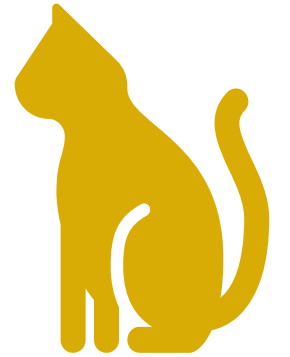# Advances in Approximate Inference

Yingzhen Li, Cheng Zhang

**Microsoft Research Cambridge**

# What is the Number?
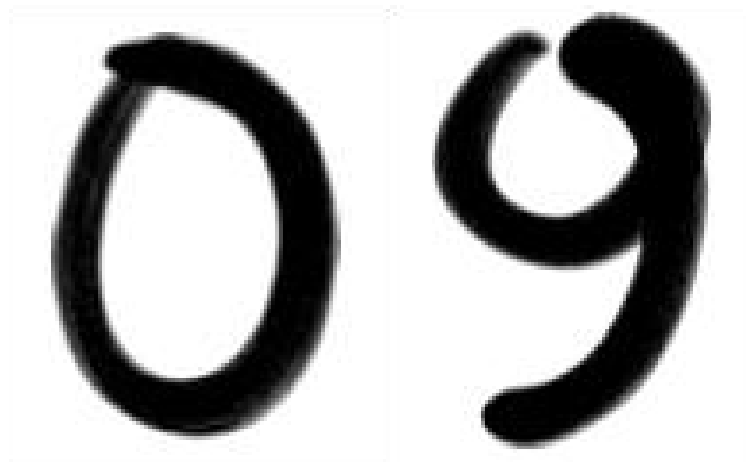
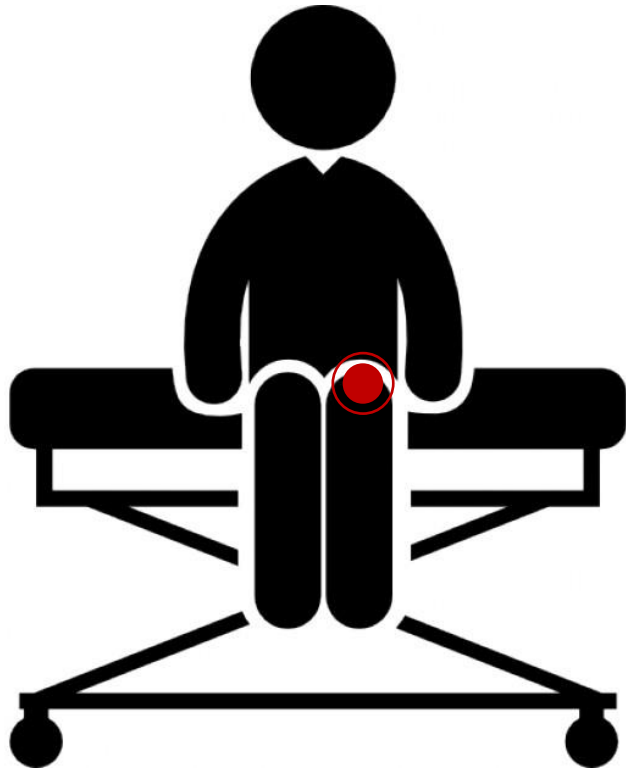07? 09?
67? 69?
…

# What is the Number?

09

# What is the Diagnosis?



Injury?
Osteoarthritis?
Neuropathic pain?
… …

# What is the Diagnosis?



Neuropathic pain
(might have spine injury)

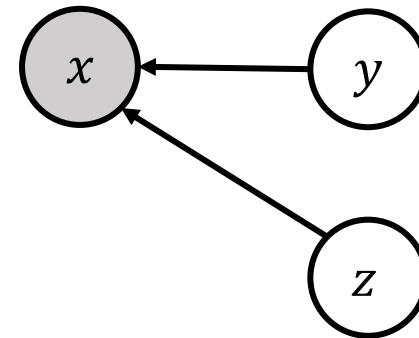# Uncertainty is Important

Bayesian ML / Probability Theory

Decision making under uncertainty

# Graphical Representation



$$p(x, y, z) = p(y)p(z)p(x|y, z)$$

# Graphical Representation



$$p(\boldsymbol{x}, \boldsymbol{y}, z) = p(z) \prod_{n}^{N} p(y_n) P(y_n | y_n, z)$$

# Graphical Representation



$$p(\boldsymbol{x}, \boldsymbol{y}, z) = p(z) \prod_{n}^{N} p(x_n) P(y_n | x_n, z)$$

$$p(\boldsymbol{x}, \boldsymbol{y}, z) = p(z) \prod_{n}^{N} p(y_n) P(x_n | y_n, z)$$

# Discriminative Model vs Generative Model



$$p(\boldsymbol{x}, \boldsymbol{y}, z) = p(z) \prod_{n}^{N} p(x_n) P(y_n | x_n, z)$$

$$p(\boldsymbol{x}, \boldsymbol{y}, z) = p(z) \prod_{n}^{N} p(y_n) P(x_n | y_n, z)$$

# Discriminative Model Example

- Bayesian Logistic Regression

| Name | A-level math score | # parents in STEM | Study STEM? |
|------|--------------------|--------------------|-------------|
| Alice | 89 | 0 | 0 (No) |
| Bob | 95 | 1 | 1 (Yes) |
| Ty | 82 | 1 | 0 (No) |
| Emma | 98 | 2 | 1 (Yes) |
| Anna | 92 | 0 | 0 (No) |
| Mo | 88 | 1 | 0 (No) |
| Li | 95 | 0 | 1 (Yes) |

$X_1$ $X_2$ $Y$

$X$

# Discriminative Model Example:

- Bayesian Logistic Regression

$$p(y = 1) = \frac{1}{1 + e^{w_0 + w_1 x_1 + w_2 x_2}}$$

| Name | A-level math score | # parents in STEM | Study STEM? |
|------|------|------|------|
| Alice | 89 | 0 | 0 (No) |
| Bob | 95 | 1 | 1 (Yes) |
| Ty | 82 | 1 | 0 (No) |
| Emma | 98 | 2 | 1 (Yes) |
| Anna | 92 | 0 | 0 (No) |
| Mo | 88 | 1 | 0 (No) |
| Li | 95 | 0 | 1 (Yes) |

$X_1$   $X_2$   $Y$

$X$

# Discriminative Model Example

- Bayesian Logistic Regression

| Name | A-level math score | # parents in STEM | Study STEM? |
|------|--------------------|-------------------|-------------|
| Alice | 89 | 0 | 0 (No) |
| Bob | 95 | 1 | 1 (Yes) |
| Ty | 82 | 1 | 0 (No) |
| Emma | 98 | 2 | 1 (Yes) |
| Anna | 92 | 0 | 0 (No) |
| Mo | 88 | 1 | 0 (No) |
| Li | 95 | 0 | 1 (Yes) |

$X_1 \qquad X_2 \qquad\qquad Y$

$X$

$$p(y = 1) = \frac{1}{1 + e^{w_0 + w_1 x_1 + w_2 x_2}}$$

$w_k$  $x_n$

$K$

$y_n$

$N$

[ 97, 1 ]

# Discriminative Model Example

Computation Graph



$$p(y = 1) = \frac{1}{1 + e^{w_0 + w_1 x_1 + w_2 x_2}}$$



[ 97, 1]

# Discriminative Model Example

Computation Graph

$$p(y = 1) = \frac{1}{1 + e^{w_0 + w_1 x_1 + w_2 x_2}}$$

$$\boldsymbol{W}^T X$$

$$\boldsymbol{W}^T \Phi(X),$$

[ 97, 1]

# Discriminative Model Example:

Computation Graph

$$p(y = 1) = \frac{1}{1 + \boxed{e^{w_0 + w_1 x_1 + w_2 x_2}}}$$

$\boldsymbol{W}^T X$

$\boldsymbol{W}^T \Phi(X),$

$x_1 \quad x_2$

$y$

$x_1 \quad x_2$

$y$

$K$ $w$ $x$

$y$
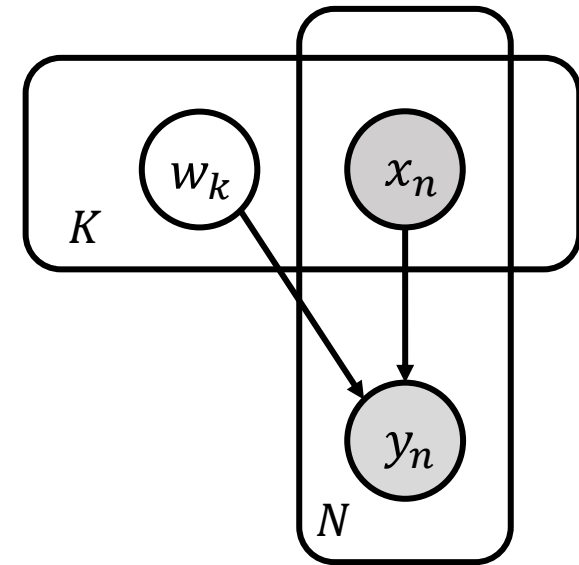
$N$

[ 97, 1]

# Discriminative Model Example

Computation Graph

$$p(y = 1) = \frac{1}{1 + \boxed{e^{w_0 + w_1 x_1 + w_2 x_2}}}$$

$$\boldsymbol{W}^T X$$

$$\boldsymbol{W}^T \Phi(X),$$



Gaussian Processes

Radford Neal's PhD thesis (1994)

[ 97, 1]

# Generative Model Example

- Topic Model

# Generative Model Example:

- Latent Dirichlet Allocation

| ID | topic | neural | distribution | ... |
|----|-------|--------|--------------|-----|
| 1  | 15    | 2      | 19           | ... |
| 2  | 1     | 13     | 21           | ... |
| 3  | 0     | 16     | 1            | ... |

Topic distribution per document:
e.g. 30% "topic model", 40% "natural language processing", 30% "interpretability"



Word distribution per topic:
e.g. under "topic model": "Dirichlet" 2%, "topic" 4%, "Categorical" 1.5%, .....

Blei et al. Latent Dirichlet Allocation. JMLR 2003.

# Generative Model Example



Topics                                          Documents

User embedding                                  Movie rating

Underlying health conditions                    Symptoms

# How to infer the unknowns?

# The Central Computation for Inference

- Inference: infer the <span style="color:red">unknowns</span>
  - Unobserved/latent variables in the model
  - Quantities depending on the latent variables in the model

# The Central Computation for Inference

- Inference: infer the <span style="color:red">unknowns</span>
  - Unobserved/latent variables in the model
  - Quantities depending on the latent variables in the model

$$\int \underbrace{F(\theta)}_{\substack{\text{integrand} \\ \text{function}}} \underbrace{\overbrace{\pi(\theta)}^{\text{prob. density}} d\theta}_{\substack{\text{probability} \\ \text{measure}}}$$

Random variable (unobserved)

(For discrete probability measures, integration becomes discrete sum.)

# Bayesian Inference

$\pi(\theta) = p(\theta|data)$

$$P(\theta \mid data) = \frac{P(\theta)P(data \mid \theta)}{P(data)}$$

- $P(\theta)$: prior distribution
- $P(data \mid \theta)$: likelihood of $\theta$ given $data$
- $P(\theta \mid data)$: posterior distribution of $\theta$ given $data$
- $P(data)$: marginal likelihood/model evidence

$$P(data) = \int P(\theta)P(data \mid \theta)$$

# Computation Challenge

- The central equation for inference:

$$\int F(\theta)\,\pi(\theta)d\theta$$

"What is the prediction distribution of the test output given a test input?"

$F(\theta) = p(y|x,\theta),\ \pi(\theta) = p(\theta\,|\,D),$
$D =$ observed datapoints

# Computation Challenge

- The central equation for inference:

$$\int F(\theta)\,\pi(\theta)d\theta$$

"What is the mean of this distribution?"

$F(\theta) = \theta, \pi(\theta)$ can be complicated and high dimensional

# Computation Challenge

- The central equation for inference:

$$\int F(\theta)\,\pi(\theta)\,d\theta$$

"What is the probability of generating this image?"

$$F(\theta) = \delta(NN(\theta) = x_0),\ \pi(\theta) = N(0, I)$$



$\theta \sim \pi(\theta)$

$x_0$

# Computation Challenge

- The central equation for inference:

$$\int F(\theta)\, \pi(\theta)\, d\theta$$

"What is the weather forecast for tomorrow?"

Answering this in a Bayesian way:
$\theta$: forecasting simulator settings
$D$: historical weather record
$F(\theta) = Simulator(\theta), \pi(\theta) = p(\theta \mid D)$

Nature laughs at the difficulties of integration.

--Pierre-Simon Laplace

Gordon and Sorkin. *The Armchair Science Reader.* New York 1959

# Integration in Bayesian Computation

The probability distribution $\pi(\theta)$ is intractable



$$\int F(\theta) \pi(\theta) d\theta$$

Approximate Inference

(in a strict sense)

This tutorial

# Integration in Bayesian Computation



The integrand $F(\theta)$ is intractable

Implicit Models

Bayesian Optimisation, Probabilistic Numerics

$$\int F(\theta) \pi(\theta) d\theta$$

The probability distribution $\pi(\theta)$ is intractable

Approximate Inference

(in a strict sense)

This tutorial

# Integration in Bayesian Computation



The integrand $F(\theta)$ is intractable

Implicit Models

Bayesian Optimization, Probabilistic Numerics

The probability distribution $\pi(\theta)$ is intractable

Approximate Inference

(in a strict sense)

This tutorial

$$\int F(\theta)\pi(\theta)d\theta$$

Both $F(\theta)$ and $\pi(\theta)$ are intractable

Approximate Bayesian Computation

# Approximate Inference

- Central task: approximate $\pi(\theta)$

$$q(\theta) \approx \pi(\theta)$$

(Assumed $\int F(\theta)q(\theta)d\theta$ can be computed or approximated efficiently.)

# Approximate Inference

- Central task: approximate $\pi(\theta)$

$$q(\theta) \approx \pi(\theta)$$

Approximate distribution design

Explicit distributions        Implicit distributions

# Approximate Inference

- Central task: approximate $\pi(\theta)$

$$q(\theta) \approx \pi(\theta)$$

Approximate distribution design

Algorithm for fitting $q(\theta)$ to $\pi(\theta)$



$$\min \; Loss(q(\theta), \pi(\theta))$$

Explicit distributions

Implicit distributions

Optimisation-based approaches

Sampling-based approaches

# Tutorial Outline

## Basics

Probabilistic modelling

Approximate inference

Variational inference

## Advances

Scalable variational inference

Monte Carlo techniques

Amortized inference

$q$ distribution design

Optimization objective design

## Applications

Bayesian neural networks

Partially observed VAEs

Future challenges

# Bayesian Inference

$$P(\theta \mid D) = \frac{P(\theta)P(D \mid \theta)}{P(D)}$$

- $P(\theta)$: prior
- $P(D \mid \theta)$: likelihood
- $P(\theta \mid D)$: posterior
- $P(D)$: marginal

# Variational Inference (VI)

The posterior

The variational distribution

$$p(\theta|D) = p(D|\theta)p(\theta)/p(D)$$

$$q_\phi(\theta)$$

# Inference as Optimization

$p(\theta|D)$             $q_\phi(\theta)$

**Kullback-Leibler (KL) divergence**

# Kullback-Leibler Divergence

$$KL[q(\theta)||p(\theta)] = -\int q(\theta)\log\frac{p(\theta)}{q(\theta)}\,d\theta = E_{q(\theta)}[\log\frac{p(\theta)}{q(\theta)}]$$

- When $p = q$, KL is 0
- Otherwise, KL > 0
- It measures how similar are these two distributions

# Let's Derive the Objective of VI

- Minimize $KL[q(\theta)||p(\theta|D)]$

$$KL[q(\theta)||p(\theta|D)] = -\mathrm{E}_{\mathrm{q}(\theta)}\left[\log\frac{\mathrm{p}(\theta|D)}{q(\theta)}\right]$$

# Let's Derive the Objective of VI

- Minimize $KL[q(\theta)||p(\theta|D)]$

$$KL[q(\theta)||p(\theta|D)] = -E_{q(\theta)}\left[\log\frac{p(\theta|D)}{q(\theta)}\right]$$

$$= -E_{q(\theta)}\left[\log\frac{p(\theta,D)}{p(D)q(\theta)}\right] = -E_{q(\theta)}\left[\log\frac{p(\theta,D)}{q(\theta)} - \log p(D)\right]$$

# Let's Derive the Objective of VI

- Minimize $KL[q(\theta)||p(\theta|D)]$

$$KL[q(\theta)||p(\theta|D)] = -E_{q(\theta)}\left[\log\frac{\mathrm{p}(\theta|D)}{q(\theta)}\right]$$

$$= -E_{q(\theta)}\left[\log\frac{\mathrm{p}(\theta,D)}{p(D)q(\theta)}\right] = -E_{q(\theta)}\left[\log\frac{\mathrm{p}(\theta,D)}{q(\theta)} - \log p(D)\right]$$

$$= \boxed{\log p(D)} - E_{q(\theta)}\left[\log\frac{\mathrm{p}(\theta,D)}{q(\theta)}\right]$$

Model Evidence

# Let's Derive the Objective of VI

**Minimize** $KL[q(\theta)||p(\theta|D)]$

$$KL[q(\theta)||p(\theta|D)] = \log p(D) - E_{q(\theta)}\left[\log\frac{p(\theta,D)}{q(\theta)}\right]$$

Maximize $E_{q(\theta)}\left[\log\frac{p(\theta,D)}{q(\theta)}\right]$

# Let's Derive the Objective of VI

**Minimize** $KL[q(\theta)||p(\theta|D)]$

$$KL[q(\theta)||p(\theta|D)] = \boxed{\log p(D)} - E_{q(\theta)}\left[\log\frac{p(\theta,D)}{q(\theta)}\right]$$

Model Evidence

Maximize $L = E_{q(\theta)}\left[\log\frac{p(\theta,D)}{q(\theta)}\right]$

Evidence Lower Bound (ELBO)

$\log p(D)$ ————————————————

$KL[q(\theta)||p(\theta|D)]$

ELBO ————————————————

"Model Evidence = ELBO + KL"

# Alternative Derivation

Let's start with the model evidence

$$\log p(D)$$

# Alternative Derivation

$$\log p(D) = \log \int p(\theta, D) \; d\theta$$

# Alternative Derivation

$$\log p(D) = \log \int p(\theta, D) \; d\theta$$

$$= \log \int \frac{\mathrm{p}(\theta, D) \, q(\theta)}{q(\theta)} \, d\theta$$

$$= \log E_{q(\theta)}\left[\frac{p(\theta, D)}{q(\theta)}\right]$$

# Jensen's Inequality



If f is a convex function, then
$$f(E[X]) \leq E[f(X)]$$

If f is a concave function, then
$$f(E[X]) \geq E[f(X)]$$

Image: Wikipedia

# Alternative Derivation

$$\log p(D) = \log \int p(\theta, D) \ d\theta$$

$$= \log \int \frac{p(\theta, D) \ q(\theta)}{q(\theta)} \ d\theta$$

$$= \log E_{q(\theta)}[\frac{p(\theta, D)}{q(\theta)}]$$

Jensen's inequality $\longrightarrow$

$$\geq E_{q(\theta)} \left[ \log \frac{p(\theta, D)}{q(\theta)} \right]$$

Log is a concave function, then
$$f(E[X]) \geq E[\ f(X)]$$

# Alternative Derivation

Model Evidence

$$\log p(D) = \log \int p(\theta, D) \, d\theta$$

$$= \log \int \frac{p(\theta, D) \, q(\theta)}{q(\theta)} \, d\theta$$

$$= \log E_{q(\theta)}\left[\frac{p(\theta, D)}{q(\theta)}\right]$$

$$\geq E_{q(\theta)}\left[\log \frac{p(\theta, D)}{q(\theta)}\right]$$

Evidence Lower Bound (ELBO)

# Variational Inference (VI)

The posterior

$$p(\theta|D) = p(D|\theta)p(\theta)/p(D)$$

The variational distribution

$$q_\phi(\theta)$$

$$L = E_{q_{\phi(\theta)}}\left[\log\frac{p(D,\theta)}{q_\phi(\theta)}\right] = \log p(D) - KL[\,q_\phi(\theta)||p(\theta)]$$



$q \in Q$

$q^*(\theta)$

$p(\theta|D)$

# Mean-field Variational Inference

- A type of choices of the variational distribution
- The name origins in the mean field theory of physics
- The variational distribution factorizes

$$q_\phi(\boldsymbol{\theta}) = \prod_{i=1}^{K} q_{\phi_i}(\theta_i)$$

Opper and Saad, eds. Advanced mean field methods: Theory and practice. MIT press, 2001.

# A Gaussian Example



$$p(\mathbf{z}) = N(z|\mu, \Lambda^{-1})$$

$$q(z) = q(z_1)q(z_2)$$

Bishop (2006). Pattern recognition and machine learning. Springer.

# Mean-field Variational Inference

ELBO

Fully Factorized Variational Distribution

$$L = E_{q(\theta)} \left[ \log \frac{\mathrm{p}(\theta, D)}{q(\theta)} \right] \qquad \longleftarrow \qquad q(\boldsymbol{\theta}) = \prod_{i=1}^{K} q_{\phi_i}(\theta_i)$$

Jordan et al. An introduction to variational methods for graphical models. Machine learning 37.2 (1999): 183-233.

# Mean-field Variational Inference

ELBO

Fully Factorized Variational Distribution

$$L = E_{q(\theta)}\left[\log\frac{p(\theta,D)}{q(\theta)}\right] \qquad \longleftarrow \qquad q(\boldsymbol{\theta}) = \prod_{i=1}^{K} q_{\phi_i}(\theta_i)$$

$$L = \int q(\theta_j)E_{q(\theta_{\neg j})}[\log p(\theta_j, D \mid \theta_{\neg j})]d\theta_j - \int q(\theta_j)\log p(\theta_j)\, d\theta_j + c_j$$

Jordan et al. An introduction to variational methods for graphical models. Machine learning 37.2 (1999): 183-233.

# Mean-field Variational Inference

ELBO

Fully Factorized Variational Distribution

$$L = E_{q(\theta)}\left[\log\frac{p(\theta, D)}{q(\theta)}\right] \longleftarrow q(\boldsymbol{\theta}) = \prod_{i=1}^{K} q_{\phi_i}(\theta_i)$$

$$L = \int q(\theta_j) E_{q(\theta_{\neg j})}[\log p(\theta_j, D \mid \theta_{\neg j})]d\theta_j - \int q(\theta_j)\log p(\theta_j)\, d\theta_j + c_j$$

$$q^*(\theta_j) \propto \exp(E_{q(\theta_{\neg j})}[\log p(\theta_j, D \mid \theta_{\neg j})])$$

Jordan et al. An introduction to variational methods for graphical models. Machine learning 37.2 (1999): 183-233.

# Part II: Advances

- **Scalable variational inference**

- **Monte Carlo methods**

- **Amortized inference**

- Approximate distribution design

- Optimization objective design

# Stochastic Variational Inference



$$\mathrm{p}(\theta, \boldsymbol{\xi}, \boldsymbol{x}) = p(\theta) \prod_{i=1}^{M} p(\xi_i|\theta) p(x_i|\xi_i, \theta)$$

# Stochastic Variational Inference



$$\mathrm{p}(\theta,\boldsymbol{\xi},\boldsymbol{x}) = p(\theta) \prod_{i=1}^{M} p(\xi_i|\theta)p(x_i|\xi_i,\theta)$$

$$L = E_q \left[ \log \frac{\mathrm{p}(\theta,\boldsymbol{\xi},\boldsymbol{x})}{q(\theta,\xi)} \right]$$

$$= E_q \left[ \log \frac{p(\theta) \prod_{i=1}^{M} p(\xi_i|\theta)p(x_i|\xi_i,\theta)}{q(\theta) \prod_{i=1}^{M} q(\xi_i)} \right]$$

$$= E_q \left[ \log \frac{p(\theta)}{q(\theta)} \right] + \sum_{i=1}^{M} E_q \left[ \log \frac{p(\xi_i|\theta)p(x_i|\xi_i)}{q(\xi_i)} \right]$$

- O(M) time to compute in each update iteration
- M can be extremely large
- Even one iteration might not be affordable

# Stochastic Variational Inference



$$p(\theta, \boldsymbol{\xi}, \boldsymbol{x}) = p(\theta) \prod_{i=1}^{M} p(\xi_i|\theta) p(x_i|\xi_i, \theta)$$

$$L = E_q\left[\log\frac{p(\theta)}{q(\theta)}\right] + \sum_{i=1}^{M} E_q\left[\log\frac{p(\xi_i|\theta)p(x_i|\xi_i)}{q(\xi_i)}\right]$$

stochastic approximation
with $S \ll M$

$$\hat{L} = E_q\left[\log\frac{p(\theta)}{q(\theta)}\right] + \frac{M}{S}\sum_{i=1}^{S} E_q\left[\log\frac{p(\xi_i|\theta)p(x_i|\xi_i)}{q(\xi_i)}\right]$$

Computational complexity: $O(M) \rightarrow O(S)$

Hoffman et al. Stochastic Variational Inference. JMLR 2013.

# How Stochastic Gradient Works



$\nabla F(x_i)$ gradient of each single data point $x_i$

$E_x[\nabla F(x)]$ batch gradient considering all data points

# How Stochastic Gradient Works



$\nabla F(x_i)$ gradient of each single data point $x_i$

$E_x[\nabla F(x)]$ batch gradient considering all $M = 10$ data points

$\frac{M}{S}\sum_{s=1}^{S}\nabla F(x_s)$ mini$-$batch gradient/stochastic gradient estimated using S=3 data points

# How Stochastic Gradient Works



$\nabla F(x_i)$ gradient of each single data point $x_i$

$E_x[\nabla F(\boldsymbol{x})]$ batch gradient considering all $M = 10$ data points

$\frac{M}{S}\sum_{s=1}^{S}\nabla F(x_s)$ mini$-$batch gradient/stochastic gradient estimated using S=3 data points

Zhang. et.al. "Determinantal Point Processes for Mini-Batch Diversification." UAI, 2017

# Stochastic Variational Inference

$$L = E_q\left[\log\frac{p(\theta)}{q(\theta)}\right] + \sum_{i=1}^{M} E_q\left[\log\frac{p(\xi_i|\theta)p(x_i|\xi_i)}{q(\xi_i)}\right] \xrightarrow{\text{gradient}} \nabla L = \nabla E_q\left[\log\frac{p(\theta)}{q(\theta)}\right] + \sum_{i=1}^{M} E_q\left[\nabla\log\frac{p(\xi_i|\theta)p(x_i|\xi_i)}{q(\xi_i)}\right]$$

Stochastic approximation
with $S \ll M$

Stochastic approximation
with $S \ll M$

$$\hat{L} = E_q\left[\log\frac{p(\theta)}{q(\theta)}\right] + \frac{M}{S}\sum_{i=1}^{S} E_q\left[\log\frac{p(\xi_i|\theta)p(x_i|\xi_i)}{q(\xi_i)}\right] \xrightarrow{\text{gradient}} \nabla\hat{L} = \nabla E_q\left[\log\frac{p(\theta)}{q(\theta)}\right] + \frac{M}{S}\sum_{i=1}^{S} E_q\left[\nabla\log\frac{p(\xi_i|\theta)p(x_i|\xi_i)}{q(\xi_i)}\right]$$

Stochastic Gradient

Hoffman et al. Stochastic Variational Inference. JMLR 2013.

Nature laughs at the difficulties of integration.

--Pierre-Simon Laplace

Gordon and Sorkin. *The Armchair Science Reader*. New York 1959

# Monte Carlo Approximation

- To approximate: $E_{p(x)}[f(x)]$

- MC Approximation:
  1. Sample $x_1, x_2, \ldots., x_K \sim p(x)$
  2. Evaluate $f(x_i)$ for each sample
  3. Compute $E[f(x)] \approx \frac{1}{K}\sum_{i=1}^{K} f(x_i)$

  <span style="color:red">Unbiased Monte Carlo estimate</span>

# Log-derivative Trick

$$\nabla_\theta \log p_\theta(x)$$

# Log-derivative Trick

$$\nabla_\theta \log p_\theta(x) = \frac{1}{p_\theta(x)} \nabla p_\theta(X)$$

# Log-derivative Trick

$$\nabla_\theta \log p_\theta(x) = \frac{1}{p_\theta(x)} \nabla p_\theta(X)$$

$$\nabla_\theta p_\theta(x) = p_\theta(x) \, \nabla_\theta \log p_\theta(x)$$

# REINFORCE Gradients

ELBO  $L = E_{q_\phi(\theta)}\left[\log\dfrac{\mathrm{p}(\theta,D)}{q_\phi(\theta)}\right]$

Gradient of the ELBO

$$\nabla_\phi L = \nabla_\phi E_{q_\phi(\theta)}\left[\log\dfrac{\mathrm{p}(\theta,D)}{q_\phi(\theta)}\right] = \int \nabla_\phi\left\{ q_\phi(\theta)\log\dfrac{\mathrm{p}(\theta,D)}{q_\phi(\theta)}\right\} d\theta$$

Glynn (1990). Likelihood ratio gradient estimation for stochastic systems. Communications of the ACM, 33(10), 75–84.
Williams (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. Machine Learning, 8(3-4), 229–256.
Fu (2006). Gradient estimation. Handbooks in Operations Research and Management Science, 13, 575–616.

# REINFORCE Gradients

ELBO $L = E_{q_\phi(\theta)}\left[\log\frac{p(\theta, D)}{q_\phi(\theta)}\right]$

Gradient of the ELBO

$$\nabla_\phi L = \nabla_\phi E_{q_\phi(\theta)}\left[\log\frac{p(\theta, D)}{q_\phi(\theta)}\right] = \int \nabla_\phi\left\{ q_\phi(\theta)\log\frac{p(\theta, D)}{q_\phi(\theta)}\right\} d\theta$$

$$= \int \nabla_\phi q_\phi(\theta)\log\frac{p(\theta, D)}{q_\phi(\theta)} d\theta + \int q_\phi(\theta)\nabla_\phi\log\frac{p(\theta, D)}{q_\phi(\theta)} d\theta$$

Glynn (1990). Likelihood ratio gradient estimation for stochastic systems. Communications of the ACM, 33(10), 75–84.

Williams (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. Machine Learning, 8(3-4), 229–256.

Fu (2006). Gradient estimation. Handbooks in Operations Research and Management Science, 13, 575–616.

# REINFORCE Gradients

$$\nabla q_\phi(\theta) = q_\phi(\theta)\, \nabla_\phi \log q_\phi(\theta)$$

ELBO $\quad L = E_{q_\phi(\theta)}\left[\log\frac{\mathrm{p}(\theta,D)}{q_\phi(\theta)}\right]$

Gradient of the ELBO

$$\nabla_\phi L = \nabla_\phi E_{q_\phi(\theta)}\left[\log\frac{\mathrm{p}(\theta,D)}{q_\phi(\theta)}\right] = \int \nabla_\phi\left\{ q_\phi(\theta)\log\frac{\mathrm{p}(\theta,D)}{q_\phi(\theta)}\right\} d\theta$$

$$= \int \nabla_\phi q_\phi(\theta)\log\frac{\mathrm{p}(\theta,D)}{q_\phi(\theta)} d\theta + \int q_\phi(\theta)\, \nabla_\phi\log\frac{\mathrm{p}(\theta,D)}{q_\phi(\theta)} d\theta$$

$$= \int q_\phi(\theta)\nabla_\phi\log q_\phi(\theta)\ \log\frac{\mathrm{p}(\theta,D)}{q_\phi(\theta)} d\theta - \int \nabla_\phi q_\phi(\theta) d\theta$$

Glynn (1990). Likelihood ratio gradient estimation for stochastic systems. Communications of the ACM, 33(10), 75–84.
Williams (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. Machine Learning, 8(3-4), 229–256.
Fu (2006). Gradient estimation. Handbooks in Operations Research and Management Science, 13, 575–616.

# REINFORCE Gradients

ELBO $\quad L = E_{q_\phi(\theta)}\left[\log\frac{p(\theta,D)}{q_\phi(\theta)}\right]$

Gradient of the ELBO

$$\nabla_\phi L = \nabla_\phi E_{q_\phi(\theta)}\left[\log\frac{p(\theta,D)}{q_\phi(\theta)}\right] = \int \nabla_\phi\left\{ q_\phi(\theta)\log\frac{p(\theta,D)}{q_\phi(\theta)}\right\}d\theta$$

$$= \int \nabla_\phi q_\phi(\theta)\log\frac{p(\theta,D)}{q_\phi(\theta)}d\theta + \int q_\phi(\theta)\nabla_\phi\log\frac{p(\theta,D)}{q_\phi(\theta)}d\theta$$

$$= \int q_\phi(\theta)\nabla_\phi\log q_\phi(\theta)\ \log\frac{p(\theta,D)}{q_\phi(\theta)}d\theta - \int \nabla_\phi q_\phi(\theta)d\theta$$

$$= \nabla_\phi\int q_\phi(\theta)d\theta = \nabla_\phi 1 = 0$$

$$\nabla_\phi\log\frac{p(\theta,D)}{q_\phi(\theta)}$$

$$= \frac{q_\phi(\theta)}{p(\theta,D)}\left(-\frac{p(\theta,D)}{q_\phi(\theta)^2}\right)\nabla q_\phi(\theta)$$

$$= -\frac{\nabla_\phi q_\phi(\theta)}{q_\phi(\theta)}$$

Glynn (1990). Likelihood ratio gradient estimation for stochastic systems. Communications of the ACM, 33(10), 75–84.

Williams (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. Machine Learning, 8(3-4), 229–256.

Fu (2006). Gradient estimation. Handbooks in Operations Research and Management Science, 13, 575–616.

# REINFORCE Gradients

ELBO  $L = E_{q_\phi(\theta)}\left[\log\frac{p(\theta,D)}{q_\phi(\theta)}\right]$

Gradient of the ELBO

$$\nabla_\phi L = \nabla_\phi E_{q_\phi(\theta)}\left[\log\frac{p(\theta,D)}{q_\phi(\theta)}\right] = \int \nabla_\phi\left\{q_\phi(\theta)\log\frac{p(\theta,D)}{q_\phi(\theta)}\right\}d\theta$$

$$= \int \nabla_\phi q_\phi(\theta)\log\frac{p(\theta,D)}{q_\phi(\theta)}d\theta + \int q_\phi(\theta)\nabla_\phi\log\frac{p(\theta,D)}{q_\phi(\theta)}d\theta$$

$$= \int q_\phi(\theta)\nabla_\phi\log q_\phi(\theta)\ \log\frac{p(\theta,D)}{q_\phi(\theta)}d\theta - \int \nabla_\phi q_\phi(\theta)d\theta$$

$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \underline{\phantom{xxxxxxxxxxxxx}}$$

$$= \ \nabla_\phi\int q_\phi(\theta)d\theta = \nabla_\phi 1 = 0$$

$$= E_{q_\phi(\theta)}\left[\underbrace{\nabla_\phi\log q_\phi(\theta)}_{\text{Score function}}\log\frac{p(\theta,D)}{q_\phi(\theta)}\right]$$

Glynn (1990). Likelihood ratio gradient estimation for stochastic systems. Communications of the ACM, 33(10), 75–84.
Williams (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. Machine Learning, 8(3-4), 229–256.
Fu (2006). Gradient estimation. Handbooks in Operations Research and Management Science, 13, 575–616.

# BBVI

ELBO $\quad L = E_{q_\phi(\theta)} \left[ \overset{f(\theta)}{\boxed{\log \frac{p(\theta, D)}{q_\phi(\theta)}}} \right]$

- To approximate: $E[f(x)]$
- MC Approximation:
  1. Sample $x_1, x_2, \ldots, x_K \sim p(x)$
  2. Evaluate $f(x_i)$ for each sample
  3. Compute $E[f(x)] \approx \frac{1}{K} \sum_{i=1}^{K} f(x_i)$

Gradient of the ELBO

$$\nabla_\phi L = E_{q_\phi(\theta)} \left[ \underbrace{\nabla_\phi \log q_\phi(\theta)}_{\text{Score function}} \overset{f(\theta)}{\log \frac{p(\theta, D)}{q_\phi(\theta)}} \right]$$

Ranganath et al. Black box variational inference. AISTATS 2014
Glynn (1990). Likelihood ratio gradient estimation for stochastic systems. Communications of the ACM, 33(10), 75–84.
Williams (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. Machine Learning, 8(3-4), 229–256.
Fu (2006). Gradient estimation. Handbooks in Operations Research and Management Science, 13, 575–616.

# BBVI

### Gradient of the ELBO

$$\nabla_\phi L = E_{q_\phi(\theta)} \left[ \nabla_\phi \log q_\phi(\theta) \log \frac{p(\theta, D)}{q_\phi(\theta)} \right]$$

- To approximate: $E[f(x)]$
- MC Approximation:
  1. Sample $x_1, x_2, \ldots, x_K \sim p(x)$
  2. Evaluate $f(x_i)$ for each sample
  3. Compute $E[f(x)] \approx \frac{1}{K} \sum_{i=1}^{K} f(x_i)$

1. Sample $\theta_1, \theta_2, \ldots, \theta_K \sim q_\phi(\theta)$

Ranganath et al. Black box variational inference. AISTATS 2014

# BBVI

### Gradient of the ELBO

$$\nabla_\phi L = E_{q_\phi(\theta)}\left[\boxed{\nabla_\phi \log q_\phi(\theta) \log \frac{p(\theta,D)}{q_\phi(\theta)}}\right]$$

- To approximate: $E[f(x)]$
- MC Approximation:
  1. Sample $x_1, x_2, \ldots, x_K \sim p(x)$
  2. Evaluate $f(x_i)$ for each sample
  3. Compute $E[f(x)] \approx \frac{1}{K}\sum_{i=1}^{K} f(x_i)$

1. Sample $\theta_1, \theta_2, \ldots, \theta_K \sim q_\phi(\theta)$

2. Evaluate $\boxed{\nabla_\phi \log q_\phi(\theta) \log \frac{p(\theta,D)}{q_\phi(\theta)}}$ for each sample

Ranganath et al. Black box variational inference. AISTATS 2014

# BBVI

### Gradient of the ELBO

$$\nabla_\phi L = E_{q_\phi(\theta)} \left[ \boxed{\nabla_\phi \log q_\phi(\theta) \log \frac{\mathrm{p}(\theta, D)}{q_\phi(\theta)}} \right]$$

- To approximate: $E[f(x)]$
- MC Approximation:
  1. Sample $x_1, x_2, \ldots, x_K \sim p(x)$
  2. Evaluate $f(x_i)$ for each sample
  3. Compute $E[f(x)] \approx \frac{1}{K} \sum_{i=1}^{K} f(x_i)$

1. Sample $\theta_1, \theta_2, \ldots, \theta_K \sim q_\phi(\theta)$

2. Evaluate $\boxed{\nabla_\phi \log q_\phi(\theta) \log \frac{\mathrm{p}(\theta, D)}{q_\phi(\theta)}}$ for each sample

3. The approximated gradient is:

$$\nabla_\phi \hat{L} = \frac{1}{K} \sum_{i=1}^{K} \nabla_\phi \log q_\phi(\theta_i) \log \frac{\mathrm{p}(\theta_i, D)}{q_\phi(\theta_i)}$$

Ranganath et al. Black box variational inference. AISTATS 2014

# Black-box Variational Inference (BBVI)



Go beyond conjugate exponential family

# Reparameterization Trick

Express $q_\phi(\theta)$ as $\epsilon \sim r(\epsilon), \; \theta = g(\epsilon, \phi)$

Kingma and Welling. Auto-encoding variational bayes. ICLR 2014
Salimans and Knowles. Fixed-Form Variational Posterior Approximation through Stochastic Linear Regression. Bayesian Analysis 2013

# Reparameterization Trick

Express $q_\phi(\theta)$ as $\epsilon \sim r(\epsilon),\ \theta = g(\epsilon, \phi)$

$$\theta \sim N(\mu, \sigma^2)$$
$$\epsilon \sim N(0,1), \theta = \mu + \sigma\epsilon$$

Kingma and Welling. Auto-encoding variational bayes. ICLR 2014
Salimans and Knowles. Fixed-Form Variational Posterior Approximation through Stochastic Linear Regression. Bayesian Analysis 2013

# Reparameterization Trick

Express $q_\phi(\theta)$ as $\epsilon \sim r(\epsilon),\ \theta = g(\epsilon, \phi)$

$$\theta \sim N(\mu, \sigma^2)$$
$$\epsilon \sim N(0,1), \theta = \mu + \sigma\epsilon$$

ELBO $\quad L = E_{q_\phi(\theta)}\left[\log\frac{\mathrm{p}(\theta,D)}{q_\phi(\theta)}\right] = E_{r(\epsilon)}\left[\log\frac{\mathrm{p}(g(\epsilon,\phi),D)}{q_\phi(g(\epsilon,\phi))}\right]$

Kingma and Welling. Auto-encoding variational bayes. ICLR 2014
Salimans and Knowles. Fixed-Form Variational Posterior Approximation through Stochastic Linear Regression. Bayesian Analysis 2013

# Reparameterization Trick

Express $q_\phi(\theta)$ as $\epsilon \sim r(\epsilon),\ \theta = g(\epsilon, \phi)$

$$\theta \sim N(\mu, \sigma^2)$$
$$\epsilon \sim N(0,1), \theta = \mu + \sigma\epsilon$$

ELBO $\ L = E_{q_\phi(\theta)}\left[\log\frac{p(\theta,D)}{q_\phi(\theta)}\right] = E_{r(\epsilon)}\left[\log\frac{p(g(\epsilon,\phi),D)}{q_\phi(g(\epsilon,\phi))}\right]$

Gradient $\ \nabla_\phi L = \nabla_\phi E_{r(\epsilon)}\left[\log\frac{p(g(\epsilon,\phi),D)}{q_\phi(g(\epsilon,\phi))}\right]$

Kingma and Welling. Auto-encoding variational bayes. ICLR 2014
Salimans and Knowles. Fixed-Form Variational Posterior Approximation through Stochastic Linear Regression. Bayesian Analysis 2013

# Reparameterization Trick

Express $q_\phi(\theta)$ as $\epsilon \sim r(\epsilon), \ \theta = g(\epsilon, \phi)$

$$\theta \sim N(\mu, \sigma^2)$$
$$\epsilon \sim N(0,1), \theta = \mu + \sigma\epsilon$$

ELBO $\quad L = E_{q_\phi(\theta)}\left[\log \frac{p(\theta, D)}{q_\phi(\theta)}\right] = E_{r(\epsilon)}\left[\log \frac{p(g(\epsilon, \phi), D)}{q_\phi(g(\epsilon, \phi))}\right]$

Gradient $\quad \nabla_\phi L = \nabla_\phi E_{r(\epsilon)}\left[\log \frac{p(g(\epsilon, \phi), D)}{q_\phi(g(\epsilon, \phi))}\right] = E_{r(\epsilon)}\left[\nabla_\phi \log \frac{p(g(\epsilon, \phi), D)}{q_\phi(g(\epsilon, \phi))}\right]$

Kingma and Welling. Auto-encoding variational bayes. ICLR 2014
Salimans and Knowles. Fixed-Form Variational Posterior Approximation through Stochastic Linear Regression. Bayesian Analysis 2013

# Reparameterization Trick

Express $q_\phi(\theta)$ as $\epsilon \sim r(\epsilon),\ \theta = g(\epsilon, \phi)$

$$\theta \sim N(\mu, \sigma^2)$$
$$\epsilon \sim N(0,1), \theta = \mu + \sigma\epsilon$$

ELBO $\quad L = E_{q_\phi(\theta)}\left[\log \frac{\mathrm{p}(\theta, D)}{q_\phi(\theta)}\right] = E_{r(\epsilon)}\left[\log \frac{\mathrm{p}(g(\epsilon, \phi), D)}{q_\phi(g(\epsilon, \phi))}\right]$

Gradient $\quad \nabla_\phi L = \nabla_\phi E_{r(\epsilon)}\left[\log \frac{\mathrm{p}(g(\epsilon, \phi), D)}{q_\phi(g(\epsilon, \phi))}\right] = E_{r(\epsilon)}\left[\nabla_\phi \log \frac{\mathrm{p}(g(\epsilon, \phi), D)}{q_\phi(g(\epsilon, \phi))}\right]$

$$\nabla_\phi \hat{L} = \frac{1}{K}\sum_{k=1}^{K} \nabla_\phi \log \frac{\mathrm{p}(g(\epsilon_k, \phi), D)}{q_\phi(g(\epsilon_k, \phi))}, \epsilon_k \sim r(\epsilon)$$

# Reparameterization Trick

Express $q_\phi(\theta)$ as $\epsilon \sim r(\epsilon),\ \theta = g(\epsilon, \phi)$

$$\theta \sim N(\mu, \sigma^2)$$
$$\epsilon \sim N(0,1), \theta = \mu + \sigma\epsilon$$

ELBO $\quad L = E_{q_\phi(\theta)}\left[\log \frac{p(\theta, D)}{q_\phi(\theta)}\right] = E_{r(\epsilon)}\left[\log \frac{p(g(\epsilon,\phi), D)}{q_\phi(g(\epsilon,\phi))}\right]$

Gradient $\quad \nabla_\phi L = \nabla_\phi E_{r(\epsilon)}\left[\log \frac{p(g(\epsilon,\phi), D)}{q_\phi(g(\epsilon,\phi))}\right] = E_{r(\epsilon)}\left[\nabla_\phi \log \frac{p(g(\epsilon,\phi), D)}{q_\phi(g(\epsilon,\phi))}\right]$

$$\nabla_\phi \hat{L} = \frac{1}{K}\sum_{k=1}^{K}\nabla_\phi \log \frac{p(g(\epsilon_k, \phi), D)}{q_\phi(g(\epsilon_k, \phi))}, \epsilon_k \sim r(\epsilon)$$

# Reparameterization Trick

Express $q_\phi(\theta)$ as $\epsilon \sim r(\epsilon), \; \theta = g(\epsilon, \phi)$

$$\theta \sim N(\mu, \sigma^2)$$
$$\epsilon \sim N(0,1), \theta = \mu + \sigma\epsilon$$

ELBO $\; L = E_{q_\phi(\theta)}\left[\log\frac{\text{p}(\theta,D)}{q_\phi(\theta)}\right] = E_{r(\epsilon)}\left[\log\frac{\text{p}(g(\epsilon,\phi),D)}{q_\phi(g(\epsilon,\phi))}\right]$

Gradient $\; \nabla_\phi L = \nabla_\phi E_{r(\epsilon)}\left[\log\frac{\text{p}(g(\epsilon,\phi),D)}{q_\phi(g(\epsilon,\phi))}\right] = E_{r(\epsilon)}\left[\nabla_\phi \log\frac{\text{p}(g(\epsilon,\phi),D)}{q_\phi(g(\epsilon,\phi))}\right]$

$\nabla_\phi \hat{L} = \frac{1}{K}\sum_{k=1}^{K}\nabla_\phi \log\frac{\text{p}(g(\epsilon_k,\phi),D)}{q_\phi(g(\epsilon_k,\phi))} \; , \epsilon_k \sim r(\epsilon)$

# Variance Reduction Techniques in MCVI

- When non-differentiable, falls back to REINFORCE gradient

High variance!

# Variance Reduction Techniques in MCVI

- When non-differentiable, falls back to REINFORCE gradient

High variance!



- Solutions to high variance REINFORCE gradients:
  - Low variance unbiased estimators with control variates
  - Biased estimators to enable reparam. trick (potentially low variance)

# Variance Reduction Techniques in MCVI

- Control variate method:
  - Assume we want to estimate with MC simulation

$$E_{q(\theta)}[F(\theta)] \approx \frac{1}{K}\sum_{k}^{K} F(\theta_k), \qquad \theta_k \sim q(\theta)$$



$F(\theta)$

# Variance Reduction Techniques in MCVI

- Control variate method:
  - Assume we want to estimate with MC simulation

$$E_{q(\theta)}[F(\theta)] \approx \frac{1}{K}\sum_k^K F(\theta_k), \qquad \theta_k \sim q(\theta)$$

  - Control variate: define a control function $G(\theta)$ satisfying:
    - $V_{q(\theta)}[G(\theta)] < \infty$
    - Known or fast computable $E_{q(\theta)}[G(\theta)]$

$F(\theta)$

$G(\theta)$

# Variance Reduction Techniques in MCVI

- Control variate method:
  - Then define the new MC estimator

$$E_{q(\theta)}[F(\theta)] \approx \frac{1}{K}\sum_{k}^{K} \hat{F}(\theta_k), \qquad \theta_k \sim q(\theta),$$



$$\hat{F}(\theta) \qquad = \qquad F(\theta) \qquad - \qquad G(\theta) \qquad + \qquad E_{q(\theta)}[G(\theta)]$$

$$V_{q(\theta)}\big[\hat{F}(\theta)\big] = V_{q(\theta)}[F(\theta)] + \underline{V_{q(\theta)}[G(\theta)] - 2\,Cov_{q(\theta)}[F(\theta), G(\theta)]}$$

$< 0$ if $F$ and $G$ are strongly and positively correlated

# Variance Reduction Techniques in MCVI

- Application to REINFORCE gradient:

  - $F(\theta) = \log \dfrac{p(D,\theta)}{q_\phi(\theta)} \nabla_\phi \log q_\phi(\theta)$

    $\coloneqq f(\theta)$

- Define $G(\theta) = g(\theta) \nabla_\phi \log q_\phi(\theta)$

# Variance Reduction Techniques in MCVI

- Application to REINFORCE gradient:

  - $F(\theta) = \log \dfrac{p(D,\theta)}{q_\phi(\theta)} \nabla_\phi \log q_\phi(\theta)$

    $:= f(\theta)$

  - Define $G(\theta) = g(\theta) \nabla_\phi \log q_\phi(\theta)$

  - Variance reduced gradient: $\hat{F}(\theta) = \big(f(\theta) - g(\theta)\big) \nabla_\phi \log q_\phi(\theta) + E_{q_\phi(\theta)}[G(\theta)]$

    $:= \Delta(\theta)$

# Variance Reduction Techniques in MCVI

- Application to REINFORCE gradient:

  - $F(\theta) = \log \frac{p(D,\theta)}{q_\phi(\theta)} \nabla_\phi \log q_\phi(\theta)$

    $:= f(\theta)$

  - Define $G(\theta) = g(\theta) \nabla_\phi \log q_\phi(\theta)$

  - Variance reduced gradient: $\hat{F}(\theta) = \big(f(\theta) - g(\theta)\big) \nabla_\phi \log q_\phi(\theta) + E_{q_\phi(\theta)}[G(\theta)]$

    $:= \Delta(\theta)$

  - "Baseline" approach:

    $g(\theta) = b$

    $\Rightarrow E_{q(\theta)}[G(\theta)] = b E_{q(\theta)}\big[\nabla_\phi \log q(\theta)\big] = b \nabla_\phi \int q_\phi(\theta) d\theta = b \nabla_\phi 1 = 0$

    (log-derivative trick)

    $\Rightarrow \hat{F}(\theta) = \Delta(\theta) \nabla_\phi \log q_\phi(\theta)$

Ranganath et al. Black Box Variational Inference. AISTATS 2014
Mnih and Gregor. Neural Variational Inference and Learning in Belief Networks. ICML 2014

# Variance Reduction Techniques in MCVI

- Application to REINFORCE gradient:

  - $F(\theta) = \log \frac{p(D,\theta)}{q_\phi(\theta)} \nabla_\phi \log q_\phi(\theta)$    • Define $G(\theta) = g(\theta)\nabla_\phi \log q_\phi(\theta)$

    $\underbrace{\phantom{\log \frac{p(D,\theta)}{q_\phi(\theta)}}}_{:= f(\theta)}$

  - Variance reduced gradient: $\hat{F}(\theta) = \underbrace{\big(f(\theta) - g(\theta)\big)}_{:= \Delta(\theta)}\nabla_\phi \log q_\phi(\theta) + E_{q_\phi(\theta)}[G(\theta)]$

- "Baseline" approach:

  $g(\theta) = b$

  $\Rightarrow E_{q(\theta)}[G(\theta)] = \underbrace{b E_{q(\theta)}[\nabla_\phi \log q(\theta)] = 0}_{\text{log-derivative trick}}$

  $\Rightarrow \hat{F}(\theta) = \Delta(\theta)\nabla_\phi \log q_\phi(\theta)$



$b$ fitted by minimising either $V_{q(\theta)}[\hat{F}(\theta)]$ or $E_{q(\theta)}[\Delta(\theta)^2]$

Ranganath et al. Black Box Variational Inference. AISTATS 2014
Mnih and Gregor. Neural Variational Inference and Learning in Belief Networks. ICML 2014

# Variance Reduction Techniques in MCVI

- Application to REINFORCE gradient:

  - $F(\theta) = \log \frac{p(D,\theta)}{q_\phi(\theta)} \nabla_\phi \log q_\phi(\theta)$

    $:= f(\theta)$

  - Define $G(\theta) = g(\theta)\nabla_\phi \log q_\phi(\theta)$

  - Variance reduced gradient: $\hat{F}(\theta) = \big(f(\theta) - g(\theta)\big)\nabla_\phi \log q_\phi(\theta) + E_{q_\phi(\theta)}[G(\theta)]$

    $:= \Delta(\theta)$

  - "Taylor expansion" approach (e.g. 1st order):

    $g(\theta) = f(\theta_0) + \nabla_{\theta_0} f(\theta_0)(\theta - \theta_0)$

Paisley et al. Variational Bayesian Inference with Stochastic Search. ICML 2012
Gu et al. MuProp: Unbiased Backpropagation for Stochastic Neural Networks. ICLR 2016

# Variance Reduction Techniques in MCVI

- Application to REINFORCE gradient:

  - $F(\theta) = \log \frac{p(D,\theta)}{q_\phi(\theta)} \nabla_\phi \log q_\phi(\theta)$         · Define $G(\theta) = g(\theta) \nabla_\phi \log q_\phi(\theta)$

    $:= f(\theta)$

  - Variance reduced gradient: $\hat{F}(\theta) = \big(f(\theta) - g(\theta)\big) \nabla_\phi \log \mathrm{q}_\phi(\theta) + E_{\mathrm{q}_\phi(\theta)}[G(\theta)]$

    $:= \Delta(\theta)$

- "Taylor expansion" approach (e.g. 1st order):

  $= 0$ (log-derivative trick)

  $g(\theta) = f(\theta_0) + \nabla_{\theta_0} f(\theta_0)(\theta - \theta_0)$         $\Rightarrow E_{q(\theta)}[G(\theta)] = \big(f(\theta_0) - \nabla_{\theta_0} f(\theta_0)\theta_0\big) E_{q(\theta)}\big[\nabla_\phi \log q(\theta)\big]$

  $+ \nabla_{\theta_0} f(\theta_0) E_{q(\theta)}\big[\theta \nabla_\phi \log q(\theta)\big]$

  $= \nabla_{\theta_0} f(\theta_0) \nabla_\phi E_{q(\theta)}[\theta]$ (log-derivative trick)

Paisley et al. Variational Bayesian Inference with Stochastic Search. ICML 2012
Gu et al. MuProp: Unbiased Backpropagation for Stochastic Neural Networks. ICLR 2016

# Variance Reduction Techniques in MCVI

- Application to REINFORCE gradient:

  - $F(\theta) = \log \frac{p(D,\theta)}{q_\phi(\theta)} \nabla_\phi \log q_\phi(\theta)$

    $:= f(\theta)$

  - Define $G(\theta) = g(\theta)\nabla_\phi \log q_\phi(\theta)$

  - Variance reduced gradient: $\hat{F}(\theta) = \big(f(\theta) - g(\theta)\big)\nabla_\phi \log q_\phi(\theta) + E_{q_\phi(\theta)}[G(\theta)]$

    $:= \Delta(\theta)$

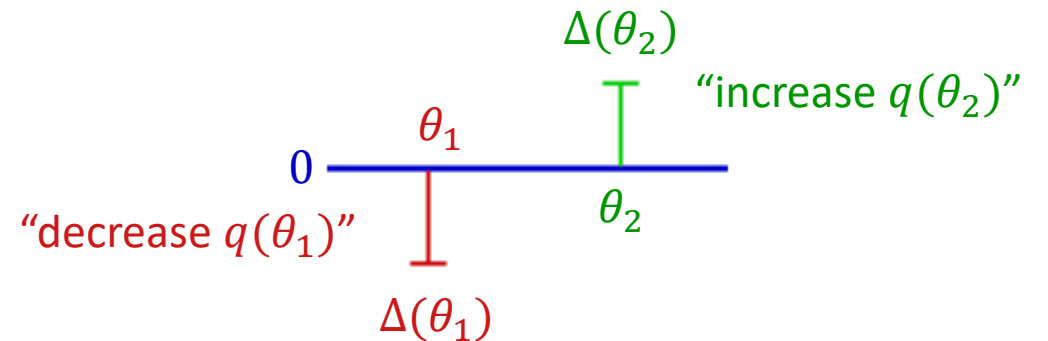- "Taylor expansion" approach (e.g. 1st order):

  $= 0$ (log-derivative trick)

  $g(\theta) = f(\theta_0) + \nabla_{\theta_0} f(\theta_0)(\theta - \theta_0)$
  $\Rightarrow E_{q(\theta)}[G(\theta)] = \big(f(\theta_0) - \nabla_{\theta_0} f(\theta_0)\theta_0\big) E_{q(\theta)}\big[\nabla_\phi \log q(\theta)\big]$

  $+\nabla_{\theta_0} f(\theta_0) E_{q(\theta)}\big[\theta \nabla_\phi \log q(\theta)\big]$

  $= \nabla_{\theta_0} f(\theta_0) \nabla_\phi E_{q(\theta)}[\theta]$ (log-derivative trick)

  $\Rightarrow \hat{F}(\theta) = \Delta(\theta)\nabla_\phi \log q_\phi(\theta) + \nabla_{\theta_0} f(\theta_0)\nabla_\phi E_{q(\theta)}[\theta]$

  $:= \nabla_\phi f(\theta_0)$ if $\theta_0 = E_{q(\theta)}[\theta]$

Paisley et al. Variational Bayesian Inference with Stochastic Search. ICML 2012
Gu et al. MuProp: Unbiased Backpropagation for Stochastic Neural Networks. ICLR 2016

# Variance Reduction Techniques in MCVI

- Gumbel-Softmax trick
    - Biased gradient estimator
    - Empirically found to have smaller variance



Categorical distribution:

$$p(y = k) = \pi_k, \sum_k \pi_k = 1$$

# Variance Reduction Techniques in MCVI

- Gumbel-Softmax trick
  - Biased gradient estimator
  - Empirically found to have smaller variance

a) Categorical

expectation

b)

sample

category

Gumbel trick to sample $y$:
$$y = \arg max \ [g_k + \log \pi_k],$$
$$g_k \sim Gumbel(0, 1)$$

# Variance Reduction Techniques in MCVI

- Gumbel-Softmax trick
  - Biased gradient estimator
  - Empirically found to have smaller variance



Gumbel trick to sample $y$:
$$y = \boxed{\arg max} \, [g_k + \log \pi_k],$$
$$g_k \sim Gumbel(0, 1)$$

replace with softmax

Jang et al. Categorical Reparameterization with Gumbel-Softmax. ICLR 2017
Maddison et al. The Concrete Distribution: A Continuous Relaxation of Discrete Random Variables. ICLR 2017

# Variance Reduction Techniques in MCVI

- Gumbel-Softmax trick
  - Biased gradient estimator
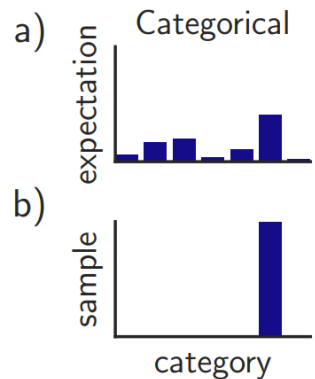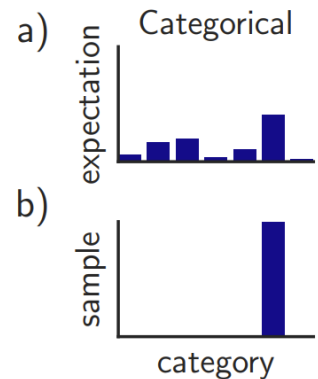  - Empirically found to have smaller variance



Concrete distribution/Gumbel-Softmax trick:
sample the "soft vector" (instead of one-hot encoding of $y$)

$$[y_1, \dots, y_K] = softmax([\frac{(g_1 + \log \pi_1)}{\tau}, \dots, \frac{(g_K + \log \pi_K)}{\tau}])$$

Jang et al. Categorical Reparameterization with Gumbel-Softmax. ICLR 2017
Maddison et al. The Concrete Distribution: A Continuous Relaxation of Discrete Random Variables. ICLR 2017

# Variance Reduction Techniques in MCVI



For an incomplete list of variance reduced gradient estimators, see http://yingzhenli.net/home/en/?page_id=1262

# Latent Variable Model

# Deep Latent Variable Model

# Amortized Inference



$\phi$  parameter for variational distribution
$\theta$  decoder parameter

# Amortized Inference



$\phi$  parameter for variational distribution
$\theta$  decoder parameter

# Amortized Inference



$$L = \log p(\boldsymbol{x}) - KL[q(\boldsymbol{z})||p(\boldsymbol{z}|\boldsymbol{x})]$$

$$L_{amortized} = \log p(\boldsymbol{x}) - KL[q(\boldsymbol{z}|\boldsymbol{x})||p(\boldsymbol{z}|\boldsymbol{x})]$$

# Variational Auto-Encoders (VAE)



$$L = \log p(\boldsymbol{x}) - KL[q(\boldsymbol{z})||p(\boldsymbol{z}|\boldsymbol{x})]$$

$$L_{amortized} = \log p(\boldsymbol{x}) - KL[q(\boldsymbol{z}|\boldsymbol{x})||p(\boldsymbol{z}|\boldsymbol{x})]$$

Encoder/
inference network

Decoder/
generator

# Variational Auto-Encoders (VAE)



$$L = \log p(\boldsymbol{x}) - KL[q(\boldsymbol{z})||p(\boldsymbol{z}|\boldsymbol{x})]$$

$$L_{amortized} = \log p(\boldsymbol{x}) - KL[q(\boldsymbol{z}|\boldsymbol{x})||p(\boldsymbol{z}|\boldsymbol{x})]$$

$\mu, \sigma$

$\epsilon \sim N(0,1)$

$z \sim N(\mu, \sigma^2)$
$\epsilon \sim N(0,1), \theta = \mu + \sigma\epsilon$

Kingma and Welling. Auto-encoding variational bayes. ICLR 2014.
Rezende et al. Stochastic backpropagation and approximate inference in deep generative models. ICML 2014.

# Variational Auto-Encoders (VAE)



$$L_{amortized} = \log p(\boldsymbol{x}) - KL[q(\boldsymbol{z}|\boldsymbol{x})||p(\boldsymbol{z}|\boldsymbol{x})]$$
$$= E_{\boldsymbol{z} \sim q(\boldsymbol{z}|\boldsymbol{x})}[\log \mathrm{p}_{\theta}(\boldsymbol{x}|\boldsymbol{z})] - KL[q(\boldsymbol{z}|\boldsymbol{x})||p(\boldsymbol{z}|\boldsymbol{x})]$$

Kingma and Welling. Auto-encoding variational bayes. ICLR 2014.
Rezende et al. Stochastic backpropagation and approximate inference in deep generative models. ICML 2014.

# How to apply amortization to other inference methods?

# Amortized Inference: Further Examples

- Amortized SMC



Find the optimal proposal distribution for each sequence $\{x_{1:T}\}$ & each time step

amortize

Explicitly parameterise & optimise $(x_{1:T}, z_{1:t}) \rightarrow$ proposal dist. for $z_t$

Naesseth et al. Variational Sequential Monte Carlo. AISTATS 2018
Maddison et al. Filtering Variational Objectives. NeurIPS 2017
Le et al. Auto-encoding Sequential Monte Carlo. ICLR 2018

# Amortized Inference: Further Examples

- Amortized MCMC

Hoffman. Learning Deep Latent Gaussian Models with Markov Chain Monte Carlo. ICML 2017
**Li** et al. Approximate Inference with Amortised MCMC. ICML 2017 Workshop on Implicit Models

# Amortized Inference: Further Examples

- Amortized Monte Carlo integration

Goal: estimate with
importance sampling

$$E_{p(z|x)}[F_\eta(z)]$$

amortize

Find the optimal proposal
distributions for each $(x, \eta)$ pair

Explicitly parameterise & optimise
$(x, \eta) \rightarrow$ proposal dist. for $z$

Goliński et al. Amortized Monte Carlo Integration. ICML 2019

# Amortized Inference: Limitations

- Amortised approximate posteriors in practice are sub-optimal



- The "refinement" idea:
  - Initialise $q(z|x) = N(z; \mu, \sigma^2)$ with the amortised solution $\mu \leftarrow \mu_\phi(x), \sigma \leftarrow \sigma_\phi(x)$
  - Then run $T$ more VI gradient steps to update $\mu, \sigma$

Cremer et al. Inference Suboptimality in Variational Autoencoders. ICML 2018
Marino et al. Iterative Amortized Inference. ICML 2018
Kim et al. Semi-Amortized Variational Autoencoders. ICML 2018

# Part II: Advances

- Scalable variational inference

- Monte Carlo methods

- Amortized inference

- **Approximate distribution design**

- **Optimization objective design**

# Designing $q$ Distributions


Structured approximations


Normalizing flows


Auxiliary variables & mixture distributions


Implicit approximate posteriors

# Structured Approximations

- introduce dependencies between random variables for $q$:

Hidden Markov Model

Mean-field approximation

Exact posterior $p(z \mid x)$

$z_i \not\perp z_j \mid x$

$$q(z) = \prod_i q(z_i)$$

# Structured Approximations

- introduce dependencies between random variables for $q$:



| Hidden Markov Model | Mean-field approximation | Structured approximation |
|---|---|---|
| Exact posterior $p(z \mid x)$ $z_i \not\perp z_j \mid x$ | $q(z) = \prod_i q(z_i)$ | $q(z) = \prod_s q(z_s)$ $q(z_s) = q(\{z_i\}_{i \in s})$ |

<span style="color:red">Main design question: the grouping and conditional dependency structure</span>

# Structured Approximations

- Auto-regressive distributions (as a specific dependency structure)



Hidden Markov Model

Exact posterior $p(z \mid x)$

$z_i \not\perp z_j \mid x$

Structured approximation

$q(z) = \prod_s q(z_s)$
$q(z_s) = q(\{z_i\}_{i \in s})$

Auto-regressive approximation

$q(z) = \prod_i q(z_i \mid z_{<i})$

$q(z_1 \mid z_{<1}) = q(z_1)$

Main design question: the ordering of the latent variables

# Normalizing Flows

- Change-of-variable formula:
  - $x$ is a random variable with probability density function (PDF) $p_X(x)$
  - $y = f(x)$ is an invertible mapping

# Normalizing Flows

- Change-of-variable formula:
  - $x$ is a random variable with probability density function (PDF) $p_X(x)$
  - $y = f(x)$ is an invertible mapping
  - The probability mass is preserved, and the PDF for $y = f(x)$ satisfies

$$p_Y(y)dy = p_X(x)dx$$

prob. mass of region around $y$        prob. mass of region around $x$

# Normalizing Flows

- Change-of-variable formula:
  - $x$ is a random variable with probability density function (PDF) $p_X(x)$
  - $y = f(x)$ is an invertible mapping
  - The probability mass is preserved, and the PDF for $y = f(x)$ satisfies

$$p_Y(y)dy = p_X(x)dx$$

<span style="color:red">prob. mass of region around $y$</span>      <span style="color:blue">prob. mass of region around $x$</span>

$$p_Y(y) = p_X(x)|\det(\frac{dx}{dy})|$$

$$p_X(x) = p_Y(y)|\det(\frac{dy}{dx})|$$



$dx$            $f$            $dy$

$x$                    $y$            hypercube approximation
(exact when $\text{vol}(dx) \to 0$)

# Normalizing Flows

- Variational inference with Normalizing flow
  - Assume $q_0(z_0) = N(z_0; 0, I)$
  - Define $z = f_\phi(z_0)$ where $f_\phi(\cdot)$ is an invertible mapping parameterized by $\phi$

$$q(z) = q_0(z_0) \left| \det\left(\frac{dz}{dz_0}\right) \right|^{-1} \qquad \text{with } z_0 = f_\phi^{-1}(z)$$

(change of variable: $q(z)dz = q_0(z_0)dz_0$)

Rezende and Mohamed. Variational Inference with Normalizing Flows. ICML 2015

# Normalizing Flows

- Variational inference with Normalizing flow
  - Assume $q_0(z_0) = N(z_0; 0, I)$
  - Define $z = f_\phi(z_0)$ where $f_\phi(\cdot)$ is an invertible mapping parameterized by $\phi$

$$q(z) = q_0(z_0) |\det\left(\frac{dz}{dz_0}\right)|^{-1} \qquad \text{with } z_0 = f_\phi^{-1}(z)$$

  - Fit $q(z)$ to $p(x \mid z)$ with VI:

$$L(q(z)) = E_{q(z)}[\log p(x \mid z) + \log p(z) - \log q(z)]$$

Rezende and Mohamed. Variational Inference with Normalizing Flows. ICML 2015

# Normalizing Flows

- Variational inference with Normalizing flow
  - Assume $q_0(z_0) = N(z_0; 0, I)$
  - Define $z = f_\phi(z_0)$ where $f_\phi(\cdot)$ is an invertible mapping parameterized by $\phi$

$$q(z) = q_0(z_0)|\det\left(\frac{dz}{dz_0}\right)|^{-1} \qquad \text{with } z_0 = f_\phi^{-1}(z)$$

  - Fit $q(z)$ to $p(x \mid z)$ with VI:

$$L\big(q(z)\big) = E_{q(z)}[\log p(x \mid z) + \log p(z) - \log q(z)] \qquad \textcolor{red}{\text{by def. of } q(z)}$$
$$= E_{q(z)}\left[\log p(x, z) - \boxed{\log q_0(z_0 = f_\phi^{-1}(z))|\det\left(\frac{dz}{dz_0}\right)|^{-1}}\right]$$

Rezende and Mohamed. Variational Inference with Normalizing Flows. ICML 2015

# Normalizing Flows

- Variational inference with Normalizing flow
  - Assume $q_0(z_0) = N(z_0; 0, I)$
  - Define $z = f_\phi(z_0)$ where $f_\phi(\cdot)$ is an invertible mapping parameterized by $\phi$

$$q(z) = q_0(z_0)|\det\left(\frac{dz}{dz_0}\right)|^{-1} \qquad \text{with} \ \ z_0 = f_\phi^{-1}(z)$$

  - Fit $q(z)$ to $p(x \mid z)$ with VI:

by def. of $q(z)$

$$L(q(z)) = E_{q(z)}[\log p(x \mid z) + \log p(z) - \log q(z)]$$

$$= E_{q(z)}\left[\log p(x, \ z) - \boxed{\log q_0(z_0 = f_\phi^{-1}(z))|\det\left(\frac{dz}{dz_0}\right)|^{-1}}\right]$$

$$= E_{q_0(z_0)}\left[\log p(x, f_\phi(z_0)) - \log q_0(z_0) + \log|\det\left(\frac{df_\phi}{dz_0}\right)|\right]$$

reparam. trick:
$z \sim q(z) \Leftrightarrow z_0 \sim q_0(z_0), z = f_\phi(z_0)$

Rezende and Mohamed. Variational Inference with Normalizing Flows. ICML 2015

# Normalizing Flows

- Variational inference with Normalizing flow
  - Assume $q_0(z_0) = N(z_0; 0, I)$
  - Define $z = f_\phi(z_0)$ where $f_\phi(\cdot)$ is an invertible mapping parameterized by $\phi$

$$q(z) = q_0(z_0) | \det\left(\frac{dz}{dz_0}\right)|^{-1} \qquad \text{with} \ \ z_0 = f_\phi^{-1}(z)$$

  - Fit $q(z)$ to $p(x \mid z)$ with VI:

$$
\begin{aligned}
L\big(q(z)\big) &= E_{q(z)}[\log p(x \mid z) + \log p(z) - \log q(z)] \qquad \textcolor{red}{\text{by def. of } q(z)} \\
&= E_{q(z)}\left[\log p(x, \ z) - \boxed{\log q_0(z_0 = f_\phi^{-1}(z)) | \det\left(\frac{dz}{dz_0}\right)|^{-1}}\right] \\
&= E_{q_0(z_0)}\left[\log p(x, f_\phi(z_0)) - \log q_0(z_0) + \log | \det\left(\frac{df_\phi}{dz_0}\right)|\right]
\end{aligned}
$$

  - Computing ELBO requires $\log | \det\left(\frac{df_\phi}{dz_0}\right)|$

<span style="color:blue">reparam. trick:</span>
<span style="color:blue">$z \sim q(z) \Leftrightarrow z_0 \sim q_0(z_0), z = f_\phi(z_0)$</span>

Rezende and Mohamed. Variational Inference with Normalizing Flows. ICML 2015

# Normalizing Flows

- Variational inference with Normalizing flow
  - Idea: define $f_\phi$ such that $\log |\det \left(\frac{df_\phi}{dz_0}\right)|$ is easy to compute!
    - Chain simple invertible mappings together to make a flexible mapping



$$f_\phi = f_K \circ f_{K-1} \circ \cdots \circ f_1, f_k(\cdot) := f_{\phi_k}(\cdot), \phi = \{\phi_k\}_{k=1}^K$$

  - For each simple mapping, hopefully the Jacobian log-determinant is easy to compute

$$\Rightarrow \log |\det \left(\frac{df_\phi}{dz_0}\right)| = \sum_{k=1}^K \log |\det(\frac{dz_k}{dz_{k-1}})|$$

Rezende and Mohamed. Variational Inference with Normalizing Flows. ICML 2015

# Normalizing Flows

- Goal: construct $f_k$ to enable fast compute of $\log|\det(\frac{dz_k}{dz_{k-1}})|$
  - Example (RealNVP): $y := f_{\phi_k}(x)$ computed as follows



Dinh et al. Density Estimation using Real NVP. ICLR 2017

# Normalizing Flows

- Goal: construct $f_k$ to enable fast compute of $\log|\det(\frac{dz_k}{dz_{k-1}})|$
  - Example (RealNVP): $y := f_{\phi_k}(x)$ computed as follows



split: $x = [x_1, x_2]$

Dinh et al. Density Estimation using Real NVP. ICLR 2017

# Normalizing Flows

- Goal: construct $f_k$ to enable fast compute of $\log|\det(\frac{dz_k}{dz_{k-1}})|$

  - Example (RealNVP): $y \coloneqq f_{\phi_k}(x)$ computed as follows

identity mapping: $y_1 = x_1$



split: $x = [x_1, x_2]$

Dinh et al. Density Estimation using Real NVP. ICLR 2017

# Normalizing Flows

- Goal: construct $f_k$ to enable fast compute of $\log |\det(\frac{dz_k}{dz_{k-1}})|$
  - Example (RealNVP): $y \coloneqq f_{\phi_k}(x)$ computed as follows

identity mapping: $y_1 = x_1$



affine transform: $y_2 = x_2 \odot \exp\big(s(x_1)\big) + t(x_1)$
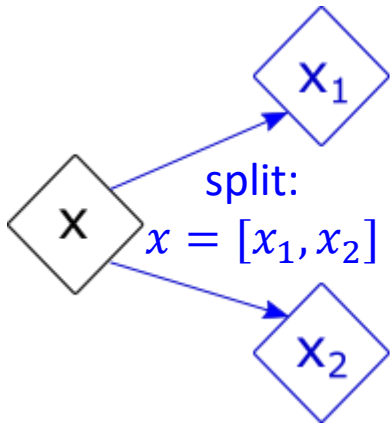
Dinh et al. Density Estimation using Real NVP. ICLR 2017

# Normalizing Flows

- Goal: construct $f_k$ to enable fast compute of $\log |\det(\frac{dz_k}{dz_{k-1}})|$
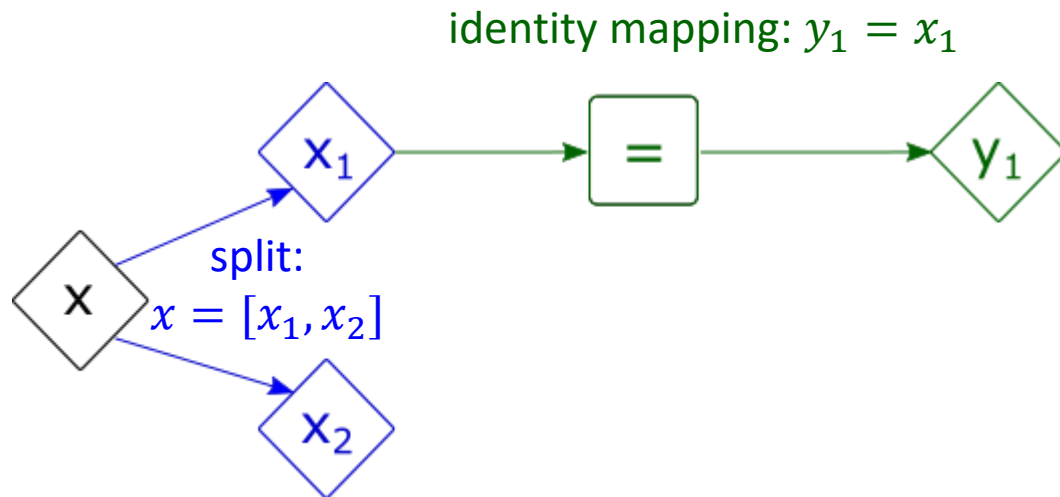  - Example (RealNVP): $y := f_{\phi_k}(x)$ computed as follows

identity mapping: $y_1 = x_1$



split:
$x = [x_1, x_2]$

merge:
$y = [y_1, y_2]$

affine transform: $y_2 = x_2 \odot \exp(s(x_1)) + t(x_1)$

Dinh et al. Density Estimation using Real NVP. ICLR 2017

# Normalizing Flows

- Goal: construct $f_k$ to enable fast compute of $\log|\det(\frac{dz_k}{dz_{k-1}})|$
  - Example (RealNVP): $y := f_{\phi_k}(x)$ computed as follows

identity mapping: $y_1 = x_1$
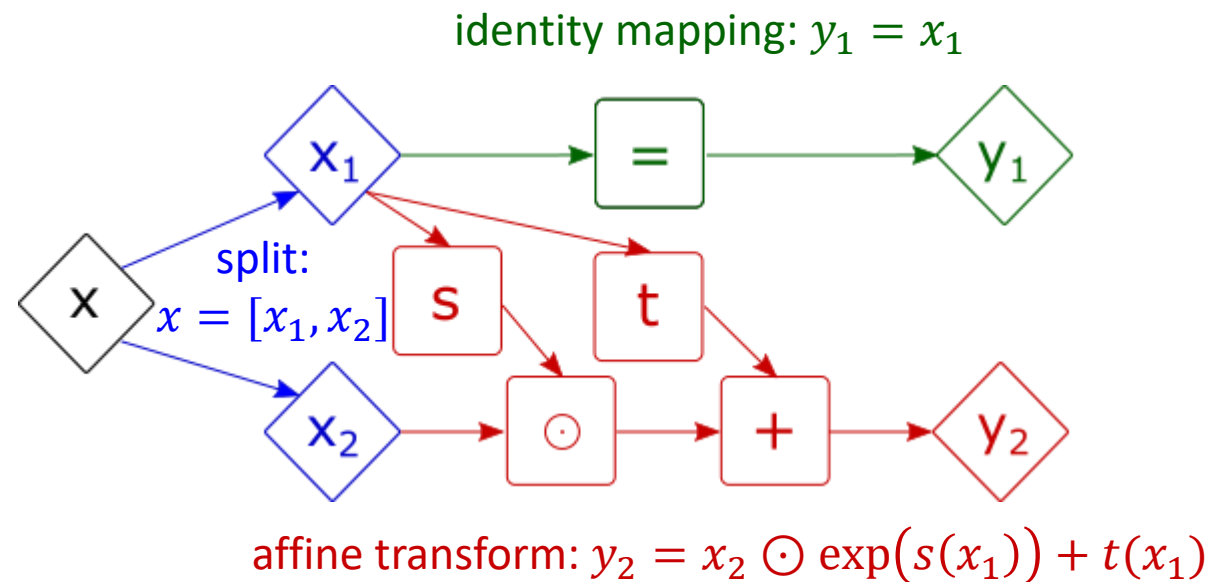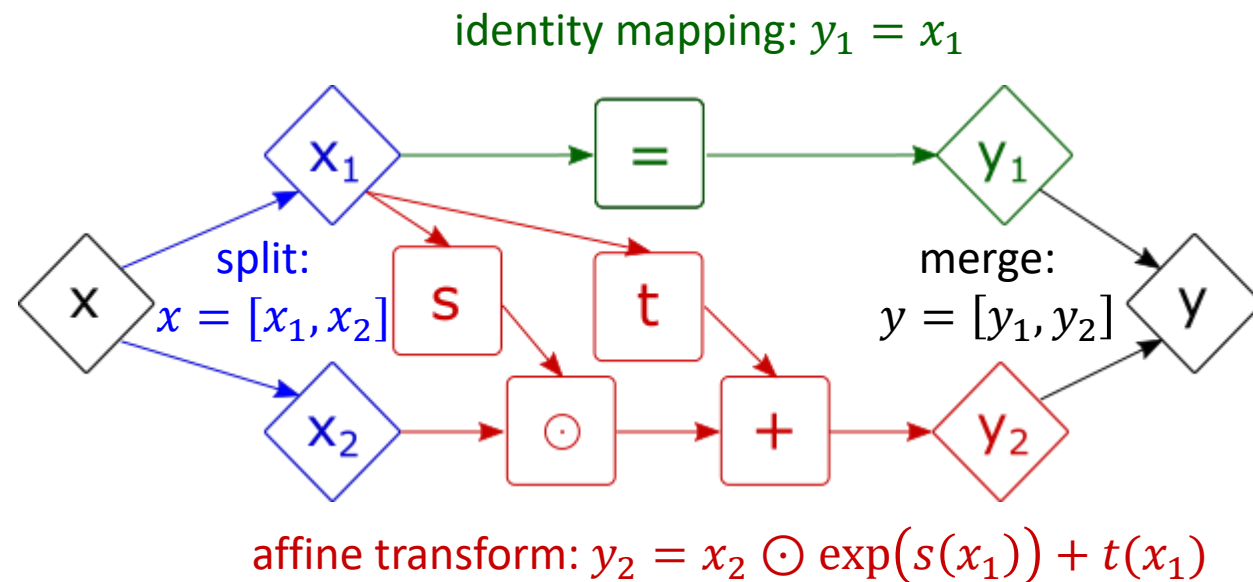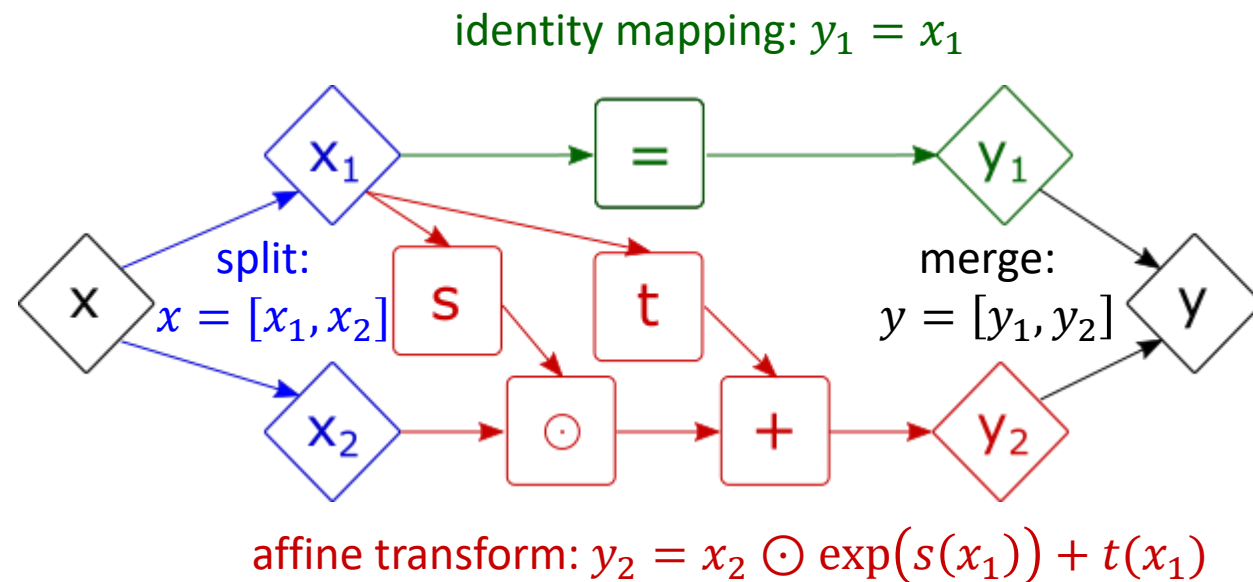


split:
$x = [x_1, x_2]$

merge:
$y = [y_1, y_2]$

affine transform: $y_2 = x_2 \odot \exp(s(x_1)) + t(x_1)$

Jacobian:
$$\frac{df_{\phi_k}}{dx} = \begin{pmatrix} I & 0 \\ dy_2/dx_1 & diag(\exp(s(x_1))) \end{pmatrix}$$

Log-determinant of Jacobian:
$$\Rightarrow \log\left|\det\left(\frac{df_\phi}{dx}\right)\right| = \sum_i s(x_1)_i$$

Dinh et al. Density Estimation using Real NVP. ICLR 2017

# Auxiliary Variables & Mixture Distributions

- Construct $q(\theta)$ as a (hierarchical) mixture distribution

$$q(\theta) = \int q(\theta \,|\, a)\, q(a)\, da$$

  - $a$ is the auxiliary variable used to enrich the approximate posterior

# Auxiliary Variables & Mixture Distributions

- Construct $q(\theta)$ as a (hierarchical) mixture distribution
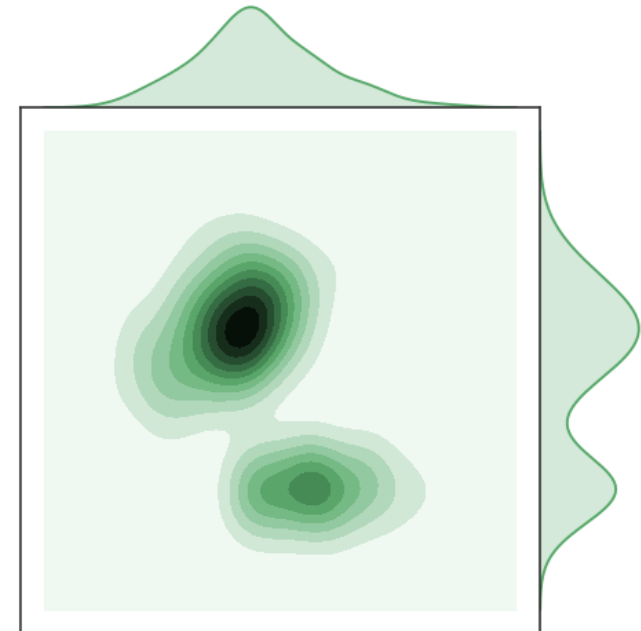
$$q(\theta) = \int q(\theta\,|a)\, q(a)\, da$$

  - $a$ is the auxiliary variable used to enrich the approximate posterior

  - Example: Mixture of Gaussians

$$a \sim q(a) = Categorical(\pi_1, \dots, \pi_K)$$
$$\theta \sim q(\theta\,|a) = N(\theta; m_a, \Sigma_a)$$

<span style="color:red">Can be very flexible with many components!</span>

# Auxiliary Variables & Mixture Distributions

- Construct $q(\theta)$ as a (hierarchical) mixture distribution

$$q(\theta) = \int q(\theta \mid a) \, q(a) \, da$$

- $a$ is the auxiliary variable used to enrich the approximate posterior
- Now the variational lower-bound becomes intractable:

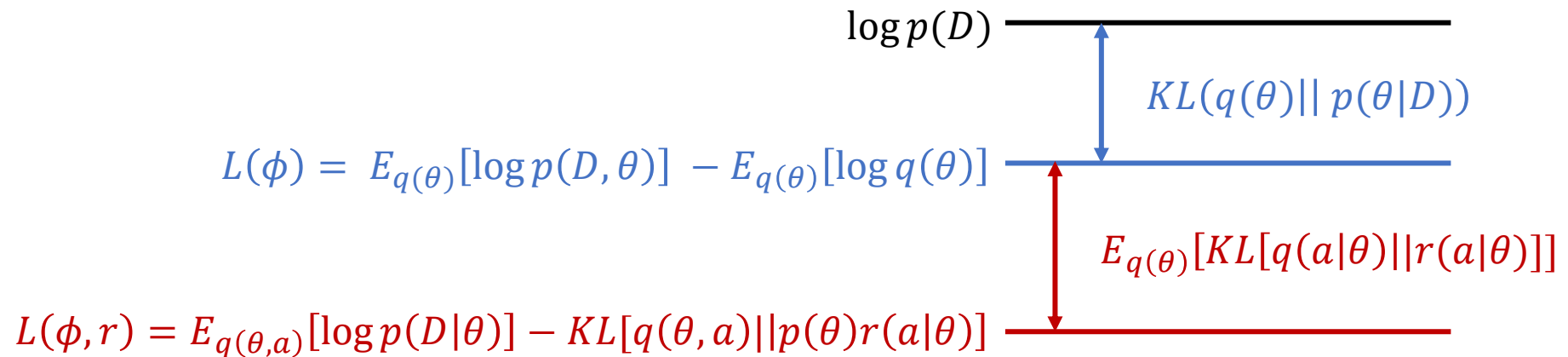$$L(\phi) = \underline{E_{q(\theta)}[\log p(D, \theta)]} - E_{q(\theta)}[\underline{\log q(\theta)}]$$

Estimated by Monte Carlo:
$a_k \sim q(a), \theta_k \sim q(\theta \mid a_k)$

Intractable density
$q(\theta) = \int q(\theta \mid a) q(a) \, da$

# Auxiliary Variables & Mixture Distributions

- Solution: introducing an auxiliary variational lower-bound $L(\phi, r)$ with an auxiliary distribution $r(a|\theta)$:

$$\log p(D)$$

$$KL(q(\theta) \| p(\theta|D))$$

$$L(\phi) = E_{q(\theta)}[\log p(D, \theta)] - E_{q(\theta)}[\log q(\theta)]$$

$$E_{q(\theta)}[KL[q(a|\theta)\|r(a|\theta)]]$$

$$L(\phi, r) = E_{q(\theta, a)}[\log p(D|\theta)] - KL[q(\theta, a)\|p(\theta)r(a|\theta)]$$

Agakov and Barber. An Auxiliary Variational Method. ICONIP 2004
Salimans et al. Markov Chain Monte Carlo and Variational Inference: Bridging the Gap. ICML 2015
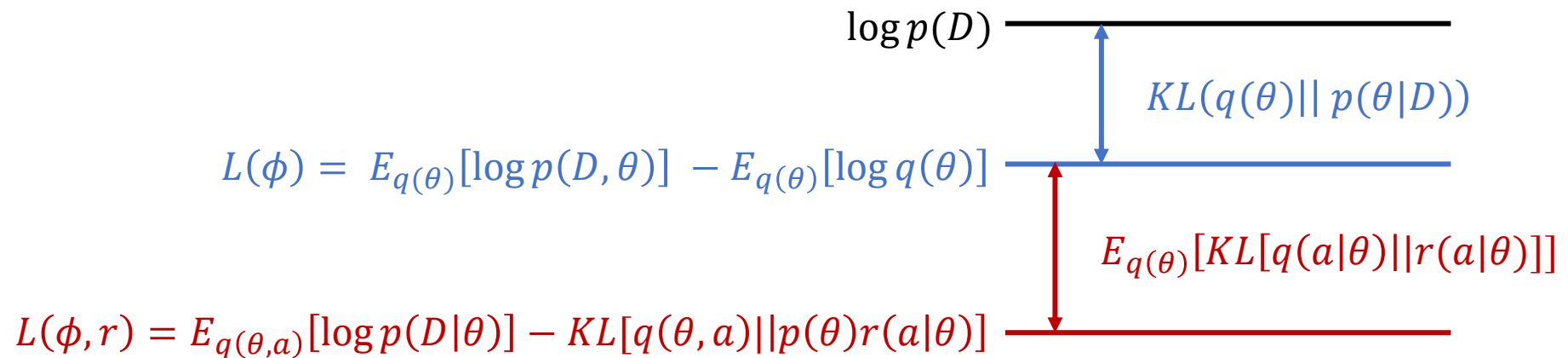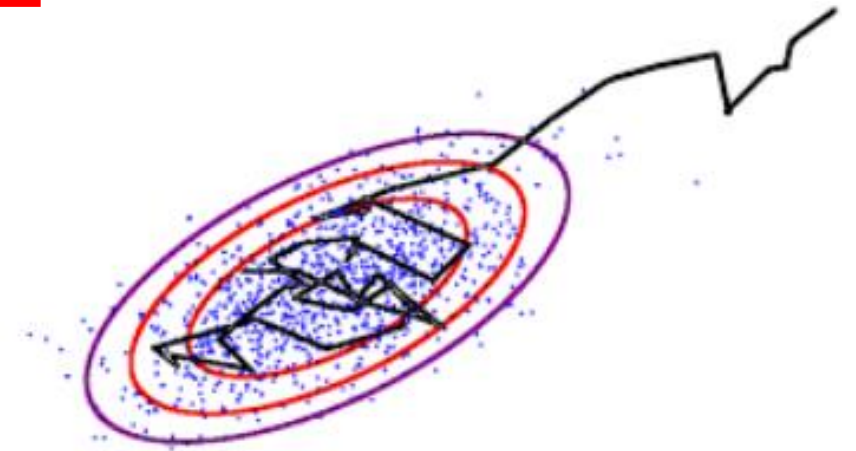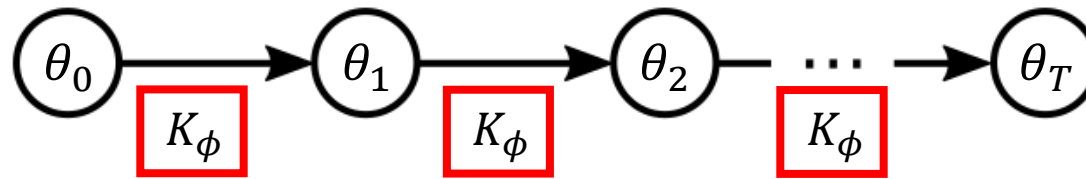Ranganath et al. Hierarchical Variational Models. ICML 2016

# Auxiliary Variables & Mixture Distributions

- Solution: introducing an auxiliary variational lower-bound $L(\phi, r)$ with an auxiliary distribution $r(a|\theta)$:

$$\log p(D) \quad\underline{\hspace{6cm}}$$

$$KL(q(\theta) \| p(\theta|D))$$

$$L(\phi) = E_{q(\theta)}[\log p(D, \theta)] - E_{q(\theta)}[\log q(\theta)]$$

$$E_{q(\theta)}[KL[q(a|\theta) \| r(a|\theta)]]$$

$$L(\phi, r) = E_{q(\theta, a)}[\log p(D|\theta)] - KL[q(\theta, a) \| p(\theta) r(a|\theta)]$$

- Optimize $r(a|\theta)$ to close the gap!
- $L(\phi, r)$ estimated by Monte Carlo: $a_k \sim q(a), \theta_k \sim q(\theta \mid a_k)$

Agakov and Barber. An Auxiliary Variational Method. ICONIP 2004
Salimans et al. Markov Chain Monte Carlo and Variational Inference: Bridging the Gap. ICML 2015
Ranganath et al. Hierarchical Variational Models. ICML 2016
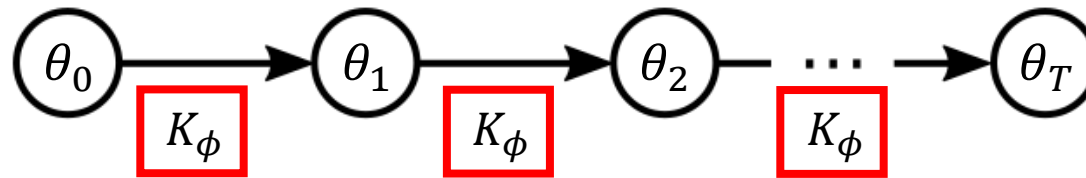
# Auxiliary Variables & Mixture Distributions

- Hierarchical mixture distributions for $q(\theta, a)$
  - VI-MCMC hybrid: build $q(\theta)$ with a Markov Chain:

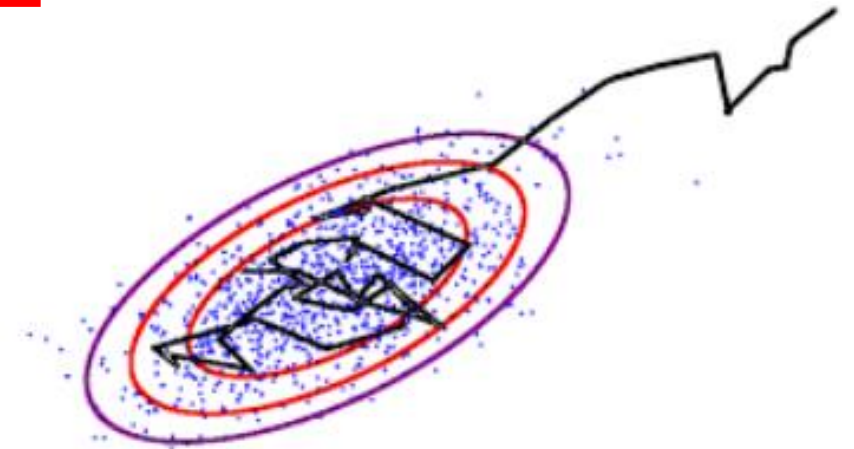# Auxiliary Variables & Mixture Distributions

- Hierarchical mixture distributions for $q(\theta, a)$
  - VI-MCMC hybrid: build $q(\theta)$ with a Markov Chain:



learn the transition kernel with VI:

$$\theta := \theta^T, a = \{\theta^{0:T-1}\}$$

$$q(\theta^T) = \int q_0(\theta^0) \prod_{t=1}^{T} K_\phi(\theta^t | \theta^{t-1}) d\theta^{0:T-1}$$

Salimans et al. Markov Chain Monte Carlo and Variational Inference: Bridging the Gap. ICML 2015
Huang et al. Improving Explorability in Variational Inference with Annealed Variational Objectives. NeurIPS 2018

# Implicit Approximate Posteriors

- Two quantities computed in (approximate) Bayesian inference:

approximate Bayesian predictive

$$p(y^*|x^*, D) \approx E_{q(\theta)}[p(y^*|x^*, \theta)]$$

approximate posterior moments

$$E_{q(\theta)}[F(\theta)]$$

# Implicit Approximate Posteriors

- Two quantities computed in (approximate) Bayesian inference:

approximate Bayesian predictive

$$p(y^*|x^*, D) \approx E_{q(\theta)}[p(y^*|x^*, \theta)]$$

$$\approx \frac{1}{K}\sum_k^K p(y^*|x^*, \theta_k), \ \theta_k \sim q(\theta)$$

approximate posterior moments

$$E_{q(\theta)}[F(\theta)]$$

$$\approx \frac{1}{K}\sum_k^K F(\theta_k), \ \theta_k \sim q(\theta)$$

Computed with Monte Carlo estimates

Only require fast sampling from $q$!
(no need for analytic form of the $q$ distribution)

Mohamed and Lakshminarayanan. Learning in Implicit Generative Models. arXiv 2016
**Li** and Liu. Wild Variational Inference. AABI 2016
Huszár. Variational Inference using Implicit Distributions. arXiv 2017

# Implicit Approximate Posteriors

- Two quantities computed in (approximate) Bayesian inference:

approximate Bayesian predictive

$$p(y^*|x^*, D) \approx E_{q(\theta)}[p(y^*|x^*, \theta)]$$

$$\approx \frac{1}{K} \sum_k^K p(y^*|x^*, \theta_k), \ \theta_k \sim q(\theta)$$

approximate posterior moments

$$E_{q(\theta)}[F(\theta)]$$

$$\approx \frac{1}{K} \sum_k^K F(\theta_k), \ \theta_k \sim q(\theta)$$

**Computed with Monte Carlo estimates**

Only require fast sampling from $q$!
(no need for analytic form of the $q$ distribution)



implicit distributions

Mohamed and Lakshminarayanan. Learning in Implicit Generative Models. arXiv 2016
**Li** and Liu. Wild Variational Inference. AABI 2016
Huszár. Variational Inference using Implicit Distributions. arXiv 2017

# Implicit Approximate Posteriors

$$L(\phi) = \underbrace{E_{q(\theta)}[\log p(D|\theta)]}_{\text{estimated by Monte Carlo}} - E_{q(\theta)}[\log \underbrace{\frac{p(\theta)}{q(\theta)}}_{}]$$

estimated by Monte Carlo

intractable
($q$ density unknown)

Mescheder et al. Adversarial Variational Bayes: Unifying Variational Autoencoders and Generative Adversarial Networks. ICML 2017

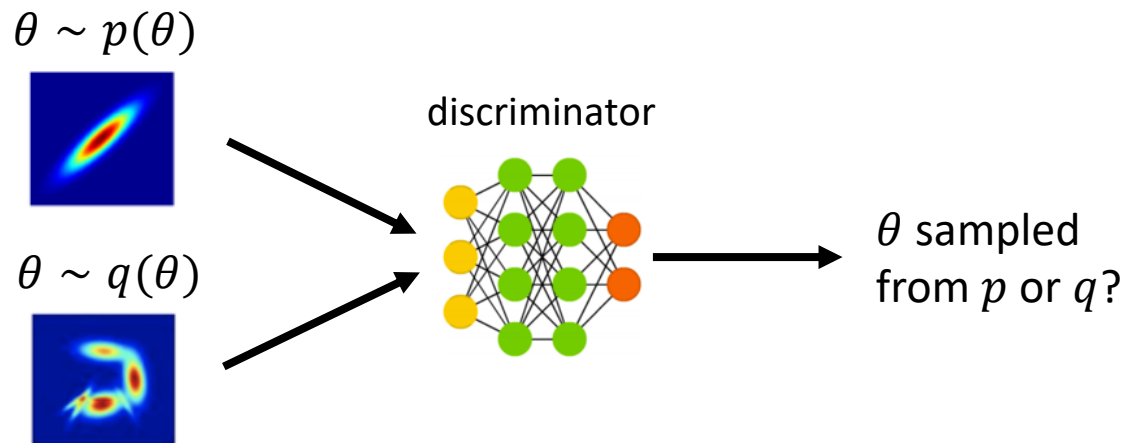Tran et al. Hierarchical Implicit Models and Likelihood-Free Variational Inference. NeurIPS 2017
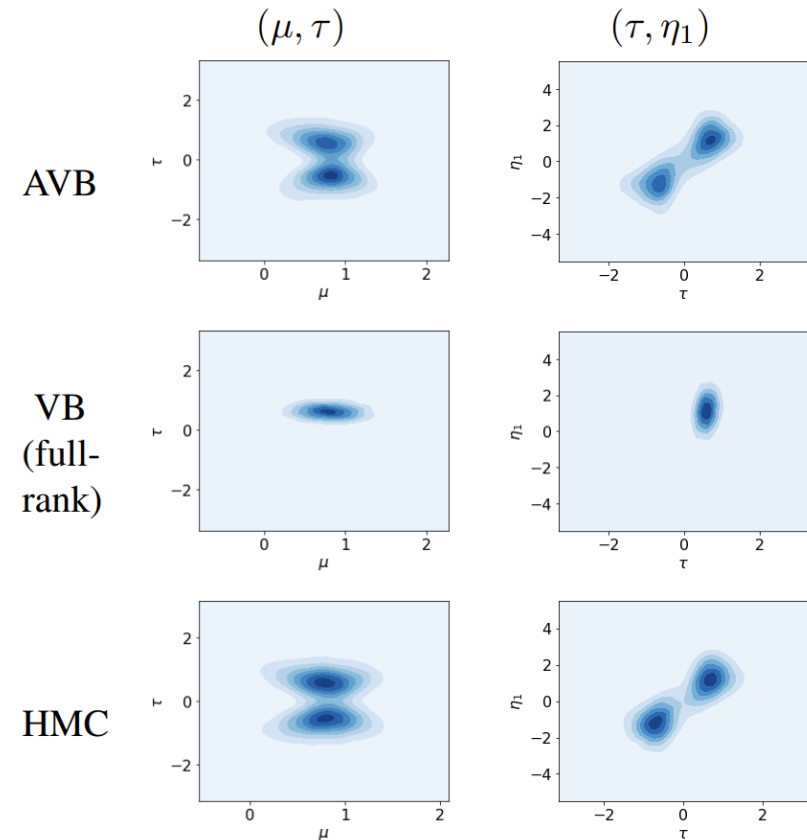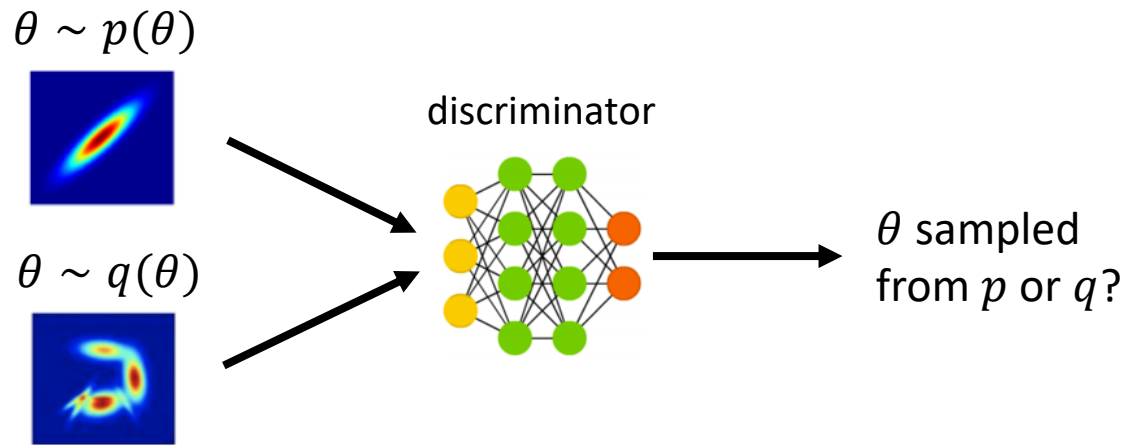
**Li** and Turner. Gradient Estimators for Implicit Models. ICLR 2018

Yin and Zhou. Semi-Implicit Variational Inference. ICML 2018

# Implicit Approximate Posteriors

$$L(\phi) = \underline{E_{q(\theta)}[\log p(D|\theta)]} - E_{q(\theta)}\left[\boxed{\log \frac{p(\theta)}{q(\theta)}}\right]$$

estimated by Monte Carlo

Approximated by using a discriminator (AVB):

$\theta \sim p(\theta)$



discriminator

$\theta \sim q(\theta)$

$\theta$ sampled from $p$ or $q$?

Mescheder et al. Adversarial Variational Bayes: Unifying Variational Autoencoders and Generative Adversarial Networks. ICML 2017
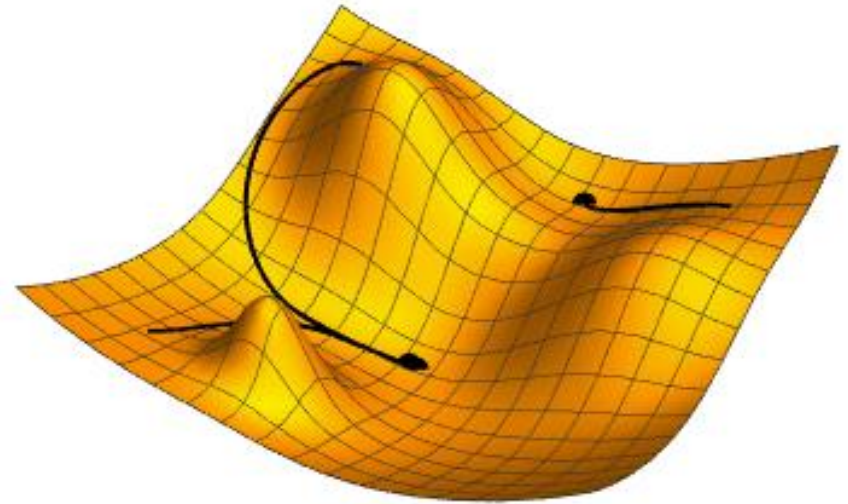Tran et al. Hierarchical Implicit Models and Likelihood-Free Variational Inference. NeurIPS 2017
**Li** and Turner. Gradient Estimators for Implicit Models. ICLR 2018
Yin and Zhou. Semi-Implicit Variational Inference. ICML 2018
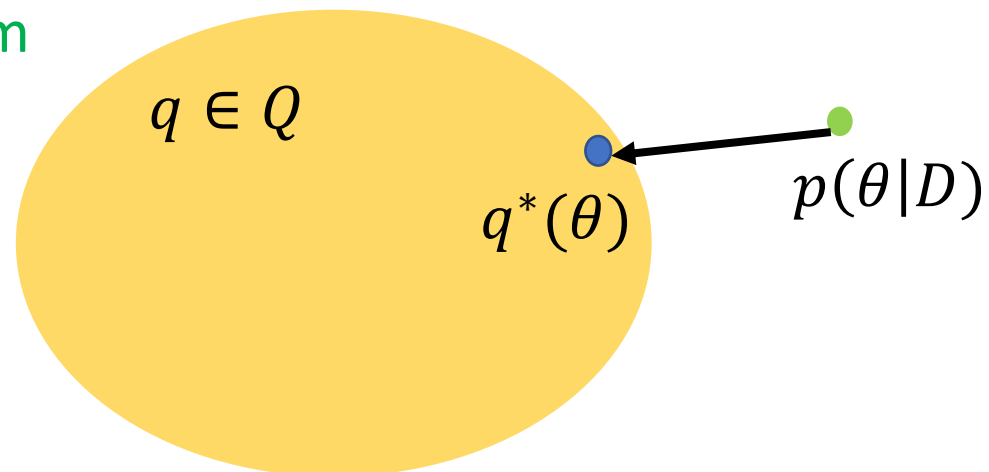
# Implicit Approximate Posteriors

$$L(\phi) = \underline{E_{q(\theta)}[\log p(D|\theta)]} - E_{q(\theta)}\left[\log \boxed{\frac{p(\theta)}{q(\theta)}}\right]$$

estimated by Monte Carlo

Approximated by using a discriminator (AVB):

$\theta \sim p(\theta)$

discriminator

$\theta \sim q(\theta)$

$\theta$ sampled from $p$ or $q$?



$(\mu, \tau)$   $(\tau, \eta_1)$

AVB

VB (full-rank)

HMC

Mescheder et al. Adversarial Variational Bayes: Unifying Variational Autoencoders and Generative Adversarial Networks. ICML 2017
Tran et al. Hierarchical Implicit Models and Likelihood-Free Variational Inference. NeurIPS 2017
**Li** and Turner. Gradient Estimators for Implicit Models. ICLR 2018
Yin and Zhou. Semi-Implicit Variational Inference. ICML 2018

# Objective Functions

For fitting the approximate posterior

# VI Ingredients

$$L = E_{\theta \sim q_\phi} \left[ \log \frac{\boxed{\mathrm{p}(\mathrm{D}, \theta)}}{\boxed{q_\phi(\theta)}} \right] = \log p(D) - \boxed{KL[\, q_\phi || p \,]}$$

p: your model design
q: your choice of variational distribution, e.g. mean field, flow based
KL: defines the algorithm



$q \in Q$

$q^*(\theta)$

$p(\theta|D)$

# Does It Work?



VI underestimate the uncertainty

Bishop (2006). Pattern recognition and machine learning. Springer.

# (Rényi) α-Divergence

$\alpha > 0, \alpha \neq 1$

$$D_\alpha[p||q] = \frac{1}{\alpha - 1}\log \int p(\theta)^\alpha q^{1-\alpha} d\,\theta$$

$\alpha = 1$

$$D_1[p||q] = \lim_{\alpha \to 1} D_{\alpha(p|q)} = KL(p||q)$$

# VI with α-Divergence

ELBO

$$L = E_{\theta \sim q_\phi} \left[ \log \frac{\mathrm{p}(\mathrm{D}, \theta)}{q_\phi(\theta)} \right] = \log p(D) - \boxed{KL[\, q_\phi || p \,]}$$

Variational Rényi bound:

$$L_\alpha = \frac{1}{1 - \alpha} E_{\theta \sim q_\phi} \left[ \left( \log \frac{\mathrm{p}(\mathrm{D}, \theta)}{q_\phi(\theta)} \right)^{1-\alpha} \right] = \log p(D) - \boxed{D_\alpha[\, q_\phi || p \,]}$$

$$\lim_{\alpha \to 1} L_\alpha = L$$

**Li** and Turner. Rényi Divergence Variational Inference. NeurIPS 2016

Dieng et al. Variational Inference via χ-Upper Bound Minimization. NeurIPS 2017

Minka, Tom. Divergence measures and message passing. Technical report, Microsoft Research, 2005.

# Does It Work?

**Li** and Turner. Rényi Divergence Variational Inference. NeurIPS 2016
Dieng et al. Variational Inference via χ-Upper Bound Minimization. NeurIPS 2017
Minka, Tom. Divergence measures and message passing. Technical report, Microsoft Research, 2005.
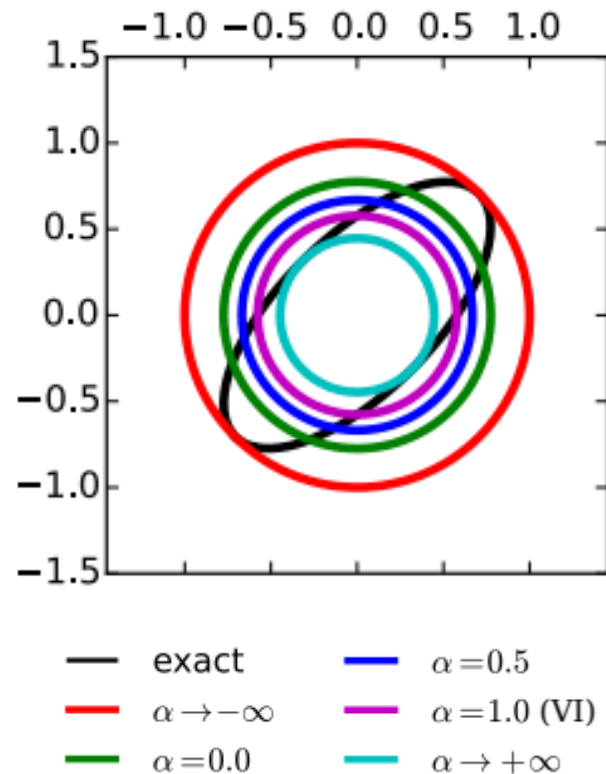
# Does It Work?



How to choose alpha?

**Li** and Turner. Rényi Divergence Variational Inference. NeurIPS 2016

Dieng et al. Variational Inference via χ-Upper Bound Minimization. NeurIPS 2017

Minka, Tom. Divergence measures and message passing. Technical report, Microsoft Research, 2005.

# Does It Work?



How to choose alpha?

$$L_\alpha = \frac{1}{1-\alpha} E_{\theta \sim q_\phi} \left[ \left( \log \frac{p(D, \theta)}{q_\phi(\theta)} \right)^{1-\alpha} \right]$$

**Li** and Turner. Rényi Divergence Variational Inference. NeurIPS 2016

Dieng et al. Variational Inference via χ-Upper Bound Minimization. NeurIPS 2017

Minka, Tom. Divergence measures and message passing. Technical report, Microsoft Research, 2005.

# Does It Work?



How to choose alpha?

$$L_\alpha = \frac{1}{1-\alpha} E_{\theta \sim q_\phi} \left[ \left( \log \frac{p(D, \theta)}{q_\phi(\theta)} \right)^{1-\alpha} \right]$$

Too small or too big alpha leads to extremely big variances

**Li** and Turner. Rényi Divergence Variational Inference. NeurIPS 2016

Dieng et al. Variational Inference via χ-Upper Bound Minimization. NeurIPS 2017

Minka, Tom. Divergence measures and message passing. Technical report, Microsoft Research, 2005.

# Revisiting Perturbation Theory for VI

$$V(x, z) \equiv \log q_\lambda(z) - \log p(x, z)$$

$$\log p(x) = \log\left( E_{z \sim q_\lambda}\left[ \frac{p(x, z)}{q_\lambda(z)} \right] \right) = \log\left( E_{z \sim q_\lambda}\left[ e^{-\beta V(x, z)} \right] \right) \Big|_{\beta=1}$$

Bamler*, **Zhang*** et al. Perturbative Black Box Variational Inference. NeurIPS 2017

# Revisiting Perturbation Theory for VI

$$V(x, z) \equiv \log q_\lambda(z) - \log p(x, z)$$

$$\log p(x) = \log \left( E_{z \sim q_\lambda} \left[ \frac{p(x, z)}{q_\lambda(z)} \right] \right) = \log \left( E_{z \sim q_\lambda} \left[ e^{-\beta V(x,z)} \right] \right) \Big|_{\beta=1}$$
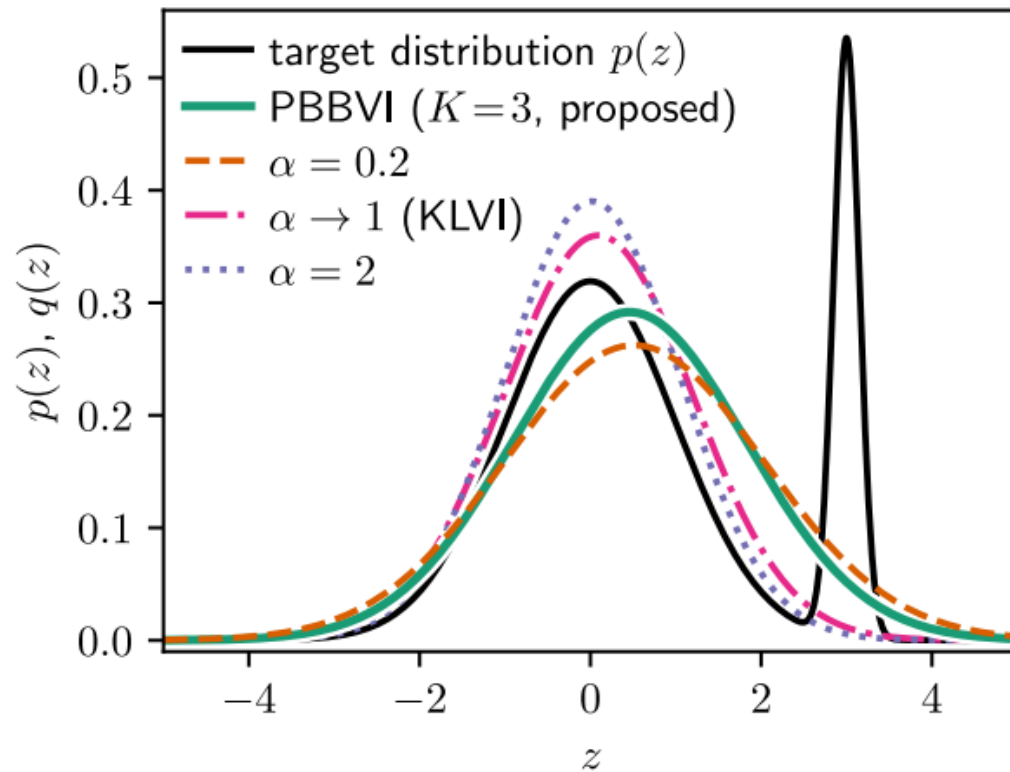
Taylor expansion around $\beta = 1$ :

$$\log p(x) \approx \boxed{E_{q_\lambda}[-V]} + \frac{1}{2} \left[ (V - E_{q_\lambda}[-V])^2 \right] - \frac{1}{3!} \left[ (V - E_{q_\lambda}[-V])^3 \right]$$
$$+ \frac{1}{4!} \left[ (V - E_{q_\lambda}[-V])^4 \right] - \ldots$$

Bamler*, **Zhang*** et al. Perturbative Black Box Variational Inference. NeurIPS 2017

# Revisiting Perturbation Theory for VI

$$V(x,z) \equiv \log q_\lambda(z) - \log p(x,z)$$

$$\log p(x) = \log\left(E_{z \sim q_\lambda}\left[\frac{p(x,z)}{q_\lambda(z)}\right]\right) = \log\left(E_{z \sim q_\lambda}\left[e^{-\beta V(x,z)}\right]\right)\Big|_{\beta=1}$$

Taylor expansion around $\beta = 1$ :

$$\log p(x) \approx \boxed{E_{q_\lambda}[-V]} + \frac{1}{2}\left[(V - E_{q_\lambda}[-V])^2\right] - \frac{1}{3!}\left[(V - E_{q_\lambda}[-V])^3\right]$$
$$+ \frac{1}{4!}\left[(V - E_{q_\lambda}[-V])^4\right] - \ \dots$$

$$E_{q_\lambda}[-V(x,z)] = E_{q_\lambda}[\log p(x,z) - \log q_\lambda(z)]$$

Bamler*, **Zhang*** et al. Perturbative Black Box Variational Inference. NeurIPS 2017

# Revisiting Perturbation Theory for VI

$$V(x,z) \equiv \log q_\lambda(z) - \log p(x,z)$$

$$\log p(x) = \log\left(E_{z \sim q_\lambda}\left[\frac{p(x,z)}{q_\lambda(z)}\right]\right) = \log\left(E_{z \sim q_\lambda}\left[e^{-\beta V(x,z)}\right]\right)\Big|_{\beta=1}$$

Taylor expansion around $\beta = 1$ :

$$\log p(x) \approx \boxed{E_{q_\lambda}[-V]} + \frac{1}{2}\left[(V - E_{q_\lambda}[-V])^2\right] - \frac{1}{3!}\left[(V - E_{q_\lambda}[-V])^3\right]$$
$$+ \frac{1}{4!}\left[(V - E_{q_\lambda}[-V])^4\right] - \dots$$

$$\boxed{E_{q_\lambda}[-V(x,z)] = E_{q_\lambda}[\log p(x,z) - \log q_\lambda(z)]}$$

Truncation at any odd number term provides a bound.

Bamler*, **Zhang*** et al. Perturbative Black Box Variational Inference. NeurIPS 2017

# Behaviour of PBBVI



- Better uncertainty estimation than KLVI

- Better bias-variance trade-off comparing to α-VI

Bamler*, **Zhang**\* et al. Perturbative Black Box Variational Inference. NeurIPS 2017

# Behaviour of PBBVI



- Better uncertainty estimation than KLVI

- Better bias-variance trade-off comparing to α-VI

Where do we truncate?
Is it flexible enough?

Bamler*, **Zhang**\* et al. Perturbative Black Box Variational Inference. NeurIPS 2017

# F-Divergence

$$D_f[p||q_\phi] = E_{\theta \sim q_\phi}[f\left(\frac{p(\theta)}{q_\phi(\theta)}\right) - f(1)]$$

# F-Divergence

$$D_f[p||q_\phi] = E_{\theta \sim q_\phi}\left[f\left(\frac{p(\theta)}{q_\phi(\theta)}\right) - f(1)\right]$$

$$f(t) = -\log t \qquad \longrightarrow \qquad KL(q||p)$$

$$f(t) = t\log t \qquad \longrightarrow \qquad KL(p||q)$$

$$f(t) = \frac{t^\alpha}{\alpha(\alpha - 1)} \qquad \longrightarrow \qquad D_\alpha(p||q)$$

Wang et.al. Variational Inference with Tail-adaptive f-Divergence. NeurIPS 2018
Wan et.al. f-Divergence Variational Inference. NeurIPS 2020

# Integral Probability Metric (IPM)

- Using a test function to describe difference:

$$D[q(z), p(z|x)] = \sup_{f \in F} |E_{q(z)}[f(z)] - E_{p(z|x)}[f(z)]|$$



Figure adapted, source: Dougal Sutherland

Gorham and Mackey. Measuring Sample Quality with Stein's Method. NeurIPS 2015
Ranganath et al. Operator Variational Inference. NeurIPS 2016
Liu and Wang. Stein Variational Gradient Descent: A General Purpose Bayesian Inference Algorithm. NeurIPS 2016

# Integral Probability Metric (IPM)

- Using a test function to describe difference:

$$D[q(z), p(z|x)] = \sup_{f \in F} |E_{q(z)}[f(z)] - E_{p(z|x)}[f(z)]|$$

- Stein discrepancy: only requires $z \sim q(z)$ and
$\nabla_z \log p(z|x) = \nabla_z \log p(z, x)$

$$S[q(z), p(z|x)] = \sup_{f \in F_q} |E_{q(z)}[\nabla_z \log p(z, x)^\top f(z) + \nabla_z^\top f(z)]|$$



Figure adapted, source: Dougal Sutherland

Gorham and Mackey. Measuring Sample Quality with Stein's Method. NeurIPS 2015
Ranganath et al. Operator Variational Inference. NeurIPS 2016
Liu and Wang. Stein Variational Gradient Descent: A General Purpose Bayesian Inference Algorithm. NeurIPS 2016

# Looking Back

VI with f divergence

VI with IPM

Perturbative VI

VI with Stein discrepancy

VI with α divergence

ELBO with KL

# How to Choose the Inference Algorithm?



VI with f divergence

VI with IPM

Perturbative VI

VI with α divergence

VI with Stein discrepancy

ELBO with KL

Choose divergence by meta-learning!

Zhang et al. Meta-Learning for Variational Inference. AABI 2019

# Improved Monte Carlo Bounds

- Importance weighted auto-encoder (IWAE) bound:

$$L_K(\phi) = E_{z_1,\ldots,z_K \sim\, q(z)}\left[\log\frac{1}{K}\sum_{k=1}^{K}\frac{p(x,z_k)}{q(z_k)}\right]$$

Importance sampling estimate of $p(x)$

Burda et al. Importance Weighted Auto-encoders. ICLR 2016
Naesseth et al. Variational Sequential Monte Carlo. AISTATS 2018
Maddison et al. Filtering Variational Objectives. NeurIPS 2017
Le et al. Auto-encoding Sequential Monte Carlo. ICLR 2018
Masrani et al. The Thermodynamic Variational Objective. NeurIPS 2019

# Improved Monte Carlo Bounds

- Importance weighted auto-encoder (IWAE) bound:

$$L_K(\phi) = E_{z_1,\ldots,z_K \sim q(z)}\left[\log \frac{1}{K}\sum_{k=1}^{K}\frac{p(x,z_k)}{q(z_k)}\right]$$

Importance sampling estimate of $p(x)$

$\log p(x)$ ——————

$L(\phi)$ ——————
$(K = 1)$

Burda et al. Importance Weighted Auto-encoders. ICLR 2016
Naesseth et al. Variational Sequential Monte Carlo. AISTATS 2018
Maddison et al. Filtering Variational Objectives. NeurIPS 2017
Le et al. Auto-encoding Sequential Monte Carlo. ICLR 2018
Masrani et al. The Thermodynamic Variational Objective. NeurIPS 2019

# Improved Monte Carlo Bounds

- Importance weighted auto-encoder (IWAE) bound:

$$L_K(\phi) = E_{z_1,\ldots,z_K \sim q(z)} \left[ \log \frac{1}{K} \sum_{k=1}^{K} \frac{p(x, z_k)}{q(z_k)} \right]$$

Importance sampling estimate of $p(x)$

$\log p(x)$ _____

$L_{k'}(\phi)$ _____
$(K = k' \geq k)$

increase $K$

$L_k(\phi)$ _____
$(K = k > 1)$

$L(\phi)$ _____
$(K = 1)$

Burda et al. Importance Weighted Auto-encoders. ICLR 2016
Naesseth et al. Variational Sequential Monte Carlo. AISTATS 2018
Maddison et al. Filtering Variational Objectives. NeurIPS 2017
Le et al. Auto-encoding Sequential Monte Carlo. ICLR 2018
Masrani et al. The Thermodynamic Variational Objective. NeurIPS 2019

# Improved Monte Carlo Bounds

- Importance weighted auto-encoder (IWAE) bound:

$$L_K(\phi) = E_{z_1,\ldots,z_K \sim q(z)}\left[\log \frac{1}{K}\sum_{k=1}^{K}\frac{p(x,z_k)}{q(z_k)}\right]$$

<span style="color:red">Importance sampling estimate of $p(x)$</span>

$\log p(x)$
$(K \rightarrow \infty)$

$L_{k'}(\phi)$
$(K = k' \geq k)$

increase $K$

$L_k(\phi)$
$(K = k > 1)$

$L(\phi)$
$(K = 1)$

Burda et al. Importance Weighted Auto-encoders. ICLR 2016
Naesseth et al. Variational Sequential Monte Carlo. AISTATS 2018
Maddison et al. Filtering Variational Objectives. NeurIPS 2017
Le et al. Auto-encoding Sequential Monte Carlo. ICLR 2018
Masrani et al. The Thermodynamic Variational Objective. NeurIPS 2019

# Improved Monte Carlo Bounds

- Importance weighted auto-encoder (IWAE) bound:

$$L_K(\phi) = E_{z_1,\ldots,z_K \sim q(z)} \left[ \log \frac{1}{K} \sum_{k=1}^{K} \frac{p(x, z_k)}{q(z_k)} \right]$$

Importance sampling estimate of $p(x)$

$\log p(x)$
$(K \to \infty)$

$L_{k'}(\phi)$
$(K = k' \geq k)$

$L_k(\phi)$
$(K = k > 1)$

$L(\phi)$
$(K = 1)$

increase $K$

Optimize
$q(z) \to p(z|x)$
to close the gap

Burda et al. Importance Weighted Auto-encoders. ICLR 2016
Naesseth et al. Variational Sequential Monte Carlo. AISTATS 2018
Maddison et al. Filtering Variational Objectives. NeurIPS 2017
Le et al. Auto-encoding Sequential Monte Carlo. ICLR 2018
Masrani et al. The Thermodynamic Variational Objective. NeurIPS 2019

# Improved Monte Carlo Bounds

- Constructing lower-bounds from an estimator $R$ of the marginal:

$$E_{q(h)}[R(h, x)] = p(x) \quad \Rightarrow \quad \underline{E_{q(h)}[\log R(h, x)] \leq \log p(x)}$$

<span style="color:red">Jensen's inequality</span>

Domke and Sheldon. Divide and Couple: Using Monte Carlo Variational Objectives for Posterior Approximation. NeurIPS 2019

# Improved Monte Carlo Bounds

- Constructing lower-bounds from an estimator $R$ of the marginal:

$$E_{q(h)}[R(h,x)] = p(x) \quad \Rightarrow \quad E_{q(h)}[\log R(h,x)] \leq \log p(x)$$

<span style="color:red">Jensen's inequality</span>

- Variational lower-bound: $h = z, \ R(z,x) = \frac{p(x,z)}{q(z)}$

- IWAE bound: $h = (z_1, \ldots, z_K), R(h,x) = \frac{1}{K}\sum_{k=1}^{K}\frac{p(x,z_k)}{q(z_k)}$

Domke and Sheldon. Divide and Couple: Using Monte Carlo Variational Objectives for Posterior Approximation. NeurIPS 2019

# Improved Monte Carlo Bounds

- Constructing lower-bounds from an estimator $R$ of the marginal:

$$E_{q(h)}[R(h,x)] = p(x) \quad \Rightarrow \quad E_{q(h)}[\log R(h,x)] \leq \log p(x)$$

Jensen's inequality

  - Variational lower-bound: $h = z, \; R(z,x) = \frac{p(x,z)}{q(z)}$

  - IWAE bound: $h = (z_1, \ldots, z_K), R(h,x) = \frac{1}{K}\sum_{k=1}^{K}\frac{p(x,z_k)}{q(z_k)}$

- Fit $q$ using existing Monte Carlo estimators of $p(x)$

  - Example: antithetic sampling with Gaussian $q(z)$:

$$R(z,x) = \frac{p(x,z)+p(x,T(z))}{2q(z)}, \quad T(z) = \mu_q - (z - \mu_q)$$

★ MC sample
★ Antithetic sample

Domke and Sheldon. Divide and Couple: Using Monte Carlo Variational Objectives for Posterior Approximation. NeurIPS 2019

# Improved Monte Carlo Bounds

- Constructing lower-bounds from an estimator $R$ of the marginal:

$$E_{q(h)}[R(h,x)] = p(x) \quad \Rightarrow \quad \underline{E_{q(h)}[\log R(h,x)] \leq \log p(x)}$$

<span style="color:red">Jensen's inequality</span>

- Variational lower-bound: $h = z$, $R(z,x) = \frac{p(x,z)}{q(z)}$

- IWAE bound: $h = (z_1, \ldots, z_K)$, $R(h,x) = \frac{1}{K}\sum_{k=1}^{K}\frac{p(x,z_k)}{q(z_k)}$

- Fit $q$ using existing Monte Carlo estimators of $p(x)$

- Example: antithetic sampling with Gaussian $q(z)$:

$$R(z,x) = \frac{p(x,z)+p(x,T(z))}{2q(z)}, \quad T(z) = \mu_q - (z - \mu_q)$$



★ MC sample
★ Antithetic sample

Variational inference

$\mathbb{E}\log R = -0.237$
— $p(z,x)$
— $q(z)$, naive

Antithetic sampling

$\mathbb{E}\log R' = -0.060$
— $p(z,x)$
— $q(z)$, antithetic

Domke and Sheldon. Divide and Couple: Using Monte Carlo Variational Objectives for Posterior Approximation. NeurIPS 2019

# Free-energy as an Objective

- Bethe free-energy & message passing:



- Both $q$ and the inference algorithm are defined by the *factor graph*
- Optimal $q$ achieved at the fixed point of the *Bethe free energy*

Wainwright and Jordan. Graphical Models, Exponential Families, and Variational Inference. 2008.

**Li** and Turner. A Unifying Approximate Inference Framework from Variational Free Energy Relaxation. AABI 2016

# Landscape of Advances

q design

e.g. mean-field: $q(\theta) = \prod_i q(\theta_i)$

objective design

variational lower-bound:
$$L(\phi) = E_{q(\theta)}[\log p(D|\theta)] - KL[q(\theta)||p(\theta)]$$

# Landscape of Advances

mini-batch training:
$$\log p(D|\theta) \approx \frac{N}{M}\log p(B|\theta), B \sim D^M$$

$q$ design

stochastic optimization

scale-up

objective design

# Landscape of Advances



requires fast sampling of
$\theta \sim q(\theta)$

$\log p(D|\theta) \approx \dfrac{N}{M} \log p(B|\theta), B \sim D^M$

enable

Monte Carlo estimates

$q$ design

(another type of MC estimate)
$B \sim D^M$

stochastic optimization

estimate
$\theta \sim q(\theta)$

scale-up

objective design

$E_{q(\theta)}\left[\log \dfrac{p(D,\theta)}{q(\theta)}\right] \approx \dfrac{1}{K}\sum_{k}^{K} \log \dfrac{p(D,\theta_k)}{q(\theta_k)}, \theta_k \sim q(\theta)$

# Landscape of Advances

REINFORCE gradient:
$$\nabla_\phi E_{q_\phi(\theta)}[F(\theta)] = E_{q_\phi(\theta)}[F(\theta)\nabla_\phi \log q_\phi(\theta)]$$

reparam. trick:
$$\nabla_\phi E_{q_\phi(\theta)}[F(\theta)] = E_{p(\epsilon)}[\nabla_\phi F(g_\phi(\epsilon))]$$

# Landscape of Advances

# Landscape of Advances

# Landscape of Advances

# Landscape of Advances

# Landscape of Advances

# Part III: Applications

- Bayesian neural networks

- Generative models for decision making

- Future directions

# Why Estimating Uncertainty in DL?

- Models are often over-parameterised
  - E.g. BERT, GPT-3 in NLP
  - E.g. ResNet-152 for vision tasks

# Why Estimating Uncertainty in DL?

- Models are often over-parameterised
  - E.g. BERT, GPT-3 in NLP
  - E.g. ResNet-152 for vision tasks

ResNet-152          BERT                    GPT-3
(~60 million)   (~110 million)        (~175 billion)

# Why Estimating Uncertainty in DL?

- Models are often over-parameterised
  - E.g. BERT, GPT-3 in NLP
  - E.g. ResNet-152 for vision tasks
- Multiple parameter settings can fit the same data
  - They might provide different predictions on test data



ResNet-152     BERT                    GPT-3
(~60 million)  (~110 million)          (~175 billion)

# Why Estimating Uncertainty in DL?

- Critical tasks need uncertainty estimates to assist decision making
  - Inform end users when uncertain, for safe decision making



Healthcare AI

# Why Estimating Uncertainty in DL?

- Critical tasks need uncertainty estimates to assist decision making
    - Inform end users when uncertain, for safe decision making



Healthcare AI



Autonomous driving

# Bayesian Neural Network (BNN) 101

Classifying different types of animals:

- $x$: input image; $y$: output label
- Build a neural network with parameters $\theta$:
$$p(y|x,\theta) = softmax(f_\theta(x))$$

# Bayesian Neural Network (BNN) 101

Classifying different types of animals:
- $x$: input image; $y$: output label
- Build a neural network with parameters $\theta$:
    $$p(y|x,\theta) = softmax(f_\theta(x))$$



A typical neural network (with non-linearity $g(\cdot)$):

$$f_\theta(x) = W^L g\big(W^{L-1}\, g\big(\dots g(W^1\, x + b^1)\big) + b^{L-1}\big) + b^L,$$

$$h^l = g(W^l\, h^{l-1} + b^l),\; h^1 = g(W^1\, x + b^1).$$

Neural network parameters: $\theta = \{W^l, b^l\}_{l=1}^{L}$

# Bayesian Neural Network (BNN) 101

Classifying different types of animals:
- $x$: input image; $y$: output label
- Build a neural network with parameters $\theta$:
$$p(y|x,\theta) = softmax(f_\theta(x))$$



Typical deep learning solution:
- Optimize $\theta$ to obtain a point estimates (MLE):

$$\theta^* = argmax \ \log p(D \mid \theta),$$
$$\log p(D \mid \theta) = \sum_{n=1}^{N} \log p(y_n \mid x_n, \theta), D = \{(x_n, y_n)\}_{n=1}^{N}$$

- Prediction: using $p(y^* \mid x^*, \theta^*)$

# Bayesian Neural Network (BNN) 101

Classifying different types of animals:
- $x$: input image; $y$: output label
- Build a neural network with parameters $\theta$:

$$p(y|x, \theta) = softmax(f_\theta(x))$$



$p(\text{``cat''}|x, \mathcal{D})$

Bayesian solution:
- Put a prior $p(\theta)$ on network parameters $\theta$, e.g. Gaussian prior

$$p(\theta) = N(\theta; 0, \sigma^2 I)$$

- Compute the posterior distribution $p(\theta \mid D)$:

$$p(\theta \mid D) \propto p(D \mid \theta) \, p(\theta)$$

- Bayesian predictive inference:

$$p(y^* \mid x^*, D) = E_{p(\theta \mid D)}[p(y^* \mid x^*, \theta)]$$

# Bayesian Neural Network (BNN) 101

Classifying different types of animals:
- $x$: input image; $y$: output label
- Build a neural network with parameters $\theta$:
$$p(y|x,\theta) = softmax(f_\theta(x))$$



$p(\text{"cat"}|x,\mathcal{D})$

**Approximate (Bayesian) inference solution:**
- Exact posterior intractable, use approximate posterior:
$$q(\theta) \approx p(\theta \mid D)$$
- Approximate Bayesian predictive inference:
$$p(y^* \mid x^*, D) \approx E_{q(\theta)}[p(y^* \mid x^*, \theta)]$$
- Monte Carlo approximation:
$$p(y^* \mid x^*, D) \approx \frac{1}{K} \sum_{k=1}^{K} p(y^* \mid x^*, \theta_k), \quad \theta_k \sim q(\theta)$$

# Bayesian Neural Network (BNN) 101



Prediction on in-distribution data:
ensemble over networks, using weights sampled from $q(\theta)$

# Bayesian Neural Network (BNN) 101



Prediction on OOD/noisy/adversarial data:

Disagreement (i.e. uncertainty) exists over networks sampled from $q(\theta)$

# Bayesian Neural Network (BNN) 101



Prediction on OOD/noisy/adversarial data when $q(\theta)$ is over-confident:
Return confidently wrong answers (close to point estimate)

# Bayesian Neural Network (BNN) 101



Prediction on in-distribution data when $q(\theta)$ is under-confident:
Low accuracy in prediction tasks (less desirable)

# Approximate Inference in BNNs

- Key steps of approximate inference in BNNs
  1. Construct the $q(\theta) \approx p(\theta \mid D)$ distribution
     - Simple distributions: e.g. Mean-field Gaussian
     - Structured approximations, e.g. low-rank Gaussians
     - Others (non-Gaussian)

# Approximate Inference in BNNs

- Key steps of approximate inference in BNNs
  1. Construct the $q(\theta) \approx p(\theta \mid D)$ distribution
     - Simple distributions: e.g. Mean-field Gaussian
     - Structured approximations, e.g. low-rank Gaussians
     - Others (non-Gaussian)
  2. Fit the $q(\theta)$ distribution
     - E.g. with variational inference

# Approximate Inference in BNNs

- Key steps of approximate inference in BNNs
    1. Construct the $q(\theta) \approx p(\theta \mid D)$ distribution
        - Simple distributions: e.g. Mean-field Gaussian
        - Structured approximations, e.g. low-rank Gaussians
        - Others (non-Gaussian)
    2. Fit the $q(\theta)$ distribution
        - E.g. with variational inference
    3. Compute prediction with Monte Carlo approximations

# Approximate Inference in BNNs

- Step 1: construct the $q(\theta) \approx p(\theta \mid D)$ distribution
  - Example: Mean-field Gaussian distribution:

$$q(\theta) = \prod_{l=1}^{L} q(W^l)\, q(b^l)$$

$$q(W_l) = \prod_{ij} q(W_{ij}^l), \qquad q(W_{ij}^l) = N(W_{ij}^l; M_{ij}^l, V_{ij}^l)$$

$$q(b^l) = \prod_i q(b_i^l), \quad q(b_i^l) = N(b_i^l; m_i^l, v_i^l)$$

  - Variational parameters: $\phi = \left\{ M_{ij}^l, \log V_{ij}^l, m_i^l, \log v_i^l \right\}_{l=1}^{L}$

Blundell et al. Weight Uncertainty in Neural Networks. ICML 2015

# Approximate Inference in BNNs

- Step 2: fit the $q(\theta)$ distribution:
  - Variational inference: $\phi^* = argmax \, L(\phi)$
$$L(\phi) = E_{q(\theta)}[\log p(D \mid \theta)] - KL[q(\theta) \mid\mid p(\theta)]$$

Blundell et al. Weight Uncertainty in Neural Networks. ICML 2015

# Approximate Inference in BNNs

- Step 2: fit the $q(\theta)$ distribution:
  - Variational inference: $\phi^* = argmax\ L(\phi)$
  $$L(\phi) = E_{q(\theta)}[\log p(D \mid \theta)] - KL[q(\theta) \mid\mid p(\theta)]$$
  - First scalable technique: <span style="color:red">Stochastic optimization</span>
    - i.i.d. assumption of data: $\log p(D \mid \theta) = \sum_{n=1}^{N} \log p(y_n \mid x_n, \theta)$
    - Enable mini-batch training with $\{(x_m, y_m)\} \sim D^M$ :
    $$L(\phi) \approx \frac{N}{M} \sum_{m=1}^{M} E_{q(\theta)}[\log p(y_m \mid x_m, \theta)] - KL[q(\theta) \mid\mid p(\theta)]$$

Blundell et al. Weight Uncertainty in Neural Networks. ICML 2015

# Approximate Inference in BNNs

- Step 2: fit the $q(\theta)$ distribution:
  - Variational inference: $\phi^* = argmax\ L(\phi)$
    $$L(\phi) = E_{q(\theta)}[\log p(D \mid \theta)] - KL[q(\theta) \mid\mid p(\theta)]$$
  - First scalable technique: <span style="color:red">Stochastic optimization</span>
    - i.i.d. assumption of data: $\log p(D \mid \theta) = \sum_{n=1}^{N} \log p(y_n \mid x_n, \theta)$
    - Enable mini-batch training with $\{(x_m, y_m)\} \sim D^M$ :
    $$L(\phi) \approx \boxed{\frac{N}{M}} \sum_{m=1}^{M} E_{q(\theta)}[\log p(y_m \mid x_m, \theta)] - KL[q(\theta) \mid\mid p(\theta)]$$

<span style="color:red">reweighting to ensure calibrated posterior concentration</span>

Blundell et al. Weight Uncertainty in Neural Networks. ICML 2015

# Approximate Inference in BNNs

- Step 2: fit the $q(\theta)$ distribution:
  - 2nd scalable technique: Monte Carlo sampling
    - $E_{q(\theta)}[\log p(y \mid x, \theta)]$ intractable even with Gaussian $q(\theta)$
    - Solution: Monte Carlo estimate:

$$E_{q(\theta)}[\log p(y \mid x, \theta)] \approx \frac{1}{K} \sum_{k}^{K} \log p(y \mid x, \theta_k), \qquad \theta_k \sim q(\theta)$$

Blundell et al. Weight Uncertainty in Neural Networks. ICML 2015

# Approximate Inference in BNNs

- Step 2: fit the $q(\theta)$ distribution:
  - 2nd scalable technique: <span style="color:red">Monte Carlo sampling</span>
    - $E_{q(\theta)}[\log p(y \mid x, \theta)]$ intractable even with Gaussian $q(\theta)$
    - <span style="color:red">Solution: Monte Carlo estimate:</span>

$$E_{q(\theta)}[\log p(y \mid x, \theta)] \approx \frac{1}{K}\sum_{k}^{K} \log p(y \mid x, \theta_k), \qquad \theta_k \sim q(\theta)$$

   - <span style="color:red">Reparameterization trick</span> to sample mean-field Gaussians:
$$\theta_k \sim q(\theta) \Leftrightarrow \theta_k = m_\theta + \sigma_\theta \, \epsilon_k, \ \epsilon_k \sim N(0, I)$$



Blundell et al. Weight Uncertainty in Neural Networks. ICML 2015

# Approximate Inference in BNNs

- Step 2: fit the $q(\theta)$ distribution:
  - 2nd scalable technique: Monte Carlo sampling
    - $E_{q(\theta)}[\log p(y \mid x, \theta)]$ intractable even with Gaussian $q(\theta)$
    - Solution: Monte Carlo estimate:

$$E_{q(\theta)}[\log p(y \mid x, \theta)] \approx \frac{1}{K} \sum_k^K \log p(y \mid x, \theta_k), \qquad \theta_k \sim q(\theta)$$

   - Reparameterization trick to sample mean-field Gaussians:

$$\theta_k \sim q(\theta) \Leftrightarrow \theta_k = m_\theta + \sigma_\theta \, \epsilon_k, \; \epsilon_k \sim N(0, I)$$

$$\Rightarrow E_{q(\theta)}[\log p(y \mid x, \theta)] \approx \frac{1}{K} \sum_k^K \log p(y \mid x, m_\theta + \sigma_\theta \epsilon_k), \epsilon_k \sim N(0, I)$$



Blundell et al. Weight Uncertainty in Neural Networks. ICML 2015

# Approximate Inference in BNNs

- Combining both steps:

$$L(\phi) \approx \frac{N}{M} \sum_{m=1}^{M} \frac{1}{K} \sum_{k=1}^{K} \log p(y_m \mid x_m, \theta_k) - \underline{KL[q(\theta) \mid\mid p(\theta)]}, \theta_k \sim q(\theta)$$

<span style="color:red">analytic between two Gaussians</span>
(if not, can also be estimated with Monte Carlo)

Blundell et al. Weight Uncertainty in Neural Networks. ICML 2015

# Approximate Inference in BNNs

- Step 3: compute prediction with Monte Carlo approximations:

$$p(y^* \mid x^*, D) \approx \frac{1}{K} \sum_{k=1}^{K} p(y^* \mid x^*, \theta_k), \quad \underline{\theta_k \sim q(\theta)}$$

<span style="color:red">Mean-field Gaussian case:
$$\theta_k = m_\theta + \sigma_\theta \, \epsilon_k, \ \epsilon_k \sim N(0, I)$$</span>



Blundell et al. Weight Uncertainty in Neural Networks. ICML 2015

# Applications of BNNs: Image Segmentation



(a) Input Image (b) Ground Truth (c) Semantic Segmentation (d) Aleatoric Uncertainty (e) Epistemic Uncertainty

Kendall and Gal. What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision? NeurIPS 2017

# Applications of BNNs: Super Resolution



Tanno et al. Uncertainty Quantification in Deep Learning for Safer Neuroimage Enhancement. Neuroimage 2020

# Applications of BNNs: Continual Learning



$t - 1$      $t$      $t + 1$

cat or dog?     dog or monkey?     monkey or owl?

learn

memorise      transfer

Update posterior belief with the current task

Posterior from the previous tasks as prior for the current task

$$L_{VCL}^t\big(q_t(\theta)\big) = E_{q_{t(\theta)}}[\log p(D_t \mid \theta)] \; - KL[\,q_t(\theta)\,||\,q_{t-1}(\theta)\,]$$

Nguyen et al. Variational Continual Learning. ICLR 2018
Pan et al. Continual Deep Learning by Functional Regularisation of Memorable Past. NeurIPS 2020

# Recent Progress in BNNs: Inference

SGD: $\quad \theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \widetilde{U}(\theta_t)$

SGLD: $\quad \theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \widetilde{U}(\theta_t) + \sqrt{2\eta}\epsilon, \qquad \epsilon \sim N(0, I)$

Stochastic gradient MCMC

Li et al. Preconditioned Stochastic Gradient Langevin Dynamics for Deep Neural Networks. AAAI 2016
Zhang et al. Cyclical Stochastic Gradient MCMC for Bayesian Deep Learning. ICLR 2020

# Recent Progress in BNNs: Inference

SGD: $\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \tilde{U}(\theta_t)$

SGLD: $\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \tilde{U}(\theta_t) + \sqrt{2\eta}\epsilon, \qquad \epsilon \sim N(0, I)$

Stochastic gradient MCMC



Monte Carlo dropout

Gal and Ghahramani. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. ICML 2016

# Recent Progress in BNNs: Inference

SGD:  $\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \tilde{U}(\theta_t)$

SGLD:  $\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \tilde{U}(\theta_t) + \sqrt{2\eta}\epsilon,$   $\epsilon \sim N(0, I)$

Stochastic gradient MCMC



Monte Carlo dropout



Deterministic approximations

Hernandez-Lobato and Adams. Probabilistic Backpropagation for Scalable Learning of Bayesian Neural Networks. ICML 2015
Wu et al. Deterministic Variational Inference for Robust Bayesian Neural Networks. ICLR 2019

# Recent Progress in BNNs: Inference

SGD:    $\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \widetilde{U}(\theta_t)$

SGLD:   $\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \widetilde{U}(\theta_t) + \sqrt{2\eta}\epsilon, \qquad \epsilon \sim N(0, I)$

Stochastic gradient MCMC



Deterministic approximations



Monte Carlo dropout



Function space approximate inference

Ma et al. Variational Implicit Processes. ICML 2019

Sun et al. Functional Variational Bayesian Neural Networks. ICLR 2019

# Recent Progress in BNNs: Theory



## Connections to GPs:

- BNN with very wide hidden layers
  $\approx$ Gaussian process
- Width limit convergence: in both
  prior (Neal's result) and posterior

Neal. Bayesian Learning for Neural Networks. PhD Thesis, 1996
Matthews et al. Gaussian Process Behaviour in Wide Deep Neural Networks. ICLR 2018
Lee et al. Deep Neural Networks as Gaussian Processes. ICLR 2018
Hron et al. Exact posterior distributions of wide Bayesian neural networks. 2020

# Recent Progress in BNNs: Theory



## Connections to GPs:

- BNN with very wide hidden layers ≈ Gaussian process
- Width limit convergence: in both prior (Neal's result) and posterior

## Approx. vs exact inference:

- Theoretical limitation of MFVI in shallow BNNs with ReLU activations
- Empirically deep BNNs with MVFI still fails in certain cases

Foong et al. On the Expressiveness of Approximate Inference in Bayesian Neural Networks. NeurIPS 2020

Farquhar et al. Liberty or Depth: Deep Bayesian Neural Nets Do Not Need Complex Weight Posterior Approximations. NeurIPS 2020

# Dynamic Information Acquisition



Doctor Li

Ma et al. Eddi: Efficient dynamic discovery of high-value information with partial VAE. ICML 2019

# Dynamic Information Acquisition



Neu: Neuropathic pain
Noc: Nociceptive pain
Pys: Psychiatric pain

Ma et al. Eddi: Efficient dynamic discovery of high-value information with partial VAE. ICML 2019

# Dynamic Information Acquisition



Step 1: Missing Value Prediction with Deep Generative Model

Ma et al. Eddi: Efficient dynamic discovery of high-value information with partial VAE. ICML 2019

# Dynamic Information Acquisition



Step 1: Missing Value Prediction with Deep Generative Model

Step 2: Active element-wise information acquisition

Ma et al. Eddi: Efficient dynamic discovery of high-value information with partial VAE. ICML 2019

# A Deep Generative Model



Variational Auto-encoder (VAE)

$$L_{amortized} = \log p(\boldsymbol{x}) - KL\left(q(\boldsymbol{z}\,|\,\boldsymbol{x}\,)|p(\boldsymbol{z}|\boldsymbol{x})\right)$$

Encoder/
inference network

Decoder/
generator

# Challenges



- Utilizing Deep Learning for scalable inference
- Cannot handle partial observation

# VAE to Partial VAE



Encoder/
inference network

Decoder/
generator

Ma et al. Eddi: Efficient dynamic discovery of high-value information with partial VAE. ICML 2019

# VAE to Partial VAE



$X$ is fully observed.

$X$ is partially observed.

**We aim to infer the missing values $X_U$ from the observed values $X_O$**

Ma et al. Eddi: Efficient dynamic discovery of high-value information with partial VAE. ICML 2019

# VAE to Partial VAE



$$L_{amortized} = \log p(\boldsymbol{x}) - KL[q(\boldsymbol{z}|\boldsymbol{x})||p(\boldsymbol{z}|\boldsymbol{x})]$$
$$= E_{\boldsymbol{z} \sim q(\boldsymbol{z}|\boldsymbol{x})}[\log p_\theta(\boldsymbol{x}|\boldsymbol{z})] - KL[q(\boldsymbol{z}|\boldsymbol{x})||p(\boldsymbol{z})]$$

$$L_{amortized} = \log p(\boldsymbol{x_o}) - KL[q(\boldsymbol{z}|\boldsymbol{x_o})||p(\boldsymbol{z}|\boldsymbol{x_o})]$$
$$= E_{\boldsymbol{z} \sim q(\boldsymbol{z}|\boldsymbol{x_o})}[\log p_\theta(\boldsymbol{x}|\boldsymbol{z})] - KL[q(\boldsymbol{z}|\boldsymbol{x_o})||p(\boldsymbol{z})]$$

The ELBO still holds.
The challenge is how to design an inference net.

Ma et al. Eddi: Efficient dynamic discovery of high-value information with partial VAE. ICML 2019

# VAE to Partial VAE



Encoder/
inference network

Decoder/
generator

$$\mathbf{c}(\mathbf{x}_O) := g(\vec{h}(\mathbf{s}_1), \vec{h}(\mathbf{s}_2), ..., \vec{h}(\mathbf{s}_O))$$

**Set Encoder**

Qi et al. Pointnet: Deep learning on point sets for 3d classification and segmentation. CVPR 2017
Zaheer et al. Deep sets. NeurIPS 2017

# VAE to Partial VAE



Encoder/
inference network

Decoder/
generator

$$\mathbf{c}(\mathbf{x}_O) := g(\vec{h}(\mathbf{s}_1), \vec{h}(\mathbf{s}_2), ..., \vec{h}(\mathbf{s}_O))$$

**Set Encoder**

| | A | B | C | D | E | F | |
|---|---|---|---|---|---|---|---|
| | | cough | cold | temperature | marlaria | age | ... |
| | Alice | 1 | 1 | | | 25 | |
| | Mary | | | | 0 | | |
| | Kevin | 0 | | | | 47 | |

Ma et al. Eddi: Efficient dynamic discovery of high-value information with partial VAE. ICML 2019

# Dynamic Information Acquisition

Yes No | Yes No | Yes No | Neu Noc Pys

Inflammation | Lower Back Pain | Have pats

Step 1: Missing Value Prediction with Deep Generative Model

Step 2: Active element-wise information acquisition

Ma et al. Eddi: Efficient dynamic discovery of high-value information with partial VAE. ICML 2019

# Actively Select the Next Variable



$X_i = ?$

$X_i = !$

$$i^* = \arg\max_{i \in U} R(i | \boldsymbol{x_o})$$

Ma et al. Eddi: Efficient dynamic discovery of high-value information with partial VAE. ICML 2019

# Actively Select the Next Variable



$$R(i, x_o) = E_{x_i \sim p(x_i|x_o)}[KL[p(x_\phi| \, x_i, x_o)||p(x_\phi|x_o)]]$$

X_i=?

X_i=!

$$i^* = \arg\max_{i \in U} R(i|\boldsymbol{x_o})$$

EDDI

User

Ma et al. Eddi: Efficient dynamic discovery of high-value information with partial VAE. ICML 2019

# Actively Select the Next Variable



$$R(i, x_o) = E_{x_i \sim p(x_i|x_o)}[KL[p(x_\phi| \; x_i, x_o)||p(x_\phi|x_o)]]$$

X$_i$=?

Intractable

User

EDDI

Efficient Approximation

X$_i$=!

$$i^* = \arg\max_{i \in U} R(i|\boldsymbol{x_o})$$

Ma et al. Eddi: Efficient dynamic discovery of high-value information with partial VAE. ICML 2019

# Predicting a House's Price



Kim (customer)

How far away is the supermarket?

Ty (broker)

MODEL

Boston - Media Home Value

Refresh

GROUND TRUTH

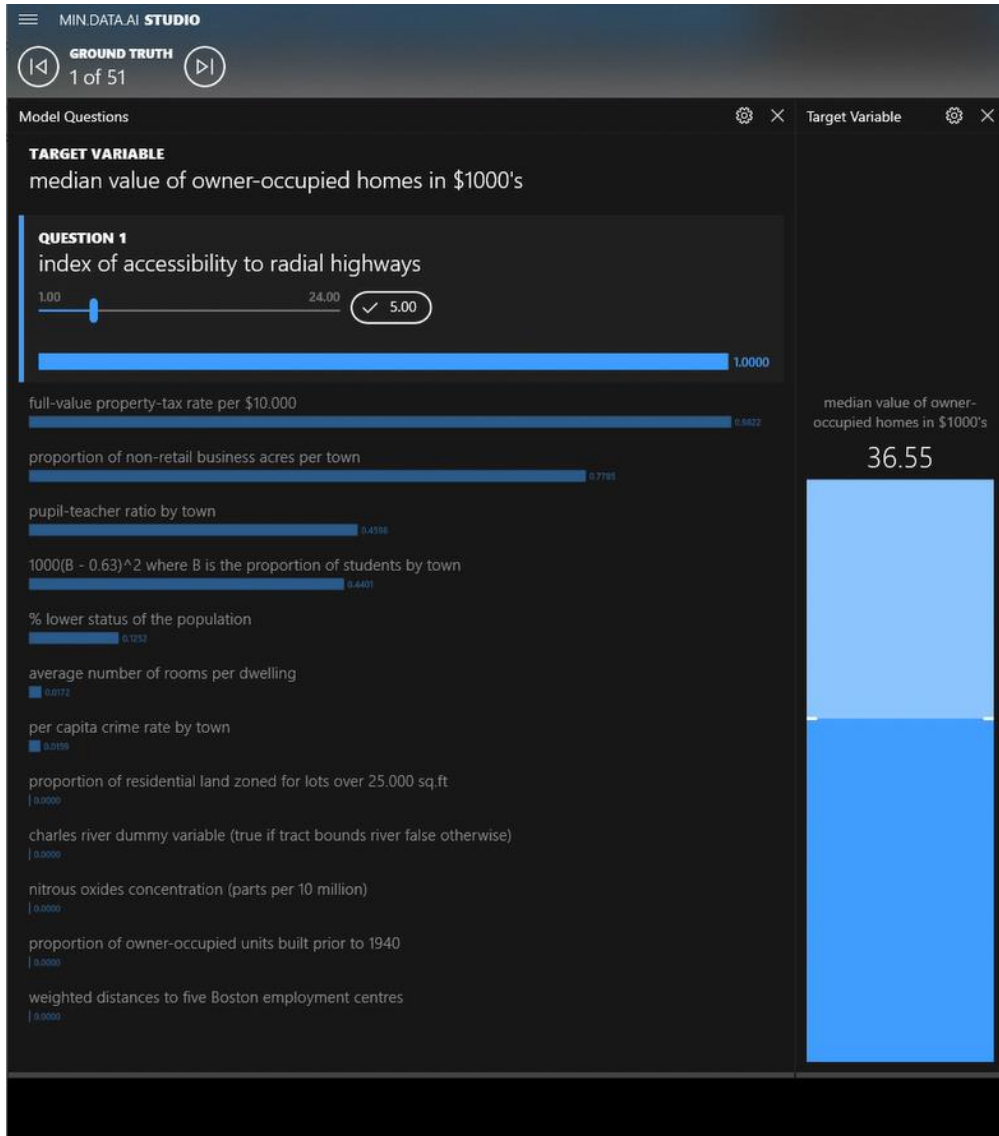C:\Users\mgrayson\Desktop\min-data-ai\boston\test.csv

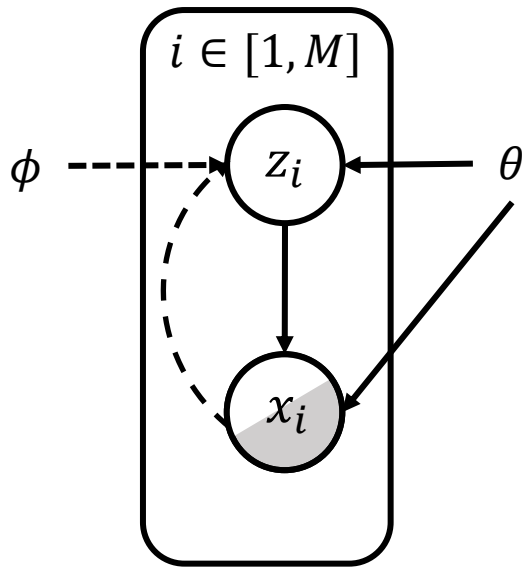Browse...

PANELS

✓ Model questions

⤬ Random questions

Tools

**Model Questions**

**TARGET VARIABLE**
median value of owner-occupied homes in $1000's

**QUESTION 1**
index of accessibility to radial highways

1.00    24.00    ✓ 5.00

1.0000

full-value property-tax rate per $10.000
0.9822

proportion of non-retail business acres per town
0.7785

pupil-teacher ratio by town
0.4556

1000(B - 0.63)^2 where B is the proportion of students by town
0.4401

% lower status of the population
0.1252

average number of rooms per dwelling
0.0172

per capita crime rate by town
0.0159

proportion of residential land zoned for lots over 25.000 sq.ft
0.0000

charles river dummy variable (true if tract bounds river false otherwise)
0.0000

nitrous oxides concentration (parts per 10 million)
0.0000

proportion of owner-occupied units built prior to 1940
0.0000

weighted distances to five Boston employment centres
0.0000

**Target Variable**

median value of owner-occupied homes in $1000's

36.55

- List of questions from survey (e.g. US veteran) for mental health monitoring

- List of technical questions from interview for recruiting
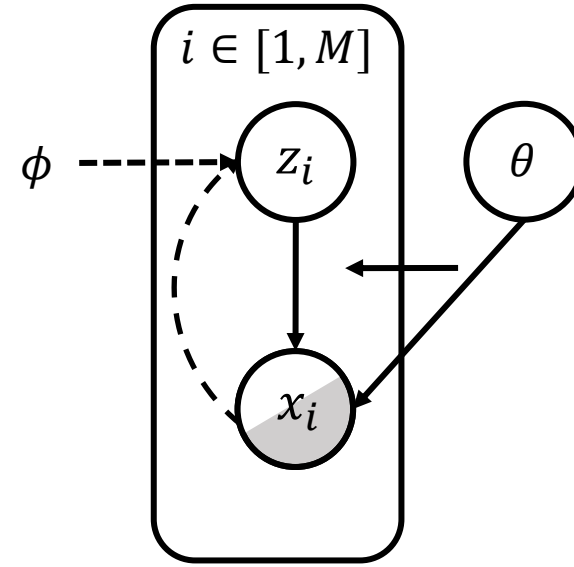
- List of medical tests for diagnosis

- … …

# Does it work when there is few training data?

# Partial Amortized Bayesian Deep Latent Gaussian Model (PA-BELGAM)



- Point estimate of global parameter $\theta$

$$p(\boldsymbol{x_o}, \boldsymbol{z}) = \prod_{i=1}^{N} \prod_{d \in O_i} p(x_{i,d}|z_i)p(z_i)$$

- Stochastic variable $\theta$

$$p(\boldsymbol{x_o}, \theta, \boldsymbol{z}) = p(\theta) \prod_{i=1}^{N} \prod_{d \in O_i} p(x_{i,d}|z_i, \theta)p(z_i)$$

Gong et al. Icebreaker: Element-wise efficient information acquisition with a Bayesian deep latent gaussian model. NeurIPS 2019

# Partial Amortized Bayesian Deep Latent Gaussian Model (PA-BELGAM)

Amortized inference for local latent variables

SGHMC for global latent variables



$$q(\theta, \boldsymbol{z} \,|\, \boldsymbol{x_o}) \approx q(\theta | \boldsymbol{x_o}) q_\phi\,(\boldsymbol{z} | \boldsymbol{x_o})$$

Gong et al. Icebreaker: Element-wise efficient information acquisition with a Bayesian deep latent gaussian model. NeurIPS 2019

# When Data are Heterogenous



Ma et.al. VAEM: a Deep Generative Model for Heterogeneous Mixed Type Data. NeurIPS 2020

# When Robustness is Needed



Causal Reasoning

Deep Casual
Manipulation
Augmented Model

**Zhang** et.al. A Causal View on Robustness of Neural Networks. NeurIPS 2020
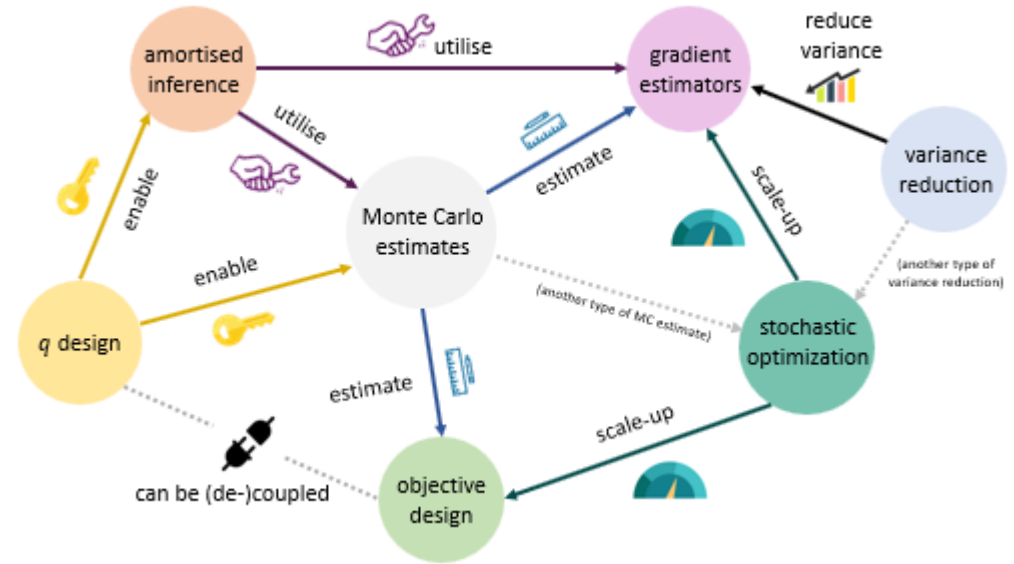
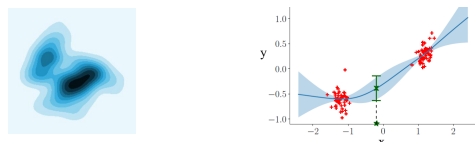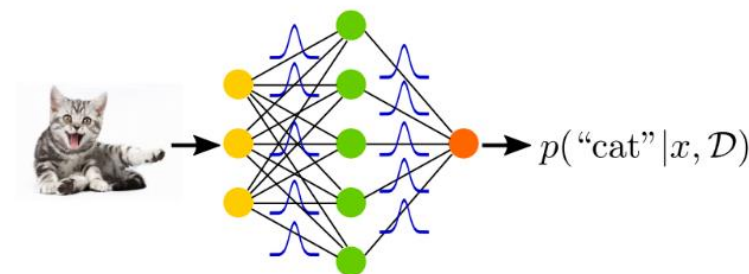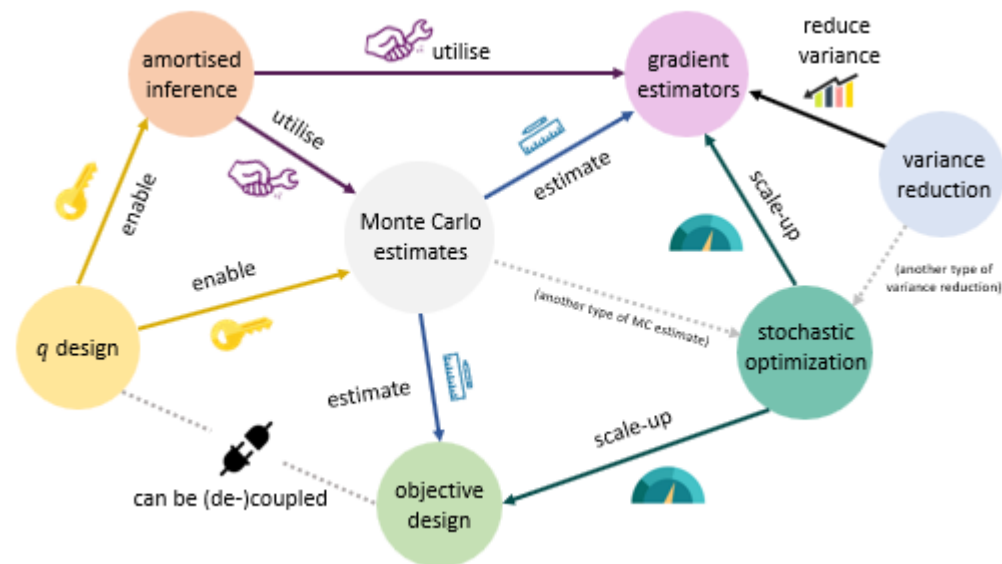# Summary

The probability distribution $\pi(\theta)$ is intractable



Approximate Inference
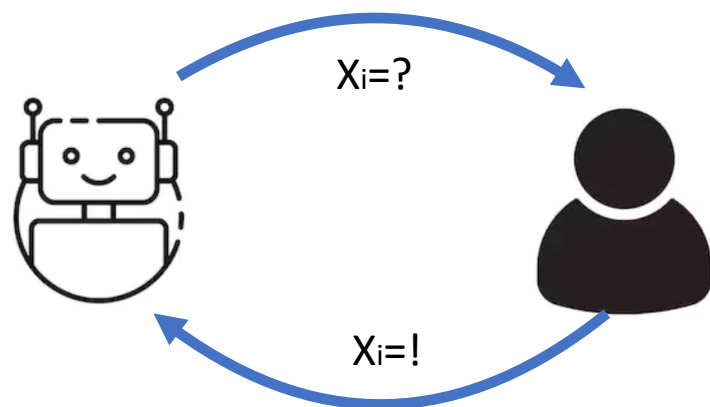
$$\int F(\theta)\,\pi(\theta)\,d\theta$$

# Summary

The probability distribution $\pi(\theta)$ is intractable

Approximate Inference

$$\int F(\theta)\,\pi(\theta)\,d\theta$$

# Summary

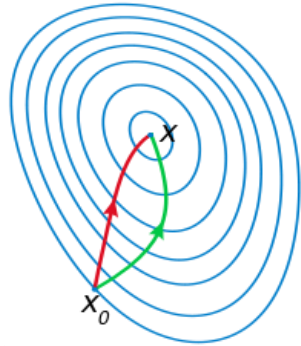The probability distribution $\pi(\theta)$ is intractable



Approximate Inference

$$\int F(\theta)\, \pi(\theta)\, d\theta$$



$X_i = ?$

$X_i = !$

amortised inference

utilise

gradient estimators

reduce variance

utilise

Monte Carlo estimates

estimate

variance reduction

enable

key

enable

$q$ design

scale-up

(another type of variance reduction)

estimate

(another type of MC estimate)

stochastic optimization

can be (de-)coupled

objective design

scale-up

$p(\text{"cat"}|x, \mathcal{D})$

# Future Directions: Methodology

Better optimization

Zhang et al. Noisy natural gradient as variational inference. ICML 2018
Khan et al. Fast and Scalable Bayesian Deep Learning by Weight-Perturbation in Adam. ICML 2018

# Future Directions: Methodology

Better optimization

Combined approaches

*q* distribution design
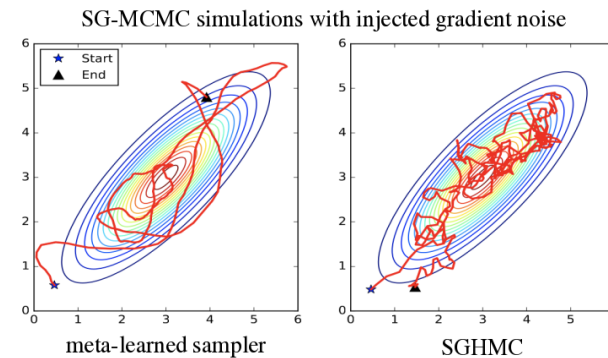objective design
amortised inference
scalable inference

...

rejection sampling
importance sampling
SMC, MCMC, Quasi MC
...



*q* distribution design
objective design
amortised inference
scalable inference
...

Burda et al. Importance Weighted Auto-encoders. ICLR 2016
Domke and Sheldon. Divide and Couple: Using Monte Carlo Variational Objectives for Posterior Approximation. NeurIPS 2019

# Future Directions: Methodology

### Better optimization

*q* distribution design
objective design
amortised inference
scalable inference

...

### Combined approaches

rejection sampling
importance sampling
SMC, MCMC, Quasi MC
...



*q* distribution design
objective design
amortised inference
scalable inference
...

### Meta-learning inference algorithms



SG-MCMC simulations with injected gradient noise

meta-learned sampler          SGHMC

Gong et al. Meta-learning for Stochastic Gradient MCMC. ICLR 2019
Zhang et al. Meta-Learning for Variational Inference. AABI 2019

# Future Directions: Error Analyses

Errors in inference

$$D[q(\theta)||p(\theta|D)] = ?$$

$$D[q(y^*|x^*)||p(y^*|x^*,D)] = ?$$

$$= \int p(y^*|x^*,\theta) \, q(\theta) d\theta$$



Analysis needed for deep probabilistic models!
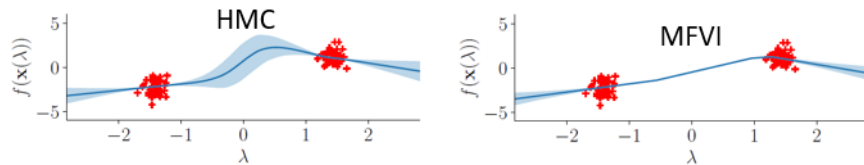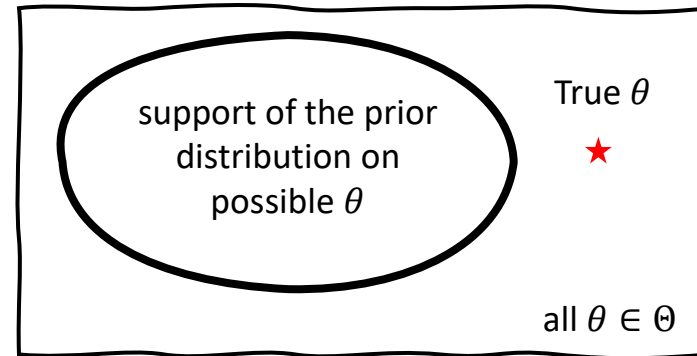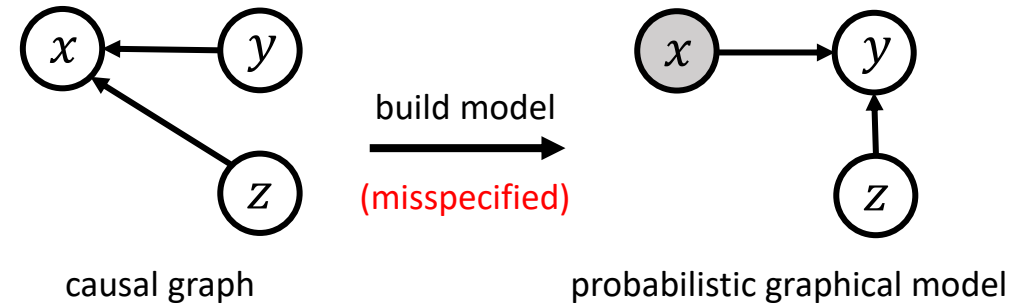
- Optimization error
- Approximation gap

Foong et al. On the Expressiveness of Approximate Inference in Bayesian Neural Networks. NeurIPS 2020

# Future Directions: Error Analyses

## Errors in inference

$$D[q(\theta)||p(\theta|D)] = ?$$

$$D[q(y^*|x^*)||p(y^*|x^*, D)] = ?$$

$$= \int p(y^*|x^*, \theta) \, q(\theta) d\theta$$



Analysis needed for deep probabilistic models!

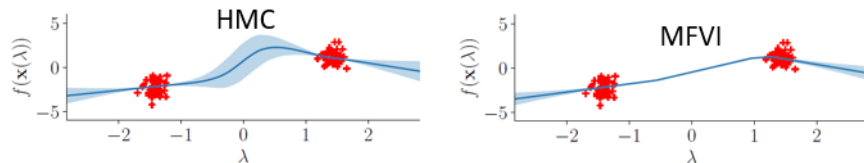- Optimization error
- Approximation gap

## Model misspecification



causal graph

build model

(misspecified)

probabilistic graphical model



support of the prior distribution on possible $\theta$

True $\theta$

★

all $\theta \in \Theta$

Wang and Blei. Variational Bayes under Model Misspecification. NeurIPS 2019

# Future Directions: Error Analyses

## Errors in inference

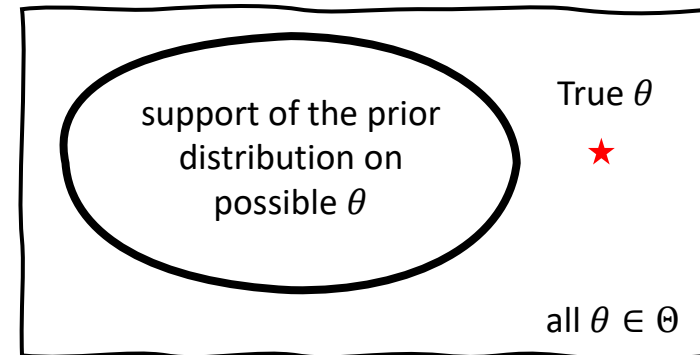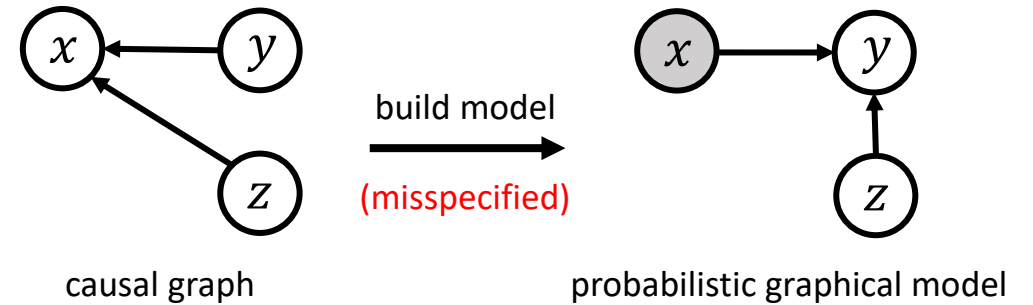$$D[q(\theta)||p(\theta|D)] = ?$$

$$D[\underline{q(y^*|x^*)}||p(y^*|x^*,D)] = ?$$

$$= \int p(y^*|x^*,\theta)\, q(\theta) d\theta$$



Analysis needed for <span style="color:red">deep</span> probabilistic models!
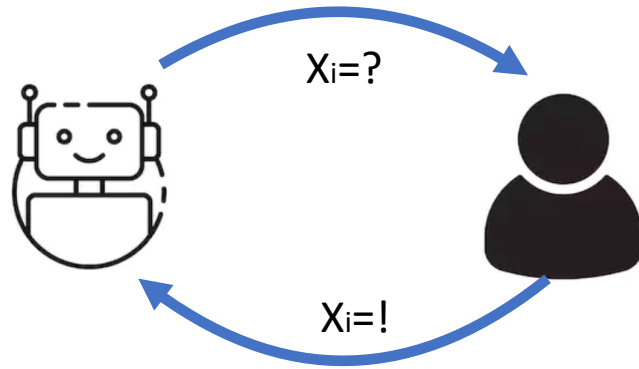- Optimization error
- Approximation gap

## Model misspecification



causal graph          probabilistic graphical model



Separation of inference & modelling?

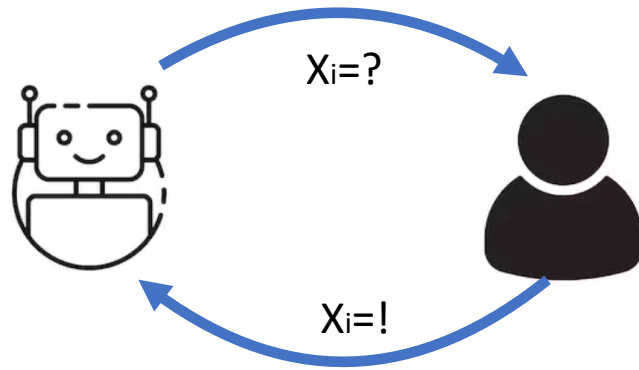Wang and Blei. Variational Bayes under Model Misspecification. NeurIPS 2019

# Future Directions: Applications

Uncertainty estimation



$X_i=?$

$X_i=!$

# Future Directions: Applications

Uncertainty estimation



Model selection & averaging



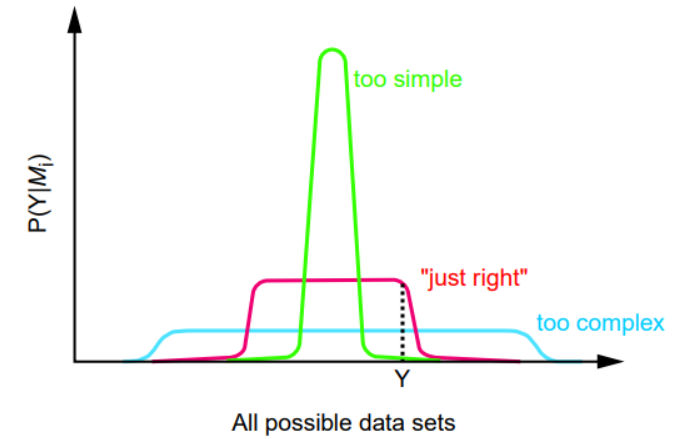Rasmussen and Ghahramani. Occam's Razor. NeurIPS 2001.
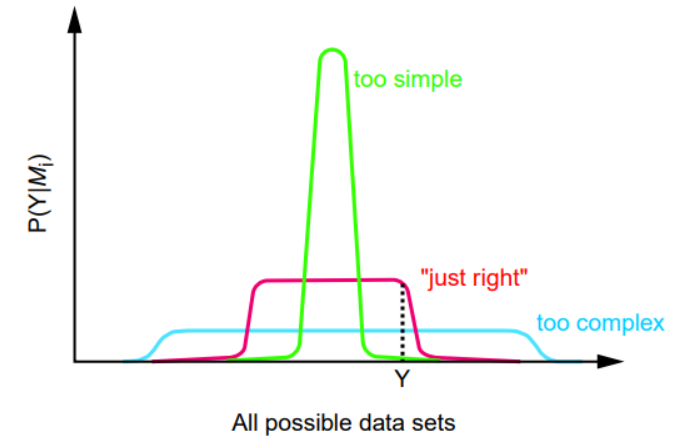
# Future Directions: Applications

Uncertainty estimation



Model selection & averaging



Causal reasoning



patient + drug A → recovery

"what if the patient was treated with drug B?"

# Thank You!

Questions? Ask at:
liyzhen2@gmail.com (Yingzhen Li)
Cheng.Zhang@microsoft.com (Cheng Zhang)