

5. Scalable Bayesian methods.

1) Probabilistic PCA.

~~Classical PCA~~: Classical PCA:

$$p(x, z | \theta) = \prod_{i=1}^n \mathcal{N}(x_i | V z_i, \sigma^2 I) \mathcal{N}(z_i | 0, I),$$

$\mathbb{R}^D \quad D \gg d$

$$\theta = [V, \sigma^2]$$

Can use EM to find $\arg \max_{\theta} p(x_{tr} | \theta)$. Why?

- EM updates have compl. $O(nDd)$, analytical solution - $O(nD^2)$
- Can process missing parts in x (or account some known z_j)
- Can determine if $p(\theta)$ is established.
- Can be extended to more general models.

Mixture of PCA:

$$p(x, z, \tau | \theta) = \prod_{i=1}^n p(x_i | t_i, z_i, \theta) p(z_i | \theta) p(t_i | \theta) = \prod_{i=1}^n \mathcal{N}(x_i | V_{t_i} z_i, \sigma_{t_i}^2 I) \cdot \mathcal{N}(z_i | 0, I) \cdot \pi_{t_i}.$$

~~Can be~~ Delivers non-linear dimension reduction.

$$E\text{-step: } q(z, \tau) = \prod_{i=1}^n \frac{\mathcal{N}(x_i | V_{t_i} z_i, \sigma_{t_i}^2 I) \mathcal{N}(z_i | 0, I) \pi_{t_i}}{\int_{t_i, z} \dots}$$

$$M\text{-step: } \mathbb{E}_{z, \tau} \log p(x, z, \tau | \theta) \rightarrow \max_{\theta}$$

2) Non-linear model and VAE.

How to account non-linear subspaces?

Use: $p(x, z | \theta) = \underbrace{p(x | z, \theta)}_{\text{generator (decoder)}} \underbrace{p(z | \theta)}_{\text{prior}}$

Idea: use Neural Network.

$$p(x, z | \theta) = \prod_{i=1}^n p(x_i | z_i, \theta) p(z_i) = \prod_{i=1}^n \left(\prod_{j=1}^D \mathcal{N}(x_{ij} | \mu_j(z_i), \sigma_j^2(z_i)) \right) \mathcal{N}(z_i | 0, I)$$

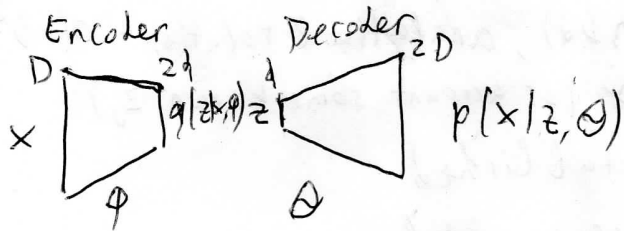
where μ_j and σ_j^2 are calculated by FFN.

$\int p(x_i | z_i, \theta) p(z_i) dz_i$ is intractable, so we couldn't compute MAP estimation ~~and~~ or simple E-step \Rightarrow we use Variational-EM

$$q(z_i | x_i, \phi) \approx \prod_{j=1}^d \mathcal{N}(z_{ij} | \mu_j(x_i), \sigma_j^2(x_i)),$$

$\approx p(z_i | x_i, \theta)$ where the mean and variance are given by another ~~FFN~~ FFN parametrized by ϕ . (Encoder).

General structure.



$$\text{ELBO: } \hat{\mathcal{L}}_V(\phi, \theta) = \int q(z | x, \phi) \log \frac{p(x | z, \theta) p(z)}{q(z | x, \phi)} dz \Rightarrow \max_{\phi, \theta}$$

$$3) \nabla_{\theta} \hat{\mathcal{L}}_V(\phi, \theta) = \nabla_{\theta} \sum_{i=1}^n \int q(z_i | x_i, \phi) \log \frac{p(x_i | z_i, \theta) p(z_i)}{q(z_i | x_i, \phi)} dz_i \approx$$

$$\stackrel{\text{mini-batching}}{\approx} n \int q(z_i | x_i, \phi) \nabla_{\theta} \log p(x_i | z_i, \theta) dz_i, z_i \sim \mathcal{U}\{1, \dots, n\}$$

$$\approx n \nabla_{\theta} \log p(x_i | z_i^*, \theta), z_i^* \sim q(z_i | x_i, \phi)$$

Monte-Carlo estimation.

4) $q(z_i | x_i, \phi)$ depends on $\phi \Rightarrow$ we have problems with $\nabla_{\phi} \hat{\mathcal{L}}_V(\phi, \theta)$ estimation.

$$\nabla_{\phi} \hat{\mathcal{L}}_V(\phi, \theta) = \sum_{i=1}^n$$

Consider differentiation of $\mathbb{E}_{y|x} h(x, y)$ in details:

$$\frac{\partial}{\partial x} \int p(y|x) h(x, y) dy = \int p(y|x) \frac{\partial}{\partial x} h(x, y) dy + \underbrace{\int h(x, y) \frac{\partial}{\partial x} p(y|x) dy}_{\text{not an expectation.}}$$

To deal with second term

we need log-derivative trick: $\frac{\partial}{\partial x} p(y|x) = p(y|x) \frac{\partial}{\partial x} \log p(y|x)$

which yields to:

$$\frac{\partial}{\partial x} \int p(y|x) h(x,y) dy = \int p(y|x) \frac{\partial}{\partial x} h(x,y) dy + \int p(y|x) h(x,y) \frac{\partial}{\partial x} \log p(y|x) dy \approx \frac{\partial}{\partial x} h(x, y_0) + h(x, y_0) \frac{\partial}{\partial x} \log p(y_0|x),$$

$y_0 \sim p(y|x)$ stochastic gradient of $\mathbb{E}_{y|x} h(x,y)$

Overall, we have now:

$$\hat{L}_v(\phi, \theta) \approx \int q(z|x, \phi) \log p(x|z, \theta) dz - D_{KL}(q(z|x, \phi) \| p(z))$$

can be computed and diff. in closed form

$$\underbrace{n \log p(x|z_i^*, \theta)}_{\text{too big at early steps}} \frac{\partial}{\partial \phi} \log q(z_i^*|x_i, \phi)$$

Variance is too big

5) To deal with big variance:

1) REINFORCE - we baseline:

consider $b(z_i, \phi)$ s.t. $\int q(z_i|\phi) b(z_i, \phi) dz_i = 0$,

e.g. $b(z_i, \phi) = b(\phi) \frac{\partial}{\partial \phi} \log q(z_i|\phi)$ - score function;

then: $\frac{\partial}{\partial \phi} \int q(z_i|x_i, \phi) \log p(x_i|z_i, \theta) dz_i =$

$$= \frac{\partial}{\partial \phi} \int q(z_i|x_i, \phi) (\log p(x_i|z_i, \theta) - b(z_i, \phi)) dz_i \approx$$

$$\approx \left(\log p(x_i|z_i^*, \theta) - b(\phi) \right) \frac{\partial}{\partial \phi} \log q(z_i^*|x_i, \phi), \text{ where } z_i^* \sim q(z_i|x_i, \phi)$$

we can make $b(\phi) = b(\phi, x_i)$

if $b(\phi, x_i)$ is close to $\mathbb{E}_{z_i} \log p(x_i|z_i, \theta)$ then it may reduce the variance

practical only for small z

2) Reparameterization trick ('24):

~~Express $y = g(x, \epsilon)$ and~~

To ~~differe~~ get stochastic gradient for $\frac{\partial}{\partial x} \int p(y|x) h(x, y) dy$

express $y = g(x, \epsilon)$, ϵ - random variable. Then:

$$\int p(y|x) h(x, y) dx = \int r(\epsilon) h(x, g(x, \epsilon)) d\epsilon, \text{ and stochastic different.}$$

is simply:

$$\frac{\partial}{\partial x} \int p(y|x) h(x, y) dx = \frac{\partial}{\partial x} \int r(\epsilon) h(x, g(x, \epsilon)) d\epsilon \approx$$

$$\approx \frac{d}{dx} h(x, g(x, \epsilon^1)), \text{ where } \epsilon^1 \sim r(\epsilon)$$

Exampler of reparameterization:

$N(x \mu, \Sigma)$	$p(x y)$	$r(\epsilon)$	$g(y, \epsilon)$
$N(x \mu, \Sigma)$	$N(\epsilon 0, I)$	$x = A\epsilon + \mu, \text{ where } A A^T = \Sigma$	
$\epsilon(x \lambda)$	$U(\epsilon 0, 1)$	$x = -\frac{\log \epsilon}{\lambda}$	
$G(x \beta, 1)$	$G(\epsilon \beta, 1)$	$x = \beta \epsilon$	

Discrete ~~is~~ distrib. cannot be reparameterized. Finally, we have:

$$\frac{\partial}{\partial \phi} \int q(z_i|x_i, \phi) \log p(x_i|z_i, \theta) = \frac{\partial}{\partial \phi} \int r(\epsilon) \log p(x_i|g(\epsilon, x_i, \phi), \theta) d\epsilon \approx$$

$$\approx \frac{\partial}{\partial \phi} \log p(x_i|g(\epsilon^1, x_i, \phi), \theta), \epsilon^1 \sim r(\epsilon)$$

5) VAE: final algorithm.

① Input: Training data x , dimension of latent space d .

②. Pick random $i \sim U\{1, \dots, n\}$ and compute stochastic gradients of ELBO:

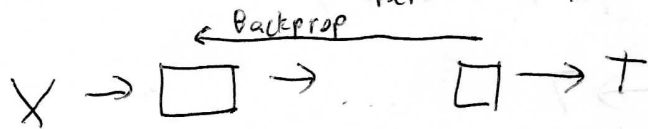
$$\text{- w.r.t. } \theta: \nabla_{\theta} \hat{\mathcal{L}}_V(\phi, \theta) \approx n \frac{\partial}{\partial \theta} \log p(x_i|z_i^*, \theta), z_i^* \sim q(z_i|x_i, \phi)$$

$$\text{- w.r.t. } \phi: \nabla_{\phi} \hat{\mathcal{L}}_V(\phi, \theta) \approx n \left[\frac{\partial}{\partial \phi} \log p(x_i|g(\epsilon^1, x_i, \phi), \theta) - \frac{\partial}{\partial \phi} D_{KL}(q(z_i|x_i, \phi) \| p(z_i)) \right]$$

$\epsilon^1 \sim r(\epsilon)$

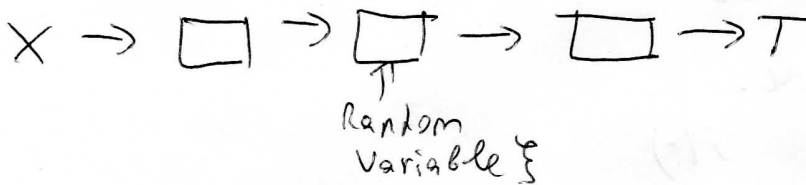
③. Update θ and ϕ using stochastic gradients.

2) Automatic differentiation.



$$\theta_{ML} = \arg \max_{\theta} p(T_{tr} | X_{tr}, \theta)$$

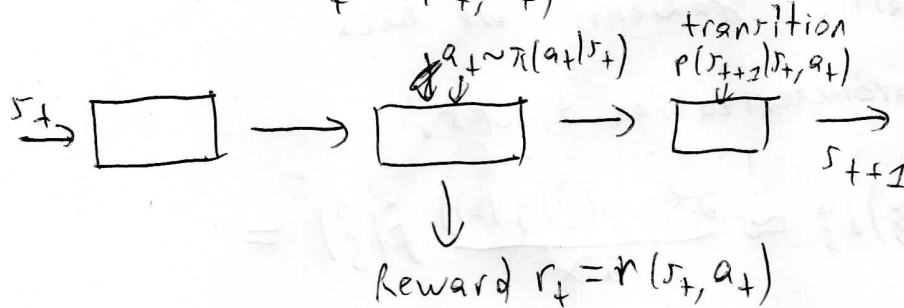
Stochastic computational graph



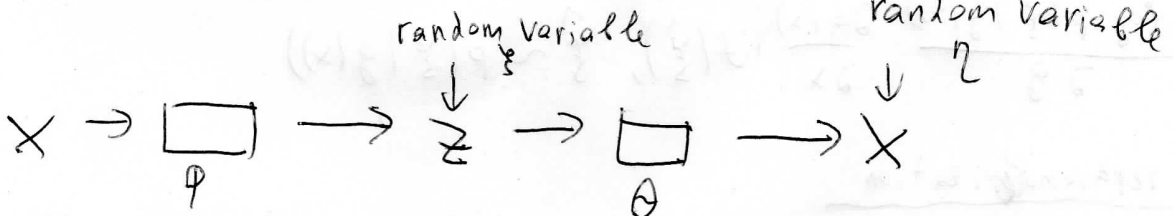
RL: states: $r_t \sim p(r_t | s_{t-1}, a_{t-1})$,

actions: $a_t \sim \pi(a_t | r_t)$

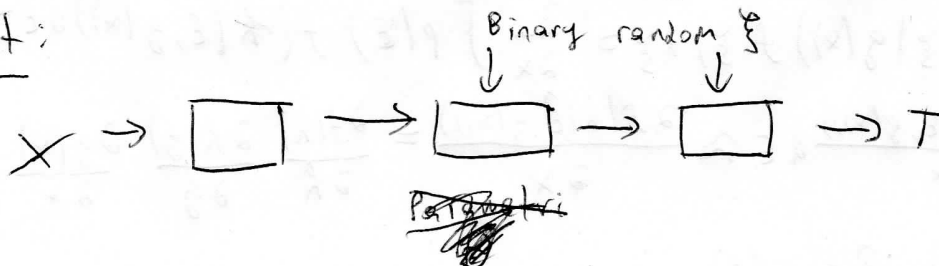
rewards: $r_t = r(s_t, a_t)$



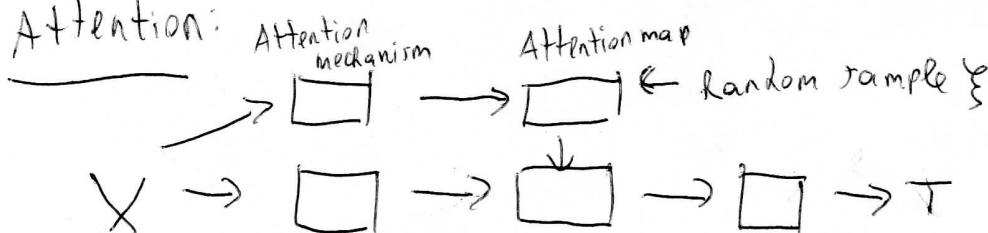
VAEs:



Dropout:



Attention:



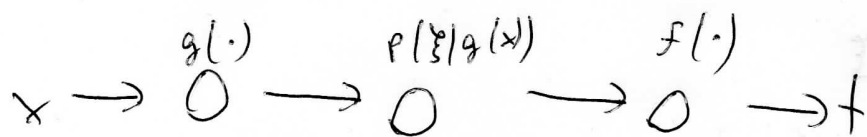
Bayesian regularization:

$$p(\tau, \theta | x) = p(\tau | x, \theta) p(\theta)$$



all steps are stochastic.

With CG: use chain rule.



$$t = \mathbb{E}_z f(z) = \int p(z | g(x)) f(z) dz$$

To propagate the (stochastic) gradients we need

log-derivative or reparameterization tricks.

with log-derivative:

$$\begin{aligned} \frac{\partial t}{\partial x} &= \int \frac{\partial p(z | g(x))}{\partial x} f(z) dz \approx \frac{\partial \log p(z^1 | g(x))}{\partial x} f(z^1) = \\ &= \frac{\partial \log p(z^1 | g(x))}{\partial g} \frac{\partial g(x)}{\partial x} f(z^1), \quad z^1 \sim p(z | g(x)). \end{aligned}$$

with reparameterization:

$$\begin{aligned} \frac{\partial t}{\partial x} &= \frac{\partial}{\partial x} \int p(z | g(x)) f(z) dz = \frac{\partial}{\partial x} \int p(\epsilon) f(h(\epsilon, g(x))) d\epsilon = \\ &= \int p(\epsilon) \frac{\partial f(h(\epsilon, g(x)))}{\partial x} d\epsilon \approx \frac{\partial f(h(\hat{\epsilon}, g(x)))}{\partial x} = \frac{\partial f(h)}{\partial h} \frac{\partial h(g)}{\partial g} \frac{\partial g(x)}{\partial x}, \end{aligned}$$

$$g \Rightarrow h(g) = h(\hat{\epsilon}, g), \quad \hat{\epsilon} \sim p(\epsilon).$$