

Rare-event estimation for Bayesian predictive rollouts: rollout Monte Carlo vs. posterior sampling (Rao–Blackwellization)

Vinayak Pathak

Abstract

We study the probability of a rare event A defined on the future outputs of a Bayesian predictive model after conditioning on observed data $y_{1:n}$. Two Monte Carlo strategies are compared: (i) direct sampling of future rollouts from the posterior predictive and counting hits of A (an indicator estimator), and (ii) sampling parameters θ from the posterior and computing $f(\theta) = \mathbb{P}(A \mid \theta, y_{1:n})$ analytically (a Rao–Blackwellized estimator). We give a self-contained multiplicative Chernoff bound for bounded random variables, derive high-probability relative-error sample complexity for both methods, and characterize when an exponential gap is possible. We also record extensions beyond iid Bernoulli to multinomials and HMM/state-space models.

1 Bayesian predictive distributions and autoregressive rollouts

1.1 General Bayesian setup

Let $\theta \in \Theta$ be a parameter with prior $\pi_0(d\theta)$. Let $Y_{1:N}$ be observations with likelihood

$$p_\theta(y_{1:N}) = p(y_{1:N} \mid \theta).$$

After observing $D := y_{1:n}$, the posterior is

$$\pi(d\theta \mid D) \propto p_\theta(D) \pi_0(d\theta).$$

Fix a future horizon $m := N - n$ and consider the future block $Y_{n+1:n+m}$.

Theorem 1 (Posterior predictive mixture / marginalization identity). *For any measurable set B in the space of length- m sequences,*

$$\mathbb{P}(Y_{n+1:n+m} \in B \mid D) = \int \mathbb{P}_\theta(Y_{n+1:n+m} \in B \mid D) \pi(d\theta \mid D).$$

Proof. By the law of total probability w.r.t. θ ,

$$\mathbb{P}(Y_{n+1:n+m} \in B \mid D) = \int \mathbb{P}(Y_{n+1:n+m} \in B \mid \theta, D) \pi(d\theta \mid D),$$

and $\mathbb{P}(Y_{n+1:n+m} \in B \mid \theta, D) = \mathbb{P}_\theta(Y_{n+1:n+m} \in B \mid D)$ by definition of the model under θ . \square

Remark 1 (Autoregressive sampling vs. block sampling). If one sequentially samples

$$Y_{n+t} \sim \mathbb{P}(\cdot \mid D, Y_{n+1:n+t-1}), \quad t = 1, 2, \dots, m,$$

then by the chain rule the joint distribution of $Y_{n+1:n+m}$ equals the posterior predictive $\mathbb{P}(\cdot \mid D)$. Thus “autoregressive rollout” (sampling from one-step conditionals and feeding back) and “block sampling” from $\mathbb{P}(Y_{n+1:n+m} \mid D)$ are equivalent ways to generate future sequences.

2 Rare events and two estimators

Definition 1 (Rare-event probability). Fix an event A measurable w.r.t. the future block $Y_{n+1:n+m}$. Define

$$q := \mathbb{P}(A \mid D).$$

For each θ , define the conditional event probability

$$f(\theta) := \mathbb{P}_\theta(A \mid D).$$

Corollary 1 (Rare-event probability as a posterior expectation).

$$q = \mathbb{E}_{\theta \sim \pi(\cdot \mid D)}[f(\theta)].$$

Proof. Apply Theorem 1 with $B = A$. □

2.1 Estimator 1: rollout hit-rate

Sample R independent rollouts $Y_{n+1:n+m}^{(r)} \sim \mathbb{P}(\cdot \mid D)$ and define

$$I_r := \mathbf{1}\{Y_{n+1:n+m}^{(r)} \in A\}, \quad \hat{q}_{\text{roll}} := \frac{1}{R} \sum_{r=1}^R I_r.$$

Then $I_r \sim \text{Bernoulli}(q)$ and $\mathbb{E}[\hat{q}_{\text{roll}}] = q$.

2.2 Estimator 2: posterior sampling + analytic $f(\theta)$ (Rao–Blackwell)

Sample $\theta_1, \dots, \theta_M \stackrel{\text{iid}}{\sim} \pi(\cdot \mid D)$ and compute

$$X_i := f(\theta_i), \quad \hat{q}_{\text{post}} := \frac{1}{M} \sum_{i=1}^M X_i.$$

Then $\mathbb{E}[\hat{q}_{\text{post}}] = q$. In what follows we assume $f(\theta)$ can be computed exactly (or to negligible error) given θ .

Lemma 1 (Variance decomposition / Rao–Blackwell). *Let $I := \mathbf{1}\{A\}$ be the indicator of A under a single posterior predictive draw. Then*

$$\text{Var}(I) = \mathbb{E}[\text{Var}(I \mid \theta)] + \text{Var}(\mathbb{E}[I \mid \theta]),$$

and since $\mathbb{E}[I \mid \theta] = f(\theta)$ we have

$$\text{Var}(f(\theta)) \leq \text{Var}(I) = q(1 - q).$$

Proof. This is the law of total variance applied to (I, θ) , using $f(\theta) = \mathbb{P}(A \mid \theta, D) = \mathbb{E}[I \mid \theta]$. □

3 A Chernoff bound for bounded random variables (with proof)

We will use a multiplicative Chernoff bound that holds for *any* independent bounded random variables in $[0, 1]$, not just Bernoullis.

Lemma 2 (MGF domination by a Bernoulli). *Let $Z \in [0, 1]$ be a random variable with $\mathbb{E}[Z] = \mu$. Then for any $\lambda \in \mathbb{R}$,*

$$\mathbb{E}[e^{\lambda Z}] \leq \exp(\mu(e^\lambda - 1)).$$

Proof. The function $x \mapsto e^{\lambda x}$ is convex, so for $x \in [0, 1]$,

$$e^{\lambda x} \leq (1 - x)e^0 + xe^\lambda = 1 + x(e^\lambda - 1).$$

Taking expectations and using $\mathbb{E}[Z] = \mu$ gives

$$\mathbb{E}[e^{\lambda Z}] \leq 1 + \mu(e^\lambda - 1) \leq \exp(\mu(e^\lambda - 1)),$$

where the last inequality uses $1 + u \leq e^u$. □

Theorem 2 (Multiplicative Chernoff bound for unit-interval variables). *Let Z_1, \dots, Z_M be independent random variables with $Z_i \in [0, 1]$ and $\mathbb{E}[Z_i] = \mu$. Let $\bar{Z} := \frac{1}{M} \sum_{i=1}^M Z_i$. Then for any $\rho \in (0, 1)$,*

$$\mathbb{P}(\bar{Z} \geq (1 + \rho)\mu) \leq \exp\left(-\mu M((1 + \rho) \ln(1 + \rho) - \rho)\right) \leq \exp\left(-\frac{\rho^2 \mu M}{3}\right),$$

and

$$\mathbb{P}(\bar{Z} \leq (1 - \rho)\mu) \leq \exp\left(-\mu M(\rho + (1 - \rho) \ln(1 - \rho))\right) \leq \exp\left(-\frac{\rho^2 \mu M}{2}\right).$$

Consequently,

$$\mathbb{P}(|\bar{Z} - \mu| \geq \rho\mu) \leq 2 \exp\left(-\frac{\rho^2 \mu M}{3}\right).$$

Proof. Let $S := \sum_{i=1}^M Z_i$, so $\mathbb{E}[S] = \mu M$.

Upper tail. For $\lambda > 0$, Markov's inequality gives

$$\mathbb{P}(S \geq (1 + \rho)\mu M) = \mathbb{P}(e^{\lambda S} \geq e^{\lambda(1 + \rho)\mu M}) \leq \frac{\mathbb{E}[e^{\lambda S}]}{e^{\lambda(1 + \rho)\mu M}}.$$

By independence and Lemma 2,

$$\mathbb{E}[e^{\lambda S}] = \prod_{i=1}^M \mathbb{E}[e^{\lambda Z_i}] \leq \prod_{i=1}^M \exp(\mu(e^\lambda - 1)) = \exp(\mu M(e^\lambda - 1)).$$

Thus

$$\mathbb{P}(S \geq (1 + \rho)\mu M) \leq \exp(\mu M(e^\lambda - 1) - \lambda(1 + \rho)\mu M).$$

Optimize over $\lambda > 0$; the minimizer is $\lambda^* = \ln(1 + \rho)$, yielding

$$\mathbb{P}(S \geq (1 + \rho)\mu M) \leq \exp\left(-\mu M((1 + \rho) \ln(1 + \rho) - \rho)\right).$$

To get the simpler $\rho^2/3$ bound for $\rho \in (0, 1)$, note that for $\rho \in [0, 1]$,

$$(1 + \rho) \ln(1 + \rho) - \rho \geq \frac{\rho^2}{3}.$$

A short proof: define $g(\rho) = (1 + \rho) \ln(1 + \rho) - \rho - \rho^2/3$. Then $g(0) = 0$ and

$$g'(\rho) = \ln(1 + \rho) - \frac{2\rho}{3}.$$

Since $\ln(1 + \rho)$ is concave, it lies above the chord from $(0, 0)$ to $(1, \ln 2)$: $\ln(1 + \rho) \geq \rho \ln 2$ for $\rho \in [0, 1]$. Because $\ln 2 > 2/3$, we get $g'(\rho) \geq \rho(\ln 2 - 2/3) \geq 0$, hence $g(\rho) \geq 0$.

Lower tail. Similarly, for $\lambda < 0$,

$$\mathbb{P}(S \leq (1 - \rho)\mu M) = \mathbb{P}(e^{\lambda S} \geq e^{\lambda(1-\rho)\mu M}) \leq \frac{\mathbb{E}[e^{\lambda S}]}{e^{\lambda(1-\rho)\mu M}} \leq \exp(\mu M(e^\lambda - 1) - \lambda(1 - \rho)\mu M).$$

Optimize over $\lambda < 0$; the minimizer is $\lambda^* = \ln(1 - \rho)$, giving

$$\mathbb{P}(S \leq (1 - \rho)\mu M) \leq \exp\left(-\mu M(\rho + (1 - \rho) \ln(1 - \rho))\right).$$

To get the simpler $\rho^2/2$ bound, define $h(\rho) = \rho + (1 - \rho) \ln(1 - \rho) - \rho^2/2$. Then $h(0) = 0$ and

$$h'(\rho) = -\ln(1 - \rho) - \rho \geq 0$$

because $-\ln(1 - \rho) \geq \rho$ for $\rho \in [0, 1)$ (equivalently $\ln(1 - \rho) \leq -\rho$). Hence $h(\rho) \geq 0$.

Finally, combine the two tails and use the weaker constant 3 to obtain the stated two-sided bound. \square

Corollary 2 (Relative-error sample size for bounded variables). *Under the conditions of Theorem 2, for $\rho \in (0, 1)$ and $\delta \in (0, 1)$,*

$$M \geq \frac{3}{\rho^2 \mu} \ln \frac{2}{\delta} \implies \mathbb{P}(|\bar{Z} - \mu| \leq \rho \mu) \geq 1 - \delta.$$

4 Sample complexity: rollouts vs posterior sampling

Throughout this section, fix $\rho \in (0, 1)$ and $\delta \in (0, 1)$.

4.1 Estimating q via rollouts

Theorem 3 (Rollout estimator sample complexity). *Let $I_r = \mathbf{1}\{A\}$ from an iid posterior-predictive rollout. Then for*

$$R \geq \frac{3}{\rho^2 q} \ln \frac{2}{\delta},$$

we have $\mathbb{P}(|\hat{q}_{\text{roll}} - q| \leq \rho q) \geq 1 - \delta$.

Proof. Apply Corollary 2 to $Z_r = I_r \in [0, 1]$ with mean $\mu = q$. \square

4.2 Estimating q via posterior sampling and analytic $f(\theta)$

Define

$$b := \sup_{\theta \in \Theta} f(\theta) \in (0, 1].$$

Since $f(\theta)$ is a probability, $b \leq 1$, but for many “thin” events $b \ll 1$.

Theorem 4 (Posterior-sampling estimator sample complexity). *Assume $0 \leq f(\theta) \leq b$ for all θ , and define $X_i = f(\theta_i)$ with $\theta_i \stackrel{iid}{\sim} \pi(\cdot | D)$. Then for*

$$M \geq \frac{3b}{\rho^2 q} \ln \frac{2}{\delta},$$

we have $\mathbb{P}(|\hat{q}_{\text{post}} - q| \leq \rho q) \geq 1 - \delta$.

Proof. Let $Z_i := X_i/b \in [0, 1]$. Then $\mathbb{E}[Z_i] = \mathbb{E}[X_i]/b = q/b$ and $\hat{q}_{\text{post}} = b \bar{Z}$. Moreover,

$$\frac{|\hat{q}_{\text{post}} - q|}{q} = \frac{|b\bar{Z} - b\mathbb{E}[Z]|}{b\mathbb{E}[Z]} = \frac{|\bar{Z} - \mathbb{E}[Z]|}{\mathbb{E}[Z]}.$$

Apply Corollary 2 with $\mu = q/b$. □

Corollary 3 (Improvement factor in sample count). *Comparing Theorems 3 and 4, the posterior method reduces the required number of Monte Carlo samples by a factor*

$$\frac{R}{M} \approx \frac{1}{b}.$$

Thus a large gap is possible exactly when $b = \sup_{\theta} f(\theta)$ is very small.

4.3 Tightness: a worst-case posterior can saturate the bound

Proposition 1 (Worst-case tightness via a two-point posterior). *Suppose there exist θ_0, θ_1 such that $f(\theta_0) = 0$ and $f(\theta_1) = b$. For any $q \in (0, b)$, define a posterior supported on $\{\theta_0, \theta_1\}$ by*

$$\pi(\theta = \theta_1 | D) = \frac{q}{b}, \quad \pi(\theta = \theta_0 | D) = 1 - \frac{q}{b}.$$

Then $Z = f(\theta)/b$ is exactly Bernoulli(q/b). In particular, estimating q by posterior sampling is information-theoretically as hard as estimating the mean of a Bernoulli(q/b), so $M = \Omega\left(\frac{b}{\rho^2 q} \log \frac{1}{\delta}\right)$ samples are necessary (up to constants).

Proof. Under this posterior, $f(\theta)$ equals b with probability q/b and 0 otherwise, hence $Z = f(\theta)/b$ is Bernoulli(q/b). Standard lower bounds for Bernoulli mean estimation (e.g. Le Cam’s method on μ vs. $(1 + \rho)\mu$) imply the stated necessity; the upper bound in Theorem 4 matches this scaling up to constants. □

5 “Seeing one hit” vs. “estimating q ”

If the goal is to *observe* at least one realization in A (rather than estimate q), then nothing beats the $1/q$ barrier.

Proposition 2 (Samples needed to see at least one event). *If $Y^{(1)}, \dots, Y^{(R)}$ are iid draws from the posterior predictive and $\mathbb{P}(A \mid D) = q$, then*

$$\mathbb{P}(\exists r : Y^{(r)} \in A) = 1 - (1 - q)^R.$$

Thus to see at least one hit with probability $\geq 1 - \delta$, it suffices (and is necessary up to constants for small q) that

$$R \gtrsim \frac{1}{q} \ln \frac{1}{\delta}.$$

Moreover, sampling $\theta \sim \pi(\cdot \mid D)$ and then sampling $Y \sim p_\theta(\cdot \mid D)$ produces exactly one posterior-predictive draw, so posterior sampling does not change this hit complexity unless one uses analytic $f(\theta)$ (i.e. one is no longer waiting for a literal hit).

6 Bernoulli strings: two instructive extremes

Let $Y_t \in \{0, 1\}$ and parameter $\theta = p \in [0, 1]$. Let $m = N - n$ be the horizon.

6.1 Event $A = \text{all ones (no sample-count gain)}$

If $A = \{Y_{n+1} = \dots = Y_{n+m} = 1\}$, then $f(p) = p^m$ and

$$b = \sup_{p \in [0, 1]} p^m = 1,$$

so Corollary 3 gives no sample-count improvement in the worst case: M and R both scale like $\Theta(\frac{1}{q})$ for fixed (ρ, δ) .

6.2 Event $A = \text{one fixed mixed string (exponential gain possible)}$

Let $A = \{Y_{n+1:n+m} = s\}$ for a prespecified string $s \in \{0, 1\}^m$ with k ones and $m - k$ zeros ($1 \leq k \leq m - 1$). Then

$$f(p) = p^k (1 - p)^{m-k}, \quad b = \max_{p \in [0, 1]} f(p) = \left(\frac{k}{m}\right)^k \left(\frac{m-k}{m}\right)^{m-k}.$$

Hence the improvement factor is

$$\frac{1}{b} = \left(\frac{m}{k}\right)^k \left(\frac{m}{m-k}\right)^{m-k} = \exp\left(m H(k/m)\right),$$

where $H(x) = -x \ln x - (1 - x) \ln(1 - x)$ is the binary entropy (natural logs). In particular, for $k = m/2$ one gets $b = 2^{-m}$ and the improvement factor 2^m .

7 Beyond Bernoulli: multinomials, HMMs, and state-space models

7.1 Multinomial (Dirichlet–Categorical) generalization

Let $Y_t \in \{1, \dots, K\}$, $\theta = \pi \in \Delta^{K-1}$, and $p_\pi(y_{1:N}) = \prod_{t=1}^N \pi_{y_t}$ (iid categorical). For a prespecified length- m string s with counts c_1, \dots, c_K (so $\sum_j c_j = m$),

$$f(\pi) = \mathbb{P}_\pi(Y_{n+1:n+m} = s \mid D) = \prod_{j=1}^K \pi_j^{c_j},$$

and

$$b = \sup_{\pi \in \Delta^{K-1}} \prod_{j=1}^K \pi_j^{c_j} = \prod_{j=1}^K \left(\frac{c_j}{m} \right)^{c_j}.$$

If $c_j = m/K$ (balanced), then $b = K^{-m}$ and the improvement factor is K^m .

7.2 Hidden Markov models: analytic $f(\theta)$ via forward DP and exponential b

Consider an HMM with finite latent states $X_t \in \{1, \dots, S\}$ and observations $Y_t \in \mathcal{Y}$. A parameter θ specifies an initial distribution $\pi_\theta(x_1)$, transition matrix $T_\theta(x' | x)$, and emission probabilities $E_\theta(y | x)$. Given θ ,

$$\mathbb{P}_\theta(y_{1:m}) = \sum_{x_{1:m}} \pi_\theta(x_1) \prod_{t=2}^m T_\theta(x_t | x_{t-1}) \prod_{t=1}^m E_\theta(y_t | x_t).$$

For a *fixed future string* $s \in \mathcal{Y}^m$ conditional on observed prefix $D = y_{1:n}$,

$$f(\theta) = \mathbb{P}_\theta(Y_{n+1:n+m} = s | D)$$

is computable by the standard forward algorithm: compute the filtered distribution over X_n from D , then propagate m steps multiplying by emissions for the fixed s . This costs $O(mS^2)$ time for discrete HMMs.

The key quantity for sample complexity is $b = \sup_\theta f(\theta)$. A simple sufficient condition for exponential smallness of b is a uniform emission peak bound.

Lemma 3 (Uniform per-step peak bound implies $b \leq \eta^m$). *Assume there exists $\eta \in (0, 1)$ such that for all θ , all states x , and all symbols $y \in \mathcal{Y}$,*

$$E_\theta(y | x) \leq \eta.$$

Then for any fixed length- m observation string $s \in \mathcal{Y}^m$ and any prefix D ,

$$\sup_\theta \mathbb{P}_\theta(Y_{n+1:n+m} = s | D) \leq \eta^m.$$

Hence $b \leq \eta^m$ and the rollout-vs-posterior improvement factor is at least η^{-m} .

Proof. Fix θ and condition on the prefix D . For each $t = 1, \dots, m$,

$$\mathbb{P}_\theta(Y_{n+t} = s_t | D, Y_{n+1:n+t-1} = s_{1:t-1}) = \sum_x \mathbb{P}_\theta(X_{n+t} = x | D, s_{1:t-1}) E_\theta(s_t | x) \leq \max_x E_\theta(s_t | x) \leq \eta.$$

Multiplying the conditional probabilities via the chain rule gives

$$\mathbb{P}_\theta(Y_{n+1:n+m} = s | D) \leq \eta^m.$$

Taking \sup_θ yields the claim. □

Remark 2 (State-space models). For continuous-observation state-space models, the event “ $Y_{n+1:n+m}$ equals an exact real-valued trajectory” typically has probability 0. A direct analogue is to take A to be a small neighborhood (e.g. an ε -tube) around a target trajectory, or to discretize/quantize observations. When the one-step observation likelihood/density is uniformly bounded above, an analogue of Lemma 3 typically yields $b \leq (\text{const} \cdot \varepsilon)^m$ and therefore an exponential separation in m .

8 Markov-chain experiment

We repeated the Bernoulli experiment for a two-state Markov chain with parameters $p = \mathbb{P}(x_n = 1 \mid x_{n-1} = 0)$ and $q = \mathbb{P}(x_n = 0 \mid x_{n-1} = 1)$. The prior is independent Beta(1, 1) on both p and q , and the initial state is sampled from a Bernoulli(1/2). The Bayes predictor updates the transition counts $(n_{01}, n_{00}, n_{10}, n_{11})$ and uses the posterior means $\mathbb{E}[p \mid D]$ and $\mathbb{E}[q \mid D]$ (conditioning on the last bit).

Transformer and training. We trained a small decoder-only transformer with $d_{\text{model}} = 32$, $n_{\text{layers}} = 2$, $n_{\text{heads}} = 4$, and $d_{\text{mlp}} = 64$, sequence length 256, batch size 64, and 12,000 total steps (initial 4,000 plus 8,000 continuation), with AdamW (lr = $3 \cdot 10^{-4}$, 400 warmup steps). Smaller configs (e.g. $d_{\text{model}} = 16$, $n_{\text{layers}} = 1$ or 2) produced noticeably worse Bayes alignment; the above was the smallest configuration with stable posterior-sampling diagnostics.

Model vs. Bayes prediction. On 100 random contexts of length 60, the model achieved $\text{MSE} = 1.53 \times 10^{-2}$ vs. Bayes $\text{MSE} = 7.54 \times 10^{-3}$ for next-bit prediction, a ratio of 2.03. Figure 1 summarizes the comparison.

Rare-event estimation. For an alternating target string of length 60 and a context of length 60, the analytical probability was $\mathbb{P}(\text{target} \mid D) = 4.91 \times 10^{-16}$. Posterior sampling with 80 samples and rollout length 400 estimated $1.99 \times 10^{-16} \pm 1.70 \times 10^{-16}$. By contrast, naive rollout would require on the order of $1/q \approx 2.0 \times 10^{15}$ rollouts to see a single hit in expectation.

Diagnostics. Across 15 random trials, the correlation between $\log_{10}(\text{true})$ and $\log_{10}(\text{estimate})$ was 0.957, and the posterior-mean correlations were 0.898 for p and 0.929 for q . Figures 2–3 show representative plots.

9 Executive summary of the main message

Let $q = \mathbb{P}(A \mid D) = \mathbb{E}_{\theta \sim \pi(\cdot \mid D)}[f(\theta)]$ with $f(\theta) = \mathbb{P}_{\theta}(A \mid D)$.

- **If you must literally wait for a hit of A in simulation:** you need $\Theta(\frac{1}{q} \log \frac{1}{\delta})$ posterior-predictive draws to see one hit with high probability.
- **If you only need to estimate q and can compute $f(\theta)$:** rollout MC needs $R = \Theta(\frac{1}{\rho^2 q} \log \frac{1}{\delta})$ samples, while posterior sampling needs

$$M = \Theta\left(\frac{b}{\rho^2 q} \log \frac{1}{\delta}\right), \quad b = \sup_{\theta} f(\theta).$$

Thus the **sample-count improvement factor** is $\Theta(1/b)$, which can be exponential in m for “thin” events where b decays exponentially (e.g. fixed strings in multinomials, noisy HMMs).

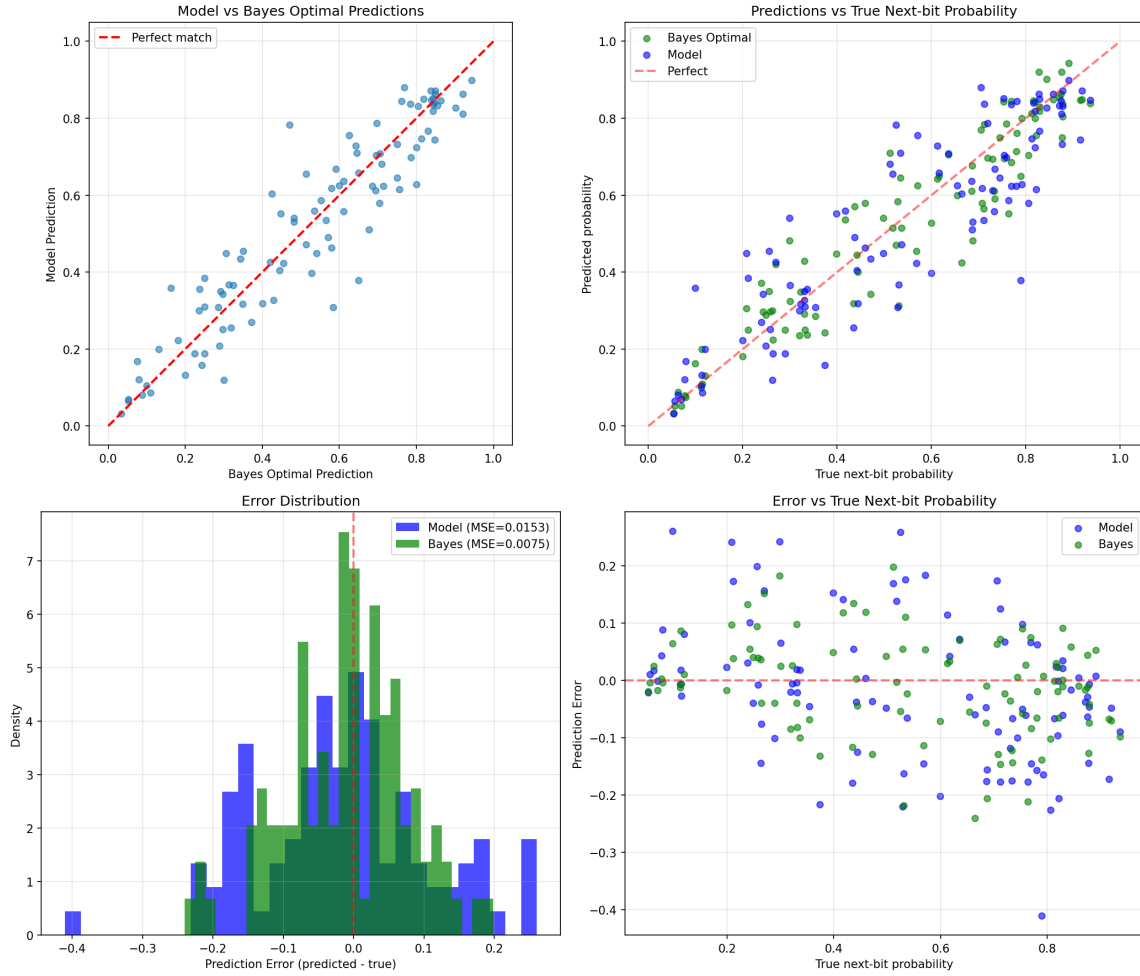


Figure 1: Markov transformer vs. Bayes predictor for next-bit probabilities (100 contexts, length 60).

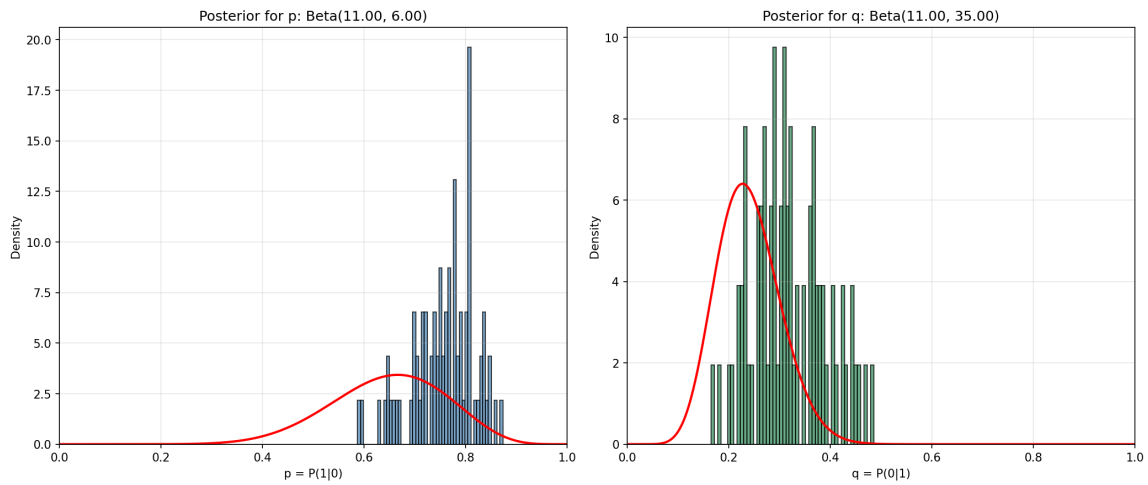


Figure 2: Posterior sampling check: histograms of inferred (p, q) vs. analytical Beta posteriors.

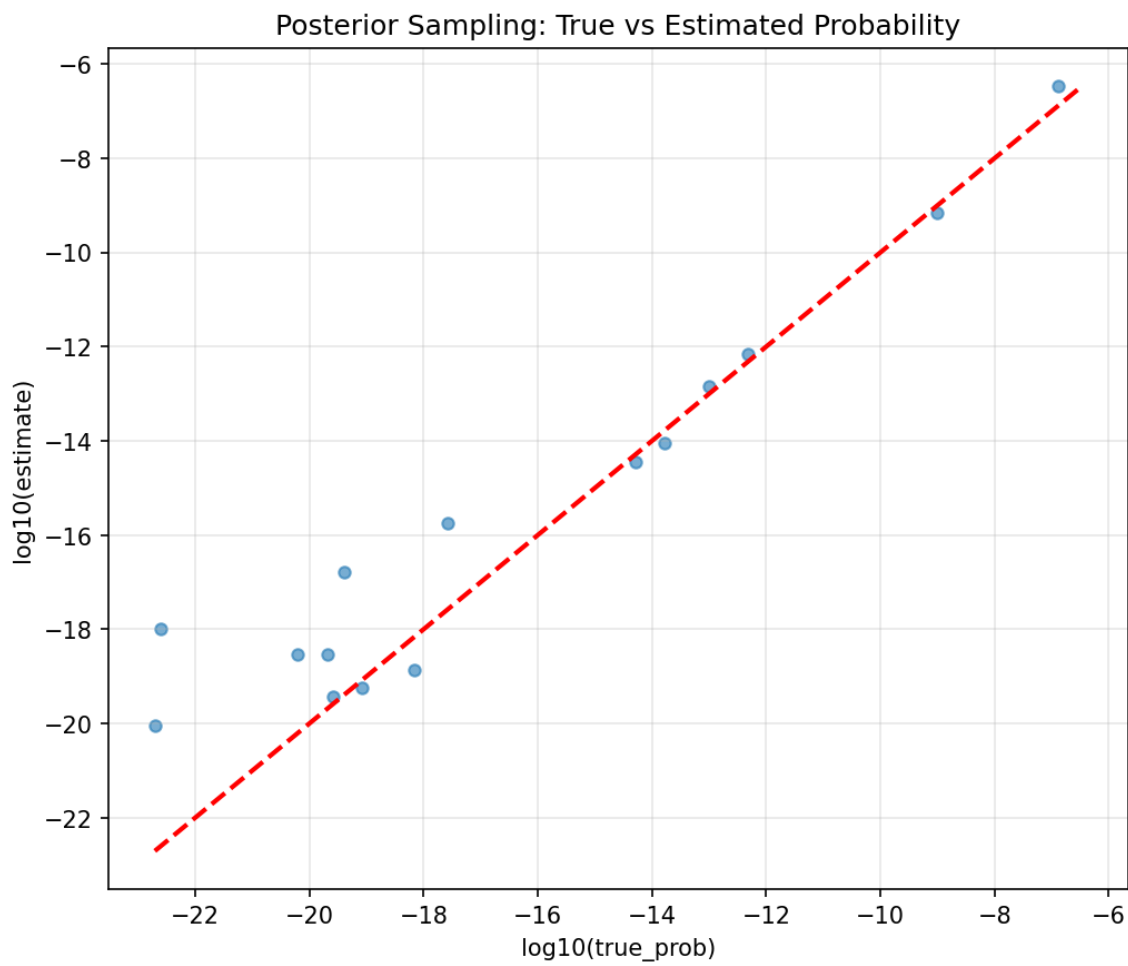


Figure 3: Diagnostics: $\log_{10}(\text{true probability})$ vs. $\log_{10}(\text{posterior-sampling estimate})$ over 15 trials.