

Robust Uncertainty Quantification for LLM Evaluations

Sam Bowyer, Laurence Aitchison, and Desi R. Ivanova



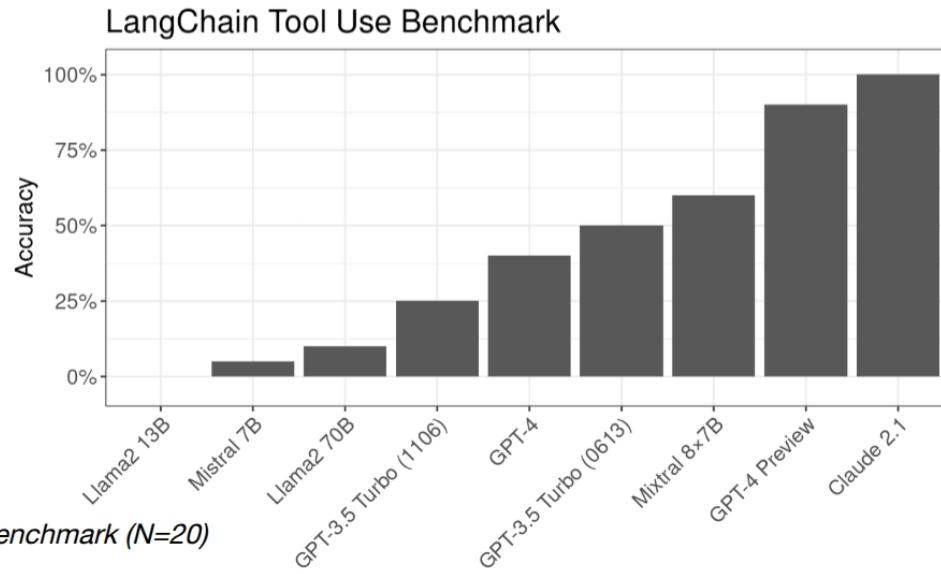
ICML 2025 Spotlight Position Paper

Position: Don't Use the CLT in LLM Evals With Fewer Than a Few Hundred Datapoints

Motivation

Motivation

- Error bars are important for interpreting evals.

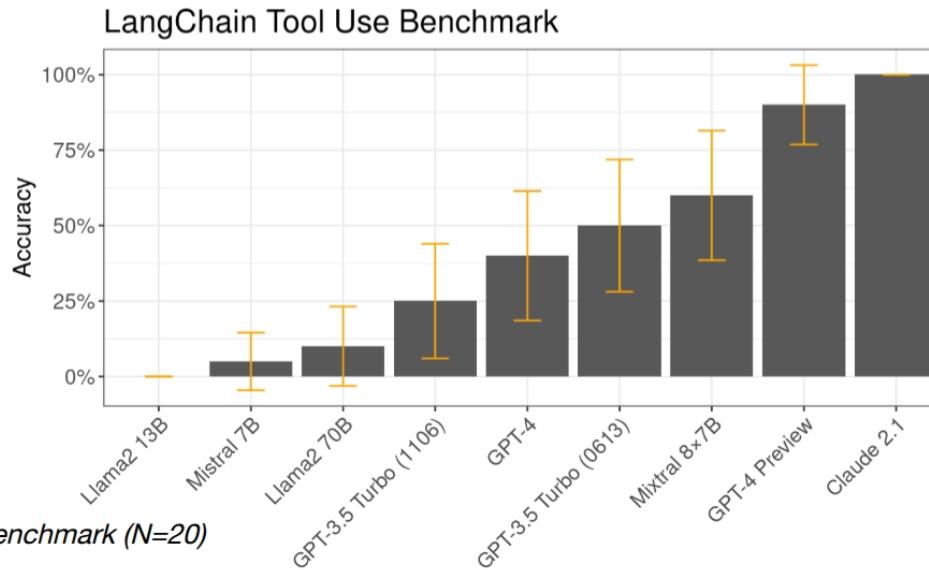


Motivation

CLT-based CI at confidence level $1 - \alpha$ for binary data $X_i \sim \text{Bernoulli}(\theta)$:

$$\text{CI}_{1-\alpha}(\theta) = \bar{X} \pm z_{\alpha/2} \sqrt{\frac{\bar{X}(1 - \bar{X})}{N}}$$

- Error bars are important for interpreting evals.
- The CLT is the most common method for computing error bars, but it's often unwise (assumes large N).

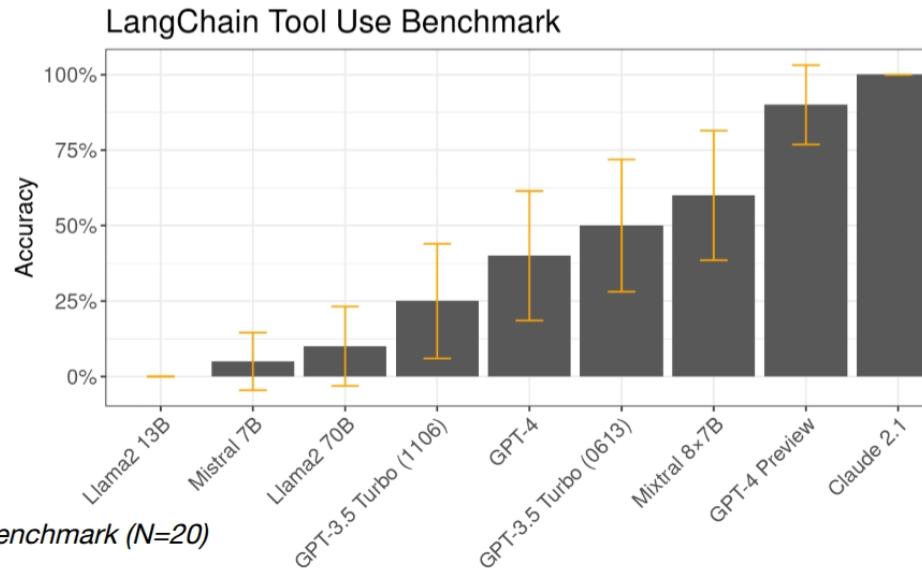


Motivation

CLT-based CI at confidence level $1 - \alpha$ for binary data $X_i \sim \text{Bernoulli}(\theta)$:

$$\text{CI}_{1-\alpha}(\theta) = \bar{X} \pm z_{\alpha/2} \sqrt{\frac{\bar{X}(1 - \bar{X})}{N}}$$

- Error bars are important for interpreting evals.
- The CLT is the most common method for computing error bars, but it's often unwise (assumes large N).
- Error bars can collapse to zero-width or extend past [0, 1].

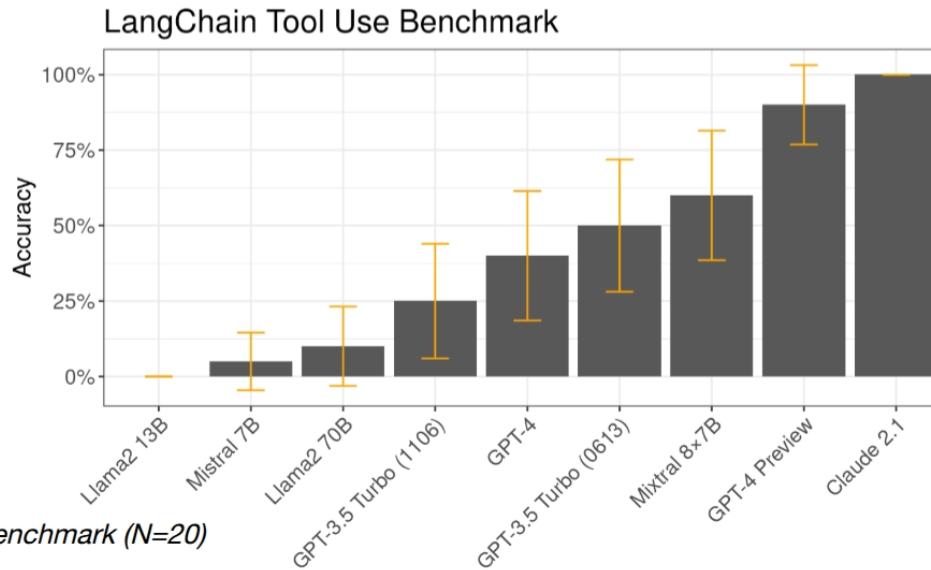


Motivation

CLT-based CI at confidence level $1 - \alpha$ for binary data $X_i \sim \text{Bernoulli}(\theta)$:

$$\text{CI}_{1-\alpha}(\theta) = \bar{X} \pm z_{\alpha/2} \sqrt{\frac{\bar{X}(1 - \bar{X})}{N}}$$

- Error bars are important for interpreting evals.
- The CLT is the most common method for computing error bars, but it's often unwise (assumes large N).
- Error bars can collapse to zero-width or extend past [0, 1].
- Smaller, more intricate and expensive LLM benchmarks are becoming increasingly common.



Bayesian Alternative – Beta-Bernoulli Model

Treat the data as IID Bernoulli with a uniform prior on the parameter θ .

$$\theta \sim \text{Beta}(1, 1) = \text{Uniform}[0, 1]$$

$$y_i \sim \text{Bernoulli}(\theta) \text{ for } i = 1, \dots, N$$

Bayesian Alternative – Beta-Bernoulli Model

Treat the data as IID Bernoulli with a uniform prior on the parameter θ .

$$\theta \sim \text{Beta}(1, 1) = \text{Uniform}[0, 1]$$

$$y_i \sim \text{Bernoulli}(\theta) \text{ for } i = 1, \dots, N$$

We say y_i is correct if $y_i = 1$ and incorrect if $y_i = 0$. (Think of θ as the probability of correctness.)

Bayesian Alternative – Beta-Bernoulli Model

Treat the data as IID Bernoulli with a uniform prior on the parameter θ .

$$\theta \sim \text{Beta}(1, 1) = \text{Uniform}[0, 1]$$

$$y_i \sim \text{Bernoulli}(\theta) \text{ for } i = 1, \dots, N$$

We say y_i is correct if $y_i = 1$ and incorrect if $y_i = 0$. (Think of θ as the probability of correctness.)

$$p(\theta|y_{1:N}) = \text{Beta}\left(1 + \sum_{i=1}^N y_i, 1 + \sum_{i=1}^N (1 - y_i)\right)$$

Bayesian Alternative – Beta-Bernoulli Model

Treat the data as IID Bernoulli with a uniform prior on the parameter θ .

$$\theta \sim \text{Beta}(1, 1) = \text{Uniform}[0, 1]$$

$$y_i \sim \text{Bernoulli}(\theta) \text{ for } i = 1, \dots, N$$

We say y_i is correct if $y_i = 1$ and incorrect if $y_i = 0$. (Think of θ as the probability of correctness.)

$$p(\theta|y_{1:N}) = \text{Beta}\left(1 + \sum_{i=1}^N y_i, 1 + \sum_{i=1}^N (1 - y_i)\right)$$

Obtain quantile-based Bayesian *credible intervals* for θ from the **closed form posterior** (with confidence level $1 - \alpha$).

Bayesian Alternative – Beta-Bernoulli Model

Treat the data as IID Bernoulli with a uniform prior on the parameter θ .

$$\theta \sim \text{Beta}(1, 1) = \text{Uniform}[0, 1]$$

$$y_i \sim \text{Bernoulli}(\theta) \text{ for } i = 1, \dots, N$$

We say y_i is correct if $y_i = 1$ and incorrect if $y_i = 0$. (Think of θ as the probability of correctness.)

$$p(\theta|y_{1:N}) = \text{Beta}\left(1 + \sum_{i=1}^N y_i, 1 + \sum_{i=1}^N (1 - y_i)\right)$$

Obtain quantile-based Bayesian *credible intervals* for θ from the **closed form posterior** (with confidence level $1 - \alpha$).

```
# y is a length N binary "eval" vector
S, N = y.sum(), len(y) # total successes & questions

# Bayesian Credible interval
posterior = scipy.stats.beta(1+S, 1+(N-S))
bayes_ci = posterior.interval(confidence=0.95)
```

Frequentist Alternatives

Frequentist Alternatives

- Wilson score interval

Frequentist Alternatives

- Wilson score interval
 - Based on the normal approximation to the binomial distribution (but **not** the CLT).

Frequentist Alternatives

- Wilson score interval
 - Based on the normal approximation to the binomial distribution (but **not** the CLT).
- Clopper-Pearson exact interval

Frequentist Alternatives

- Wilson score interval
 - Based on the normal approximation to the binomial distribution (but **not** the CLT).
- Clopper-Pearson exact interval
 - 'Worst-case' approach (very conservative method; guaranteed to never under-cover).

Frequentist Alternatives

- Wilson score interval
 - Based on the normal approximation to the binomial distribution (but **not** the CLT).
- Clopper-Pearson exact interval
 - 'Worst-case' approach (very conservative method; guaranteed to never under-cover).

```
# y is a length N binary "eval" vector
S, N = y.sum(), len(y) # total successes & questions
result = scipy.stats.binomtest(k=S, n=N)

# 95% Wilson score interval and Clopper-Pearson exact interval
wilson_ci = result.proportion_ci("wilson", 0.95)
cp_ci = result.proportion_ci("exact", 0.95)
```

Interval Comparison Simulations

We have to rely on synthetic eval data so that we *know* the true parameter θ .

Interval Comparison Simulations

We have to rely on synthetic eval data so that we *know* the true parameter θ .

- Draw $\theta \sim \text{Uniform}[0, 1]$.

Interval Comparison Simulations

We have to rely on synthetic eval data so that we *know* the true parameter θ .

- Draw $\theta \sim \text{Uniform}[0, 1]$.
- Draw $N \in \{3, 10, 30, 100\}$ IID Bernoulli datapoints with parameter θ .

Interval Comparison Simulations

We have to rely on synthetic eval data so that we *know* the true parameter θ .

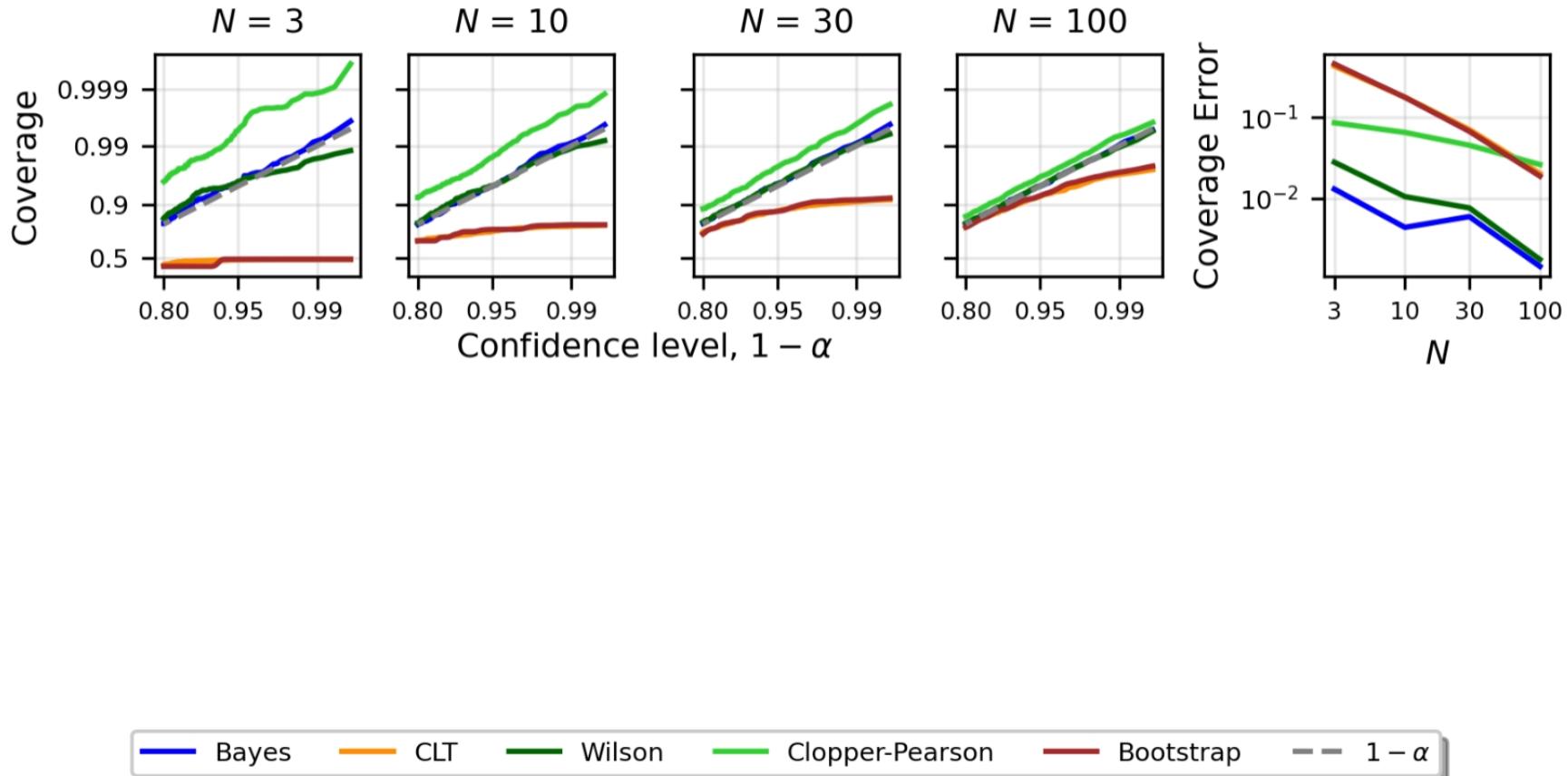
- Draw $\theta \sim \text{Uniform}[0, 1]$.
- Draw $N \in \{3, 10, 30, 100\}$ IID Bernoulli datapoints with parameter θ .
- Construct intervals with various methods for various $1 - \alpha$ confidence levels.

Interval Comparison Simulations

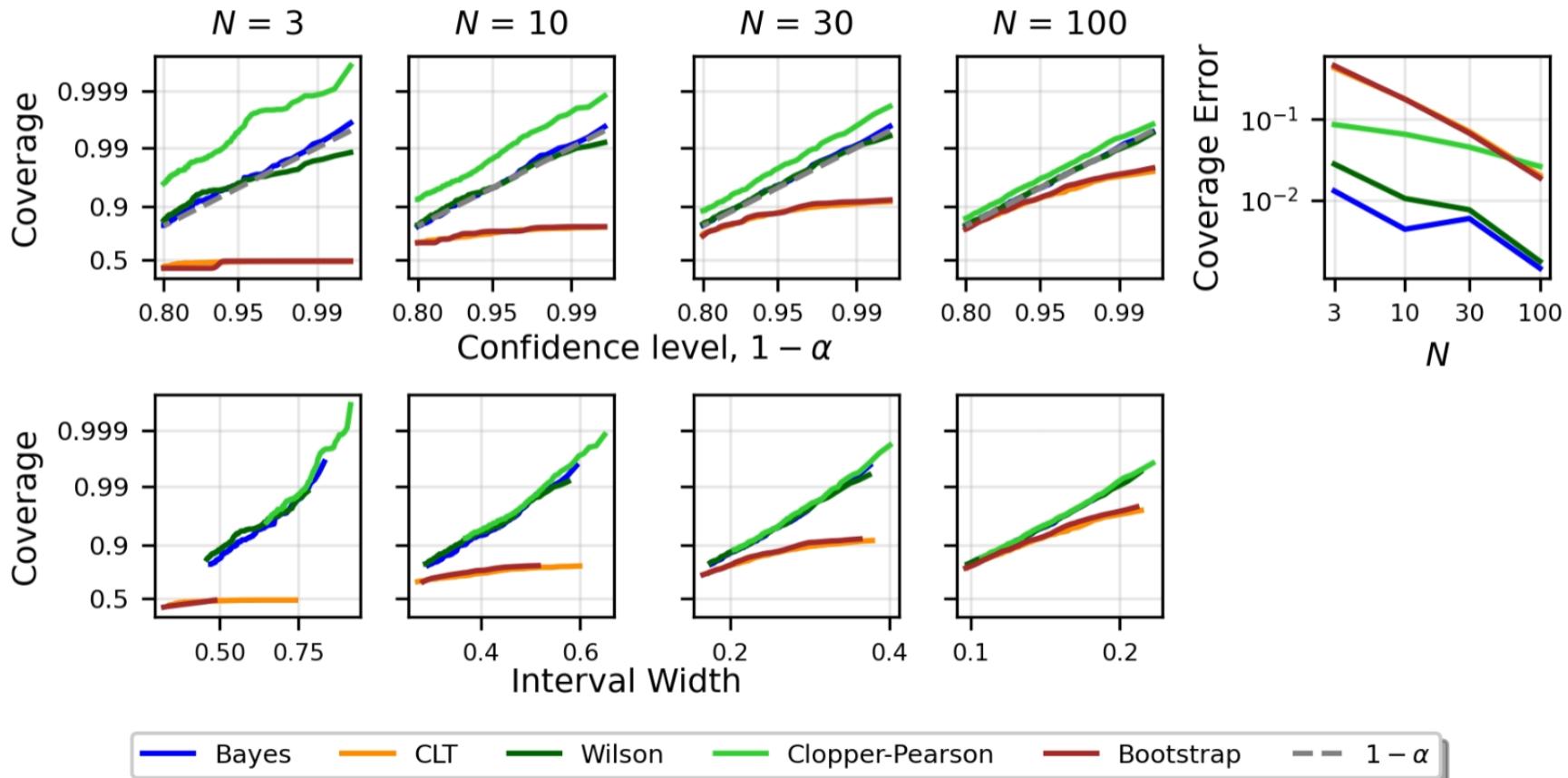
We have to rely on synthetic eval data so that we *know* the true parameter θ .

- Draw $\theta \sim \text{Uniform}[0, 1]$.
- Draw $N \in \{3, 10, 30, 100\}$ IID Bernoulli datapoints with parameter θ .
- Construct intervals with various methods for various $1 - \alpha$ confidence levels.
- Repeat many times and calculate the true coverage and width of the intervals.

IID Questions Setting



IID Questions Setting



Other Eval Settings

Other Eval Settings

Clustered Questions

Other Eval Settings

Clustered Questions

Instead of N IID questions, we have T tasks, each with N_t IID questions.

Other Eval Settings

Clustered Questions

Instead of N IID questions, we have T tasks, each with N_t IID questions.

Comparisons Between Two Models, θ_A and θ_B

Other Eval Settings

Clustered Questions

Instead of N IID questions, we have T tasks, each with N_t IID questions.

Comparisons Between Two Models, θ_A and θ_B

- **Independent Comparisons:** Using N_A , N_B , \bar{y}_A , and \bar{y}_B .
- **Paired Comparisons:** Using $N_A = N_B$, $\{y_{A;i}\}_{i=1}^N$, and $\{y_{B;i}\}_{i=1}^N$.

Other Eval Settings

Clustered Questions

Instead of N IID questions, we have T tasks, each with N_t IID questions.

Comparisons Between Two Models, θ_A and θ_B

- **Independent Comparisons:** Using N_A , N_B , \bar{y}_A , and \bar{y}_B .
- **Paired Comparisons:** Using $N_A = N_B$, $\{y_{A;i}\}_{i=1}^N$, and $\{y_{B;i}\}_{i=1}^N$.

Metrics that aren't simple averages of binary results (e.g. F1 score)

Other Eval Settings

Clustered Questions

Instead of N IID questions, we have T tasks, each with N_t IID questions.

Comparisons Between Two Models, θ_A and θ_B

- **Independent Comparisons:** Using N_A , N_B , \bar{y}_A , and \bar{y}_B .
- **Paired Comparisons:** Using $N_A = N_B$, $\{y_{A;i}\}_{i=1}^N$, and $\{y_{B;i}\}_{i=1}^N$.

Metrics that aren't simple averages of binary results (e.g. F1 score)

Also, what if the prior is mismatched? (i.e. $\theta \sim \text{Uniform}[0, 1]$)

Conclusion

Advice to practitioners who might not be so familiar with stats:

Conclusion

Advice to practitioners who might not be so familiar with stats:

- Use Bayesian Beta-Bernoulli or Wilson Score intervals.

Conclusion

Advice to practitioners who might not be so familiar with stats:

- Use Bayesian Beta-Bernoulli or Wilson Score intervals.
- It's not hard (use `scipy` or `bayes_evals`).

Conclusion

Advice to practitioners who might not be so familiar with stats:

- Use Bayesian Beta-Bernoulli or Wilson Score intervals.
- It's not hard (use `scipy` or `bayes_evals`).
- It's safer than CLT-based methods.

Conclusion

Advice to practitioners who might not be so familiar with stats:

- Use Bayesian Beta-Bernoulli or Wilson Score intervals.
- It's not hard (use `scipy` or `bayes_evals`).
- It's safer than CLT-based methods.
- It's still cheap for large N .

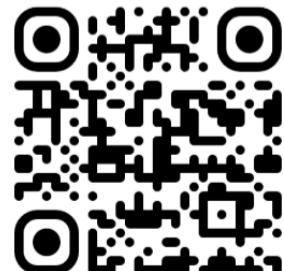
Conclusion

Advice to practitioners who might not be so familiar with stats:

- Use Bayesian Beta-Bernoulli or Wilson Score intervals.
- It's not hard (use `scipy` or `bayes_evals`).
- It's safer than CLT-based methods.
- It's still cheap for large N .

Paper

<https://arxiv.org/pdf/2503.01747>



bayes_evals package

https://github.com/sambowyer/bayes_evals

