	164	1.1
siqa	1954	0.93
humaneval+	164	0.50

signal-to-noise

by model class as well



When everyone is too busy to look into the details of the eval because it's too complicated and 100k tokens, bad things can happen

<https://github.com/SWE-bench/SWE-bench/issues/465>

@UniverseFly found it, several others from the team investigated, as mentioned in the issue

Paired comparisons

