

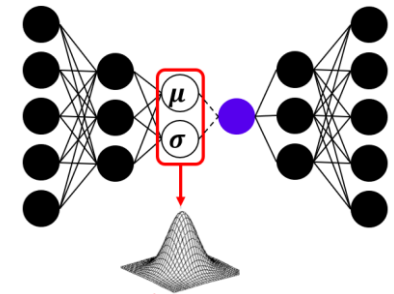
Probabilistic Models for Sequential Data

杨帆

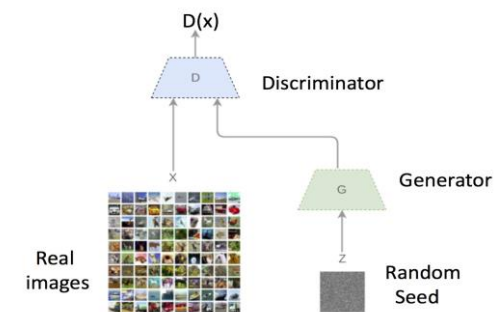
2018年7月19日

Motivations For This Talk (1)

- Some trends in machine learning: make ML to be
 - Robust to uncertain and adversarial inputs
 - Unsupervised, semi-supervised or self-supervised
- Key tools: ***probabilistic generative models***
 - Explicitly model the uncertainty in data
 - Recent advances make them scalable to large datasets and complex data distributions
 - E.g., two popular frameworks: VAE and GAN



Variational Autoencoders



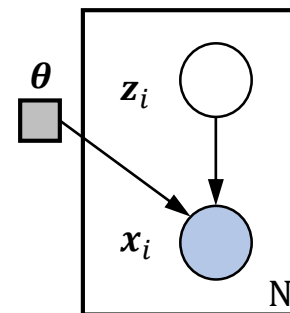
Generative Adversarial Networks

Motivations For This Talk (2)

- Last time we introduced the **Variational Autoencoders**
 - They assume a deep latent gaussian generative model

$$\mathbf{z}_i \sim \mathcal{N}(\mathbf{0}, \mathbb{I}), \quad i = 1, \dots, N,$$
$$\mathbf{x}_i | \mathbf{z}_i \sim \mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{z}_i), \boldsymbol{\sigma}_{\boldsymbol{\theta}}^2(\mathbf{z}_i)), \quad i = 1, \dots, N.$$

where $\boldsymbol{\mu}_{\boldsymbol{\theta}}(\cdot)$ and $\boldsymbol{\sigma}_{\boldsymbol{\theta}}(\cdot)$ are neural networks



- For sequential data, we may simply plug CNN or RNN and play
 - However, deep latent gaussian models are not so natural for **sequential data**
- Today we will introduce several probabilistic models specially designed for sequential data
 - Enable us do prediction and inference in a principled way

Outline

- Short (Re)Introduction to VI and VAE
- State Space Models and Stochastic RNNs
- Variational Inference of SSM
- (Variational) Sequential Monte Carlo
- Conclusions

Outline

- **Short (Re)Introduction to VI and VAE**
- State Space Models and Stochastic RNNs
- Variational Inference of SSM
- (Variational) Sequential Monte Carlo
- Conclusions

Some Clarifications

- **Maximum Likelihood Estimation (MLE)**

- θ are treated as unknown static parameters

$$\theta_{MLE} = \operatorname{argmax}_{\theta} p_{\theta}(\mathcal{D})$$

- **Maximum A Posterior estimation (MAP)**

- θ are treated as latent variables and we have the prior $p(\theta)$

$$\theta_{MAP} = \operatorname{argmax}_{\theta} p(\theta|\mathcal{D})$$

where

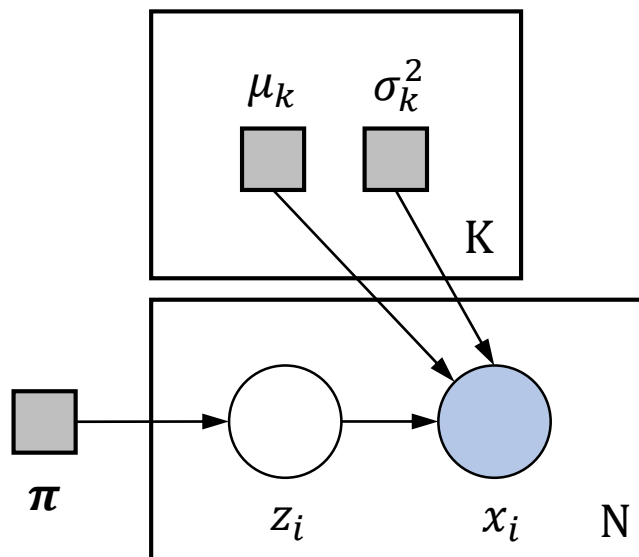
$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})}$$

- **Nonparametric Bayes: Gaussian process, Dirichlet process, ...**

- a class of models in which the number of parameters grows with sample size

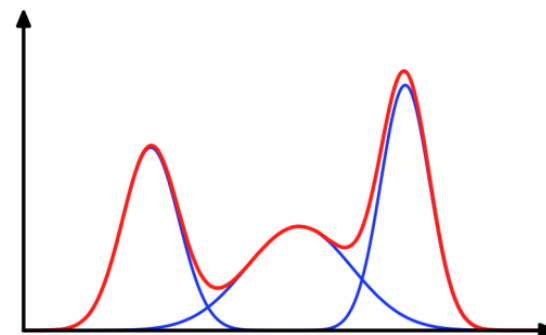
Probabilistic Generative Models

- Example: Gaussian Mixture Models



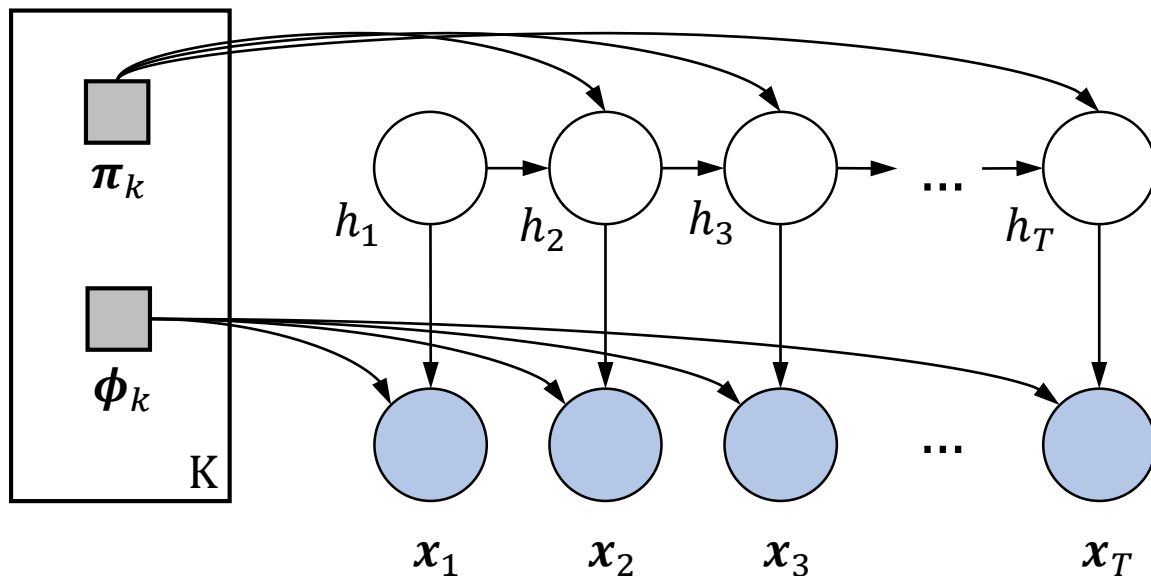
$$z_i \sim \text{categorical}(\pi_1, \dots, \pi_K), \quad i = 1, \dots, N,$$

$$x_i | z_i \sim \mathcal{N}(\mu_{z_i}, \sigma_{z_i}^2) \quad i = 1, \dots, N.$$



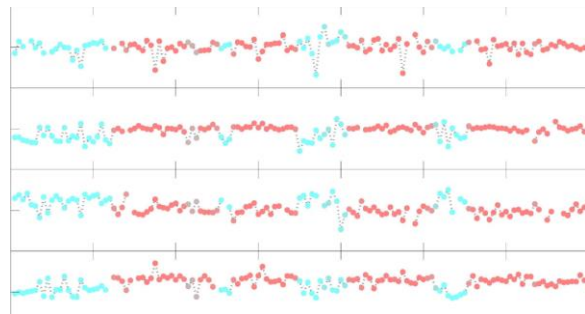
Probabilistic Generative Models

- Example: Hidden Markov Models



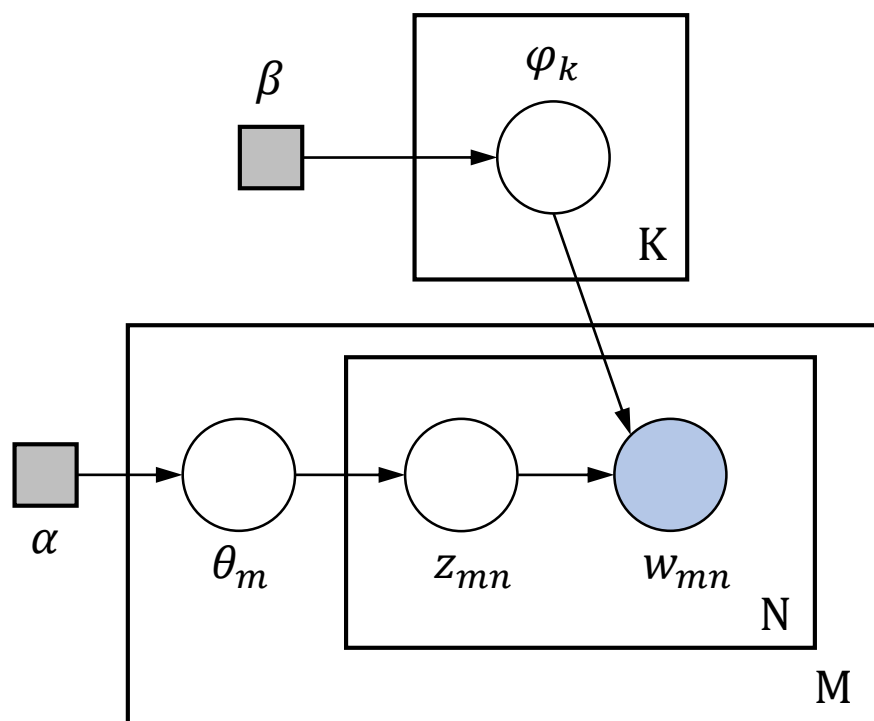
$$h_t \sim \text{categorical}(\pi_{h_{t-1}})$$

$$x_t | h_t \sim F(\phi_{h_t})$$



Probabilistic Generative Models

- Example: Latent Dirichlet Allocation (fully Bayesian)

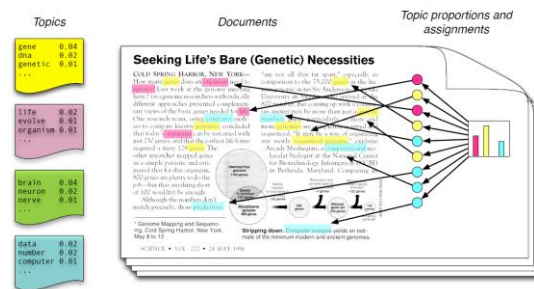


$$\varphi_k \sim \text{Dir}(\beta), \quad k = 1, \dots, K,$$

$$\theta_m \sim \text{Dir}(\alpha), \quad m = 1, \dots, M,$$

$$z_{mn} | \theta_m \sim \text{categorical}(\theta_m),$$

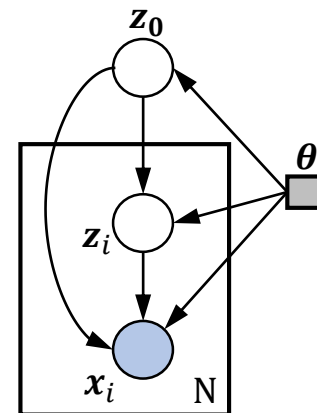
$$w_{mn} | \varphi, z_{mn} \sim \text{categorical}(\varphi_{z_{mn}})$$



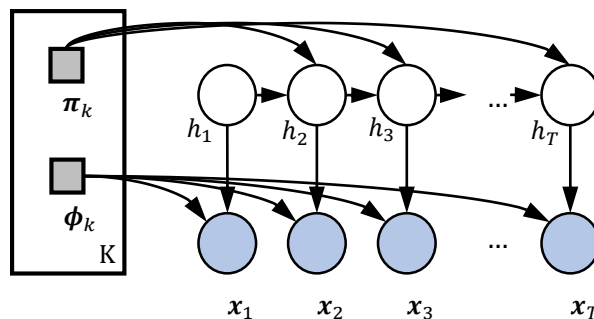
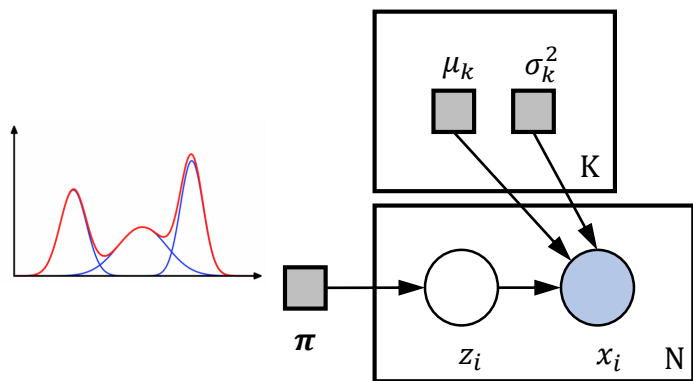
Basic Tasks in Probabilistic Generative Models

- **Parameter Learning:** fit the model to the dataset
 - **Maximum Likelihood Estimation** for the parameters θ

$$\theta_{MLE} = \operatorname{argmax}_{\theta} p_{\theta}(\mathbf{x}) \quad p_{\theta}(\mathbf{z}|\mathbf{x}) = \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{x})}$$



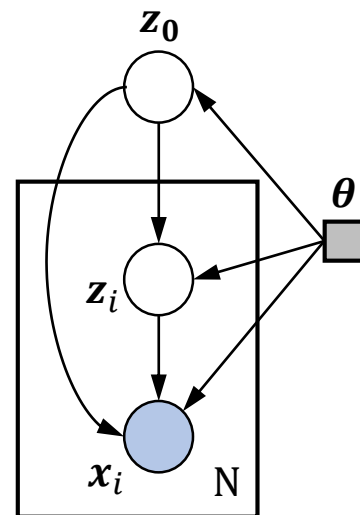
- **Inference:** compute unknown probability distributions
 - **Marginal** distribution of observations \mathbf{x}
 - **Posterior** distribution of latent variable \mathbf{z}



The Difficulty of Learning and Inference

- **Posterior** distribution of latent variables $\mathbf{z} = \{\mathbf{z}_i\}$, $i = 0 \dots N$
- **Marginal** distribution of observations $\mathbf{x} = \{\mathbf{x}_i\}$, $i = 1 \dots N$
- MLE for model parameters θ : $\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} p_{\theta}(\mathbf{x})$

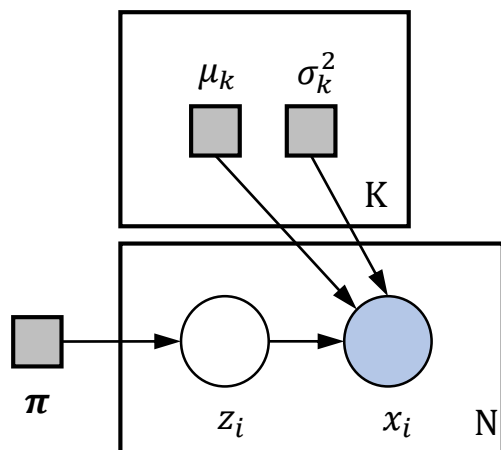
$$p_{\theta}(\mathbf{z}|\mathbf{x}) = \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{x})} = \frac{\overbrace{p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z})}^{\text{likelihood prior}}}{\underbrace{\int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}}_{\text{marginal/evidence}}}$$



For many models, the integral (or sum) is **intractable**:

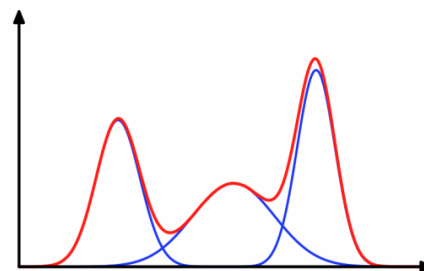
- Unavailable in closed form, or
- Requires exponential time to compute.

A Not-So-Good Example: Univariate GMM



$$z_i \sim \text{categorical}(\pi_1, \dots, \pi_K), \quad i = 1, \dots, N,$$

$$x_i | z_i \sim \mathcal{N}(\mu_{z_i}, \sigma_{z_i}^2) \quad i = 1, \dots, N.$$



$$p_{\theta}(\mathbf{x}, \mathbf{z}) = \prod_{i=1}^N p_{\theta}(x_i, z_i), \quad \text{where } \theta = \{\pi_k, \mu_k, \sigma_k^2: k = 1 \dots K\}$$

$$\begin{aligned} p_{\theta}(\mathbf{x}) &= \sum_{\mathbf{z}} p_{\theta}(\mathbf{x}, \mathbf{z}) = \prod_{i=1}^N \sum_{z_i=1}^K p_{\theta}(z_i) p_{\theta}(x_i | z_i) \\ &= \prod_{i=1}^N \sum_{k=1}^K \pi_k \cdot \mathcal{N}(x_i; \mu_k, \sigma_k^2) \end{aligned}$$

$O(K^N)$

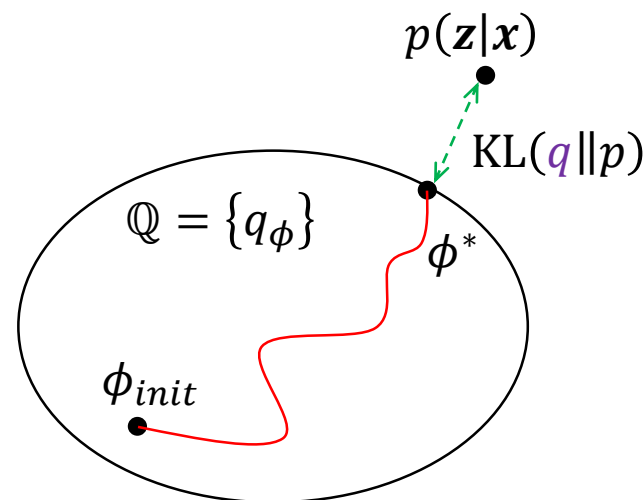
Variational Inference

- Task: computing the posterior of latent variables \mathbf{z} : $p(\mathbf{z}|\mathbf{x})$
- **Variational Inference** turns inference into an **optimization** problem
 - Set up a **family** of approximate densities $\mathbb{Q} = \{q_\phi\}$ over the latent variables \mathbf{z}
 - Find the member $q^*(\mathbf{z}|\mathbf{x}) \in \mathbb{Q}$ that is closest to the exact posterior
- Specifically, given a family of distributions $\mathbb{Q} = \{q_\phi\}$, try to find:

$$q^*(\mathbf{z}|\mathbf{x}) = \arg \min_{q_\phi \in \mathbb{Q}} \text{KL}(q_\phi(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}|\mathbf{x}))$$

equivalent to find:

$$\begin{aligned} \phi^* &= \arg \min_{\phi} \text{KL}(q_\phi(\mathbf{z}) \| p(\mathbf{z}|\mathbf{x})) \\ &= \arg \max_{\phi} \text{ELBO}(q_\phi) \end{aligned}$$



Maximizing the ELBO

- Variational Inference: maximizing the **E**vidence **L**ower **B**ound
 - Maximizing the ELBO is equivalent to minimizing the KL divergence

$$\log p(\mathbf{x}) = \downarrow \text{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x})) + \text{ELBO}(q_\phi) \uparrow$$

$$\geq \text{ELBO}(q_\phi) = \mathbb{E}_{q_\phi} \left[\log \frac{p(\mathbf{z}, \mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]$$

\downarrow
 $\int q_\phi(\mathbf{z}|\mathbf{x})(...)d\mathbf{z}$

- We want to optimize the ELBO using stochastic gradient ascent
- The key is how to take the gradient ∇_ϕ of the **intractable** expectation
 - The general REINFORCE gradient provides an estimator (with high variance)

$$\nabla \mathbb{E}_p[f(\dots)] = \mathbb{E}_p[\nabla \log p(x) f(\dots)] \approx \sum_{i=1}^M \nabla \log p(x_i) f(\dots), \quad x_i \sim p(x)$$

Another Usage of the ELBO

- Another task is finding the MLE for θ : $\theta_{MLE} = \operatorname{argmax}_{\theta} p_{\theta}(\mathbf{x})$
- Recall that $p_{\theta}(\mathbf{x})$ is generally **intractable** for complex models
- Instead, we can maximize the ELBO, which is a lower bound of $p_{\theta}(\mathbf{x})$

$$\log p_{\theta}(\mathbf{x}) = \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) + \mathbf{ELBO}(q_{\phi})$$

$$\geq \mathbf{ELBO}(q_{\phi}) = \mathbb{E}_{q_{\phi}} \left[\log \frac{p_{\theta}(\mathbf{z}, \mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]$$

- So the ELBO provides an optimization target for both ϕ and θ
 - Variational parameters ϕ and model/generative parameters θ

$$\mathcal{L}(\phi, \theta) = \mathbb{E}_{q_{\phi}} [\log p_{\theta}(\mathbf{z}, \mathbf{x}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})]$$

- We want to take the gradients $\nabla_{\phi, \theta} \mathcal{L}$ and optimize these parameters

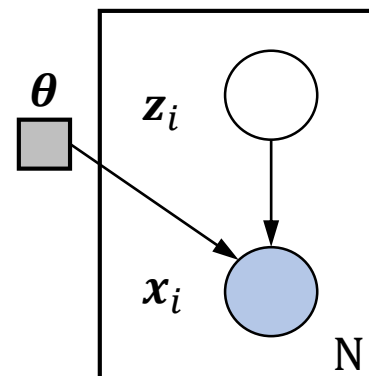
Variational Autoencoder: Model and ELBO

- VAE assumes a **Deep Latent Gaussian Model**

$$\mathbf{z}_i \sim \mathcal{N}(\mathbf{0}, \mathbb{I}), \quad i = 1, \dots, N,$$

$$\mathbf{x}_i | \mathbf{z}_i \sim \mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{z}_i), \boldsymbol{\sigma}_{\boldsymbol{\theta}}^2(\mathbf{z}_i)), \quad i = 1, \dots, N.$$

where $\boldsymbol{\mu}_{\boldsymbol{\theta}}(\cdot)$ and $\boldsymbol{\sigma}_{\boldsymbol{\theta}}(\cdot)$ are neural networks



- Here the global ELBO can be factorized into many local ELBOs
 - One ELBO for each individual data point

$$\mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\theta}) = \sum_{i=1}^N \mathbb{E}_{q_{\boldsymbol{\phi}_i}} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{z}_i, \mathbf{x}_i)}{q_{\boldsymbol{\phi}_i}(\mathbf{z}_i | \mathbf{x}_i)} \right] = \sum_{i=1}^N \mathcal{L}_i(\boldsymbol{\phi}_i, \boldsymbol{\theta})$$

where $\boldsymbol{\phi} = \{\boldsymbol{\phi}_i : i = 1 \dots N\}$

VAE: Optimization of the ELBO (1)

$$\mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\theta}) = \sum_{i=1}^N \mathcal{L}_i(\boldsymbol{\phi}_i, \boldsymbol{\theta}) = \sum_{i=1}^N \mathbb{E}_{q_{\boldsymbol{\phi}_i}} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{z}_i, \mathbf{x}_i)}{q_{\boldsymbol{\phi}_i}(\mathbf{z}_i | \mathbf{x}_i)} \right]$$

$$\nabla_{\boldsymbol{\phi}_i} \mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\phi}_i} \mathcal{L}_i(\boldsymbol{\phi}_i, \boldsymbol{\theta}), \quad i = 1 \dots N$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\theta}) \approx \frac{N}{S} \sum_{s=1}^S \nabla_{\boldsymbol{\theta}} \mathcal{L}_{i_s}(\boldsymbol{\phi}_{i_s}, \boldsymbol{\theta})$$

- They introduced a low-variance gradient estimator for the ELBO
 - Called *reparameterization trick*, or *path-wise gradient estimator*

$$\begin{aligned} \nabla_{\boldsymbol{\theta}, \boldsymbol{\phi}_i} \mathbb{E}_{q_{\boldsymbol{\phi}_i}} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{z}_i, \mathbf{x}_i)}{q_{\boldsymbol{\phi}_i}(\mathbf{z}_i | \mathbf{x}_i)} \right] &= \mathbb{E}_{\mathcal{N}(\boldsymbol{\epsilon})} \left[\nabla_{\boldsymbol{\theta}, \boldsymbol{\phi}_i} \log \frac{p_{\boldsymbol{\theta}}(t(\boldsymbol{\epsilon}, \boldsymbol{\phi}_i), \mathbf{x}_i)}{q_{\boldsymbol{\phi}_i}(t(\boldsymbol{\epsilon}, \boldsymbol{\phi}_i) | \mathbf{x}_i)} \right] \\ &\approx \frac{1}{M} \sum_{k=1}^M \nabla_{\boldsymbol{\theta}, \boldsymbol{\phi}_i} \log \frac{p_{\boldsymbol{\theta}}(t(\boldsymbol{\epsilon}_k, \boldsymbol{\phi}_i), \mathbf{x}_i)}{q_{\boldsymbol{\phi}_i}(t(\boldsymbol{\epsilon}_k, \boldsymbol{\phi}_i) | \mathbf{x}_i)}, \quad \boldsymbol{\epsilon}_k \sim \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{1}) \end{aligned}$$

Assume:

$$\mathbf{z}_i \sim q_{\boldsymbol{\phi}_i}(\mathbf{z}_i | \mathbf{x}_i)$$

\Leftrightarrow

$$\mathbf{z}_i = t(\boldsymbol{\epsilon}, \boldsymbol{\phi}_i, \mathbf{x}_i),$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{1})$$

VAE: Optimization of the ELBO (2)

$$\mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\theta}) = \sum_{i=1}^N \mathcal{L}_i(\boldsymbol{\phi}_i, \boldsymbol{\theta}) = \sum_{i=1}^N \mathbb{E}_{q_{\boldsymbol{\phi}_i}} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{z}_i, \mathbf{x}_i)}{q_{\boldsymbol{\phi}_i}(\mathbf{z}_i | \mathbf{x}_i)} \right]$$

$$\nabla_{\boldsymbol{\phi}_i} \mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\phi}_i} \mathcal{L}_i(\boldsymbol{\phi}_i, \boldsymbol{\theta}), \quad i = 1 \dots N$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\theta}) \approx \frac{N}{S} \sum_{s=1}^S \nabla_{\boldsymbol{\theta}} \mathcal{L}_{i_s}(\boldsymbol{\phi}_{i_s}, \boldsymbol{\theta})$$

- There are N local variational parameters $\boldsymbol{\phi}_i$ to be optimized
 - For large datasets N is large, this is too expensive
- VAE uses a shared **inference network** to predict $\boldsymbol{\phi}_i$ from \mathbf{x}_i instead
 - $g(\mathbf{x}_i) = \boldsymbol{\phi}_i$
 - Also called *recognition network*
 - this technique is known as **A**mortized **V**ariational **I**nference

VAE: Optimization of the ELBO (3)

$$\mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\theta}) = \sum_{i=1}^N \mathcal{L}_i(g_{\boldsymbol{\phi}}(\mathbf{x}_i), \boldsymbol{\theta}) = \sum_{i=1}^N \mathbb{E}_{q_{g_{\boldsymbol{\phi}}(\mathbf{x}_i)}} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{z}_i, \mathbf{x}_i)}{q_{g_{\boldsymbol{\phi}}(\mathbf{x}_i)}(\mathbf{z}_i | \mathbf{x}_i)} \right]$$

$$\nabla_{\boldsymbol{\phi}} \mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\theta}) \approx \frac{N}{S} \sum_{s=1}^S \nabla_{\boldsymbol{\phi}} \mathcal{L}_{i_s}(g_{\boldsymbol{\phi}}(\mathbf{x}_{i_s}), \boldsymbol{\theta})$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\theta}) \approx \frac{N}{S} \sum_{s=1}^S \nabla_{\boldsymbol{\theta}} \mathcal{L}_{i_s}(g_{\boldsymbol{\phi}}(\mathbf{x}_{i_s}), \boldsymbol{\theta})$$

- We treat the parameters of the inference network g as $\boldsymbol{\phi}$
 - Instead of $\boldsymbol{\phi} = \{\boldsymbol{\phi}_i : i = 1 \dots N\}$
- The reparameterization trick gives the gradients $\nabla_{\boldsymbol{\phi}_i}$ for $\boldsymbol{\phi}_i$
- Since $g_{\boldsymbol{\phi}}(\mathbf{x}_i) = \boldsymbol{\phi}_i$, $\nabla_{\boldsymbol{\phi}_i}$ can be propagated back into $g_{\boldsymbol{\phi}}$

Contributions of VAE

- Proposed the reparameterization tricks (a.k.a. path-wise estimator)
 - which yield a low-variance gradient estimator for variational inference
- Introduced the inference network (a.k.a. *recognition network*)
 - This scheme is also called *amortized variational inference*
 - In traditional VI, we have to optimize local variational parameters for every datapoint
 - In amortized VI, we optimize shared parameters of the inference network instead

$$\mathbf{z}_i | \mathbf{x}_i \sim q_\phi(\mathbf{z}_i | \mathbf{x}_i) = \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\sigma}_i^2) = \mathcal{N}(\boldsymbol{\mu}_\phi(\mathbf{x}_i), \boldsymbol{\sigma}_\phi^2(\mathbf{x}_i) \mathbb{I})$$

where $\boldsymbol{\mu}_\phi(\cdot)$ and $\boldsymbol{\sigma}_\phi(\cdot)$ are neural networks

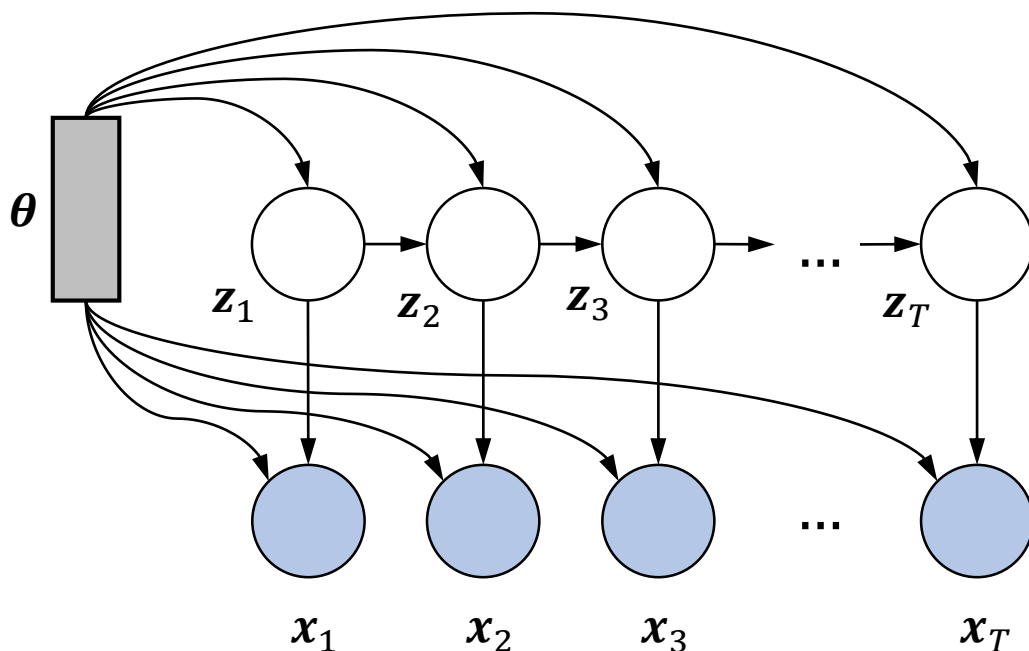
- Co-training of variational parameters and generative parameters

Outline

- Short (Re)Introduction to VI and VAE
- **State Space Models and Stochastic RNNs**
- Variational Inference of SSM
- (Variational) Sequential Monte Carlo
- Conclusions

State Space Models

- **SSM** are a general class of probabilistic sequential models
 - Usually we assume that the Markovian property holds for state transitions
 - Non-Markovian transitions can be convert to Markovian transitions by memorization



$$z_1 \sim \pi_{\theta}(\cdot)$$

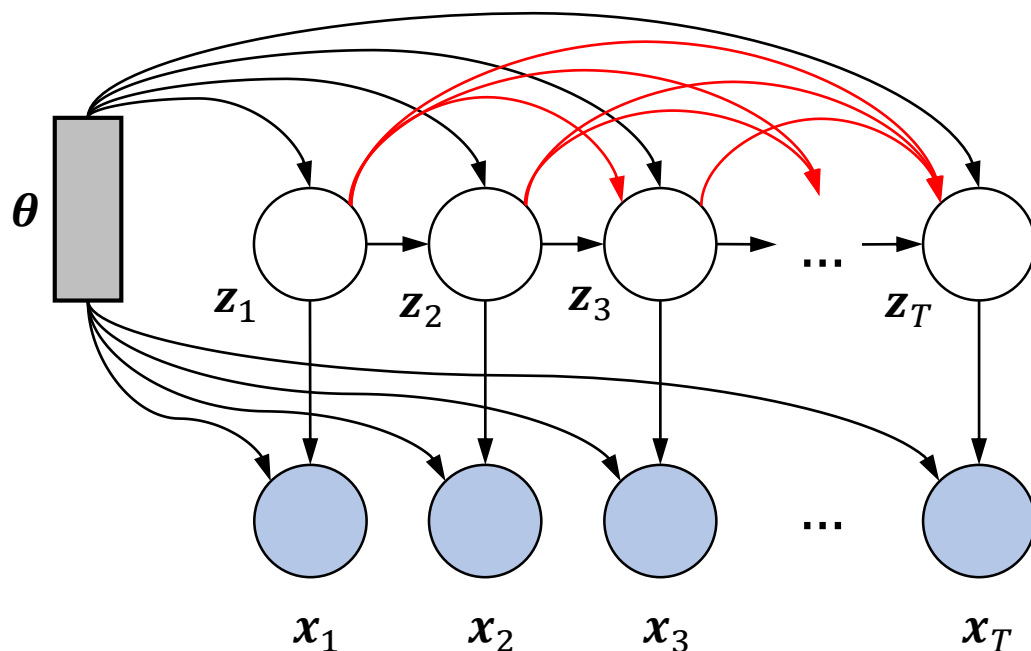
$$z_t | z_{t-1} \sim f_{\theta}(z_{t-1})$$

$$x_t | z_t \sim g_{\theta}(z_t)$$

where $t = 2 \dots T$

State Space Models: Non-Markovian Variant

- **SSM** are a general class of probabilistic sequential models
 - Usually we assume that the Markovian property holds for state transitions
 - Non-Markovian transitions can be convert to Markovian transitions by memorization
 - However, explicitly modelling the **non-Markovian** property can be beneficial



$$z_1 \sim \pi_{\theta}(\cdot)$$

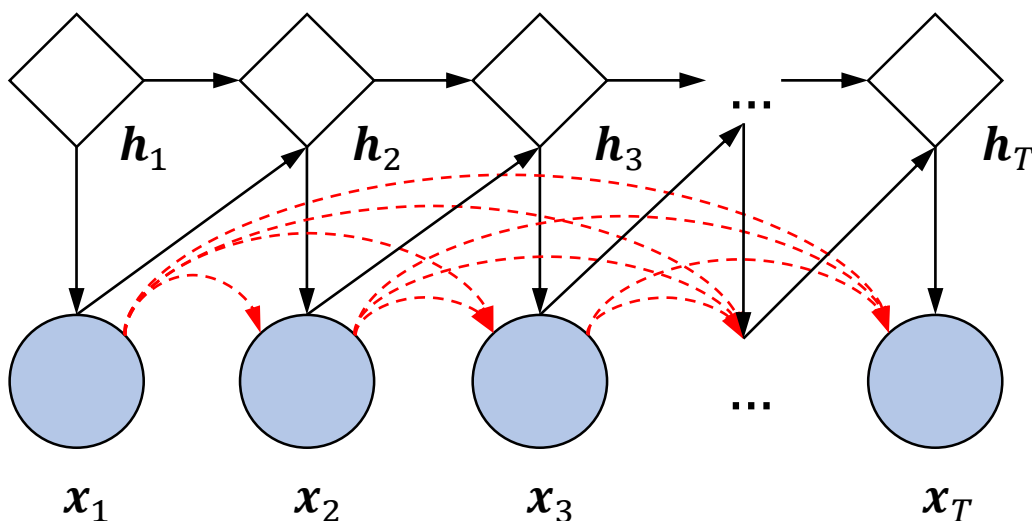
$$z_t | z_{1:t-1} \sim f_{\theta}(z_{1:t-1})$$

$$x_t | z_t \sim g_{\theta}(z_t)$$

where $t = 2 \dots T$

Recurrent Neural Network

- In general, RNN is a deterministic model
- However in language modeling, RNN is exactly a probabilistic model
 - The output at each timestep is a softmax, measured by cross-entropy loss
- A.k.a. **autoregressive** models, in which there are no **latent** variables
 - order sensitive, not-so-good at modeling stochastic transitions



$$x_1 \sim g_\theta(h_1)$$
$$h_t = \text{LSTM}_\theta(h_{t-1}, x_{t-1})$$

$$x_t | x_{1:t-1} \sim g_\theta(\cdot | h_t)$$

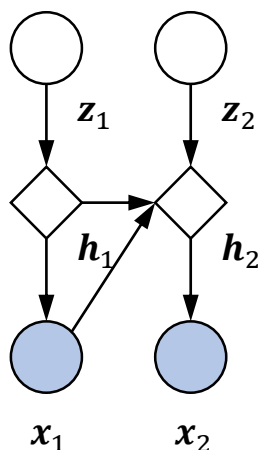
where $t = 2 \dots T$

the actual conditional dependences
are shown in red dashed lines

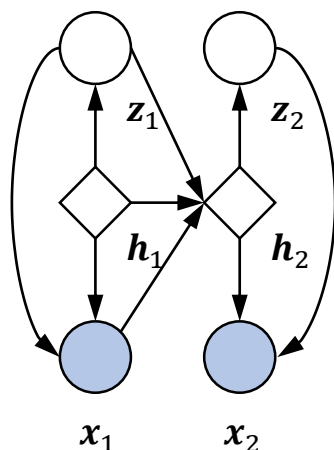
Stochastic RNNs

- Some work focused on injecting random variables into RNNs

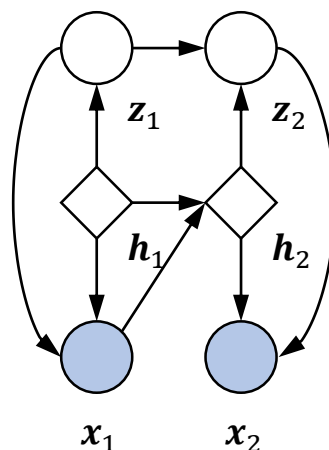
STORN [Arxiv 2014]
Noisin [ICML 2018]



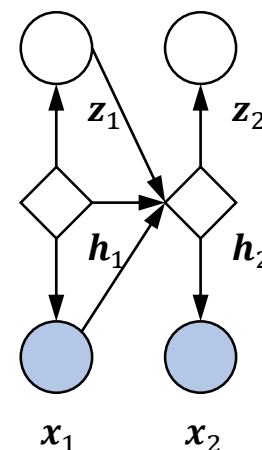
VRNN [NIPS 2015]



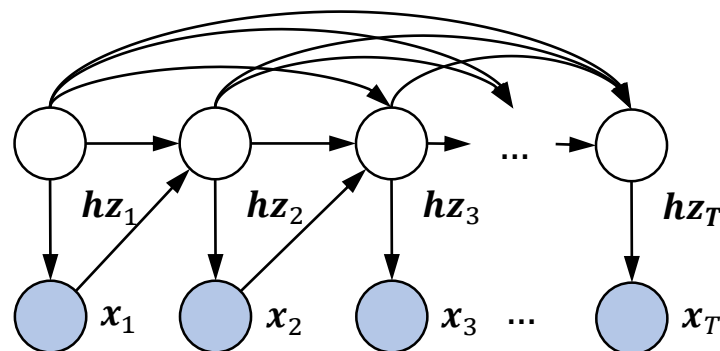
SRNN [NIPS 2016]



Z-Forcing [NIPS 2017]



After merging
the deterministic & stochastic state
at each time step,
almost all of them resemble this:



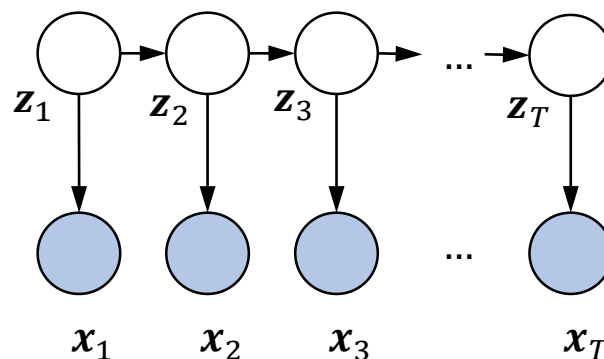
Two Elegant Models that I Like

- Deep Markov Model

$$\mathbf{z}_1 \sim \pi_{\theta}(\cdot)$$

$$\mathbf{z}_t | \mathbf{z}_{t-1} \sim \mathcal{N}(\mu_{\theta}(\mathbf{z}_{t-1}), \sigma_{\theta}^2(\mathbf{z}_{t-1}))$$

$$\mathbf{x}_t | \mathbf{z}_t \sim g_{\theta}(\cdot | \mathbf{z}_t)$$



where $\mu_{\theta}(\cdot)$ and $\sigma_{\theta}^2(\cdot)$ are (gated) neural networks

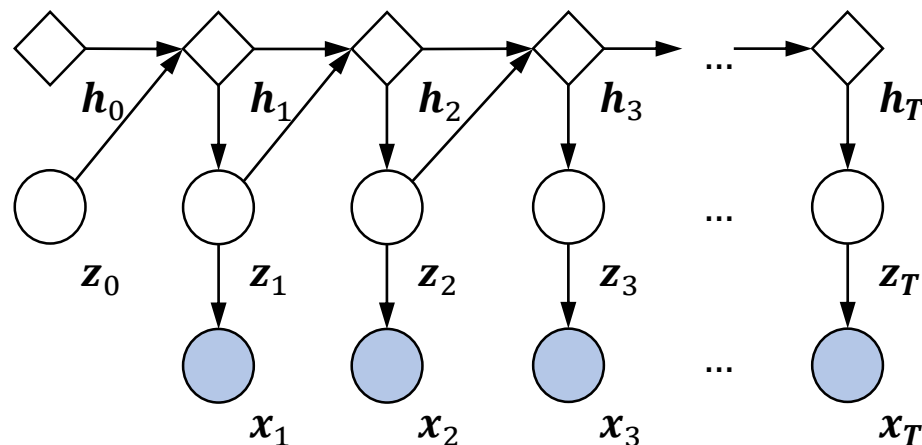
- LSTM State Space Model

$$\mathbf{z}_0 \sim \pi_{\theta}(\cdot)$$

$$\mathbf{h}_t = \text{LSTM}_{\theta}(\mathbf{h}_{t-1}, \mathbf{z}_{t-1})$$

$$\mathbf{z}_t | \mathbf{z}_{1:t-1} \sim f_{\theta}(\cdot | \mathbf{h}_t)$$

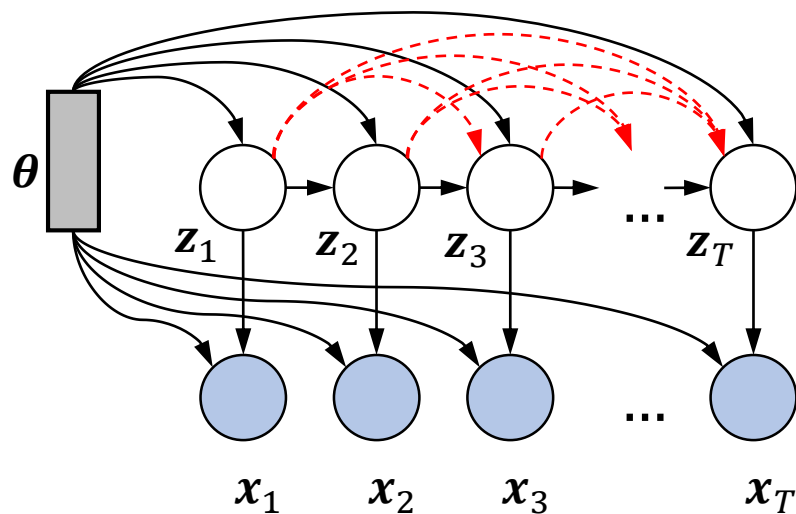
$$\mathbf{x}_t | \mathbf{z}_t \sim g_{\theta}(\cdot | \mathbf{z}_t)$$



Structured Inference Networks for Nonlinear State Space Models. Krishnan, et al. AAAI 2017
State Space LSTM Models with Particle MCMC Inference. Zheng, et al. Arxiv 2017

Using Probabilistic Sequential Models

- Given these models, our basic tasks are:
 - Model learning: MLE for the model parameters θ : $\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} p_{\theta}(x)$
 - Bayesian inference: posterior distribution of latent variables z : $p_{\theta}(z_{1:T} | x_{1:T})$
 - Prediction: conditional distribution of future states: $p_{\theta}(z_{T+1:T+L} | x_{1:T})$
 - And future observations: $p_{\theta}(x_{T+1:T+L} | x_{1:T})$
- Next: do learning and inference using
 - Variational Inference
 - (Variational) Sequential Monte Carlo

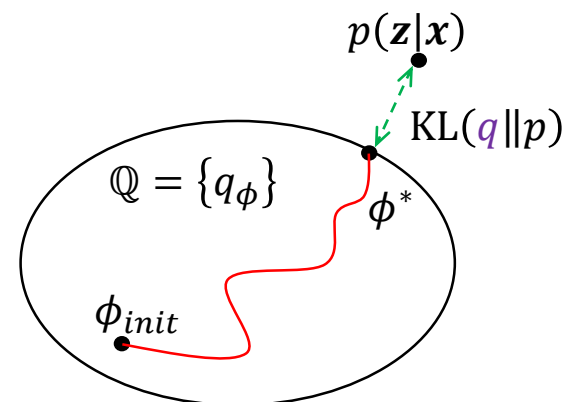
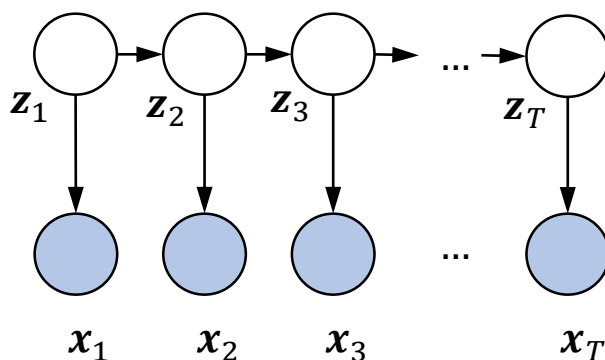


Outline

- Short (Re)Introduction to VI and VAE
- State Space Models and Stochastic RNNs
- **Variational Inference of SSM**
- (Variational) Sequential Monte Carlo
- Conclusions

Variational Inference of SSM

- Recall that in VI, we need to design a variational family $\mathbb{Q} = \{q_\phi(\mathbf{z}|\mathbf{x})\}$
 - Or equivalently, a family of inference networks in VAE
 - A good variational family should be expressive enough
 - To include the true posterior or some element close enough to the posterior
 - We can look at the factorized form of true posterior and imitate it

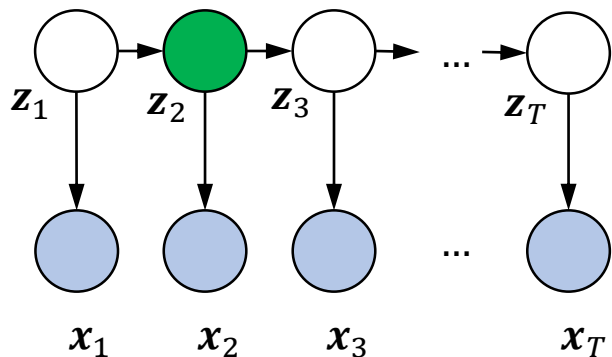


$$p_\theta(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = \pi_\theta(\mathbf{z}_1) g_\theta(\mathbf{x}_1 | \mathbf{z}_1) \prod_{t=2}^T f_\theta(\mathbf{z}_t | \mathbf{z}_{t-1}) g_\theta(\mathbf{x}_t | \mathbf{z}_t)$$

$$p_\theta(\mathbf{z}_{1:T} | \mathbf{x}_{1:T}) = p_\theta(\mathbf{z}_1 | \mathbf{x}_{1:T}) \prod_{t=2}^T p_\theta(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{x}_{t:T})$$

True Posterior for SSMs

- Using the standard **d-separation** criterion for Bayesian networks:



$$\mathbf{z}_{t+1} \perp \mathbf{z}_{1:t-1} \mid \mathbf{z}_t$$

$$\mathbf{z}_t \perp \mathbf{x}_{1:t-1} \mid \mathbf{z}_{t-1}$$

$$p_{\theta}(\mathbf{z}_{1:T} \mid \mathbf{x}_{1:T})$$

$$= p_{\theta}(\mathbf{z}_1 \mid \mathbf{x}_{1:T}) \prod_{t=2}^T p_{\theta}(\mathbf{z}_t \mid \mathbf{z}_{1:t-1}, \mathbf{x}_{1:T})$$

$$= p_{\theta}(\mathbf{z}_1 \mid \mathbf{x}_{1:T}) \prod_{t=2}^T p_{\theta}(\mathbf{z}_t \mid \mathbf{z}_{t-1}, \mathbf{x}_{1:T})$$

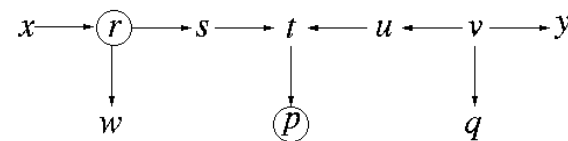
$$= p_{\theta}(\mathbf{z}_1 \mid \mathbf{x}_{1:T}) \prod_{t=2}^T p_{\theta}(\mathbf{z}_t \mid \mathbf{z}_{t-1}, \mathbf{x}_{t:T})$$

information from the **future** is crucial

Rule 1: x and y are d-connected if there is an unblocked path between them.

Rule 2: x and y are d-connected, conditioned on a set Z of nodes, if there is a collider-free path between x and y that traverses no member of Z .

Rule 3: If a collider is a member of the conditioning set Z , or has a descendant in Z , then it no longer blocks any path that traces this collider.



Inference Network for SSM

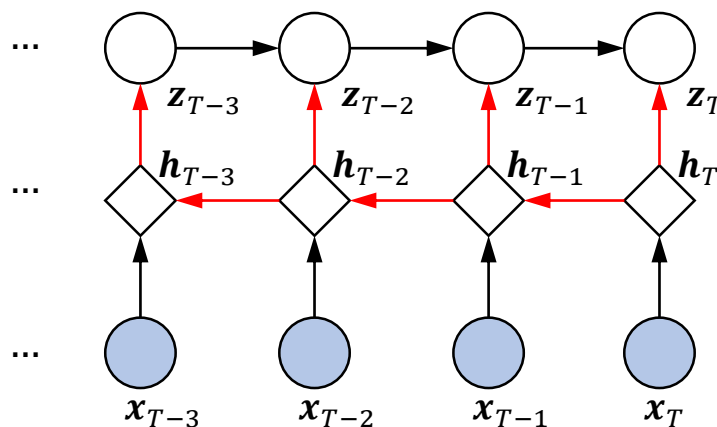
- Let the form of $q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$ be exactly the same as $p_{\theta}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$:

$$p_{\theta}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) = p_{\theta}(\mathbf{z}_1|\mathbf{x}_{1:T}) \prod_{t=2}^T p_{\theta}(\mathbf{z}_t|\mathbf{z}_{t-1}, \mathbf{x}_{t:T})$$

$$q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) = q_{\phi}(\mathbf{z}_1|\mathbf{x}_{1:T}) \prod_{t=2}^T q_{\phi}(\mathbf{z}_t|\mathbf{z}_{t-1}, \mathbf{x}_{t:T})$$

$$\text{where } q_{\phi}(\mathbf{z}_t|\mathbf{z}_{t-1}, \mathbf{x}_{t:T}) = \mathcal{N}(\mu_{\phi}(\mathbf{z}_{t-1}, \mathbf{x}_{t:T}), \sigma_{\phi}^2(\mathbf{z}_{t-1}, \mathbf{x}_{t:T}))$$

- How to implement $q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$? Well, it looks like a **Backward RNN** !
 - $\mu_{\phi}(\cdot)$ and $\sigma_{\phi}^2(\cdot)$ use backward RNNs as their backbones



Learning and Inference

- We have constructed the variational family backbone by RNNs
- So we have the ELBO:

$$\text{ELBO}(q_\phi) = \mathbb{E}_{q_\phi} \left[\log \frac{p_\theta(\mathbf{z}, \mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]$$

- Now we can perform stochastic gradient ascent in both ϕ and θ
 - Using the gradient estimator provided by reparameterization tricks
 - The ELBO can be factorized to reduce the variance of gradient estimator
- Given the parameters θ and ϕ , we can:
 - Generating new sequences (using generative parameters θ)
 - Performing inference (using variational parameters ϕ)
 - Predicting the future from current observations (using both ϕ and θ)

Outline

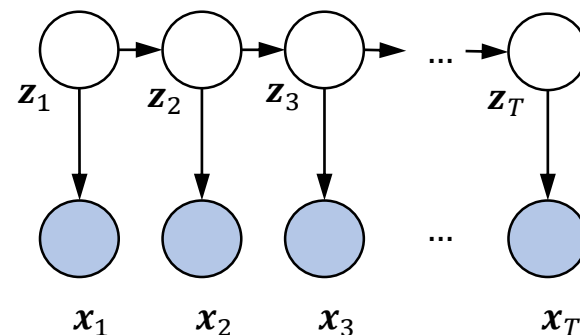
- Short (Re)Introduction to VI and VAE
- State Space Models and Stochastic RNNs
- Variational Inference of SSM
- **(Variational) Sequential Monte Carlo**
- Conclusions

Sequential Monte Carlo

- **SMC** methods are standard tools for nonlinear dynamical systems
 - a.k.a. **P**article **F**ilters
 - Widely used in many fields, including physics, engineering, robotics, finance, ...

- Consider the state space model

$$p_{\theta}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = \mu_{\theta}(\mathbf{z}_1) g_{\theta}(\mathbf{x}_1 | \mathbf{z}_1) \prod_{t=2}^T f_{\theta}(\mathbf{z}_t | \mathbf{z}_{t-1}) g_{\theta}(\mathbf{x}_t | \mathbf{z}_t)$$



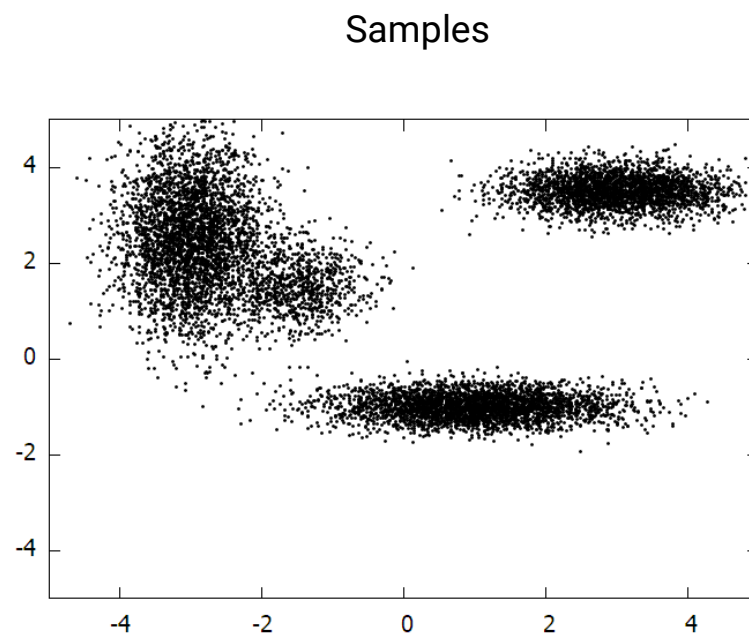
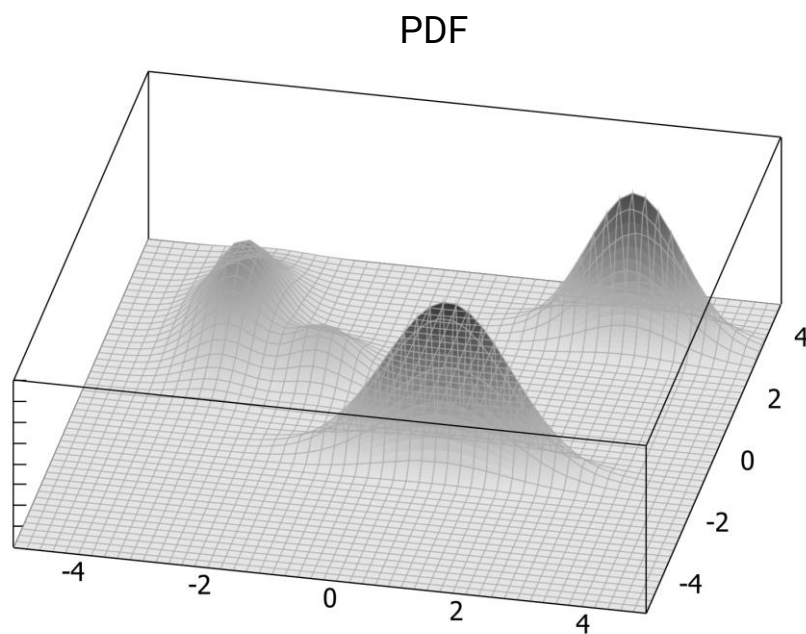
- PFs are mainly used for solving the **filtering** & **smoothing** problem
 - Filtering: the posterior distribution of current state: $\{p_{\theta}(\mathbf{z}_t | \mathbf{x}_{1:t})\}_{t \geq 1}$
 - Smoothing: the posterior distribution of past states: $\{p_{\theta}(\mathbf{z}_{1:t} | \mathbf{x}_{1:t})\}_{t \geq 1}$
 - or $p_{\theta}(\mathbf{z}_{m:n} | \mathbf{x}_{1:t}), 1 \leq m \leq n < t$
 - While SMCs can be used in more general case
 - e.g., non-Markovian SSM
- (most literature writes $\mathbf{x} \rightarrow \mathbf{y}$ instead of $\mathbf{z} \rightarrow \mathbf{x}$)

Monte Carlo Basics (1)

- Monte Carlo approximation of the intractable density π

$$\mathbb{E}_{\pi}[f(x)] = \int \pi(x)f(x) dx \approx \mathbb{E}_{\hat{\pi}}[f(x)] = \frac{1}{N} \sum_{i=1}^N f(x_i), \quad x_i \sim \pi(x), i = 1 \dots N$$

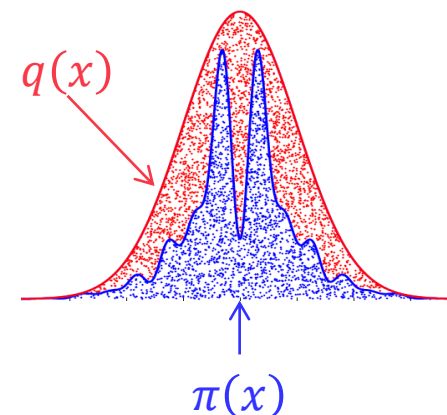
- i.e., approximate π by its empirical distribution $\hat{\pi}$: $\hat{\pi}(x) = \frac{1}{N} \sum_{i=1}^N \delta_{x_i}(x)$



Monte Carlo Basics (2)

- MC assumes that it is possible to generate $x_i \sim \pi(x), i = 1 \dots N$
 - i.e., to generate i.i.d. samples from the target distribution
 - However, we only have access to uniform PRNGs
 - For some simple distributions, this can be done by inverting the CDF
 - Inverting the CDF: $u \sim \mathbf{U}(0, 1)$, then $F^{-1}(u) \sim \pi(x)$
- For distributions w/o tractable CDF, **rejection sampling** can be used

```
given proposal  $q(x)$  s.t.  $\pi(x) \leq Cq(x)$ :  
for  $i = 1 \dots N$ :  
  repeat  
    sample  $\tilde{x} \sim q(x)$   
    sample  $u \sim \mathbf{U}(0, 1)$   
  until  $u \leq \pi(\tilde{x})/Cq(\tilde{x})$   
  set  $x_i = \tilde{x}$ 
```



- the *acceptance probability* decays exponentially as the dimension increases

Importance Sampling

$$\begin{aligned}\mathbb{E}_\pi[f(x)] &= \int \pi(x)f(x) dx = \int q(x) \frac{\pi(x)}{q(x)} f(x) dx = \mathbb{E}_q \left[\frac{\pi(x)}{q(x)} f(x) \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \omega(x_i) f(x_i), \quad \omega(x) = \frac{\pi(x)}{q(x)}, \quad x_i \sim q(x), i = 1 \dots N\end{aligned}$$

- Importance Sampling also uses a proposal distribution that is easy to sample
- Rather than discarding some samples, IS makes use of all samples but assigns individual **importance weights** to them
- How about if we can only evaluate $\tilde{\pi}(x)$ rather than $\pi(x) = \frac{\tilde{\pi}(x)}{\int \tilde{\pi}(x) dx}$?
 - Let $\tilde{\omega}(x) = \tilde{\pi}(x)/q(x)$, $Z = \int \tilde{\pi}(x) dx$. Then:

$$Z = \int \tilde{\pi}(x) dx = \int q(x) \frac{\tilde{\pi}(x)}{q(x)} dx = \mathbb{E}_q[\tilde{\omega}(x)] \approx \frac{1}{N} \sum_{i=1}^N \tilde{\omega}(x_i), \quad x_i \sim q(x)$$

$$\frac{1}{N} \omega(x_i) = \frac{1}{N} \frac{\pi(x_i)}{q(x_i)} = \frac{1}{N} \frac{1}{Z} \frac{\tilde{\pi}(x_i)}{q(x_i)} \approx \frac{\tilde{\omega}(x_i)}{\sum_{k=1}^N \tilde{\omega}(x_k)} \triangleq w_i, \quad \mathbb{E}_\pi[f(x)] \approx \sum_{i=1}^N w_i f(x_i)$$

self-normalized importance sampling estimator

Sequential Importance Sampling (1)

- Aim: approximate a sequence of target distributions $\{\pi_t(x_{1:t})\}_{t \geq 1}$
 - $\pi_1(x_1), \pi_2(x_{1:2}), \pi_3(x_{1:3}), \dots$
 - e.g., in state space models: $p_\theta(z_1|x_1), p_\theta(z_{1:2}|x_{1:2}), p_\theta(z_{1:3}|x_{1:3}), \dots$

- IS: approximate $\mathbb{E}_{\pi_t}(f_t)$ by N weighted samples $\{x_{1:t}^{(i)}: i = 1 \dots N\}$:

$$\mathbb{E}_{\pi_t}(f_t) \approx \sum_{i=1}^N w_t^{(i)} f_t(x_{1:t}^{(i)})$$
$$Z_t = \int \tilde{\pi}_t(x_{1:t}) dx_{1:t} \approx \frac{1}{N} \sum_{i=1}^N \tilde{w}(x_{1:t}^{(i)})$$

- To generate such weighted samples, we need to sampling from a sequence of proposal distributions: $\{q_t(x_{1:t})\}_{t \geq 1}$
 - Randomly sampling from high-dimensional distributions can be inefficient
 - Requiring numerous samples to achieve desired accuracy & computationally expensive
- Basic idea of SIS: reuse $\{x_{1:\underline{t}}^{(i)}: i = 1 \dots N\}$ to build $\{x_{1:\underline{t}+1}^{(i)}: i = 1 \dots N\}$

Sequential Importance Sampling (2)

- **SIS** uses a structured proposal to reuse samples:

$$\begin{aligned} q_t(x_{1:t+1}) &= q_t(x_{1:t})q_{t+1}(x_{t+1}|x_{1:t}) \\ &= q_1(x_1)q_2(x_2|x_1)q_3(x_3|x_{1:2}) \dots q_{t+1}(x_{t+1}|x_{1:t}) \end{aligned}$$

- Given $x_{1:t}^{(i)} \sim q_t(x_{1:t})$, sampling $x_{t+1}^{(i)} \sim q_{t+1}(x_{t+1}|x_{1:t}^{(i)})$ to obtain $x_{1:t+1}^{(i)}$:

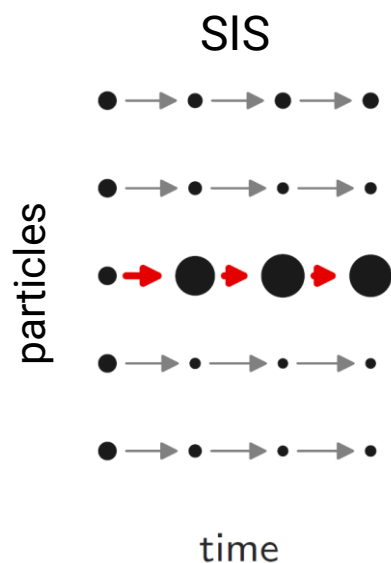
$$x_{1:t+1}^{(i)} = \langle x_{1:t}^{(i)}, x_{t+1}^{(i)} \rangle$$

- And update the unnormalized weights as:

$$\begin{aligned} \tilde{\omega}_{t+1}(x_{1:t+1}) &= \frac{\tilde{\pi}_{t+1}(x_{1:t+1})}{q_{t+1}(x_{1:t+1})} = \frac{\tilde{\pi}_{t+1}(x_{1:t+1})}{q_t(x_{1:t})q_{t+1}(x_{t+1}|x_{1:t})} \\ &= \frac{\tilde{\pi}_t(x_{1:t})}{q_t(x_{1:t})} \frac{\tilde{\pi}_{t+1}(x_{1:t+1})}{\tilde{\pi}_t(x_{1:t})q_{t+1}(x_{t+1}|x_{1:t})} = \tilde{\omega}_t(x_{1:t}) \frac{\tilde{\pi}_{t+1}(x_{1:t+1})}{\tilde{\pi}_t(x_{1:t})q_{t+1}(x_{t+1}|x_{1:t})} \end{aligned}$$

Sequential Importance Sampling (3)

- The overall SIS algorithm



input: number of samples N ,
 unnormalized target distributions $\{\tilde{\pi}_t(x_{1:t})\}_{t \geq 1}$,
 proposals $q_1(\cdot)$ and $\{q_t(\cdot | x_{1:t-1})\}_{t > 1}$

for $i = 1 \dots N$:

sample $x_1^{(i)} \sim q_1(\cdot)$, set $\tilde{w}_1(x_1^{(i)}) = \frac{\tilde{\pi}_1(x_1^{(i)})}{q_1(x_1^{(i)})}$

for $t = 2 \dots T$:

for $i = 1 \dots N$:

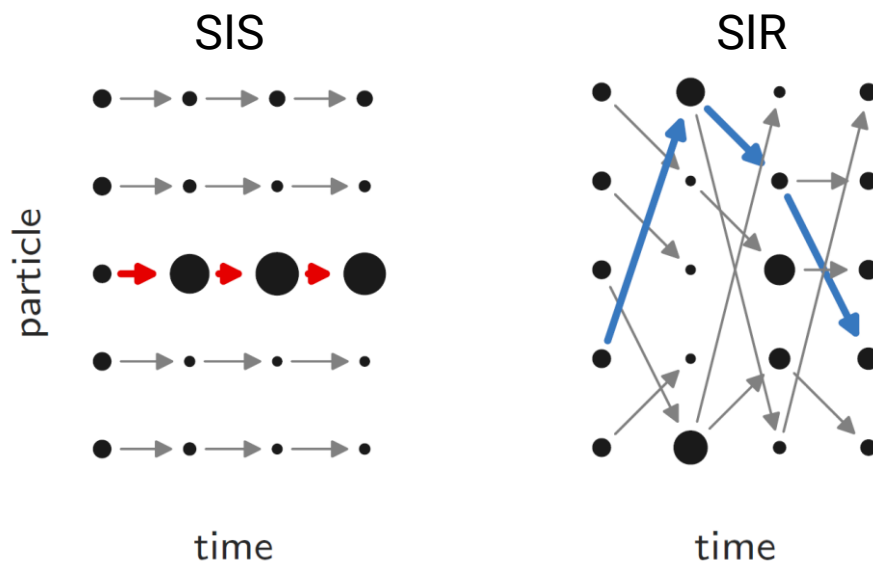
sample $x_t^{(i)} \sim q_t(\cdot | x_{1:t-1}^{(i)})$, set $x_{1:t}^{(i)} = \langle x_{1:t-1}^{(i)}, x_t^{(i)} \rangle$

set $\tilde{w}_t(x_{1:t}^{(i)}) = \tilde{w}_{t-1}(x_{1:t-1}^{(i)}) \frac{\tilde{\pi}_t(x_{1:t}^{(i)})}{\tilde{\pi}_{t-1}(x_{1:t-1}^{(i)}) q_t(x_t^{(i)} | x_{1:t-1}^{(i)})}$

output: $\{(x_{1:T}^{(i)}, w_T^{(i)}) : i = 1 \dots N\}$, $\hat{Z}_T = \frac{1}{N} \sum_{i=1}^N \tilde{w}_T(x_{1:T}^{(i)})$

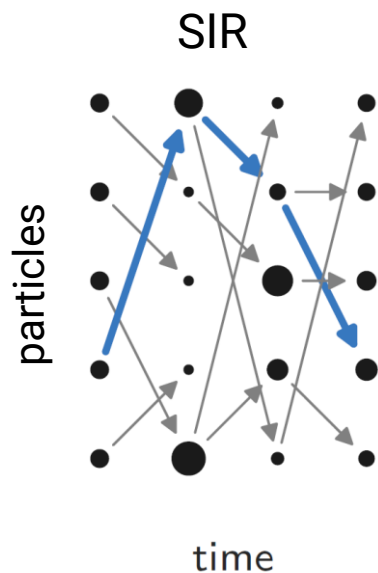
Sequential Importance Resampling

- Problem of SIS: the variance of $\{\bar{\omega}_t^{(i)}\}$ tends to grow as t increases
 - All the mass will concentrate on a few samples
 - Converge to $\bar{\omega}_t^{(i)} = 1$ for some i and others being zero
- Resampling to rescue: **Sequential Importance Resampling**
 - Using resampling to focus on promising samples/trajectories/particles



Sequential Monte Carlo

- The general SMC algorithm



input: number of samples N ,
 unnormalized target distributions $\{\tilde{\pi}_t(x_{1:t})\}_{t \geq 1}$,
 proposals $q_1(\cdot)$ and $\{q_t(\cdot | x_{1:t-1})\}_{t > 1}$

for $i = 1 \dots N$:

sample $x_1^{(i)} \sim q_1(\cdot)$, set $\tilde{\omega}_1(x_1^{(i)}) = \frac{\tilde{\pi}_1(x_1^{(i)})}{q_1(x_1^{(i)})}$, $w_1^{(i)} = \frac{\tilde{\omega}_1(x_1^{(i)})}{\sum_{k=1}^N \tilde{\omega}_1(x_1^{(k)})}$

for $t = 2 \dots T$:

for $i = 1 \dots N$:

sample ancestor index $a_{t-1}^{(i)} \sim \text{categorical}(\cdot | w_{t-1}^{(i)}, \dots, w_{t-1}^{(N)})$

sample $x_t^{(i)} \sim q_t(\cdot | x_{1:t-1}^{a_{t-1}^{(i)}})$, set $x_{1:t}^{(i)} = \langle x_{1:t-1}^{a_{t-1}^{(i)}}, x_t^{(i)} \rangle$

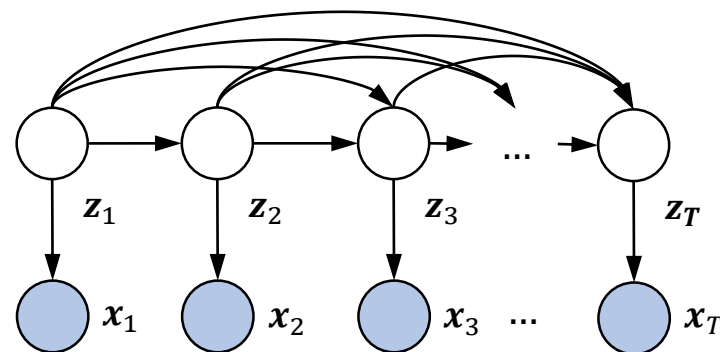
set $w_t^{(i)} \propto \tilde{\omega}_t(x_{1:t}^{(i)}) = \tilde{\pi}_t(x_{1:t}^{(i)}) / [\tilde{\pi}_{t-1}(x_{1:t-1}^{a_{t-1}^{(i)}}) q_t(x_t^{(i)} | x_{1:t-1}^{a_{t-1}^{(i)}})]$

output: $\{(x_{1:T}^{(i)}, w_T^{(i)}) : i = 1 \dots N\}$, $\hat{Z}_T = \prod_{t=1}^T \left[\frac{1}{N} \sum_{i=1}^N \tilde{\omega}_t(x_{1:t}^{(i)}) \right]$

Sequential Monte Carlo for SSM

- Consider the non-Markovian SSM

$$p_{\theta}(\mathbf{z}_{1:T}, \mathbf{x}_{1:T}) = \mu_{\theta}(\mathbf{z}_1) g_{\theta}(\mathbf{x}_1 | \mathbf{z}_1) \prod_{t=2}^T f_{\theta}(\mathbf{z}_t | \mathbf{z}_{1:t-1}) g_{\theta}(\mathbf{x}_t | \mathbf{z}_t)$$



- We set the unnormalized target distributions for SMC to be:

$$\pi_t(\cdot) = p_{\theta}(\mathbf{z}_{1:t} | \mathbf{x}_{1:t}) = \frac{p_{\theta}(\mathbf{z}_{1:t}, \mathbf{x}_{1:t})}{p_{\theta}(\mathbf{x}_{1:t})} = \frac{\tilde{\pi}_t(\cdot)}{Z_t}$$

- Now:

$$\tilde{\pi}_1(\cdot) = \mu_{\theta}(\mathbf{z}_1) g_{\theta}(\mathbf{x}_1 | \mathbf{z}_1)$$

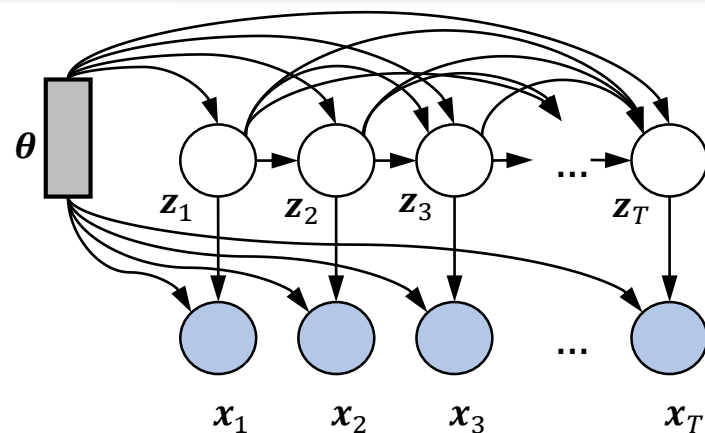
$$\frac{\tilde{\pi}_t(\cdot)}{\tilde{\pi}_{t-1}(\cdot)} = \frac{p_{\theta}(\mathbf{z}_{1:t}, \mathbf{x}_{1:t})}{p_{\theta}(\mathbf{z}_{1:t-1}, \mathbf{x}_{1:t-1})} = f_{\theta}(\mathbf{z}_t | \mathbf{z}_{1:t-1}) g_{\theta}(\mathbf{x}_t | \mathbf{z}_t)$$

- And the proposals are structured as: $q_t(\cdot | \mathbf{x}_{1:t-1}) = q_t(\mathbf{z}_t | \mathbf{z}_{1:t-1}, \mathbf{x}_{1:t})$

SMC for Parameter Learning

- Given the model

$$p_{\theta}(\mathbf{z}_{1:T}, \mathbf{x}_{1:T}) = \mu_{\theta}(\mathbf{z}_1) g_{\theta}(\mathbf{x}_1 | \mathbf{z}_1) \prod_{t=2}^T f_{\theta}(\mathbf{z}_t | \mathbf{z}_{1:t-1}) g_{\theta}(\mathbf{x}_t | \mathbf{z}_t)$$



- We want to learn the MLE for θ : $\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} p_{\theta}(\mathbf{x}_{1:T})$
- Two basic approaches: gradient ascent & **Expectation-Maximization**

(1) Gradient ascent: optimizing the likelihood estimated by SMC

$$\hat{p}_{\theta}(\mathbf{x}_{1:T}) = \prod_{t=1}^T \left[\frac{1}{N} \sum_{i=1}^N \tilde{\omega}_t^{(i)} \right]$$

(2) EM: using the SMC *smoothing* approximation at the E step

$$\mathbb{E}_{p_{\theta}(\mathbf{z}_{1:T} | \mathbf{x}_{1:T})} [\mathcal{L}(\theta, \mathbf{z}_{1:T}, \mathbf{x}_{1:T})] \approx \sum_{i=1}^N w_i \mathcal{L}(\theta, \mathbf{z}_{1:T}^{(i)}, \mathbf{x}_{1:T})$$

ELBO and Monte Carlo Objectives

- Recall the ELBO in VI:

$$\begin{aligned}\mathbf{ELBO}(q_\phi) &= \mathbb{E}_{q_\phi} \left[\log \frac{p_\theta(\mathbf{z}, \mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \leq \log \mathbb{E}_{q_\phi} \left[\frac{p_\theta(\mathbf{z}, \mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \\ &= \log \int q_\phi(\mathbf{z}|\mathbf{x}) \frac{p_\theta(\mathbf{z}, \mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} = \log p_\theta(\mathbf{x})\end{aligned}$$

- Note that $\mathbb{E}_{q_\phi} \left[\frac{p_\theta(\mathbf{z}, \mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] = p_\theta(\mathbf{x})$
 - i.e., $\frac{p_\theta(\mathbf{z}, \mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})}$ is an unbiased estimator of $p_\theta(\mathbf{x})$
- Monte Carlo Objectives** generalize the ELBO
 - Given an unbiased estimator $\hat{p}_\theta(\mathbf{x})$ of $p_\theta(\mathbf{x})$: $\mathbb{E}[\hat{p}_\theta(\mathbf{x})] = p_\theta(\mathbf{x})$, the MCO is:

$$\mathbf{MCO}(\cdot) = \mathbb{E}[\log \hat{p}_\theta(\mathbf{x})] \leq \log \mathbb{E}[\hat{p}_\theta(\mathbf{x})] = \log p_\theta(\mathbf{x})$$

Variational SMC (1)

- **Monte Carlo Objectives:** given $\mathbb{E}[\hat{p}_\theta(\mathbf{x})] = p_\theta(\mathbf{x})$

$$\mathbf{MCO}(\cdot) = \mathbb{E}[\log \hat{p}_\theta(\mathbf{x})] \leq \log \mathbb{E}[\hat{p}_\theta(\mathbf{x})] = \log p_\theta(\mathbf{x})$$

- Recall that SMC provides an estimator for the margin distribution:

$$\hat{p}_\theta(\mathbf{x}_{1:T}) = \prod_{t=1}^T \left[\frac{1}{N} \sum_{i=1}^N \tilde{\omega}_t^{(i)} \right]$$

- It can be shown that $\hat{p}_\theta(\mathbf{x}_{1:T})$ is unbiased

$$\mathbb{E}_{Q_{\text{SMC}}}[\hat{p}_\theta(\mathbf{x}_{1:T})] = p_\theta(\mathbf{x}_{1:T}), \quad \text{where } Q_{\text{SMC}}(x_{1:T}^{1:N}, a_{1:T-1}^{1:N}) =$$

$$\left(\prod_{i=1}^N q_{1,\phi}(x_1^{(i)}) \right) \prod_{t=2}^T \prod_{i=1}^N q_{t,\phi}(x_t^{(i)} | x_{1:t-1}^{a_{t-1}^{(i)}}) \cdot \text{categorical}(a_{t-1}^{(i)} | w_{t-1}^{1:N})$$

Variational SMC (2)

- So we can construct a MCO using SMC:

$$\begin{aligned}\mathbf{MCO}_{\text{SMC}}(\boldsymbol{\theta}, \boldsymbol{\phi}, \mathbf{x}_{1:T}) &= \mathbb{E}_{Q_{\text{SMC}}}[\log \hat{p}_{\boldsymbol{\theta}}(\mathbf{x}_{1:T})] \\ &= \mathbb{E}_{Q_{\text{SMC}}} \left[\prod_{t=1}^T \left(\frac{1}{N} \sum_{i=1}^N \tilde{\omega}_t(\mathbf{x}_{1:t}^{(i)}) \right) \right]\end{aligned}$$

- $\boldsymbol{\theta}$ are the model parameters, $\boldsymbol{\phi}$ are parameters of the proposal distribution
- Now by optimizing this objective with respect to $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$, we can learn the model parameters and the proposals together
 - This objective is (much) tighter than ELBO: $\mathbf{MCO}_{\text{SMC}} \rightarrow \log p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T})$ as $N \rightarrow \infty$
 - In practice, it leads to better models and faster convergence
 - The learned proposals can be very useful

Outline

- Short (Re)Introduction to VI and VAE
- State Space Models and Stochastic RNNs
- Variational Inference of SSM
- (Variational) Sequential Monte Carlo
- **Conclusions**

Backup Slides

