

Jin Gu

Department of Automation, Tsinghua University

Email: jgu@tsinghua.edu.cn

Phone: (010) 62794294-866

Chapter 9

Inference as Optimization:

Structured Variational Inference

2021 Fall

Jin Gu (古槿)

Outlines

- Variational inferences
 - Basic idea
 - Evidence lower bound
- Mean field variational inferences
 - Basic methods
 - Example for Ising model
- Structured variational inferences
- Deep variational autoencoder

Chapter 9 Variational Inference & MAP Inference

Textbook2

Chapter 21 Variational Inference

Textbook1

Chapter 11.3.6 Variational Analysis of Belief Propagation

Chapter 11.5 Structured Variational Approximations

Chapter 13 MAP Inference

Other references

[1] Doersch C. **Tutorial on variational autoencoders.**

<https://arxiv.org/abs/1606.05908>

[2] Kingma DP, Welling M. **Auto-encoding variational Bayes.**

<https://arxiv.org/abs/1312.6114>

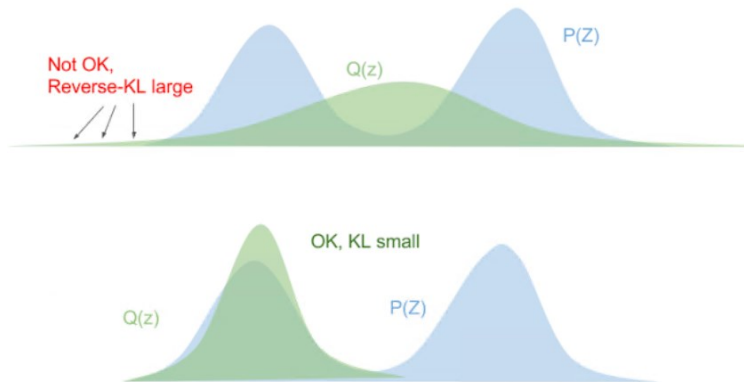
[3] Blei MN, Ng AY and Jordan MI. **Latent Dirichlet Allocation.** *Journal of Machine Learning Research* 2003, 3:993-1022.

What is Variational Inference?

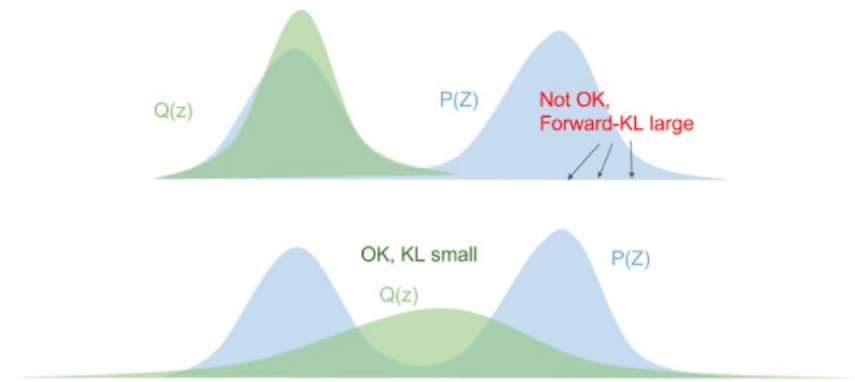
- To infer $P(X)$, find another distribution $Q(X)$ which can approximate it
 - Target distribution: $P(X)$
 - Proposal distribution: $Q(X)$
 - Q is restricted to a family of distribution with simple form
 - KL divergence: $D_{KL}(Q||P)$ or $D_{KL}(P||Q)$
 - $D_{KL}(Q||P) = E_{X \sim Q} \left(\ln \left(\frac{Q}{P} \right) \right)$
 - Aim: $\min_Q D_{KL}(Q||P)$
- Structured variational inference
 - Define $Q(X)$ based on the graph with simple structure
 - E.g. A graph with no edge (*mean field algorithm*)

Comments On the Two Divergences

- Assume $Q(X)$ is restricted as Gaussians
 - Reverse KL (I-projection): $D_{KL}(Q||P)$
 - Forward KL (M-projection): $D_{KL}(P||Q)$



Reverse KL: Q should model one of the peaks in $P > 0$



Forward KL: Q should cover $P > 0$

Q : Why *reverse* KL is commonly used?

Calculations in Variational Inferences

- I-projection
 - General conditional target distribution $P(X|Z=z)$

$$D_{KL}(Q||P) = E_{X \sim Q} \left(\ln \frac{Q(X)}{P(X|z)} \right) \quad P(X|z) = P(X, z)/P(z)$$

$$D_{KL}(Q||P) = E_{X \sim Q} \left(\ln \frac{Q(X)}{P(X, z)} \right) + \ln P(z)$$

$$\ln P(z) - D_{KL}(Q||P) = -E_{X \sim Q} \left(\ln \frac{Q(X)}{P(X, z)} \right) = E_{X \sim Q} (\ln P(X, z)) + H(Q(X))$$

- Two key questions:
 - **Q1**: How to choose the family of the proposal distribution?
 - **Q2**: How to maximize the above energy functional?

Exact Inference as Optimization

- Energy functional

$$D_{KL}(Q||P_{\Phi}) = \ln Z - \left(H(Q(X)) + \sum_{\phi \in \Phi} E_Q(\phi) \right)$$

– The second term is energy functional $F[\tilde{P}_{\Phi}, Q]$

- For a clique tree with clique C_i , we have a set of beliefs \mathbf{Q} ($\beta_i, \mu_{i,j}$ - not calibrated). Its energy functional:

$$\tilde{P}_{\Phi} = \prod_i \psi_i \quad (\psi_i \text{ is the initial factor for each clique})$$

$$\tilde{F}[\tilde{P}_{\Phi}, \mathbf{Q}] = \sum_i H_{\beta_i}(C_i) - \sum_{(i-j) \in \mathcal{E}} H_{\mu_{i,j}}(S_{i,j}) + \sum_i E_{\beta_i}[\ln \psi_i]$$

Exact Inference as Optimization

- If \mathbf{Q} is calibrated as \mathbf{Q}^* , we can conclude that
$$\tilde{F}[\tilde{P}_\Phi, \mathbf{Q}^*] = \max_{\mathbf{Q}} \tilde{F}[\tilde{P}_\Phi, \mathbf{Q}] = \ln Z$$
- Because the **distribution is invariant for any calibrated beliefs**, the relative entropy is minimized as 0
- Transform SP/BU algorithm as optimization

Ctree-Optimize-KL:

Find $\mathbf{Q} = \{\beta_i : i \in \mathcal{V}_T\} \cup \{\mu_{i,j} : (i,j) \in \mathcal{E}_T\}$
maximizing $-D(\mathbf{Q} \| P_\Phi)$
subject to

$$\mu_{i,j}[s_{i,j}] = \sum_{\mathbf{c}_{i-j}} \beta_i(\mathbf{c}_i) \quad \forall (i,j) \in \mathcal{E}_T, \forall s_{i,j} \in \text{Val}(\mathbf{S}_{i,j})$$

$$\sum_{\mathbf{c}_i} \beta_i(\mathbf{c}_i) = 1 \quad \forall i \in \mathcal{V}_T.$$

Energy Functional

- Key is to *find a Q* to *minimize the KL distance* that is equal to *maximize the energy functional*
 $L(Q) = E_{X \sim Q}(\ln P(X, z)) + H(Q)$
 - Or minimize $J(Q) = -L(Q)$
- Normally, Q is restricted to a distribution family
- The first term: $E_{X \sim Q}(\ln P(X, z))$
 - Q should mimic P
- The second term: $H(Q)$
 - Q should diffuse (reduce overfitting)

$$D_{KL}(Q || P_{\Phi}) = \ln Z - (H(Q(X)) + \sum_{\phi \in \Phi} E_Q(\phi))$$

Evidence Lower Bound (ELBO)

- Because KL divergence is non-negative, the energy functional is bound by $\ln P(z)$

$$\begin{aligned} L(Q) &= E_{X \sim Q}(\ln P(X, z)) + H(Q) \\ &= \ln P(z) - D_{KL}(Q || P) \leq \ln P(z) \end{aligned}$$

- So, we also call $L(Q)$ as ELBO
- Find a distribution Q to maximize the ELBO

What to **LEARN** in this lecture?

- Simplest case: treat all variables independently (**no edge**) in the proposal distribution Q

Mean field inference (平均場)

- Extended to general structures of proposal distribution Q (structured VI)
- Use deep NNs to do variational inferences

Mean Field Variational Inference

- Now, we need to calculate (infer)
 - $p(x|z)$ [*Note: p or P are treated as the same here*]
 - But, it is usually hard to calculate...
- We want to use a simplified distribution q to mimic the distribution p
- The **simplest** distribution family is in *fully factorized* format

$$q(x; \theta) = \prod_i q_i(x_i)$$

Mean Field Variational Inference

- Recall the energy functional

$$L(q) = E_{x \sim q}(\ln p(x, z)) + H(q)$$

- For fully factorized $q = \prod_i q_i(x_i)$

$$E_{x \sim q}(\ln p(x, z)) = \sum_x \left[\prod_i q_i(x_i) \right] \ln p(x, z)$$

$$H(q) = \sum_x \left[\prod_i q_i(x_i) \right] \left[- \sum_i \ln q_i(x_i) \right]$$

- So finally, we get

$$L(q) = \sum_x \left[\prod_i q_i(x_i) \right] \left[\ln p(x, z) - \sum_k \ln q_k(x_k) \right]$$

Mean Field Variational Inference

- We can derive the energy function for each q_j ,
assuming the other terms q_{-j} are fixed

$$\begin{aligned} L(q_j) &= \sum_x \prod_i q_i(x_i) \left[\log \tilde{p}(x) - \sum_k \log q_k(x_k) \right] \\ &= \sum_{x_j} \sum_{x_{-j}} q_j(x_j) \prod_{i \neq j} q_i(x_i) \left[\log \tilde{p}(x) - \sum_k \log q_k(x_k) \right] \\ &= \sum_{x_j} q_j(x_j) \sum_{x_{-j}} \prod_{i \neq j} q_i(x_i) \log \tilde{p}(x) \\ &\quad - \sum_{x_j} q_j(x_j) \sum_{x_{-j}} \prod_{i \neq j} q_i(x_i) \left[\sum_{k \neq j} \log q_k(x_k) + \log q_j(x_j) \right] \\ &= \sum_{x_j} q_j(x_j) \log f_j(x_j) - \sum_{x_j} q_j(x_j) \log q_j(x_j) + \text{const} \end{aligned}$$

$$\log f_j(x_j) \triangleq \sum_{x_{-j}} \prod_{i \neq j} q_i(x_i) \log \tilde{p}(x) = \mathbb{E}_{-q_j}[\log \tilde{p}(x)]$$

Mean Field Variational Inference

- Maximize the energy $L(q_j)$

$$L(q_j) = \sum_{x_j} q_j(x_j) \log \left[\frac{\exp \{E_{-q_j}[\ln \tilde{p}(x)]\}}{q_j(x_j)} \right] + c$$

- Recall the KL divergence

$$D_{KL}(q || \tilde{p}) = - \sum q \ln \frac{\tilde{p}}{q}$$

$q^* \propto \tilde{p}$ will minimize the divergence

KEY: $L(q_j)$ follows a form of KL divergence

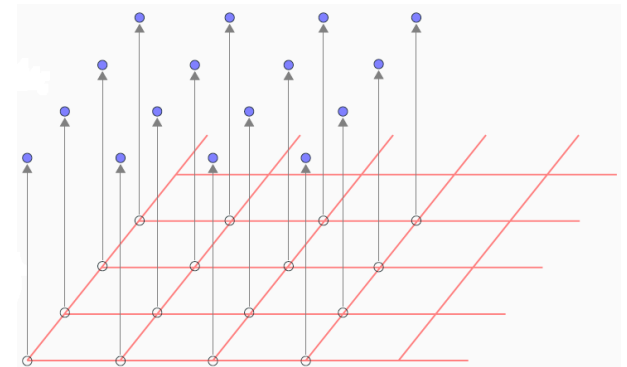
- **We can easily see the solution:**

$$q_j(x_j) \propto \exp \{E_{-q_j}[\ln \tilde{p}(x)]\}$$

An Example for Ising Model

- The hidden variables x_i are binary and z_i are always observed. (Please note *the edges are directed* from x to z)
- Aim: calculate $E_{-q_j}[\ln \tilde{p}(x)]$
- The target distribution P
 $p(x|z) \propto p(z|x)p(x)$
 - $p(x) = \frac{1}{Z} e^{\sum_{i,j} w_{ij} x_i x_j} = \frac{1}{Z} e^{-E_0(x)}$
 - $p(z|x) = e^{\sum_i \ln p(z_i|x_i)} = e^{\sum_i L_i(x_i)}$
- We can rewrite as

$$\ln \tilde{p}(x) = \sum_{i,j} w_{ij} x_i x_j + \sum_i L_i(x_i) + \text{const}$$



$$q_j(x_j) \propto \exp \left\{ E_{-q_j}[\ln \tilde{p}(x)] \right\}$$

Mean Field for Ising Model

- Set the family of the variational distribution $q = \prod_i q(x_i; \theta_i)$.
 - Each $q(x_i; \theta_i)$ is a *binary* distribution
 - $q_i(x_i = 1) = \theta_i$ and $q_i(x_i = -1) = 1 - \theta_i$
- NOTE: the parameter θ_i is not determined
- KEY: variational inference here is to **find the optimal parameter θ_i^*** to maximize the energy functional (then we get q^*)
- AIM: inference in q^* rather than in p

$$q_j(x_j) \propto \exp \left\{ E_{-q_j} [\ln \tilde{p}(x)] \right\}$$

Mean Field for Ising Model

- Recall the previous result

$$\ln q_j^*(x_j) = E_{-q_j}[\ln \tilde{p}(x)] + \text{const}$$

Please remember the above optimal solution for mean-field inference

- For Ising model

$$\ln \tilde{p}(x) = \sum_{i,j} w_{ij} x_i x_j + \sum_i L_i(x_i) + \text{const}$$

$$\begin{aligned} E_{-q_j}[\ln \tilde{p}(x)] &= E_{-q_j} \left(\sum_{i,j} w_{ij} x_i x_j + \sum_i L_i(x_i) \right) + c \\ &= x_j \sum_{i \in \text{Neighbor}(j)} [w_{ij}(2\theta_i - 1)] + L_j(x_j) + c \end{aligned}$$

- Set $m_j = \sum_{i \in \text{Neighbor}(j)} [w_{ij}(2\theta_i - 1)]$

$$q_j(x_j) \propto \exp \{ E_{-q_j}[\ln \tilde{p}(x)] \}$$

Mean Field for Ising Model

- So we can get the optimal solution in this step

$$\ln q^*(x_j; \theta_j) = x_j m_j + L_j(x_j) + \text{const}$$

- Obviously, the optimal parameter

$$\theta_j^* = \frac{\exp[m_j + L_j(x_j=1)]}{\exp[m_j + L_j(x_j=1)] + \exp[-m_j + L_j(x_j=-1)]}$$

- Here $L_j(x_j) = \ln p(z_j | x_j)$

$$m_j = \sum_{i \in \text{Neighbor}(j)} [w_{ij} (2\theta_i - 1)]$$

- We need to update the parameters by iteration, because we need to specify θ_{-j}

$$q_j(x_j) \propto \exp \{E_{-q_j}[\ln \tilde{p}(x)]\}$$

Mean Field Inference for De-noising

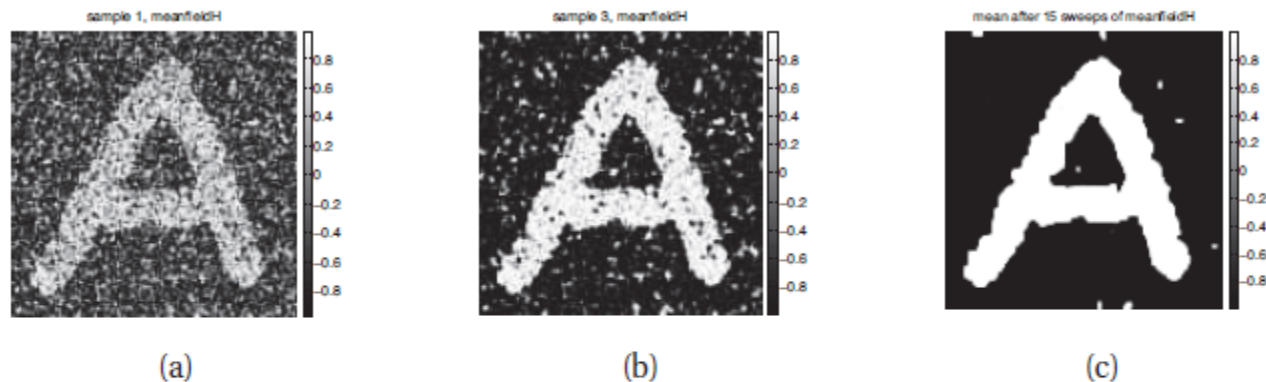
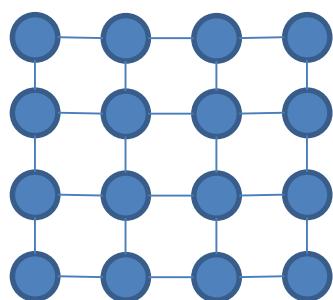


Figure 21.3 Example of image denoising using mean field (with parallel updates and a damping factor of 0.5). We use an Ising prior with $W_{ij} = 1$ and a Gaussian noise model with $\sigma = 2$. We show the results after 1, 3 and 15 iterations across the image. Compare to Figure 24.1. Figure generated by `isingImageDenoiseDemo`.

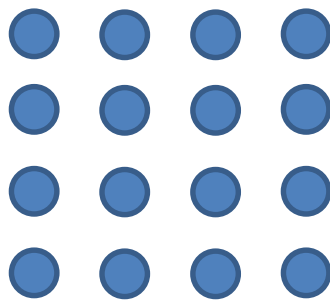
Infer the probability of the hidden variables according to q rather than p . Sometimes, the variational probability is “better” than the target probability: *de-noising* or *smoothing* effects.

General Structured Variational Inference

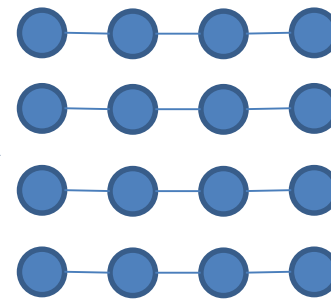
- Graph for mean field inference
 - No edge is in the variational distribution family
- General structured variational inference
 - We can use any simpler structure to do inference



$p(x|z)$



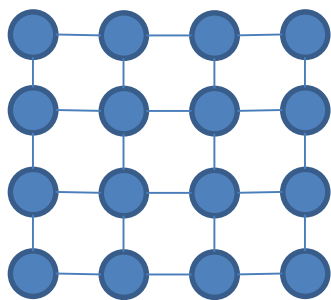
Mean field



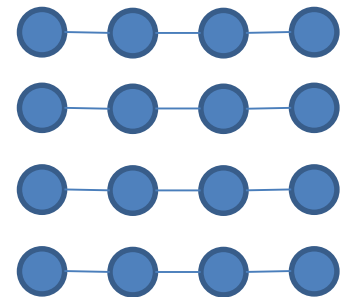
Linear chain

General Structured Variational Inference

- If we use linear chain for Ising model
- Recall: *find a Q* to *minimize the KL distance* that is equal to *maximize the energy functional*
$$L(q) = E_{x \sim q}(\ln p) + H(q)$$
- We can set $q \propto \prod_l q_l(x_l)$



$p(x|z)$



Linear chain

$$L(q) = E_{x \sim q}(\ln p) + H(q)$$

- $q_i(x_i) = \frac{1}{Z_i} e^{\sum_j \beta_{j,j+1} x_j x_{j+1}}$

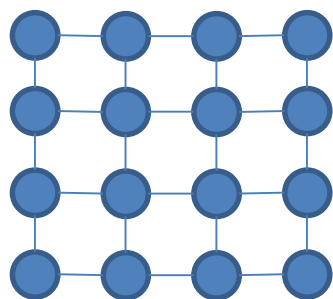
Each chain has a separate variational Gibbs distribution.

- $p = \frac{1}{Z} e^{\sum_{i,j} w_{i,j} x_i x_j}$

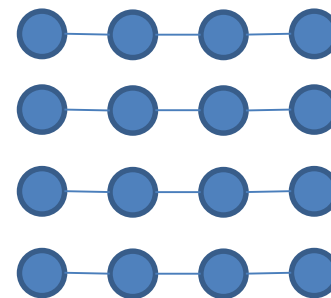
- $E_{x \sim q}(\ln p) = \sum_x (-\ln Z + \sum_{i,j} w_{ij} x_i x_j) \prod_i q_i$

- $H(q) = \sum_x (-\sum_i \ln q_i) \prod_i q_i$

Maximize the energy functional **chain by chain!**



$p(x|z)$



Linear chain

Variational Inference for LDA

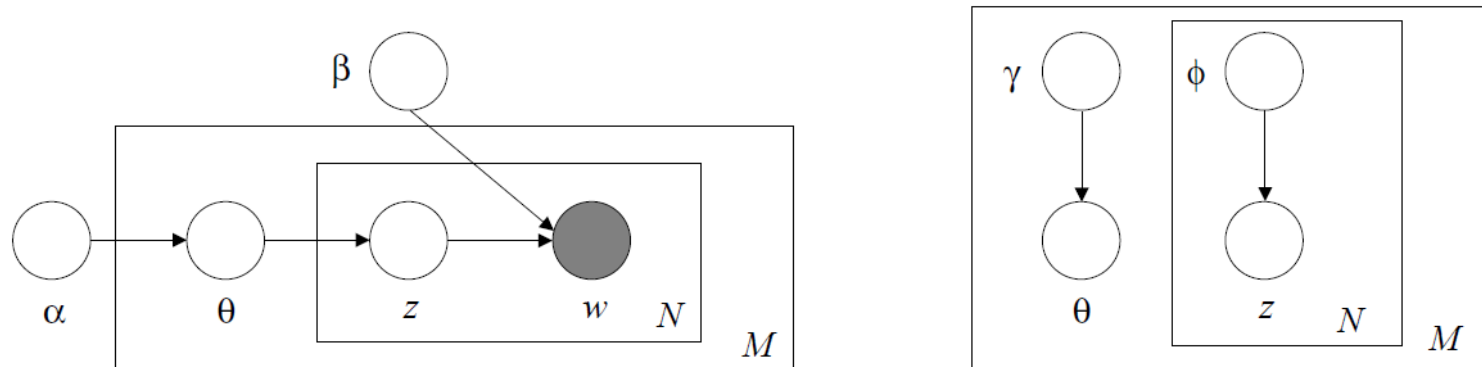


Figure 5: (Left) Graphical model representation of LDA. (Right) Graphical model representation of the variational distribution used to approximate the posterior in LDA.

$$p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta) = p(\theta | \alpha) \prod_{n=1}^N p(z_n | \theta) p(w_n | z_n, \beta)$$

$$q(\theta, \mathbf{z} | \gamma, \phi) = q(\theta | \gamma) \prod_{n=1}^N q(z_n | \phi_n)$$

NOTE: **variables** should be the same and **parameters** are different

Variables: θ and \mathbf{z}
Parameters: γ and ϕ_n

General Structured Variational Inference

- For a general variational structures, we can finally derive:

$$\psi_j(\mathbf{c}_j) \propto \exp \left\{ E_Q \left[\ln \tilde{P}_\Phi \mid \mathbf{c}_j \right] - \sum_{k \neq j} E_Q [\ln \psi_k \mid \mathbf{c}_j] \right\}$$

- Comments:
 - The energy functional can be separated into maximal cliques on q
 - The update also depends on the expectation

Stochastic Gradient for Parameter Update

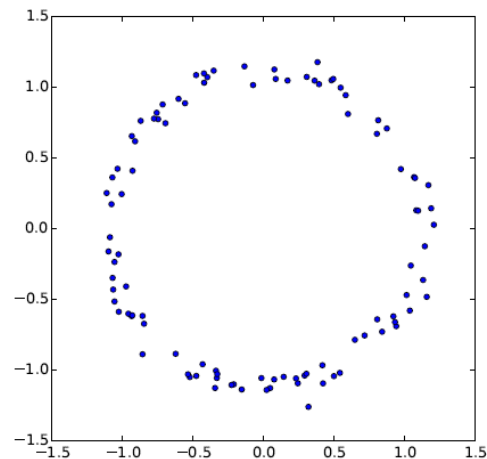
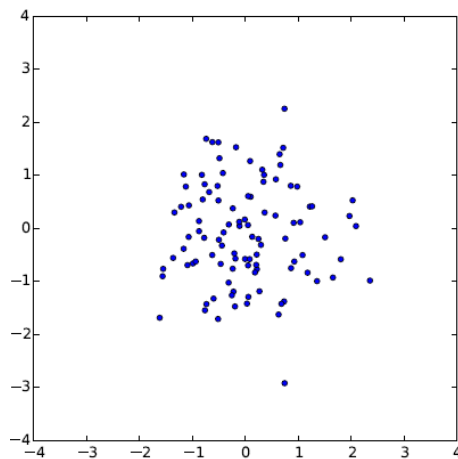
- In many cases, we cannot derive the close-form solution as the mean field inference
- Many heuristic optimization methods can be used to optimize the target function
 - Remember the belief propagation on cluster graph
 - Stochastic gradient is commonly used
- Stochastic gradient (for maximizing)
$$L(q) = E_{x \sim q}(\ln p) + H(q)$$
 - Calculate the partial derivate $\frac{\partial L(q)}{\partial \theta_i}$
 - Then, update $\theta_i^{t+1} = \theta_i^t + \delta \frac{\partial L(q)}{\partial \theta_i^t}$

Limitation of Variational Inference

- The variational inference distribution family is hard to choose
 - Simple family cannot capture the complex features in target distribution
 - Complex family is hard for computation and easy for over-fitting

NNs for Variational Inference

- ▶ Assume the observed data $X \in R^N$ are generated from a hidden subspace $Z \in R^L$ as standard Gaussian
- ▶ In most cases, the posterior $P(Z|X)$ is hard to infer



$$g(z) = z/10 + z/||z||$$

NNs for Variational Inference

- ▶ Another distribution $Q(Z|X)$ can be used to approximate $P(Z|X)$ with respect to the minimization of KL-divergence

$$\begin{aligned} & \log P(\mathbf{X}) - D[Q(\mathbf{z}|\mathbf{X})||P(\mathbf{z}|\mathbf{X})] \\ &= E_{\mathbf{z} \sim Q}[\log P(\mathbf{X}|\mathbf{z})] - D[Q(\mathbf{z}|\mathbf{X})||P(\mathbf{z})] \end{aligned}$$

First term: maximize the data generation probability

Second term: minimize the KL-divergence between standard Gaussian and the variational distribution

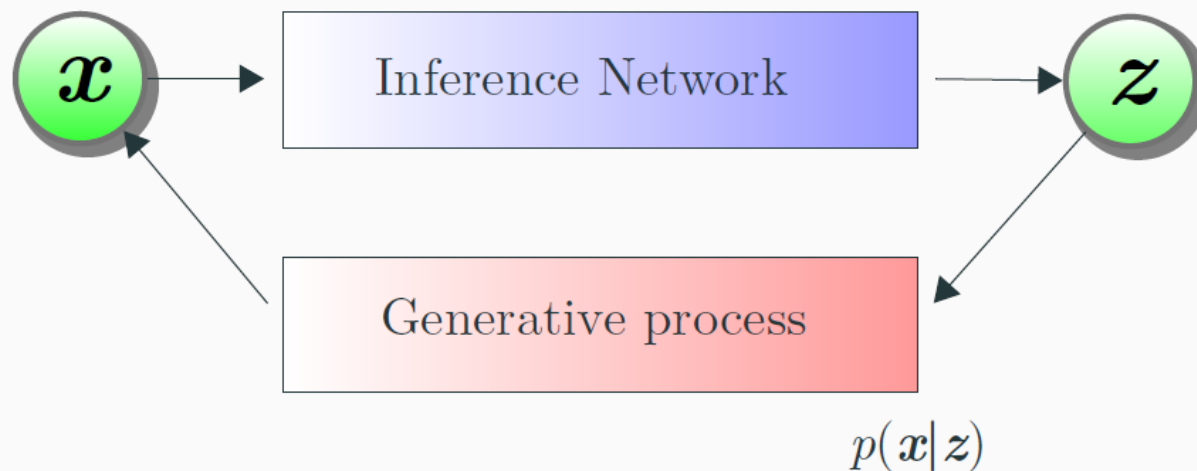
- ▶ The two terms can be treated as the “decoding” (from \mathbf{z} to \mathbf{X}) and “encoding” (from \mathbf{X} to \mathbf{z}) process. Two deep neural networks can be used to approximate these terms

NNs for Variational Inference

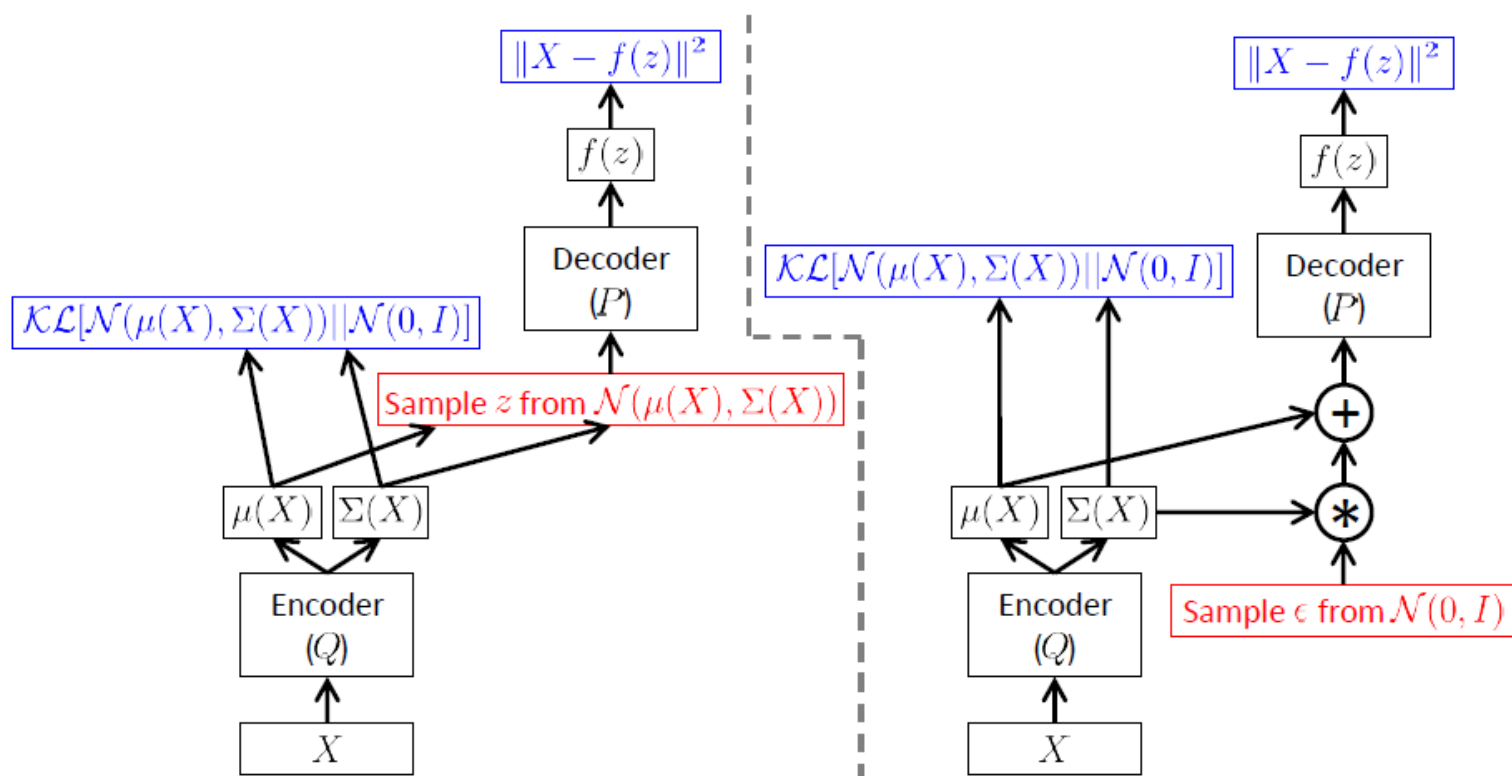
VAE idea: Use neural networks to approximate both **variational** and **generative** parameters.

$$\begin{aligned}\mathcal{L}(\theta; \mathbf{x}) &= \mathbb{E}_{z \sim q(z|\mathbf{x}; \theta)} \log p(\mathbf{x}, z) + \mathbb{H}(q(z|\mathbf{x}; \theta)) \\ &= \mathbb{E}_{z \sim q(z|\mathbf{x}; \theta)} \log p(\mathbf{x}|z) - \mathbb{KL}(q(z|\mathbf{x}; \theta) \parallel p(z))\end{aligned}$$

$$q(z|\mathbf{x}) = \mathcal{N}(f_{\mu}(\mathbf{x}), f_{\sigma}(\mathbf{x}))$$

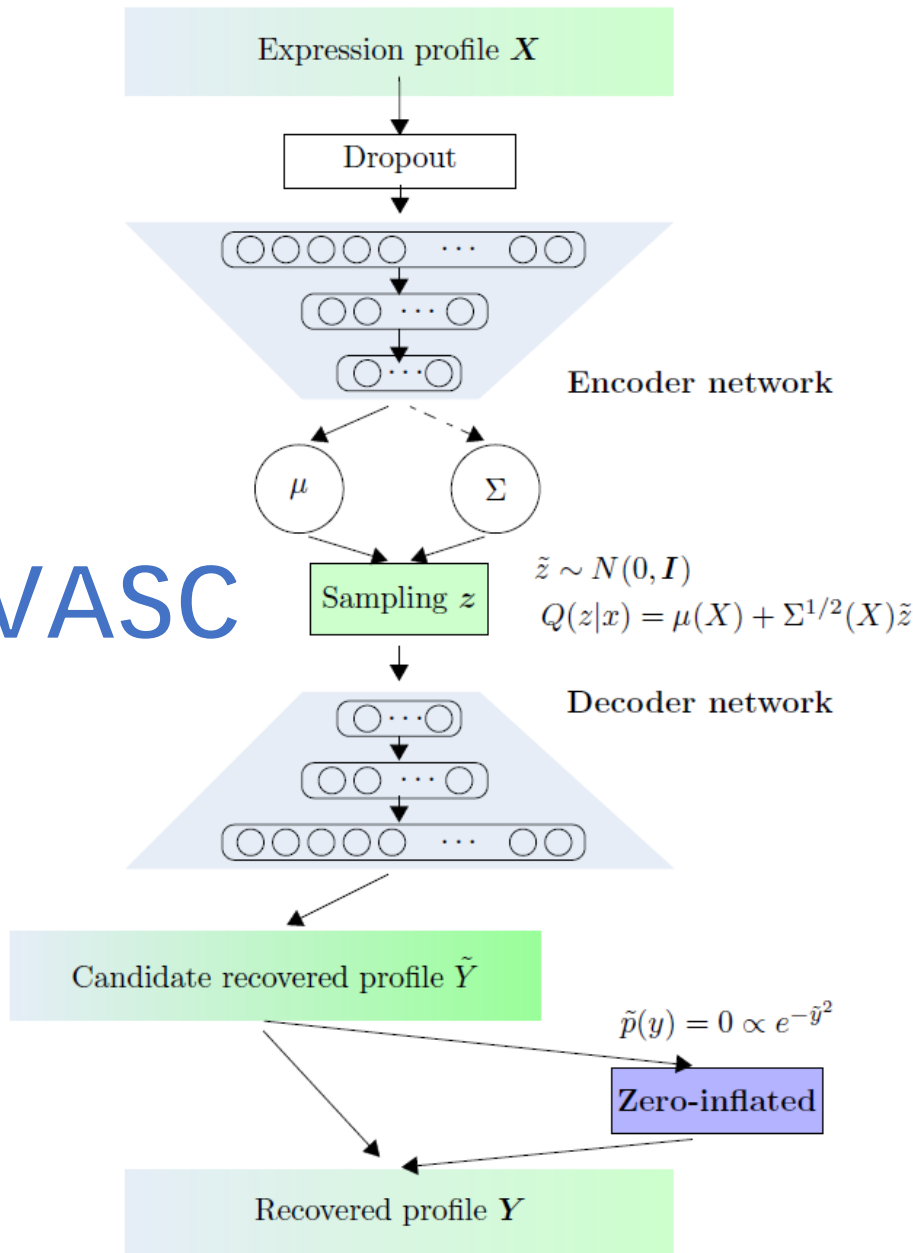


NNs for Variational Inference



Re-parameterization trick for learning

VASC



Input: scRNA-seq FPKM matrix

Additional random dropout!!

Three-layer encoder NNs:
L1 norm in first layer
ReLU activation

Target: Low-dimension
representation by sampling

Three-layer decoder NNs:
ReLU activation

Model dropouts by Gumbel
distribution instead of “hard” ZI

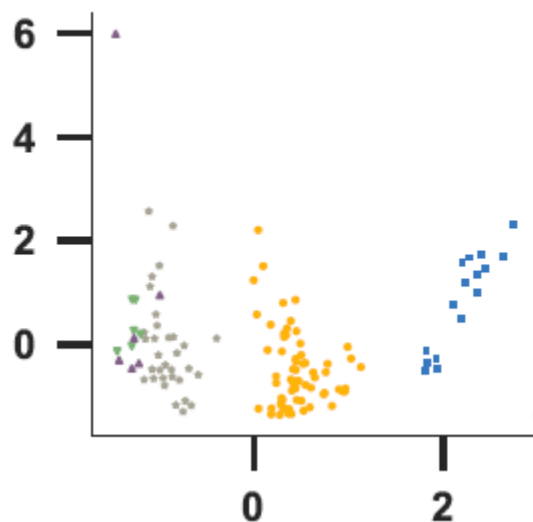
Output: recovered FPKM matrix

$$\text{Loss}(X, Y) = \text{binary_entropy}(X, Y) + \text{KL}(Q(z|X) || P(z))$$

Advantages of VAE

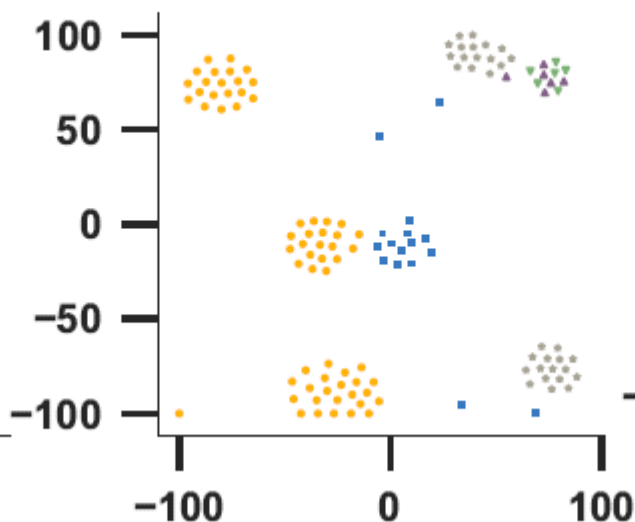
- ▶ **Unequal** structures in encoding & decoding networks
 - ▶ VAE is doing inference rather than autoencoding
 - ▶ Encoding/decoding NNs are used to mimic the two terms in the variational energy functional
 - ▶ So, we can freely design the topology of two NNs
- ▶ Dimension reduction
 - ▶ Use the inferred mean value calculated by the encoding network $\mu(X)$
- ▶ A generative model
 - ▶ Fill the missing values and generate new data

ZIFA



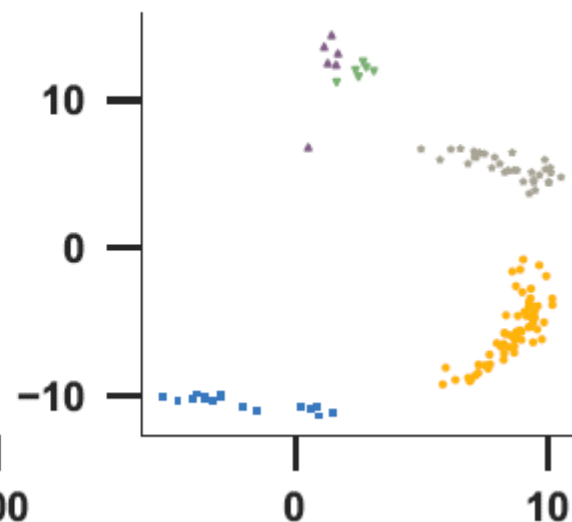
■ 2cell ● 4cell ★ 8cell ▼ 16cell ▲ blast

SIMLR

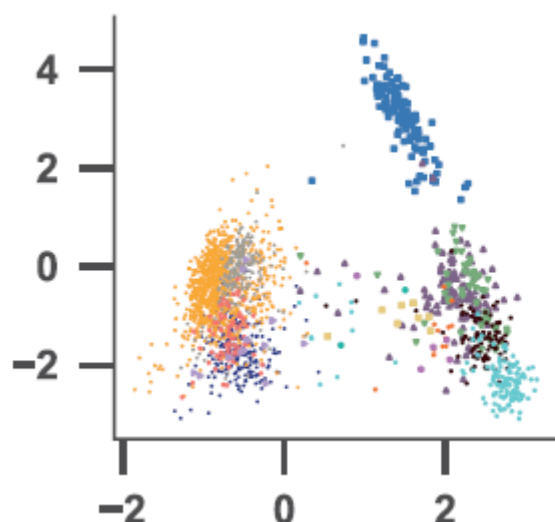


■ 2cell ● 4cell ★ 8cell ▼ 16cell ▲ blast

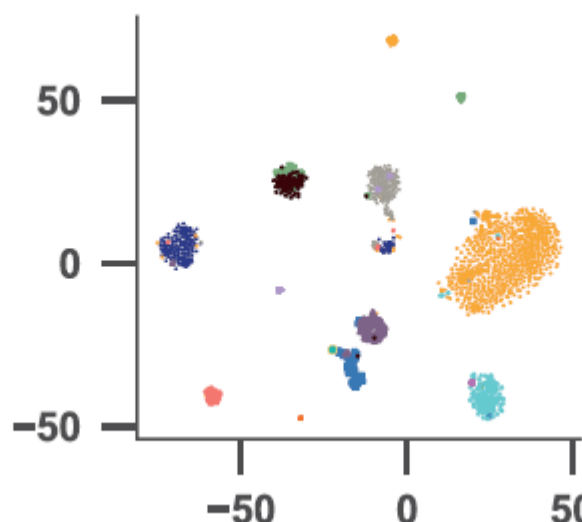
VASC



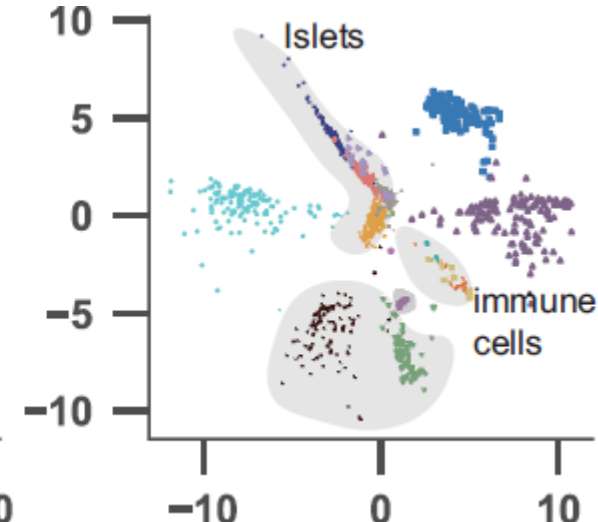
■ 2cell ● 4cell ★ 8cell ▼ 16cell ▲ blast



◄ alpha ● beta ★ delta ◄ gamma ◄ epsilon ◄ acinar ◄ ductal ▼ activated_stellate
◄ quiescent_stellate ◄ endothelial ◄ macrophage ◄ mast ◄ t_cell ◄ schwann

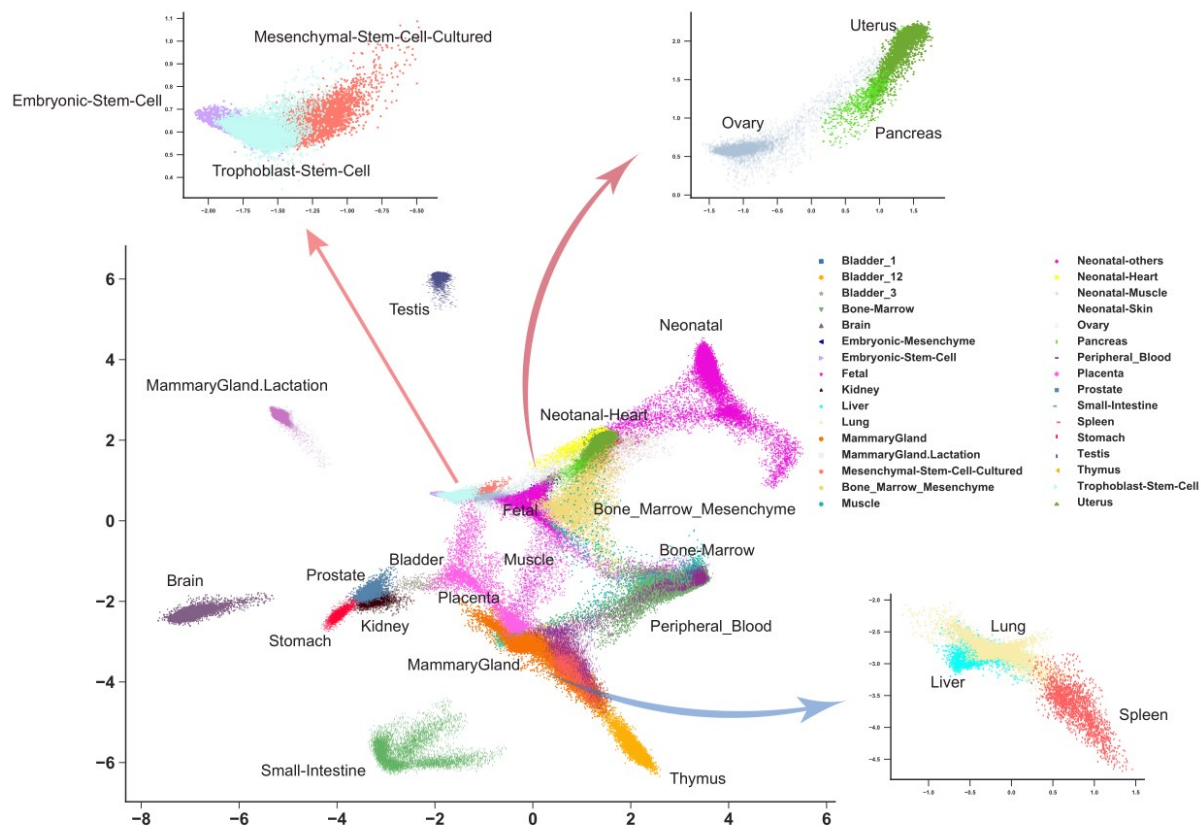


◄ alpha ● beta ★ delta ◄ gamma ◄ epsilon ◄ acinar ◄ ductal ▼ activated_stellate
◄ quiescent_stellate ◄ endothelial ◄ macrophage ◄ mast ◄ t_cell ◄ schwann

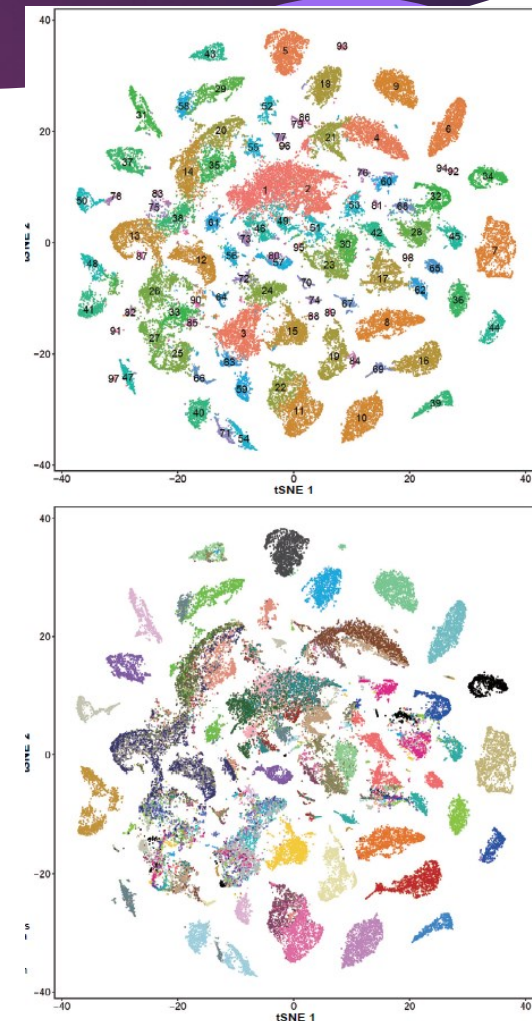


◄ alpha ● beta ★ delta ◄ gamma ◄ epsilon ◄ acinar ◄ ductal ▼ activated_stellate
◄ quiescent_stellate ◄ endothelial ◄ macrophage ◄ mast ◄ t_cell ◄ schwann

Test on 10X & Microwell datasets



t-SNE (right) only consider the local similarity
VASC preserves the topology of cell expression patterns



Summary

- Variational inference is to find a simpler distribution q to mimic the target distribution p . Then, you can do easy inference in q .
- Variational inference is an *optimization* process to find a variational distribution q from a distribution family by *maximizing an energy functional*.
- Then variational distribution is usually simpler than the target distribution. So, it can be used for de-noising or dimension reduction.

Comments: From Inference to Learning

- Variational inference can be regarded as a kind of learning: find the optimal parameters by maximizing an energy functional

$$- \max_q L(q) \quad L(q) = E_{x \sim q}(\ln p(x|z)) + H(q)$$

- Learning as MLE: find the optimal parameters by maximizing likelihood function given data

$$- \max_{\theta} L(D; \theta)$$

The End of Chapter 9

Structured variational inference is
commonly used in PGMs

MAP Inference

- Probability inference
 - Compute the distribution $P(Y)$ or $P(Y|\mathbf{e})$
- MAP (maximum a posterior) inference
 - Compute the optimal assignment or configuration
 - $\text{MAP}(Y|\mathbf{e}) = \arg \max_{\mathbf{y} \in \text{Val}(Y)} P(\mathbf{y}, \mathbf{e})$
- The max-marginal of a function f
 - $\text{MaxMarg}_f(\mathbf{y}) = \max_{\mathbf{y} \in \text{Val}(Y)} f(\mathbf{y})$
 - $\text{MaxMarg}_F(\mathbf{y}) = \max_{\mathbf{y} \in \text{Val}(Y)} \sum_{\mathbf{X}=\mathbf{x}} F(\mathbf{y}, \mathbf{X})$

Compute $\underset{Y}{\operatorname{argmax}} P(X, Y|\theta)$

Viterbi algorithm Revisited

- Known transition matrix and emission matrix
- Infer the *maximum* probability STATE series
- The probability at time 0 with observation x_0

$$\text{For } Y_1 = i: \quad \delta_{1,i} = \pi_i e_{i,x_1}$$

- The probability at time 1 with observation x_1

$$\delta_{2,i} = e_{i,x_2} \max_{y_1=1,\dots} (\pi_i e_{i,x_2} \times t_{y_1,i}) = e_{i,x_2} \max_{\dots} (\delta_{1,y_1} \times t_{y_1,i})$$

$$\phi_2(i) = \arg \max_{y_1=1,\dots} (\delta_{1,y_1} \times t_{y_1,i})$$

Compute $\underset{Y}{\operatorname{argmax}} P(X, Y|\theta)$

Viterbi algorithm Revisited

- The probability at time t with observation x_t

$$\delta_{t,i} = e_{i,x_t} \max_{y_{t-1}=1,\dots} (\delta_{t-1,y_{t-1}} \times t_{y_{t-1},i})$$

$$\phi_t(i) = \arg \max_{y_{t-1}=1,\dots} (\delta_{t-1,y_{t-1}} \times t_{y_{t-1},i})$$

- For the last observation

$$\delta_{T,i} = e_{i,x_T} \max_{y_{T-1}=1,\dots} (\delta_{T-1,y_{T-1}} \times t_{y_{T-1},i})$$

$$\phi_T(i) = \arg \max_{y_{T-1}=1,\dots} (\delta_{T-1,y_{T-1}} \times t_{y_{T-1},i})$$

$$y_T^* = \arg \max_{y_T=1,\dots} (\delta_{T,y_T})$$

Compute $\underset{Y}{\operatorname{argmax}} P(X, Y|\theta)$

Viterbi algorithm Revisited

- After y_t is inferred, trace back to get other y_i

$$y_T^* = \underset{y_T=1,\dots}{\operatorname{argmax}} \left(\delta_{T,y_T} \right)$$

$$\phi_t(i) = \underset{y_{t-1}=1,\dots}{\operatorname{argmax}} \left(\delta_{t-1,y_{t-1}} \times t_{y_{t-1},i} \right)$$

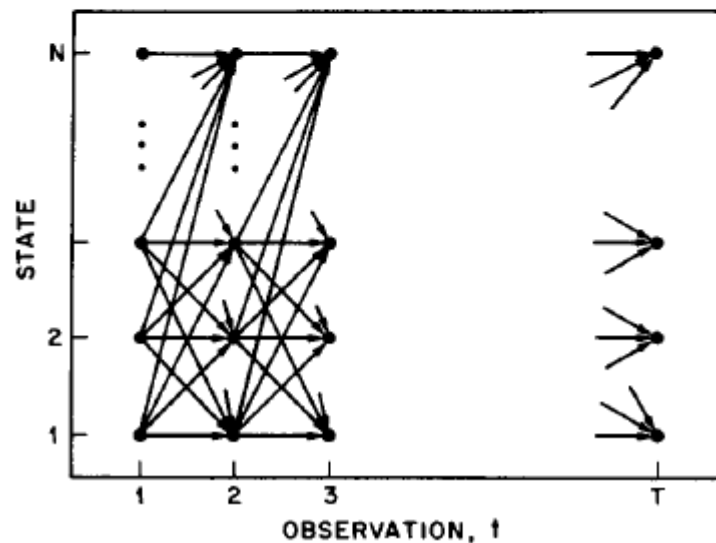
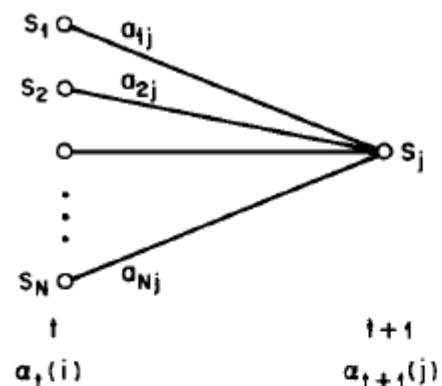
- We can infer other y_i

$$y_{T-1}^* = \phi_T \left(i = y_T^* \right) = \underset{y_{T-1}=1,\dots}{\operatorname{argmax}} \left(\delta_{T-1,y_{T-1}} \times t_{y_{T-1},i} \right)$$

$$y_{t-1}^* = \phi_t \left(i = y_t^* \right) = \underset{y_{t-1}=1,\dots}{\operatorname{argmax}} \left(\delta_{t-1,y_{t-1}} \times t_{y_{t-1},i} \right)$$

Viterbi Algorithm Revisited

- Do variable elimination and record the paths which maximizing the joint distribution for the first k time slices after k -th run of elimination
- After all query variables are eliminated, **trace back** the path to get the configuration which maximizes the joint distribution
- *Could we generalize Viterbi algorithm to Bayesian Networks and Markov Networks?*



Move to *Learning*

- Representation $P \Leftrightarrow \{P, G\}$

- Inference $P(Y | E = e, \theta)$

- Learning
$$\max_{\theta} P(\mathbf{x}[1], \mathbf{x}[2], \dots | \theta)$$
$$P(\theta | \mathbf{x}[1], \mathbf{x}[2], \dots)$$