

Rfigre

# On Variational Bayesian Methods

## DRAFT V10

David Ewing

1st January 2026

### Abstract

Probabilistic machine learning often involves models with hidden (latent) variables. The main goal is to uncover the latent structure  $Z$  that explains the observed data  $X$ . In many practical settings, however, the exact posterior distribution  $p(z \mid x)$  cannot be computed in closed form or evaluated efficiently. Variational Inference (VI) addresses this challenge by reframing posterior inference as an optimisation problem. Instead of working with the true posterior, a tractable family of distributions  $q_\nu(z)$ , is chosen, and parameterised by variational parameters  $\nu$ . A member of this family that is closest to the true posterior, is chosen, typically measured using Kullback–Leibler divergence. This optimisation is equivalent to maximising the Evidence Lower Bound (ELBO), which provides a lower bound on the log marginal likelihood of the data. In this paper, I outline the core ideas behind VI, explain the role of the ELBO, and show how common choices of variational families—such as fixed-form and mean-field approximations—fit into the broader framework.

## 1 Probabilistic Modelling and the Inference Problem

### Bayesian and frequentist goals (why the choice matters)

Both *Bayesian* and *frequentist* statistics are valid approaches to inference, and each targets different outputs. The choice is driven by the goal of the analysis.

- **Frequentist goal:** estimate unknown parameters (treated as fixed constants) and quantify uncertainty via repeated-sampling properties of the estimator (for example, standard errors and confidence intervals).
- **Bayesian goal:** represent uncertainty about unknown quantities using probability distributions, producing a posterior distribution over latent variables and parameters (for example, posterior means, credible intervals, and posterior predictive distributions).

A simple example illustrates the distinction. Suppose  $Y_1, \dots, Y_n \sim N(\mu, \sigma^2)$  with known  $\sigma^2$ . A frequentist analysis reports an estimate  $\hat{\mu} = \bar{Y}$  and a 95% confidence interval  $\bar{Y} \pm 1.96 \sigma / \sqrt{n}$ , where the 95% refers to long-run coverage of the interval construction under repeated sampling. A Bayesian analysis places a prior on  $\mu$  (for example,  $\mu \sim N(\mu_0, \tau_0^2)$ ) and then reports the posterior distribution  $p(\mu \mid y)$ ; a 95% credible interval refers to posterior probability mass given the model, prior, and observed data.

Probabilistic modelling provides a systematic way to connect our assumptions about the world to the data observed. A model is specified in terms of:

- *Observed variables*  $X$ , representing the data;
- *Latent variables* or *parameters*  $Z$ , representing unknown structure or quantities of interest;

- A joint distribution  $p(x, z)$  that encodes our assumptions and prior information.

The central inferential object is the *posterior distribution*

$$p(z | x) = \frac{p(x, z)}{p(x)}, \quad p(x) = \int p(x, z) dz,$$

which tells us how plausible different configurations of  $Z$  are, given the observed data  $X$ .

In simple models,  $p(z | x)$  can be derived analytically. However, in many modern applications (hierarchical models, complex latent-variable models, deep generative models, and so on) the marginal likelihood  $p(x)$  is intractable: the integral cannot be evaluated exactly and is expensive to approximate directly. This makes exact posterior inference infeasible, and motivates approximate methods such as Markov chain Monte Carlo (MCMC) and Variational Inference.

## 2 Approximate Bayesian Inference in Practice

### Variational inference as Bayesian inference (and how it differs from other Bayesian workflows)

Up to this point, I have described probabilistic modelling in terms of a joint distribution  $p(x, z)$  and an inferential target that is naturally expressed as a distribution. In Bayesian work, that target is the posterior distribution  $p(z | x)$ . In realistic models, this posterior is rarely available in closed form, so we turn to *approximate Bayesian inference* methods.

There are several common ways to approximate the posterior:

- **Simulation-based Bayes (MCMC):** construct a Markov chain with stationary distribution  $p(z | x)$  and use draws to approximate expectations and uncertainty. This is widely treated as a reference standard for accuracy, but can be slow in high dimensions or hierarchical settings.
- **Deterministic local approximations (Laplace):** approximate the posterior by a Gaussian expansion around a mode (often the MAP), giving fast approximate uncertainty when the posterior is close to normal, but potentially misleading summaries when the posterior is skewed or multi-modal.
- **Optimisation-based Bayes (Variational Inference):** choose a tractable family  $\mathcal{Q}$  and optimise within that family to obtain  $q_{\nu^*}(z)$  as a proxy for  $p(z | x)$ . This is typically much faster than MCMC, but (especially under mean-field factorisations) can underestimate uncertainty.

In other words: VI is not an alternative to Bayesian inference; it is one way of doing Bayesian inference when exact posterior computation is infeasible. The difference is not the inferential target (still a posterior), but the computational route taken to approximate it.

### Convergence properties and computational differences

These three methods differ not only in speed but also in their convergence guarantees and computational patterns:

**MCMC** iteratively samples from the posterior. At each iteration  $t$ , the algorithm evaluates the likelihood  $p(x | z^{(t)})$  and prior  $p(z^{(t)})$  for a proposed state  $z^{(t)}$ . Under suitable conditions (ergodicity), the Markov chain converges to the true posterior distribution:

$$q_{\text{MCMC}}(z) \xrightarrow[t \rightarrow \infty]{} p(z | x).$$

In the limit, MCMC provides unbiased estimates of posterior expectations. In practice, convergence diagnostics (such as  $\hat{R}$  and effective sample size) help assess whether the chain has run long enough.

**Laplace approximation** computes the posterior mode (MAP) once:

$$z^* = \arg \max_z \log p(x, z),$$

and then approximates the posterior as Gaussian:

$$p(z \mid x) \approx N(z^*, H^{-1}),$$

where  $H$  is the negative Hessian of  $\log p(x, z)$  evaluated at  $z^*$ . This requires computing the likelihood and its derivatives at a single point (the MAP), making it extremely fast. However, the approximation is only accurate when the posterior is close to Gaussian; for skewed or multimodal posteriors, Laplace can be misleading.

**Variational Inference** optimises over a family of distributions  $q_\nu(z)$ . Unlike MCMC, which evaluates the likelihood at many sampled points, VI computes (or estimates via sampling) expectations of the log-likelihood under the current approximation:

$$\mathbb{E}_{q_\nu}[\log p(x, z)].$$

The algorithm adjusts  $\nu$  to maximise the ELBO, effectively minimising  $\text{KL}(q_\nu \| p(z \mid x))$ . The solution

$$q_{\nu^*}(z) = \arg \min_{q \in \mathcal{Q}} \text{KL}(q \| p(z \mid x))$$

is the best approximation within the chosen family  $\mathcal{Q}$ . However, if the true posterior lies outside  $\mathcal{Q}$  (as it typically does for restrictive families like mean-field), then even the optimal  $q_{\nu^*}$  will be systematically biased.

## Asymptotic behavior as sample size increases

All three methods improve as the dataset size  $n \rightarrow \infty$ , but in different ways:

**MCMC:** The posterior concentrates around the true parameter value at rate  $O(1/\sqrt{n})$  (Bernstein-von Mises theorem). MCMC draws from this concentrating distribution, so posterior summaries converge to the truth. Monte Carlo error (from finite chain length) decreases with more iterations, independently of  $n$ .

**Laplace:** As  $n \rightarrow \infty$ , under regularity conditions, the posterior becomes approximately Gaussian (Bernstein-von Mises). Thus, the Laplace approximation becomes asymptotically exact:

$$\text{KL}(p(z \mid x) \| N(z^*, H^{-1})) \xrightarrow{n \rightarrow \infty} 0.$$

For large datasets, Laplace can be highly accurate at a fraction of the computational cost of MCMC.

**Variational Inference:** The posterior also concentrates as  $n \rightarrow \infty$ . If the variational family  $\mathcal{Q}$  is flexible enough to contain the limiting Gaussian, then

$$\text{KL}(q_{\nu^*} \| p(z \mid x)) \xrightarrow{n \rightarrow \infty} 0.$$

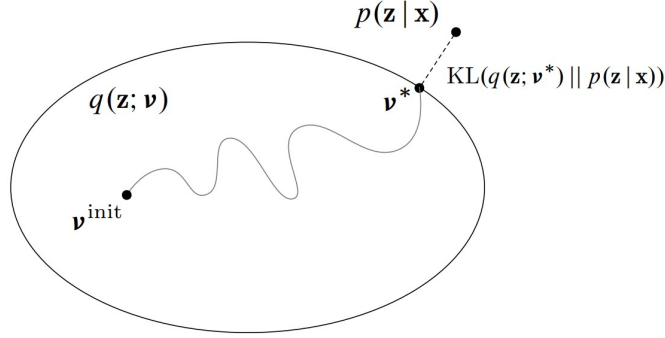
In this asymptotic regime, VI becomes accurate. However, for finite  $n$ , the gap between  $q_{\nu^*}$  and  $p(z \mid x)$  depends on how restrictive the family  $\mathcal{Q}$  is. Mean-field approximations, which enforce independence among parameters, often remain systematically biased even for moderately large  $n$ , particularly for variance components in hierarchical models.

**Summary:** MCMC converges to the true posterior (up to Monte Carlo error). Laplace and VI both produce approximations that become asymptotically accurate as  $n \rightarrow \infty$ , provided the posterior becomes Gaussian and (for VI) the variational family is sufficiently flexible. For finite samples, MCMC is the gold standard, but Laplace and VI trade off accuracy for computational speed. VI's systematic bias is most problematic in hierarchical models with restrictive mean-field factorisations.

### 3 Variational Inference as Optimisation

Variational Inference approaches posterior inference through optimisation rather than sampling. The idea is to replace the exact posterior  $p(z | x)$  with a simpler, tractable distribution  $q_\nu(z)$  drawn from a chosen family. Figure 1 provides the standard picture: attention is restricted to a variational family  $q(z; \nu)$  and then optimise  $\nu$  so that  $q(z; \nu)$  is close (in KL divergence) to the true posterior  $p(z | x)$ .

#### Variational Inference



- VI turns **inference into optimization**.
- Posit a **variational family** of distributions over the latent variables,

$$q(z; \nu)$$

- Fit the **variational parameters**  $\nu$  to be close (in KL) to the exact posterior.  
(There are alternative divergences, which connect to algorithms like EP, BP and others.)

Figure 1: Variational inference as optimisation over variational parameters  $\nu$  within a variational family  $q(z; \nu)$ , choosing  $\nu^*$  to minimise  $\text{KL}(q(z; \nu) || p(z | x))$  [?].

The ellipse represents all families of tractable approximations, indexed by the variational parameters  $\nu$ . Starting from an initial value  $\nu^{\text{init}}$ , an optimisation routine moves through parameter space (the grey path) to reach  $\nu^*$ , the best approximation available within the family. The true posterior  $p(z | x)$  sits outside the ellipse because it is generally too complex to belong to the variational family, and the dashed segment indicates the remaining discrepancy measured by  $\text{KL}(q(z; \nu^*) || p(z | x))$ . In the next subsection we define the family  $\mathcal{Q}$  formally and make precise how  $\nu^*$  is determined.

#### 3.1 The Variational Family

We first specify a family of candidate distributions

$$\mathcal{Q} = \{q_\nu(z) : \nu \in \mathcal{V}\},$$

where  $\nu$  denotes the collection of *variational parameters*. This family might be defined, for example, by:

- A parametric form (e.g. multivariate Gaussian with mean vector and covariance matrix);
- A factorisation assumption (e.g. a product of simpler distributions across components of  $Z$ );
- Some combination of structural and parametric choices.

The goal of VI is then to choose

$$\nu^* = \arg \min_{\nu} \text{KL}(q_\nu(z) || p(z | x)),$$

and to use  $q_{\nu^*}(z)$  as our approximation to the posterior. Once  $\nu^*$  has been found, this approximate posterior can be used for point estimates, interval estimates, prediction, and other downstream tasks.

### 3.2 KL Divergence and the ELBO

The Kullback–Leibler (KL) divergence between  $q_\nu(z)$  and  $p(z \mid x)$  is defined as

$$\text{KL}(q_\nu(z) \parallel p(z \mid x)) = \mathbb{E}_{q_\nu} \left[ \log \frac{q_\nu(z)}{p(z \mid x)} \right].$$

Directly minimising this quantity is usually not possible, because it depends on the intractable marginal likelihood  $p(x)$ . Instead, we work with the *Evidence Lower Bound* (ELBO), defined by

$$\mathcal{L}(\nu) = \mathbb{E}_{q_\nu} [\log p(x, z)] - \mathbb{E}_{q_\nu} [\log q_\nu(z)].$$

It can be shown that

$$\log p(x) = \mathcal{L}(\nu) + \text{KL}(q_\nu(z) \parallel p(z \mid x)),$$

so that for fixed data  $x$ , maximising  $\mathcal{L}(\nu)$  is equivalent to minimising the KL divergence. The ELBO thus serves as a surrogate objective that we can evaluate (or approximate) using only  $p(x, z)$  and  $q_\nu(z)$ .

The ELBO has a useful interpretation as a balance between two terms:

- An *expected log joint* term,  $\mathbb{E}_{q_\nu} [\log p(x, z)]$ , which encourages  $q_\nu(z)$  to place mass on configurations of  $z$  that explain the data well.
- An *entropy* term,  $-\mathbb{E}_{q_\nu} [\log q_\nu(z)]$ , which encourages  $q_\nu(z)$  to be diffuse and to avoid collapsing onto a single point.

Optimising the ELBO therefore trades off goodness-of-fit against complexity, in a way that is closely related to (but distinct from) classical regularisation ideas.

### 3.3 Different Views of the ELBO

Although the ELBO is a single mathematical quantity, it can be written and interpreted in several equivalent ways. Different algorithms tend to emphasise different views. It is helpful to keep three of them in mind.

#### Projection view (KL divergence)

We can write,

$$\log p(x) = \mathcal{L}(\nu) + \text{KL}(q_\nu(z) \parallel p(z \mid x)).$$

For fixed data  $x$ , the term  $\log p(x)$  is a constant. Maximising the ELBO is therefore exactly the same as minimising  $\text{KL}(q_\nu(z) \parallel p(z \mid x))$ .

In this view, variational inference is a projection operation: we project the true posterior onto the chosen family of approximations, and the discrepancy is measured by the KL divergence from  $q_\nu$  to  $p(\cdot \mid x)$ .

#### Energy–entropy view

The ELBO can also be written as

$$\mathcal{L}(\nu) = \mathbb{E}_{q_\nu} [\log p(x, z)] - \mathbb{E}_{q_\nu} [\log q_\nu(z)].$$

The first term,  $\mathbb{E}_{q_\nu}[\log p(x, z)]$ , is sometimes called an “energy” term: it rewards placing probability mass on configurations of  $z$  that explain the data and respect the prior. The second term,  $-\mathbb{E}_{q_\nu}[\log q_\nu(z)]$ , is the entropy of  $q_\nu$ , which rewards spread or uncertainty in the approximation.

Optimising the ELBO can then be seen as a trade-off:

- improve the fit to the data and prior by increasing the expected log joint;
- avoid collapsing  $q_\nu$  onto a single point by maintaining entropy.

This energy–entropy view is especially useful when thinking about coordinate-ascent algorithms and mean-field approximations.

### Likelihood–regularisation view

A third, equally valid, way to arrange the same terms is

$$\mathcal{L}(\nu) = \mathbb{E}_{q_\nu}[\log p(x | z)] - \text{KL}(q_\nu(z) \| p(z)).$$

Here the ELBO looks like:

- an expected log-likelihood term,  $\mathbb{E}_{q_\nu}[\log p(x | z)]$ ;
- minus a KL penalty that keeps  $q_\nu(z)$  close to the prior  $p(z)$ .

Variational inference resembles a regularised fitting problem: we try to choose  $q_\nu$  so that it explains the data well on average, but we pay a regularisation cost if  $q_\nu$  moves too far from the prior.

This perspective connects naturally to fixed-form and black box variational methods, where the ELBO is treated as an objective to optimise numerically [?, ?].

### A note on the direction of KL

The divergence that appears is  $\text{KL}(q_\nu(z) \| p(z | x))$ . This “exclusive” direction penalises  $q_\nu$  for placing mass in regions where the true posterior density is low.

One practical consequence, which will matter later, is that minimising  $\text{KL}(q_\nu \| p)$  tends to produce approximations that concentrate on one main mode of the posterior and avoid spreading mass into low-density regions. This often leads to underestimation of posterior uncertainty, especially in simple mean-field approximations [?, ?].

Other divergences and other directions of KL are possible, but in this introductory note we focus on this standard choice, because it underpins most of the mean-field and fixed-form methods discussed next.

### Why these views matter for algorithms

Although the ELBO is the same object in all three views, different formulations support different algorithmic strategies:

- The projection view emphasises VI as “KL minimisation” and is useful for high-level understanding.
- The energy–entropy view fits naturally with coordinate-ascent and mean-field factorisations, where we update one factor at a time to improve the balance between fit and entropy.

- The likelihood–regularisation view fits naturally with fixed-form and black box VI, where we treat the ELBO as a regularised objective and optimise it with gradient-based methods [?, ?].

*In latent-variable models, the ELBO is engineered so that maximising it is exactly the same as minimising  $\text{KL}(q_\nu(z) \parallel p(z \mid x))$ ; outside that setting, KL minimisation need not come with an “ELBO” attached.*

In later sections, when we focus on mean-field VI and fixed-form VI, these interpretations will reappear in slightly different guises.

## 4 Variational Families: Fixed-Form and Mean-Field (Preview)

A key modelling choice in VI is how we define the family  $\mathcal{Q}$ . Two broad patterns that appear repeatedly in the literature are:

### 4.1 Fixed-Form Variational Inference

In *fixed-form* variational inference, we choose a specific parametric family (for example, multivariate Gaussian distributions) and constrain  $q_\nu(z)$  to lie within that family for all models under consideration. The focus is then on developing general optimisation methods that can handle any model  $p(x, z)$  while keeping  $q_\nu(z)$  in this shared, tractable form [?, ?].

### 4.2 Mean-Field Variational Inference

In *mean-field* variational inference, we assume that the latent variables factorise under  $q_\nu$ :

$$q_\nu(z) = \prod_j q_{\nu_j}(z_j),$$

where  $z = (z_1, \dots, z_J)$  is partitioned into components and each  $q_{\nu_j}(z_j)$  is a simpler distribution. This assumption is often unrealistic as a literal description of the true posterior, but it makes the optimisation problem more tractable and can lead to closed-form coordinate ascent updates in many models [?, ?].

An important limitation of mean-field approximations concerns their performance in hierarchical models with variance components. This under-dispersion phenomenon is discussed in detail in Section 6.

A detailed comparison of fixed-form, mean-field, and more structured approximations (including their strengths and limitations) will be developed in later sections. Here we simply note that the definition of  $\mathcal{Q}$  is central: it encodes both the computational tractability and the expressive power of the variational approximation.

## 5 Optimisation

Maximising the ELBO is achieved by numerical optimisation. Common strategies include:

- *Coordinate ascent* algorithms, particularly in conjugate-exponential models, where updates for each factor or parameter have closed-form expressions.



- *Gradient-based* methods, which use derivatives of  $\mathcal{L}(\nu)$  with respect to  $\nu$  and can be combined with modern optimisation techniques (such as variants of stochastic gradient descent) [?, ?].

The ELBO is typically not a convex function of  $\nu$ , so optimisation procedures may converge to local optima rather than a unique global solution. Nevertheless, in many applications the resulting approximate posteriors are accurate enough to support interpretation and prediction, especially when exact inference is infeasible.

Modern probabilistic programming frameworks and automatic differentiation tools make it increasingly straightforward to implement gradient-based VI for complex models, without having to derive analytic updates by hand. More advanced schemes—such as black box variational inference and reparameterisation-based stochastic gradients—build on the same core ideas presented here [?, ?].

## 6 Under-Dispersion in Variational Bayes

Variational inference provides fast approximate posterior inference, but the quality of the approximation depends critically on the choice of variational family  $\mathcal{Q}$ . One of the most important systematic failures of VI occurs in hierarchical models, where mean-field approximations tend to produce *under-dispersed* posteriors for variance components. This section explains the phenomenon, why it occurs, and what it means for practical inference.

### 6.1 The Under-Dispersion Phenomenon

Variational Bayes approximations are not equally accurate for all types of parameters. In general:

**Location parameters** (such as regression coefficients  $\beta$ , group means  $\mu$ , or random intercepts  $u_j$ ) tend to have:

- Posterior means that are reasonably accurate (close to MCMC estimates).
- Posterior variances that may be slightly underestimated, but often acceptable for practical purposes.

**Scale parameters** (such as standard deviations  $\sigma$ , variance components  $\sigma_u^2$ , or precision parameters  $\tau$ ) tend to have:

- Posterior distributions that are systematically *under-dispersed*: too narrow, with variance substantially smaller than the true posterior.
- Mass concentrated on smaller values, underestimating the true uncertainty about variability.

This asymmetry is especially pronounced in hierarchical or random-effects models, where variance components control the distribution of lower-level parameters. The under-dispersion of variance components has cascading effects: it leads to *over-shrinkage* of random effects and overconfident predictions.

#### Quantifying under-dispersion: the variance ratio diagnostic

One way to quantify under-dispersion is to compare posterior variances from variational Bayes with those from a reference MCMC fit. For each parameter  $\theta$ , we can form the variance ratio

$$\text{VR}(\theta) = \frac{\text{Var}_{\text{VB}}(\theta)}{\text{Var}_{\text{MCMC}}(\theta)}.$$

Values close to 1.0 indicate similar levels of uncertainty, while values substantially below 1.0 signal that the variational posterior is under-dispersed [?, ?].

#### Typical variance ratio values:

- For regression coefficients  $\beta$ : VR  $\approx$  0.8–0.95 (mild under-dispersion).
- For residual variance  $\sigma^2$ : VR  $\approx$  0.6–0.8 (moderate under-dispersion).
- For variance components  $\sigma_u^2$  in hierarchical models: VR  $\approx$  0.3–0.7 (severe under-dispersion).

In hierarchical models, variance ratios for variance components are often well below one, and this aligns with the visual impression of over-shrinkage in the random effects. Such variance-ratio summaries provide a compact way to report how far a variational approximation departs from a more accurate MCMC reference.

## 6.2 Why Mean-Field VI Causes Under-Dispersion in Hierarchical Models

The under-dispersion phenomenon is not inherent to variational inference itself, but arises specifically from the *mean-field factorisation* assumption when applied to hierarchical models. This subsection explains the mechanism.

### The hierarchical model structure

Consider a simple random-intercept model:

$$y_{ij} = \mu + u_j + \varepsilon_{ij}, \quad u_j \sim N(0, \sigma_u^2), \quad \varepsilon_{ij} \sim N(0, \sigma^2),$$

where  $y_{ij}$  is the  $i$ -th observation in group  $j$ ,  $u_j$  is the random intercept for group  $j$ , and  $\sigma_u^2$  is the variance component controlling between-group variability. The parameters of interest are:

- The random intercepts  $\mathbf{u} = (u_1, \dots, u_J)$  (location parameters).
- The variance component  $\sigma_u^2$  (scale parameter, hyper-parameter).

The true posterior  $p(\mathbf{u}, \sigma_u^2 \mid \mathbf{y})$  exhibits strong dependence between  $\mathbf{u}$  and  $\sigma_u^2$ :

- If  $\sigma_u^2$  is large, the data support larger deviations  $|u_j|$  from zero.
- If  $\sigma_u^2$  is small, the posterior for each  $u_j$  is pulled tightly toward zero (shrinkage).

This correlation is encoded in the joint posterior. MCMC correctly captures this dependence by sampling  $(\mathbf{u}, \sigma_u^2)$  jointly.

### Mean-field factorisation breaks the dependence

Mean-field VI imposes the factorisation

$$q(\mathbf{u}, \sigma_u^2) = q(\mathbf{u}) q(\sigma_u^2),$$

forcing the variational distributions for  $\mathbf{u}$  and  $\sigma_u^2$  to be independent. During coordinate-ascent optimisation:

1. The algorithm updates  $q(\mathbf{u})$  given the current  $q(\sigma_u^2)$ , averaging over the variational uncertainty in  $\sigma_u^2$ .
2. Then it updates  $q(\sigma_u^2)$  given the current  $q(\mathbf{u})$ , averaging over the variational distribution of  $\mathbf{u}$ .

Because  $q(\mathbf{u})$  and  $q(\sigma_u^2)$  cannot coordinate, the algorithm systematically underestimates the posterior variance of  $\sigma_u^2$ . Intuitively:

- $q(\mathbf{u})$  sees an averaged (smoothed) version of  $\sigma_u^2$ , not the full range of plausible values.
- $q(\sigma_u^2)$  must explain the observed variability in  $\mathbf{u}$  using only the mean of  $q(\mathbf{u})$ , losing information about the uncertainty in  $\mathbf{u}$ .
- The factorisation prevents the two distributions from "communicating" about their joint uncertainty.

The result:  $q(\sigma_u^2)$  is *under-dispersed*—its posterior variance is too small, and the mass is pulled towards smaller values compared to a well-calibrated MCMC analysis [?, ?].

### Consequences: over-shrinkage and overconfidence

In practice, this under-dispersion at the variance level appears as *over-shrinkage*: the random effects are pulled too tightly towards the group mean, so the fitted model looks more confident and more homogeneous than it should. Specifically:

- Because  $q(\sigma_u^2)$  underestimates the true variability, the algorithm "believes" that groups are more similar than they actually are.
- This leads to excessive shrinkage of individual random effects  $u_j$  toward zero.
- The model appears to have high precision (narrow posterior intervals), but this precision is spurious—it reflects the mean-field constraint, not the information in the data.

This effect can be subtle in the marginal posteriors for individual random effects  $u_j$  (where posterior means may still be reasonable), but it becomes very visible when we compare posterior distributions for  $\sigma_u^2$  itself. A typical plot shows:

- **MCMC posterior for  $\sigma_u^2$ :** Wider distribution, capturing genuine uncertainty about between-group variability.
- **Mean-field VB posterior for  $\sigma_u^2$ :** Narrow spike, overconfident estimate that underestimates the true range of plausible values.

### Why variance components are especially affected

Variance components are *hyper-parameters*—parameters that appear in the priors of other parameters. In the random-intercept model,  $\sigma_u^2$  governs the distribution of  $u_j$ :

$$u_j \sim N(0, \sigma_u^2).$$

This hierarchical structure creates strong posterior dependence between  $\sigma_u^2$  and  $\mathbf{u}$ , which mean-field factorisations cannot capture. As a result:

- Under-dispersion is most severe for variance components (and other hyper-parameters).

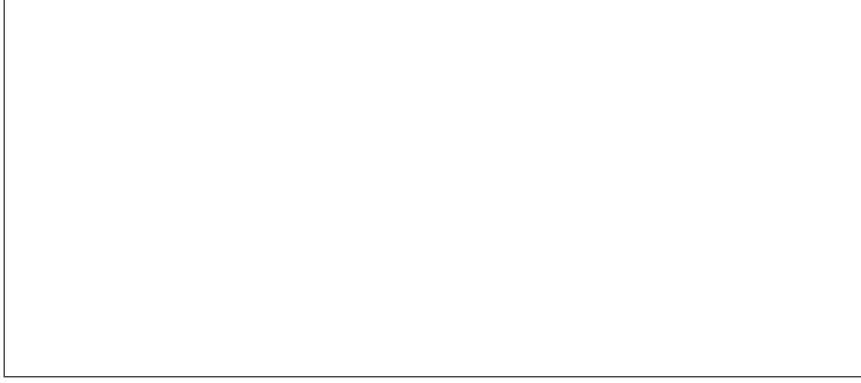


Figure 2: example Venn Diagram

- Location parameters like  $u_j$  or  $\beta$  are less affected, though they still suffer from indirect bias due to the misestimated  $\sigma_u^2$ .

This is the flagship demonstration of mean-field VI's limitations in hierarchical settings, and motivates the use of more flexible variational families or MCMC when accurate variance component estimation is critical.

### Alternative approaches

To mitigate under-dispersion in hierarchical models, one can:

- Use a less restrictive variational family (e.g., full-form VI with correlated  $q(\mathbf{u}, \sigma_u^2)$ ) that preserves the dependence between random effects and variance components.
- Apply structured mean-field approximations that group strongly correlated parameters into joint blocks.
- Fall back to MCMC for hierarchical models where accurate variance component estimation is critical.
- Use variance ratio diagnostics (VR) to quantify the severity of under-dispersion and decide whether the VI approximation is acceptable.

The choice depends on the trade-off between computational speed and accuracy: mean-field VI is fast but systematically biased for variance components; MCMC is slow but asymptotically exact; full-form VI offers a middle ground with better accuracy at moderate computational cost.

# APPENDIX

## A Bayesian vs Frequentist Perspectives

This appendix clarifies key conceptual differences between Bayesian and frequentist approaches to statistical inference, particularly as they relate to variational methods.

### A.1 Posterior Distributions

**Bayesian perspective:** Parameters  $\theta$  are treated as random variables. The posterior distribution  $p(\theta \mid \text{data})$  represents the probability distribution of  $\theta$  after observing the data. This distribution quantifies uncertainty by specifying how likely different parameter values are given the observed data.

**Frequentist perspective:** Parameters  $\theta$  are fixed but unknown constants. There is no probability distribution over parameters. A frequentist estimates  $\theta$  (for example,  $\hat{\theta}$  via maximum likelihood) and may construct confidence intervals, but does not assign probabilities to parameter values. A 95% confidence interval does not mean “ $\theta$  has 95% probability of lying in this interval”; rather, it means “if we repeated this procedure many times, 95% of such intervals would contain the true  $\theta$ .”

### A.2 Prior Distributions

**Bayesian perspective:** Before observing data, prior beliefs or knowledge about  $\theta$  are encoded in a prior distribution  $p(\theta)$ . This prior may be informative (reflecting strong prior knowledge) or vague (reflecting minimal assumptions). Bayes’ theorem updates the prior with data via

$$p(\theta \mid \text{data}) \propto p(\text{data} \mid \theta) p(\theta).$$

**Frequentist perspective:** No prior distributions are used. Inference is based solely on the likelihood  $p(\text{data} \mid \theta)$  and the data. While frequentist methods may employ regularisation (which has similar mathematical effects to Bayesian priors), this is not framed as encoding beliefs about parameters before observing data.

### A.3 Latent Variables and Parameters

Both paradigms use latent (unobserved) variables  $Z$  that influence the data. Examples include cluster memberships in mixture models or true abilities in item-response models.

**Frequentist approach:** Latent variables  $Z$  are random, but parameters  $\theta$  remain fixed. The likelihood is obtained by integrating over  $Z$ :

$$L(\theta) = \int p(\text{data}, Z \mid \theta) dZ.$$

One then estimates  $\hat{\theta} = \arg \max_{\theta} L(\theta)$ , treating  $\theta$  as a fixed unknown.

**Bayesian approach:** Both  $Z$  and  $\theta$  are random. A joint posterior distribution  $p(\theta, Z \mid \text{data})$  is obtained, quantifying uncertainty about both latent variables and parameters simultaneously.

### A.4 Approximating Posteriors

**Bayesian goal:** The posterior  $p(\theta \mid \text{data})$  is the inferential target, but it is often intractable. Approximation methods include:

- Markov chain Monte Carlo (MCMC), which generates samples from  $p(\theta \mid \text{data})$ ;
- Variational inference, which finds a tractable distribution  $q(\theta)$  that approximates  $p(\theta \mid \text{data})$  [?, ?].

**Frequentist goal:** The objective is to estimate a point value  $\hat{\theta}$  and quantify sampling uncertainty (how  $\hat{\theta}$  would vary across hypothetical repeated samples). There is no notion of “approximating a distribution over  $\theta$ ,” because parameters are not viewed as random.

## A.5 Random vs Fixed Parameters

**Frequentist view:** Parameters  $\theta$  are fixed but unknown constants. Randomness arises only from sampling: different data sets yield different estimates  $\hat{\theta}$ , but the true  $\theta$  does not vary. Standard errors describe how  $\hat{\theta}$  varies across hypothetical repeated samples.

**Bayesian view:** Parameters  $\theta$  are random variables with probability distributions. Uncertainty about  $\theta$  is represented by a prior distribution before data and a posterior distribution after data. The posterior distribution directly encodes degrees of belief about different parameter values.

## A.6 Hierarchical Variance Components

Both frameworks use hierarchical (multilevel) models in which observations are grouped (for example, students within schools). Group-level effects often vary randomly.

**Frequentist approach (mixed models):**

- Random effects  $u_j \sim N(0, \sigma^2)$  for group  $j$  are treated as random.
- The variance component  $\sigma^2$  is a *fixed parameter* to be estimated (for example, via restricted maximum likelihood).
- Inference yields a point estimate  $\hat{\sigma}^2$  and possibly a confidence interval.

**Bayesian approach:**

- Random effects  $u_j \sim N(0, \sigma^2)$  for group  $j$ .
- The variance component  $\sigma^2$  also receives a prior distribution  $p(\sigma^2)$ .
- Inference yields a posterior distribution  $p(\sigma^2 \mid \text{data})$ , providing full uncertainty quantification.

In the context of mean-field variational Bayes, the approximate posterior  $q(\sigma^2)$  for the variance component is often under-dispersed (too narrow) relative to the true posterior obtained via MCMC [?, ?]. This under-dispersion manifests as overconfidence in the estimated variance and can lead to excessive shrinkage of random effects towards their mean. A frequentist analysis would yield a different point estimate but would not frame the issue as “posterior distribution is too narrow,” because frequentist inference does not produce posterior distributions in the first place.

## A.7 Philosophical Underpinnings

**Frequentist:** Probability describes long-run frequency in repeated experiments. Parameters are not assigned probabilities; only data (and functions of data) have probability distributions under repeated sampling.

**Bayesian:** Probability quantifies degree of belief about uncertain quantities, including parameters. Both data and parameters can have probability distributions, and inference updates beliefs via Bayes’ theorem.

These foundational differences explain why concepts central to variational inference—such as approximating posterior distributions and quantifying parameter uncertainty through probability distributions—are specific to the Bayesian paradigm and may appear unfamiliar from a frequentist perspective.