Machine Learning Reading Group

Presenter: Daniel long

Paper #1: Auto-encoding Variational Bayes (Kingma et al. (2014)) **Paper #2**: Variational Inference with Normalizing Flows (Rezende et al. (2015))

October 1, 2021

Overview

1. Background

Bayesian Inference/Latent variable modeling Variational Inference

- 2. Overview of contributions
- Paper #1
 Reparameterization trick
 Stochastic Gradient VB Estimators
 Auto-encoding VB Algorithm

Variational Auto-Encoder

- Paper #2 Normalizing flows VI algorithm
- 5. Comparison of papers
- 6. Related work

Overview

1. Background

Bayesian Inference/Latent variable modeling Variational Inference

- 2. Overview of contributions
- Paper #1 Reparameterization trick Stochastic Gradient VB Estimators Auto-encoding VB Algorithm Variational Auto-Encoder
- 4. Paper #2 Normalizing flows VI algorithm
- 5. Comparison of papers
- 6. Related work

Consider a latent variable model of data $\mathbf{x} = \{x_i\}$ and latent variables $\mathbf{z} = \{z_i\}$ with joint density

 $p(\mathbf{z}, \mathbf{x}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z}).$

Consider a latent variable model of data $\mathbf{x} = \{x_i\}$ and latent variables $\mathbf{z} = \{z_i\}$ with joint density

$$p(\mathbf{z}, \mathbf{x}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$$

Example: Gaussian Mixture model

$$p(z_i = k) = \pi_k, \ k = 1, \dots, K$$
$$x_i | z_i = k \sim N(\mu_k, \sigma_k^2)$$

Consider a latent variable model of data $\mathbf{x} = \{x_i\}$ and latent variables $\mathbf{z} = \{z_i\}$ with joint density

$$p(\mathbf{z}, \mathbf{x}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$$

Example: Gaussian Mixture model

$$p(z_i = k) = \pi_k, \ k = 1, \dots, K$$
$$x_i | z_i = k \sim N(\mu_k, \sigma_k^2)$$

Objective: Compute/sample from the posterior distribution $p_{\theta}(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$.

• $p(\mathbf{x})$ is usually difficult/impossible to compute.

Consider a latent variable model of data $\mathbf{x} = \{x_i\}$ and latent variables $\mathbf{z} = \{z_i\}$ with joint density

$$p(\mathbf{z}, \mathbf{x}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$$

Example: Gaussian Mixture model

$$p(z_i = k) = \pi_k, \ k = 1, \dots, K$$
$$x_i | z_i = k \sim N(\mu_k, \sigma_k^2)$$

Objective: Compute/sample from the posterior distribution $p_{\theta}(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$.

• $p(\mathbf{x})$ is usually difficult/impossible to compute.

Two solutions

1. MCMC: Construct an ergodic Markov chain whose stationary distribution is $p_{\theta}(\mathbf{z}|\mathbf{x})$.

Consider a latent variable model of data $\mathbf{x} = \{x_i\}$ and latent variables $\mathbf{z} = \{z_i\}$ with joint density

$$p(\mathsf{z},\mathsf{x}) = p(\mathsf{z})p(\mathsf{x}|\mathsf{z})$$

Example: Gaussian Mixture model

$$p(z_i = k) = \pi_k, \ k = 1, \dots, K$$
$$x_i | z_i = k \sim N(\mu_k, \sigma_k^2)$$

Objective: Compute/sample from the posterior distribution $p_{\theta}(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$.

• $p(\mathbf{x})$ is usually difficult/impossible to compute.

Two solutions

- 1. MCMC: Construct an ergodic Markov chain whose stationary distribution is $p_{\theta}(\mathbf{z}|\mathbf{x})$.
- 2. Variational Inference (VI): Compute an approximation of $p_{\theta}(\mathbf{z}|\mathbf{x})$ by optimizing over a family of approx. distributions.

Objective: Obtain approximate posterior (or recognition model)

$$q_{\phi}^{*}(\mathsf{z}|\mathsf{x}) = rgmin_{q \in \mathcal{Q}} D_{\mathcal{KL}}(q_{\phi}(\mathsf{z}|\mathsf{x})||p_{ heta}(\mathsf{z}|\mathsf{x})),$$

- Q is a (specified) variational family (e.g. mean-field)
- ϕ contain the variational parameters; θ contain the model parameters.

Objective: Obtain approximate posterior (or recognition model)

$$q_{\phi}^{*}(\mathsf{z}|\mathsf{x}) = rgmin_{q \in \mathcal{Q}} D_{\mathcal{KL}}(q_{\phi}(\mathsf{z}|\mathsf{x})||p_{ heta}(\mathsf{z}|\mathsf{x})),$$

- Q is a (specified) variational family (e.g. mean-field)
- ϕ contain the variational parameters; θ contain the model parameters.

KL Divergence:

$$D_{\mathcal{KL}}(q_{\phi}(\mathsf{z}|\mathsf{x})||p_{ heta}(\mathsf{z}|\mathsf{x})) := \int \log\Big(rac{q_{\phi}(\mathsf{z}|\mathsf{x})}{p_{ heta}(\mathsf{z}|\mathsf{x})}\Big)q_{\phi}(\mathsf{z}|\mathsf{x})d\mathsf{z} = E_q\Big[\lograc{q_{\phi}(\mathsf{Z}|\mathsf{x})}{p_{ heta}(\mathsf{Z}|\mathsf{x})}\Big]$$

Objective: Obtain approximate posterior (or recognition model)

$$q_{\phi}^{*}(\mathsf{z}|\mathsf{x}) = rgmin_{q \in \mathcal{Q}} D_{\mathcal{KL}}(q_{\phi}(\mathsf{z}|\mathsf{x})||p_{ heta}(\mathsf{z}|\mathsf{x})),$$

- Q is a (specified) variational family (e.g. mean-field)
- ϕ contain the variational parameters; θ contain the model parameters.

KL Divergence:

$$D_{\mathit{KL}}(q_{\phi}(\mathsf{z}|\mathsf{x})||p_{ heta}(\mathsf{z}|\mathsf{x})) := \int \log\Big(rac{q_{\phi}(\mathsf{z}|\mathsf{x})}{p_{ heta}(\mathsf{z}|\mathsf{x})}\Big) q_{\phi}(\mathsf{z}|\mathsf{x}) d\mathsf{z} = E_q\Big[\lograc{q_{\phi}(\mathsf{Z}|\mathsf{x})}{p_{ heta}(\mathsf{Z}|\mathsf{x})}\Big]$$

- Asymmetric: $D_{KL}(q||p)$ (exclusive KL) $\neq D_{KL}(p||q)$ (inclusive KL).
 - VI chooses to optimize $D_{\mathcal{KL}}(q||p)$ b/c it's easier to take expectations w.r.t. q.
 - There's existing work about optimizing $D_{KL}(p||q)$ (Markovian score climbing).
- Exclusive KL: $q(x) > 0 \Rightarrow p(x) > 0$

I

• The marginal likelihood (evidence) $p(\mathbf{x})$ can be expressed as

$$\log p_{ heta}(\mathbf{x}) = D_{\mathcal{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{ heta}(\mathbf{z}|\mathbf{x})) + \underbrace{E_qig[-\log q_{\phi}(\mathbf{z}|\mathbf{x}) + \log p_{ heta}(\mathbf{x},\mathbf{z})ig]}_{\mathcal{L}(heta,\phi;\mathbf{x})} \geq \mathcal{L}(heta,\phi;\mathbf{x}),$$

where $\mathcal{L}(\theta, \phi; \mathbf{x})$ is the **evidence lower bound** (ELBO) which can also be written as

$$\mathcal{L}(heta, \phi; \mathbf{x}) = -D_{\mathcal{K}\mathcal{L}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z})) + E_{q} \big[\log p_{\theta}(\mathbf{x}|\mathbf{z})\big].$$

• The marginal likelihood (evidence) $p(\mathbf{x})$ can be expressed as

$$\log p_{\theta}(\mathbf{x}) = D_{\mathcal{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) + \underbrace{\mathcal{E}_{q}\big[-\log q_{\phi}(\mathbf{z}|\mathbf{x}) + \log p_{\theta}(\mathbf{x},\mathbf{z})\big]}_{\mathcal{L}(\theta,\phi;\mathbf{x})} \geq \mathcal{L}(\theta,\phi;\mathbf{x}),$$

where $\mathcal{L}(\theta, \phi; \mathbf{x})$ is the **evidence lower bound** (ELBO) which can also be written as

$$\mathcal{L}(\theta,\phi;\mathbf{x}) = -D_{\mathcal{K}\mathcal{L}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z})) + E_{q}\big[\log p_{\theta}(\mathbf{x}|\mathbf{z})\big].$$

• **Goal**: Optimize $\mathcal{L}(\theta, \phi; \mathbf{x})$ w.r.t. both ϕ and θ .

• The marginal likelihood (evidence) $p(\mathbf{x})$ can be expressed as

$$\log p_{\theta}(\mathbf{x}) = D_{\mathcal{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) + \underbrace{\mathcal{E}_{q}\big[-\log q_{\phi}(\mathbf{z}|\mathbf{x}) + \log p_{\theta}(\mathbf{x},\mathbf{z})\big]}_{\mathcal{L}(\theta,\phi;\mathbf{x})} \geq \mathcal{L}(\theta,\phi;\mathbf{x}),$$

where $\mathcal{L}(\theta, \phi; \mathbf{x})$ is the **evidence lower bound** (ELBO) which can also be written as

$$\mathcal{L}(heta, \phi; \mathbf{x}) = -D_{\mathcal{K}\mathcal{L}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{ heta}(\mathbf{z})) + E_{q} \big[\log p_{ heta}(\mathbf{x}|\mathbf{z})\big].$$

- **Goal**: Optimize $\mathcal{L}(\theta, \phi; \mathbf{x})$ w.r.t. both ϕ and θ .
 - Optimizing w.r.t. θ is easy
 - Unbiased gradient estimates are straightforward to obtain.

• The marginal likelihood (evidence) $p(\mathbf{x})$ can be expressed as

$$\log p_{ heta}(\mathbf{x}) = D_{\mathcal{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}))||p_{ heta}(\mathbf{z}|\mathbf{x})) + \underbrace{E_q ig[-\log q_{\phi}(\mathbf{z}|\mathbf{x}) + \log p_{ heta}(\mathbf{x}, \mathbf{z})ig]}_{\mathcal{L}(heta, \phi; \mathbf{x})} \geq \mathcal{L}(heta, \phi; \mathbf{x}),$$

where $\mathcal{L}(\theta, \phi; \mathbf{x})$ is the **evidence lower bound** (ELBO) which can also be written as

$$\mathcal{L}(\theta,\phi;\mathbf{x}) = -D_{\mathcal{K}\mathcal{L}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z})) + E_{q} \big[\log p_{\theta}(\mathbf{x}|\mathbf{z})\big].$$

- **Goal**: Optimize $\mathcal{L}(\theta, \phi; \mathbf{x})$ w.r.t. both ϕ and θ .
 - Optimizing w.r.t. θ is easy
 - Unbiased gradient estimates are straightforward to obtain.
 - Optimizing w.r.t. ϕ is harder.
 - Expectations in the ELBO are taken w.r.t. $q_{\phi}(\mathbf{z}|\mathbf{x})$, which is a function of ϕ .
 - Reparameterization trick is useful here.

Overview

1. Background

Bayesian Inference/Latent variable modeling Variational Inference

2. Overview of contributions

3. Paper #1 Reparameterization trick Stochastic Gradient VB Estimators

Auto-encoding VB Algorithm Variational Auto-Encoder

4. Paper #2 Normalizing flows VI algorithm

- 5. Comparison of papers
- 6. Related work

Overview of contributions

Paper #1: Auto-encoding Variational Bayes (Kingma et al. 2014)

- Two significant contributions
 - 1. SGVB estimator + AEVB algorithm
 - 2. Variational Auto-encoders (VAEs)

Overview of contributions

Paper #1: Auto-encoding Variational Bayes (Kingma et al. 2014)

- Two significant contributions
 - 1. SGVB estimator + AEVB algorithm
 - 2. Variational Auto-encoders (VAEs)

Paper #2: Variational Inference with Normalizing Flows (Rezende et al. 2015)

- Provides a rich, computationally-feasible approximate posterior using normalizing flows.
 - In the asymptotic regime, the space of solutions is rich enough to contain the true posterior distribution.
- Uses gradient estimator from Paper #1 (amortized VI).

Overview of contributions

Paper #1: Auto-encoding Variational Bayes (Kingma et al. 2014)

- Two significant contributions
 - 1. SGVB estimator + AEVB algorithm
 - 2. Variational Auto-encoders (VAEs)

Paper #2: Variational Inference with Normalizing Flows (Rezende et al. 2015)

- Provides a rich, computationally-feasible approximate posterior using normalizing flows.
 - In the asymptotic regime, the space of solutions is rich enough to contain the true posterior distribution.
- Uses gradient estimator from Paper #1 (amortized VI).

Topics I will not cover:

- Experiments/empirical results
- infinitesimal flows

Overview

1. Background

Bayesian Inference/Latent variable modeling Variational Inference

2. Overview of contributions

Paper #1
 Reparameterization trick
 Stochastic Gradient VB Estimators
 Auto-encoding VB Algorithm
 Variational Auto-Encoder

4. Paper #2 Normalizing flows VI algorithm

- 5. Comparison of papers
- 6. Related work

Suppose $\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})$. Then, it is often possible to express \mathbf{z} as

 $\mathbf{z} = g_{\phi}(\epsilon, \mathbf{x}),$

where $\epsilon \sim p(\epsilon)$ (a distribution that doesn't depend on ϕ).

Example: Suppose $z \sim N(\mu, \sigma^2)$. Then $z = \mu + \sigma \epsilon$, where $\epsilon \sim N(0, 1)$.

Suppose $\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})$. Then, it is often possible to express \mathbf{z} as

 $\mathbf{z} = g_{\phi}(\epsilon, \mathbf{x}),$

where $\epsilon \sim p(\epsilon)$ (a distribution that doesn't depend on ϕ).

Example: Suppose $z \sim N(\mu, \sigma^2)$. Then $z = \mu + \sigma \epsilon$, where $\epsilon \sim N(0, 1)$. Why is it useful?

• In general, for a differentiable function f,

 $abla_{\phi} E_{q_{\phi}}[f(\mathbf{z})] \neq E_{q_{\phi}}[
abla_{\phi}f(\mathbf{z})]$

Suppose $\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})$. Then, it is often possible to express \mathbf{z} as

 $\mathbf{z} = g_{\phi}(\epsilon, \mathbf{x}),$

where $\epsilon \sim p(\epsilon)$ (a distribution that doesn't depend on ϕ).

Example: Suppose $z \sim N(\mu, \sigma^2)$. Then $z = \mu + \sigma \epsilon$, where $\epsilon \sim N(0, 1)$. Why is it useful?

• In general, for a differentiable function f,

 $abla_{\phi} E_{\boldsymbol{q}_{\phi}}[f(\boldsymbol{z})] \neq E_{\boldsymbol{q}_{\phi}}[
abla_{\phi}f(\boldsymbol{z})]$

• However, with the reparameterization trick,

$$egin{aligned}
abla_{\phi} E_{q_{\phi}}[f(\mathbf{z})] &=
abla_{\phi} E_{p(\epsilon)}[f(g_{\phi}(\epsilon, \mathbf{x}))] \ &= E_{p(\epsilon)}[
abla_{\phi}f(g_{\phi}(\epsilon, \mathbf{x}))] \ &pprox rac{1}{L}\sum_{l=1}^{L}f(g_{\phi}(\epsilon^{(l)}, \mathbf{x})), ext{ where } \epsilon^{(l)} \sim p(\epsilon) \end{aligned}$$



Illustration of reparameterization trick from Kingma (2017).

Paper #1: Stochastic Gradient Variational Bayes Estimators

Using the reparameterization trick, a generic SGVB estimator of the ELBO is given by

$$ilde{\mathcal{L}}^{\mathcal{A}}(heta,\phi;\mathbf{x}) = rac{1}{L}\sum_{l=1}^{L} \Big[\log p_{ heta}(\mathbf{x},\mathbf{z}^{(l)}) - \log q_{\phi}(\mathbf{z}^{(l)}|\mathbf{x})\Big],$$

where $\mathbf{z}^{(l)} = g_{\phi}(\epsilon^{(l)}, \mathbf{x})$ and $\epsilon^{(l)} \sim p(\epsilon)$, for l = 1, ..., L.

Paper #1: Stochastic Gradient Variational Bayes Estimators

Using the reparameterization trick, a generic SGVB estimator of the ELBO is given by

$$ilde{\mathcal{L}}^{\mathcal{A}}(heta,\phi;\mathbf{x}) = rac{1}{L}\sum_{l=1}^{L} \Big[\log p_{ heta}(\mathbf{x},\mathbf{z}^{(l)}) - \log q_{\phi}(\mathbf{z}^{(l)}|\mathbf{x})\Big],$$

where
$$\mathbf{z}^{(l)} = g_{\phi}(\epsilon^{(l)}, \mathbf{x})$$
 and $\epsilon^{(l)} \sim p(\epsilon)$, for $l = 1, \dots, L$.

When $D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}))$ can be computed analytically (e.g. Gaussian case),

$$ilde{\mathcal{L}}^{B}(heta, \phi; \mathbf{x}) = -D_{\mathcal{K}L}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{ heta}(\mathbf{z})) + rac{1}{L}\sum_{l=1}^{L}\log p_{ heta}(\mathbf{x}|\mathbf{z}^{(l)})$$

• $\tilde{\mathcal{L}}^B$ typically has less variance than $\tilde{\mathcal{L}}^A$.

Paper #1: Auto-Encoding Variational Bayes Algorithm

Algorithm 1 Minibatch version of the Auto-Encoding VB (AEVB) algorithm. Either of the two SGVB estimators in section 2.3 can be used. We use settings M = 100 and L = 1 in experiments.

$$\boldsymbol{\theta}, \boldsymbol{\phi} \leftarrow \text{Initialize parameters}$$

repeat

 $\mathbf{X}^M \leftarrow \text{Random minibatch of } M \text{ datapoints (drawn from full dataset)}$

 $\epsilon \leftarrow \text{Random samples from noise distribution } p(\epsilon)$

 $\mathbf{g} \leftarrow \nabla_{\boldsymbol{\theta}, \boldsymbol{\phi}} \widetilde{\mathcal{L}}^{M}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{X}^{M}, \boldsymbol{\epsilon})$ (Gradients of minibatch estimator (8)) $\boldsymbol{\theta}, \boldsymbol{\phi} \leftarrow$ Update parameters using gradients \mathbf{g} (e.g. SGD or Adagrad [DHS10]) until convergence of parameters ($\boldsymbol{\theta}, \boldsymbol{\phi}$)

return $\boldsymbol{\theta}, \boldsymbol{\phi}$

The mini-batch estimator $\tilde{\mathcal{L}}^{M}$ is given by

$$ilde{\mathcal{L}}^{M}(heta,\phi;X^{M}) = rac{N}{M}\sum_{i=1}^{M} ilde{\mathcal{L}}(heta,\phi;\mathbf{x})$$

How is it related to auto-encoders?

An **autoencoder** is a neural network used for unsupervised learning that minimizes an objective function with the form: reconstruction error + regularizer.

Recall:

$$\tilde{\mathcal{L}}^{\mathcal{B}}(\theta,\phi;\mathbf{x}) = \underbrace{-D_{\mathcal{K}L}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}))}_{\text{regularizer}} + \underbrace{\frac{1}{L}\sum_{l=1}^{L}\log p_{\theta}(\mathbf{x}|\mathbf{z}^{(l)})}_{\text{Negative reconstruction error}}$$

VAEs are given as an example but AFAIK, this is the original paper that introduced them.

Assume **x**, **z** are continuous.

Decoder (generative model)

$$p(\mathbf{z}) = N(\mathbf{z}; 0, I)$$

$$p(\mathbf{x}|\mathbf{z}) = N(\mathbf{x}; \mu, \sigma^2 I)$$

$$\mu = W_2 h + b_2$$

$$\log \sigma^2 = W_3 h + b_3$$

$$h = \tanh(W_1 \mathbf{z} + b_1)$$

$$\theta = \{W_j, b_j : j = 1, 2, 3\}, \phi = \{\tilde{W}_j, \tilde{b}_j : j = 1, 2, 3\}$$

Encoder (inference model)

$$egin{aligned} q_{\phi}(\mathbf{z}|\mathbf{x}) &= \mathcal{N}(\mathbf{z}; ilde{\mu}, ilde{\sigma}^2 \mathbf{I}) \ & ilde{h} &= ext{tanh}(ilde{W}_1\mathbf{x} + ilde{b}_1) \ & ilde{\mu} &= ilde{W}_1 ilde{h} + ilde{b}_2 \ & ext{log}\, ilde{\sigma}^2 &= ilde{W}_3 ilde{h} + ilde{b}_3 \end{aligned}$$

VAEs are given as an example but AFAIK, this is the original paper that introduced them.

Assume **x**, **z** are continuous.

Decoder (generative model)

 $p(\mathbf{z}) = N(\mathbf{z}; 0, I) \qquad q_{\phi}(\mathbf{z}|\mathbf{x}) = N(\mathbf{z}; \tilde{\mu}, \tilde{\sigma}^2 I) \\ p(\mathbf{x}|\mathbf{z}) = N(\mathbf{x}; \mu, \sigma^2 I) \qquad \tilde{h} = \tanh(\tilde{W}_1 \mathbf{x} + \tilde{b}_1) \\ \mu = W_2 h + b_2 \qquad \tilde{\mu} = \tilde{W}_1 \tilde{h} + \tilde{b}_2 \\ \log \sigma^2 = W_3 h + b_3 \qquad \log \tilde{\sigma}^2 = \tilde{W}_3 \tilde{h} + \tilde{b}_3 \\ h = \tanh(W_1 \mathbf{z} + b_1)$

Encoder (inference model)

 $\theta = \{W_j, b_j : j = 1, 2, 3\}, \phi = \{\tilde{W}_j, \tilde{b}_j : j = 1, 2, 3\}$

- VAEs are a non-linear generalization of probabilistic PCA, where $\mu = \mathbf{W}\mathbf{z}$.
 - In this case, the evidence has an analytical form: $p(\mathbf{x}) = N(0, \mathbf{W}\mathbf{W}' + \sigma^2 I)$.

VAEs are given as an example but AFAIK, this is the original paper that introduced them.

Assume **x**, **z** are continuous.

Decoder (generative model)

 $p(\mathbf{z}) = N(\mathbf{z}; 0, I) \qquad q_{\phi}(\mathbf{z}|\mathbf{x}) = N(\mathbf{z}; \tilde{\mu}, \tilde{\sigma}^2 I) \\ p(\mathbf{x}|\mathbf{z}) = N(\mathbf{x}; \mu, \sigma^2 I) \qquad \tilde{h} = \tanh(\tilde{W}_1 \mathbf{x} + \tilde{b}_1) \\ \mu = W_2 h + b_2 \qquad \tilde{\mu} = \tilde{W}_1 \tilde{h} + \tilde{b}_2 \\ \log \sigma^2 = W_3 h + b_3 \qquad \log \tilde{\sigma}^2 = \tilde{W}_3 \tilde{h} + \tilde{b}_3 \\ h = \tanh(W_1 \mathbf{z} + b_1)$

Encoder (inference model)

 $\theta = \{W_j, b_j : j = 1, 2, 3\}, \phi = \{\tilde{W}_j, \tilde{b}_j : j = 1, 2, 3\}$

- VAEs are a non-linear generalization of probabilistic PCA, where $\mu = \mathbf{W}\mathbf{z}$.
 - In this case, the evidence has an analytical form: $p(\mathbf{x}) = N(0, \mathbf{WW}' + \sigma^2 I)$.

To sample from $q_{\phi}(\mathbf{z}|\mathbf{x})$,

- 1. Sample $\epsilon^{(\ell)} \sim N(0, I)$, $\ell = 1, \dots, L$
- 2. Compute $\mathbf{z}^{(\ell)} = \tilde{\mu} + \tilde{\sigma} \odot \epsilon^{(\ell)}$, where \odot denotes element-wise product.

In the case where p(z) and $q_{\phi}(z|\mathbf{x})$ are Gaussian, $-D_{KL}(q_{\phi}(z||\mathbf{x})||p(z))$ can be computed analytically.

$$-D_{\mathcal{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) = rac{1}{2}\sum_{j=1}^{J}\left[1+\log(ilde{\sigma}_{j}^{2})- ilde{\mu}_{j}^{2}- ilde{\sigma}_{j}^{2}
ight]$$

In the case where p(z) and $q_{\phi}(z|\mathbf{x})$ are Gaussian, $-D_{KL}(q_{\phi}(z||\mathbf{x})||p(z))$ can be computed analytically.

$$- \mathcal{D}_{\mathcal{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) = rac{1}{2}\sum_{j=1}^{J} \left[1 + \log(ilde{\sigma}_{j}^{2}) - ilde{\mu}_{j}^{2} - ilde{\sigma}_{j}^{2}
ight]$$

The SGVB estimator for VAEs is given by,

$$\tilde{\mathcal{L}}^{B}(\theta,\phi;\mathbf{x}) = -D_{\mathcal{K}L}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) + \frac{1}{L}\sum_{\ell=1}^{L}\log p_{\theta}(\mathbf{x}|\mathbf{z}^{(\ell)})$$
(1)

where $\mathbf{z}^{(\ell)} = \tilde{\mu} + \tilde{\sigma} \odot \epsilon^{(\ell)}$ and $\epsilon^{(\ell)} \sim N(0, I)$.

In the case where p(z) and $q_{\phi}(z|\mathbf{x})$ are Gaussian, $-D_{KL}(q_{\phi}(z||\mathbf{x})||p(z))$ can be computed analytically.

$$- \mathcal{D}_{\mathcal{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) = rac{1}{2}\sum_{j=1}^{J} \left[1 + \log(ilde{\sigma}_{j}^{2}) - ilde{\mu}_{j}^{2} - ilde{\sigma}_{j}^{2}
ight]$$

The SGVB estimator for VAEs is given by,

$$\tilde{\mathcal{L}}^{B}(\theta,\phi;\mathbf{x}) = -D_{\mathcal{K}L}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) + \frac{1}{L}\sum_{\ell=1}^{L}\log p_{\theta}(\mathbf{x}|\mathbf{z}^{(\ell)})$$
(1)

where $\mathbf{z}^{(\ell)} = \tilde{\mu} + \tilde{\sigma} \odot \epsilon^{(\ell)}$ and $\epsilon^{(\ell)} \sim N(0, I)$.

• We can optimize (1) using the AEVB algorithm to obtain estimates of θ and ϕ .

Overview

1. Background

Bayesian Inference/Latent variable modeling Variational Inference

- 2. Overview of contributions
- Paper #1
 Reparameterization trick
 Stochastic Gradient VB Estimators
 Auto-encoding VB Algorithm
 Variational Auto-Encoder
- Paper #2 Normalizing flows VI algorithm
- 5. Comparison of papers
- 6. Related work

- In paper #1, a Gaussian distribution with diagonal covariance was used as the approximate posterior.
- How can we obtain a more flexible family of variational distributions?
 - Solution: Normalizing flows
- **Normalizing flows** transform simple distributions (e.g. Gaussian) through a sequence of invertible mappings into rich complex distributions.

- In paper #1, a Gaussian distribution with diagonal covariance was used as the approximate posterior.
- How can we obtain a more flexible family of variational distributions?
 - Solution: Normalizing flows
- **Normalizing flows** transform simple distributions (e.g. Gaussian) through a sequence of invertible mappings into rich complex distributions.

Review: Change of variables

Let $f : \mathbb{R}^d \to \mathbb{R}^d$ be an invertible, smooth function with inverse f^{-1} and $\mathbf{z} \sim q(\mathbf{z})$. Then $\mathbf{z}' = f(\mathbf{z})$ has the distribution:

$$q(\mathbf{z}') = q(\mathbf{z}) \Big| \mathsf{det} rac{\partial f^{-1}}{\partial \mathbf{z}'} \Big| = q(\mathbf{z}) \Big| \mathsf{det} rac{\partial f}{\partial \mathbf{z}} \Big|^{-1}$$

Suppose $\mathbf{z}_0 \sim q_0(\mathbf{z}_0)$ and $\mathbf{z}_K = f_K \circ \ldots f_2 \circ f_1(\mathbf{z}_0)$, where f_1, \ldots, f_K is a sequence of mappings. Then the log density of \mathbf{z}_K is given by

$$\log q_{\mathcal{K}}(\mathsf{z}_{\mathcal{K}}) = \log q_0(\mathsf{z}_0) - \sum_{k=1}^{\mathcal{K}} \log \det \Big| rac{\partial f_k}{\partial \mathsf{z}_k} \Big|$$

Suppose $\mathbf{z}_0 \sim q_0(\mathbf{z}_0)$ and $\mathbf{z}_K = f_K \circ \ldots f_2 \circ f_1(\mathbf{z}_0)$, where f_1, \ldots, f_K is a sequence of mappings. Then the log density of \mathbf{z}_K is given by

$$\log q_{\mathcal{K}}(\mathsf{z}_{\mathcal{K}}) = \log q_0(\mathsf{z}_0) - \sum_{k=1}^{\mathcal{K}} \log \det \Bigl| rac{\partial f_k}{\partial \mathsf{z}_k} \Bigr|$$

By the law of the unconscious statistician, for any function h,

$$E_{q_{\mathcal{K}}}[h(\mathbf{z})] = E_{q_0}[h(f_{\mathcal{K}} \circ \cdots \circ f_1(\mathbf{z}_0))]$$

• There are many good choices for the sequence of invertible transformations. Paper #2 uses planar and radial flows.

Planar flows are transformations of the form

$$f(\mathbf{z}) = \mathbf{z} + \mathbf{u}h(\mathbf{w}'\mathbf{z} + b),$$

where $\lambda = {\mathbf{w} \in \mathbb{R}^D, \mathbf{u} \in \mathbb{R}^D, b \in \mathbb{R}}$ are free parameters and *h* is a smooth function with derivative *h'*.

Planar flows are transformations of the form

$$f(\mathbf{z}) = \mathbf{z} + \mathbf{u}h(\mathbf{w}'\mathbf{z} + b),$$

where $\lambda = {\mathbf{w} \in \mathbb{R}^D, \mathbf{u} \in \mathbb{R}^D, b \in \mathbb{R}}$ are free parameters and *h* is a smooth function with derivative *h'*.

• The logdet-Jacobian term can computed in O(D) time:

$$\det \left| \frac{\partial f}{\partial \mathbf{z}} \right| = \left| \det(\mathbf{I} + \mathbf{u}\psi(\mathbf{z})') \right| = |\mathbf{1} + \mathbf{u}'\psi(\mathbf{z})|,$$

$$\psi(\mathbf{z}) = h'(\mathbf{w}'\mathbf{z} + b)\mathbf{w}.$$

Planar flows are transformations of the form

$$f(\mathbf{z}) = \mathbf{z} + \mathbf{u}h(\mathbf{w}'\mathbf{z} + b),$$

where $\lambda = {\mathbf{w} \in \mathbb{R}^D, \mathbf{u} \in \mathbb{R}^D, b \in \mathbb{R}}$ are free parameters and *h* is a smooth function with derivative *h'*.

• The logdet-Jacobian term can computed in O(D) time:

$$\det \left| \frac{\partial f}{\partial \mathbf{z}} \right| = \left| \det(\mathbf{I} + \mathbf{u}\psi(\mathbf{z})') \right| = |\mathbf{1} + \mathbf{u}'\psi(\mathbf{z})|,$$

$$\psi(\mathbf{z}) = h'(\mathbf{w}'\mathbf{z} + b)\mathbf{w}.$$

• The log-density of $z_{\mathcal{K}} = f_{\mathcal{K}} \circ \cdots \circ f_1(\mathbf{z})$, where $f_k = \mathbf{z} + \mathbf{u}_k h(\mathbf{w}'_k \mathbf{z} + b_k)$:

$$\log q_{\mathcal{K}}(\mathsf{z}_{\mathcal{K}}) = \log q_0(\mathsf{z}) - \sum_{k=1}^{\mathcal{K}} \log |1 + \mathsf{u}_k' \psi_k(\mathsf{z}_k)|,$$

Radial flows are transformations of the form:

$$f(\mathbf{z}) = \mathbf{z} + \beta h(\alpha, r)(\mathbf{z} - \mathbf{z}_0),$$

where $r = ||\mathbf{z} - \mathbf{z}_0||$; $h(\alpha, r) = 1/(\alpha + r)$; and $\lambda = {\mathbf{z}_0 \in \mathbb{R}^D, \alpha \in \mathbb{R}, \beta \in \mathbb{R}}$ are parameters.

• The logdet-Jacobian can also be computed in linear time:

$$\det \left| \frac{\partial f}{\partial \mathbf{z}} \right| = \left[1 + \beta h(\alpha, r) \right]^{d-1} \left[1 + \beta h(\alpha, r) + h'(\alpha, r) r \right]^{d-1}$$



Figure 1 in Rezende et al. (2015) showing effects planar/radial flows on two distributions.

Flow-based ELBO (negative free energy bound)

Suppose the approximate posterior is parameterized with a (planar) flow of length K i.e. $q_{\phi}(\mathbf{z}|\mathbf{x}) := q_{K}(\mathbf{z}_{K})$.

Then the ELBO is given by

$$\begin{split} \mathcal{L}(\theta,\phi;\mathbf{x}) &= E_{q_{\phi}(\mathbf{z}|\mathbf{x})} \Big[-\log q_{\phi}(\mathbf{z}|\mathbf{x}) + \log p(\mathbf{x},\mathbf{z}) \Big] \\ &= E_{q_{0}(z_{0})} \Big[-\log q_{K}(\mathbf{z}_{K}) + \log p(\mathbf{x},\mathbf{z}_{K}) \Big] \\ &= -E_{q_{0}(z_{0})} \big[\log q_{0}(z_{0}) \big] + E_{q_{0}(z_{0})} \big[\log p(\mathbf{x},\mathbf{z}_{K}) \big] + E_{q_{0}(z_{0})} \Big[\sum_{k=1}^{K} \log |1 + \mathbf{u}_{k}'\psi_{k}(\mathbf{z}_{k}) \Big] \end{split}$$

-

Flow-based ELBO (negative free energy bound)

Suppose the approximate posterior is parameterized with a (planar) flow of length K i.e. $q_{\phi}(\mathbf{z}|\mathbf{x}) := q_{K}(\mathbf{z}_{K})$.

Then the ELBO is given by

$$\begin{aligned} \mathcal{L}(\theta,\phi;\mathbf{x}) &= E_{q_{\phi}(\mathbf{z}|\mathbf{x})} \Big[-\log q_{\phi}(\mathbf{z}|\mathbf{x}) + \log p(\mathbf{x},\mathbf{z}) \Big] \\ &= E_{q_{0}(z_{0})} \Big[-\log q_{K}(\mathbf{z}_{K}) + \log p(\mathbf{x},\mathbf{z}_{K}) \Big] \\ &= -E_{q_{0}(z_{0})} \Big[\log q_{0}(z_{0}) \Big] + E_{q_{0}(z_{0})} \Big[\log p(\mathbf{x},\mathbf{z}_{K}) \Big] + E_{q_{0}(z_{0})} \Big[\sum_{k=1}^{K} \log |1 + \mathbf{u}_{k}^{\prime} \psi_{k}(\mathbf{z}_{k}) \Big] \end{aligned}$$

- Paper #2 also constructs a recognition model (inference network) using a deep neural network that maps observations x to the parameters of q₀ = N(μ, σ²) and flow parameters λ.
 - They don't give explicit details about how they do this like paper #1 does.
 - It could be interesting to work out the explicit form of the ELBO here (*if time permits*).

Algorithm 1 Variational Inf. with Normalizing Flows

Parameters: ϕ variational, θ generative while not converged do $\mathbf{x} \leftarrow \{\text{Get mini-batch}\}$ $\mathbf{z}_0 \sim q_0(\bullet | \mathbf{x})$ $\mathbf{z}_K \leftarrow f_K \circ f_{K-1} \circ \ldots \circ f_1(\mathbf{z}_0)$ $\mathcal{F}(\mathbf{x}) \approx \mathcal{F}(\mathbf{x}, \mathbf{z}_K)$ $\Delta \theta \propto -\nabla_{\theta} \mathcal{F}(\mathbf{x})$ $\Delta \phi \propto -\nabla_{\phi} \mathcal{F}(\mathbf{x})$ end while

- Similiar in nature to the AEVB algorithm
- The authors don't make this explicit in the paper MC estimates are still necessary (I think?)

Overview

1. Background

Bayesian Inference/Latent variable modeling Variational Inference

- 2. Overview of contributions
- Paper #1
 Reparameterization trick
 Stochastic Gradient VB Estimators
 Auto-encoding VB Algorithm
 Variational Auto-Encoder

4. Paper #2 Normalizing flows VI algorithm

5. Comparison of papers

6. Related work

- Both papers focus on different problems in VI.
 - Paper #1 focuses on estimating the ELBO and its gradient.
 - Paper #2 focuses on on coming up with a flexible variational family.
- In paper #2, the initial density is q₀ = N(μ, σ²), where μ and σ are parameterized with neural networks. This is similiar to paper #1 but they use this as the approximate posterior, whereas this is transformed through normalizing flows in paper #2.
- Paper #1 gives a specific example where AEVB is applied (VAEs) whereas paper #2 talks about their algorithm more generally.
- Your thoughts?

Overview

1. Background

Bayesian Inference/Latent variable modeling Variational Inference

- 2. Overview of contributions
- Paper #1
 Reparameterization trick
 Stochastic Gradient VB Estimators
 Auto-encoding VB Algorithm
 Variational Auto-Encoder
- 4. Paper #2 Normalizing flows VI algorithm
- 5. Comparison of papers
- 6. Related work

- Wake-sleep algorithm (Hinton et al. (1995))
 - Applicable to same general class of continuous latent variable models as AEVB and also discrete latent variables.
 - Drawback: concurrent optimization of two objective functions ⇒ optimization of (a bound of) the evidence.
- Non-linear Independent Components Estimation (NICE) (Dinh et al. (2014))
 - Transformations are neural networks with easy to compute inverses
- Hamiltonian variational approx. (HVI) (Salimans et al. (2015))
 - infinitesimal volume-preserving flow
 - Elegant but not as computationally efficient