

7 Proofs for Main Framework (Sec. 3)

Lemma 1. Let ω be a random variable with distribution $Q(\omega)$ and let $R(\omega)$ be a positive estimator such that $\mathbb{E}_{Q(\omega)} R(\omega) = p(x)$. Then

$$P^{MC}(\omega, x) = Q(\omega)R(\omega)$$

is an unnormalized distribution over ω with normalization constant $p(x)$ and $R(\omega) = P^{MC}(\omega, x)/Q(\omega)$ for $Q(\omega) > 0$. Furthermore as defined above,

$$\log p(x) = \mathbb{E}_{Q(\omega)} \log R(\omega) + \text{KL} [Q(\omega) \| P^{MC}(\omega|x)]. \quad (3)$$

Proof. Since $P^{MC}(\omega, x) \geq 0$ and $P^{MC}(x) = \int P^{MC}(\omega, x) d\omega = \mathbb{E}_{Q(\omega)} R(\omega) = p(x)$, it is a valid distribution. Thus, one can apply the standard ELBO decomposition to $Q(\omega)$ and $P^{MC}(\omega, x)$. But since $R = P^{MC}/Q$, it follows that $\mathbb{E}_{Q(\omega)} \log (P^{MC}(\omega, x)/Q(\omega)) = \mathbb{E}_{Q(\omega)} \log R(\omega)$. \square

Theorem 2. Suppose that $R(\omega)$ and $a(z|\omega)$ are a valid estimator-coupling pair under $Q(\omega)$. Then,

$$Q(z, \omega) = Q(\omega)a(z|\omega), \quad (5)$$

$$P^{MC}(z, \omega, x) = Q(\omega)R(\omega)a(z|\omega), \quad (6)$$

are valid distributions, $P^{MC}(z, x) = p(z, x)$, and

$$\log p(x) = \mathbb{E}_{Q(\omega)} \log R(\omega) + \text{KL} [Q(z) \| p(z|x)] + \text{KL} [Q(\omega|z) \| P^{MC}(\omega|z, x)]. \quad (7)$$

Proof. First, note that

$$\begin{aligned} P^{MC}(z, x) &= \int P^{MC}(z, \omega, x) d\omega \\ &= \int Q(\omega)R(\omega)a(z|\omega) d\omega \\ &= \mathbb{E}_{Q(\omega)} R(\omega)a(z|\omega) \\ &= p(z, x), \end{aligned}$$

so $P^{MC}(z, \omega, x)$ is a valid augmentation of $p(z, x)$.

Next, observe for P^{MC} and Q as defined,

$$\frac{P^{MC}(z, \omega, x)}{Q(z, \omega)} = R(\omega).$$

Applying the ELBO decomposition from Eq. (1) to $Q(z, \omega)$ and $P^{MC}(z, \omega, x)$ we get that

$$\log P^{MC}(x) = \mathbb{E}_{Q(z, \omega)} \left[\log \frac{P^{MC}(z, \omega, x)}{Q(z, \omega)} \right] + \text{KL} [Q(z, \omega) \| P^{MC}(z, \omega|x)].$$

Using the observations above and the chain rule of KL-divergence means that

$$\begin{aligned} \log p(x) &= \mathbb{E}_{Q(\omega)} \log R(\omega) + \text{KL} [Q(z, \omega) \| P^{MC}(z, \omega|x)] \\ &= \mathbb{E}_{Q(\omega)} \log R(\omega) + \text{KL} [Q(z) \| P^{MC}(z|x)] + \text{KL} [Q(\omega|z) \| P^{MC}(\omega|z, x)] \\ &= \mathbb{E}_{Q(\omega)} \log R(\omega) + \text{KL} [Q(z) \| p(z|x)] + \text{KL} [Q(\omega|z) \| P^{MC}(\omega|z, x)]. \end{aligned}$$

\square

Claim 5. Suppose that $Q(T(\omega)) = Q(\omega)$. Then, the antithetic estimator

$$R(\omega) = \frac{p(\omega, x) + p(T(\omega), x)}{2Q(\omega)}$$

and the coupling distribution

$$\begin{aligned} a(z|\omega) &= \pi(\omega) \delta(z - \omega) + (1 - \pi(\omega)) \delta(z - T(\omega)), \\ \pi(\omega) &= \frac{p(\omega, x)}{p(\omega, x) + p(T(\omega), x)}. \end{aligned}$$

form a valid estimator / coupling pair under $Q(\omega)$.

Proof.

$$\begin{aligned} & \mathbb{E}_{Q(\omega)} R(\omega) a(z|\omega) \\ &= \mathbb{E}_{Q(\omega)} \frac{p(\omega, x) + p(T(\omega), x)}{2Q(\omega)} (\pi(\omega) \delta(z - \omega) + (1 - \pi(\omega)) \delta(z - T(\omega))) \\ &= \mathbb{E}_{Q(\omega)} \frac{p(\omega, x) + p(T(\omega), x)}{2Q(\omega)} \left(\frac{p(\omega, x)}{p(\omega, x) + p(T(\omega), x)} \delta(z - \omega) \right. \\ & \quad \left. + \frac{p(T(\omega), x)}{p(\omega, x) + p(T(\omega), x)} \delta(z - T(\omega)) \right) \\ &= \mathbb{E}_{Q(\omega)} \frac{1}{2Q(\omega)} (p(\omega, x) \delta(z - \omega) + p(T(\omega), x) \delta(z - T(\omega))) \\ &= \mathbb{E}_{Q(\omega)} \frac{1}{2} \left(\frac{1}{Q(\omega)} p(\omega, x) \delta(z - \omega) + \frac{1}{Q(\omega)} p(T(\omega), x) \delta(z - T(\omega)) \right) \\ &= \mathbb{E}_{Q(\omega)} \frac{1}{2} \left(\frac{1}{Q(\omega)} p(\omega, x) \delta(z - \omega) + \frac{1}{Q(T(\omega))} p(T(\omega), x) \delta(z - T(\omega)) \right) \quad (10) \\ &= \mathbb{E}_{Q(\omega)} \frac{1}{2} \left(\frac{1}{Q(\omega)} p(\omega, x) \delta(z - \omega) + \frac{1}{Q(\omega)} p(\omega, x) \delta(z - \omega) \right) \quad (11) \\ &= \mathbb{E}_{Q(\omega)} \left(\frac{1}{Q(\omega)} p(\omega, x) \delta(z - \omega) \right) \\ &= \int (p(\omega, x) \delta(z - \omega)) d\omega \\ &= p(z, x) \end{aligned}$$

Here, [Eq. \(10\)](#) follows from the fact that $Q(T(\omega)) = Q(\omega)$ while [Eq. \(11\)](#) follows from the fact that $T(\omega)$ is equal in distribution to ω when $\omega \sim Q$. \square

8 Measure-Theoretic Details

The content of this section draws from [14, 15]. We do not use sans-serif font in this section.

8.1 Measures, KL, ELBO

Let (Ω, \mathcal{A}) be a measurable space and Q and P be two measures over it. Write $Q \ll P$ when Q is absolutely continuous with respect to P , i.e. when $P(A) = 0 \Rightarrow Q(A) = 0$. Whenever $Q \ll P$, there exists measurable $f : \Omega \rightarrow \mathbb{R}$ such that

$$Q(A) = \int_A f dP.$$

The function f is the *Radon-Nikodym derivative*, denoted as $f = \frac{dQ}{dP}$. Write $Q \sim P$ when $Q \ll P$ and $P \ll Q$; in this case $\frac{dQ}{dP} = \left(\frac{dP}{dQ}\right)^{-1}$ Q -a.e.

For two probability measures $Q \ll P$, the KL-divergence is

$$\text{KL}[Q\|P] = \int \log\left(\frac{dQ}{dP}\right) Q(d\omega) = \mathbb{E}_{Q(\omega)} \log \frac{dQ}{dP}.$$

For a probability measure Q and measure \hat{P} (not necessarily a probability measure) with $Q \ll \hat{P}$, the evidence lower bound or “ELBO” is

$$\text{ELBO}[Q\|\hat{P}] = -\mathbb{E}_Q \log \frac{dQ}{d\hat{P}}.$$

When $Q \sim \hat{P}$, we can equivalently write $\text{ELBO}[Q\|\hat{P}] = \mathbb{E}_Q \log \frac{d\hat{P}}{dQ}$.

Let (Z, \mathcal{B}) be a measurable space. Let $P_{z,x}$ be an unnormalized distribution over z representing the joint distribution over (z, x) for a fixed x . Write either $P_{z,x}(B)$ or $P_{z,x}(z \in B)$ for the measure of $B \in \mathcal{B}$. Define

$$p(x) = P_{z,x}(Z)$$

to be the total measure or the normalization constant of $P_{z,x}$, and write $P_{z|x}(z \in B) := P_{z,x}(z \in B)/p(x)$ for the corresponding normalized measure, which represents the conditional distribution of z given x . Henceforth, x will *always* denote a fixed constant, and, for any u , the measure $P_{u,x}$ is unnormalized with total measure $p(x)$.

The following gives a measure-theoretic version of the “ELBO decomposition” from Eq. (1).

Lemma 6. *Given a probability measure Q and a measure $P_{z,x}$ on (Z, \mathcal{B}) , whenever $Q \ll P_{z,x}$ we have the following “ELBO decomposition”:*

$$\log p(x) = \text{ELBO}[Q\|P_{z,x}] + \text{KL}[Q\|P_{z|x}].$$

Proof. It is easy to check that $\frac{dQ}{dP_{z|x}} = p(x) \frac{dQ}{dP_{z,x}}$.³ Then

$$\begin{aligned} \text{KL}[Q\|P_{z|x}] &= \mathbb{E}_Q \log \frac{dQ}{dP_{z|x}} = \mathbb{E}_Q \log \left(p(x) \frac{dQ}{dP_{z,x}} \right) \\ &= \log p(x) + \mathbb{E}_Q \log \frac{dQ}{dP_{z,x}} = \log p(x) - \text{ELBO}[Q\|P_{z,x}]. \end{aligned}$$

Rearranging, we see the ELBO decomposition. □

³ $\int_A p(x) \frac{dQ}{dP_{z,x}} dP_{z|x} = \int_A \frac{dQ}{dP_{z,x}} dP_{z,x} = Q(z \in A)$.

8.2 Conditional, Marginal, and Joint Distributions

Standard Borel and product spaces We will assume that each relevant measure space is a *standard Borel space*, that is, isomorphic to a Polish space (a separable complete metric space) with the Borel σ -algebra. Standard Borel spaces capture essentially all spaces that arise in practice in probability theory [15]. Let (Ω, \mathcal{A}) and (Z, \mathcal{B}) be standard Borel spaces. The *product space* $(\Omega \times Z, \mathcal{A} \otimes \mathcal{B})$ is the measurable space on $\Omega \times Z$ with $\mathcal{A} \otimes \mathcal{B} = \{A \times B : A \in \mathcal{A}, B \in \mathcal{B}\}$, and is also a standard Borel space.

Conditional distributions We require tools to augment a distribution with a new random variable and define the conditional distribution of one random variable with respect to another. We begin with a Markov kernel, which we will use to augment a distribution P_ω with a new random variable to obtain a joint distribution $P_{\omega,z}$.

Formally, a *Markov kernel* [15, Def. 8.24] from (Ω, \mathcal{A}) to (Z, \mathcal{B}) is a mapping $a(B|\omega)$ that satisfies:

1. For fixed ω , $a(B|\omega)$ is a probability measure on (Z, \mathcal{B}) .
2. For fixed B , $a(B|\omega)$ is an \mathcal{A} -measurable function of ω .

Let P_ω be a measure on (Ω, \mathcal{A}) and $a(B|\omega)$ a Markov kernel from (Ω, \mathcal{A}) to (Z, \mathcal{B}) . These define a unique measure $P_{\omega,z}$ over the product space defined as

$$P_{\omega,z}(\omega \in A, z \in B) = \int_A a(z \in B|\omega) P_\omega(d\omega),$$

such that if P_ω is a probability measure, then $P_{\omega,z}$ is also a probability measure [15, Cor. 14.23].

Alternately, we may have a joint distribution $P_{\omega,z}$ (a measure on the product space $(\Omega \times Z, \mathcal{A} \otimes \mathcal{B})$) and want to define the marginals and conditionals. The *marginal distribution* P_z is the measure on (Z, \mathcal{B}) with $P_z(z \in B) = P_{\omega,z}(\omega \in \Omega, z \in B)$, and the marginal P_ω is defined analogously. Since the product space is standard Borel [15, Thm 14.8], there exists a *regular conditional distribution* $P_{\omega|z}(\omega \in A|z)$ [15, Def. 8.27, Thm. 8.36], which is a Markov kernel (as above) and satisfies the following for all $A \in \mathcal{A}, B \in \mathcal{B}$:

$$P_{\omega,z}(\omega \in A, z \in B) = \int_B P_{\omega|z}(\omega \in A|z) P_z(dz).$$

The regular conditional distribution is unique up to null sets of P_z .

The conditional distribution $P_{z|\omega}$ is defined analogously.

8.3 KL Chain Rule

Let $P_{\omega,z}$ and $Q_{\omega,z}$ be two probability measures on the standard Borel product space $(\Omega \times Z, \mathcal{A} \otimes \mathcal{B})$ with $Q_{\omega,z} \ll P_{\omega,z}$. The *conditional KL-divergence* $\text{KL}[Q_{\omega|z}||P_{\omega|z}]$ is defined⁴ as [14, Ch. 5.3]

$$\text{KL}[Q_{\omega|z}||P_{\omega|z}] = \mathbb{E}_{Q_{\omega,z}} \left(\frac{dQ_{\omega|z}}{dP_{\omega|z}} \right),$$

where $\frac{dQ_{\omega|z}}{dP_{\omega|z}}(\omega|z) = \left(\frac{dQ_{\omega,z}}{dP_{\omega,z}}(\omega, z) \right) \left(\frac{dP_z}{dP_z}(z) \right)^{-1}$ when $\frac{dP_z}{dP_z}(z) > 0$ and 1 otherwise. When all densities exist, $\frac{dQ_{\omega|z}}{dP_{\omega|z}}(\omega|z) = \frac{q(\omega|z)}{p(\omega|z)}$. Under the same conditions as above, we have the *chain rule for KL-divergence* [14, Lem. 5.3.1]

$$\text{KL}[Q_{\omega,z}||P_{\omega,z}] = \text{KL}[Q_\omega||P_\omega] + \text{KL}[Q_{\omega|z}||P_{\omega|z}] = \text{KL}[Q_z||P_z] + \text{KL}[Q_{z|\omega}||P_{z|\omega}].$$

⁴While this (standard) notation for the divergence refers to “ $Q_{\omega|z}$ ” it is a function of the joint $Q_{\omega,z}$ and similarly for $P_{\omega,z}$.

8.4 Our Results

Now consider a strictly positive estimator $R(\omega)$ over probability space $(\Omega, \mathcal{A}, Q_\omega)$ such that $\mathbb{E}_{Q_\omega} R = \int R dQ_\omega = p(x)$. We wish to define $P_{\omega,x}^{\text{MC}}$ so that $\frac{dP_{\omega,x}^{\text{MC}}}{dQ_\omega} = R$, to justify interpreting $\mathbb{E}_{Q_\omega} \log R$ as an ELBO. This is true when $R = \frac{dP_{\omega,x}^{\text{MC}}}{dQ_\omega}$ is the Radon-Nikodym derivative, i.e., a change of measure from Q_ω to $P_{\omega,x}^{\text{MC}}$, and is strictly positive. This leads to the definition

$$P_{\omega,x}^{\text{MC}}(\omega \in A) = \int_A R dQ_\omega.$$

Lemma 7. *Let $R(\omega)$ be an almost-everywhere positive random variable on $(\Omega, \mathcal{A}, Q_\omega)$ with $\mathbb{E}_{Q_\omega} R = p(x)$ and define $P_{\omega,x}^{\text{MC}}(\omega \in A) = \int_A R dQ_\omega$. The ELBO decomposition applied to Q_ω and $P_{\omega,x}^{\text{MC}}$ gives:*

$$\log p(x) = \mathbb{E}_{Q_\omega} \log R + \text{KL} [Q_\omega \| P_{\omega|x}^{\text{MC}}].$$

Proof. By construction, $R = \frac{dP_{\omega,x}^{\text{MC}}}{dQ_\omega}$ and $P_{\omega,x}^{\text{MC}} \sim Q_\omega$, since R is positive Q -a.e. Therefore $\mathbb{E}_{Q_\omega} \log R = \mathbb{E}_{Q_\omega} \log \frac{dP_{\omega,x}^{\text{MC}}}{dQ_\omega} = \text{ELBO} [Q \| P_{\omega,x}^{\text{MC}}]$, where the final equality uses the definition of the ELBO for the case when $P_{\omega,x}^{\text{MC}} \sim Q_\omega$. Now apply [Lem. 6](#) and the fact that $\text{ELBO} [Q \| P_{\omega,x}^{\text{MC}}] = \mathbb{E}_{Q_\omega} \log R$. \square

[Lem. 7](#) provides distributions Q_ω and $P_{\omega,x}^{\text{MC}}$ so that $\mathbb{E}_{Q_\omega} \log R = \text{ELBO} [Q_\omega \| P_{\omega,x}^{\text{MC}}]$, which justifies maximizing the likelihood bound $\mathbb{E}_{Q_\omega} \log R$ as minimizing the KL-divergence from Q_ω to the “target” $P_{\omega|x}^{\text{MC}}$. However, neither distribution contains the random variable z from the original target distribution $P_{z|x}$, so the significance of [Lem. 7](#) on its own is unclear. We now describe a way to couple $P_{\omega,x}^{\text{MC}}$ to the original target distribution using a Markov kernel $a(z \in B|\omega)$.

Definition 8. A valid *estimator-coupling pair* with respect to target distribution $P_{z,x}$ is an estimator $R(\omega)$ on probability space $(\Omega, \mathcal{A}, Q_\omega)$ and Markov kernel $a(z \in B|\omega)$ from (Ω, \mathcal{A}) to (Z, \mathcal{B}) such that:

$$\mathbb{E}_{Q_\omega} R(\omega) a(z \in B|\omega) = P_{z,x}(z \in B).$$

Lemma 9. *Assume $R(\omega)$ and $a(z \in B|\omega)$ are a valid estimator-coupling pair with respect to target $P_{z,x}$, and define*

$$P_{\omega,z,x}^{\text{MC}}(\omega \in A, z \in B) = \int_A a(z \in B|\omega) R(\omega) Q_\omega(d\omega).$$

Then $P_{\omega,z,x}^{\text{MC}}$ admits $P_{z,x}$ as a marginal, i.e., $P_{z,x}^{\text{MC}}(z \in B) = P_{z,x}(z \in B)$.

Proof. We have

$$\begin{aligned} P_{z,x}^{\text{MC}}(z \in B) &= P_{\omega,z,x}^{\text{MC}}(\omega \in \Omega, z \in B) \\ &= \int_\Omega a(z \in B|\omega) R(\omega) Q_\omega(d\omega). \\ &= \mathbb{E}_{Q_\omega} R(\omega) a(z \in B|\omega) \\ &= P_{z,x}(z \in B). \end{aligned}$$

The second line uses the definition of $P_{\omega,z,x}^{\text{MC}}$. The last line uses the definition of a valid estimator-coupling pair. \square

Theorem 10. *Let $P_{z,x}$ be an unnormalized distribution with normalization constant $p(x)$. Assume $R(\omega)$ and $a(z \in B|\omega)$ are a valid estimator-coupling pair with respect to $P_{z,x}$. Define $P_{\omega,z,x}^{\text{MC}}$ as in [Lem. 9](#) and define $Q_{\omega,z}(\omega \in A, z \in B) = \int_A a(z \in B|\omega) Q_\omega(d\omega)$. Then*

$$\log p(x) = \mathbb{E}_{Q_\omega} \log R + \text{KL} [Q_z \| P_{z|x}] + \text{KL} [Q_{\omega|z} \| P_{\omega|z,x}^{\text{MC}}].$$

Proof. From [Lem. 7](#), we have

$$\log p(x) = \mathbb{E}_{Q_\omega} \log R + \text{KL} [Q_\omega \| P_{\omega|x}^{\text{MC}}].$$

We will show by two applications of the KL chain rule that the second term can be expanded as

$$\text{KL} [Q_\omega \| P_{\omega|x}^{\text{MC}}] = \text{KL} [Q_z \| P_{z|x}] + \text{KL} [Q_{\omega|z} \| P_{\omega|z,x}^{\text{MC}}], \quad (12)$$

which will complete the proof.

We first apply the KL chain rule as follows:

$$\text{KL} [Q_{\omega,z} \| P_{\omega,z|x}^{\text{MC}}] = \text{KL} [Q_\omega \| P_{\omega|x}^{\text{MC}}] + \underbrace{\text{KL} [Q_{z|\omega} \| P_{z|\omega,x}^{\text{MC}}]}_{=0}. \quad (13)$$

We now argue that the second term is zero, as indicated in the equation. Note from above that $\frac{dP_{\omega,x}^{\text{MC}}}{dQ_\omega} = R$. It is also true that $\frac{dP_{\omega,z,x}^{\text{MC}}}{dQ_{\omega,z}} = R$. To see this, observe that

$$\begin{aligned} \int_{A \times B} R(\omega) Q_{\omega,z}(d\omega, dz) &= \int_A \left(\int_B R(\omega) a(z \in dz|\omega) \right) Q_\omega(d\omega) \\ &= \int_A R(\omega) \left(\int_B a(z \in dz|\omega) \right) Q_\omega(d\omega) \\ &= \int_A R(\omega) a(z \in B|\omega) Q_\omega(d\omega) \\ &= P_{\omega,z,x}^{\text{MC}}(\omega \in A, z \in B). \end{aligned}$$

The first equality above uses a version of Fubini's theorem for Markov kernels [[?](#), Thm. 14.29]. Because $P_{\omega,x}^{\text{MC}} \sim Q_\omega$ it also follows that $\frac{dQ_\omega}{dP_{\omega,x}^{\text{MC}}} = \frac{dQ_{\omega,z}}{dP_{\omega,z,x}^{\text{MC}}} = \frac{1}{R}$. Since the normalized distributions $P_{\omega|x}^{\text{MC}}$ and $P_{\omega,z|x}^{\text{MC}}$ differ from the unnormalized counterparts by the constant factor $p(x)$, it is straightforward to see that $\frac{dQ_\omega}{dP_{\omega|x}^{\text{MC}}} = \frac{dQ_{\omega,z}}{dP_{\omega,z|x}^{\text{MC}}} = \frac{p(x)}{R}$.⁵ This implies that $\frac{dP_{\omega,z,x}^{\text{MC}}}{dQ_{z|\omega}} = 1$ a.e., which in turn implies that the conditional divergence $\text{KL} [Q_{z|\omega} \| P_{z|\omega,x}^{\text{MC}}]$ is equal to zero.

We next apply the chain rule the other way and use the fact that $P_{z|x}^{\text{MC}} = P_{z|x}$ ([Lem. 9](#)) to see that:

$$\text{KL} [Q_{\omega,z} \| P_{\omega,z|x}^{\text{MC}}] = \text{KL} [Q_z \| P_{z|x}^{\text{MC}}] + \text{KL} [Q_{\omega|z} \| P_{\omega|z,x}^{\text{MC}}] = \text{KL} [Q_z \| P_{z|x}] + \text{KL} [Q_{\omega|z} \| P_{\omega|z,x}^{\text{MC}}]. \quad (14)$$

Combining [Eq. \(13\)](#) and [Eq. \(14\)](#) we get [Eq. \(12\)](#), as desired. \square

⁵ $\int_A \frac{p(x)}{R} dP_{\omega|x}^{\text{MC}} = \int_A \frac{1}{R} dP_{\omega,x}^{\text{MC}} = Q_\omega(\omega \in A)$, and similarly for $Q_{\omega,z}$ and $P_{\omega,z,x}^{\text{MC}}$.

9 Specific Variance Reduction Techniques

9.1 IID Mean

As a simple example, consider the IID mean. Suppose $R_0(\omega)$ and $a_0(z|\omega)$ are valid under Q_0 . If we define

$$Q(\omega_1, \dots, \omega_M, m) = \frac{1}{M} \prod_{m=1}^M Q_0(\omega_m)$$

(with $\omega_1, \dots, \omega_M \sim Q_0$ i.i.d. and m uniform on $\{1, \dots, M\}$) then this satisfies the condition of [Thm. 3](#) that $\omega_m \sim Q_0$. Thus we can define R and a as in [Eq. \(8\)](#) and [Eq. \(9\)](#), to get that

$$\begin{aligned} R(\omega_1, \dots, \omega_M, m) &= R_0(\omega_m) \\ a(z|\omega_1, \dots, \omega_M, m) &= a_0(z|\omega_m) \end{aligned}$$

are a valid estimator-coupling pair under Q . Note that $Q(m|\omega_1, \dots, \omega_M) = \frac{1}{M}$, so if we apply [Thm. 4](#) to marginalize out m , we get that

$$\begin{aligned} R(\omega_1, \dots, \omega_M) &= \mathbb{E}_{Q(m|\omega_1, \dots, \omega_M)} R(\omega_1, \dots, \omega_M, m) \\ &= \frac{1}{M} \sum_{m=1}^M R(\omega_1, \dots, \omega_M, m) \\ &= \frac{1}{M} \sum_{m=1}^M R_0(\omega_m) \\ a(z|\omega_1, \dots, \omega_M) &= \frac{1}{R(\omega_1, \dots, \omega_M)} \mathbb{E}_{Q(m|\omega_1, \dots, \omega_M)} [R(\omega_1, \dots, \omega_M, m) a(z|\omega_1, \dots, \omega_M, m)] \\ &= \frac{1}{R(\omega_1, \dots, \omega_M)} \frac{1}{M} \sum_{m=1}^M [R(\omega_1, \dots, \omega_M, m) a(z|\omega_1, \dots, \omega_M, m)] \\ &= \frac{1}{\frac{1}{M} \sum_{m=1}^M R_0(\omega_m)} \frac{1}{M} \sum_{m=1}^M [R_0(\omega_m) a_0(z|\omega_m)] \\ &= \frac{\sum_{m=1}^M [R_0(\omega_m) a_0(z|\omega_m)]}{\sum_{m=1}^M R_0(\omega_m)}. \end{aligned}$$

These are exactly the forms for $R(\cdot)$ and $a(z|\cdot)$ shown in the table.

9.2 Stratified Sampling

As another example, take stratified sampling. The estimator-coupling pair can be derived similarly to with the i.i.d. mean. For simplicity, we assume here one sample in each strata ($N_m = 1$). Suppose $\Omega_1 \dots \Omega_M$ partition the state-space and define

$$Q(\omega_1, \dots, \omega_M, m) = \frac{1}{M} \prod_{k=1}^M \frac{Q_0(\omega_k) I(\omega_k \in \Omega_m)}{\mu(k)} \times \mu(m), \quad \mu(m) = \mathbb{E}_{Q_0(\omega)} I(\omega \in \Omega_m).$$

This again satisfies the condition of [Thm. 3](#), so [Eq. \(8\)](#) and [Eq. \(9\)](#) give that

$$\begin{aligned} R(\omega_1, \dots, \omega_M, m) &= R_0(\omega_m) \\ a(z|\omega_1, \dots, \omega_M, m) &= a_0(z|\omega_m) \end{aligned}$$

is a valid estimator-coupling pair with respect to Q . Note that $Q(m|\omega_1, \dots, \omega_M) = \mu(m)$, so if we apply [Thm. 4](#) to marginalize out m , we get that

$$\begin{aligned}
R(\omega_1, \dots, \omega_M) &= \mathbb{E}_{Q(m|\omega_1, \dots, \omega_M)} R(\omega_1, \dots, \omega_M, m) \\
&= \sum_{m=1}^M \mu(m) R(\omega_1, \dots, \omega_M, m) \\
&= \sum_{m=1}^M \mu(m) R_0(\omega_m) \\
a(z|\omega_1, \dots, \omega_M) &= \frac{1}{R(\omega_1, \dots, \omega_M)} \mathbb{E}_{Q(m|\omega_1, \dots, \omega_M)} [R(\omega_1, \dots, \omega_M, m) a(z|\omega_1, \dots, \omega_M, m)] \\
&= \frac{1}{R(\omega_1, \dots, \omega_M)} \sum_{m=1}^M \mu(m) R(\omega_1, \dots, \omega_M, m) a(z|\omega_1, \dots, \omega_M, m) \\
&= \frac{\sum_{m=1}^M \mu(m) R_0(\omega_m) a_0(z|\omega_m)}{\sum_{m=1}^M \mu(m) R_0(\omega_m)}.
\end{aligned}$$

Again, this is the form shown in the table.

10 Proofs for Deriving Couplings (Sec. 4)

Theorem 3. Suppose that $R_0(\omega)$ and $a_0(z|\omega)$ are a valid estimator-coupling pair under $Q_0(\omega)$. Let $Q(\omega_1, \dots, \omega_M, m)$ be any distribution such that if $(\omega_1, \dots, \omega_M, m) \sim Q$, then $\omega_m \sim Q_0$. Then,

$$R(\omega_1, \dots, \omega_M, m) = R_0(\omega_m) \quad (8)$$

$$a(z|\omega_1, \dots, \omega_M, m) = a_0(z|\omega_m) \quad (9)$$

are a valid estimator-coupling pair under $Q(\omega_1, \dots, \omega_M, m)$.

Proof. Substitute the definitions of R and a to get that

$$\begin{aligned} \mathbb{E}_{Q(\omega_1, \dots, \omega_M, m)} R(\omega_1, \dots, \omega_M, m) a(z|\omega_1, \dots, \omega_M, m) &= \mathbb{E}_{Q(\omega_1, \dots, \omega_M, m)} R_0(\omega_m) a_0(z|\omega_m) \\ &= \mathbb{E}_{Q_0(\omega)} R_0(\omega) a_0(z|\omega) \\ &= p(z, x), \end{aligned}$$

which is equivalent to the definition of R and a being a valid estimator-coupling pair. The second line follows from the assumption on $Q(\omega_1, \dots, \omega_M, m)$. \square

Theorem 4. Suppose that $R_0(\omega, \nu)$ and $a_0(z|\omega, \nu)$ are a valid estimator-coupling pair under $Q_0(\omega, \nu)$. Then

$$\begin{aligned} R(\omega) &= \mathbb{E}_{Q_0(\nu|\omega)} R_0(\omega, \nu), \\ a(z|\omega) &= \frac{1}{R(\omega)} \mathbb{E}_{Q_0(\nu|\omega)} [R_0(\omega, \nu) a_0(z|\omega, \nu)], \end{aligned}$$

are a valid estimator-coupling pair under $Q(\omega) = \int Q_0(\omega, \nu) d\nu$.

Proof. Substitute the definition of a to get that

$$\begin{aligned} \mathbb{E}_{Q(\omega)} R(\omega) a(z|\omega) &= \mathbb{E}_{Q(\omega)} R(\omega) \frac{1}{R(\omega)} \mathbb{E}_{Q_0(\nu|\omega)} [R_0(\omega, \nu) a_0(z|\omega, \nu)] \\ &= \mathbb{E}_{Q_0(\omega, \nu)} [R_0(\omega, \nu) a_0(z|\omega, \nu)] \\ &= p(z, x), \end{aligned}$$

which is equivalent to R and a being a valid estimator-coupling pair under $R(\omega)$. The last line follows from the fact that R_0 and a_0 are a valid estimator-coupling pair under Q_0 . \square

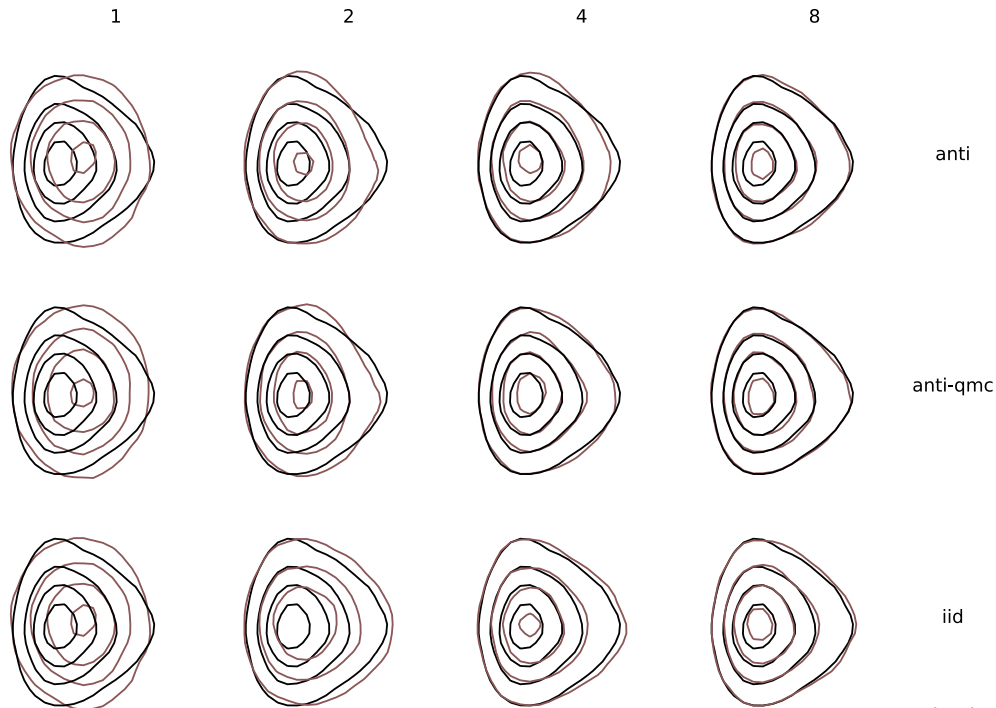


Figure 8: The target density $p(z|x)$ and the approximation $Q(z|x)$ produced by various sampling methods (row) with various M (columns). The dark curves show isocontours of kernel density estimate for samples generated using Stan and projected to the first two principal components. The darker curves show isocontours for the process applied to samples from $Q(z|x)$. Antithetic sampling is visibly (but subtly) better than iid for $M = 2$ while the combination of quasi-Monte Carlo and antithetic sampling is (still more subtly) best for $M = 8$.

11 Additional Experimental Results

[Fig. 9](#) shows additional aggregate statistics of ELBO and posterior variance error for different methods across model from the Stan library.

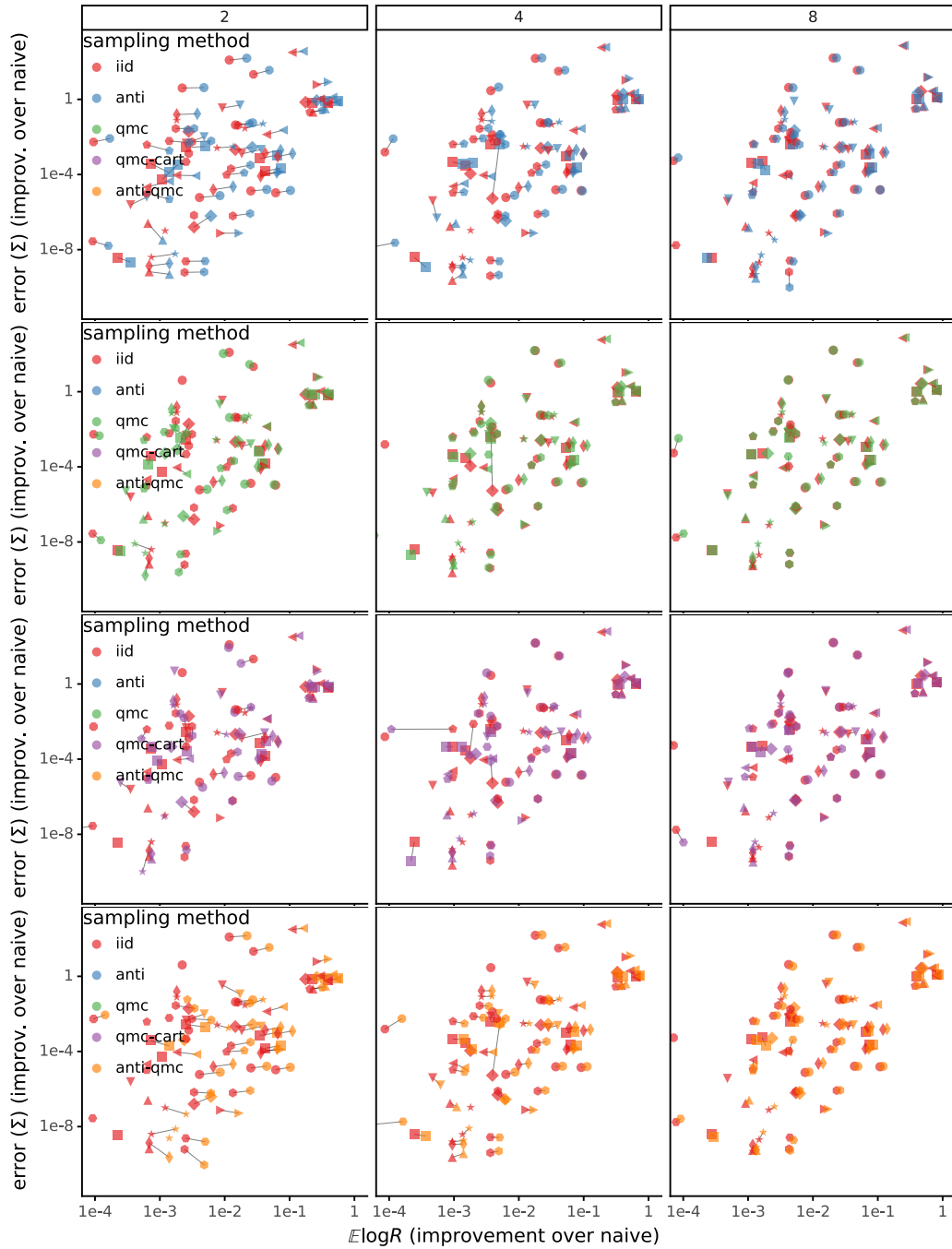


Figure 9: **How much do methods improve over naive VI in likelihood bound (x-axis) and in estimating posterior variance (y-axis)?** Each point corresponds to a model from the Stan library, with a random shape. Each plot compares iid sampling against some other strategy. From top, these are antithetic sampling (anti), Quasi-Monte Carlo, either using an elliptical mapping (qmc) or a Cartesian mapping (qmc-cart), and antithetic sampling after an elliptical mapping (anti-qmc). The columns correspond to using $M = 2, 4$ and 8 samples for each estimate. **Conclusions:** Improvements in ELBO and error are correlated. Improvements converge to those of iid for larger M , as all errors decay towards zero. Different sampling methods are best on different datasets. A few cases are not plotted where the measured “improvement” was negative (if naive VI has near-zero error, or due to local optima).

12 Full Results For All Models

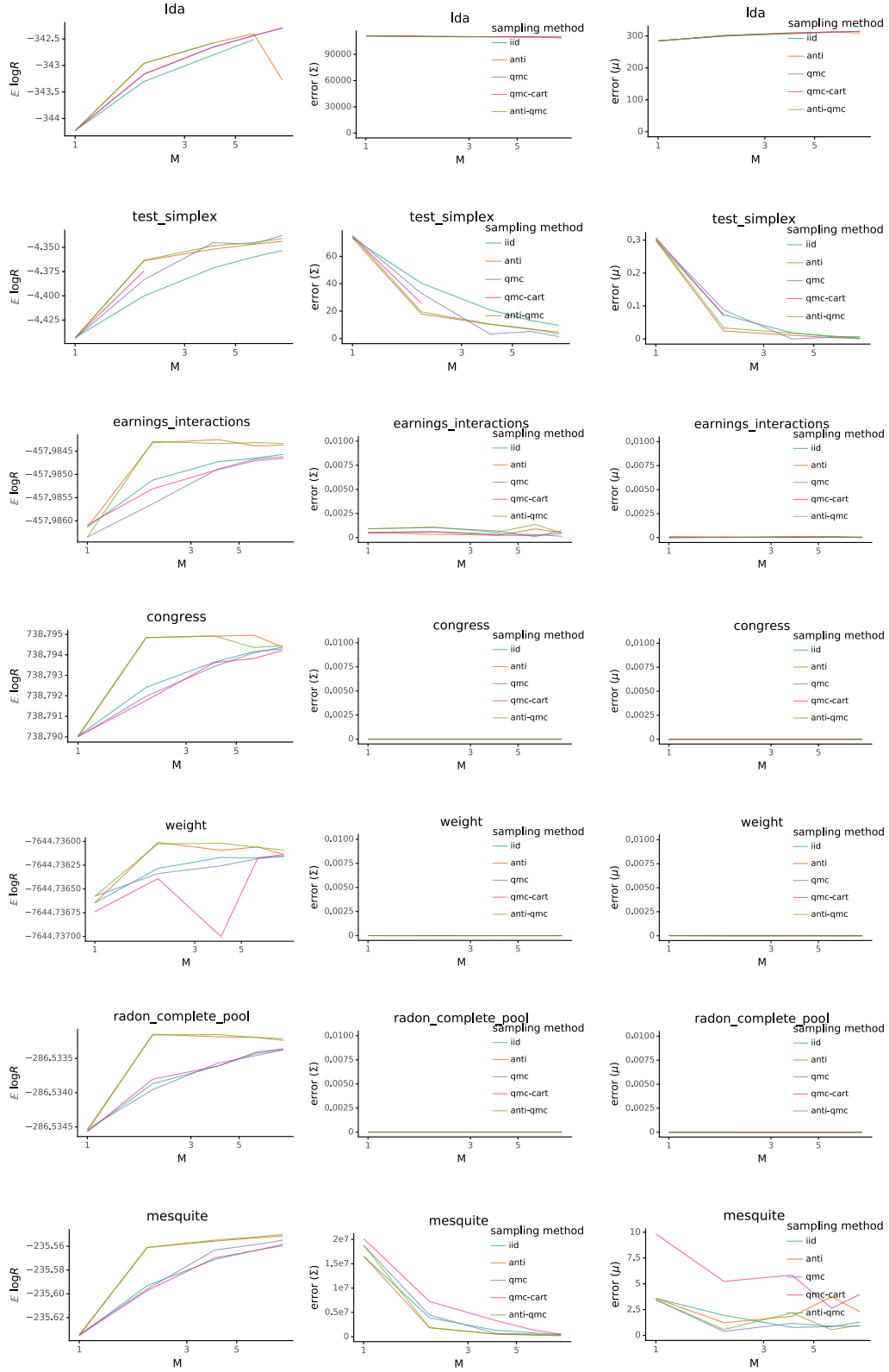


Figure 10: Across all models, improvements in likelihood bounds correlate strongly with improvements in posterior accuracy. Better sampling methods can improve both.

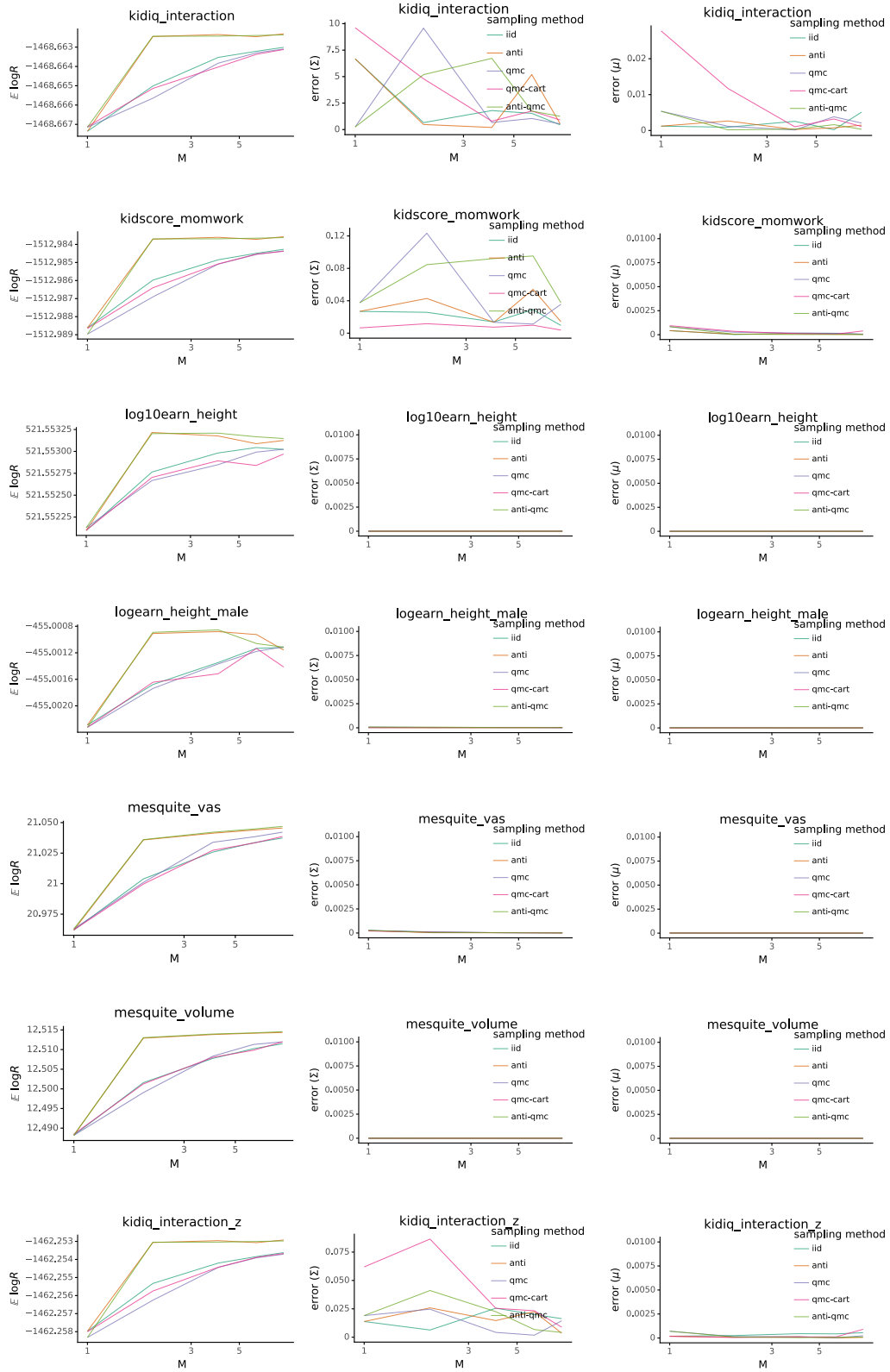


Figure 11: Across all models, improvements in likelihood bounds correlate strongly with improvements in posterior accuracy. Better sampling methods can improve both.

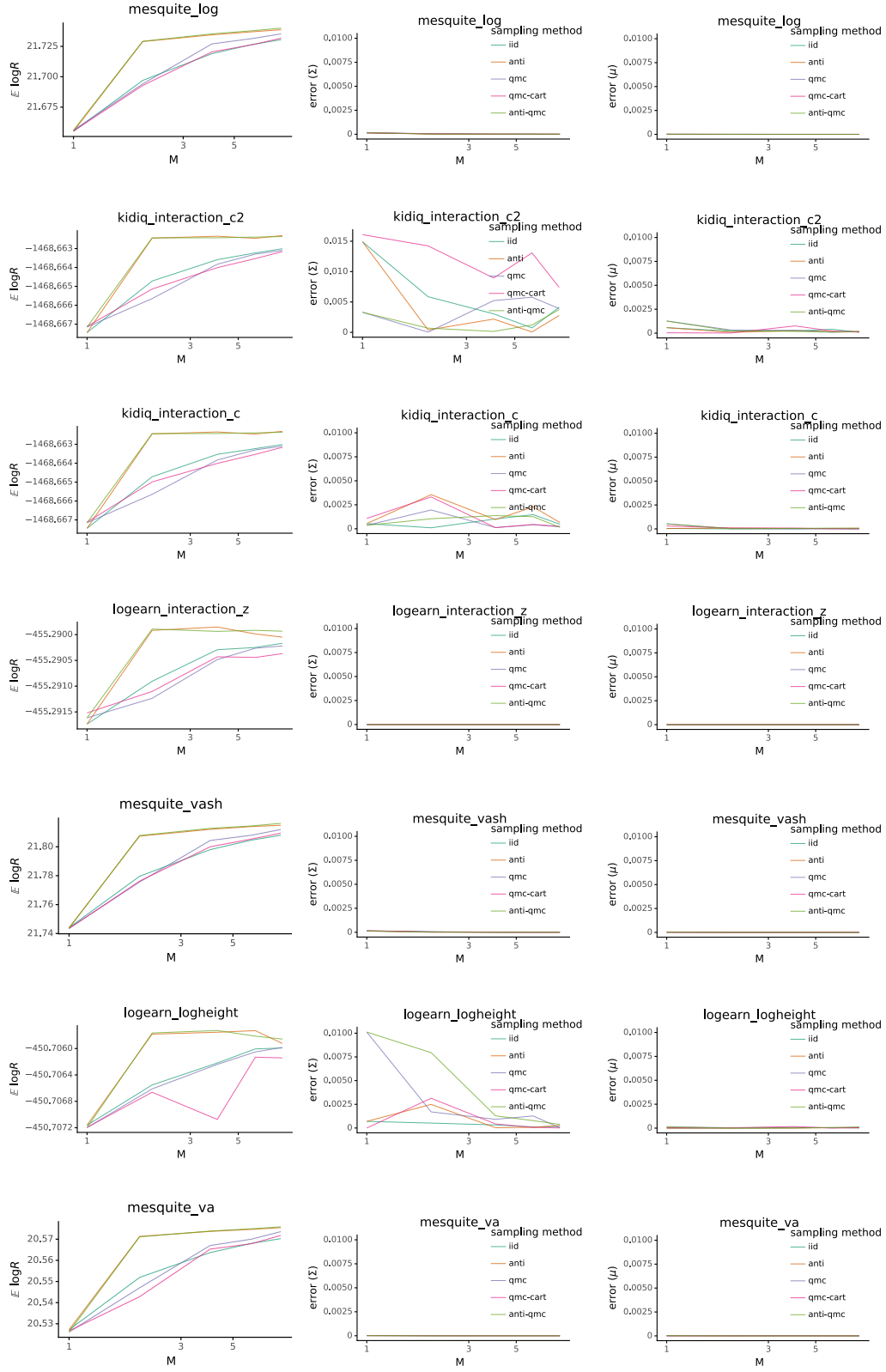


Figure 12: Across all models, improvements in likelihood bounds correlate strongly with improvements in posterior accuracy. Better sampling methods can improve both.

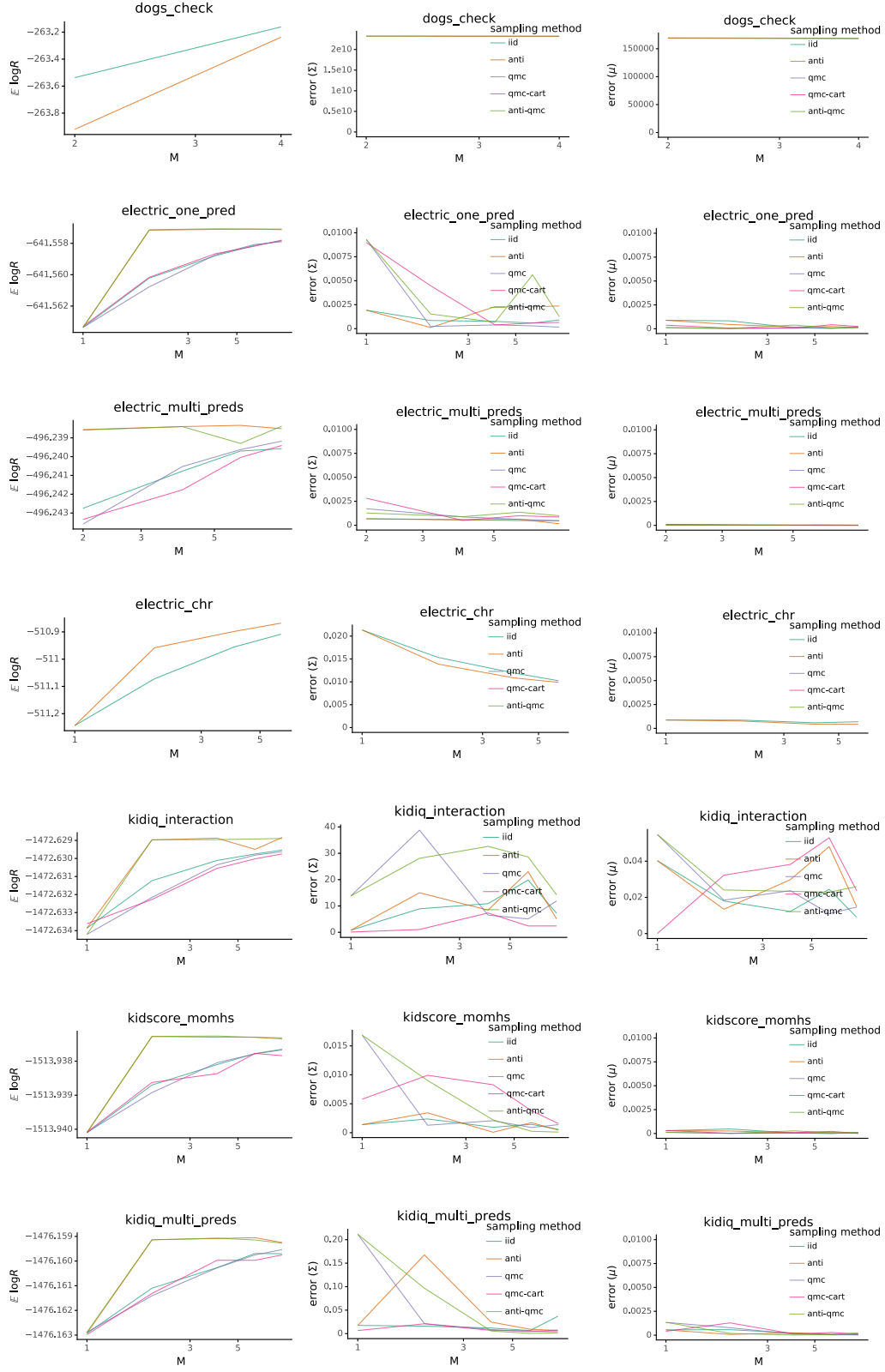


Figure 13: Across all models, improvements in likelihood bounds correlate strongly with improvements in posterior accuracy. Better sampling methods can improve both.

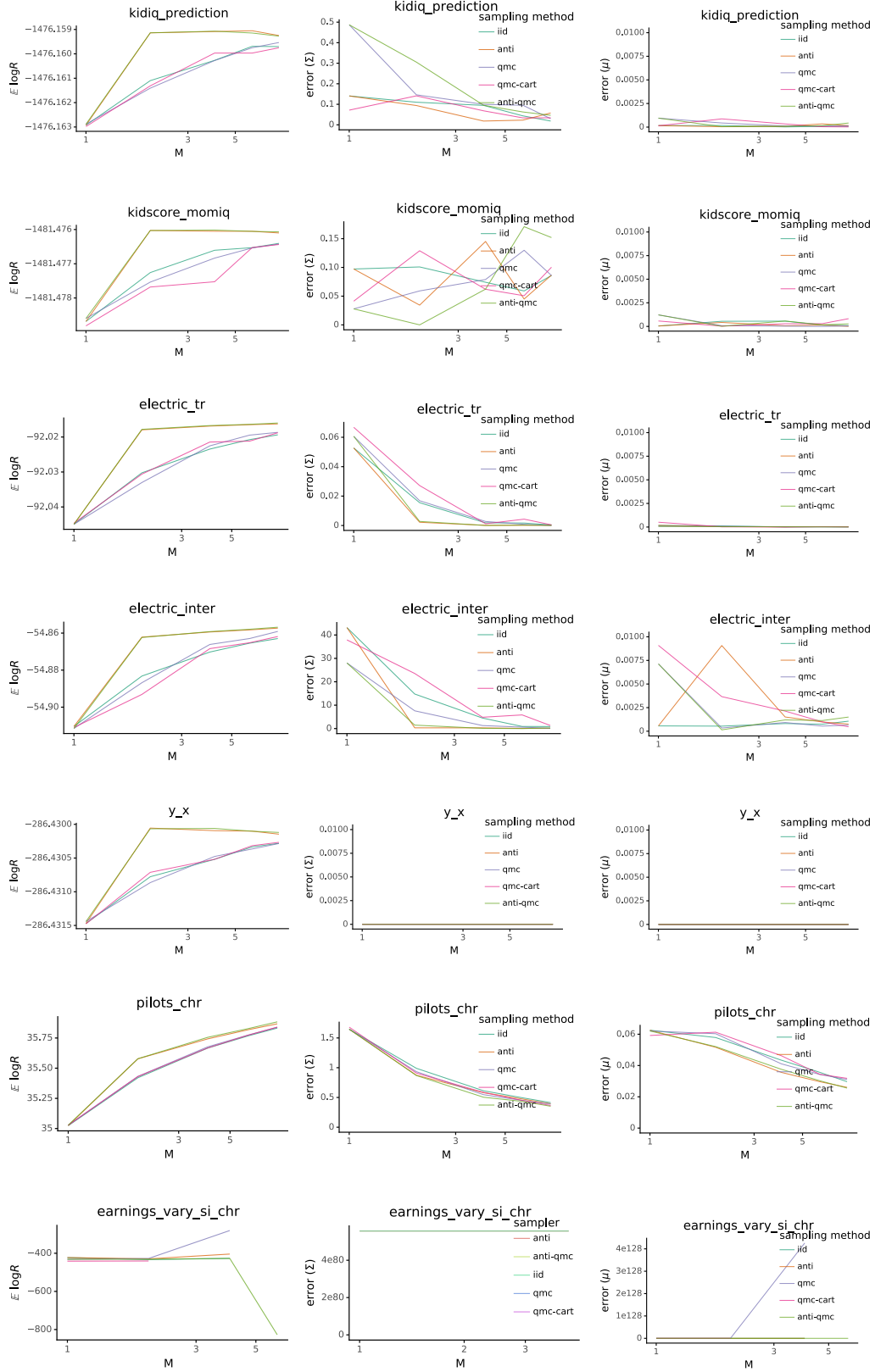


Figure 14: Across all models, improvements in likelihood bounds correlate strongly with improvements in posterior accuracy. Better sampling methods can improve both.

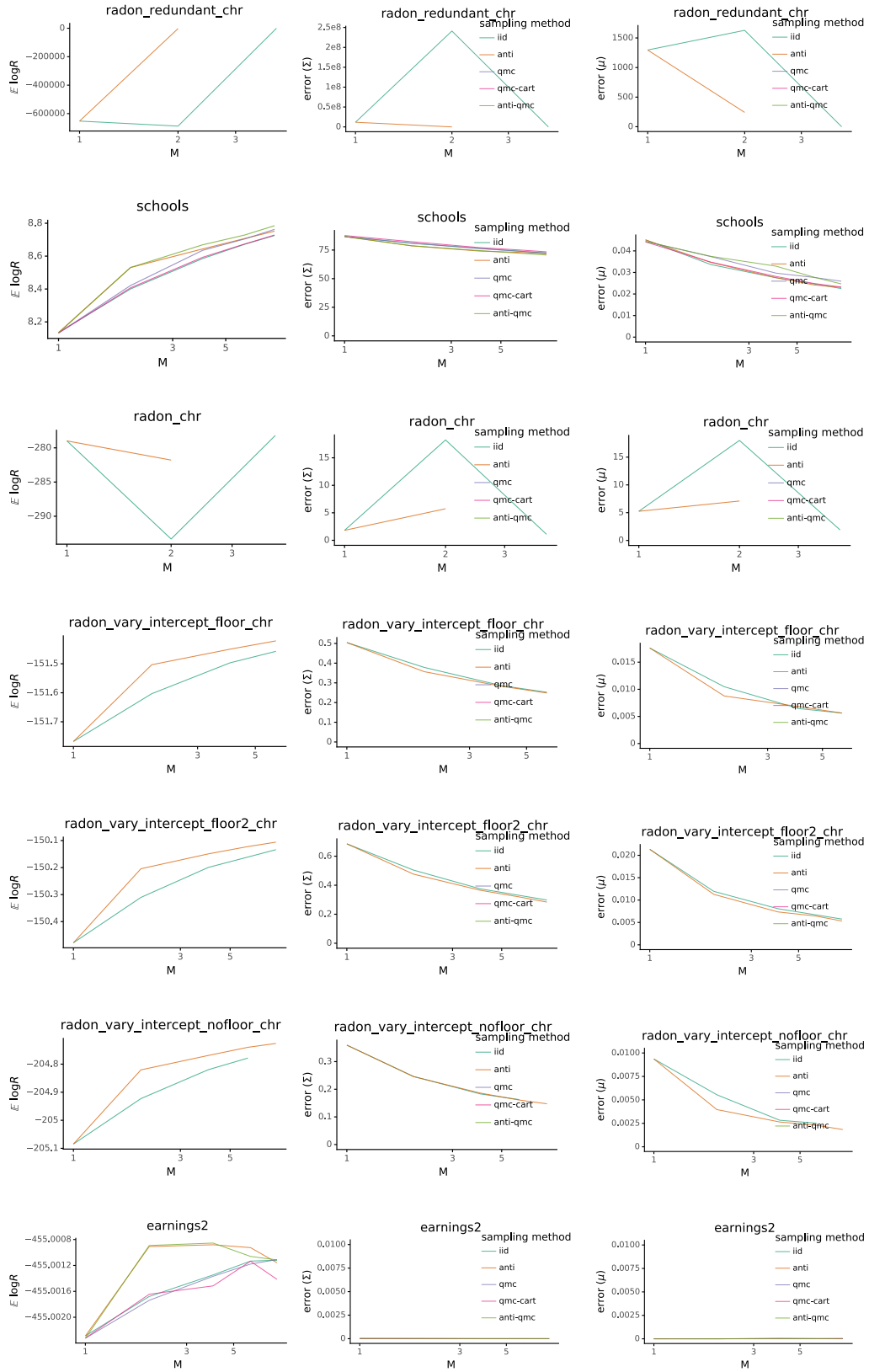


Figure 15: Across all models, improvements in likelihood bounds correlate strongly with improvements in posterior accuracy. Better sampling methods can improve both.

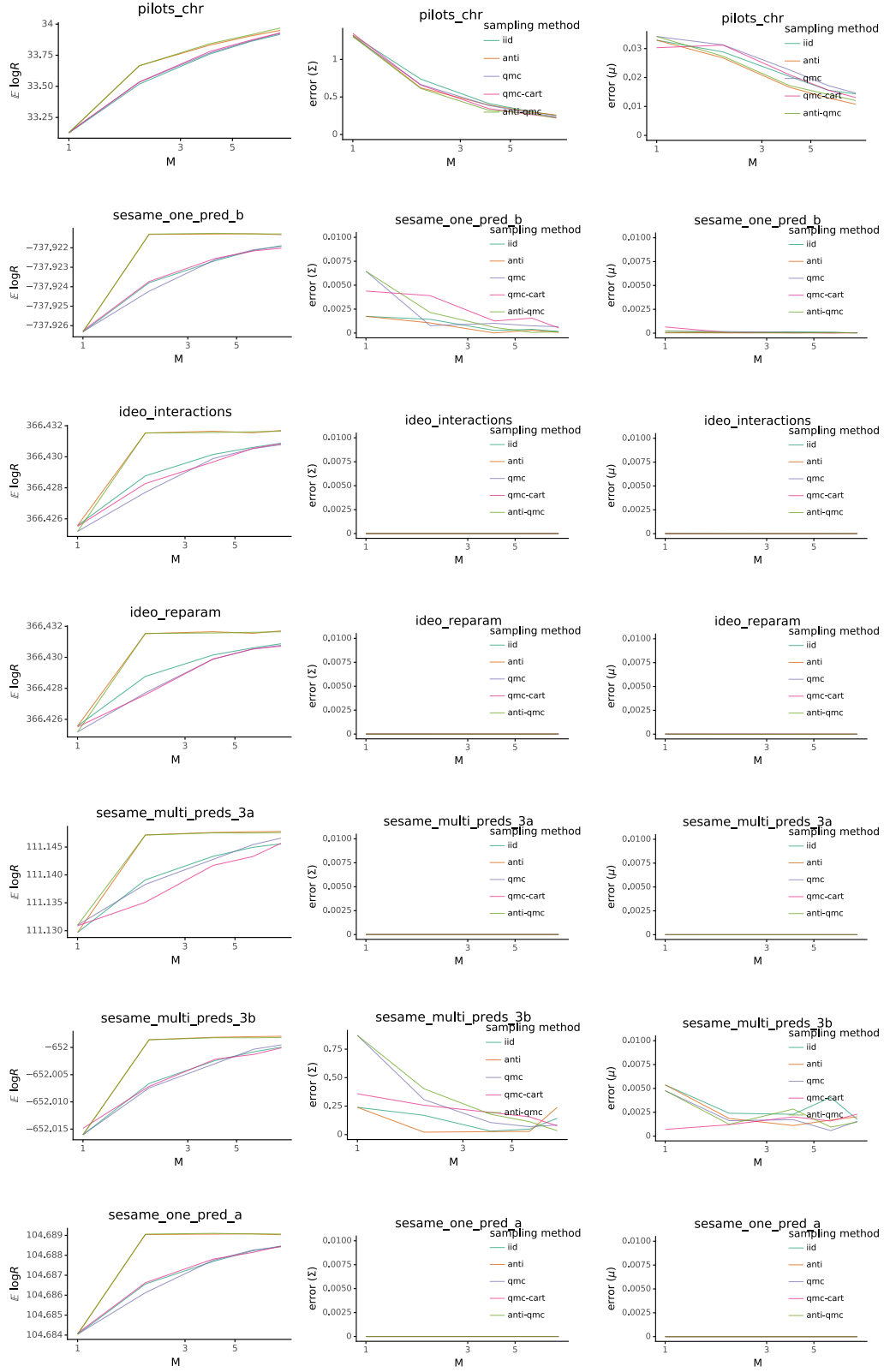


Figure 16: Across all models, improvements in likelihood bounds correlate strongly with improvements in posterior accuracy. Better sampling methods can improve both.

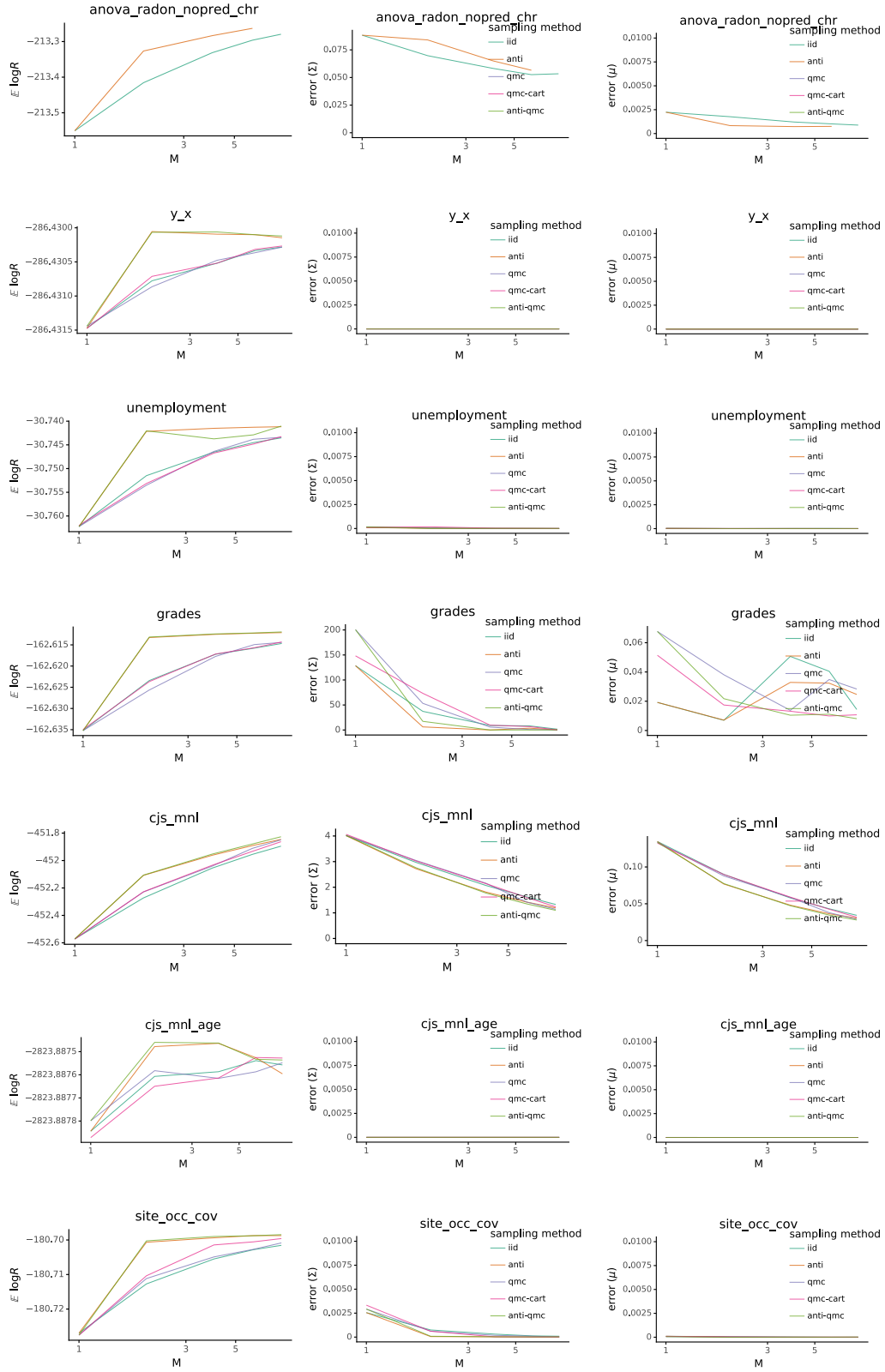


Figure 17: Across all models, improvements in likelihood bounds correlate strongly with improvements in posterior accuracy. Better sampling methods can improve both.

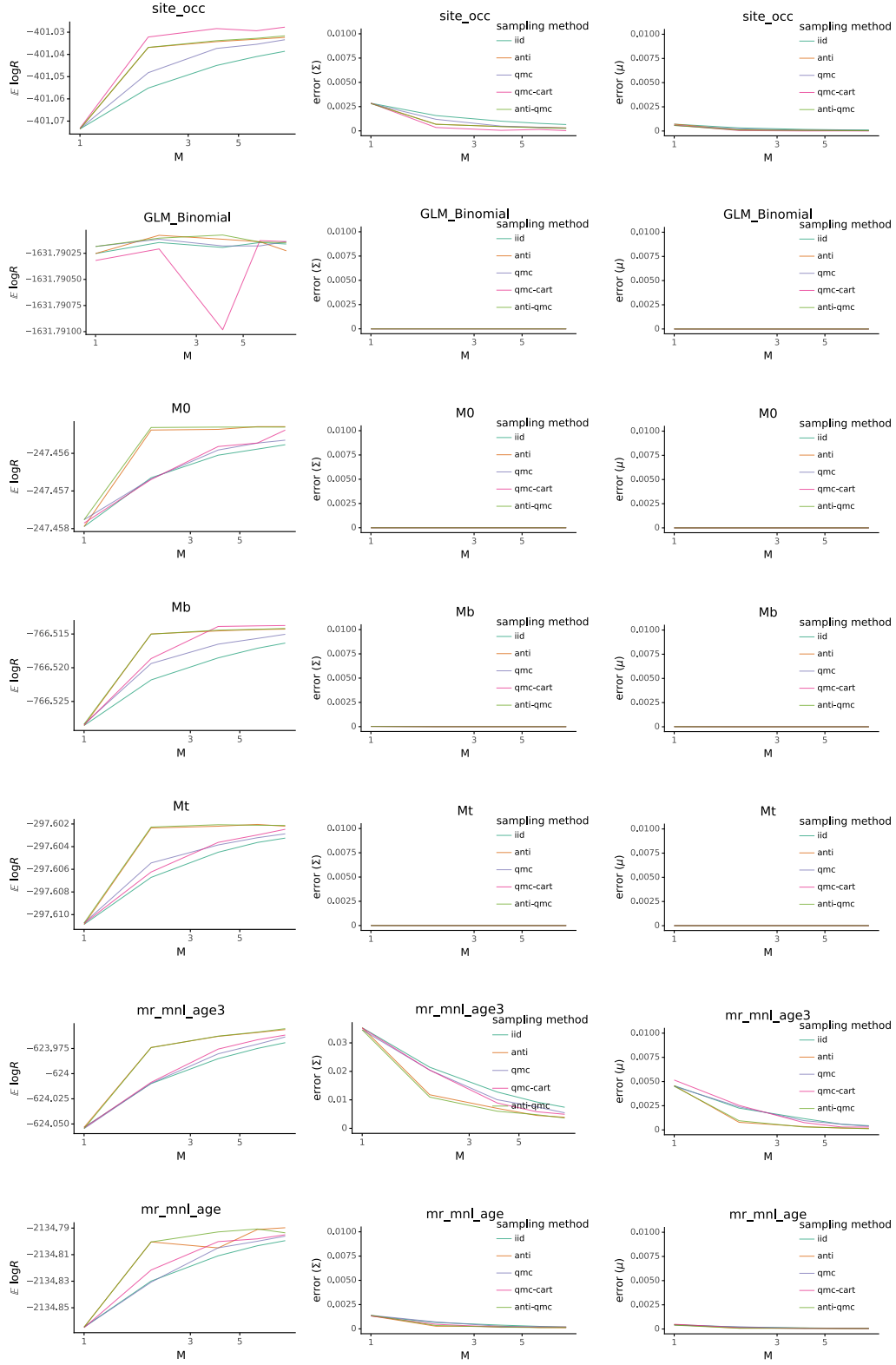


Figure 18: Across all models, improvements in likelihood bounds correlate strongly with improvements in posterior accuracy. Better sampling methods can improve both.

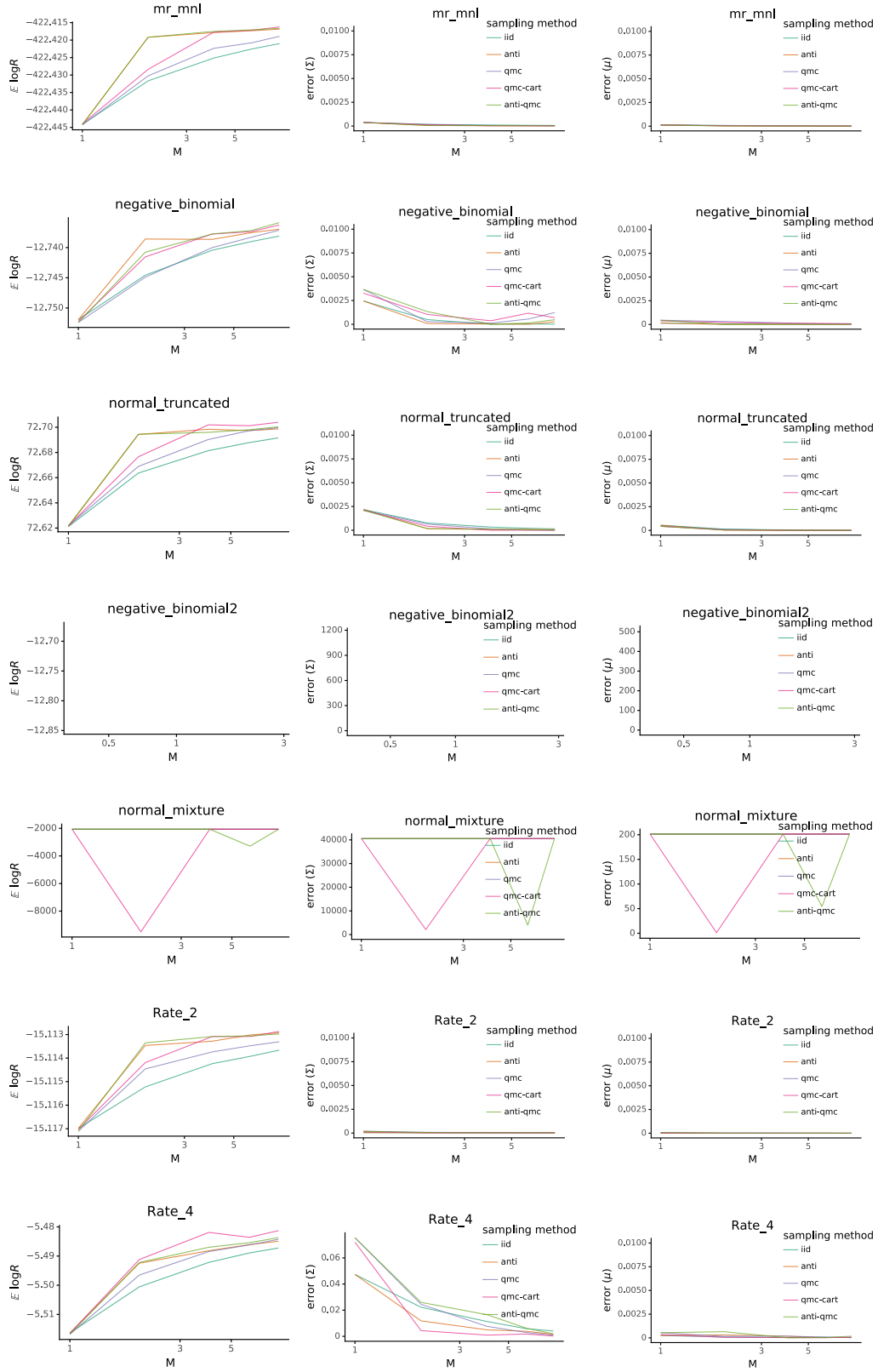


Figure 19: Across all models, improvements in likelihood bounds correlate strongly with improvements in posterior accuracy. Better sampling methods can improve both.

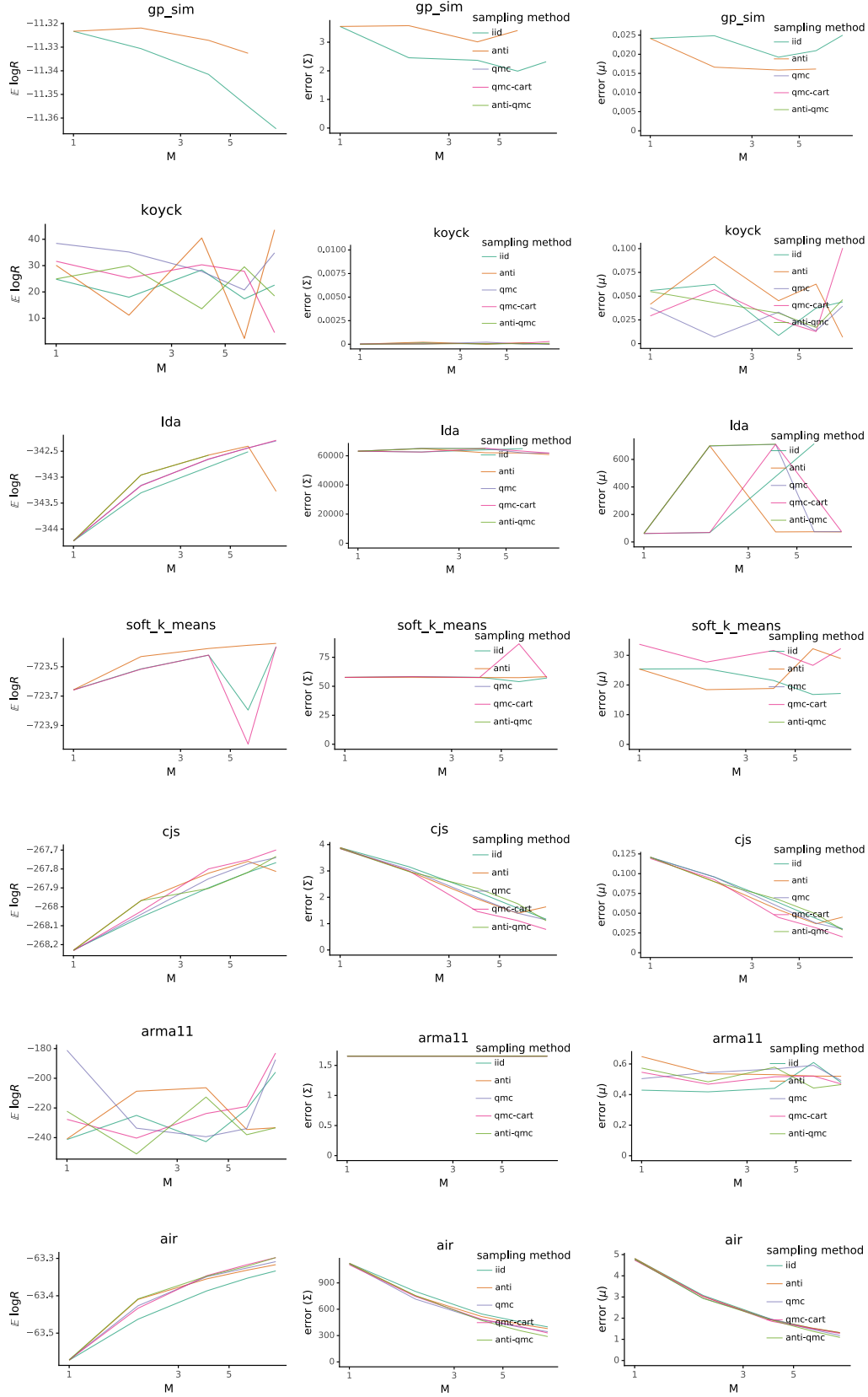


Figure 20: Across all models, improvements in likelihood bounds correlate strongly with improvements in posterior accuracy. Better sampling methods can improve both.

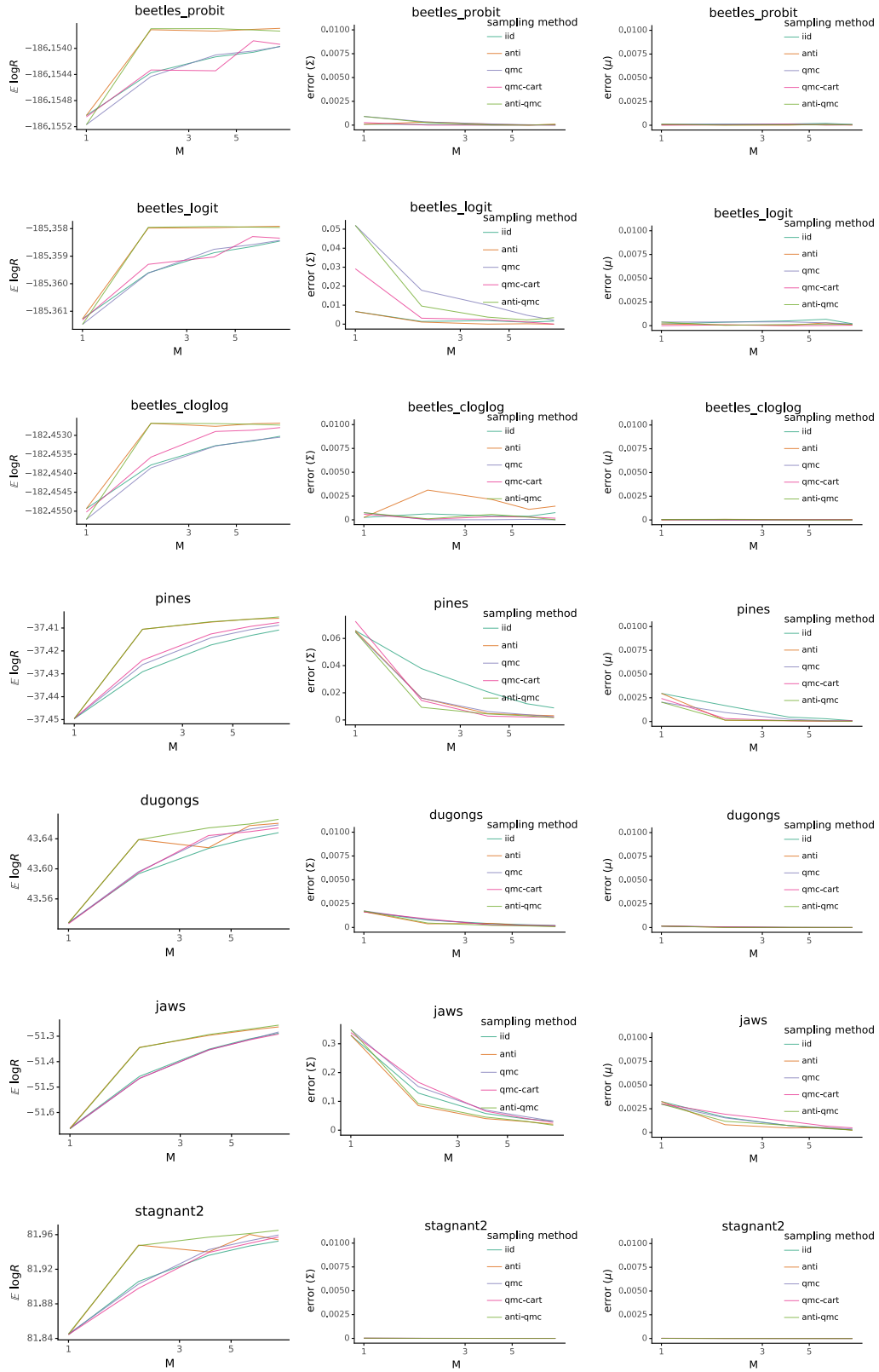


Figure 21: Across all models, improvements in likelihood bounds correlate strongly with improvements in posterior accuracy. Better sampling methods can improve both.

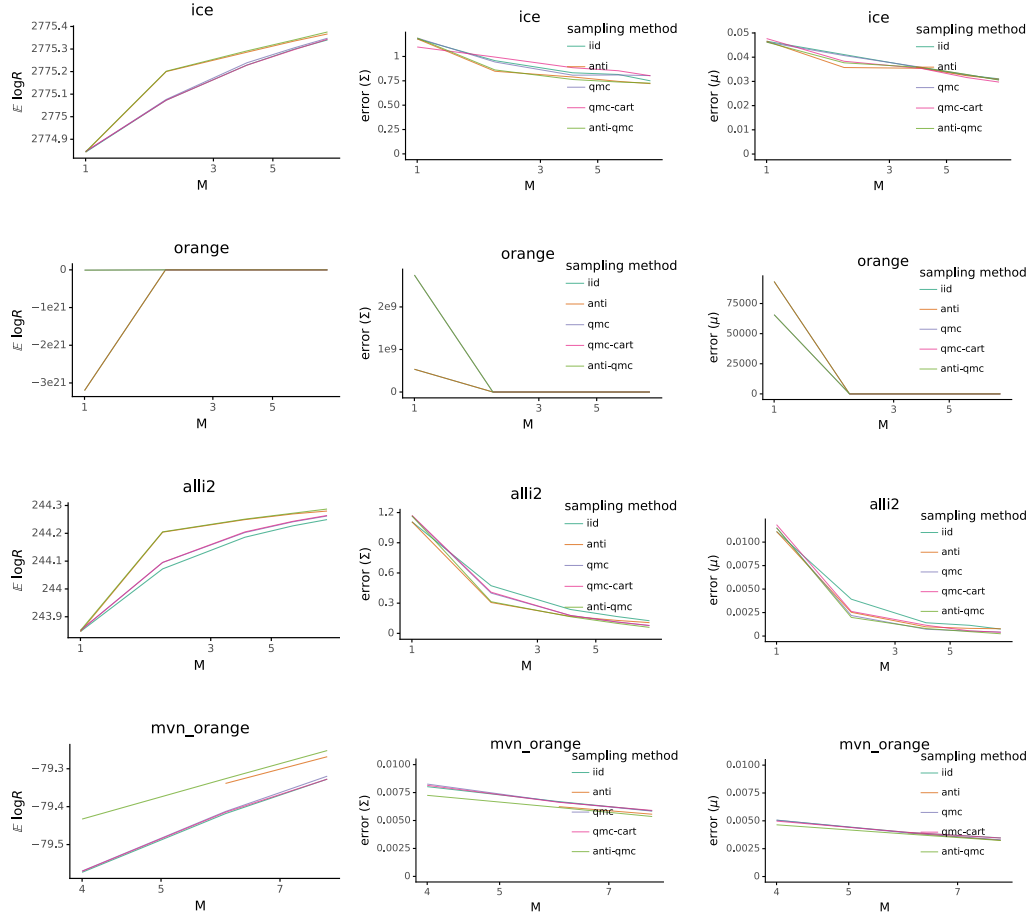


Figure 22: Across all models, improvements in likelihood bounds correlate strongly with improvements in posterior accuracy. Better sampling methods can improve both.