

Variational approximation(Continued)

Let $L(Q, \lambda) = Q(q, g) + \lambda(\sum_z q_z - 1)$, then

$$\begin{aligned}\frac{\partial L(Q, \lambda)}{\partial q_j} &= \log g(y, j) - \log q_j - 1 + \lambda \\ &\implies q_j^* \equiv \alpha g(y, j)\end{aligned}$$

$$\sum_z q_z = 1 : q_j^* = \frac{g(y, j)}{\sum_z g(y, z)}$$

Substituting q^* in the original maximization equation, we get
 $\log \sum_{z=1}^k g(y, z)$.

Hence LHS=RHS.

For all other q ,

$$\log \sum_{z=1}^k g(y, z) \geq \sum_{z=1}^k q_z \log g(y, z) - \sum_z q_z \log q_z$$

Applying variational approximation

$$\begin{aligned}\theta^{ML} &\equiv \arg \max_{\theta} \sum_{i=1}^N \log \sum_{\mathbf{z}: z_1, \dots, m} P(\mathbf{y}^i, \mathbf{z} | \theta, \mathbf{x}^i) \\ &\equiv \arg \max_{\theta} \sum_{i=1}^N \max_{q_{i,\mathbf{z}}: \sum_{\mathbf{z}} q_{i,\mathbf{z}} = 1} \sum_{\mathbf{z}} q_{i,\mathbf{z}} \log P(\mathbf{y}^i, \mathbf{z} | \theta, \mathbf{x}^i) - \sum_{\mathbf{z}} q_{i,\mathbf{z}} \log q_{i,\mathbf{z}}\end{aligned}$$

Solve by alternating optimization.

EM algorithm

E-Step: Solve for q_{iz} keeping θ fixed at θ^t

$$\begin{aligned} q_{i,z}^t &\equiv \frac{P(\mathbf{y}^i, \mathbf{z} | \theta^t, \mathbf{x}^i)}{\sum_z P(\mathbf{y}^i, \mathbf{z} | \theta^t, \mathbf{x}^i)} \\ &\equiv P(\mathbf{z} | \theta^t, \mathbf{x}^i, \mathbf{y}^i) \end{aligned}$$

Hence $q_{i,z}$ is the posterior distribution of hidden variable at time t .

M-Step: Solve for θ , keeping $q_{i,z}$ fixed to $q_{i,z}^t$. The problem becomes

$$\max_{\theta} \sum_i \sum_z q_{i,z}^t \log P(\mathbf{y}^i, \mathbf{z} | \theta, \mathbf{x}^i)$$

which is concave in θ and can be often solved in closed form.(Example:HMM)

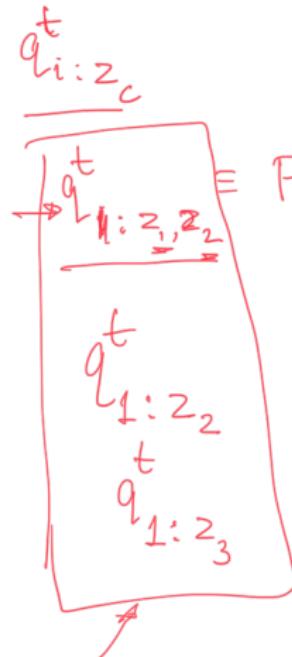
EM algorithm for graphical models

Above formulation may lead to lot of computations for calculating $q_{i,\mathbf{z}}^t$ depending on the possible values of \mathbf{z} . In graphical models, we don't need to directly compute all $q_{i,\mathbf{z}}$.

M-Step:

$$\begin{aligned} & \max_{\theta} \sum_i \sum_{\mathbf{z}} q_{i,\mathbf{z}}^t \log P(\mathbf{y}^i, \mathbf{z} | \theta, \mathbf{x}^i) \\ = & \max_{\theta} \sum_i \sum_{\mathbf{z}} q_{i,\mathbf{z}}^t (\log \exp(\sum_C F_{\theta}(\mathbf{y}_C^i, \mathbf{z}_C | \mathbf{x}^i)) - \log Z_{\theta}(\mathbf{x}^i)) \\ = & \max_{\theta} \sum_i \sum_{\mathbf{z}} q_{i,\mathbf{z}}^t (\sum_C F_{\theta}(\mathbf{y}_C^i, \mathbf{z}_C | \mathbf{x}^i)) - \sum_i \log Z_{\theta}(\mathbf{x}^i) \sum_{\mathbf{z}} q_{i,\mathbf{z}}^t \\ = & \max_{\theta} \sum_i \sum_C \sum_{\mathbf{z}_C} F_{\theta}(\mathbf{y}_C^i, \mathbf{z}_C | \mathbf{x}^i) \sum_{\mathbf{z} - \mathbf{z}_C} q_{i,\mathbf{z}}^t - \sum_i \log Z_{\theta}(\mathbf{x}^i) \\ = & \max_{\theta} \sum_i \sum_C \sum_{\mathbf{z}_C} q_{i,\mathbf{z}_C}^t F_{\theta}(\mathbf{y}_C^i, \mathbf{z}_C | \mathbf{x}^i) - \sum_i \log Z_{\theta}(\mathbf{x}^i) \end{aligned}$$

Example: CRFs



$i=1$

z_1	z_2	y_3	y_4	z_3	y_6
hostel	II,	IIT	Bombay, Parai, Mumbai		
x_1	x_2	x_3	x_4	x_5	x_6

ln ln city

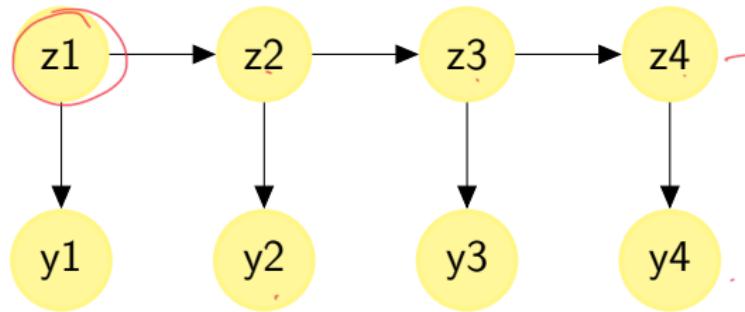
$z_1 \rightarrow z_2 \rightarrow y_3 \rightarrow y_4 \rightarrow z_3 \rightarrow y_6$

$(z_1, z_2) \quad (z_2, y_3) \quad - \quad - \quad - \quad -$

$y \in \{ \text{House\#}, \text{Inst Name}, \text{Area}, \text{City}, \text{Other} \}$

Example: CRF continued

Example: HMM training



Parameters θ of the network:

- $P(z_j|z_{j-1})$, $z_j \in 1, \dots, m$ are table potentials (m^2 entries). Since they are independent of j , they can be written as $P(z|z')$. m^2
- $P(y_j|z_j)$, $y_j \in 1, \dots, n$ are table potentials (nm entries). Since they are independent of j , they can be written as $P(y|z)$.
- $P(z)$ has m entries.

For HMM, x is not present.

Example: HMM training(continued)

There are two different kinds of $F_\theta(\mathbf{y}_C^i, \mathbf{z}_C)$

- $\log P(\mathbf{y}_p^i | z_p) , C \equiv (\mathbf{y}_p, z_p)$
- $\log P(z_p | z_{p-1}) , C \equiv (z_p, z_{p-1})$

$$\begin{aligned} & \max_\theta \sum_i \sum_C \sum_{\mathbf{z}_C} q_{i, \mathbf{z}_C}^t F_\theta(\mathbf{y}_C^i, \mathbf{z}_C) \\ &= \sum_{i=1}^N \sum_{p=1}^n q_{i, z_p}^t \log P(y_{ip} | z_p) + \\ & \quad \sum_{p=2}^n q_{i, z_p, z_{p-1}}^t \log P(z_p | z_{p-1}) \end{aligned}$$

Call $P(y|z)$ as θ_{yz} and $P(z|z')$ as $\theta_{zz'}$. In the M-step, we calculate

$$\theta_{yz}^{t+1} = \frac{\sum_{i=1}^N \sum_{p=1}^n q_{i, z_p=z}^t [[y_p^i = y]]}{\sum_{i=1}^N \sum_{p=1}^n q_{i, z_p \neq z}^t}$$

$$\theta_{zz'}^{t+1} = \frac{\sum_{i=1}^N \sum_{p=2}^n q_{i, z_p=z, z_{p-1}=z'}^t}{\sum_{i=1}^N \sum_{p=2}^n q_{i, z_{p-1}=z'}^t}$$

Example: HMM training(continued)

In E-Step, we calculate:

$$q_{i,z_p=z} = P(z_p = z | y_i, \theta^t)$$

and

$$q_{i,z_p=z, z_{p-1}=z'} = P(z_p = z, z_{p-1} = z' | y_i, \theta^t)$$

Could be calculated using junction trees.

This algorithm is called **Baum Welch algorithm**.

Overall training algorithm for BNs with hidden variables

EM Algorithm

Input: Graph G , Data D with observed subset of variables \mathbf{y} and hidden variables \mathbf{z} .

Initially ($t = 0$): Assign random variables of parameters

$$\theta^t = \Pr(y_j | \mathbf{y}_{\text{Pa}(j)})^t$$

for $i = 1, \dots, T$ **do**

E-step

for $i = 1, \dots, N$ **do**

 Use inference in G to estimate conditionals

$\rightarrow q_{i,\mathbf{z}_c}^t = \Pr(\mathbf{z}_c | \mathbf{y}^i, \theta^t)$ for all variable subsets $c \subseteq (j, \text{pa}(j))$
 involving any hidden variable.

end for

M-step

When j is observed: $\Pr(y_j | \text{pa}(y_j) = \mathbf{z}_c)^t = \frac{\sum_{i=1}^N q_{i,\mathbf{z}_c}^t [[y_j^i == y_j]]}{\sum_{i=1}^N q_{i,\mathbf{z}_c}^t}$

When j is not observed:

$$\Pr(z_j | \text{pa}(y_j) = \mathbf{z}_{c-j})^t = \frac{\sum_{i=1}^N q_{i,\mathbf{z}_c}^t}{\sum_{i=1}^N \sum_{z_j=1}^m q_{i,[\mathbf{z}_{c-j}, z_j]}^t}$$

end for

Take Aways for Learning with Hidden variables

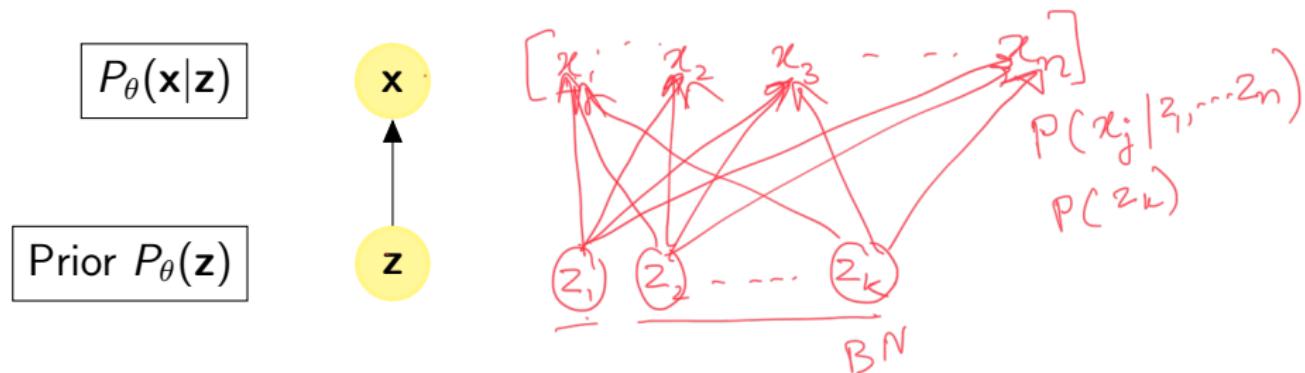
- ① When not all variables are observed in each training sample, learning of potentials expressing their interaction becomes complicated.
- ② EM algorithm is an efficient option: it alternates between estimating probability distribution over hidden variables and solving for potentials with weighted instances.
- ③ Variational approximation provides a formal justification for the EM algorithm and is more broadly applicable as we will see next.
- ④ For graphical models: the distribution over hidden variables can be factorized.

Variational Auto Encoders(VAEs)

VAEs: useful for generating high dimensional objects with latent variables. E.g. generate images with latent factor capturing its properties.

Hidden variables \underline{z} are continuous & multidimensional. $\underline{z} = z_1 \dots z_k$.
Observed variables are continuous. We denote them as $\underline{x} = x_1, \dots, x_n$ to be consistent with paper on the topic.

All \underline{z} are parents of each observed variable which are also continuous.

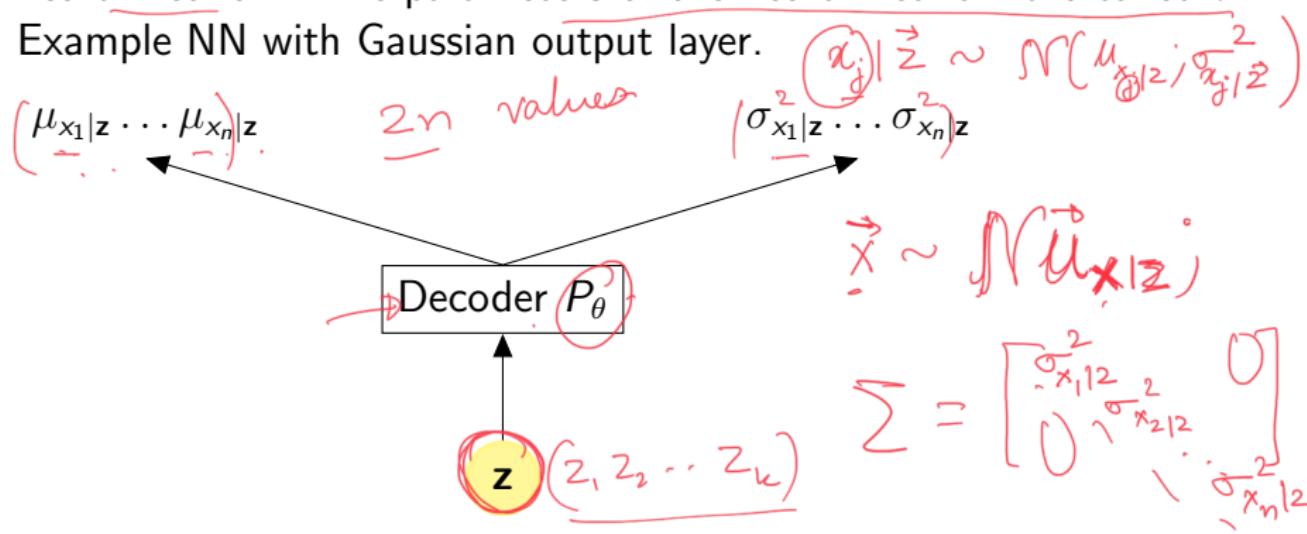


Potentials in VAE

$P(z_j)$ is usually simple like Gaussian. This provides prior distributions over the latent variables.

The $P(\mathbf{x}|\mathbf{z})$ distribution is complicated (it has to generate complex objects, example images, time series, etc). We choose to make it a neural network. The parameters of the neural network are called θ .

Example NN with Gaussian output layer.



Training VAEs

Training data $D = \{\mathbf{x}^i\}_{i=1}^N$ is observed and \mathbf{z} is hidden. Marginalize over \mathbf{z} to compute $P(\mathbf{x})$.

$$\max_{\theta} \sum_{i=1}^N \log P(\mathbf{x}^i | \theta) = \max_{\theta} \sum_{i=1}^N \log \int P_{\theta}(\mathbf{x}^i | \mathbf{z}) P_{\theta}(\mathbf{z}) d\mathbf{z}$$

$P_{\theta}(\mathbf{x}, \mathbf{z}) =$
 \mathbf{z} is continuous

$\mathbf{z} \in R^k, \mathbf{x} \in R^d$ We can apply EM algorithm to train with hidden variables.

Training VAEs(Continued)

$$\begin{aligned}\sum_{i=1}^N \log P(\mathbf{x}^i | \theta) &= \sum_{i=1}^N \log \int_z P_\theta(\mathbf{x}^i | \mathbf{z}) P_\theta(\mathbf{z}) dz \\ &= \sum_{i=1}^N \max_{q_{i,z}} \int_z (q_{i,z} \log P_\theta(\mathbf{x}^i | \mathbf{z}) P_\theta(\mathbf{z}) - q_{i,z} \log q_{i,z}) dz \\ s.t. q_{i,z} &\geq 0, \int_z q_{i,z} = 1\end{aligned}$$

But unlike in the usual applications of EM, for calculating optimal $\underline{q}_{i,z}^t$, we have to compute

$$q_{i,z}^t = \underline{P_{\theta=\theta^t}(\mathbf{z} | \mathbf{x}^i)} = \frac{\underline{P_{\theta^t}(\mathbf{x}^i | \mathbf{z}) P_{\theta^t}(\mathbf{z})}}{\int_z \underline{P_{\theta^t}(\mathbf{x}^i | \mathbf{z}) P_{\theta^t}(\mathbf{z})}}$$

But since \mathbf{z} is continuous, calculating $\underline{q}_{i,z}^t$ is difficult.

intractable
Using EM is not promising

Training VAE(Continued)

To compute $q_{i,z}$ we can use another neural network $q_\phi(z|x^i)$.

$$\begin{aligned} & \sum_{i=1}^N \max_{q_{i,z}} \int_z (q_{i,z} \log P_\theta(x^i|z) P_\theta(z) - q_{i,z} \log q_{i,z}) dz \\ &= \max_{\phi} \sum_{i=1}^N \int_z (q_\phi(z|x^i) \log P_\theta(x^i|z) - \int_z q_\phi(z|x^i) \log \frac{q_\phi(z|x^i)}{P_\theta(z)} dz) \\ &\quad \text{first term} \approx q_{i,z} \quad \text{second term} \quad \text{cosmetic} \end{aligned}$$

Estimating the first term

We use sampling to estimate the integral

$$\max_{\theta, \phi} \sum_{i=1}^N \int_{\mathbf{z}} (q_{\phi}(\mathbf{z}|\mathbf{x}^i) \log P_{\theta}(\mathbf{x}^i|\mathbf{z}) d\mathbf{z} = \max_{\theta, \phi} \sum_{i=1}^N E_{q_{\phi}(\mathbf{z}|\mathbf{x}^i)} [\log P_{\theta}(\mathbf{x}^i|\mathbf{z})]$$

(approximating expectation with sampling.)

$$\approx \max_{\theta, \phi} \sum_{i=1}^N \frac{1}{r} \sum_{j=1}^r \log P_{\theta}(\mathbf{x}^i|\mathbf{z}_j)$$

where $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_r \sim q_{\phi}(\mathbf{z}|\mathbf{x}^i)$ main constraint.

How to sample from q_{ϕ} while still allowing the gradients to flow through parameters ϕ ?

Approximating expected values with averages on a sample.

$$E_{P(x)} [g(x)] = \int p(x) g(x) dx$$

$x^1, x^2, \dots, x^r \sim P(x)$

$$\approx \frac{1}{r} \sum_{i=1}^r g(x^i)$$

Reparameterization trick

Assume $q_\phi(\mathbf{z}|\mathbf{x}^i) = \prod_k q_\phi(z_k|\mathbf{x}^i)$

Let $q_\phi(z_k|\mathbf{x}^i) \sim \mathcal{N}(\mu_{z_k|\mathbf{x}^i}, \sigma_{z_k|\mathbf{x}^i}^2)$

To sample from $\mathcal{N}(\mu_{z|\mathbf{x}^i}, \sigma_{z|\mathbf{x}^i}^2)$, we use **Reparameterization trick**.

Sample from a parameterless distribution. Here, we choose $\mathcal{N}(0, I)$.

Let the samples be $v_1 \dots v_r$

$z_j = v_j \sigma_{z|\mathbf{x}^i} + \mu_{z|\mathbf{x}^i}$ are samples from $\mathcal{N}(\mu_{z|\mathbf{x}^i}, \sigma_{z|\mathbf{x}^i}^2)$

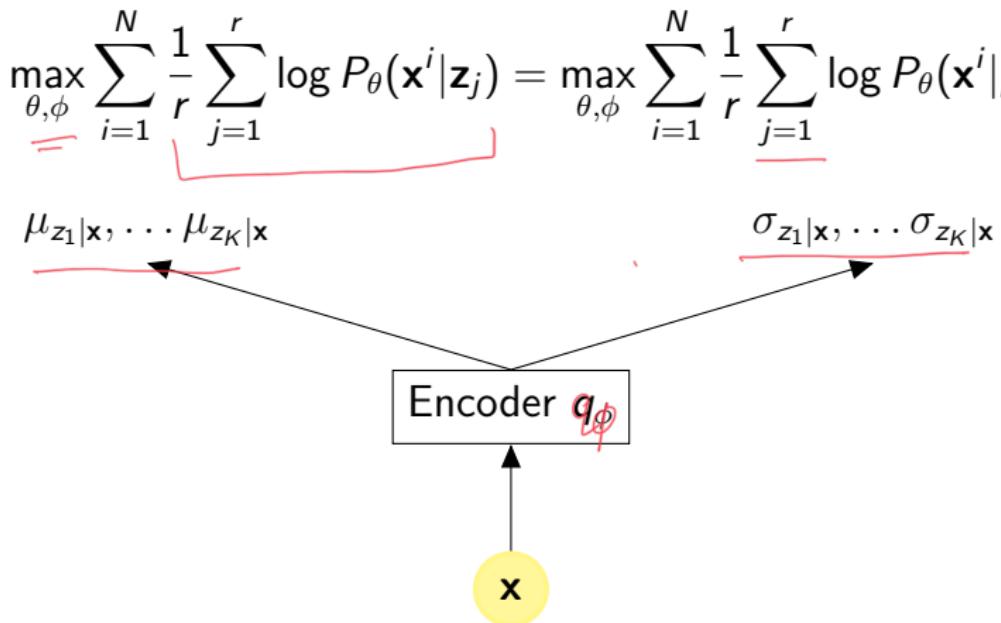
$$(z_j) \sim \mathcal{N}(\bar{\mu}_{z_j|\mathbf{x}}, \bar{\sigma}_{z_j|\mathbf{x}}^2)$$

$$v_j \sim \mathcal{N}(0, 1)$$

$$\mathcal{N}(0, 1)$$

$$v_j = \bar{\mu}_{z_j|\mathbf{x}} + \bar{\sigma}_{z_j|\mathbf{x}} z_j$$

Training VAE(Continued)



v. samples
do not
involve ϕ
 $\sim \mathcal{N}(0, I)$

Calculating second term: KL distance

$\int_z q_\phi(z|x^i) \log \frac{q_\phi(z|x^i)}{P_\theta(z)}$ is the KL distance between $q_\phi(z|x^i)$ and $P_\theta(z)$.

KL distance in closed form

Assume $P_\theta(z) = \mathcal{N}(0, 1)$, $q(z|x^i) = \mathcal{N}(\mu_i, \sigma_i^2)$

$$\begin{aligned} \text{KL}(q(z|x), P(z)) &= \int_z q(z|x) \log \frac{q(z|x^i)}{P(z)} \\ &= \int_z \mathcal{N}(\mu_i, \sigma_i^2) 0.5 \left[-\frac{(z - \mu_i)^2}{\sigma_i^2} - \log \sigma_i^2 + z^2 \right] \\ &= 0.5 \left[-\frac{\sigma_i^2}{\sigma_i^2} - \log \sigma_i^2 + \int_z z^2 \mathcal{N}(\mu_i, \sigma_i^2) \right] \\ &= 0.5[-1.0 - \log \sigma_i^2 + \mu_i^2 + \sigma_i^2] \end{aligned}$$