

# Representation Learning — Module 1

Christopher Zach

Department of Electrical Engineering  
Chalmers University

Lecture, Sep 2020

# Outline

## 1 Latent variable models

- Variational auto-encoders
- Variational NCE
- Boltzmann machines

## 2 Sparse coding and dictionary learning

## 3 Exercise 2

## Latent variable model (LVM)

A LVM is a joint probability distribution  $p_\theta(x, z)$ , where  $x$  is observed in training data but  $z$  is not. We only can match the data distribution  $p_d(x)$  (via samples) with the marginal of the model

$$p_\theta(x) = \int p_\theta(x, z) dz$$

- In most interesting cases computing the marginal  $p_\theta(x)$  is intractable

## Discussion

When do we expect exact computation of  $p_\theta(x)$  to be intractable?

# Latent variable models

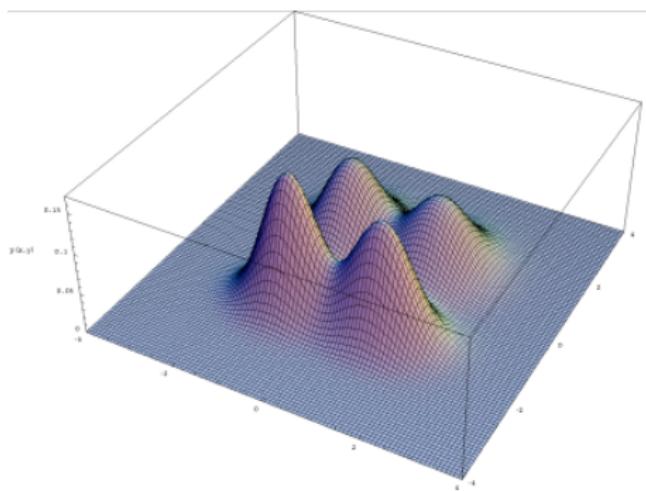
## Examples

- Mixture models

$$p_{\theta}(x|z) = p_{\theta}^z(x) \quad z \in \{1, \dots, K\} \quad p(z=k) = w_k$$

$$p_{\theta}(x) = \sum_k w_k p_{\theta}^k(x)$$

We cannot observe which mixture component  $x$  came from



# Latent variable models

Examples (all usually with intractable partition function)

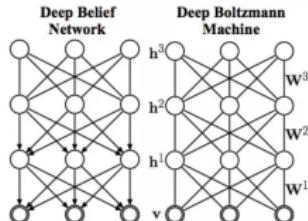
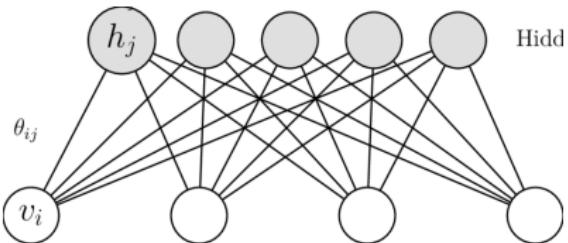
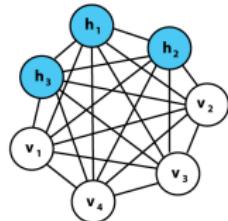
- Boltzmann machines (binary random vectors for visibles  $x$  and latent variables  $z$ )

$$p_{\theta}(x, z) = \frac{1}{Z} \exp \left( \begin{pmatrix} x \\ z \end{pmatrix}^T W \begin{pmatrix} x \\ z \end{pmatrix} + a^T x + b^T z \right) \quad \text{diag}(W) = 0$$

- Restricted Boltzmann machines

$$p_{\theta}(x, z) = \frac{1}{Z} \exp (x^T W z + a^T x + b^T z)$$

- Deep belief nets and deep Boltzmann machines



# Latent variable models

Examples (all usually with intractable partition function)

- Markov random fields / conditional random fields

$$P(z|x) = \frac{1}{Z} \exp \left( - \sum_s \phi_s(z_s | x_s) - \sum_{s \sim t} \phi_{st}(z_s, z_t) \right)$$

- Topic models

- Probabilistic latent semantic analysis (PLSA)

$$P(w, d) = P(d) \sum_c P(c|d) P(w|c)$$

- Latent Dirichlet allocation (LDA)

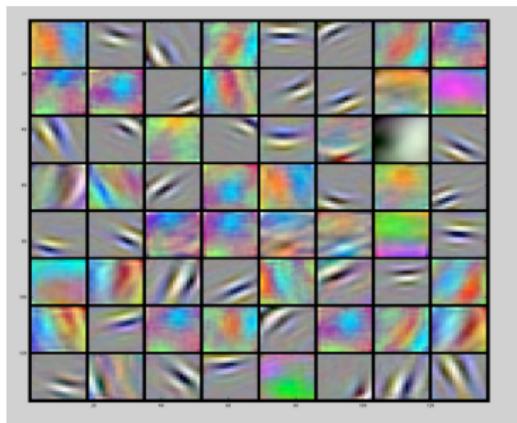
## Latent variable models

Sometimes we optimize over latent variables instead of marginalization

- Sparse dictionary learning / sparse coding

$$-\log p_\theta(x) \doteq \min_z \|x - Dz\|^2 + \lambda \|z\|_\varepsilon \quad \varepsilon \in [0, 1]$$

- Important tool to model images / image patches
- Motivated by the primary visual cortex (V1)
- Usually not interpreted as probabilistic model (but seen as EBM)



# Latent variable models

- Recall

$$p_{\theta}(x) = \int p_{\theta}(x, z) dz$$

- How can we work with latent variables?
- Expectation-maximization algorithm
  - Infer latent variables  $z$  for current parameter
  - Optimize log-likelihood with inferred  $z$
- Learn jointly a posterior  $q(z|x)$ 
  - Variational Bayes, variational inference
  - Variational auto-encoder

# Expectation maximization

- Operational description

$$\begin{aligned} Q(\theta; \theta^{(t)}) &\leftarrow \mathbb{E}_{z \sim p_{\theta^{(t)}}(z|x)} [\log p_{\theta}(x, z)] && \text{E-step} \\ \theta^{(t+1)} &\leftarrow \arg \max_{\theta'} Q(\theta; \theta^{(t)}) && \text{M-step} \end{aligned}$$

- E-step: infer latent variables using current estimate of the posterior
- M-step: maximize likelihood using the now fixed posterior
- In one line

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta'} \mathbb{E}_{z \sim p_{\theta^{(t)}}(z|x)} [\log p_{\theta}(x, z)]$$

- Think of  $x$  as complete training set
  - $\log p_{\theta}(x)$  is sum over training samples

# Expectation maximization

- Operational description

$$\begin{aligned} Q(\theta; \theta^{(t)}) &\leftarrow \mathbb{E}_{z \sim p_{\theta^{(t)}}(z|x)} [\log p_{\theta}(x, z)] && \text{E-step} \\ \theta^{(t+1)} &\leftarrow \arg \max_{\theta'} Q(\theta; \theta^{(t)}) && \text{M-step} \end{aligned}$$

- E-step: infer latent variables using current estimate of the posterior
- M-step: maximize likelihood using the now fixed posterior
- In one line

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta'} \mathbb{E}_{z \sim p_{\theta^{(t)}}(z|x)} [\log p_{\theta}(x, z)]$$

- Think of  $x$  as complete training set
  - $\log p_{\theta}(x)$  is sum over training samples

# Expectation maximization

- Operational description

$$\begin{aligned} Q(\theta; \theta^{(t)}) &\leftarrow \mathbb{E}_{z \sim p_{\theta^{(t)}}(z|x)} [\log p_{\theta}(x, z)] && \text{E-step} \\ \theta^{(t+1)} &\leftarrow \arg \max_{\theta'} Q(\theta; \theta^{(t)}) && \text{M-step} \end{aligned}$$

- E-step: infer latent variables using current estimate of the posterior
- M-step: maximize likelihood using the now fixed posterior
- In one line

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta'} \mathbb{E}_{z \sim p_{\theta^{(t)}}(z|x)} [\log p_{\theta}(x, z)]$$

- Think of  $x$  as complete training set
  - $\log p_{\theta}(x)$  is sum over training samples

# Expectation maximization

- Recall

$$Q(\theta; \theta^{(t)}) \leftarrow \mathbb{E}_{z \sim p_{\theta^{(t)}}(z|x)} [\log p_\theta(x, z)] \quad \text{E-step}$$

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta'} Q(\theta; \theta^{(t)}) \quad \text{M-step}$$

- $Q(\theta; \theta')$  is lower bound of  $\log p_\theta(x)$

$$\begin{aligned} \log p_\theta(x) &= \log \frac{p_\theta(x)p_\theta(x, z)}{p_\theta(x, z)} = \log \frac{p_\theta(x, z)}{p_\theta(z|x)} \quad \forall z : p_\theta(z|x) > 0 \\ &= \log p_\theta(x, z) - \log p_\theta(z|x) \end{aligned}$$

Take the expectation w.r.t.  $p_{\theta'}(z|x)$

$$\begin{aligned} \log p_\theta(x) &= \mathbb{E}_{z \sim p_{\theta'}(z|x)} [\log p_\theta(x, z) - \log p_\theta(z|x)] \\ &= \mathbb{E}_{z \sim p_{\theta'}(z|x)} [\log p_\theta(x, z)] - \mathbb{E}_{z \sim p_{\theta'}(z|x)} [\log p_\theta(z|x)] \\ &= Q(\theta; \theta') - \underbrace{\mathbb{E}_{z \sim p_{\theta'}(z|x)} [\log p_\theta(z|x)]}_{\leq 0} \geq Q(\theta; \theta') \end{aligned}$$

# Expectation maximization

- Recall

$$Q(\theta; \theta^{(t)}) \leftarrow \mathbb{E}_{z \sim p_{\theta^{(t)}}(z|x)} [\log p_\theta(x, z)] \quad \text{E-step}$$

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta'} Q(\theta; \theta^{(t)}) \quad \text{M-step}$$

- $Q(\theta; \theta')$  is lower bound of  $\log p_\theta(x)$

$$\begin{aligned} \log p_\theta(x) &= \log \frac{p_\theta(x)p_\theta(x, z)}{p_\theta(x, z)} = \log \frac{p_\theta(x, z)}{p_\theta(z|x)} \quad \forall z : p_\theta(z|x) > 0 \\ &= \log p_\theta(x, z) - \log p_\theta(z|x) \end{aligned}$$

Take the expectation w.r.t.  $p_{\theta'}(z|x)$

$$\begin{aligned} \log p_\theta(x) &= \mathbb{E}_{z \sim p_{\theta'}(z|x)} [\log p_\theta(x, z) - \log p_\theta(z|x)] \\ &= \mathbb{E}_{z \sim p_{\theta'}(z|x)} [\log p_\theta(x, z)] - \mathbb{E}_{z \sim p_{\theta'}(z|x)} [\log p_\theta(z|x)] \\ &= Q(\theta; \theta') - \underbrace{\mathbb{E}_{z \sim p_{\theta'}(z|x)} [\log p_\theta(z|x)]}_{\leq 0} \geq Q(\theta; \theta') \end{aligned}$$

# Expectation maximization

- Recall

$$Q(\theta; \theta^{(t)}) \leftarrow \mathbb{E}_{z \sim p_{\theta^{(t)}}(z|x)} [\log p_\theta(x, z)] \quad \text{E-step}$$

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta'} Q(\theta; \theta^{(t)}) \quad \text{M-step}$$

- $Q(\theta; \theta')$  is lower bound of  $\log p_\theta(x)$

$$\begin{aligned} \log p_\theta(x) &= \log \frac{p_\theta(x)p_\theta(x, z)}{p_\theta(x, z)} = \log \frac{p_\theta(x, z)}{p_\theta(z|x)} \quad \forall z : p_\theta(z|x) > 0 \\ &= \log p_\theta(x, z) - \log p_\theta(z|x) \end{aligned}$$

Take the expectation w.r.t.  $p_{\theta'}(z|x)$

$$\begin{aligned} \log p_\theta(x) &= \mathbb{E}_{z \sim p_{\theta'}(z|x)} [\log p_\theta(x, z) - \log p_\theta(z|x)] \\ &= \mathbb{E}_{z \sim p_{\theta'}(z|x)} [\log p_\theta(x, z)] - \mathbb{E}_{z \sim p_{\theta'}(z|x)} [\log p_\theta(z|x)] \\ &= Q(\theta; \theta') - \underbrace{\mathbb{E}_{z \sim p_{\theta'}(z|x)} [\log p_\theta(z|x)]}_{\leq 0} \geq Q(\theta; \theta') \end{aligned}$$

# Expectation maximization

- Cross entropy between  $p$  and  $q$

$$H(p, q) = -\mathbb{E}_{x \sim p} [\log q(x)] = H(p) + D_{KL}(p\|q) \geq H(p) = H(p, p)$$

- Recall

$$\begin{aligned}\log p_\theta(x) &= Q(\theta; \theta') - \underbrace{\mathbb{E}_{z \sim p_{\theta'}(z|x)} [\log p_\theta(z|x)]}_{=H(p_{\theta'}(z|x), p_\theta(z|x))}\end{aligned}$$

- Therefore

$$\begin{aligned}\log p_\theta(x) &= Q(\theta; \theta') + H(p_{\theta'}(z|x), p_\theta(z|x)) \\ \log p_{\theta'}(x) &= Q(\theta'; \theta') + H(p_{\theta'}(z|x))\end{aligned}$$

- Assume  $\theta$  maximizes  $Q(\theta; \theta')$  for given  $\theta'$

$$\begin{aligned}0 &\leq Q(\theta; \theta') - Q(\theta'; \theta') \\ &= \log p_\theta(x) - \log p_{\theta'}(x) + \underbrace{H(p_{\theta'}(z|x)) - H(p_{\theta'}(z|x), p_\theta(z|x))}_{\leq 0} \leq \log p_\theta(x) - \log p_{\theta'}(x)\end{aligned}$$

- Updating  $\theta^{(t+1)} \leftarrow \arg \max_\theta Q(\theta; \theta^{(t)})$  improves marginal likelihood

$$p_{\theta^{(t+1)}}(x) \geq p_{\theta^{(t)}}(x)$$

# Expectation maximization

- Cross entropy between  $p$  and  $q$

$$H(p, q) = -\mathbb{E}_{x \sim p} [\log q(x)] = H(p) + D_{KL}(p\|q) \geq H(p) = H(p, p)$$

- Recall

$$\log p_\theta(x) = Q(\theta; \theta') - \underbrace{\mathbb{E}_{z \sim p_{\theta'}(z|x)} [\log p_\theta(z|x)]}_{=H(p_{\theta'}(z|x), p_\theta(z|x))}$$

- Therefore

$$\begin{aligned}\log p_\theta(x) &= Q(\theta; \theta') + H(p_{\theta'}(z|x), p_\theta(z|x)) \\ \log p_{\theta'}(x) &= Q(\theta'; \theta') + H(p_{\theta'}(z|x))\end{aligned}$$

- Assume  $\theta$  maximizes  $Q(\theta; \theta')$  for given  $\theta'$

$$\begin{aligned}0 &\leq Q(\theta; \theta') - Q(\theta'; \theta') \\ &= \log p_\theta(x) - \log p_{\theta'}(x) + \underbrace{H(p_{\theta'}(z|x)) - H(p_{\theta'}(z|x), p_\theta(z|x))}_{\leq 0} \leq \log p_\theta(x) - \log p_{\theta'}(x)\end{aligned}$$

- Updating  $\theta^{(t+1)} \leftarrow \arg \max_\theta Q(\theta; \theta^{(t)})$  improves marginal likelihood

$$p_{\theta^{(t+1)}}(x) \geq p_{\theta^{(t)}}(x)$$

# Expectation maximization

- Cross entropy between  $p$  and  $q$

$$H(p, q) = -\mathbb{E}_{x \sim p} [\log q(x)] = H(p) + D_{KL}(p\|q) \geq H(p) = H(p, p)$$

- Recall

$$\log p_\theta(x) = Q(\theta; \theta') - \underbrace{\mathbb{E}_{z \sim p_{\theta'}(z|x)} [\log p_\theta(z|x)]}_{=H(p_{\theta'}(z|x), p_\theta(z|x))}$$

- Therefore

$$\begin{aligned}\log p_\theta(x) &= Q(\theta; \theta') + H(p_{\theta'}(z|x), p_\theta(z|x)) \\ \log p_{\theta'}(x) &= Q(\theta'; \theta') + H(p_{\theta'}(z|x))\end{aligned}$$

- Assume  $\theta$  maximizes  $Q(\theta; \theta')$  for given  $\theta'$

$$\begin{aligned}0 &\leq Q(\theta; \theta') - Q(\theta'; \theta') \\ &= \log p_\theta(x) - \log p_{\theta'}(x) + \underbrace{H(p_{\theta'}(z|x)) - H(p_{\theta'}(z|x), p_\theta(z|x))}_{\leq 0} \leq \log p_\theta(x) - \log p_{\theta'}(x)\end{aligned}$$

- Updating  $\theta^{(t+1)} \leftarrow \arg \max_\theta Q(\theta; \theta^{(t)})$  improves marginal likelihood

$$p_{\theta^{(t+1)}}(x) \geq p_{\theta^{(t)}}(x)$$

# Expectation maximization

## EM as alternating optimization

- Consider

$$\begin{aligned} F(q, \theta) &:= \log p_\theta(x) - D_{KL}(q \| p_\theta(\cdot|x)) \\ &= \log p_\theta(x) - \int q(z) \log \frac{q(z)}{p_\theta(z|x)} dz \\ &= \mathbb{E}_{z \sim q} \left[ \log p_\theta(x) - \log \frac{q(z)}{p_\theta(z|x)} \right] = \mathbb{E}_{z \sim q} \left[ \log \frac{p_\theta(x)p_\theta(z|x)}{q(z)} \right] \\ &= \mathbb{E}_{z \sim q} \left[ \log \frac{p_\theta(x, z)}{q(z)} \right] = \mathbb{E}_{z \sim q} [\log p_\theta(x, z)] + H(q) \end{aligned}$$

- Optimize w.r.t.  $q$  and  $\theta$  by alternating

$$q^{(t+1)} \leftarrow \arg \max_q -D_{KL}(q \| p_\theta(\cdot|x)) = p_\theta(\cdot|x)$$

$$\theta^{(t+1)} \leftarrow \arg \max_\theta \mathbb{E}_{z \sim q^{(t+1)}} [\log p_\theta(x, z)]$$

- This interpretation is useful when exact optimization w.r.t.  $q$  (E-step) is hard

# Expectation maximization

## EM as alternating optimization

- Consider

$$\begin{aligned} F(q, \theta) &:= \log p_\theta(x) - D_{KL}(q \| p_\theta(\cdot|x)) \\ &= \log p_\theta(x) - \int q(z) \log \frac{q(z)}{p_\theta(z|x)} dz \\ &= \mathbb{E}_{z \sim q} \left[ \log p_\theta(x) - \log \frac{q(z)}{p_\theta(z|x)} \right] = \mathbb{E}_{z \sim q} \left[ \log \frac{p_\theta(x)p_\theta(z|x)}{q(z)} \right] \\ &= \mathbb{E}_{z \sim q} \left[ \log \frac{p_\theta(x, z)}{q(z)} \right] = \mathbb{E}_{z \sim q} [\log p_\theta(x, z)] + H(q) \end{aligned}$$

- Optimize w.r.t.  $q$  and  $\theta$  by alternating

$$q^{(t+1)} \leftarrow \arg \max_q -D_{KL}(q \| p_\theta(\cdot|x)) = p_\theta(\cdot|x)$$

$$\theta^{(t+1)} \leftarrow \arg \max_\theta \mathbb{E}_{z \sim q^{(t+1)}} [\log p_\theta(x, z)]$$

- This interpretation is useful when exact optimization w.r.t.  $q$  (E-step) is hard

# Expectation maximization

## EM as alternating optimization

- Consider

$$\begin{aligned} F(q, \theta) &:= \log p_\theta(x) - D_{KL}(q \| p_\theta(\cdot|x)) \\ &= \log p_\theta(x) - \int q(z) \log \frac{q(z)}{p_\theta(z|x)} dz \\ &= \mathbb{E}_{z \sim q} \left[ \log p_\theta(x) - \log \frac{q(z)}{p_\theta(z|x)} \right] = \mathbb{E}_{z \sim q} \left[ \log \frac{p_\theta(x)p_\theta(z|x)}{q(z)} \right] \\ &= \mathbb{E}_{z \sim q} \left[ \log \frac{p_\theta(x, z)}{q(z)} \right] = \mathbb{E}_{z \sim q} [\log p_\theta(x, z)] + H(q) \end{aligned}$$

- Optimize w.r.t.  $q$  and  $\theta$  by alternating

$$q^{(t+1)} \leftarrow \arg \max_q -D_{KL}(q \| p_\theta(\cdot|x)) = p_\theta(\cdot|x)$$

$$\theta^{(t+1)} \leftarrow \arg \max_\theta \mathbb{E}_{z \sim q^{(t+1)}} [\log p_\theta(x, z)]$$

- This interpretation is useful when exact optimization w.r.t.  $q$  (E-step) is hard

# Expectation maximization

## EM as alternating optimization

- Consider

$$\begin{aligned} F(q, \theta) &:= \log p_\theta(x) - D_{KL}(q \| p_\theta(\cdot|x)) \\ &= \log p_\theta(x) - \int q(z) \log \frac{q(z)}{p_\theta(z|x)} dz \\ &= \mathbb{E}_{z \sim q} \left[ \log p_\theta(x) - \log \frac{q(z)}{p_\theta(z|x)} \right] = \mathbb{E}_{z \sim q} \left[ \log \frac{p_\theta(x)p_\theta(z|x)}{q(z)} \right] \\ &= \mathbb{E}_{z \sim q} \left[ \log \frac{p_\theta(x, z)}{q(z)} \right] = \mathbb{E}_{z \sim q} [\log p_\theta(x, z)] + H(q) \end{aligned}$$

- Optimize w.r.t.  $q$  and  $\theta$  by alternating

$$q^{(t+1)} \leftarrow \arg \max_q -D_{KL}(q \| p_\theta(\cdot|x)) = p_\theta(\cdot|x)$$

$$\theta^{(t+1)} \leftarrow \arg \max_\theta \mathbb{E}_{z \sim q^{(t+1)}} [\log p_\theta(x, z)]$$

- This interpretation is useful when exact optimization w.r.t.  $q$  (E-step) is hard

# Expectation maximization

EM as alternating optimization

- Consider

$$\begin{aligned} F(q, \theta) &:= \log p_\theta(x) - D_{KL}(q \| p_\theta(\cdot|x)) \\ &= \log p_\theta(x) - \int q(z) \log \frac{q(z)}{p_\theta(z|x)} dz \\ &= \mathbb{E}_{z \sim q} \left[ \log p_\theta(x) - \log \frac{q(z)}{p_\theta(z|x)} \right] = \mathbb{E}_{z \sim q} \left[ \log \frac{p_\theta(x)p_\theta(z|x)}{q(z)} \right] \\ &= \mathbb{E}_{z \sim q} \left[ \log \frac{p_\theta(x, z)}{q(z)} \right] = \mathbb{E}_{z \sim q} [\log p_\theta(x, z)] + H(q) \end{aligned}$$

- Optimize w.r.t.  $q$  and  $\theta$  by alternating

$$q^{(t+1)} \leftarrow \arg \max_q -D_{KL}(q \| p_\theta(\cdot|x)) = p_\theta(\cdot|x)$$

$$\theta^{(t+1)} \leftarrow \arg \max_\theta \mathbb{E}_{z \sim q^{(t+1)}} [\log p_\theta(x, z)]$$

- This interpretation is useful when exact optimization w.r.t.  $q$  (E-step) is hard

## Evidence lower bound (ELBO)

- One of *the* workhorses in ML
- Closely linked with EM algorithm
- ELBO: evidence lower bound

$$\begin{aligned}\log p_\theta(x) &= \log \int p_\theta(x, z) dz \\&= \log \int \frac{q(z|x)p_\theta(x, z)}{q(z|x)} dx \\&\geq \int q(z|x) \log \frac{p_\theta(x, z)}{q(z|x)} dx \\&= \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q(z|x)} \right] \\&= \underbrace{\mathbb{E}_{z \sim q(\cdot|x)} [\log p_\theta(x, z)] + H(q(\cdot|x))}_{=: \mathcal{L}}\end{aligned}$$

Def. of marginal prob.

Insert 1 =  $\frac{q(z|x)}{q(z|x)}$

Jensen ineq.

Identify expectation

Identify entropy

## Evidence lower bound (ELBO)

- One of *the* workhorses in ML
- Closely linked with EM algorithm
- ELBO: evidence lower bound

$$\begin{aligned}\log p_\theta(x) &= \log \int p_\theta(x, z) dz && \text{Def. of marginal prob.} \\ &= \log \int \frac{q(z|x)p_\theta(x, z)}{q(z|x)} dx && \text{Insert 1} = \frac{q(z|x)}{q(z|x)} \\ &\geq \int q(z|x) \log \frac{p_\theta(x, z)}{q(z|x)} dx && \text{Jensen ineq.} \\ &= \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q(z|x)} \right] && \text{Identify expectation} \\ &= \underbrace{\mathbb{E}_{z \sim q(\cdot|x)} [\log p_\theta(x, z)] + H(q(\cdot|x))}_{=: \mathcal{L}} && \text{Identify entropy}\end{aligned}$$

# Evidence lower bound (ELBO)

- One of *the* workhorses in ML
- Closely linked with EM algorithm
- ELBO: evidence lower bound

$$\begin{aligned}\log p_\theta(x) &= \log \int p_\theta(x, z) dz \\&= \log \int \frac{q(z|x)p_\theta(x, z)}{q(z|x)} dx \\&\geq \int q(z|x) \log \frac{p_\theta(x, z)}{q(z|x)} dx \\&= \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q(z|x)} \right] \\&= \underbrace{\mathbb{E}_{z \sim q(\cdot|x)} [\log p_\theta(x, z)] + H(q(\cdot|x))}_{=: \mathcal{L}}\end{aligned}$$

Def. of marginal prob.

$$\text{Insert 1} = \frac{q(z|x)}{q(z|x)}$$

Jensen ineq.

Identify expectation

Identify entropy

## Evidence lower bound (ELBO)

- One of *the* workhorses in ML
- Closely linked with EM algorithm
- ELBO: evidence lower bound

$$\begin{aligned}\log p_\theta(x) &= \log \int p_\theta(x, z) dz \\&= \log \int \frac{q(z|x)p_\theta(x, z)}{q(z|x)} dx \\&\geq \int q(z|x) \log \frac{p_\theta(x, z)}{q(z|x)} dx \\&= \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q(z|x)} \right] \\&= \underbrace{\mathbb{E}_{z \sim q(\cdot|x)} [\log p_\theta(x, z)] + H(q(\cdot|x))}_{=: \mathcal{L}}\end{aligned}$$

Def. of marginal prob.

$$\text{Insert 1} = \frac{q(z|x)}{q(z|x)}$$

Jensen ineq.

Identify expectation

Identify entropy

## Evidence lower bound (ELBO)

- One of *the* workhorses in ML
- Closely linked with EM algorithm
- ELBO: evidence lower bound

$$\begin{aligned}\log p_\theta(x) &= \log \int p_\theta(x, z) dz \\&= \log \int \frac{q(z|x)p_\theta(x, z)}{q(z|x)} dx \\&\geq \int q(z|x) \log \frac{p_\theta(x, z)}{q(z|x)} dx \\&= \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q(z|x)} \right] \\&= \underbrace{\mathbb{E}_{z \sim q(\cdot|x)} [\log p_\theta(x, z)] + H(q(\cdot|x))}_{=: \mathcal{L}}\end{aligned}$$

Def. of marginal prob.

$$\text{Insert 1} = \frac{q(z|x)}{q(z|x)}$$

Jensen ineq.

Identify expectation

Identify entropy

# Evidence lower bound (ELBO)

- ELBO  $\mathcal{L}$  can be rewritten as

$$\begin{aligned}\mathcal{L} &= \mathbb{E}_{z \sim q(\cdot|x)} [\log p_\theta(x, z)] + H(q(\cdot|x)) \\ &= \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q(z|x)} \right] \\ &= \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(z|x)p_\theta(x)}{q(z|x)} \right] \\ &= \log p_\theta(x) + \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(z|x)}{q(z|x)} \right] \\ &= \log p_\theta(x) - \underbrace{D_{KL}(q(\cdot|x) \| p_\theta(\cdot|x))}_{\geq 0} \leq \log p_\theta(x)\end{aligned}$$

- $\mathcal{L} = \log p_\theta(x)$  (ELBO is tight) iff  $q(\cdot|x) = p_\theta(\cdot|x)$  a.e.

## Evidence lower bound (ELBO)

- ELBO  $\mathcal{L}$  can be rewritten as

$$\begin{aligned}\mathcal{L} &= \mathbb{E}_{z \sim q(\cdot|x)} [\log p_\theta(x, z)] + H(q(\cdot|x)) \\ &= \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q(z|x)} \right] \\ &= \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(z|x)p_\theta(x)}{q(z|x)} \right] \\ &= \log p_\theta(x) + \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(z|x)}{q(z|x)} \right] \\ &= \log p_\theta(x) - \underbrace{D_{KL}(q(\cdot|x) \| p_\theta(\cdot|x))}_{\geq 0} \leq \log p_\theta(x)\end{aligned}$$

- $\mathcal{L} = \log p_\theta(x)$  (ELBO is tight) iff  $q(\cdot|x) = p_\theta(\cdot|x)$  a.e.

## Evidence lower bound (ELBO)

- ELBO  $\mathcal{L}$  can be rewritten as

$$\begin{aligned}\mathcal{L} &= \mathbb{E}_{z \sim q(\cdot|x)} [\log p_\theta(x, z)] + H(q(\cdot|x)) \\ &= \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q(z|x)} \right] \\ &= \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(z|x)p_\theta(x)}{q(z|x)} \right] \\ &= \log p_\theta(x) + \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(z|x)}{q(z|x)} \right] \\ &= \log p_\theta(x) - \underbrace{D_{KL}(q(\cdot|x) \| p_\theta(\cdot|x))}_{\geq 0} \leq \log p_\theta(x)\end{aligned}$$

- $\mathcal{L} = \log p_\theta(x)$  (ELBO is tight) iff  $q(\cdot|x) = p_\theta(\cdot|x)$  a.e.

## Evidence lower bound (ELBO)

- ELBO  $\mathcal{L}$  can be rewritten as

$$\begin{aligned}\mathcal{L} &= \mathbb{E}_{z \sim q(\cdot|x)} [\log p_\theta(x, z)] + H(q(\cdot|x)) \\ &= \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q(z|x)} \right] \\ &= \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(z|x)p_\theta(x)}{q(z|x)} \right] \\ &= \log p_\theta(x) + \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(z|x)}{q(z|x)} \right] \\ &= \log p_\theta(x) - \underbrace{D_{KL}(q(\cdot|x) \| p_\theta(\cdot|x))}_{\geq 0} \leq \log p_\theta(x)\end{aligned}$$

- $\mathcal{L} = \log p_\theta(x)$  (ELBO is tight) iff  $q(\cdot|x) = p_\theta(\cdot|x)$  a.e.

## Evidence lower bound (ELBO)

- ELBO  $\mathcal{L}$  can be rewritten as

$$\begin{aligned}\mathcal{L} &= \mathbb{E}_{z \sim q(\cdot|x)} [\log p_\theta(x, z)] + H(q(\cdot|x)) \\&= \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q(z|x)} \right] \\&= \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(z|x)p_\theta(x)}{q(z|x)} \right] \\&= \log p_\theta(x) + \mathbb{E}_{z \sim q(\cdot|x)} \left[ \log \frac{p_\theta(z|x)}{q(z|x)} \right] \\&= \log p_\theta(x) - \underbrace{D_{KL}(q(\cdot|x) \| p_\theta(\cdot|x))}_{\geq 0} \leq \log p_\theta(x)\end{aligned}$$

- $\mathcal{L} = \log p_\theta(x)$  (ELBO is tight) iff  $q(\cdot|x) = p_\theta(\cdot|x)$  a.e.

# EM and ELBO / variational Bayes

Are EM and VB the same?

- ELBO/VB is more general than EM (IMHO)
- We usually call it EM if the E-step is tractable
  - Clustering, mixture of distributions
- We call it variational Bayes e.g. when
  - we use a simpler parametric approximation  $q(\cdot|x)$  for the true posterior  $p_\theta(\cdot|x)$
  - we are interested in full posteriors over the parameters  $p(\theta|x)$

$$\log p(\theta|x) \geq \mathbb{E}_{\theta \sim q(\cdot|x)} [\log p(\theta, x)] + H(q(\cdot|x))$$

# Outline

- 1 Latent variable models
  - Variational auto-encoders
  - Variational NCE
  - Boltzmann machines

- 2 Sparse coding and dictionary learning

- 3 Exercise 2

## Overview

The variational auto-encoder (VAE) / auto-encoding variational Bayes

- uses a generative model for  $p_{\theta}(x, z)$
  - identifies the ELBO as sum of reconstruction and regularization terms
  - extends auto-encoders with stochastic latent variables
  - proposes a method to apply back-prop through stochastic network units
- 
- D.P. Kingma & M. Welling, "Auto-Encoding Variational Bayes"
  - Rezende, D. J., S. Mohamed, and D. Wierstra, "Stochastic back-propagation and approximate inference in deep generative models"

## Variational auto-encoders (VAE)

- Recall the ELBO (and parametrize  $q = q_\phi$  via  $\phi$ )

$$\begin{aligned}\mathcal{L} &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x, z)] + H(q_\phi(\cdot|x)) \\ &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q_\phi(z|x)} \right] \leq \log p_\theta(x)\end{aligned}$$

- Generative model for  $p_\theta(x, z)$

$$p_\theta(x, z) = p_\theta(x|z)p_z(z)$$

- $p_z$  user-specified prior over latent variables, e.g.

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

- Typical model for  $p_\theta(x|z)$

$$p_\theta(x|z) = \mathcal{N}(x; \mu_\theta(z), \Sigma_\theta(z))$$

$\mu_\theta$  and  $\Sigma_\theta$  are functions defined via NNs (with parameters  $\theta$ )

## Variational auto-encoders (VAE)

- Recall the ELBO (and parametrize  $q = q_\phi$  via  $\phi$ )

$$\begin{aligned}\mathcal{L} &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x, z)] + H(q_\phi(\cdot|x)) \\ &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q_\phi(z|x)} \right] \leq \log p_\theta(x)\end{aligned}$$

- Generative model for  $p_\theta(x, z)$

$$p_\theta(x, z) = p_\theta(x|z)p_z(z)$$

- $p_z$  user-specified prior over latent variables, e.g.

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

- Typical model for  $p_\theta(x|z)$

$$p_\theta(x|z) = \mathcal{N}(x; \mu_\theta(z), \Sigma_\theta(z))$$

$\mu_\theta$  and  $\Sigma_\theta$  are functions defined via NNs (with parameters  $\theta$ )

# Variational auto-encoders (VAE)

- Recall ELBO and the generative model

$$\mathcal{L} = \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q_\phi(z|x)} \right] \quad p_\theta(x, z) = p_\theta(x|z)p_z(z)$$

- Insert into  $\mathcal{L}$

$$\begin{aligned}\mathcal{L} &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x|z)p_z(z)}{q_\phi(z|x)} \right] \\ &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \log \frac{q_\phi(z|x)}{p_z(z)} \right] \\ &= \underbrace{\mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x|z)]}_{\text{Reconstruction term}} - \underbrace{D_{KL}(q_\phi(\cdot|x) \| p_z)}_{\text{Regularization of } q_\phi} \rightarrow \max_{\theta, \phi}\end{aligned}$$

- Encoder  $q_\phi(z|x)$
- Decoder  $p_\theta(x|z)$
- ELBO is tight if  $q_\phi(z|x) = p_\theta(z|x)$

# Variational auto-encoders (VAE)

- Recall ELBO and the generative model

$$\mathcal{L} = \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q_\phi(z|x)} \right] \quad p_\theta(x, z) = p_\theta(x|z)p_z(z)$$

- Insert into  $\mathcal{L}$

$$\begin{aligned}\mathcal{L} &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x|z)p_z(z)}{q_\phi(z|x)} \right] \\ &= \underbrace{\mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x|z)]}_{\text{Reconstruction term}} - \underbrace{\mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \log \frac{q_\phi(z|x)}{p_z(z)} \right]}_{\text{Regularization of } q_\phi} \\ &= \underbrace{\mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x|z)]}_{\text{Reconstruction term}} - D_{KL}(q_\phi(\cdot|x) \| p_z) \rightarrow \max_{\theta, \phi}\end{aligned}$$

- Encoder  $q_\phi(z|x)$
- Decoder  $p_\theta(x|z)$
- ELBO is tight if  $q_\phi(z|x) = p_\theta(z|x)$

# Variational auto-encoders (VAE)

- Recall ELBO and the generative model

$$\mathcal{L} = \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q_\phi(z|x)} \right] \quad p_\theta(x, z) = p_\theta(x|z)p_z(z)$$

- Insert into  $\mathcal{L}$

$$\begin{aligned}\mathcal{L} &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x|z)p_z(z)}{q_\phi(z|x)} \right] \\ &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \log \frac{q_\phi(z|x)}{p_z(z)} \right] \\ &= \underbrace{\mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x|z)]}_{\text{Reconstruction term}} - \underbrace{D_{KL}(q_\phi(\cdot|x) \| p_z)}_{\text{Regularization of } q_\phi} \rightarrow \max_{\theta, \phi}\end{aligned}$$

- Encoder  $q_\phi(z|x)$
- Decoder  $p_\theta(x|z)$
- ELBO is tight if  $q_\phi(z|x) = p_\theta(z|x)$

## Variational auto-encoders (VAE)

- Recall ELBO and the generative model

$$\mathcal{L} = \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{q_\phi(z|x)} \right] \quad p_\theta(x, z) = p_\theta(x|z)p_z(z)$$

- Insert into  $\mathcal{L}$

$$\begin{aligned}\mathcal{L} &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x|z)p_z(z)}{q_\phi(z|x)} \right] \\ &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \log \frac{q_\phi(z|x)}{p_z(z)} \right] \\ &= \underbrace{\mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x|z)]}_{\text{Reconstruction term}} - \underbrace{D_{KL}(q_\phi(\cdot|x) \| p_z)}_{\text{Regularization of } q_\phi} \rightarrow \max_{\theta, \phi}\end{aligned}$$

- Encoder  $q_\phi(z|x)$
- Decoder  $p_\theta(x|z)$
- ELBO is tight if  $q_\phi(z|x) = p_\theta(z|x)$

# Variational auto-encoders (VAE)

- Why “reconstruction term”

$$\mathcal{L} = \underbrace{\mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x|z)]}_{\text{Reconstruction term}} - \underbrace{D_{KL}(q_\phi(\cdot|x) \| p_z)}_{\text{Regularization of } q_\phi}$$

- With our choice of  $\log p_\theta(x|z) = \mathcal{N}(x; \mu_\theta(z), \Sigma_\theta(z))$

$$\log p_\theta(x|z) \doteq -\frac{1}{2} (x - \mu_\theta(z))^T \Sigma_\theta(z)^{-1} (x - \mu_\theta(z))$$

- Squared Mahalanobis distance between  $x$  and decoded  $\mu_\theta(z)$
- With  $\log p_\theta(x|z) = \mathcal{N}(x; \mu_\theta(z), \sigma^2 I)$

$$\log p_\theta(x|z) \doteq -\frac{\|x - \mu_\theta(z)\|^2}{2\sigma^2}$$

# Variational auto-encoders (VAE)

- Why “reconstruction term”

$$\mathcal{L} = \underbrace{\mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x|z)]}_{\text{Reconstruction term}} - \underbrace{D_{KL}(q_\phi(\cdot|x) \| p_z)}_{\text{Regularization of } q_\phi}$$

- With our choice of  $\log p_\theta(x|z) = \mathcal{N}(x; \mu_\theta(z), \Sigma_\theta(z))$

$$\log p_\theta(x|z) \doteq -\frac{1}{2} (x - \mu_\theta(z))^T \Sigma_\theta(z)^{-1} (x - \mu_\theta(z))$$

- Squared Mahalanobis distance between  $x$  and decoded  $\mu_\theta(z)$
- With  $\log p_\theta(x|z) = \mathcal{N}(x; \mu_\theta(z), \sigma^2 \mathbf{I})$

$$\log p_\theta(x|z) \doteq -\frac{\|x - \mu_\theta(z)\|^2}{2\sigma^2}$$

## Variational auto-encoders (VAE)

- Role of the “regularization term”

$$D_{KL}(q_\phi(\cdot|x) \| p_z)$$

- Favors posterior  $q_\phi(\cdot|x)$  to be “disentangled” if  $p_z(\cdot) = \mathcal{N}(\cdot; 0, I)$ 
  - Components of  $\mathcal{N}(0, I)$  random vector are independent
- Stronger disentanglement by increasing weight:  $\beta$ -VAE

$$\mathcal{L} = \mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x|z)] - \beta D_{KL}(q_\phi(\cdot|x) \| p_z)$$

for a  $\beta \geq 1$

- Higgins et al., “ $\beta$ -VAE: Learning basic visual concepts with a constrained variational framework”

# Variational auto-encoders (VAE)

- How to train a VAE?
- Given training data  $\{x_i\}$

$$J(\theta, \phi) = \frac{1}{N} \sum \left( \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)] - D_{KL}(q_\phi(\cdot|x_i) \| p_z) \right)$$

- 2nd term often has closed-form expression
  - $p_z$  and  $q(\cdot|x_i)$  are Gaussians

$$D_{KL}(\mathcal{N}_0 \| \mathcal{N}_1) \doteq \frac{1}{2} \left( \text{trace}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) + \log \frac{|\Sigma_1|}{|\Sigma_0|} \right)$$

- $p_z(z) = \mathcal{N}(z; \mathbf{0}, \mathbf{I})$

$$D_{KL}(\mathcal{N}(\mu, \Sigma) \| p_z) \doteq \frac{1}{2} \left( \text{trace}(\Sigma) + \|\mu\|^2 - \log |\Sigma| \right)$$

# Variational auto-encoders (VAE)

- How to train a VAE?
- Given training data  $\{x_i\}$

$$J(\theta, \phi) = \frac{1}{N} \sum \left( \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)] - D_{KL}(q_\phi(\cdot|x_i) \| p_z) \right)$$

- 2nd term often has closed-form expression
  - $p_z$  and  $q(\cdot|x_i)$  are Gaussians

$$D_{KL}(\mathcal{N}_0 \| \mathcal{N}_1) \doteq \frac{1}{2} \left( \text{trace}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) + \log \frac{|\Sigma_1|}{|\Sigma_0|} \right)$$

- $p_z(z) = \mathcal{N}(z; \mathbf{0}, \mathbf{I})$

$$D_{KL}(\mathcal{N}(\mu, \Sigma) \| p_z) \doteq \frac{1}{2} \left( \text{trace}(\Sigma) + \|\mu\|^2 - \log |\Sigma| \right)$$

# Variational auto-encoders (VAE)

- How to train a VAE?
- Given training data  $\{x_i\}$

$$J(\theta, \phi) = \frac{1}{N} \sum \left( \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)] - D_{KL}(q_\phi(\cdot|x_i) \| p_z) \right)$$

- 1st term

$$\nabla_\theta \frac{1}{N} \sum \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)] \quad \nabla_\phi \frac{1}{N} \sum \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)]$$

- $\nabla_\theta$ : MC estimator

$$\nabla_\theta \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)] \approx \frac{1}{T} \nabla_\theta \log p_\theta(x_i|z^{(t)}) \quad z^{(t)} \sim q_\phi(z|x_i)$$

Works okay

# Variational auto-encoders (VAE)

- How to train a VAE?
- $\nabla_\phi$ : MC estimator

$$\begin{aligned}\nabla_\phi \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)] &= \nabla_\phi \int q_\phi(z|x_i) \log p_\theta(x_i|z) dz \\ &= \int \frac{q_\phi(z|x_i)}{q_\phi(z|x_i)} \nabla_\phi q_\phi(z|x_i) \log p_\theta(x_i|z) dz \\ &= \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z) \nabla_\phi \log q_\phi(z|x_i)] \\ &\approx \frac{1}{T} \log p_\theta(x_i|z^{(t)}) \nabla_\phi \log q_\phi(z|x_i) \quad z^{(t)} \sim q_\phi(z|x_i)\end{aligned}$$

High variance in practice

# Variational auto-encoders (VAE)

- How to train a VAE?
- $\nabla_\phi$ : MC estimator

$$\begin{aligned}\nabla_\phi \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)] &= \nabla_\phi \int q_\phi(z|x_i) \log p_\theta(x_i|z) dz \\ &= \int \frac{q_\phi(z|x_i)}{q_\phi(z|x_i)} \nabla_\phi q_\phi(z|x_i) \log p_\theta(x_i|z) dz \\ &= \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z) \nabla_\phi \log q_\phi(z|x_i)] \\ &\approx \frac{1}{T} \log p_\theta(x_i|z^{(t)}) \nabla_\phi \log q_\phi(z|x_i) \quad z^{(t)} \sim q_\phi(z|x_i)\end{aligned}$$

High variance in practice

# Variational auto-encoders (VAE)

- How to train a VAE?
- $\nabla_\phi$ : MC estimator

$$\begin{aligned}\nabla_\phi \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)] &= \nabla_\phi \int q_\phi(z|x_i) \log p_\theta(x_i|z) dz \\ &= \int \frac{q_\phi(z|x_i)}{q_\phi(z|x_i)} \nabla_\phi q_\phi(z|x_i) \log p_\theta(x_i|z) dz \\ &= \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z) \nabla_\phi \log q_\phi(z|x_i)] \\ &\approx \frac{1}{T} \log p_\theta(x_i|z^{(t)}) \nabla_\phi \log q_\phi(z|x_i) \quad z^{(t)} \sim q_\phi(z|x_i)\end{aligned}$$

High variance in practice

# Variational auto-encoders (VAE)

- How to train a VAE?
- $\nabla_\phi$ : MC estimator

$$\begin{aligned}\nabla_\phi \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)] &= \nabla_\phi \int q_\phi(z|x_i) \log p_\theta(x_i|z) dz \\ &= \int \frac{q_\phi(z|x_i)}{q_\phi(z|x_i)} \nabla_\phi q_\phi(z|x_i) \log p_\theta(x_i|z) dz \\ &= \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z) \nabla_\phi \log q_\phi(z|x_i)] \\ &\approx \frac{1}{T} \log p_\theta(x_i|z^{(t)}) \nabla_\phi \log q_\phi(z|x_i) \quad z^{(t)} \sim q_\phi(z|x_i)\end{aligned}$$

High variance in practice

# Variational auto-encoders (VAE)

- How to train a VAE?
- Assume Gaussian posterior

$$q(z|x) = \mathcal{N}(z; \mu_\phi(x), \Sigma_\phi(x)) \quad \Sigma_\phi(x) = L_\phi(x)L_\phi(x)^T$$

- “Reparametrization trick”

$$\varepsilon \sim \mathcal{N}(\cdot; \mathbf{0}, \mathbf{I}) \quad z = L_\phi(x)\varepsilon + \mu_\phi(x) \sim \mathcal{N}(\cdot; \mu_\phi(x), \Sigma_\phi(x))$$

- 1st term

$$\mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)] = \mathbb{E}_{\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\log p_\theta(x_i|L_\phi(x)\varepsilon + \mu_\phi(x))]$$

- Monte-Carlo estimate

$$\nabla_\phi \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)] \approx \frac{1}{T} \nabla_\phi \log p_\theta(x_i|L_\phi(x)\varepsilon^{(t)} + \mu_\phi(x)) \quad \varepsilon^{(t)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Works fine

# Variational auto-encoders (VAE)

- How to train a VAE?
- Assume Gaussian posterior

$$q(z|x) = \mathcal{N}(z; \mu_\phi(x), \Sigma_\phi(x)) \quad \Sigma_\phi(x) = L_\phi(x)L_\phi(x)^T$$

- “Reparametrization trick”

$$\varepsilon \sim \mathcal{N}(\cdot; \mathbf{0}, \mathbf{I}) \quad z = L_\phi(x)\varepsilon + \mu_\phi(x) \sim \mathcal{N}(\cdot; \mu_\phi(x), \Sigma_\phi(x))$$

- 1st term

$$\mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)] = \mathbb{E}_{\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\log p_\theta(x_i|L_\phi(x)\varepsilon + \mu_\phi(x))]$$

- Monte-Carlo estimate

$$\nabla_\phi \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)] \approx \frac{1}{T} \nabla_\phi \log p_\theta(x_i|L_\phi(x)\varepsilon^{(t)} + \mu_\phi(x)) \quad \varepsilon^{(t)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Works fine

# Variational auto-encoders (VAE)

- How to train a VAE?
- Assume Gaussian posterior

$$q(z|x) = \mathcal{N}(z; \mu_\phi(x), \Sigma_\phi(x)) \quad \Sigma_\phi(x) = L_\phi(x)L_\phi(x)^T$$

- “Reparametrization trick”

$$\varepsilon \sim \mathcal{N}(\cdot; \mathbf{0}, \mathbf{I}) \quad z = L_\phi(x)\varepsilon + \mu_\phi(x) \sim \mathcal{N}(\cdot; \mu_\phi(x), \Sigma_\phi(x))$$

- 1st term

$$\mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)] = \mathbb{E}_{\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\log p_\theta(x_i|L_\phi(x)\varepsilon + \mu_\phi(x))]$$

- Monte-Carlo estimate

$$\nabla_\phi \mathbb{E}_{z \sim q_\phi(\cdot|x_i)} [\log p_\theta(x_i|z)] \approx \frac{1}{T} \nabla_\phi \log p_\theta(x_i|L_\phi(x)\varepsilon^{(t)} + \mu_\phi(x)) \quad \varepsilon^{(t)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Works fine

# Variational auto-encoders (VAE)

- How to train a VAE?
- “Reparametrization trick”: applicable
  - Location-scale family of distributions (e.g. Gaussian)
  - Easily invertible cdf:  $\varepsilon \sim \mathcal{U}[0, 1]$
  - Functions of easier RVs
- Caveat: requires continuous latent variables
- Reparametrization trick needed for  $\nabla_{\theta}$ ?

$$\nabla_{\theta} \mathbb{E}_{z \sim q_{\phi}(\cdot|x)} [\log p_{\theta}(x|z)] = \nabla_{\theta} \mathbb{E}_{\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\log p_{\theta}(x_i | L_{\phi}(x)\varepsilon + \mu_{\phi}(x))]$$

- No difference (besides different parametrization of  $\Sigma_{\phi}(x)$ )
- Drawing samples from  $q_{\phi}(\cdot|x)$  is equivalent to samples  $L_{\phi}(x)\varepsilon + \mu_{\phi}(x)$ ,  $\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

# Variational auto-encoders (VAE)

Connection to regular AE

- Deterministic encoder

$$z = f_\phi(x)$$

$$q_\phi(z|x) = \delta_{f_\phi(x)}(z)$$

- Gaussian decoder

$$p_\theta(x|z) = \mathcal{N}(x; g_\theta(z), \sigma^2 \mathbf{I})$$

- Gaussian prior  $p_z(z) = \mathcal{N}(z; \mathbf{0}, \mathbf{I})$

$$\begin{aligned}\log p_\theta(x) &\geq \mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x|z)] - D_{KL}(q_\phi(\cdot|x) \| p_z) \\ &= -\frac{1}{2\sigma^2} \|x - g_\theta(f_\phi(x))\|^2 - \frac{1}{2} \|f_\theta(x)\|^2\end{aligned}$$

- Functions  $f_\phi$  and  $g_\theta$  can be deep networks

# Single layer AEs

- We choose

$$f_\phi(x) = s(W^T x) \quad g_\theta(z) = Wz \quad p_\theta(x|z) = \mathcal{N}(x; g_\theta(z), I) \quad p_z(z) = \mathcal{N}(z; 0, \tau^2 I)$$

with  $\tau \gg 0$

- Recall  $s(u) = S'(u)$

## Score matching model

$$E(x) = -\log p(x) = \frac{1}{2} \|x\|^2 - \sum_k S(w_k^T x) + \log Z(W)$$

### Training objective

$$\mathbb{E}_{x \sim p_d} \left[ \frac{1}{2} \|x - Ws(W^T x)\|^2 + \frac{1}{2} \sum_k \|w_k\|^2 s'(w_k^T x) \right] \rightarrow \min_W$$

# Single layer AEs

- We choose

$$f_\phi(x) = s(W^T x) \quad g_\theta(z) = Wz \quad p_\theta(x|z) = \mathcal{N}(x; g_\theta(z), I) \quad p_z(z) = \mathcal{N}(z; 0, \tau^2 I)$$

with  $\tau \gg 0$

- Recall  $s(u) = S'(u)$

## Auto-encoding variational Bayes

$$E(x) = -\log p(x) \leq \frac{1}{2} \|x - Ws(W^T x)\|^2 + \frac{1}{2\tau^2} \|s(W^T x)\|^2$$

Training objective

$$\mathbb{E}_{x \sim p_d} \left[ \frac{1}{2} \|x - Ws(W^T x)\|^2 + \frac{1}{2\tau^2} \sum_k \|s(w_k^T x)\|^2 \right] \rightarrow \min_w$$

# Single layer VAEs

- Encoder model for discrete hiddens  $z_k \in \{0, 1\}$

$$q_\phi(z_k = 1|x) = \sigma(w_k^T x)$$

- Decoder

$$p_\theta(x|z) = \mathcal{N}(x; Wz, I) \quad \log p_\theta(x|z) \doteq -\frac{1}{2} \|x - Wz\|^2$$

- Lower bound

$$\begin{aligned}\mathcal{L} &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x|z)] - D_{KL}(q_\phi(\cdot|x) \| p_z) \\ &\doteq -\frac{1}{2\sigma^2} \|x - W\sigma(W^T x)\|^2 - D_{KL}(q_\phi(\cdot|x) \| p_z)\end{aligned}$$

- $p_{z_k} = \text{Ber}(1/2)$ :  $D_{KL}(q_\phi(\cdot|x) \| p_z) \doteq -\sum_k H(q_\phi(z_k|x))$

## Auto-encoding variational Bayes: Bernoulli hiddens

### Training objective

$$\mathbb{E}_{x \sim p_d} \left[ \frac{1}{2} \|x - Ws(W^T x)\|^2 + \sum_k H(q_\phi(z_k|x)) \right] \rightarrow \min_W$$

# Outline

## 1 Latent variable models

- Variational auto-encoders
- **Variational NCE**
- Boltzmann machines

## 2 Sparse coding and dictionary learning

## 3 Exercise 2

## Overview

Variational noise-contrastive estimation

- combines NCE with variational Bayes
- proposes a lower bound on the NCE objective
- allows unnormalized latent variable models
- Rhodes & Gutmann, "Variational Noise-Contrastive Estimation"

# Variational NCE

- Recall the NCE objective

$$J_{NCE}(\theta) = \mathbb{E}_{x \sim p_d} \left[ \log \frac{p_\theta(x)}{p_\theta(x) + \nu p_n(x)} \right] + \nu \mathbb{E}_{x \sim p_n} \left[ \log \frac{\nu p_n(x)}{p_\theta(x) + \nu p_n(x)} \right]$$

- Now

$$p_\theta(x) = \int p_\theta(x, z) dz$$

- Naively plugging into  $J$ ?

$$J_{NCE}(p_\theta) = \mathbb{E}_{x \sim p_d} \left[ \log \frac{\int p_\theta(x, z) dz}{\int p_\theta(x, z) dz + \nu p_n(x)} \right] + \nu \mathbb{E}_{x \sim p_n} \left[ \log \frac{\nu p_n(x)}{\int p_\theta(x, z) dz + \nu p_n(x)} \right]$$

- We cannot evaluate  $\int p_\theta(x, z) dz$  or  $\nabla_\theta \int p_\theta(x, z) dz$

## Variational NCE

- Recall the NCE objective

$$J_{NCE}(\theta) = \mathbb{E}_{x \sim p_d} \left[ \log \frac{p_\theta(x)}{p_\theta(x) + \nu p_n(x)} \right] + \nu \mathbb{E}_{x \sim p_n} \left[ \log \frac{\nu p_n(x)}{p_\theta(x) + \nu p_n(x)} \right]$$

- Now

$$p_\theta(x) = \int p_\theta(x, z) dz$$

- Naively plugging into  $J$ ?

$$J_{NCE}(p_\theta) = \mathbb{E}_{x \sim p_d} \left[ \log \frac{\int p_\theta(x, z) dz}{\int p_\theta(x, z) dz + \nu p_n(x)} \right] + \nu \mathbb{E}_{x \sim p_n} \left[ \log \frac{\nu p_n(x)}{\int p_\theta(x, z) dz + \nu p_n(x)} \right]$$

- We cannot evaluate  $\int p_\theta(x, z) dz$  or  $\nabla_\theta \int p_\theta(x, z) dz$

# Variational NCE

- Focus on the first term in the NCE objective

$$\begin{aligned} J_1 &:= \mathbb{E}_{x \sim p_d} \left[ \log \frac{p_\theta(x)}{p_\theta(x) + \nu p_n(x)} \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ \log \frac{1}{1 + \frac{\nu p_n(x)}{p_\theta(x)}} \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ -\log \left( 1 + \frac{\nu p_n(x)}{p_\theta(x)} \right) \right] \end{aligned}$$

- The mapping  $g : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$  with

$$g(r) = -\log \left( 1 + \frac{\nu}{r} \right)$$

is concave (why?)

# Variational NCE

- Focus on the first term in the NCE objective

$$\begin{aligned} J_1 &:= \mathbb{E}_{x \sim p_d} \left[ \log \frac{p_\theta(x)}{p_\theta(x) + \nu p_n(x)} \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ \log \frac{1}{1 + \frac{\nu p_n(x)}{p_\theta(x)}} \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ -\log \left( 1 + \frac{\nu p_n(x)}{p_\theta(x)} \right) \right] \end{aligned}$$

- The mapping  $g : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$  with

$$g(r) = -\log \left( 1 + \frac{\nu}{r} \right)$$

is concave (why?)

# Variational NCE

- Focus on the first term in the NCE objective

$$\begin{aligned} J_1 &:= \mathbb{E}_{x \sim p_d} \left[ \log \frac{p_\theta(x)}{p_\theta(x) + \nu p_n(x)} \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ \log \frac{1}{1 + \frac{\nu p_n(x)}{p_\theta(x)}} \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ -\log \left( 1 + \frac{\nu p_n(x)}{p_\theta(x)} \right) \right] \end{aligned}$$

- The mapping  $g : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$  with

$$g(r) = -\log \left( 1 + \frac{\nu}{r} \right)$$

is concave (why?)

## Variational NCE

- Focus on the first term in the NCE objective

$$\begin{aligned} J_1 &:= \mathbb{E}_{x \sim p_d} \left[ \log \frac{p_\theta(x)}{p_\theta(x) + \nu p_n(x)} \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ \log \frac{1}{1 + \frac{\nu p_n(x)}{p_\theta(x)}} \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ -\log \left( 1 + \frac{\nu p_n(x)}{p_\theta(x)} \right) \right] \end{aligned}$$

- The mapping  $g : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$  with

$$g(r) = -\log \left( 1 + \frac{\nu}{r} \right)$$

is concave (why?)

# Variational NCE

- Recall the 1st term

$$J_1 = \mathbb{E}_{x \sim p_d} \left[ -\log \left( 1 + \frac{\nu p_n(x)}{p_\theta(x)} \right) \right]$$

$$= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{p_\theta(x)}{p_n(x)} \right) \right]$$

$$= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{1}{p_n(x)} \int p_\theta(x, z) dz \right) \right]$$

$$= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{1}{p_n(x)} \int \frac{q_\phi(z|x)p_\theta(x, z)}{q_\phi(x|z)} dz \right) \right]$$

$$\geq \mathbb{E}_{x \sim p_d} \left[ \int q_\phi(z|x) g \left( \frac{1}{p_n(x)} \frac{p_\theta(x, z)}{q_\phi(x|z)} \right) dx \right]$$

$$= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ g \left( \frac{p_\theta(x, z)}{q_\phi(x|z)p_n(x)} \right) \right]$$

$$g(r) = -\log \left( 1 + \frac{\nu}{r} \right)$$

$$\text{Insert } 1 = \frac{q_\phi(z|x)}{q_\phi(z|z)}$$

Jensen ineq.

# Variational NCE

- Recall the 1st term

$$J_1 = \mathbb{E}_{x \sim p_d} \left[ -\log \left( 1 + \frac{\nu p_n(x)}{p_\theta(x)} \right) \right]$$

$$= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{p_\theta(x)}{p_n(x)} \right) \right]$$

$$= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{1}{p_n(x)} \int p_\theta(x, z) dz \right) \right]$$

$$= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{1}{p_n(x)} \int \frac{q_\phi(z|x)p_\theta(x, z)}{q_\phi(x|z)} dz \right) \right]$$

$$\geq \mathbb{E}_{x \sim p_d} \left[ \int q_\phi(z|x) g \left( \frac{1}{p_n(x)} \frac{p_\theta(x, z)}{q_\phi(x|z)} \right) dx \right]$$

$$= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ g \left( \frac{p_\theta(x, z)}{q_\phi(x|z)p_n(x)} \right) \right]$$

$$g(r) = -\log \left( 1 + \frac{\nu}{r} \right)$$

$$\text{Insert } 1 = \frac{q_\phi(z|x)}{q_\phi(z|z)}$$

Jensen ineq.

# Variational NCE

- Recall the 1st term

$$J_1 = \mathbb{E}_{x \sim p_d} \left[ -\log \left( 1 + \frac{\nu p_n(x)}{p_\theta(x)} \right) \right]$$

$$= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{p_\theta(x)}{p_n(x)} \right) \right]$$

$$= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{1}{p_n(x)} \int p_\theta(x, z) dz \right) \right]$$

$$= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{1}{p_n(x)} \int \frac{q_\phi(z|x)p_\theta(x, z)}{q_\phi(x|z)} dz \right) \right]$$

$$\geq \mathbb{E}_{x \sim p_d} \left[ \int q_\phi(z|x) g \left( \frac{1}{p_n(x)} \frac{p_\theta(x, z)}{q_\phi(x|z)} \right) dx \right]$$

$$= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ g \left( \frac{p_\theta(x, z)}{q_\phi(x|z)p_n(x)} \right) \right]$$

$$g(r) = -\log \left( 1 + \frac{\nu}{r} \right)$$

$$\text{Insert 1} = \frac{q_\phi(z|x)}{q_\phi(z|z)}$$

Jensen ineq.

# Variational NCE

- Recall the 1st term

$$\begin{aligned} J_1 &= \mathbb{E}_{x \sim p_d} \left[ -\log \left( 1 + \frac{\nu p_n(x)}{p_\theta(x)} \right) \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{p_\theta(x)}{p_n(x)} \right) \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{1}{p_n(x)} \int p_\theta(x, z) dz \right) \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{1}{p_n(x)} \int \frac{q_\phi(z|x)p_\theta(x, z)}{q_\phi(x|z)} dz \right) \right] \\ &\geq \mathbb{E}_{x \sim p_d} \left[ \int q_\phi(z|x) g \left( \frac{1}{p_n(x)} \frac{p_\theta(x, z)}{q_\phi(x|z)} \right) dx \right] \\ &= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ g \left( \frac{p_\theta(x, z)}{q_\phi(x|z)p_n(x)} \right) \right] \end{aligned}$$

$$g(r) = -\log \left( 1 + \frac{\nu}{r} \right)$$

$$\text{Insert 1} = \frac{q_\phi(z|x)}{q_\phi(x|z)}$$

Jensen ineq.

# Variational NCE

- Recall the 1st term

$$\begin{aligned} J_1 &= \mathbb{E}_{x \sim p_d} \left[ -\log \left( 1 + \frac{\nu p_n(x)}{p_\theta(x)} \right) \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{p_\theta(x)}{p_n(x)} \right) \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{1}{p_n(x)} \int p_\theta(x, z) dz \right) \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{1}{p_n(x)} \int \frac{q_\phi(z|x)p_\theta(x, z)}{q_\phi(x|z)} dz \right) \right] \\ &\geq \mathbb{E}_{x \sim p_d} \left[ \int q_\phi(z|x) g \left( \frac{1}{p_n(x)} \frac{p_\theta(x, z)}{q_\phi(x|z)} \right) dx \right] \\ &= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ g \left( \frac{p_\theta(x, z)}{q_\phi(x|z)p_n(x)} \right) \right] \end{aligned}$$

$$g(r) = -\log \left( 1 + \frac{\nu}{r} \right)$$

Insert 1 =  $\frac{q_\phi(z|x)}{q_\phi(x|z)}$

Jensen ineq.

# Variational NCE

- Recall the 1st term

$$\begin{aligned} J_1 &= \mathbb{E}_{x \sim p_d} \left[ -\log \left( 1 + \frac{\nu p_n(x)}{p_\theta(x)} \right) \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{p_\theta(x)}{p_n(x)} \right) \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{1}{p_n(x)} \int p_\theta(x, z) dz \right) \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ g \left( \frac{1}{p_n(x)} \int \frac{q_\phi(z|x)p_\theta(x, z)}{q_\phi(x|z)} dz \right) \right] \\ &\geq \mathbb{E}_{x \sim p_d} \left[ \int q_\phi(z|x) g \left( \frac{1}{p_n(x)} \frac{p_\theta(x, z)}{q_\phi(x|z)} \right) dx \right] \\ &= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ g \left( \frac{p_\theta(x, z)}{q_\phi(x|z)p_n(x)} \right) \right] \end{aligned}$$

$$g(r) = -\log \left( 1 + \frac{\nu}{r} \right)$$

Insert 1 =  $\frac{q_\phi(z|x)}{q_\phi(x|z)}$

Jensen ineq.

# Variational NCE

- Recall

$$\begin{aligned} J_1 &\geq \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ g \left( \frac{p_\theta(x, z)}{q_\phi(x|z)p_n(x)} \right) \right] \\ &= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ -\log \left( 1 + \frac{\nu q_\phi(x|z)p_n(x)}{p_\theta(x, z)} \right) \right] \quad \text{Expand } g \\ &= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{\nu p_n(x) q_\phi(x|z)} \right] \end{aligned}$$

- Lower bound of the NCE objective

$$J_{NCE} \geq \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{\nu p_n(x) q_\phi(x|z)} \right] + \nu \mathbb{E}_{x \sim p_n} \left[ \log \frac{\nu p_n(x)}{p_\theta(x) + \nu p_n(x)} \right]$$

- What should we do with the 2nd term?

- Why does the same trick not work for the 2nd term?

# Variational NCE

- Recall

$$\begin{aligned} J_1 &\geq \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ g \left( \frac{p_\theta(x, z)}{q_\phi(x|z)p_n(x)} \right) \right] \\ &= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ -\log \left( 1 + \frac{\nu q_\phi(x|z)p_n(x)}{p_\theta(x, z)} \right) \right] \quad \text{Expand } g \\ &= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{\nu p_n(x) q_\phi(x|z)} \right] \end{aligned}$$

- Lower bound of the NCE objective

$$J_{NCE} \geq \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{\nu p_n(x) q_\phi(x|z)} \right] + \nu \mathbb{E}_{x \sim p_n} \left[ \log \frac{\nu p_n(x)}{p_\theta(x) + \nu p_n(x)} \right]$$

- What should we do with the 2nd term?

- Why does the same trick not work for the 2nd term?

# Variational NCE

- Recall

$$\begin{aligned} J_1 &\geq \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ g \left( \frac{p_\theta(x, z)}{q_\phi(x|z)p_n(x)} \right) \right] \\ &= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ -\log \left( 1 + \frac{\nu q_\phi(x|z)p_n(x)}{p_\theta(x, z)} \right) \right] \quad \text{Expand } g \\ &= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{\nu p_n(x) q_\phi(x|z)} \right] \end{aligned}$$

- Lower bound of the NCE objective

$$J_{NCE} \geq \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{\nu p_n(x) q_\phi(x|z)} \right] + \nu \mathbb{E}_{x \sim p_n} \left[ \log \frac{\nu p_n(x)}{p_\theta(x) + \nu p_n(x)} \right]$$

- What should we do with the 2nd term?

- Why does the same trick not work for the 2nd term?

# Variational NCE

- Recall

$$\begin{aligned} J_1 &\geq \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ g \left( \frac{p_\theta(x, z)}{q_\phi(x|z)p_n(x)} \right) \right] \\ &= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ -\log \left( 1 + \frac{\nu q_\phi(x|z)p_n(x)}{p_\theta(x, z)} \right) \right] \quad \text{Expand } g \\ &= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{\nu p_n(x) q_\phi(x|z)} \right] \end{aligned}$$

- Lower bound of the NCE objective

$$J_{NCE} \geq \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{\nu p_n(x) q_\phi(x|z)} \right] + \nu \mathbb{E}_{x \sim p_n} \left[ \log \frac{\nu p_n(x)}{p_\theta(x) + \nu p_n(x)} \right]$$

- What should we do with the 2nd term?

- Why does the same trick not work for the 2nd term?

# Variational NCE

- Recall lower bound

$$\mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{\nu p_n(x) q_\phi(z|x)} \right] + \nu \mathbb{E}_{x \sim p_n} \left[ \log \frac{\nu p_n(x)}{p_\theta(x) + \nu p_n(x)} \right]$$

- Paper proposes to use importance sampling for the 2nd terms

$$p_\theta(x) = \int p_\theta(x, z) dz = \int \frac{q_\phi(z|x)}{q_\phi(z|x)} p_\theta(x, z) dz = \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \frac{p_\theta(x, z)}{q_\phi(z|x)} \right]$$

- We arrive at

$$\begin{aligned} J_{VNCE}(\theta, \phi) &= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{\nu p_n(x) q_\phi(z|x)} \right] \\ &\quad + \nu \mathbb{E}_{x \sim p_n} \left[ \log \frac{\nu p_n(x)}{\mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \frac{p_\theta(x, z)}{q_\phi(z|x)} \right] + \nu p_n(x)} \right] \end{aligned}$$

# Variational NCE

- Recall lower bound

$$\mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{\nu p_n(x) q_\phi(z|x)} \right] + \nu \mathbb{E}_{x \sim p_n} \left[ \log \frac{\nu p_n(x)}{p_\theta(x) + \nu p_n(x)} \right]$$

- Paper proposes to use importance sampling for the 2nd terms

$$p_\theta(x) = \int p_\theta(x, z) dz = \int \frac{q_\phi(z|x)}{q_\phi(z|x)} p_\theta(x, z) dz = \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \frac{p_\theta(x, z)}{q_\phi(z|x)} \right]$$

- We arrive at

$$\begin{aligned} J_{VNCE}(\theta, \phi) &= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{\nu p_n(x) q_\phi(z|x)} \right] \\ &\quad + \nu \mathbb{E}_{x \sim p_n} \left[ \log \frac{\nu p_n(x)}{\mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \frac{p_\theta(x, z)}{q_\phi(z|x)} \right] + \nu p_n(x)} \right] \end{aligned}$$

# Variational NCE

- Recall lower bound

$$\mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{\nu p_n(x) q_\phi(z|x)} \right] + \nu \mathbb{E}_{x \sim p_n} \left[ \log \frac{\nu p_n(x)}{p_\theta(x) + \nu p_n(x)} \right]$$

- Paper proposes to use importance sampling for the 2nd terms

$$p_\theta(x) = \int p_\theta(x, z) dz = \int \frac{q_\phi(z|x)}{q_\phi(z|x)} p_\theta(x, z) dz = \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \frac{p_\theta(x, z)}{q_\phi(z|x)} \right]$$

- We arrive at

$$\begin{aligned} J_{VNCE}(\theta, \phi) &= \mathbb{E}_{x \sim p_d, z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x, z)}{\nu p_n(x) q_\phi(z|x)} \right] \\ &\quad + \nu \mathbb{E}_{x \sim p_n} \left[ \log \frac{\nu p_n(x)}{\mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[ \frac{p_\theta(x, z)}{q_\phi(z|x)} \right] + \nu p_n(x)} \right] \end{aligned}$$

## Variational NCE

- By construction we have

$$J_{VNCE}(\theta, \phi) \leq J_{NCE}(\theta)$$

- It can be shown that for a suitable  $f$

$$J_{VNCE}(\theta, \phi) = J_{NCE}(\theta) - D_f(q_\phi(\cdot|x) \| p_\theta(\cdot|x))$$

- $D_f(p\|q) \geq 0 \dots f\text{-divergence (generalization of } D_{KL})$

$$D_f(p\|q) := \mathbb{E}_{x \sim q} \left[ f \left( \frac{p(x)}{q(x)} \right) \right]$$

$f$  is convex with  $f(0) = 1$

$$D_{KL}(p\|q) = D_f(p\|q) \quad \text{with } f(x) = x \log x$$

### Discussion

- When is  $J_{VNCE}(\theta, \phi) = J_{NCE}(\theta)$ ?
- What about unnormalized models?

## Variational NCE

- By construction we have

$$J_{VNCE}(\theta, \phi) \leq J_{NCE}(\theta)$$

- It can be shown that for a suitable  $f$

$$J_{VNCE}(\theta, \phi) = J_{NCE}(\theta) - D_f(q_\phi(\cdot|x) \| p_\theta(\cdot|x))$$

- $D_f(p\|q) \geq 0 \dots f\text{-divergence (generalization of } D_{KL})$

$$D_f(p\|q) := \mathbb{E}_{x \sim q} \left[ f \left( \frac{p(x)}{q(x)} \right) \right]$$

$f$  is convex with  $f(0) = 1$

$$D_{KL}(p\|q) = D_f(p\|q) \quad \text{with } f(x) = x \log x$$

### Discussion

- When is  $J_{VNCE}(\theta, \phi) = J_{NCE}(\theta)$ ?
- What about unnormalized models?

1

## Latent variable models

- Variational auto-encoders
- Variational NCE
- Boltzmann machines

2

## Sparse coding and dictionary learning

3

## Exercise 2

## Boltzmann machine

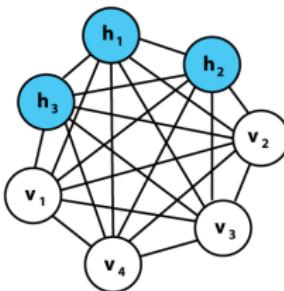
A Boltzmann machine is an energy-based model over binary variables  $x_j \in \{0, 1\}$  with

$$\log p(x) = x^T W x + b^T x + c \quad \text{diag}(W) = 0$$

- Boltzmann machine are undirected graphical models
- Strongly connected with Hopfield networks used to model associative memory
  - Clamp visible units  $x_j, j \in V$  to given values

$$p(x_{\setminus V} | x_V) = ?$$

- Intractable, since  $x_{\setminus V}$  are not independent



# Restricted Boltzmann machines

## Restricted Boltzmann machine

A restricted Boltzmann machine (RBM) is an energy-based model over two sets of binary variables  $x_j \in \{0, 1\}$  and  $z_k \in \{0, 1\}$  with

$$\log p(x, z) \doteq x^T W z + a^T x + b^T z$$

Visible units only talk to hidden units (and vice versa).

- Clamp hidden units  $z$ :

$$p(x|z) \propto \exp(x^T W z + a^T x) = \prod_j e^{x_j W_j z + a_j x_j} = \prod_j p(x_j|z)$$

$$p(x_j = 0|z) \propto e^0 = 1 \quad p(x_j = 1|z) \propto e^{W_j z + a_j}$$
$$p(x_j = 1|z) = \frac{e^{W_j z + a_j}}{1 + e^{W_j z + a_j}} = \sigma(W_j z + a_j) \quad \sigma(u) = \frac{e^u}{1 + e^u}$$

- Clamp visible units  $x$ :

$$p(z_k = 1|x) = \sigma(W^k x + b_k)$$

# Restricted Boltzmann machines

## Restricted Boltzmann machine

A restricted Boltzmann machine (RBM) is an energy-based model over two sets of binary variables  $x_j \in \{0, 1\}$  and  $z_k \in \{0, 1\}$  with

$$\log p(x, z) \doteq x^T W z + a^T x + b^T z$$

Visible units only talk to hidden units (and vice versa).

- Clamp hidden units  $z$ :

$$p(x|z) \propto \exp(x^T W z + a^T x) = \prod_j e^{x_j W_{jz} + a_j x_j} = \prod_j p(x_j|z)$$

$$p(x_j = 0|z) \propto e^0 = 1 \quad p(x_j = 1|z) \propto e^{W_{jz} + a_j}$$

$$p(x_j = 1|z) = \frac{e^{W_{jz} + a_j}}{1 + e^{W_{jz} + a_j}} = \sigma(W_{jz} + a_j) \quad \sigma(u) = \frac{e^u}{1 + e^u}$$

- Clamp visible units  $x$ :

$$p(z_k = 1|x) = \sigma(W^k x + b_k)$$

# Restricted Boltzmann machines

## Restricted Boltzmann machine

A restricted Boltzmann machine (RBM) is an energy-based model over two sets of binary variables  $x_j \in \{0, 1\}$  and  $z_k \in \{0, 1\}$  with

$$\log p(x, z) \doteq x^T W z + a^T x + b^T z$$

Visible units only talk to hidden units (and vice versa).

- Clamp hidden units  $z$ :

$$p(x|z) \propto \exp(x^T W z + a^T x) = \prod_j e^{x_j W_{jz} + a_j x_j} = \prod_j p(x_j|z)$$

$$p(x_j = 0|z) \propto e^0 = 1 \quad p(x_j = 1|z) \propto e^{W_{jz} + a_j}$$

$$p(x_j = 1|z) = \frac{e^{W_{jz} + a_j}}{1 + e^{W_{jz} + a_j}} = \sigma(W_{jz} + a_j) \quad \sigma(u) = \frac{e^u}{1 + e^u}$$

- Clamp visible units  $x$ :

$$p(z_k = 1|x) = \sigma(W^k x + b_k)$$

# Restricted Boltzmann machines

- RBMs allow to compute  $p(z|x)$  and  $p(x|z)$  easily and *exactly*
  - Visible units  $x$  are conditionally independent given hidden units  $z$
  - Hidden units  $z$  are conditionally independent given visible units  $x$
- RBMs naturally introduce layers
  - Deep Boltzmann machines are a stack of RBMs
  - Salakhutdinov & Hinton, "Deep boltzmann machines"
- Maximum likelihood training
  - Contrastive divergence and variants
  - Hinton, Training products of experts by minimizing contrastive divergence, 2002
  - We will point out possible different algorithms to train an RBM

# Restricted Boltzmann machines

We cannot compute  $p(x, z)$  (or  $p(x)$ ) exactly

- Partition function is still intractable
- But:

$$\begin{aligned} p(x) &= \frac{1}{Z} \int \exp(x^T W z + a^T x + b^T z) dz \\ &= \frac{1}{Z} \int \exp \left( \sum_k x^T W^k z_k + a^T x + b_k z_k \right) dz \\ &= \frac{\exp(a^T x)}{Z} \sum_{z_1 \in \{0,1\}, \dots, z_K \in \{0,1\}} \prod_{k=1}^K \exp(x^T W^k z_k + b_k z_k) \\ &= \frac{\exp(a^T x)}{Z} \sum_{z_1 \in \{0,1\}} \exp(x^T W^1 z_1 + b_1 z_1) \sum_{z_2 \in \{0,1\}} \exp(x^T W^2 z_2 + b_2 z_2) \cdots \\ &= \frac{\exp(a^T x)}{Z} \prod_{k=1}^K (1 + \exp(x^T W^k + b_k)) \end{aligned}$$

- Therefore

$$\log p(x) \doteq a^T x + \sum_k S(x^T W^k + b_k) \quad S(u) := \log(1 + e^u) \quad S'(u) = \sigma(u)$$

- We can use NCE to train an RBM!

- Score matching is not applicable as  $x$  is not continuous

# Restricted Boltzmann machines

We cannot compute  $p(x, z)$  (or  $p(x)$ ) exactly

- Partition function is still intractable
- But:

$$\begin{aligned} p(x) &= \frac{1}{Z} \int \exp(x^T W z + a^T x + b^T z) dz \\ &= \frac{1}{Z} \int \exp \left( \sum_k x^T W^k z_k + a^T x + b_k z_k \right) dz \\ &= \frac{\exp(a^T x)}{Z} \sum_{z_1 \in \{0,1\}, \dots, z_K \in \{0,1\}} \prod_{k=1}^K \exp(x^T W^k z_k + b_k z_k) \\ &= \frac{\exp(a^T x)}{Z} \sum_{z_1 \in \{0,1\}} \exp(x^T W^1 z_1 + b_1 z_1) \sum_{z_2 \in \{0,1\}} \exp(x^T W^2 z_2 + b_2 z_2) \cdots \\ &= \frac{\exp(a^T x)}{Z} \prod_{k=1}^K (1 + \exp(x^T W^k + b_k)) \end{aligned}$$

- Therefore

$$\log p(x) \doteq a^T x + \sum_k S(x^T W^k + b_k) \quad S(u) := \log(1 + e^u) \quad S'(u) = \sigma(u)$$

- We can use NCE to train an RBM!

- Score matching is not applicable as  $x$  is not continuous

# Restricted Boltzmann machines

We cannot compute  $p(x, z)$  (or  $p(x)$ ) exactly

- Partition function is still intractable
- But:

$$\begin{aligned} p(x) &= \frac{1}{Z} \int \exp(x^T W z + a^T x + b^T z) dz \\ &= \frac{1}{Z} \int \exp \left( \sum_k x^T W^k z_k + a^T x + b_k z_k \right) dz \\ &= \frac{\exp(a^T x)}{Z} \sum_{z_1 \in \{0,1\}, \dots, z_K \in \{0,1\}} \prod_{k=1}^K \exp(x^T W^k z_k + b_k z_k) \\ &= \frac{\exp(a^T x)}{Z} \sum_{z_1 \in \{0,1\}} \exp(x^T W^1 z_1 + b_1 z_1) \sum_{z_2 \in \{0,1\}} \exp(x^T W^2 z_2 + b_2 z_2) \cdots \\ &= \frac{\exp(a^T x)}{Z} \prod_{k=1}^K (1 + \exp(x^T W^k + b_k)) \end{aligned}$$

- Therefore

$$\log p(x) \doteq a^T x + \sum_k S(x^T W^k + b_k) \quad S(u) := \log(1 + e^u) \quad S'(u) = \sigma(u)$$

- We can use NCE to train an RBM!

- Score matching is not applicable as  $x$  is not continuous

# Restricted Boltzmann machines

We cannot compute  $p(x, z)$  (or  $p(x)$ ) exactly

- Partition function is still intractable
- But:

$$\begin{aligned} p(x) &= \frac{1}{Z} \int \exp(x^T W z + a^T x + b^T z) dz \\ &= \frac{1}{Z} \int \exp \left( \sum_k x^T W^k z_k + a^T x + b_k z_k \right) dz \\ &= \frac{\exp(a^T x)}{Z} \sum_{z_1 \in \{0,1\}, \dots, z_K \in \{0,1\}} \prod_{k=1}^K \exp(x^T W^k z_k + b_k z_k) \\ &= \frac{\exp(a^T x)}{Z} \sum_{z_1 \in \{0,1\}} \exp(x^T W^1 z_1 + b_1 z_1) \sum_{z_2 \in \{0,1\}} \exp(x^T W^2 z_2 + b_2 z_2) \cdots \\ &= \frac{\exp(a^T x)}{Z} \prod_{k=1}^K (1 + \exp(x^T W^k + b_k)) \end{aligned}$$

- Therefore

$$\log p(x) \doteq a^T x + \sum_k S(x^T W^k + b_k) \quad S(u) := \log(1 + e^u) \quad S'(u) = \sigma(u)$$

- We can use NCE to train an RBM!

- Score matching is not applicable as  $x$  is not continuous

# Restricted Boltzmann machines

We cannot compute  $p(x, z)$  (or  $p(x)$ ) exactly

- Partition function is still intractable
- But:

$$\begin{aligned} p(x) &= \frac{1}{Z} \int \exp(x^T W z + a^T x + b^T z) dz \\ &= \frac{1}{Z} \int \exp \left( \sum_k x^T W^k z_k + a^T x + b_k z_k \right) dz \\ &= \frac{\exp(a^T x)}{Z} \sum_{z_1 \in \{0,1\}, \dots, z_K \in \{0,1\}} \prod_{k=1}^K \exp(x^T W^k z_k + b_k z_k) \\ &= \frac{\exp(a^T x)}{Z} \sum_{z_1 \in \{0,1\}} \exp(x^T W^1 z_1 + b_1 z_1) \sum_{z_2 \in \{0,1\}} \exp(x^T W^2 z_2 + b_2 z_2) \cdots \\ &= \frac{\exp(a^T x)}{Z} \prod_{k=1}^K (1 + \exp(x^T W^k + b_k)) \end{aligned}$$

- Therefore

$$\log p(x) \doteq a^T x + \sum_k S(x^T W^k + b_k) \quad S(u) := \log(1 + e^u) \quad S'(u) = \sigma(u)$$

- We can use NCE to train an RBM!

- Score matching is not applicable as  $x$  is not continuous

# Restricted Boltzmann machines

We cannot compute  $p(x, z)$  (or  $p(x)$ ) exactly

- Partition function is still intractable
- But:

$$\begin{aligned} p(x) &= \frac{1}{Z} \int \exp(x^T W z + a^T x + b^T z) dz \\ &= \frac{1}{Z} \int \exp \left( \sum_k x^T W^k z_k + a^T x + b_k z_k \right) dz \\ &= \frac{\exp(a^T x)}{Z} \sum_{z_1 \in \{0,1\}, \dots, z_K \in \{0,1\}} \prod_{k=1}^K \exp(x^T W^k z_k + b_k z_k) \\ &= \frac{\exp(a^T x)}{Z} \sum_{z_1 \in \{0,1\}} \exp(x^T W^1 z_1 + b_1 z_1) \sum_{z_2 \in \{0,1\}} \exp(x^T W^2 z_2 + b_2 z_2) \cdots \\ &= \frac{\exp(a^T x)}{Z} \prod_{k=1}^K (1 + \exp(x^T W^k + b_k)) \end{aligned}$$

- Therefore

$$\log p(x) \doteq a^T x + \sum_k S(x^T W^k + b_k) \quad S(u) := \log(1 + e^u) \quad S'(u) = \sigma(u)$$

- We can use NCE to train an RBM!

- Score matching is not applicable as  $x$  is not continuous

# Contrastive divergence

- RBMs (and other models) trained with contrastive divergence
- Gradient of MLE for unnormalized latent variable models

$$p_\theta(x) = \frac{\int e^{f_\theta(x,z)} dz}{Z(\theta)}$$
$$Z(\theta) = \int e^{f_\theta(x,z)} dz dx$$

- What is

$$\nabla_\theta \log p_\theta(x) = \nabla_\theta \log \frac{\int e^{f_\theta(x,z)} dz}{Z(\theta)}$$
$$= \nabla_\theta \log \int e^{f_\theta(x,z)} dz - \nabla_\theta \log Z(\theta)?$$

# Contrastive divergence

- Let us start with the 2nd term (recall  $p_\theta(x, z) = e^{f_\theta(x, z)} / Z(\theta)$ )

$$\begin{aligned}\nabla_\theta \log Z(\theta) &= \frac{\nabla_\theta Z(\theta)}{Z(\theta)} = \frac{\nabla_\theta \int e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int \nabla_\theta e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int \frac{e^{f_\theta(x, z)}}{e^{f_\theta(x, z)}} \nabla_\theta e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int e^{f_\theta(x, z)} \left( \nabla_\theta \log e^{f_\theta(x, z)} \right) dz dx}{Z(\theta)} \\ &= \int p_\theta(x, z) \nabla_\theta f_\theta(x, z) dz dx \\ &= \mathbb{E}_{(x, z) \sim p_\theta(\cdot, \cdot)} [\nabla_\theta f_\theta(x, z)]\end{aligned}$$

# Contrastive divergence

- Let us start with the 2nd term (recall  $p_\theta(x, z) = e^{f_\theta(x, z)} / Z(\theta)$ )

$$\begin{aligned}\nabla_\theta \log Z(\theta) &= \frac{\nabla_\theta Z(\theta)}{Z(\theta)} = \frac{\nabla_\theta \int e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int \nabla_\theta e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int \frac{e^{f_\theta(x, z)}}{e^{f_\theta(x, z)}} \nabla_\theta e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int e^{f_\theta(x, z)} \left( \nabla_\theta \log e^{f_\theta(x, z)} \right) dz dx}{Z(\theta)} \\ &= \int p_\theta(x, z) \nabla_\theta f_\theta(x, z) dz dx \\ &= \mathbb{E}_{(x, z) \sim p_\theta(\cdot, \cdot)} [\nabla_\theta f_\theta(x, z)]\end{aligned}$$

# Contrastive divergence

- Let us start with the 2nd term (recall  $p_\theta(x, z) = e^{f_\theta(x, z)} / Z(\theta)$ )

$$\begin{aligned}\nabla_\theta \log Z(\theta) &= \frac{\nabla_\theta Z(\theta)}{Z(\theta)} = \frac{\nabla_\theta \int e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int \nabla_\theta e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int \frac{e^{f_\theta(x, z)}}{e^{f_\theta(x, z)}} \nabla_\theta e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int e^{f_\theta(x, z)} \left( \nabla_\theta \log e^{f_\theta(x, z)} \right) dz dx}{Z(\theta)} \\ &= \int p_\theta(x, z) \nabla_\theta f_\theta(x, z) dz dx \\ &= \mathbb{E}_{(x, z) \sim p_\theta(\cdot, \cdot)} [\nabla_\theta f_\theta(x, z)]\end{aligned}$$

# Contrastive divergence

- Let us start with the 2nd term (recall  $p_\theta(x, z) = e^{f_\theta(x, z)} / Z(\theta)$ )

$$\begin{aligned}\nabla_\theta \log Z(\theta) &= \frac{\nabla_\theta Z(\theta)}{Z(\theta)} = \frac{\nabla_\theta \int e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int \nabla_\theta e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int \frac{e^{f_\theta(x, z)}}{e^{f_\theta(x, z)}} \nabla_\theta e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int e^{f_\theta(x, z)} \left( \nabla_\theta \log e^{f_\theta(x, z)} \right) dz dx}{Z(\theta)} \\ &= \int p_\theta(x, z) \nabla_\theta f_\theta(x, z) dz dx \\ &= \mathbb{E}_{(x, z) \sim p_\theta(\cdot, \cdot)} [\nabla_\theta f_\theta(x, z)]\end{aligned}$$

# Contrastive divergence

- Let us start with the 2nd term (recall  $p_\theta(x, z) = e^{f_\theta(x, z)} / Z(\theta)$ )

$$\begin{aligned}\nabla_\theta \log Z(\theta) &= \frac{\nabla_\theta Z(\theta)}{Z(\theta)} = \frac{\nabla_\theta \int e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int \nabla_\theta e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int \frac{e^{f_\theta(x, z)}}{e^{f_\theta(x, z)}} \nabla_\theta e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int e^{f_\theta(x, z)} \left( \nabla_\theta \log e^{f_\theta(x, z)} \right) dz dx}{Z(\theta)} \\ &= \int p_\theta(x, z) \nabla_\theta f_\theta(x, z) dz dx \\ &= \mathbb{E}_{(x, z) \sim p_\theta(\cdot, \cdot)} [\nabla_\theta f_\theta(x, z)]\end{aligned}$$

# Contrastive divergence

- Let us start with the 2nd term (recall  $p_\theta(x, z) = e^{f_\theta(x, z)} / Z(\theta)$ )

$$\begin{aligned}\nabla_\theta \log Z(\theta) &= \frac{\nabla_\theta Z(\theta)}{Z(\theta)} = \frac{\nabla_\theta \int e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int \nabla_\theta e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int \frac{e^{f_\theta(x, z)}}{e^{f_\theta(x, z)}} \nabla_\theta e^{f_\theta(x, z)} dz dx}{Z(\theta)} \\ &= \frac{\int e^{f_\theta(x, z)} \left( \nabla_\theta \log e^{f_\theta(x, z)} \right) dz dx}{Z(\theta)} \\ &= \int p_\theta(x, z) \nabla_\theta f_\theta(x, z) dz dx \\ &= \mathbb{E}_{(x, z) \sim p_\theta(\cdot, \cdot)} [\nabla_\theta f_\theta(x, z)]\end{aligned}$$

# Contrastive divergence

- Now for the 1st term

$$\begin{aligned}\nabla_{\theta} \log \int e^{f_{\theta}(x, z)} dz &= \frac{\nabla_{\theta} \int e^{f_{\theta}(x, z)} dz}{\int e^{f_{\theta}(x, z)} dz} \\&= \frac{Z(\theta) \nabla_{\theta} \int e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{Z(\theta) \int \nabla_{\theta} e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{Z(\theta) \int \frac{e^{f_{\theta}(x, z)}}{e^{f_{\theta}(x, z)}} \nabla_{\theta} e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{\int p_{\theta}(x, z) \left( \nabla_{\theta} \log e^{f_{\theta}(x, z)} \right) dz}{p_{\theta}(x)} \\&= \int p_{\theta}(z|x) (\nabla_{\theta} f_{\theta}(x, z)) dz \\&= \mathbb{E}_{z \sim p_{\theta}(\cdot|x)} [\nabla_{\theta} f_{\theta}(x, z)]\end{aligned}$$

# Contrastive divergence

- Now for the 1st term

$$\begin{aligned}\nabla_{\theta} \log \int e^{f_{\theta}(x, z)} dz &= \frac{\nabla_{\theta} \int e^{f_{\theta}(x, z)} dz}{\int e^{f_{\theta}(x, z)} dz} \\&= \frac{Z(\theta) \nabla_{\theta} \int e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{Z(\theta) \int \nabla_{\theta} e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{Z(\theta) \int \frac{e^{f_{\theta}(x, z)}}{e^{f_{\theta}(x, z)}} \nabla_{\theta} e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{\int p_{\theta}(x, z) \left( \nabla_{\theta} \log e^{f_{\theta}(x, z)} \right) dz}{p_{\theta}(x)} \\&= \int p_{\theta}(z|x) (\nabla_{\theta} f_{\theta}(x, z)) dz \\&= \mathbb{E}_{z \sim p_{\theta}(\cdot|x)} [\nabla_{\theta} f_{\theta}(x, z)]\end{aligned}$$

# Contrastive divergence

- Now for the 1st term

$$\begin{aligned}\nabla_{\theta} \log \int e^{f_{\theta}(x, z)} dz &= \frac{\nabla_{\theta} \int e^{f_{\theta}(x, z)} dz}{\int e^{f_{\theta}(x, z)} dz} \\&= \frac{Z(\theta) \nabla_{\theta} \int e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{Z(\theta) \int \nabla_{\theta} e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{Z(\theta) \int \frac{\partial f_{\theta}(x, z)}{\partial f_{\theta}(x, z)} \nabla_{\theta} e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{\int p_{\theta}(x, z) \left( \nabla_{\theta} \log e^{f_{\theta}(x, z)} \right) dz}{p_{\theta}(x)} \\&= \int p_{\theta}(z|x) (\nabla_{\theta} f_{\theta}(x, z)) dz \\&= \mathbb{E}_{z \sim p_{\theta}(\cdot|x)} [\nabla_{\theta} f_{\theta}(x, z)]\end{aligned}$$

# Contrastive divergence

- Now for the 1st term

$$\begin{aligned}\nabla_{\theta} \log \int e^{f_{\theta}(x, z)} dz &= \frac{\nabla_{\theta} \int e^{f_{\theta}(x, z)} dz}{\int e^{f_{\theta}(x, z)} dz} \\&= \frac{Z(\theta) \nabla_{\theta} \int e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{Z(\theta) \int \nabla_{\theta} e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{Z(\theta) \int \frac{e^{f_{\theta}(x, z)}}{e^{f_{\theta}(x, z)}} \nabla_{\theta} e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{\int p_{\theta}(x, z) \left( \nabla_{\theta} \log e^{f_{\theta}(x, z)} \right) dz}{p_{\theta}(x)} \\&= \int p_{\theta}(z|x) (\nabla_{\theta} f_{\theta}(x, z)) dz \\&= \mathbb{E}_{z \sim p_{\theta}(\cdot|x)} [\nabla_{\theta} f_{\theta}(x, z)]\end{aligned}$$

# Contrastive divergence

- Now for the 1st term

$$\begin{aligned}\nabla_{\theta} \log \int e^{f_{\theta}(x, z)} dz &= \frac{\nabla_{\theta} \int e^{f_{\theta}(x, z)} dz}{\int e^{f_{\theta}(x, z)} dz} \\&= \frac{Z(\theta) \nabla_{\theta} \int e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{Z(\theta) \int \nabla_{\theta} e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{Z(\theta) \int \frac{e^{f_{\theta}(x, z)}}{e^{f_{\theta}(x, z)}} \nabla_{\theta} e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{\int p_{\theta}(x, z) \left( \nabla_{\theta} \log e^{f_{\theta}(x, z)} \right) dz}{p_{\theta}(x)} \\&= \int p_{\theta}(z|x) (\nabla_{\theta} f_{\theta}(x, z)) dz \\&= \mathbb{E}_{z \sim p_{\theta}(\cdot|x)} [\nabla_{\theta} f_{\theta}(x, z)]\end{aligned}$$

# Contrastive divergence

- Now for the 1st term

$$\begin{aligned}\nabla_{\theta} \log \int e^{f_{\theta}(x, z)} dz &= \frac{\nabla_{\theta} \int e^{f_{\theta}(x, z)} dz}{\int e^{f_{\theta}(x, z)} dz} \\&= \frac{Z(\theta) \nabla_{\theta} \int e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{Z(\theta) \int \nabla_{\theta} e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{Z(\theta) \int \frac{e^{f_{\theta}(x, z)}}{e^{f_{\theta}(x, z)}} \nabla_{\theta} e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{\int p_{\theta}(x, z) \left( \nabla_{\theta} \log e^{f_{\theta}(x, z)} \right) dz}{p_{\theta}(x)} \\&= \int p_{\theta}(z|x) (\nabla_{\theta} f_{\theta}(x, z)) dz \\&= \mathbb{E}_{z \sim p_{\theta}(\cdot|x)} [\nabla_{\theta} f_{\theta}(x, z)]\end{aligned}$$

# Contrastive divergence

- Now for the 1st term

$$\begin{aligned}\nabla_{\theta} \log \int e^{f_{\theta}(x, z)} dz &= \frac{\nabla_{\theta} \int e^{f_{\theta}(x, z)} dz}{\int e^{f_{\theta}(x, z)} dz} \\&= \frac{Z(\theta) \nabla_{\theta} \int e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{Z(\theta) \int \nabla_{\theta} e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{Z(\theta) \int \frac{e^{f_{\theta}(x, z)}}{e^{f_{\theta}(x, z)}} \nabla_{\theta} e^{f_{\theta}(x, z)} dz}{p_{\theta}(x)} \\&= \frac{\int p_{\theta}(x, z) \left( \nabla_{\theta} \log e^{f_{\theta}(x, z)} \right) dz}{p_{\theta}(x)} \\&= \int p_{\theta}(z|x) (\nabla_{\theta} f_{\theta}(x, z)) dz \\&= \mathbb{E}_{z \sim p_{\theta}(\cdot|x)} [\nabla_{\theta} f_{\theta}(x, z)]\end{aligned}$$

# Contrastive divergence

- Gradient of MLE for unnormalized latent variable models

$$p_\theta(x) = \frac{\int e^{f_\theta(x,z)} dz}{Z(\theta)}$$
$$Z(\theta) = \int e^{f_\theta(x,z)} dz dx$$

- What is

$$\nabla_\theta \log p_\theta(x) = \nabla_\theta \log \int e^{f_\theta(x,z)} dz - \nabla_\theta \log Z(\theta)?$$

- Answer

$$\nabla_\theta \log p_\theta(x) = \mathbb{E}_{z \sim p_\theta(\cdot|x)} [\nabla_\theta f_\theta(x, z)] - \mathbb{E}_{(x', z) \sim p_\theta} [\nabla_\theta f_\theta(x', z)]$$

- Expectation over the posterior minus expectation over the model

# Contrastive divergence

Back to RBM

- Gradient of MLE for unnormalized latent variable models

$$\nabla_{\theta} \log p_{\theta}(x) = \mathbb{E}_{z \sim p_{\theta}(\cdot|x)} [\nabla_{\theta} f_{\theta}(x, z)] - \mathbb{E}_{(x', z) \sim p_{\theta}} [\nabla_{\theta} f_{\theta}(x', z)]$$

- $\mathbb{E}_{z \sim p_{\theta}(\cdot|x)} [\nabla_{\theta} f_{\theta}(x, z)]$  can be estimated from samples

$$p_{\theta}(z|x) = \prod_k Ber(z_k; \sigma(W^k x + b_k))$$

- $\mathbb{E}_{(x', z) \sim p_{\theta}} [\nabla_{\theta} f_{\theta}(x', z)]$  much harder

- We need to draw samples from the joint
- Solution: Gibbs sampling

$$\begin{aligned} z \sim p_{\theta}(\cdot|x) &\rightarrow x' \sim p_{\theta}(\cdot|z) \rightarrow z' \sim p_{\theta}(\cdot|x') \rightarrow x'' \sim p_{\theta}(\cdot|z') \rightarrow \dots \\ &\rightarrow x^{\infty} \sim p_{\theta}(\cdot|z^{\infty}) \end{aligned}$$

- Contrastive divergence stops at  $(x', z)$  to produce a sample from  $p_{\theta}$

# Contrastive divergence

Back to RBM

- Gradient of MLE for unnormalized latent variable models

$$\nabla_{\theta} \log p_{\theta}(x) = \mathbb{E}_{z \sim p_{\theta}(\cdot|x)} [\nabla_{\theta} f_{\theta}(x, z)] - \mathbb{E}_{(x', z) \sim p_{\theta}} [\nabla_{\theta} f_{\theta}(x', z)]$$

- $\mathbb{E}_{z \sim p_{\theta}(\cdot|x)} [\nabla_{\theta} f_{\theta}(x, z)]$  can be estimated from samples

$$p_{\theta}(z|x) = \prod_k Ber(z_k; \sigma(W^k x + b_k))$$

- $\mathbb{E}_{(x', z) \sim p_{\theta}} [\nabla_{\theta} f_{\theta}(x', z)]$  much harder

- We need to draw samples from the joint
- Solution: Gibbs sampling

$$\begin{aligned} z \sim p_{\theta}(\cdot|x) &\rightarrow x' \sim p_{\theta}(\cdot|z) \rightarrow z' \sim p_{\theta}(\cdot|x') \rightarrow x'' \sim p_{\theta}(\cdot|z') \rightarrow \dots \\ &\rightarrow x^{\infty} \sim p_{\theta}(\cdot|z^{\infty}) \end{aligned}$$

- Contrastive divergence stops at  $(x', z)$  to produce a sample from  $p_{\theta}$

# Contrastive divergence

Back to RBM

- Gradient of MLE for unnormalized latent variable models

$$\nabla_{\theta} \log p_{\theta}(x) = \mathbb{E}_{z \sim p_{\theta}(\cdot|x)} [\nabla_{\theta} f_{\theta}(x, z)] - \mathbb{E}_{(x', z) \sim p_{\theta}} [\nabla_{\theta} f_{\theta}(x', z)]$$

- $\mathbb{E}_{z \sim p_{\theta}(\cdot|x)} [\nabla_{\theta} f_{\theta}(x, z)]$  can be estimated from samples

$$p_{\theta}(z|x) = \prod_k Ber(z_k; \sigma(W^k x + b_k))$$

- $\mathbb{E}_{(x', z) \sim p_{\theta}} [\nabla_{\theta} f_{\theta}(x', z)]$  much harder

- We need to draw samples from the joint
- Solution: Gibbs sampling

$$\begin{aligned} z \sim p_{\theta}(\cdot|x) &\rightarrow x' \sim p_{\theta}(\cdot|z) \rightarrow z' \sim p_{\theta}(\cdot|x') \rightarrow x'' \sim p_{\theta}(\cdot|z') \rightarrow \dots \\ &\rightarrow x^{\infty} \sim p_{\theta}(\cdot|z^{\infty}) \end{aligned}$$

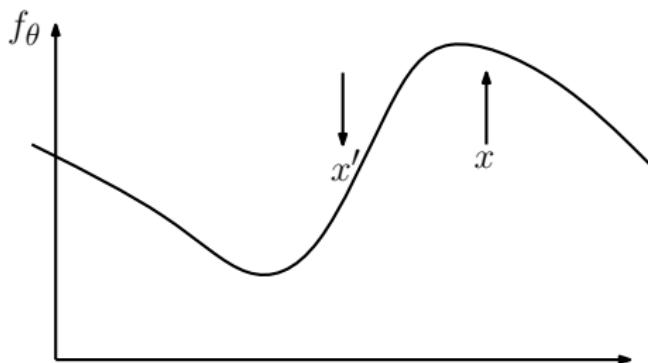
- Contrastive divergence stops at  $(x', z)$  to produce a sample from  $p_{\theta}$

# Contrastive divergence

## Learning using energy-based models

Contrastive divergence learning increases (unnormalized) log-likelihood  $f_\theta$  for given datapoints  $x$  and decreases it for nearby samples  $x'$ .

- Which other training methods share this feature?



# Gaussian-Bernoulli RBMs (GRBMs)

Real-valued visible units?

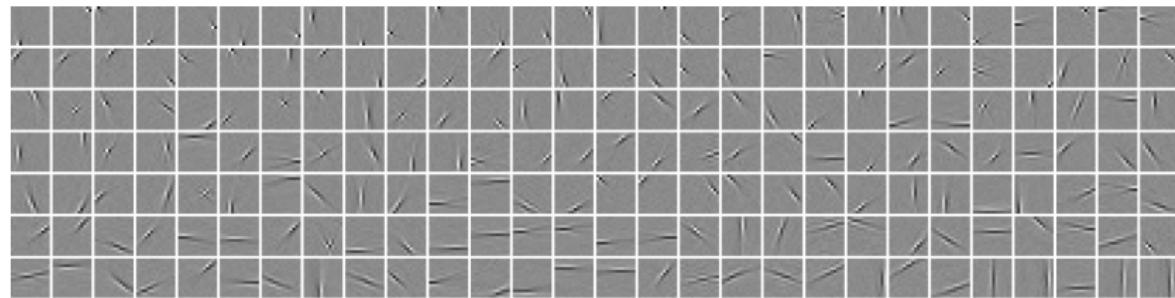
- Gaussian-Bernoulli RBM:  $x_j \in \mathbb{R}$ ,  $z_k \in \{0, 1\}$

$$\log p(x, z) \doteq -\frac{1}{2\sigma^2} \|x - a\|^2 + \frac{1}{\sigma^2} x^T W z + b^T z$$

- Hinton & Salakhutdinov, "Reducing the Dimensionality of Data with Neural Networks"

## Questions

- What is  $p(x|z)$ ?
- What is  $p(x)$ ?



Melchior et al., "Gaussian-binary restricted Boltzmann machines for modeling natural image statistics"

# Gaussian-Bernoulli RBMs (GRBMs)

Real-valued visible units?

- Gaussian-Bernoulli RBM:  $x_j \in \mathbb{R}$ ,  $z_k \in \{0, 1\}$

$$\log p(x, z) \doteq -\frac{1}{2\sigma^2} \|x - a\|^2 + \frac{1}{\sigma^2} x^T W z + b^T z$$

- Hinton & Salakhutdinov, "Reducing the Dimensionality of Data with Neural Networks"

## Questions

- What is  $p(x|z)$ ?

$$p(x|z) = \mathcal{N}(x; a + Wz, \sigma^2 I)$$

- What is  $p(x)$ ?

$$\log p(x) \doteq -\frac{\|x - a\|^2}{2\sigma^2} + \sum_k S(x^T W^k / \sigma^2 + b_k)$$

Should look familiar

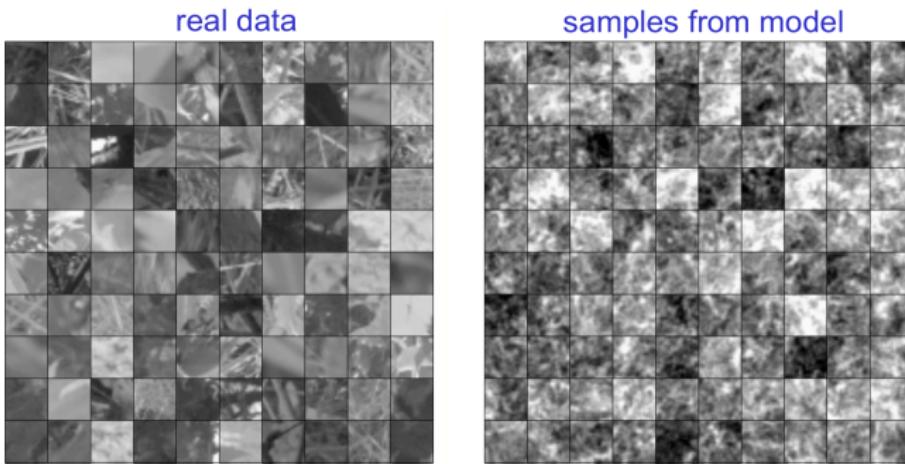
- We can use score matching now as well!

# Semi-restricted Boltzmann machines

- Semi-restricted Boltzmann machines allow coupling between visible units via lateral connections

$$\log p(x, z) \doteq -\frac{1}{2}(x - a)^T \Lambda(x - a) + b^T z + x^T \Lambda W z$$
$$p(x|z) = \mathcal{N}(x; a + Wz, \Lambda)$$

- Osindero & Hinton, "Modeling image patches with a directed hierarchy of Markov random fields"

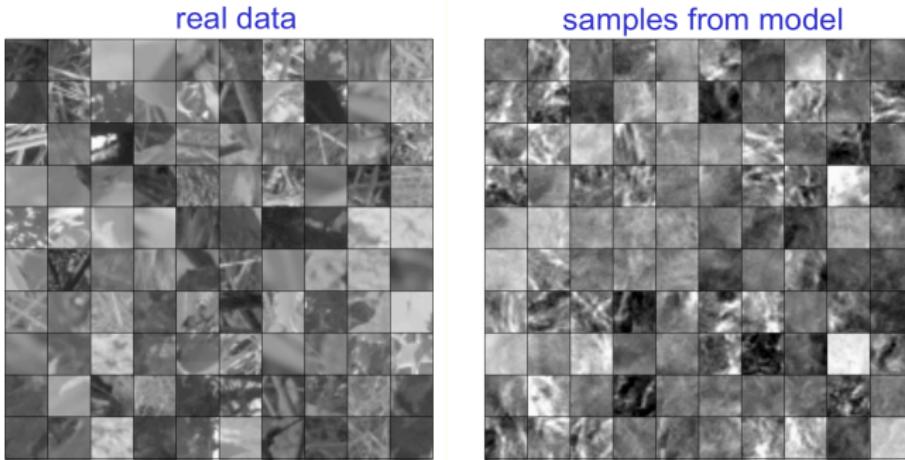


# Semi-restricted Boltzmann machines

- Semi-restricted Boltzmann machines allow coupling between visible units via lateral connections

$$\log p(x, z) \doteq -\frac{1}{2}(x - a)^T \Lambda(x - a) + b^T z + x^T \Lambda W z$$
$$p(x|z) = \mathcal{N}(x; a + Wz, \Lambda)$$

- Osindero & Hinton, "Modeling image patches with a directed hierarchy of Markov random fields"



## Mean-covariance RBM (mcRBM)

- Idea: two sets of hidden units
  - One set modulates the mean as in the GRBM
  - A 2nd set of hiddens modulates the covariance matrix
- Energy model  $\log p(x, y, z) = \log p_M(x, z) + \log p_C(x, y)$

$$\log p_M(x, z) = \frac{1}{2} \|x - a\|^2 + x^T \mathbf{W}z + b^T x \quad \text{GRBM}$$

$$\log p_C(x, y) = (x^T \mathbf{R})^2 \mathbf{P}y + c^T y \quad \text{cRBM}$$

- It can be shown that

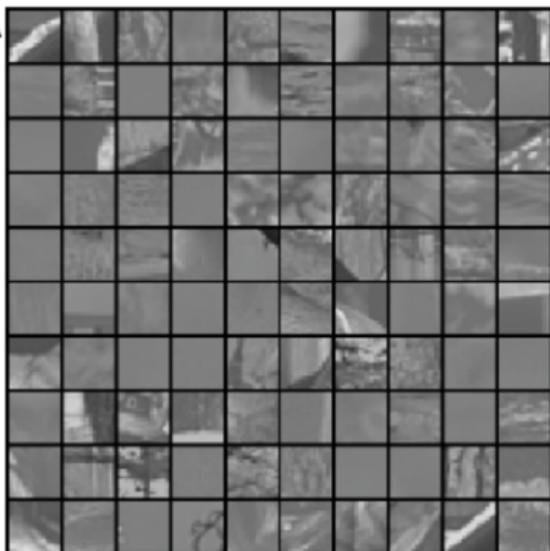
$$p(x|y, z) = \mathcal{N}(x; \Sigma \mathbf{W}z, \Sigma) \quad \Sigma := (\mathbf{R} \operatorname{diag}(-\mathbf{P}^T y) \mathbf{R}^T)^{-1}$$

- Training using contrastive divergence
- Ranzato & Hinton, "Modeling Pixel Means and Covariances Using Factorized Third-Order Boltzmann Machines"

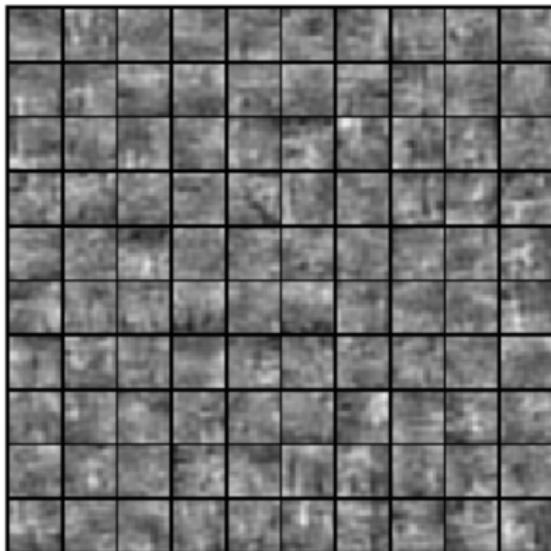
## Mean-covariance RBM (mcRBM)

- Training on natural image patches

A



C



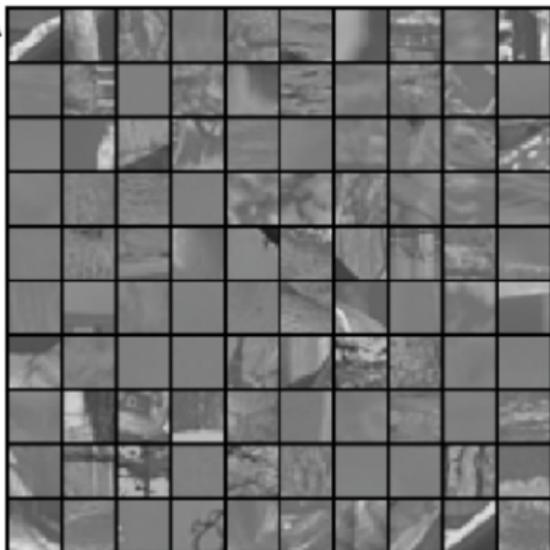
GRBM

- Ranzato & Hinton, "Modeling Pixel Means and Covariances Using Factorized Third-Order Boltzmann Machines"

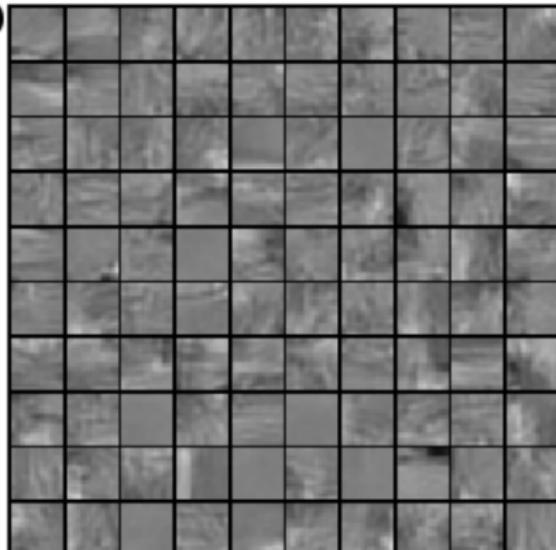
## Mean-covariance RBM (mcRBM)

- Training on natural image patches

A



D



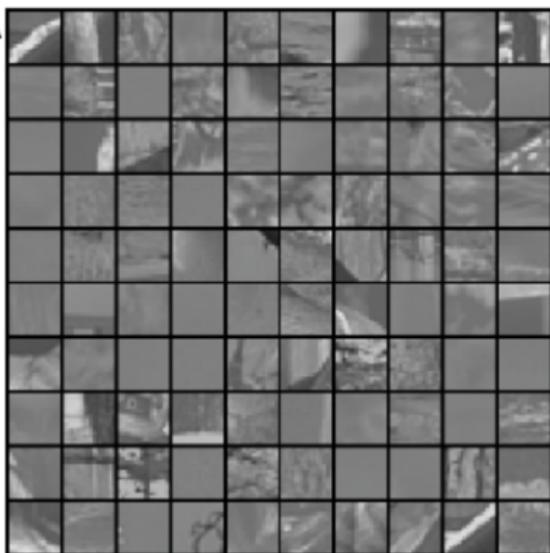
cRBM

- Ranzato & Hinton, "Modeling Pixel Means and Covariances Using Factorized Third-Order Boltzmann Machines"

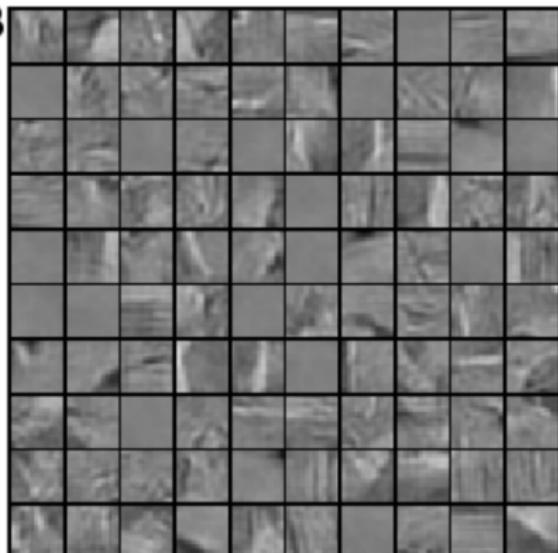
## Mean-covariance RBM (mcRBM)

- Training on natural image patches

A



B



mcRBM

- Ranzato & Hinton, "Modeling Pixel Means and Covariances Using Factorized Third-Order Boltzmann Machines"

# Higher-order Boltzmann machines

- mcRBM is 3rd order Boltzmann machine
  - 3-way multiplicative interaction  $z_k x_j x_{j'}$
- What about 3-way interaction  $z_k x_j y_n$ ? Learning image relations

$$\log p(x, y, z) \doteq \sum_k z_k x^T W_k y + b^T x + c^T y + d^T z - \frac{1}{2} \|x\|^2 - \frac{1}{2} \|y\|^2$$

$$p(x|y, z) = \mathcal{N}\left(x; b + \left(\sum_k z_k W_k\right) y, I\right)$$

$$p(y|x, z) = \mathcal{N}\left(y; c + \left(\sum_k z_k W_k^T\right) x, I\right)$$

$$p(z|x, y) = \prod_k \text{Ber}\left(z_k; \sigma(x^T W_k y + d_k)\right)$$

- Hidden units  $z$  determine the linear transformation
- Memisevic & Hinton, "Learning to represent spatial transformations with factored higher-order Boltzmann machines"

# Higher-order Boltzmann machines

- mcRBM is 3rd order Boltzmann machine
  - 3-way multiplicative interaction  $z_k x_j x_{j'}$
- What about 3-way interaction  $z_k x_j y_n$ ? Learning image relations

$$\log p(x, y, z) \doteq \sum_k z_k x^T W_k y + b^T x + c^T y + d^T z - \frac{1}{2} \|x\|^2 - \frac{1}{2} \|y\|^2$$

$$p(x|y, z) = \mathcal{N}\left(x; b + \left(\sum_k z_k W_k\right) y, I\right)$$

$$p(y|x, z) = \mathcal{N}\left(y; c + \left(\sum_k z_k W_k^T\right) x, I\right)$$

$$p(z|x, y) = \prod_k \text{Ber}\left(z_k; \sigma(x^T W_k y + d_k)\right)$$

- Hidden units  $z$  determine the linear transformation
- Memisevic & Hinton, "Learning to represent spatial transformations with factored higher-order Boltzmann machines"

# Higher-order Boltzmann machines

- Learning image relations

$$\log p(x, y, z) \doteq \sum_k z_k x^T W_k y + b^T x + c^T y + d^T z - \frac{1}{2} \|x\|^2 - \frac{1}{2} \|y\|^2$$

- Efficiency?

- $\mathcal{W} := (W_k)_{k=1}^N$  forms a 3-tensor
- Many parameters  $D \times D \times K$
- Replace by rank- $F$  approximation

$$\mathcal{W} = \sum_{f=1}^F w_f^x \otimes w_f^y \otimes w_f^z$$

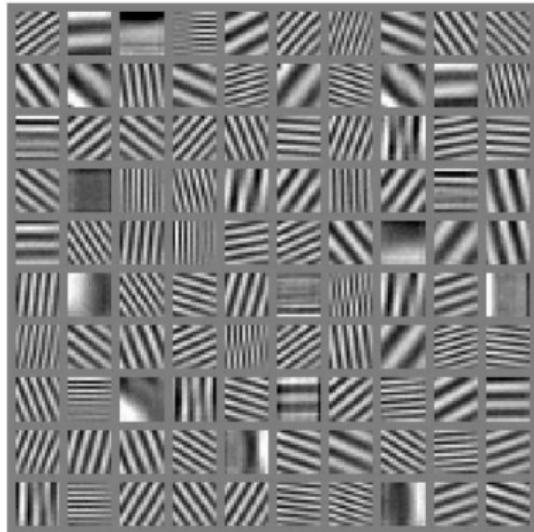
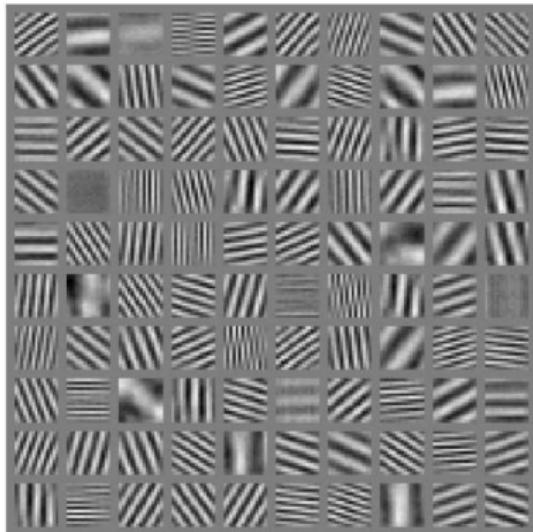
- Factored higher-order RBM

$$\log p(x, y, z) \doteq \sum_{k,f} w_{fk}^z z_k (x^T w_f^x) (y^T w_f^y) + b^T x + c^T y + d^T z - \frac{1}{2} \|x\|^2 - \frac{1}{2} \|y\|^2$$

- Memisevic & Hinton, "Learning to represent spatial transformations with factored higher-order Boltzmann machines"

# Higher-order Boltzmann machines

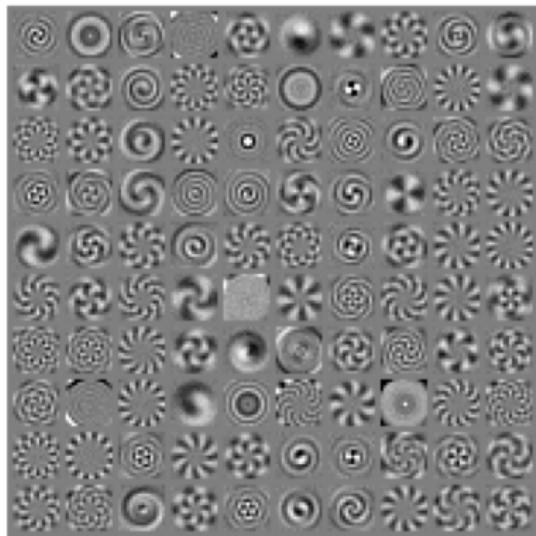
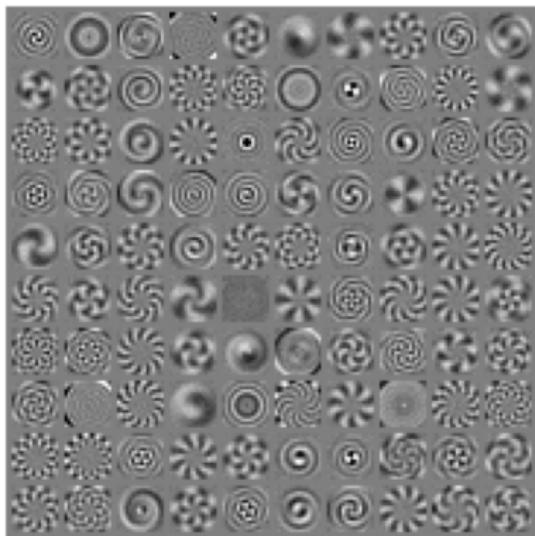
- Learning: training set of corresponding image pairs  $\{(x_i, y_i)\}$
- Filters  $w_f^x$  and  $w_f^y$  learned from data
  - Filters independent of image content
  - translated pairs



- Memisevic & Hinton, "Learning to represent spatial transformations with factored higher-order Boltzmann machines"

# Higher-order Boltzmann machines

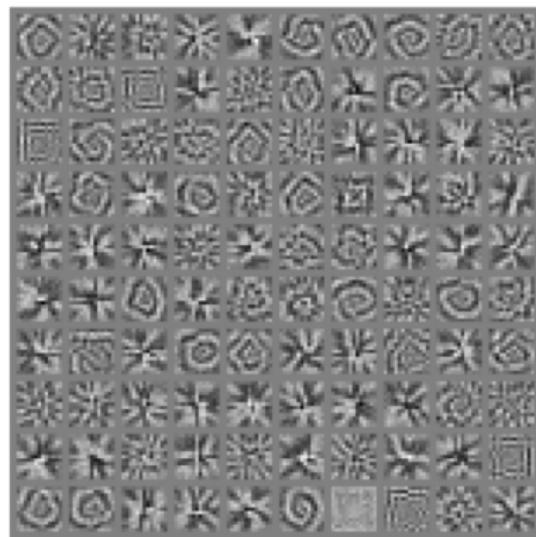
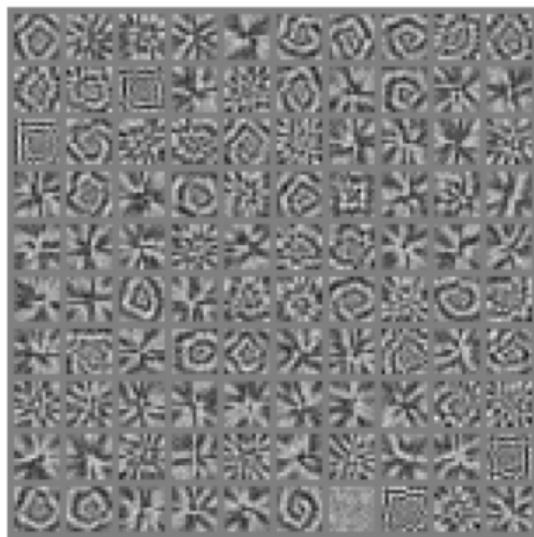
- Learning: training set of corresponding image pairs  $\{(x_i, y_i)\}$
- Filters  $w_f^x$  and  $w_f^y$  learned from data
  - Filters independent of image content
  - rotated pairs



- Memisevic & Hinton, "Learning to represent spatial transformations with factored higher-order Boltzmann machines"

# Higher-order Boltzmann machines

- Learning: training set of corresponding image pairs  $\{(x_i, y_i)\}$
- Filters  $w_f^x$  and  $w_f^y$  learned from data
  - Filters independent of image content
  - scaled pairs



- Memisevic & Hinton, "Learning to represent spatial transformations with factored higher-order Boltzmann machines"

# Higher-order Boltzmann machines

- Theory predicts how  $w_f^x$  and  $w_f^y$  should look
- Assume transformations (parametrized by  $\kappa$ ) are orthonormal

$$x = A_\kappa y \qquad \qquad y = A_\kappa^T x$$

- Permutations, rotations, translations (with wrap-around boundary conditions)
- Eigen-decomposition: complex eigenvalues with magnitude 1

$$A_\kappa = U_\kappa \Lambda_\kappa U_\kappa^T$$

- Assume transformations commute:  $A_\kappa A_{\kappa'} = A_{\kappa'} A_\kappa$ 
  - Commuting matrices share eigenspaces!

$$U_\kappa = U_{\kappa'} = U$$

- We obtain

$$x^T A_\theta y = (\underbrace{U^T x}_{x^T w^x})^T \underbrace{\Lambda_\theta}_{\text{diag}(w^z z_k)} \underbrace{U^T y}_{y^T w^y}$$

- $z$  parametrizes transformation  $\kappa$
- Memisevic, "Learning related images"

# Higher-order Boltzmann machines

- Theory predicts how  $w_f^x$  and  $w_f^y$  should look
- Assume transformations (parametrized by  $\kappa$ ) are orthonormal

$$x = A_\kappa y \qquad \qquad y = A_\kappa^T x$$

- Permutations, rotations, translations (with wrap-around boundary conditions)
- Eigen-decomposition: complex eigenvalues with magnitude 1

$$A_\kappa = U_\kappa \Lambda_\kappa U_\kappa^T$$

- Assume transformations commute:  $A_\kappa A_{\kappa'} = A_{\kappa'} A_\kappa$ 
  - Commuting matrices share eigenspaces!

$$U_\kappa = U_{\kappa'} = U$$

- We obtain

$$x^T A_\theta y = (\underbrace{U^T x}_{x^T w^x})^T \underbrace{\Lambda_\theta}_{\text{diag}(w^z z_k)} \underbrace{U^T y}_{y^T w^y}$$

- $z$  parametrizes transformation  $\kappa$
- Memisevic, "Learning related images"

# Higher-order Boltzmann machines

- Theory predicts how  $w_f^x$  and  $w_f^y$  should look
- Assume transformations (parametrized by  $\kappa$ ) are orthonormal

$$x = A_\kappa y \qquad \qquad y = A_\kappa^T x$$

- Permutations, rotations, translations (with wrap-around boundary conditions)
- Eigen-decomposition: complex eigenvalues with magnitude 1

$$A_\kappa = U_\kappa \Lambda_\kappa U_\kappa^T$$

- Assume transformations commute:  $A_\kappa A_{\kappa'} = A_{\kappa'} A_\kappa$ 
  - Commuting matrices share eigenspaces!

$$U_\kappa = U_{\kappa'} = U$$

- We obtain

$$x^T A_\theta y = (\underbrace{U^T x}_{x^T w^x})^T \underbrace{\Lambda_\theta}_{\text{diag}(w^z z_k)} \underbrace{U^T y}_{y^T w^y}$$

- $z$  parametrizes transformation  $\kappa$
- Memisevic, "Learning related images"

# Training RBMs and whitening of data

- Pixels in images are highly correlated
- Whitening aims to remove 2nd order correlation from images
  - Learned model can focus on higher-order correlations
- PCA whitening
  - Subtract empirical mean

$$\mathbf{X} \leftarrow \mathbf{X} - \frac{1}{N} \mathbf{X}\mathbf{1} \quad \mathbf{X} = [x_1, \dots, x_N]$$

- Empirical covariance matrix and eigendecomposition

$$\mathbf{C} \leftarrow \frac{1}{N-1} \mathbf{X}\mathbf{X}^T \quad \mathbf{C} = \mathbf{U}\Lambda\mathbf{U}^T$$

- Whitening step

$$\mathbf{X} \leftarrow \Lambda^{-1/2} \mathbf{U}^T \mathbf{X}$$

- Krizhevsky, "Learning Multiple Layers of Features from Tiny Images"

# Training RBMs and whitening of data

- Pixels in images are highly correlated
- Whitening aims to remove 2nd order correlation from images
  - Learned model can focus on higher-order correlations
- ZCA whitening
  - Subtract empirical mean

$$\mathbf{X} \leftarrow \mathbf{X} - \frac{1}{N} \mathbf{X}\mathbf{1} \quad \mathbf{X} = [x_1, \dots, x_N]$$

- Empirical covariance matrix and eigendecomposition

$$\mathbf{C} \leftarrow \frac{1}{N-1} \mathbf{X}\mathbf{X}^T \quad \mathbf{C} = \mathbf{U}\Lambda\mathbf{U}^T$$

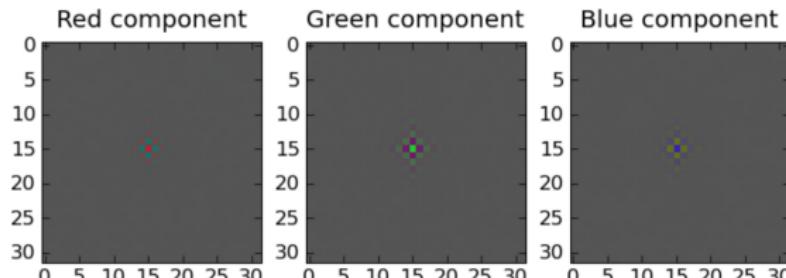
- Whitening step: also rotate back to original space

$$\mathbf{X} \leftarrow \mathbf{U}^T \Lambda^{-1/2} \mathbf{U}^T \mathbf{X}$$

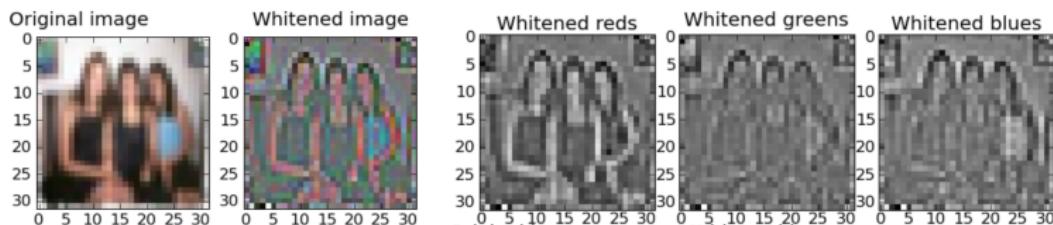
- Krizhevsky, "Learning Multiple Layers of Features from Tiny Images"

## Training RBMs and whitening of data

- Whitening filters  $\mathbf{U}^T \Lambda^{-1/2} \mathbf{U}^T$  for tiny images



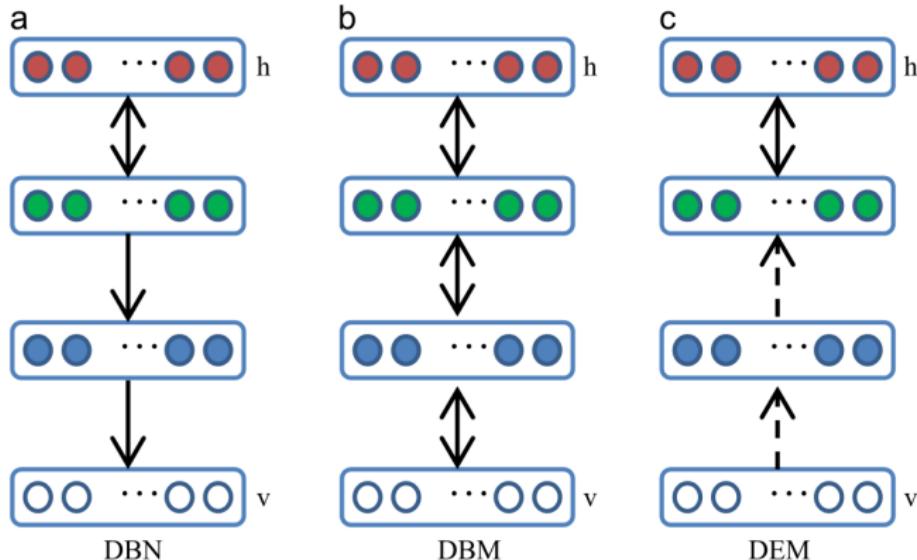
- ZCA whitening enhances edges in images



- Reason for observing Gabor filters in the 1st layer?
- PCA whitening leads to unrecognizable images
- Krizhevsky, "Learning Multiple Layers of Features from Tiny Images"

# Deep extensions of RBMs

- Deep belief networks (Hinton et al., “A Fast Learning Algorithm for Deep Belief Nets”, 2006)
- Deep Boltzmann machines (Salakhutdinov & Hinton, “Deep Boltzmann machines”, 2009)
- Deep energy models (Ngiam et al., “Learning Deep Energy Models”, 2011)



# Deep energy models (DEMs)

- Deep energy model (DEM)

$$\log p_\theta(x, z) \doteq -\frac{1}{2\sigma^2} \|x\|^2 + b^T x + c^T z + z^T W g_\theta(x) \quad z \in \{0, 1\}^K$$

- $-\|x\|^2$  makes the model (usually) normalizable
- Marginalizing over  $z$

$$p_\theta(z_k = 1|x) = s(w_k^T g_\theta(x) + c_k)$$

$$\log p_\theta(x) \doteq -\frac{1}{2\sigma^2} \|x\|^2 + b^T x + \sum_k S(w_k^T g_\theta(x) + c_k)$$

- $p_\theta(x|z)$  is difficult in general
- General DEM

$$\log p_\theta(x) \doteq -\frac{1}{2\sigma^2} \|x\|^2 + b^T x + H_\theta(x)$$

- Generalization of
  - GRBMs
  - mcRBMs
- Ngiam et al., "Learning Deep Energy Models"

## Deep energy models (DEMs)

- Deep energy model (DEM)

$$\log p_\theta(x, z) \doteq -\frac{1}{2\sigma^2} \|x\|^2 + b^T x + c^T z + z^T W g_\theta(x) \quad z \in \{0, 1\}^K$$

- $-\|x\|^2$  makes the model (usually) normalizable
- Marginalizing over  $z$

$$p_\theta(z_k = 1 | x) = s(w_k^T g_\theta(x) + c_k)$$

$$\log p_\theta(x) \doteq -\frac{1}{2\sigma^2} \|x\|^2 + b^T x + \sum_k S(w_k^T g_\theta(x) + c_k)$$

- $p_\theta(x|z)$  is difficult in general
- General DEM

$$\log p_\theta(x) \doteq -\frac{1}{2\sigma^2} \|x\|^2 + b^T x + H_\theta(x)$$

- Generalization of
  - GRBMs
  - mcRBMs
- Ngiam et al., "Learning Deep Energy Models"

## Deep energy models (DEMs)

- Deep energy model (DEM)

$$\log p_\theta(x, z) \doteq -\frac{1}{2\sigma^2} \|x\|^2 + b^T x + c^T z + z^T W g_\theta(x) \quad z \in \{0, 1\}^K$$

- $-\|x\|^2$  makes the model (usually) normalizable
- Marginalizing over  $z$

$$p_\theta(z_k = 1 | x) = s(w_k^T g_\theta(x) + c_k)$$

$$\log p_\theta(x) \doteq -\frac{1}{2\sigma^2} \|x\|^2 + b^T x + \sum_k S(w_k^T g_\theta(x) + c_k)$$

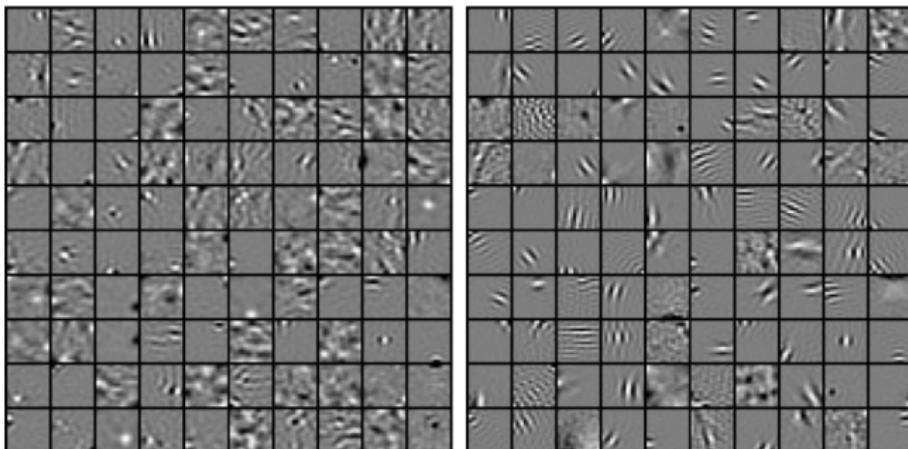
- $p_\theta(x|z)$  is difficult in general
- General DEM

$$\log p_\theta(x) \doteq -\frac{1}{2\sigma^2} \|x\|^2 + b^T x + H_\theta(x)$$

- Generalization of
  - GRBMs
  - mcRBMs
- Ngiam et al., "Learning Deep Energy Models"

## Deep energy models (DEMs)

- Learning DEM via contrastive divergence
  - Ngiam et al., "Learning Deep Energy Models", 2011

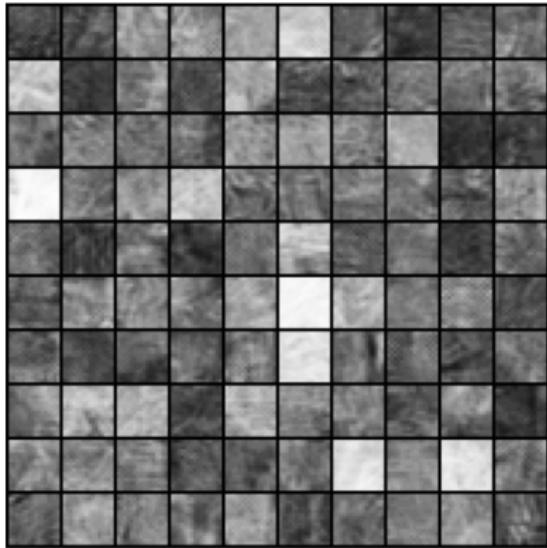
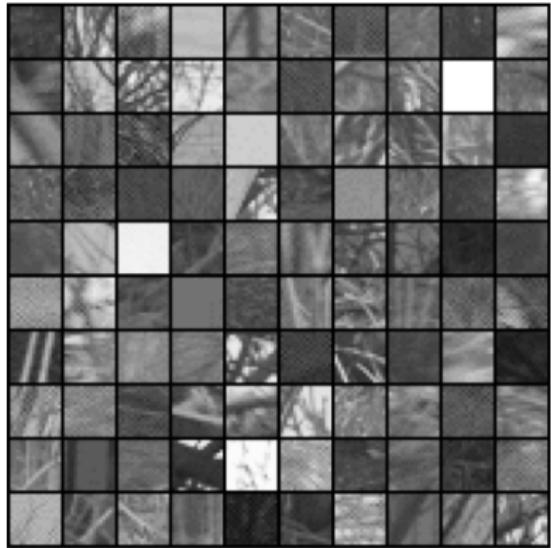


Single layer (sigmoid) DEM vs. after adding a 2nd layer

- Learning similar DEMs via score matching
  - Köster & Hyvärinen, "A Two-Layer Model of Natural Stimuli Estimated with Score Matching", 2009
  - Kingma & LeCun, "Regularized estimation of image statistics by Score Matching", 2010

## Deep energy models (DEMs)

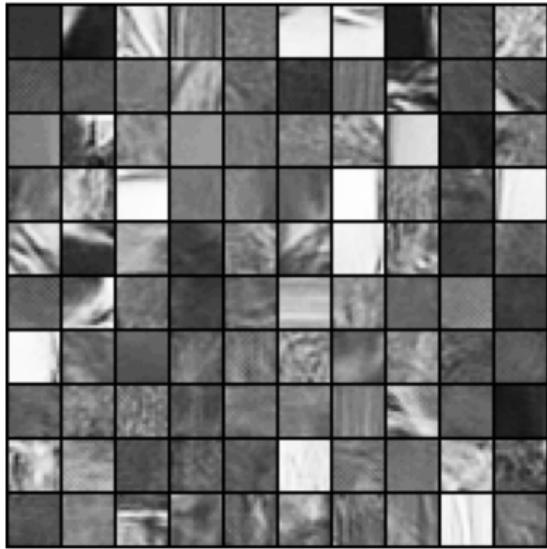
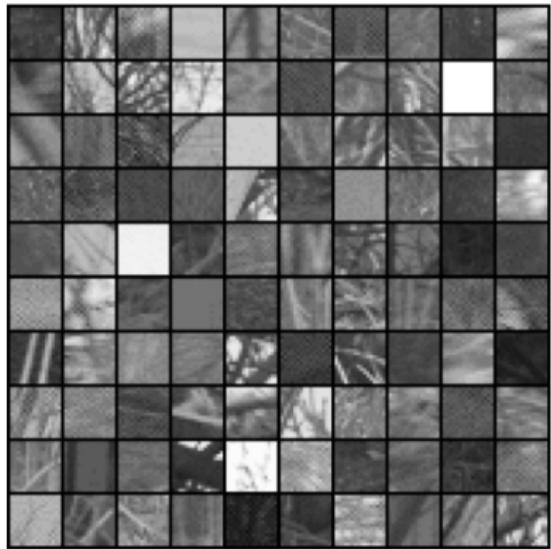
- Left: samples from training data
- Right: samples from model with 1 layer(s) in  $H_\theta$



- Ngiam et al., "Learning Deep Energy Models"

## Deep energy models (DEMs)

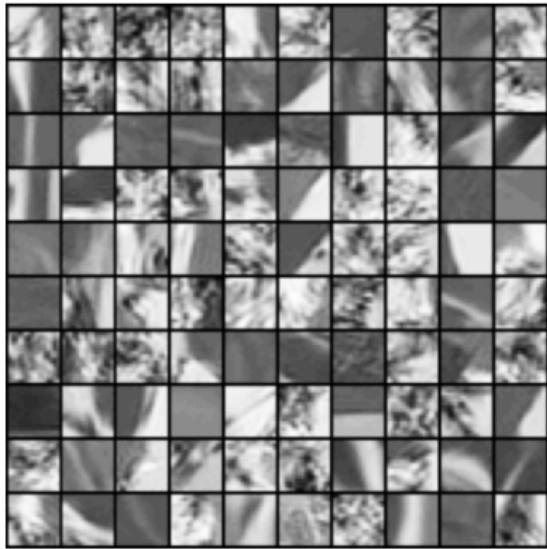
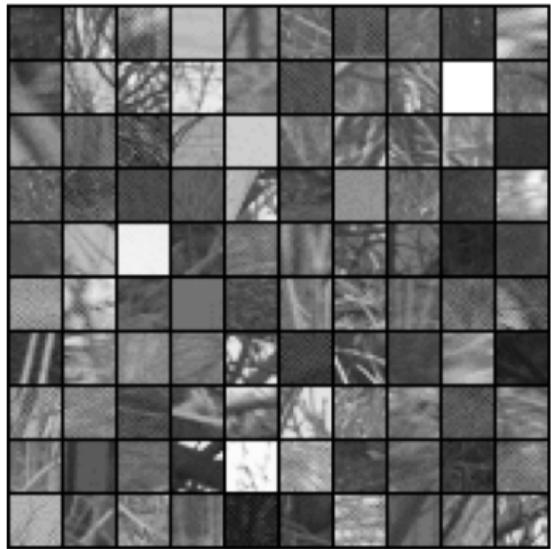
- Left: samples from training data
- Right: samples from model with 2 layer(s) in  $H_\theta$



- Ngiam et al., "Learning Deep Energy Models"

## Deep energy models (DEMs)

- Left: samples from training data
- Right: samples from model with 3 layer(s) in  $H_\theta$



- Ngiam et al., "Learning Deep Energy Models"

# Outline

## 1 Latent variable models

- Variational auto-encoders
- Variational NCE
- Boltzmann machines

## 2 Sparse coding and dictionary learning

## 3 Exercise 2

## Overview

### Sparse coding

- represents an input signal as sparse linear combination of atoms (dictionary elements)
- uses overcomplete dictionaries
- appears in
  - signal processing (wavelets)
  - statistics (LASSO)
  - compressed sensing
  - neuroscience
- Tibshirani, "Regression Shrinkage and Selection via the LASSO"
- Donoho & Elad, "Optimally sparse representation in general (nonorthogonal) dictionaries via L1 minimization"
- Olshausen & Field, "Emergence of simple-cell receptive field properties by learning a sparse code for natural images"

# Sparse coding and dictionary learning

- Main objective

$$J(\mathbf{D}, \mathbf{z}) = \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|^2 + \lambda \|\mathbf{z}\|_q \quad q \in [0, 1]$$

- Dictionary D

- Columns of D are dictionary elements / atoms
- Usually #rows(D)  $\ll$  #columns(D)

- Sparse code z: sparsity favored via  $\|\mathbf{z}\|_q$

- Sparse coding: given D and x infer

$$\mathbf{z}(\mathbf{x}) = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|^2 + \lambda \|\mathbf{z}\|_q$$

- Dictionary learning: given  $\{\mathbf{x}_i\}$  infer D s.t.  $\|\mathbf{D}(:,j)\| \leq 1$  (why?)

$$\min_{\mathbf{D}} \sum_i \left( \min_{\mathbf{z}_i} \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\mathbf{z}_i\|^2 + \lambda \|\mathbf{z}_i\|_q \right)$$

- Usually not interpreted as probabilistic model

- But: Bayesian sparse learning and automatic relevance determination

# Sparse coding and dictionary learning

- Main objective

$$J(\mathbf{D}, \mathbf{z}) = \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|^2 + \lambda \|\mathbf{z}\|_q \quad q \in [0, 1]$$

- Dictionary D

- Columns of D are dictionary elements / atoms
- Usually #rows(D)  $\ll$  #columns(D)

- Sparse code z: sparsity favored via  $\|\mathbf{z}\|_q$

- Sparse coding: given D and x infer

$$\mathbf{z}(\mathbf{x}) = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|^2 + \lambda \|\mathbf{z}\|_q$$

- Dictionary learning: given  $\{\mathbf{x}_i\}$  infer D s.t.  $\|\mathbf{D}(:,j)\| \leq 1$  (why?)

$$\min_{\mathbf{D}} \sum_i \left( \min_{\mathbf{z}_i} \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\mathbf{z}_i\|^2 + \lambda \|\mathbf{z}_i\|_q \right)$$

- Usually not interpreted as probabilistic model

- But: Bayesian sparse learning and automatic relevance determination

# Sparse coding and dictionary learning

- Main objective

$$J(\mathbf{D}, \mathbf{z}) = \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|^2 + \lambda \|\mathbf{z}\|_q \quad q \in [0, 1]$$

- Dictionary D

- Columns of D are dictionary elements / atoms
- Usually #rows(D)  $\ll$  #columns(D)

- Sparse code z: sparsity favored via  $\|\mathbf{z}\|_q$

- Sparse coding: given D and x infer

$$\mathbf{z}(\mathbf{x}) = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|^2 + \lambda \|\mathbf{z}\|_q$$

- Dictionary learning: given  $\{\mathbf{x}_i\}$  infer D s.t.  $\|\mathbf{D}(:,j)\| \leq 1$  (why?)

$$\min_{\mathbf{D}} \sum_i \left( \min_{\mathbf{z}_i} \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\mathbf{z}_i\|^2 + \lambda \|\mathbf{z}_i\|_q \right)$$

- Usually not interpreted as probabilistic model

- But: Bayesian sparse learning and automatic relevance determination

# Sparse coding and dictionary learning

- Main objective

$$J(\mathbf{D}, \mathbf{z}) = \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|^2 + \lambda \|\mathbf{z}\|_q \quad q \in [0, 1]$$

- Dictionary D

- Columns of D are dictionary elements / atoms
- Usually #rows(D)  $\ll$  #columns(D)

- Sparse code z: sparsity favored via  $\|\mathbf{z}\|_q$

- Sparse coding: given D and x infer

$$\mathbf{z}(\mathbf{x}) = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|^2 + \lambda \|\mathbf{z}\|_q$$

- Dictionary learning: given  $\{x_i\}$  infer D s.t.  $\|\mathbf{D}(:,j)\| \leq 1$  (why?)

$$\min_{\mathbf{D}} \sum_i \left( \min_{\mathbf{z}_i} \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\mathbf{z}_i\|^2 + \lambda \|\mathbf{z}_i\|_q \right)$$

- Usually not interpreted as probabilistic model

- But: Bayesian sparse learning and automatic relevance determination

# Sparse coding and dictionary learning

- In many applications  $q = 1$
- Convex but non-smooth sparse coding problem

$$\mathbf{z}(\mathbf{x}) = \arg \min_z \frac{1}{2} \|\mathbf{x} - \mathbf{Dz}\|^2 + \lambda \|\mathbf{z}\|_1$$

- Iterative Shrinkage-Thresholding Algorithm (ISTA)
  - Proximal gradient descent
- Fast Iterative Shrinkage-Thresholding Algorithm (FISTA)
  - Accelerated projected gradient descent
- Beck & Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems"

# Sparse coding and dictionary learning

- Sparse coding objective

$$\mathbf{z}(\mathbf{x}) = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|^2 + \lambda \|\mathbf{z}\|_1$$

- ISTA

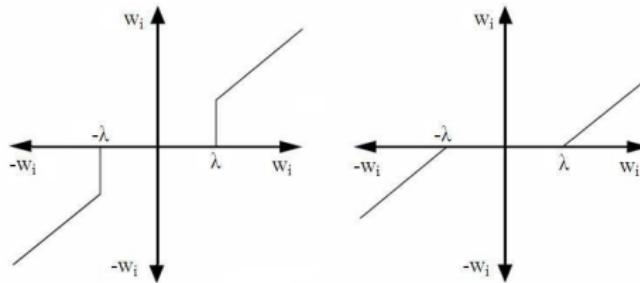
$$\tau \in (0, 1/\|\mathbf{D}^T \mathbf{D}\|_2), \mathbf{z}^{(0)} \leftarrow \mathbf{0}$$

$$\mathbf{z}^{(t+1/2)} \leftarrow \mathbf{z}^{(t)} - \tau \mathbf{D}^T (\mathbf{D}\mathbf{z}^{(t)} - \mathbf{x})$$

$$\mathbf{z}^{(t+1)} \leftarrow \arg \min_{\mathbf{z}} \lambda \|\mathbf{z}\|_1 + \frac{1}{2\tau} \|\mathbf{z} - \mathbf{z}^{(t+1/2)}\|^2$$

- Define

$$\mathcal{S}_\lambda(\mathbf{z}^0) := \arg \min_{\mathbf{z}} \lambda \|\mathbf{z}\|_1 + \frac{1}{2} \|\mathbf{z} - \mathbf{z}^0\|^2 = \begin{pmatrix} \vdots \\ [|z_j^0| - \lambda]_+ \operatorname{sgn}(z_j^0) \\ \vdots \end{pmatrix}$$



# Sparse coding and dictionary learning

- Sparse coding objective

$$\mathbf{z}(\mathbf{x}) = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|^2 + \lambda \|\mathbf{z}\|_1$$

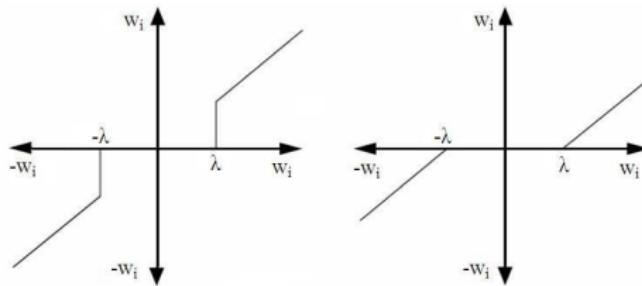
- ISTA in one line

$$\tau \in (0, 1/\|\mathbf{D}^T \mathbf{D}\|_2), \mathbf{z}^{(0)} \leftarrow \mathbf{0}$$

$$\mathbf{z}^{(t+1)} \leftarrow \mathcal{S}_{\tau\lambda} \left( \mathbf{z}^{(t)} - \tau \mathbf{D}^T (\mathbf{D}\mathbf{z}^{(t)} - \mathbf{x}) \right)$$

- Define

$$\mathcal{S}_\lambda(\mathbf{z}^0) := \arg \min_{\mathbf{z}} \lambda \|\mathbf{z}\|_1 + \frac{1}{2} \|\mathbf{z} - \mathbf{z}^0\|^2 = \begin{pmatrix} \vdots \\ [|z_j^0| - \lambda]_+ \operatorname{sgn}(z_j^0) \\ \vdots \end{pmatrix}$$

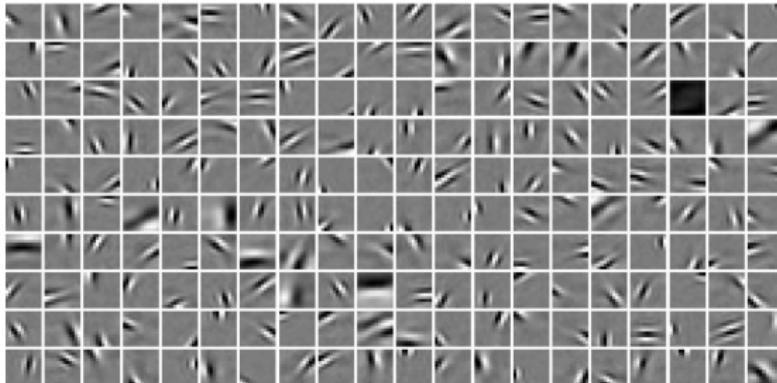


# Sparse coding and dictionary learning

- Jointly learning sparse codes and dictionary

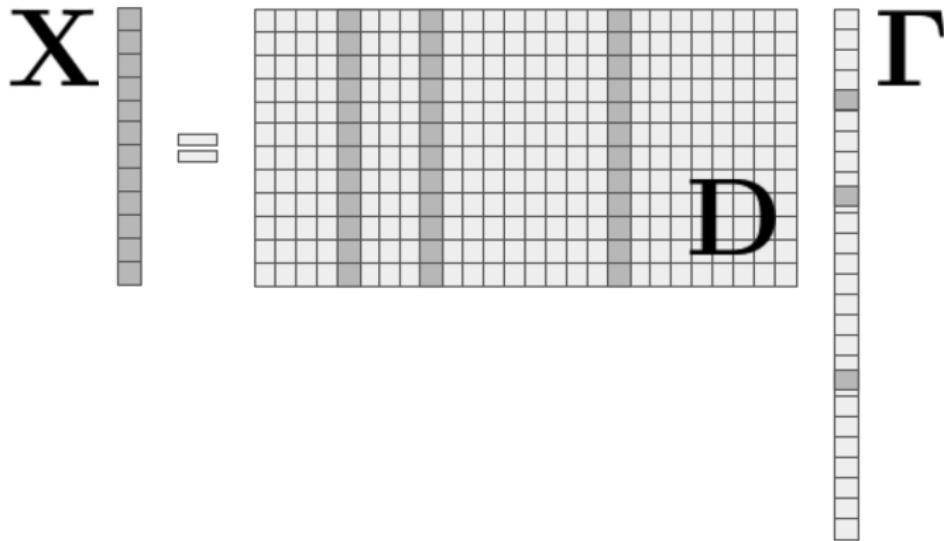
$$\min_{\mathbf{D}: \|\mathbf{D}(:,j)\|_1 \leq 1} \min_{\{\mathbf{z}_i\}} \sum_i \left( \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\mathbf{z}_i\|^2 + \lambda \|\mathbf{z}_i\|_1 \right)$$

- Block-convex problem
- E.g. minimization by alternation
  - Find  $\mathbf{z}_i$  for all  $i$  via FISTA
  - Update  $\mathbf{D}$  via convex optimization (or column-by-column in closed form)



# Convolutional sparse coding

- Going from patches to images
- Sparse coding



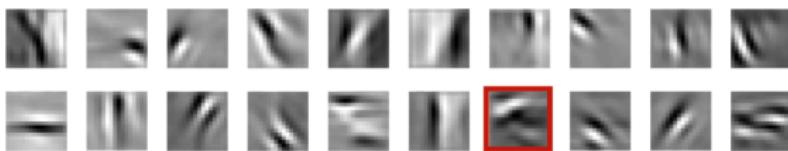
# Convolutional sparse coding

- Going from patches to images
- Convolutional sparse coding

$$\mathbf{X} = \mathbf{D} \mathbf{I}$$

The diagram illustrates the convolutional sparse coding equation. On the left, the input patch  $\mathbf{X}$  is shown as a grid of colored bars. An equals sign follows. To the right of the equals sign is a large matrix  $\mathbf{D}$ , which is also a grid of colored bars. To the right of  $\mathbf{D}$  is the output image  $\mathbf{I}$ , represented as a grid of colored bars. The matrix  $\mathbf{D}$  has a diagonal band of non-zero values, indicating the receptive fields of the neurons in the image  $\mathbf{I}$ . The input  $\mathbf{X}$  and output  $\mathbf{I}$  have vertical color bars at their boundaries, while the matrix  $\mathbf{D}$  has horizontal color bars.

# Convolutional sparse coding



- Bristow et al., "Fast Convolutional Sparse Coding"

# Deep Sparse Coding

- Hierarchical (deep) sparse coding

$$\begin{array}{ll} \text{find } \{\mathbf{z}_k\}_{k=1}^K \text{ s.t. } & \|\mathbf{x} - \mathbf{D}_1 \mathbf{z}_1\| \leq \varepsilon_1 \quad \|\mathbf{z}_1\|_1 \leq \mu_1 \\ & \|\mathbf{z}_1 - \mathbf{D}_2 \mathbf{z}_2\| \leq \varepsilon_2 \quad \|\mathbf{z}_2\|_1 \leq \mu_2 \\ & \vdots \quad \vdots \\ & \|\mathbf{z}_{K-1} - \mathbf{D}_K \mathbf{z}_K\| \leq \varepsilon_K \quad \|\mathbf{z}_K\|_1 \leq \mu_K \end{array}$$

- Papyan et al., "Convolutional Neural Networks Analyzed via Convolutional Sparse Coding"

# Deep Sparse Coding

- A layered thresholding algorithm

$$\mathbf{z}_1 \leftarrow \mathcal{S}_{\lambda_1}(\mathbf{D}_1^T \mathbf{x})$$

$$\mathbf{z}_2 \leftarrow \mathcal{S}_{\lambda_2}(\mathbf{D}_2^T \mathbf{z}_1)$$

⋮

$$\mathbf{z}_K \leftarrow \mathcal{S}_{\lambda_K}(\mathbf{D}_K^T \mathbf{z}_{K-1})$$

- One step of ISTA (with  $\mathbf{z}_k^{(0)} = 0$ )
- Starts to look like a ReLU DNN
  - Non-negative sparse coding: ReLU DNN
  - Convolutional sparse coding: ReLU CNN
- Popyan et al., “Convolutional Neural Networks Analyzed via Convolutional Sparse Coding”

# Deep Sparse Coding

- Non-negative sparse coding

$$\mathbf{z}(\mathbf{x}) = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|^2 + \lambda \mathbf{z}$$

subject to  $\mathbf{z} \geq 0$

- Proximal gradient method

$$\mathbf{z}^{(0)} \leftarrow 0$$

$$\mathbf{z}^{(t+1)} \leftarrow \mathcal{S}_{\tau\lambda}^+ \left( \mathbf{z}^{(t)} - \tau \mathbf{D}^T (\mathbf{D}\mathbf{z}^{(t)} - \mathbf{x})^2 \right)$$

$$\mathcal{S}_\lambda^+(\mathbf{z}^0) := \arg \min_{\mathbf{z} \geq 0} \lambda \mathbf{z} + \frac{1}{2} \|\mathbf{z} - \mathbf{z}^0\|^2 = \begin{pmatrix} \vdots \\ [z_j^0 - \lambda]_+ \\ \vdots \end{pmatrix}$$

- Papyan et al., "Convolutional Neural Networks Analyzed via Convolutional Sparse Coding"

# Deep Sparse Coding

- A layered thresholding algorithm for non-negative SC

$$\mathbf{z}_1 \leftarrow \mathcal{S}_{\lambda_1}^+(\mathbf{D}_1^T \mathbf{x}) = [\mathbf{D}_1^T \mathbf{x} - \lambda_1]_+$$

$$\mathbf{z}_2 \leftarrow \mathcal{S}_{\lambda_2}^+(\mathbf{D}_2^T \mathbf{z}_1) = [\mathbf{D}_2^T \mathbf{z}_1 - \lambda_2]_+$$

⋮

$$\mathbf{z}_K \leftarrow \mathcal{S}_{\lambda_K}^+(\mathbf{D}_K^T \mathbf{z}_{K-1}) = [\mathbf{D}_K^T \mathbf{z}_{K-1} - \lambda_K]_+$$

- Looks like a ReLU DNN
- Convolutional sparse coding: ReLU CNN

## Discussion

Is this analogy between deep sparse coding and DNNs too far fetched?

- Very crude approximation of ISTA
- Impact of batch normalization etc.

- Popyan et al., "Convolutional Neural Networks Analyzed via Convolutional Sparse Coding"

# Deconvolutional networks

- Layerwise convolutional sparse coding

$$\mathbf{z}_l = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{z}_{l-1} - \mathbf{D}_l * \mathbf{z}\|^2 + \lambda_l \|\mathbf{z}_l\|_1$$

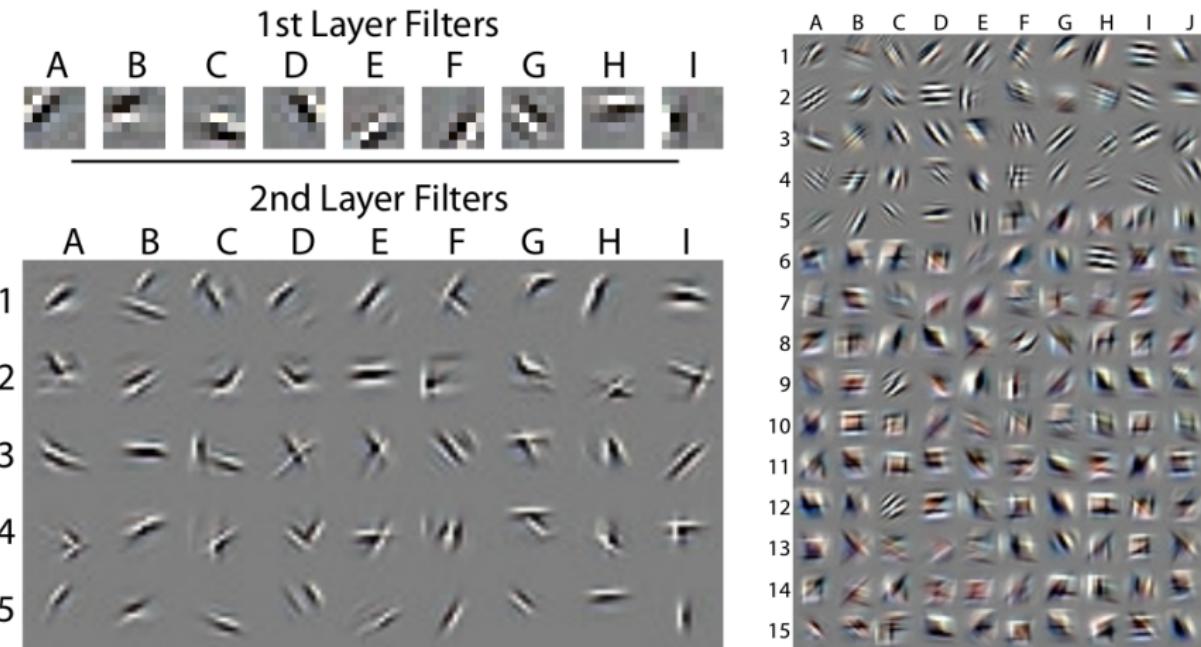
- Reconstruction operator by chaining, e.g.

$$x' = \mathbf{D}_1 * (\mathbf{D}_2 * (\mathbf{D}_3 * \mathbf{z}_3))$$

- Layerwise training

- Zeiler et al., "Deconvolutional Networks"

# Deconvolutional networks



- Zeiler et al., “Deconvolutional Networks”

# Adaptive deconvolutional networks

- Add max-pooling / unpooling
  - Store pooling index to unpool
- Reconstruct input (instead of previous layer)

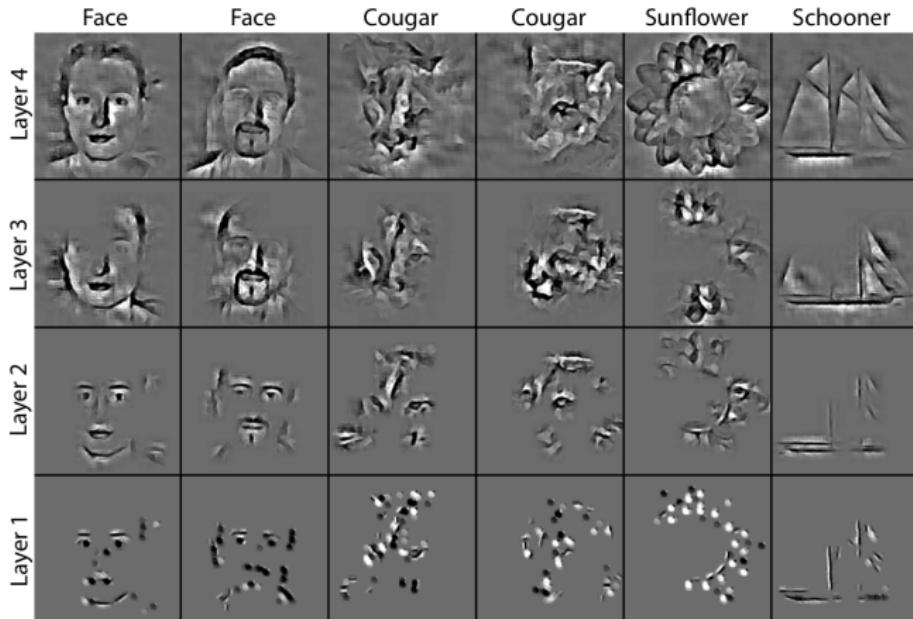
$$\mathbf{z}_l = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{x} - \mathbf{R}_l * \mathbf{z}\|^2 + \lambda_l \|\mathbf{z}_l\|_1$$

- Reconstruction operator by chaining, e.g.

$$x' = \mathbf{R}_3 \mathbf{z}_3 = \mathbf{D}_1 * \mathbf{U}_1 (\mathbf{D}_2 * \mathbf{U}_2 (\mathbf{D}_3 * \mathbf{z}_3))$$

- Layerwise training
- Zeiler et al., "Adaptive Deconvolutional Networks for Mid and High Level Feature Learning"

# Adaptive deconvolutional networks



- Zeiler et al., "Adaptive Deconvolutional Networks for Mid and High Level Feature Learning"

# Outline

## 1 Latent variable models

- Variational auto-encoders
- Variational NCE
- Boltzmann machines

## 2 Sparse coding and dictionary learning

## 3 Exercise 2

## Exercise 2: learning a two layer DEM on *your* image patches

### Exercise

- Extract 50.000 random  $28 \times 28$  natural image patches from a set of images
  - Holiday pictures, public domain pictures
- Convert patches to grayscale  $\in [0, 1]^{28 \times 28}$
- Compute a (constrained) Gaussian representing this data
  - Computing the empirical mean
  - Using the method from exercise 1 to estimate a constrained precision matrix  $\Lambda$
  - Use the method of your choice (SM, NCE, cNCE)
  - We will optionally whiten the patches using this estimate in the later steps
- Define  $S(u) := \log(1 + e^u)$
- Attention: numerically robust way to compute  $S$  and  $s = S'$

$$S(u) = \begin{cases} \log(1 + e^u) & \text{if } u \leq 0 \\ u + \log(1 + e^{-u}) & \text{if } u \geq 0 \end{cases} \quad s(u) = \begin{cases} \frac{e^u}{1+e^u} & \text{if } u \leq 0 \\ \frac{1}{1+e^{-u}} & \text{if } u \geq 0 \end{cases}$$

## Exercise 2: learning a two layer DEM on *your* image patches

### Exercise

- Train a 2-layer DEM of the following form (using a method of your choice)

$$\log p_\theta(x) \doteq -\frac{1}{2\sigma^2} \|x\|^2 + b^T x + \sum_{k=1}^K S(w_k^T g_\theta(x) + c_k)$$
$$g_\theta(x) = s(Vx) \quad \text{single layer sigmoid NN}$$

- Parameters:  $W, V, b, c$  e.g.  $K = 64, V \in \mathbb{R}^{64 \times 784}$
- Choose  $\sigma \in \{1, 0.1\}$
- Run on whitened and non-whitened data
- Visualize filters  $v_j$  (rows in  $V$ )
- Is your model able to distinguish between
  - Hold-out natural patches
  - Samples generated by the fitted Gaussian ( $\mathcal{N}(\mathbf{0}, I)$  in the whitened version)
  - MNIST digits
- Report mean and std. deviation of  $\log p_\theta$
- Bonus: generate modes of  $p_\theta$  by maximizing  $\log p_\theta(x)$ 
  - Initial  $x$ : sample from fitted Gaussian