

Importance Weighting and Variational Inference

Justin Domke and Daniel Sheldon, University of Massachusetts, Amherst

- Variational autoencoders can use **importance weighting** for **better likelihood bounds**.

- But how to apply to “**pure probabilistic**” variational inference (VI)?

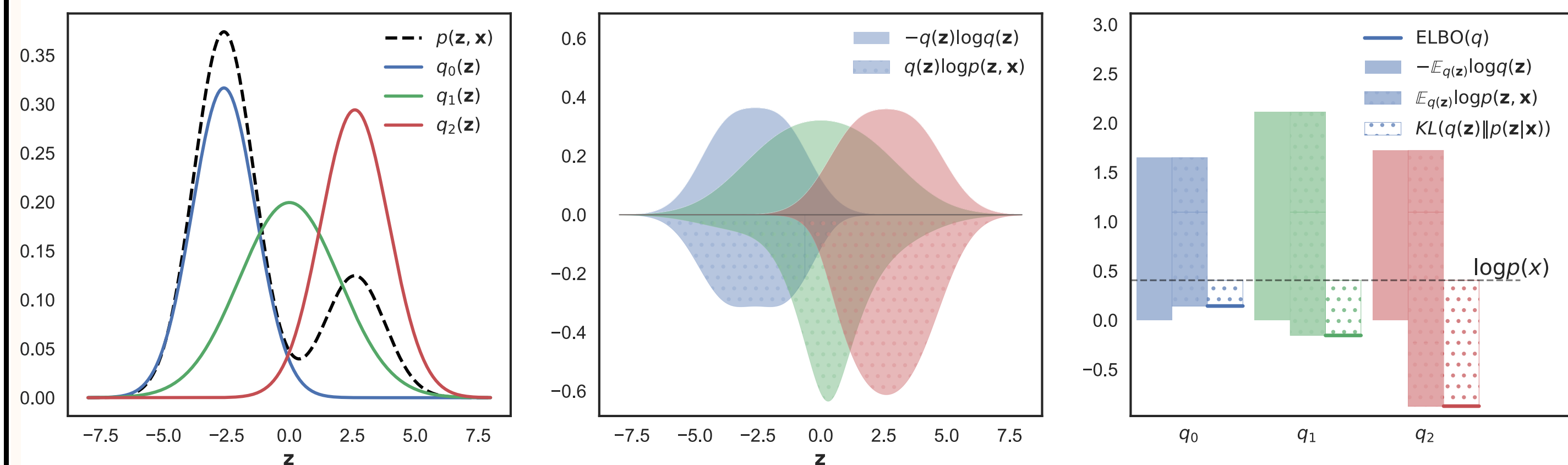
- We show that using importance-weighting is equivalent to **traditional VI** on **augmented distributions**. This informs test-time inference and clarifies looseness of existing bounds.

- Investigate VI on **elliptical distributions** via an “inverse CDF trick”.

1. The ELBO Decomposition

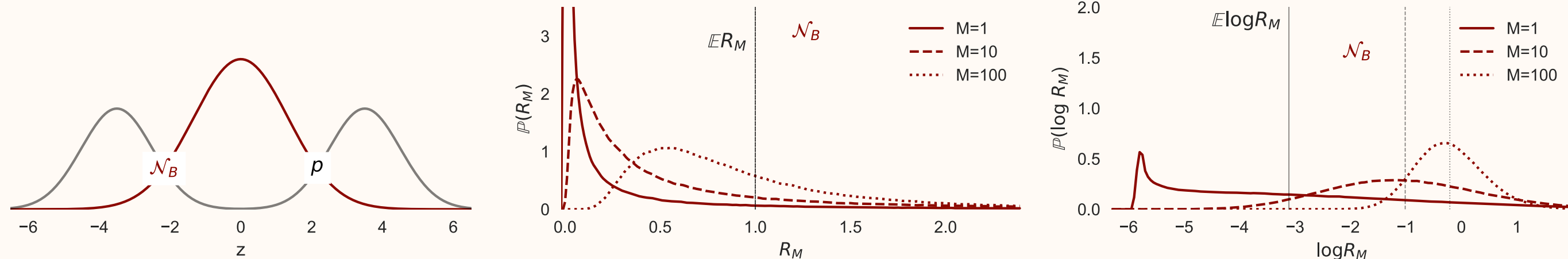
$$\log p(\mathbf{x}) = \underbrace{\mathbb{E}_{\mathbf{z} \sim q} \log \frac{p(\mathbf{z}, \mathbf{x})}{q(\mathbf{z})}}_{\text{ELBO}(q||p)} + KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x}))$$

- Any q gives $\text{ELBO} \leq \log p(\mathbf{x})$.
- Looseness is KL-divergence.



2. Importance Weighting

- For any $R > 0$ with $\mathbb{E} R = p(\mathbf{x})$: $\log p(\mathbf{x}) = \underbrace{\mathbb{E} \log R}_{\text{bound}} + \underbrace{\mathbb{E} \log \frac{p(\mathbf{x})}{R}}_{\text{looseness}}$.
- Traditional VI: $R = \frac{p(\mathbf{z}, \mathbf{x})}{q(\mathbf{z})}$, $\mathbf{z} \sim q$.
- Better bound: Average i.i.d. samples: $R_M = \frac{1}{M} \sum_{m=1}^M \frac{p(\mathbf{z}_m, \mathbf{x})}{q(\mathbf{z}_m)}$, $\mathbf{z}_m \sim q$



IWAEs: [Burda et al., 2015]

- $p(\mathbf{z}, \mathbf{x}) = \text{model}(\mathbf{z}, \mathbf{x})$
- Input \mathbf{x} (dataset)
- Maximize $\mathbb{E} \log R_M$ w.r.t. p and q
- Use p

Importance Weighted VI (IWVI):

- $p(\mathbf{z}, \mathbf{x}) = (\text{Some model})$
- Input \mathbf{x} (evidence)
- Maximize $\mathbb{E} \log R_M$ w.r.t. q
- Use q

Good-old-fashioned VI:

$$\log p(\mathbf{x}) = \text{ELBO}(q(\mathbf{z})||p(\mathbf{z}, \mathbf{x})) + KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x}))$$

- Learning:** $\text{ELBO} \leq \log p(\mathbf{x})$
- Inference:** $\mathbb{E}_{p(\mathbf{z}|\mathbf{x})} t(\mathbf{z}) \approx \mathbb{E}_{q(\mathbf{z})} t(\mathbf{z})$

IWVI:

$$\log p(\mathbf{x}) = \mathbb{E} \log R_M + \mathbb{E} \log \frac{p(\mathbf{x})}{R_M}$$

- Learning:** $\mathbb{E} \log R_M \leq \log p(\mathbf{x})$
- Inference:** ???

3. Main Technical Results

Summary:

- Theorem 1: For augmented p_M / q_M , IWVI minimizes $KL(q_M(\mathbf{z}_{1:M})||p_M(\mathbf{z}_{1:M}|\mathbf{x}))$.
- Theorem 2: That is exactly $\underbrace{KL(q_M(\mathbf{z}_1)||p(\mathbf{z}_1|\mathbf{x}))}_{\text{what we care about}} + \underbrace{KL(q_M(\mathbf{z}_{2:M})||q(\mathbf{z}_{2:M}))}_{\text{other stuff}}$.
- Theorem 3: When M is large that is approximately $\frac{1}{M} \frac{\mathbb{V}[R]}{2p(\mathbf{x})}$.

3.1 Theorem 1: IVWI is Normal VI on Augmented Distributions

Definition of $q_M(\mathbf{z}_{1:M})$:

- Draw $\hat{\mathbf{z}}_1, \hat{\mathbf{z}}_2, \dots, \hat{\mathbf{z}}_M$ independently from $q(\mathbf{z})$.
- Choose $m \in \{1, \dots, M\}$ with prob $\mathbb{P}(m) \propto \frac{p(\hat{\mathbf{z}}_m, \mathbf{x})}{q(\hat{\mathbf{z}}_m)}$.
- Set $\mathbf{z}_1 = \hat{\mathbf{z}}_m$ and $\mathbf{z}_{2:M} = \hat{\mathbf{z}}_{-m}$

(Self-normalized importance sampling for $\hat{\mathbf{z}}_m$; also keep and relabel unselected $\hat{\mathbf{z}}_i$)

Definition of $p_M(\mathbf{z}_{1:M})$:

(One sample from p and $M-1$ “dummy” samples from q)

- $p_M(\mathbf{z}_{1:M}, \mathbf{x}) = p(\mathbf{z}_1, \mathbf{x})q(\mathbf{z}_2) \dots q(\mathbf{z}_M)$.

Previously known [Bachman and Precup, 2015, Cremer et al., 2017, Naesseth et al., 2018, Le et al., 2018]: $\log p(\mathbf{x}) \geq \text{ELBO}(q_M(\mathbf{z}_1)||p_M(\mathbf{z}_1, \mathbf{x})) \geq \mathbb{E} \log R_M$.

Our Result: $\log p(\mathbf{x}) = \mathbb{E} \log R_M + KL(q_M(\mathbf{z}_{1:M})||p_M(\mathbf{z}_{1:M}|\mathbf{x}))$.

Thus, approximate test integrals as

$$\mathbb{E}_{p(\mathbf{z}|\mathbf{x})} t(\mathbf{z}) = \mathbb{E}_{p_M(\mathbf{z}_1|\mathbf{x})} t(\mathbf{z}_1) \approx \mathbb{E}_{q_M(\mathbf{z}_1)} t(\mathbf{z}_1).$$

3.2 Theorem 2: IVWI is Tightening an Upper Bound

Result:

$$\underbrace{KL(q_M(\mathbf{z}_{1:M})||p_M(\mathbf{z}_{1:M}|\mathbf{x}))}_{\text{what IWVI minimizes}} = \underbrace{KL(q_M(\mathbf{z}_1)||p(\mathbf{z}_1|\mathbf{x}))}_{\text{what we care about}} + \underbrace{KL(q_M(\mathbf{z}_{2:M})||q(\mathbf{z}_{2:M}))}_{\text{other stuff}}$$

Proof: KL chain rule + definition of p_M .

If you will use normalized importance sampling, IWVI truly optimizes a bound.

3.3 Theorem 3: Asymptotics

Result: If $\mathbb{E} |R - p(\mathbf{x})|^{2+\alpha} < \infty$ for $\alpha > 0$ and $\limsup_{M \rightarrow \infty} \mathbb{E} \frac{1}{R_M} < \infty$,

$$\lim_{M \rightarrow \infty} M (\log p(\mathbf{x}) - \mathbb{E} \log R_M) = \frac{\mathbb{V}[R_1]}{2p(\mathbf{x})}.$$

(Maddison et al. [2017] showed for $\alpha = 4$)

Non-Proof: CLT + 2nd-order delta-method:

$$M (\log p(\mathbf{x}) - \mathbb{E} \log R_M) \xrightarrow{d} \frac{\mathbb{V}[R_1]}{2p(\mathbf{x})} \chi_1^2$$

Problem: $X_M \xrightarrow{d} X$ does not imply $\mathbb{E}[X_M] \rightarrow \mathbb{E}[X]$.

Proof: Long. Broadly follows Maddison et al. [2017] to bound higher terms in a Taylor expansion. Biggest technical innovation is using Marcinkiewicz–Zygmund inequality to bound sample moments from true moments.

4. Elliptical Distributions

- For $M = 1$, IWVI minimizes KL .
- For M large, IWVI minimizes $\mathbb{V}[R]$ (Equiv. χ^2 divergence).
- This is **mode finding**.
- This is **mode spanning**.

Suggests we want **different tail behavior** as M changes.

Given some spherically symmetric distribution g_ν , an “Elliptical” distribution is

$$q(\mathbf{z}|\mu, \Sigma, \nu) = \frac{1}{|\Sigma|^{1/2}} g_\nu \left((\mathbf{z} - \mu)^\top \Sigma^{-1} (\mathbf{z} - \mu) \right).$$

- Fit μ, Σ as “Normal”. (Reparameterization trick with $g_\nu(\epsilon)$ as base density).
- Fit ν by backpropagating through inverse CDF of $\|\epsilon\|$, $\epsilon \sim g_\nu$.
- No inverse-CDF? Sample $\epsilon \sim g_\nu$, then “pretend”: (Same idea at this conference: *Implicit Reparameterization Gradients*, Figurnov et al.)

$$\nabla_\nu F_\nu^{-1}(u) = -\frac{\nabla_\nu F_\nu(\|\epsilon\|)}{\nabla_\nu F_\nu(\|\epsilon\|)}, \quad u = F_\nu(\|\epsilon\|).$$

5. Experiments

Variational Families:

- IWVI : Gaussians
- E-IWVI: Student-T with ν deg. of freedom.

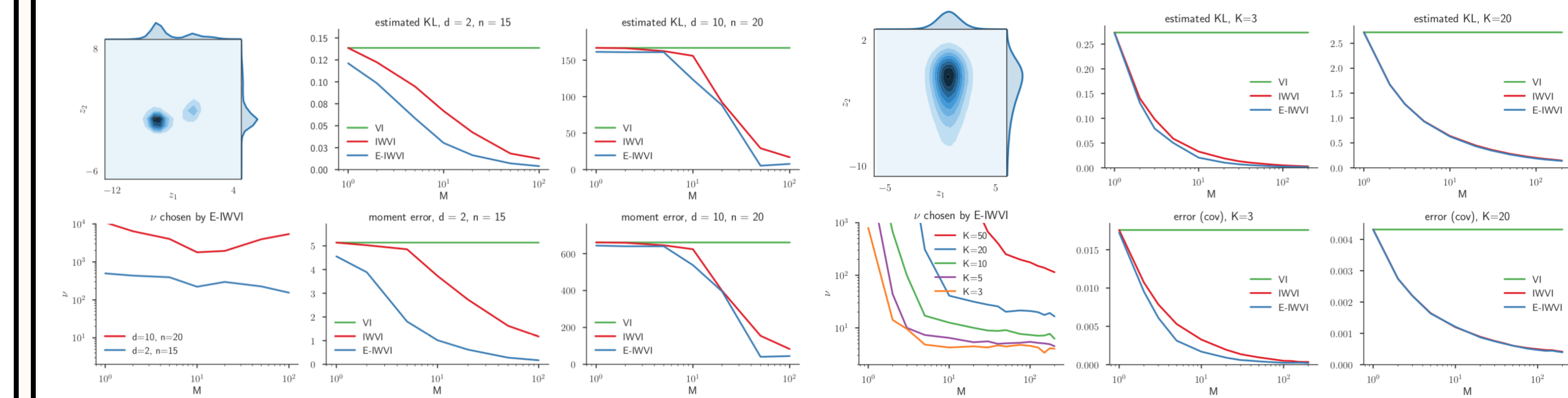
Error metrics:

- $KL(q(\mathbf{z})||p(\mathbf{z}))$
- $\mathbb{C}[p(\mathbf{z})]$ vs. $\mathbb{C}[q_M(\mathbf{z}_1)]$

Clutter model: [Minka, 2001]

Random Dirichlets:

- $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}, 0, 100I)$ - hidden location
- Sample $\alpha_1, \dots, \alpha_K \sim \text{Gamma}(10)$
- $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ noisy observations
- $p(\mathbf{z})$ is density of $\text{StickBreak}(\theta)$, $\theta \sim \text{Dirichlet}(\theta|\alpha)$
- $-p(\mathbf{x}_i|\mathbf{z}) = \frac{1}{4}\mathcal{N}(\mathbf{x}_i|\mathbf{z}_i, I) + \frac{3}{4}\mathcal{N}(\mathbf{x}_i|0, 10I)$



Logistic Regression (Cauchy Prior)

