

A Statistical Approach to Language Model Evaluations

fzeng, 2024-11-18

Challenges with LLM Evals

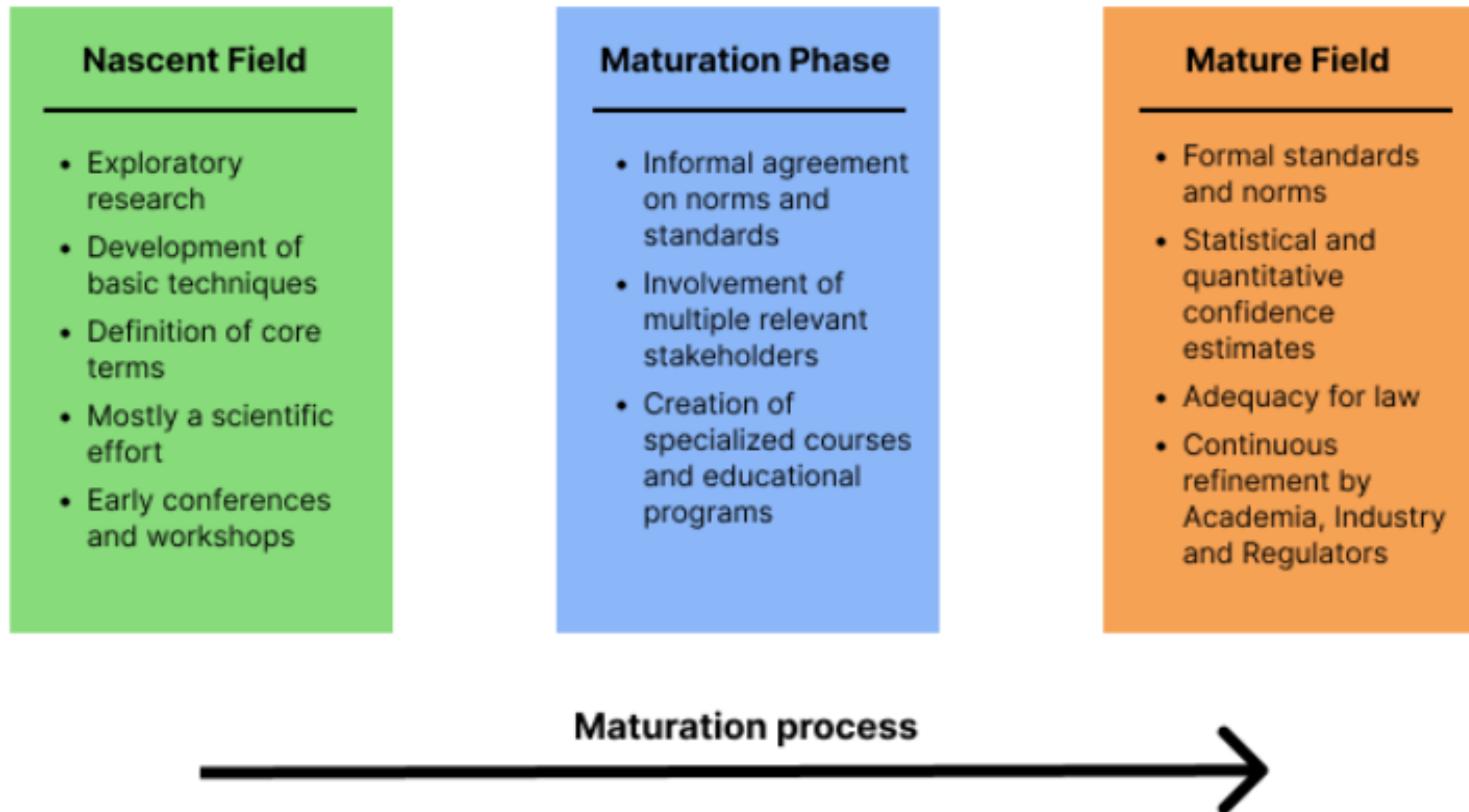
- LLM highly sensitive to prompts ([Liang et al., 2022](#); [Mizrahi et al., 2023](#); [Scalar et al., 2023](#); [Weber et al., 2023](#), [Bsharat et al., 2023](#))
- Several widely used open-source LLMs extremely sensitive to subtle changes in prompt formatting in few-shot settings, with performance differences of up to 76 accuracy points ([Scalar et al., 2023](#))
- Changing the options from (A) to (1) or changing the parentheses from (A) to [A], or adding an extra space between the option and the answer can lead to a ~5 percentage point change in accuracy on the evaluation ([Anthropic](#))
- Tipping a language model 300K for a better solution" leads to increased capabilities ([Bsharat et al., 2023](#))

Papers often don't report standard errors

| Category | Benchmark | Llama 3 8B | Gemma 2 9B | Mistral 7B | Llama 3 70B | Mixtral 8x22B | GPT 3.5 Turbo | Llama 3 405B | Nemotron 4 340B | GPT-4 (0125) | GPT-4o | Claude 3.5 Sonnet |
|--------------|------------------------|-------------|-------------------|------------|-------------|---------------|---------------|--------------|-------------------|--------------|--------------|-------------------|
| General | MMLU (5-shot) | 69.4 | 72.3 | 61.1 | 83.6 | 76.9 | 70.7 | 87.3 | 82.6 | 85.1 | 89.1 | 89.9 |
| | MMLU (0-shot, CoT) | 73.0 | 72.3 [△] | 60.5 | 86.0 | 79.9 | 69.8 | 88.6 | 78.7 [◁] | 85.4 | 88.7 | 88.3 |
| | MMLU-Pro (5-shot, CoT) | 48.3 | – | 36.9 | 66.4 | 56.3 | 49.2 | 73.3 | 62.7 | 64.8 | 74.0 | 77.0 |
| | IFEval | 80.4 | 73.6 | 57.6 | 87.5 | 72.7 | 69.9 | 88.6 | 85.1 | 84.3 | 85.6 | 88.0 |
| Code | HumanEval (0-shot) | 72.6 | 54.3 | 40.2 | 80.5 | 75.6 | 68.0 | 89.0 | 73.2 | 86.6 | 90.2 | 92.0 |
| | MBPP EvalPlus (0-shot) | 72.8 | 71.7 | 49.5 | 86.0 | 78.6 | 82.0 | 88.6 | 72.8 | 83.6 | 87.8 | 90.5 |
| Math | GSM8K (8-shot, CoT) | 84.5 | 76.7 | 53.2 | 95.1 | 88.2 | 81.6 | 96.8 | 92.3 [◇] | 94.2 | 96.1 | 96.4 [◇] |
| | MATH (0-shot, CoT) | 51.9 | 44.3 | 13.0 | 68.0 | 54.1 | 43.1 | 73.8 | 41.1 | 64.5 | 76.6 | 71.1 |
| Reasoning | ARC Challenge (0-shot) | 83.4 | 87.6 | 74.2 | 94.8 | 88.7 | 83.7 | 96.9 | 94.6 | 96.4 | 96.7 | 96.7 |
| | GPQA (0-shot, CoT) | 32.8 | – | 28.8 | 46.7 | 33.3 | 30.8 | 51.1 | – | 41.4 | 53.6 | 59.4 |
| Tool use | BFCL | 76.1 | – | 60.4 | 84.8 | – | 85.9 | 88.5 | 86.5 | 88.3 | 80.5 | 90.2 |
| | Nexus | 38.5 | 30.0 | 24.7 | 56.7 | 48.5 | 37.2 | 58.7 | – | 50.3 | 56.1 | 45.7 |
| Long context | ZeroSCROLLS/QuALITY | 81.0 | – | – | 90.5 | – | – | 95.2 | – | 95.2 | 90.5 | 90.5 |
| | InfiniteBench/En.MC | 65.1 | – | – | 78.2 | – | – | 83.4 | – | 72.1 | 82.5 | – |
| | NIH/Multi-needle | 98.8 | – | – | 97.5 | – | – | 98.1 | – | 100.0 | 100.0 | 90.8 |
| Multilingual | MGSM (0-shot, CoT) | 68.9 | 53.2 | 29.9 | 86.9 | 71.1 | 51.4 | 91.6 | – | 85.9 | 90.5 | 91.6 |

| Model | HumanEval | HumanEval+ | MBPP | MBPP EvalPlus (base) |
|-------------------|------------------|------------------|------------------|----------------------|
| Llama 3 8B | 72.6 ±6.8 | 67.1 ±7.2 | 60.8 ±4.3 | 72.8 ±4.5 |
| Gemma 2 9B | 54.3 ±7.6 | 48.8 ±7.7 | 59.2 ±4.3 | 71.7 ±4.5 |
| Mistral 7B | 40.2 ±7.5 | 32.3 ±7.2 | 42.6 ±4.3 | 49.5 ±5.0 |
| Llama 3 70B | 80.5 ±6.1 | 74.4 ±6.7 | 75.4 ±3.8 | 86.0 ±3.5 |
| Mixtral 8×22B | 75.6 ±6.6 | 68.3 ±7.1 | 66.2 ±4.1 | 78.6 ±4.1 |
| GPT-3.5 Turbo | 68.0 ±7.1 | 62.8 ±7.4 | 71.2 ±4.0 | 82.0 ±3.9 |
| Llama 3 405B | 89.0 ±4.8 | 82.3 ±5.8 | 78.8 ±3.6 | 88.6 ±3.2 |
| GPT-4 | 86.6 ±5.2 | 77.4 ±6.4 | 80.2 ±3.5 | 83.6 ±3.7 |
| GPT-4o | 90.2 ±4.5 | 86.0 ±5.3 | 81.4 ±3.4 | 87.8 ±3.3 |
| Claude 3.5 Sonnet | 92.0 ±4.2 | 82.3 ±5.8 | 76.6 ±3.7 | 90.5 ±3.0 |
| Nemotron 4 340B | 73.2 ±6.8 | 64.0 ±7.3 | 75.4 ±3.8 | 72.8 ±4.5 |

Science of Evals still young



Paper recommendations

1. Computing standard errors of the mean using the Central Limit Theorem
2. When questions are drawn in related groups, computing clustered standard errors
3. Reducing variance by resampling answers and by analyzing next-token probabilities
4. When two models are being compared, conducting statistical inference on the question level paired differences, rather than the population-level summary statistics
5. Using power analysis to determine whether an eval (or a random subsample) is capable of testing a hypothesis of interest

Standard Errors of the Mean

Standard Error of the Mean

- Some notation:
- For some question i in the dataset,

$$\underbrace{s_i}_{\text{score of question } i} = \underbrace{x_i}_{\text{conditional mean on question } i} + \underbrace{\epsilon_i}_{\text{conditional variance on question } i}$$

- Can also talk about any question in the dataset unconditionally: $s = x + \epsilon$

- Mean of scores: $\bar{s} = \frac{1}{n} \sum_i s_i$

Standard Error of the Mean

- Our scores can come from any distribution; how can we say anything about error bounds if we don't know this distribution?
- CLT to the rescue!
- CLT says mean of i.i.d random variables with finite mean and variance converges can be approximated with standard normal

Central Limit Theorem: Let Y_1, Y_2, \dots, Y_n be independent and identically distributed random variables with $E(Y_i) = \mu$ and $V(Y_i) = \sigma^2 < \infty$. Define

$$U_n = \frac{\sum_{i=1}^n Y_i - n\mu}{\sigma\sqrt{n}} = \frac{\bar{Y} - \mu}{\sigma/\sqrt{n}} \quad \text{where } \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i.$$

Then the distribution function of U_n converges to the standard normal distribution function as $n \rightarrow \infty$. That is,

$$\lim_{n \rightarrow \infty} P(U_n \leq u) = \int_{-\infty}^u \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt \quad \text{for all } u.$$

Standard Error of the Mean

- So the estimate of our mean can be transformed into a standard normal
- We can then also get unbiased estimator of sample variance:

$$\text{Var}(s) = \frac{1}{n-1} \sum_i (s_i - \bar{s})^2$$

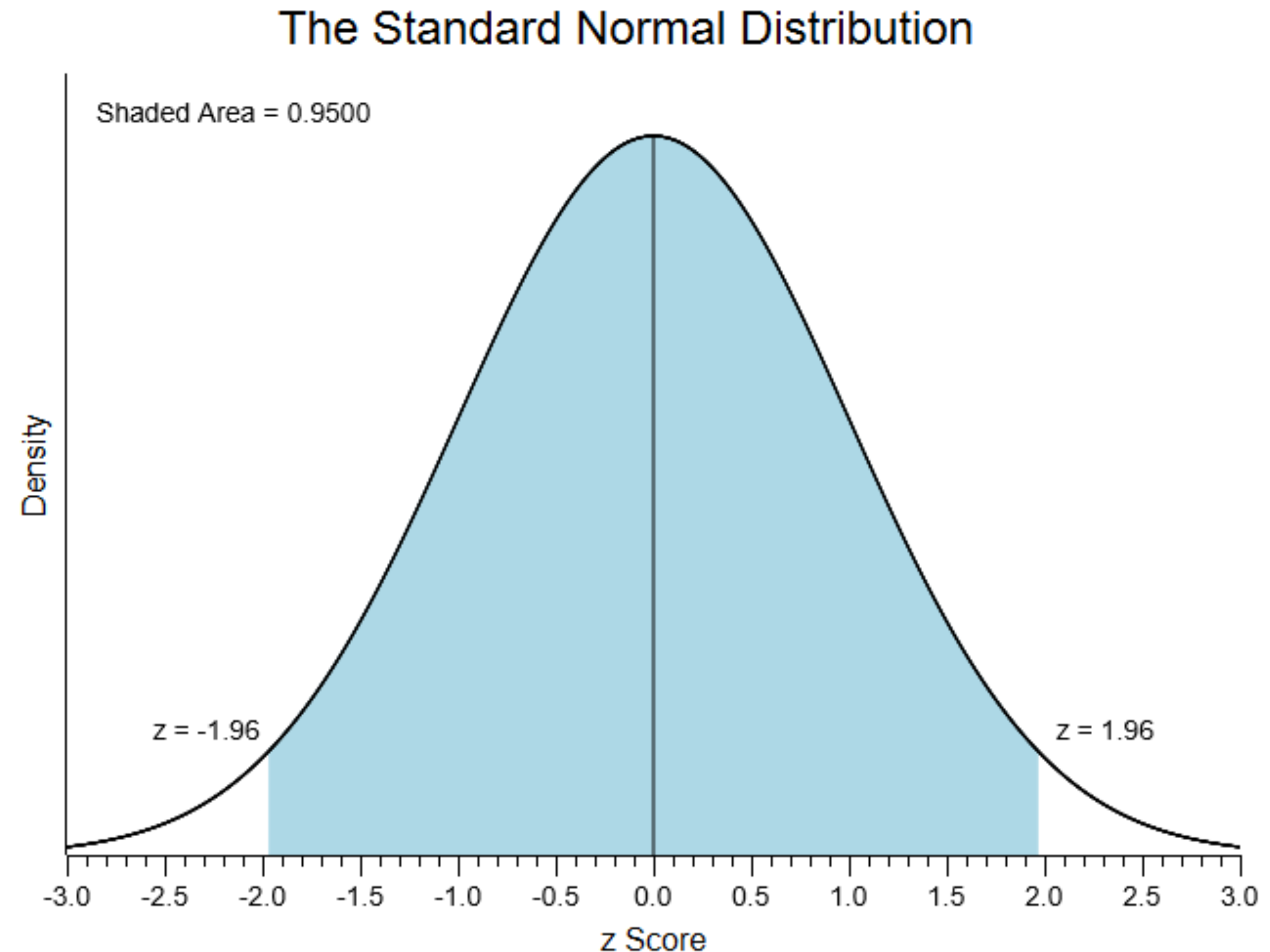
- For n samples, by linearity of variance we recover

$$\text{SE}_{\text{C.L.T.}} = \sqrt{\text{Var}(s)/n} = \sqrt{\left(\frac{1}{n-1} \sum_i (s_i - \bar{s})^2 \right) / n} \quad (1)$$

Standard Error of the Mean

- Using maximum likelihood estimator (MLE), we declare \bar{s} to be the estimate of population mean, and draw a 95% confidence interval around it (1.96 sigma)
- Recovers Eq (3):

$$CI_{95\%} = \bar{s} \pm 1.96 \times SE_{C.L.T.}$$



Clustered Standard Errors

Clustered Standard Error

- CLT requires i.i.d assumption
- Some datasets are clearly not i.i.d
- MGSM (Multilingual Grade-School Math):
 - 2500 grade-school math questions
 - But really: 250 questions translated into 10 different languages
 - 250 clusters of 10

Clustered Standard Error

$$SE_{C.L.T.} = \sqrt{\text{Var}(s)/n}$$

- Why does it fail if observations not i.i.d?
 - “Effective” number of observations much fewer than 2500, probably more like 250
- Case 1: observation in each cluster is iid (implies 0 covariance)
- Then

$$SE_{\text{clustered}} = \sqrt{SE_{C.L.T.}^2 + \underbrace{\frac{1}{n^2} \sum_c \sum_i \sum_{j \neq i} (s_{i,c} - \bar{s}) (s_{j,c} - \bar{s})}_{=0}}$$

Clustered Standard Error

- Case 2: observation in each cluster perfectly correlated
- Then

$$SE_{\text{clustered}} = \sqrt{SE_{\text{C.L.T.}}^2 + \frac{1}{n^2} \sum_c \sum_i \sum_{j \neq i} (s_{i,c} - \bar{s})^2}$$

and you add back variance contributions within each cluster

$$SE_{\text{C.L.T.}} = \sqrt{\text{Var}(s)/n} = \sqrt{\left(\frac{1}{n-1} \sum_i (s_i - \bar{s})^2 \right) / n}$$

Recommendation for reporting errors

| | # Questions | # Clusters | “Galleon” | “Dreadnought” |
|--------|-------------|------------|-----------------|-----------------|
| DROP | 9,622 | 588 | 87.1 (0.8) | 83.1 (0.9) |
| RACE-H | 3,498 | 1,045 | 91.5% (0.5%) | 82.9% (0.7%) |
| MGSM | 2,500 | 250 | 75.3% (1.6%) | 78.0% (1.5%) |

Table 3: We suggest including the cluster count alongside the question count when reporting cluster-adjusted standard errors (fictional models and numbers).

| | $SE_{\text{clustered}}$ | $SE_{\text{C.L.T.}}$ | Ratio |
|--------|-------------------------|----------------------|-------|
| DROP | (1.34) | (0.44) | 3.05 |
| RACE-H | (0.51%) | (0.46%) | 1.10 |
| MGSM | (1.62%) | (0.86%) | 1.88 |

Table 4: Clustered and naive standard errors computed on two popular evals using Anthropic models (non-fictional numbers). Analyzing the same data, clustered standard errors can be over 3X larger than naive standard errors.

Variance Reduction

Variance Reduction

- $\text{Var}(\hat{\mu}) = \text{Var}\left(\frac{1}{n} \sum_i s_i\right) = \text{Var}(s)/n$
- Increase number of samples directly reduces variance
- But we still have another trick..

Law of Total Variance

- This is tricky to get intuition on
- $\text{Var}(Y) = E[\text{Var}(Y | X)] + \text{Var}(E[Y | X])$
- Example: Y is dog's weight, X is breed
- First term: avg of variance of weight within each breed (within-group variance)
- Second term: variance of avg of each breed (between-group variance)

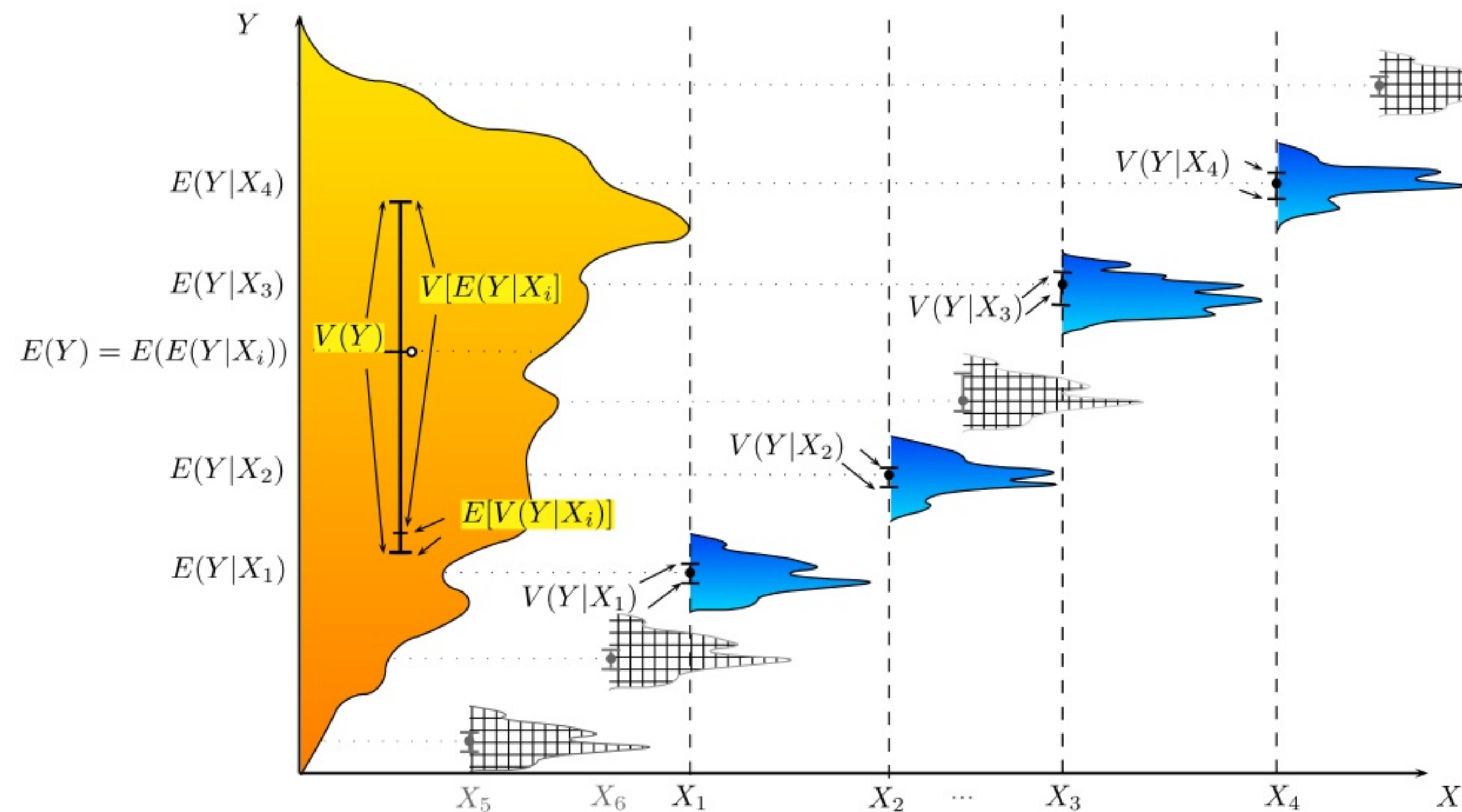


Figure 3: ANOVA : very good fit

Variance Reduction

- FYI: I don't like their notation for this part, very imprecise

- $$\text{Var}(s) = \underbrace{\text{Var}(\mathbb{E}[x_i | i])}_{\text{variance in scores across different questions}} + \underbrace{\mathbb{E}[\text{Var}(x_i | i)]}_{\text{variance in scores from answering the same question across different attempts}}$$

- Let's consider resampling
- Resampling won't help the first term - this is inherent in the distribution of questions
- But it can help to decrease the second term: sampling n times & taking mean will reduce it by n
- Increasing n is economical until the point that second term is same size as first term (then first term dominates)

Variance Reduction

- $$\text{Var}(s) = \underbrace{\text{Var}(\mathbb{E}[x_i | i])}_{\text{variance in scores across different questions}} + \underbrace{\mathbb{E}[\text{Var}(x_i | i)]}_{\text{variance in scores from answering the same question across different attempts}}$$

- Tempting thing to eliminate second term: set temp=0

3.3 Don't touch the thermostat!

It may be tempting to reduce the “sampling temperature” [10] of the model in order to reduce (or eliminate) the conditional variance. However, we advise against this practice, unless the purpose is to study the model at the new temperature. Besides altering the model’s behavior, adjusting the sampling temperature may simply shift the conditional variance (which can be mitigated using the two techniques above) into the variance of the conditional means (which cannot), or else reduce conditional variance by injecting bias into the estimator. Two short examples will

- illustrate these points.

- Their example: setting T=0 increases first term

Variance Reduction

- $$\text{Var}(s) = \underbrace{\text{Var}(\mathbb{E}[x_i | i])}_{\text{variance in scores across different questions}} + \underbrace{\mathbb{E}[\text{Var}(x_i | i)]}_{\text{variance in scores from answering the same question across different attempts}}$$
- For problems where you can use model logprobs to get probability of correct answer (i.e true/false qn), here's another trick:
- Instead of sampling the answer token & giving a binary score, return the probability of correct answer as score
- Then second term becomes 0 😊

To be continued next session...