

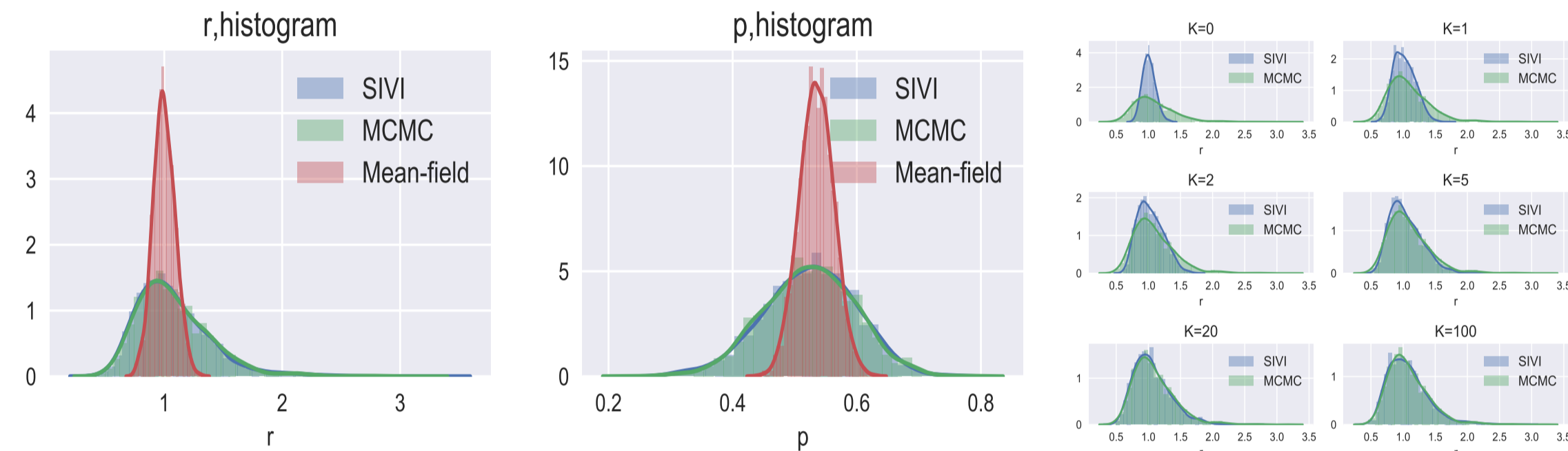
Variational inference

Find $Q(z) \in \mathcal{Q}$ to maximize evidence lower bound (ELBO)

$$\text{ELBO} = \mathcal{L}(Q) = \mathbf{E}_Q[\ln P(\mathbf{x}, \mathbf{z})] - \mathbf{E}_Q[\ln Q(\mathbf{z})] \\ = \ln p(\mathbf{x}) - \text{KL}(Q(z) || P(z|\mathbf{x}))$$

- Optimizing $Q(z)$ is considered as approximating posterior inference;
- Mean-field VI makes a fully factorized assumption as $Q(z) = \prod_{i=1}^K q_{\phi_i}(z_i)$;
- Capturing latent dependencies is crucial for correct uncertainty estimation.

Negative binomial $x_i \stackrel{iid}{\sim} \text{NB}(r, p)$, $r \sim \text{Gamma}(a, 1/b)$, $p \sim \text{Beta}(\alpha, \beta)$



Semi-implicit model

Implicit model consists of a source of randomness $q(\epsilon)$ and a deterministic transform $T_\phi: \mathbb{R}^p \rightarrow \mathbb{R}^d$

$$z = T_\phi(\epsilon), \quad \epsilon \sim q(\epsilon)$$

For an implicit distribution, it is often easy to generate random samples from it but intractable to calculate its probability density function, making variational inference difficult

$$q_\phi(z) = \frac{\partial}{\partial z_1} \dots \frac{\partial}{\partial z_d} \int_{T_\phi(\epsilon) \leq z} q(\epsilon) d\epsilon$$

Semi-implicit model is a two-stage model

$$z \sim q(z|\psi), \quad \psi \sim q_\phi(\psi)$$

- The first layer distribution $q(z|\psi)$ is explicit, while the mixing distribution $q_\phi(\psi)$ is allowed to be implicit;
- The marginal distribution $h_\phi(z)$ is used as variational distribution

$$\mathcal{H} = \left\{ h_\phi(z) : h_\phi(z) = \mathbf{E}_{\psi \sim q_\phi(\psi)} [q(z|\psi)] = \int_{\psi} \left[\prod_{k=1}^K q(z_k|\psi_k) \right] q_\phi(\psi) d\psi \right\}$$

- The components of z are conditionally independent but marginally dependent;
- It is evident that $q(z|\psi) \in \mathcal{Q} \subseteq \mathcal{H}$, i.e., \mathcal{H} forms an expansion;
- Semi-implicit distribution $h_\phi(z)$ achieves a balance between expressiveness and tractability.

Lower and upper bound of ELBO

Optimize $\text{ELBO} = \mathbf{E}_{h_\phi(z)}[\ln p(\mathbf{x}, \mathbf{z}) - \ln h_\phi(z)]$ for SIVI is generally intractable if $h_\phi(z) = \mathbf{E}_{q_\phi(\psi)} q(z|\psi)$ is not analytic

- KL convexity and Jensen's inequality lead to an ELBO lower bound:

$$\mathcal{L}(q(z|\psi), q_\phi(\psi)) = \mathbf{E}_{\psi \sim q_\phi(\psi)} \mathbf{E}_{z \sim q(z|\psi)} \log \frac{p(\mathbf{x}, \mathbf{z})}{q(z|\psi)} \\ = -\mathbf{E}_{\psi \sim q_\phi(\psi)} \text{KL}(q(z|\psi) || p(z|\mathbf{x})) + \log p(\mathbf{x}) \\ \leq -\text{KL}(\mathbf{E}_{\psi \sim q_\phi(\psi)} q(z|\psi) || p(z|\mathbf{x})) + \log p(\mathbf{x}) = \mathcal{L} = \mathbf{E}_{z \sim h_\phi(z)} \log \frac{p(\mathbf{x}, \mathbf{z})}{h_\phi(z)}$$

- Using the concavity of the logarithmic function, we have $\log h_\phi(z) \geq \mathbf{E}_{\psi \sim q_\phi(\psi)} \log q(z|\psi)$ and hence an ELBO upper bound:

$$\bar{\mathcal{L}}(q(z|\psi), q_\phi(\psi)) = \mathbf{E}_{\psi \sim q_\phi(\psi)} \mathbf{E}_{z \sim h_\phi(z)} \log \frac{p(\mathbf{x}, \mathbf{z})}{q(z|\psi)} \geq \mathcal{L}$$

- Note there is a subtle but critical difference between \mathcal{L} and $\bar{\mathcal{L}}$
- Direct optimizing on \mathcal{L} will result in degeneracy; namely, $q_\phi(\psi) \rightarrow \delta_{\psi^*}$

$$\mathcal{L}(q(z|\psi), q_\phi(\psi)) \leq -\mathbf{E}_{z \sim q(z|\psi^*)} \log \frac{q(z|\psi^*)}{p(\mathbf{x}, \mathbf{z})},$$

with $\psi^* = \text{argmax}_\psi -\text{KL}(q(z|\psi) || p(\mathbf{x}, \mathbf{z}))$

Asymptotically exact surrogate ELBOs

Add regularization as $\mathcal{L}_K = \mathcal{L} + B_K$

$$B_K = \mathbf{E}_{\psi, \psi^{(1)}, \dots, \psi^{(K)} \sim q_\phi(\psi)} \text{KL}(q(z|\psi) || \tilde{h}_K(z))$$

The regularized surrogate ELBO can be expressed as

$$\mathcal{L}_K = \mathbf{E}_{\psi \sim q_\phi(\psi)} \mathbf{E}_{z \sim q(z|\psi)} \mathbf{E}_{\psi^{(1)}, \dots, \psi^{(K)} \sim q_\phi(\psi)} \log \frac{p(\mathbf{x}, \mathbf{z})}{\frac{1}{K+1} [q(z|\psi) + \sum_{k=1}^K q(z|\psi^{(k)})]}$$

The Jensen gap can also be narrowed from upper side by $\bar{\mathcal{L}}_k = \bar{\mathcal{L}} - A_k$

$$\bar{\mathcal{L}}_K = \mathbf{E}_{\psi \sim q_\phi(\psi)} \mathbf{E}_{z \sim q(z|\psi)} \mathbf{E}_{\psi^{(1)}, \dots, \psi^{(K)} \sim q_\phi(\psi)} \log \frac{p(\mathbf{x}, \mathbf{z})}{\frac{1}{K} \sum_{k=1}^K q(z|\psi^{(k)})}$$

Property: Surrogate ELBOs

The regularized lower bound \mathcal{L}_K is an asymptotically exact ELBO that satisfies $\mathcal{L}_0 = \mathcal{L}$ and $\lim_{K \rightarrow \infty} \mathcal{L}_K = \mathcal{L}$. The regularized upper bound satisfies $\bar{\mathcal{L}}_1 = \bar{\mathcal{L}}$, $\bar{\mathcal{L}}_{K+1} \leq \bar{\mathcal{L}}_K$, and $\lim_{K \rightarrow \infty} \bar{\mathcal{L}}_K = \mathcal{L}$.

For non-reparameterizable but conjugate model, the gradient can be expressed as

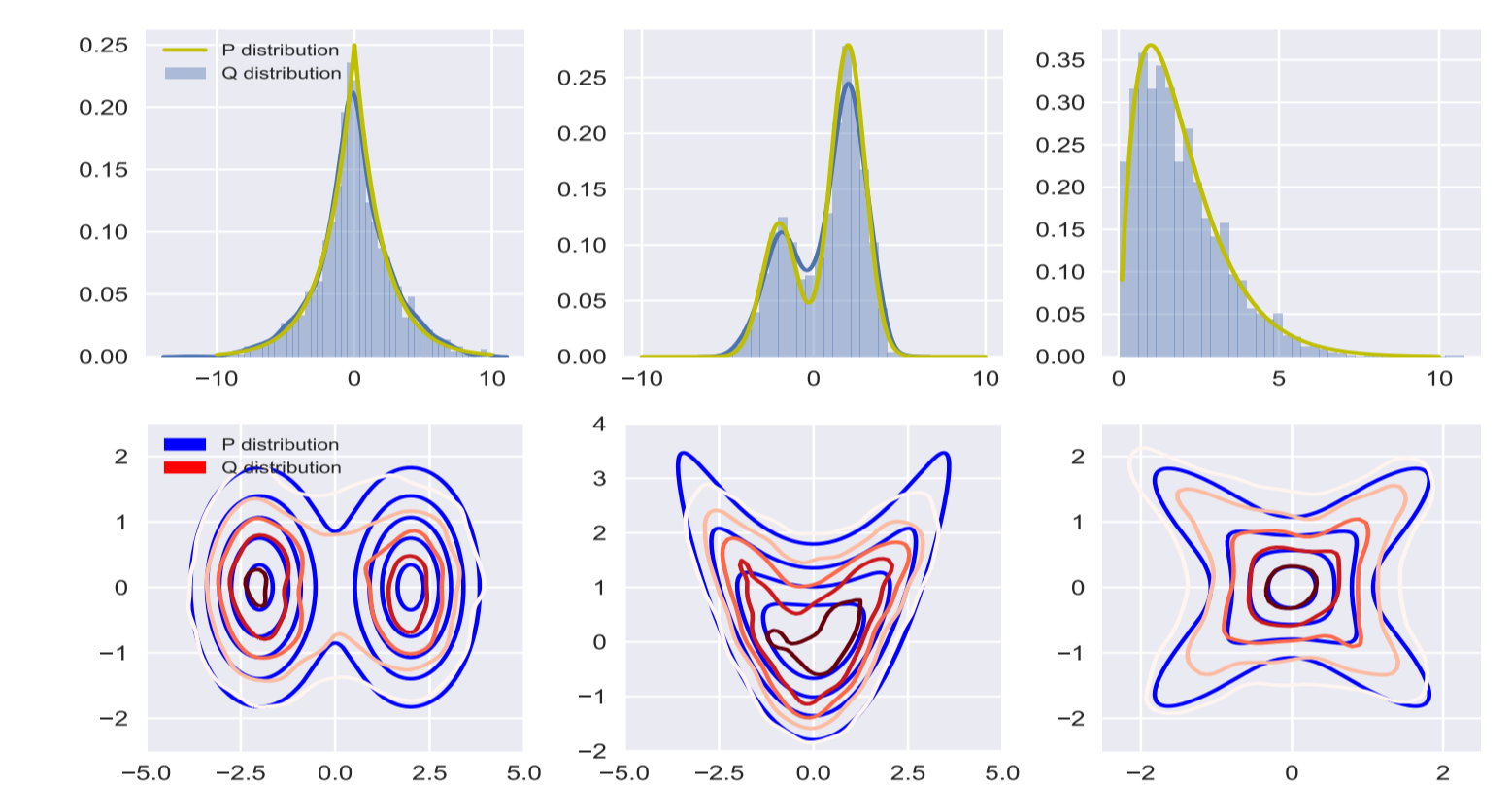
$$\nabla_\phi \mathcal{L}_K \approx \frac{1}{J} \sum_{j=1}^J \left\{ -\nabla_\phi \mathbf{E}_{z \sim q_\xi(z|T_\phi(\epsilon_j))} [\log \frac{q_\xi(z|T_\phi(\epsilon_j))}{p(\mathbf{x}, \mathbf{z})}] \right. \\ \left. + \nabla_\phi \log r_{\xi, \phi}(z_j, \epsilon_j, \epsilon^{(1:K)}) \right. \\ \left. + [\nabla_\phi \log q_\xi(z_j | T_\phi(\epsilon_j))] \log r_{\xi, \phi}(z_j, \epsilon_j, \epsilon^{(1:K)}) \right\},$$

$$r_{\xi, \phi}(z, \epsilon, \epsilon^{(1:K)}) = q_\xi(z | T_\phi(\epsilon)) / \left[\frac{q_\xi(z | T_\phi(\epsilon)) + \sum_{k=1}^K q_\xi(z | T_\phi(\epsilon^{(k)}))}{K+1} \right]$$

Experiments

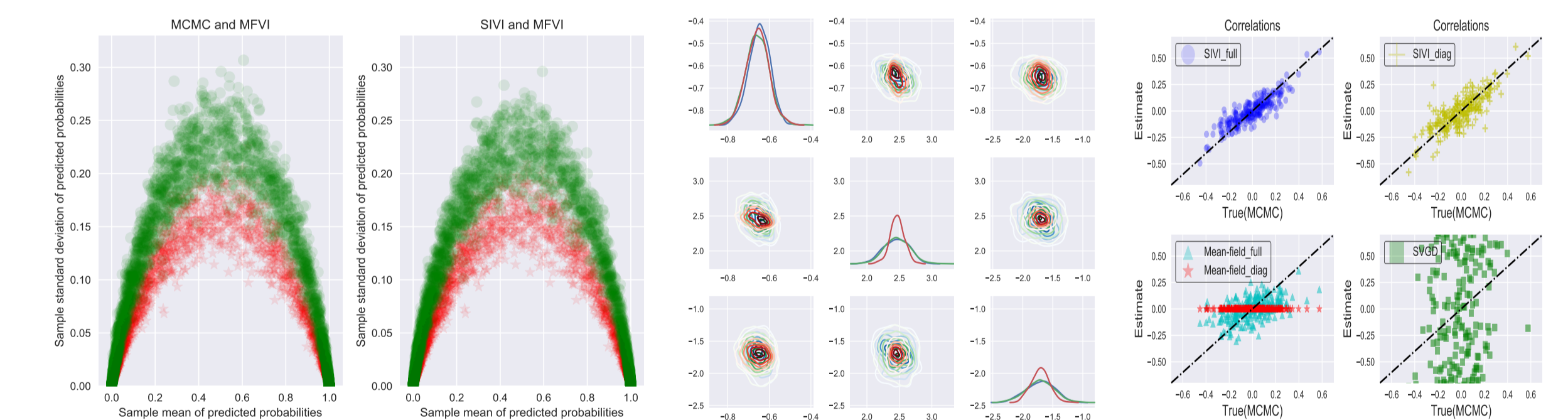
- Toy examples (capturing skewness, kurtosis, and multimodality)

$$h(z) = \mathbf{E}_{\psi \sim q(\psi)} q(z|\psi), \quad q(z|\psi) = (\log) \text{Normal}(\psi, 0.1)$$



- Bayesian Logistic regression

$$y_i \sim \text{Bernoulli}[(1 + e^{-x_i^T \beta})^{-1}], \quad \beta \sim \mathcal{N}(\mathbf{0}, \alpha^{-1} \mathbf{I}_{V+1}) \\ q(\beta | \psi) = \mathcal{N}(\psi, \Sigma), \quad \psi \sim q_\phi(\psi)$$



- Variational autoencoder (VAE)

Inject random noise at M different stochastic layers. Let $h(z|\mathbf{x}) = \int q(z|\mathbf{x}, \epsilon) q(\epsilon) d\epsilon$

$$\ell_t = T_t(\ell_{t-1}, \epsilon_t, \mathbf{x}; \phi), \quad \epsilon_t \sim q_t(\epsilon), \quad t = 1, \dots, M, \\ \mu(\mathbf{x}, \phi) = f(\ell_M, \mathbf{x}; \phi), \quad \Sigma(\mathbf{x}, \phi) = g(\ell_M, \mathbf{x}; \phi), \\ q_\phi(z|\mathbf{x}, \mu, \Sigma) = \mathcal{N}(\mu(\mathbf{x}, \phi), \Sigma(\mathbf{x}, \phi)),$$

Results below from Mescheder et al. (2017)		
VAE + IAF (Kingma et al., 2016)		$\approx 84.9 \pm 0.3$
Auxiliary VAE (Maaloe et al., 2016)		$\approx 83.8 \pm 0.3$
AVB + AC		$\approx 83.7 \pm 0.3$
SIVI (3 stochastic layers)		$= 84.07$
SIVI (3 stochastic layers) + IW($\bar{K} = 10$)		$= 83.25$

Full version at <https://arxiv.org/abs/1805.11183>

