
Image Generation

Slides adapted from S. Lazebnik, and others

Image generation

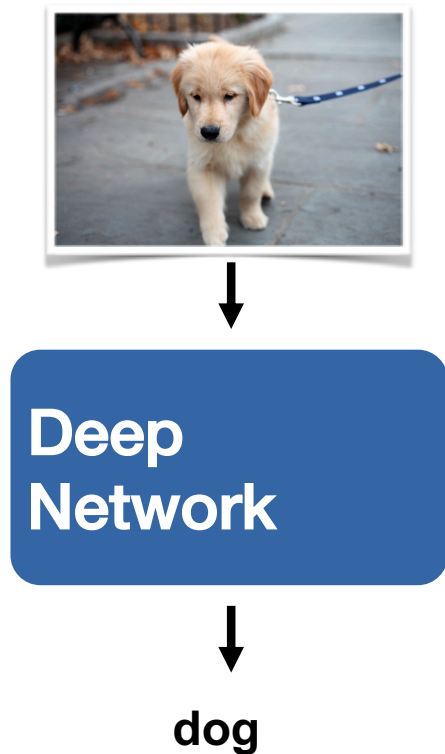
Input: Image

- High dimensional
- Structured

Output: Label

- Low dimensional

Easy



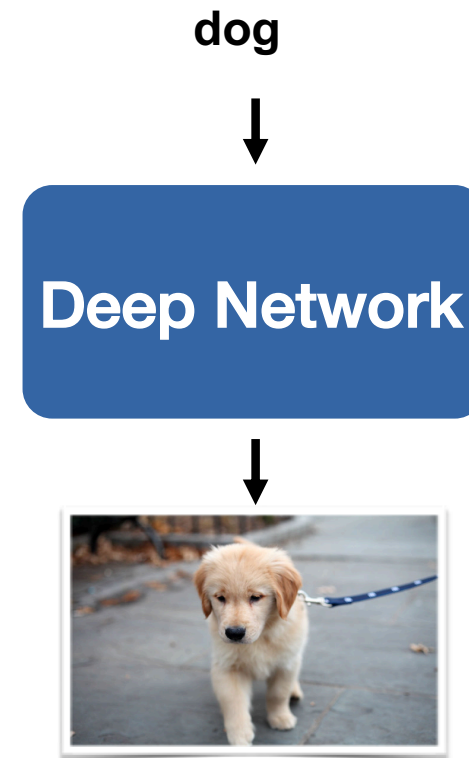
Input: Label / nothing

- Low dimensional

Output: Image

- High dimensional
- Many possibilities

Very hard



Boltzman Machines

Generative graphical models

Graphical model (Markov Random Field)

- Nodes are random variables
- Edges compatibility constraints
- Boltzman Machines – general undirected graphs
- nodes – binary random variables – edges (weights)
- Graph defines so energy function

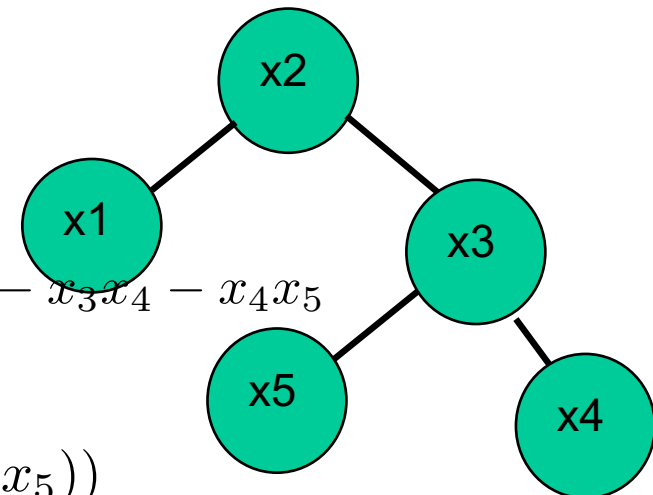
$$E(x_1, x_2, x_3, x_4, x_5) = - \sum_{i=1}^5 b_i x_i + \sum_i \sum_j x_i w_{ij} x_j$$

$$E(x_1, x_2, x_3, x_4, x_5) = x_1 + x_2 + 2x_3 - x_1x_2 - x_2x_3 - x_3x_4 - x_4x_5$$

- Associated probability distribution

$$P(x_1, x_2, x_3, x_4, x_5) = \frac{1}{Z} \exp(-E(x_1, x_2, x_3, x_4, x_5))$$

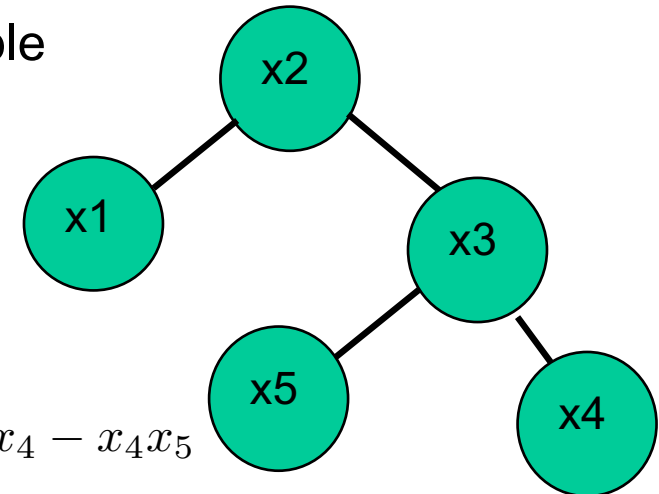
- Structured probabilistic models



Boltzman Machine

Inference Tasks

- What is the probability of particular configuration
- What the the maximum energy
- What is the most probable configuration
- Combinatorial problem (evaluate over all possible
- 5 dim binary vectors 2^5
- use sampling



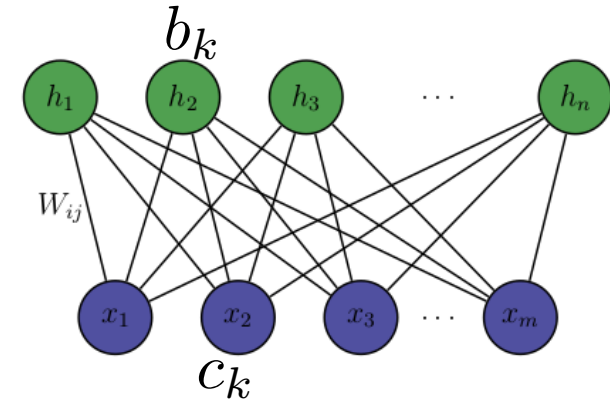
$$E(x_1, x_2, x_3, x_4, x_5) = - \sum_{i=1}^5 b_i x_i + \sum_i \sum_j x_i w_{ij} x_j$$

$$E(x_1, x_2, x_3, x_4, x_5) = x_1 + x_2 + 2x_3 - x_1x_2 - x_2x_3 - x_3x_4 - x_4x_5$$

$$P(x_1, x_2, x_3, x_4, x_5) = \frac{1}{Z} \exp(-E(x_1, x_2, x_3, x_4, x_5))$$

Restricted Boltzman machines

- Graph is bipartite
- \mathbf{x} set of observable nodes
- \mathbf{h} set of hidden nodes
- \mathbf{c} , \mathbf{b} are bias vectors
- Define probability distribution over \mathbf{x} and \mathbf{h}
- Z partition function intractable



$$E(\mathbf{x}, \mathbf{h}) = -\mathbf{h}^T \mathbf{W} \mathbf{x} - \mathbf{c}^T \mathbf{x} - \mathbf{b}^T \mathbf{h}$$

$$p(\mathbf{x}, \mathbf{h}) = \exp(-E(\mathbf{x}, \mathbf{h})) / Z$$

- Some intuition how parameters affect the probability of configuration

Restricted Boltzman machines

Factor view, probability as product of factors

$$p(\mathbf{x}, \mathbf{h}) = \exp(-E(\mathbf{x}, \mathbf{h})) / Z$$

$$p(\mathbf{x}, \mathbf{h}) = \exp(-(-\mathbf{h}^T \mathbf{W} \mathbf{x} - \mathbf{c}^T \mathbf{x} - \mathbf{b}^T \mathbf{h})) / Z$$

$$p(\mathbf{x}, \mathbf{h}) = \exp(\mathbf{h}^T \mathbf{W} \mathbf{x}) \exp(\mathbf{c}^T \mathbf{x}) \exp(\mathbf{b}^T \mathbf{h}) / Z$$

Inference in RBM's

$$p(\mathbf{x}, \mathbf{h}) = \exp(\mathbf{h}^T \mathbf{W} \mathbf{x}) \exp(\mathbf{c}^T \mathbf{x}) \exp(\mathbf{b}^T \mathbf{h}) / Z$$

Computation of probability – intractable (approximate)

Conditional inferences (conditional probabilities) possible

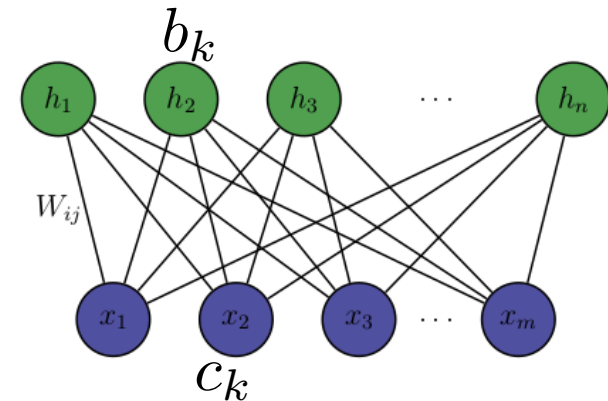
- similar to outputs on NN

$$p(\mathbf{h}|\mathbf{x}) = \prod_j p(h_j|\mathbf{x})$$

$$p(h_j = 1|\mathbf{x}) = \frac{1}{1 + \exp(-(-b_j + \mathbf{W}_j \mathbf{x}))}$$

$$p(\mathbf{x}|\mathbf{h}) = \prod_k p(x_k|\mathbf{h})$$

$$p(x_k = 1|\mathbf{h}) = \frac{1}{1 + \exp(-(-c_k + \mathbf{h}^T \mathbf{W}_k))}$$



This can be derived from

$$p(\mathbf{h}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{h})}{\sum_{\mathbf{h}'} p(\mathbf{x}, \mathbf{h}')}$$

Partition function cancels out

RBM's

$$p(\mathbf{x}, \mathbf{h}) = \exp(\mathbf{h}^T \mathbf{W} \mathbf{x}) \exp(\mathbf{c}^T \mathbf{x}) \exp(\mathbf{b}^T \mathbf{h}) / Z$$

Free energy

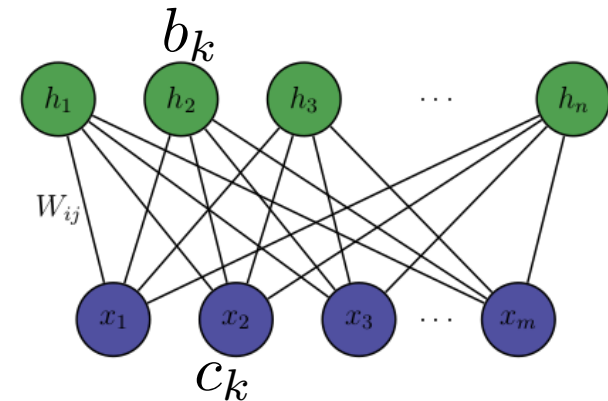
How to model

$$p(\mathbf{x}) = \sum_{\mathbf{h} \in \{0,1\}^H} p(\mathbf{x}, \mathbf{h})$$

... after some derivations

$$= \exp(\mathbf{c}^T \mathbf{x} + \sum_{i=1}^H \log(1 + \exp(b_j + \mathbf{W}_j \mathbf{x}))) / Z$$

$$p(\mathbf{x}) = \exp(-F(\mathbf{x})) / Z$$



This can be derived from

$$p(\mathbf{h}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{h})}{\sum_{h'} p(\mathbf{x}, \mathbf{h}')}$$

Partition function cancels out

Training RBM's

$$p(\mathbf{x}, \mathbf{h}) = \exp(\mathbf{h}^T \mathbf{W} \mathbf{x}) \exp(\mathbf{c}^T \mathbf{x}) \exp(\mathbf{b}^T \mathbf{h}) / Z$$

$$p(\mathbf{x}) = \exp(-F(\mathbf{x})) / Z$$

How to train RBM that models well the data
 Minimize average negative log likelihood

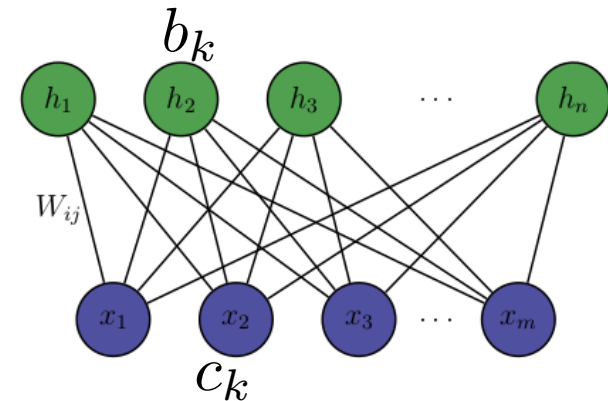
$$\frac{1}{T} \sum_t -\log p(\mathbf{x}^{(t)})$$

$$\frac{\partial -\log p(\mathbf{x}^{(t)})}{\partial \theta} = \underbrace{\mathbb{E}_{\mathbf{h}} \left[\frac{\partial E(\mathbf{x}^{(t)}, \mathbf{h})}{\partial \theta} \mid \mathbf{x}^{(t)} \right]}_{\text{positive phase}} - \underbrace{\mathbb{E}_{\mathbf{x}, \mathbf{h}} \left[\frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta} \right]}_{\text{negative phase}}$$

Dahl, Adams, Larochelle Training Restricted Boltzmann Machines on Word Observations, ICML 2012

Slides from Hugo Larochelle

http://info.usherbrooke.ca/hlarochelle/neural_networks/content.html



This can be derived from

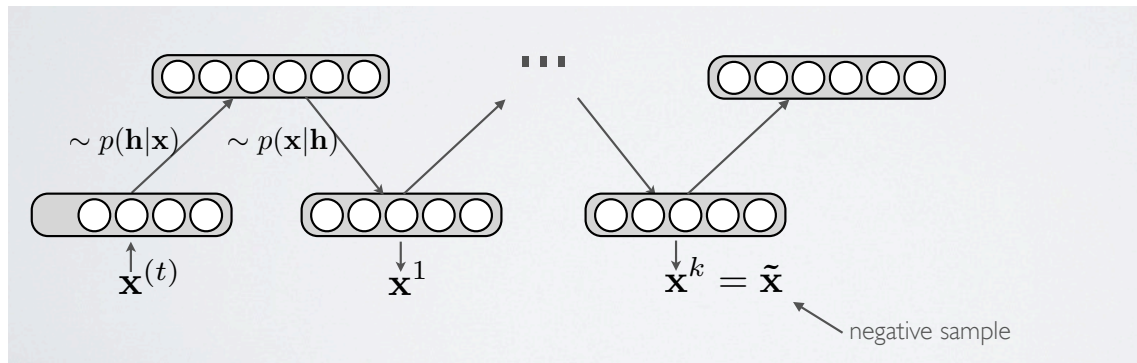
$$p(\mathbf{h} | \mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{h})}{\sum_{h'} p(\mathbf{x}, \mathbf{h}')}$$

Partition function cancels out

Contrastive Divergence

Idea

1. Replace expectation by point estimate $\tilde{\mathbf{x}}$
2. Obtain the point estimate by Gibbs sampling $\tilde{\mathbf{x}}$
3. Start sampling chain at $\mathbf{x}^{(t)}$

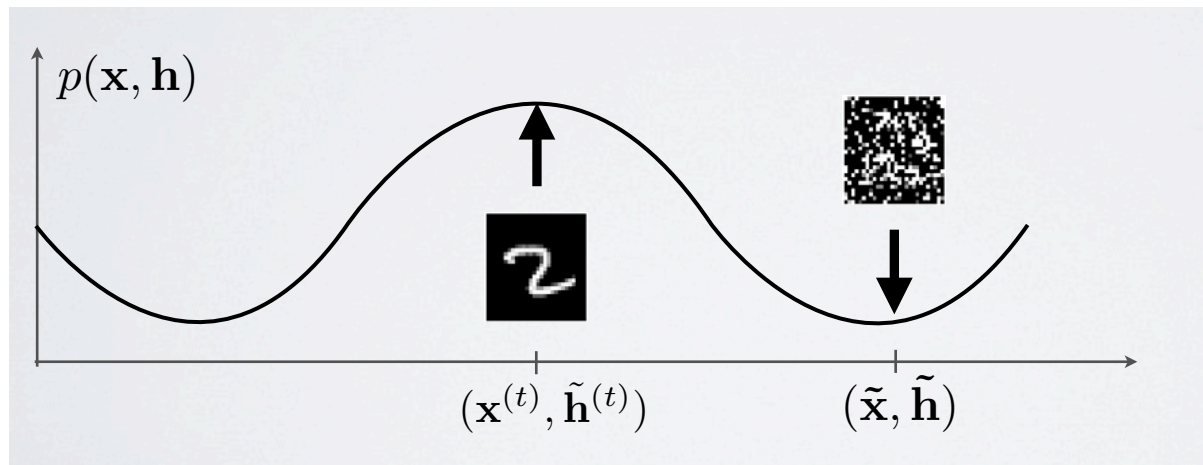
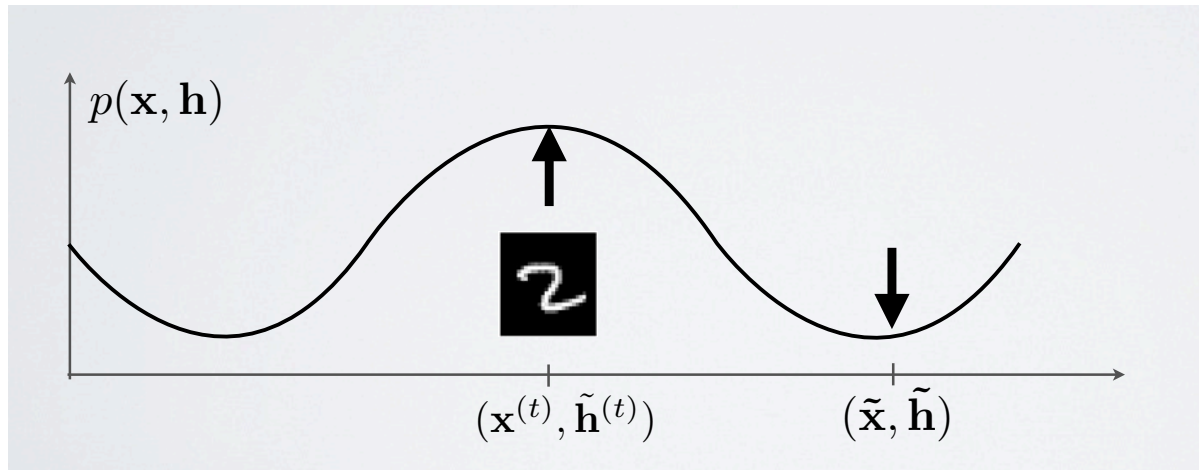


Slides from Hugo Larochelle

http://info.usherbrooke.ca/hlarochelle/neural_networks/content.html

Contrastive divergence training

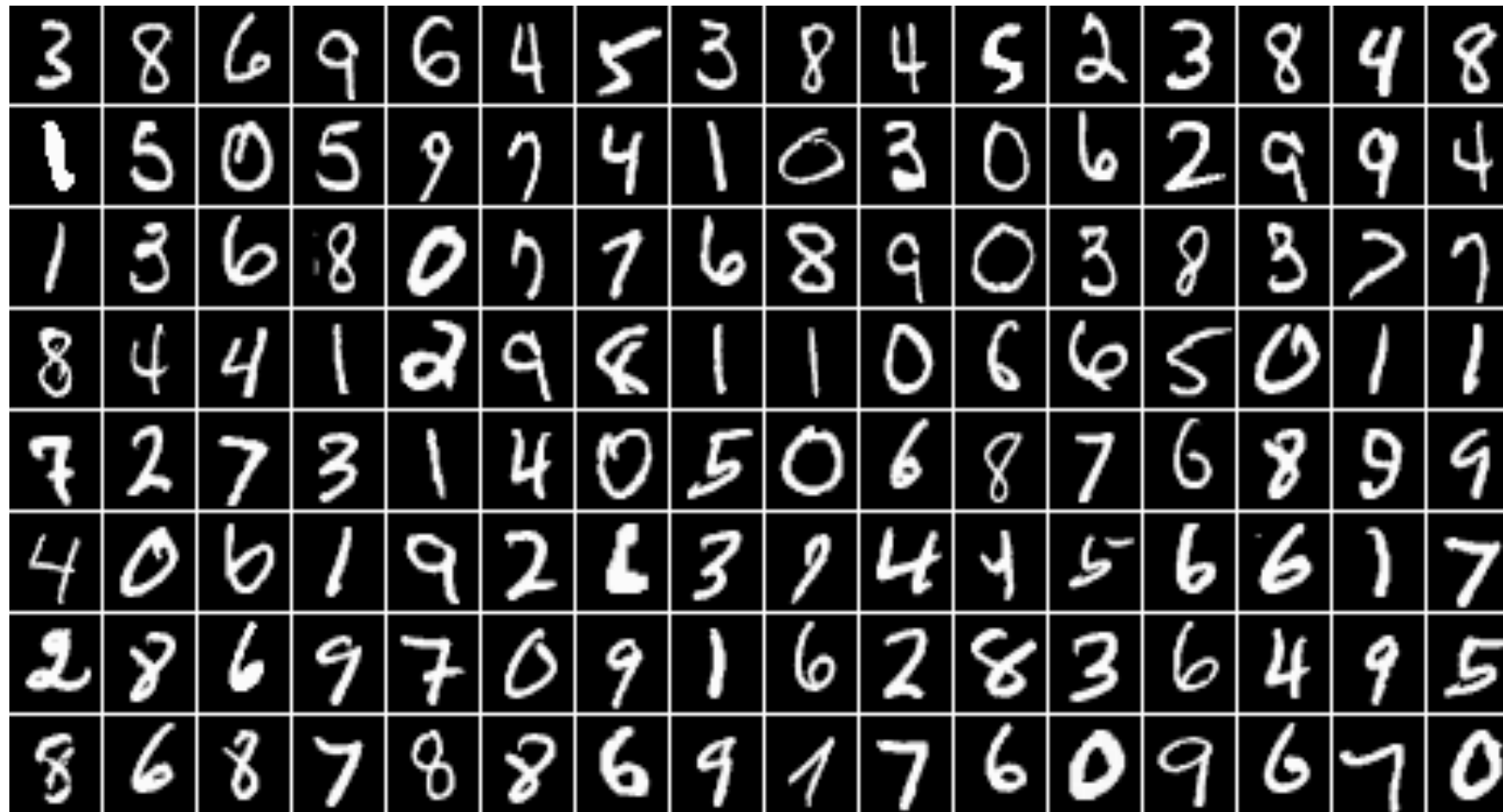
As you proceed sampled values are going to look more
Training examples



Slides from Hugo Larochelle

http://info.usherbrooke.ca/hlarochelle/neural_networks/content.html

Example

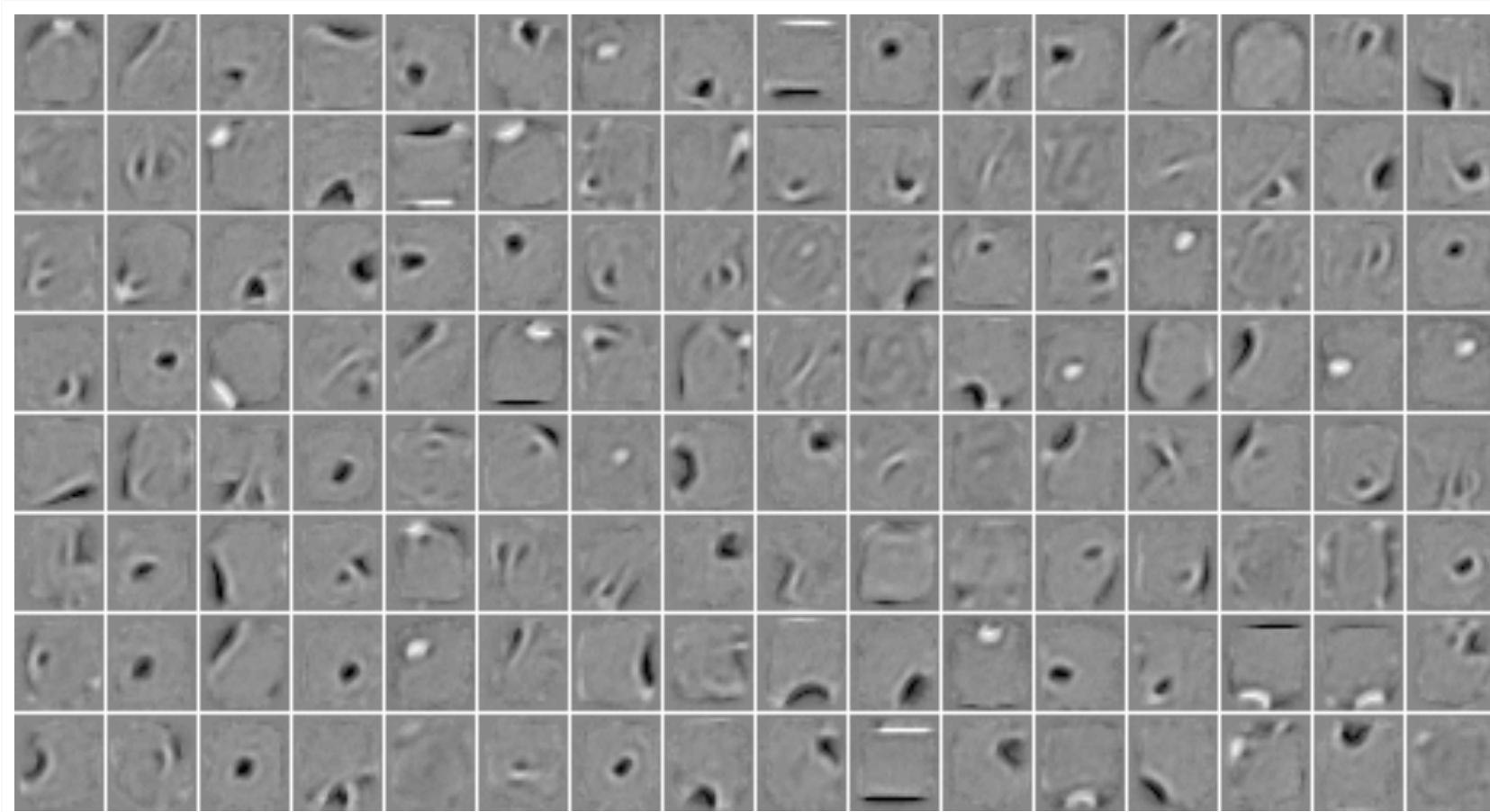


Larochelle et al JMLR 2009

Slides from Hugo Larochelle

http://info.usherbrooke.ca/hlarochelle/neural_networks/content.html

Learned filters

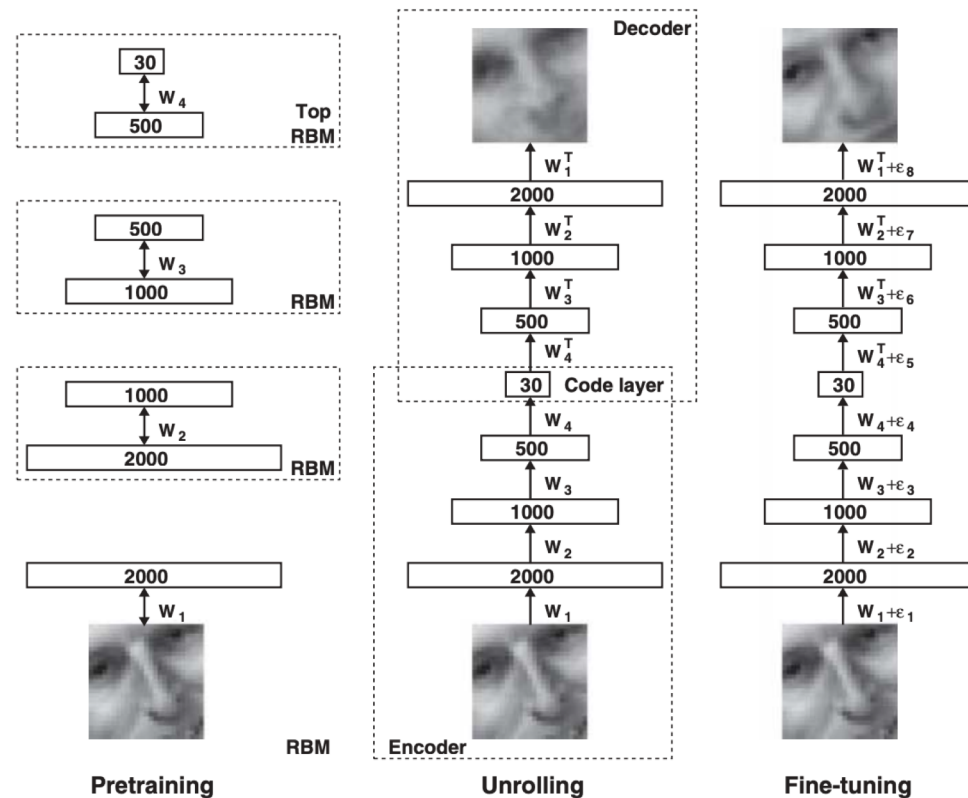


Larochelle et al JMLR 2009

Nonlinear dimensionality reduction

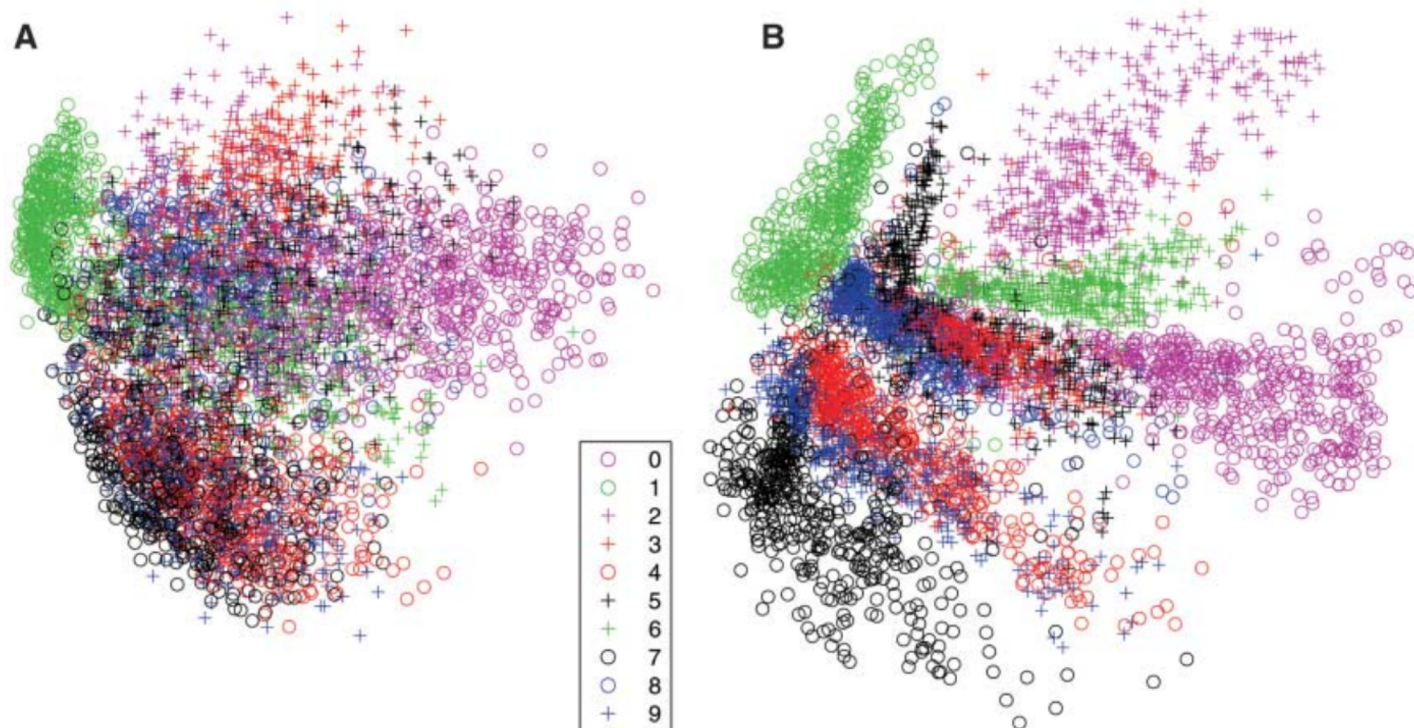
- extensions to discrete random variables
- deeper models, stack of RBM's

Reducing the Dimensionality of Data with Neural Networks, Hinton and Salakhutdinov, Science, 2006



Nonlinear Dimensionality reduction

Fig. 3. (A) The two-dimensional codes for 500 digits of each class produced by taking the first two principal components of all 60,000 training images. (B) The two-dimensional codes found by a 784-1000-500-250-2 autoencoder. For an alternative visualization, see (8).



*Reducing the Dimensionality of Data with Neural Networks,
Hinton and Salakhutdinov, Science, 2006*

Autoencoder

Feedforward NN

Tries to reconstruct input

The dim of output is the same as input

For binary inputs sigm

Single layer NN encoder

$$\mathbf{h}(\mathbf{x}) = \text{sigm}(\mathbf{b} + \mathbf{W}\mathbf{x})$$

Decoder

$$\hat{\mathbf{x}} = \text{sigm}(\mathbf{c} + \mathbf{W}^*\mathbf{h}(\mathbf{x}))$$

Sometimes weights are tied

$$\mathbf{W}^* = \mathbf{W}^T$$

Hidden representation maintains everything about input – idea compression of data



Encoder



Bottleneck

Decoder



Loss function

For binary inputs

$$l(f(\mathbf{x})) = - \sum_k (x_k \log(\hat{x}_k) + (1 - x_k) \log(1 - \hat{x}_k))$$

For real valued inputs

$$l(f(\mathbf{x})) = \frac{1}{2} \sum_k (\hat{x}_k - x_k)^2$$

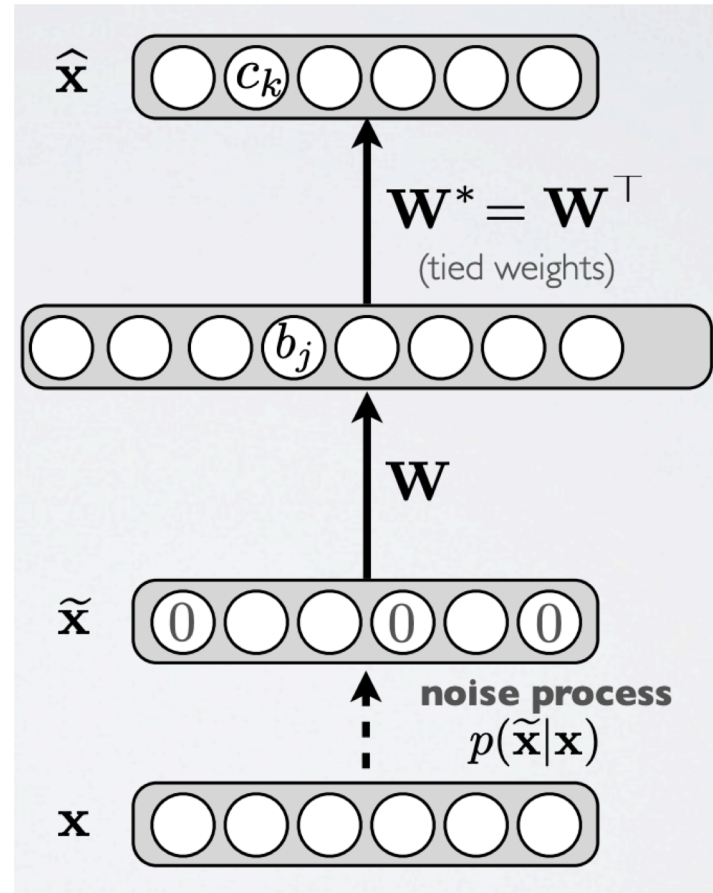
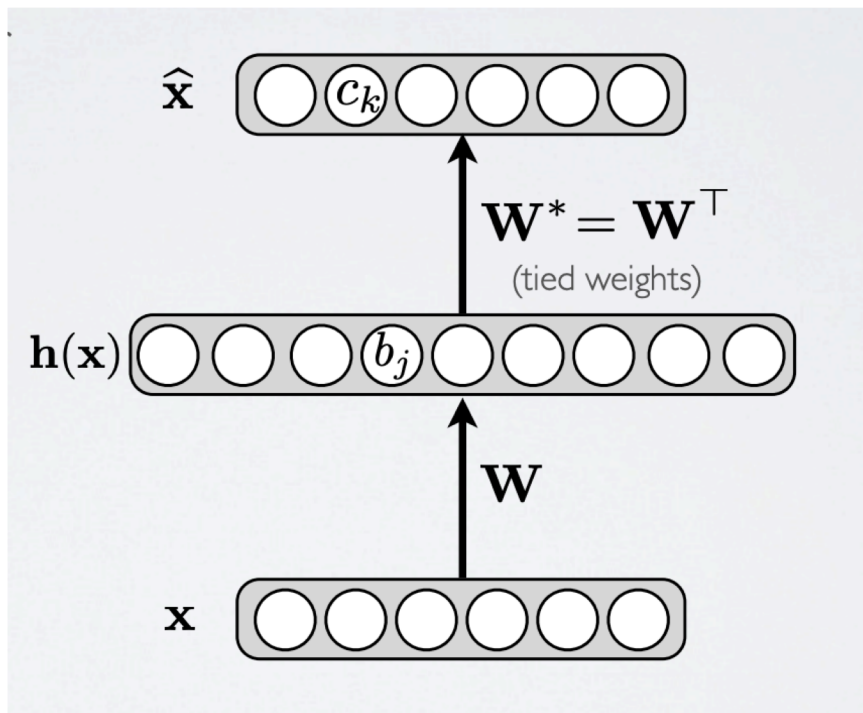
Gradient simple form – regular backprop

Choices

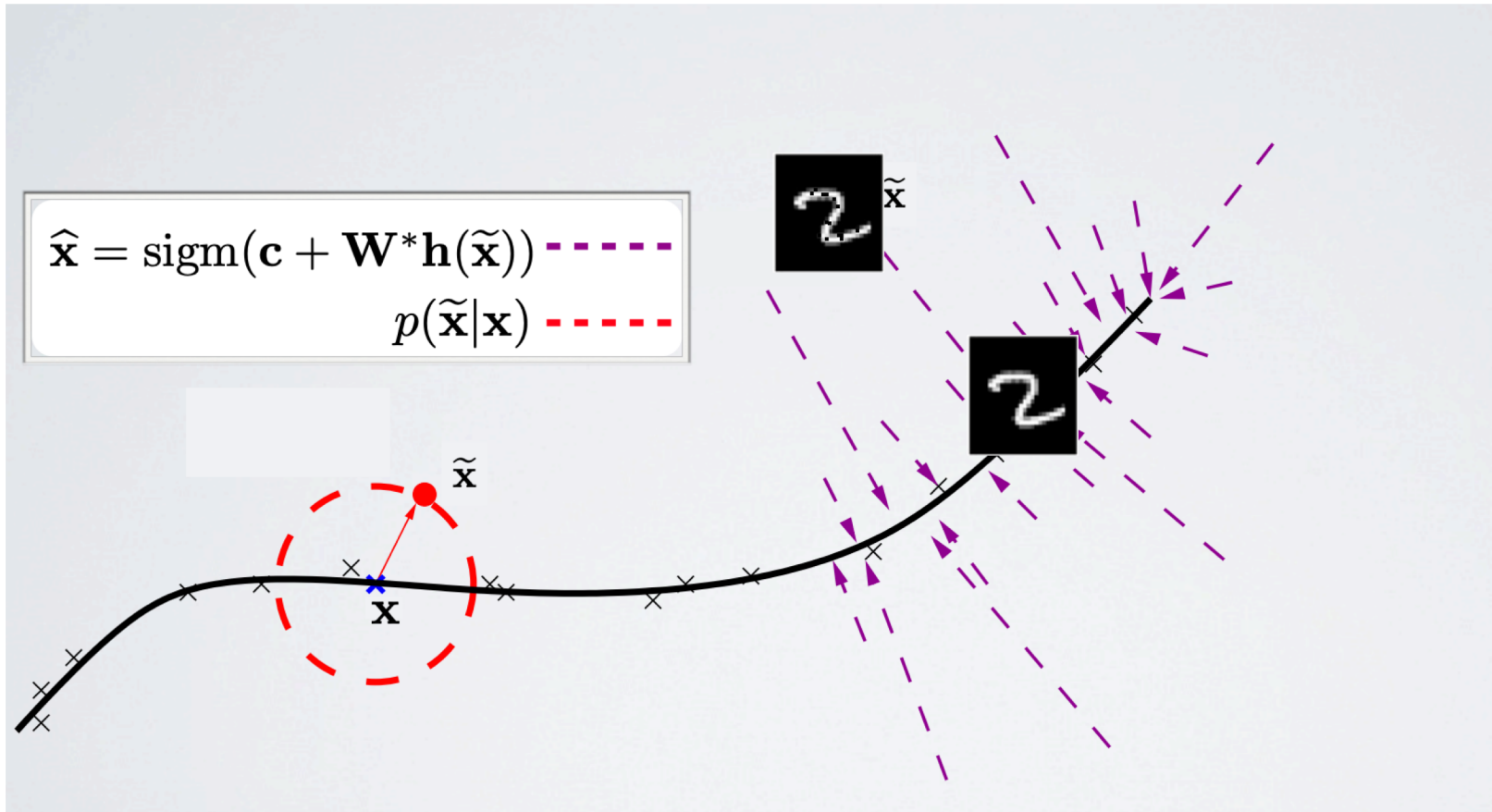
- hidden layers undercomplete, compression for the training distribution, not so good for other examples
- hidden layers overcomplete, does not extract useful representation

Denoising Autoencoder

Idea: train with the noise corrupted samples –
Learn how to denoise, learn how to denoise



Denoising autoencoder



Slides from Hugo Larochelle

http://info.usherbrooke.ca/hlarochelle/neural_networks/content.html

Variational Autoencoder

- Encoder $P(\mathbf{x}) = \int P(\mathbf{x}|\mathbf{z})P(\mathbf{z})d\mathbf{z}$

- $q(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu_{\theta}(\mathbf{x}), \sigma_{\theta}^2(\mathbf{x})\mathbf{I})$

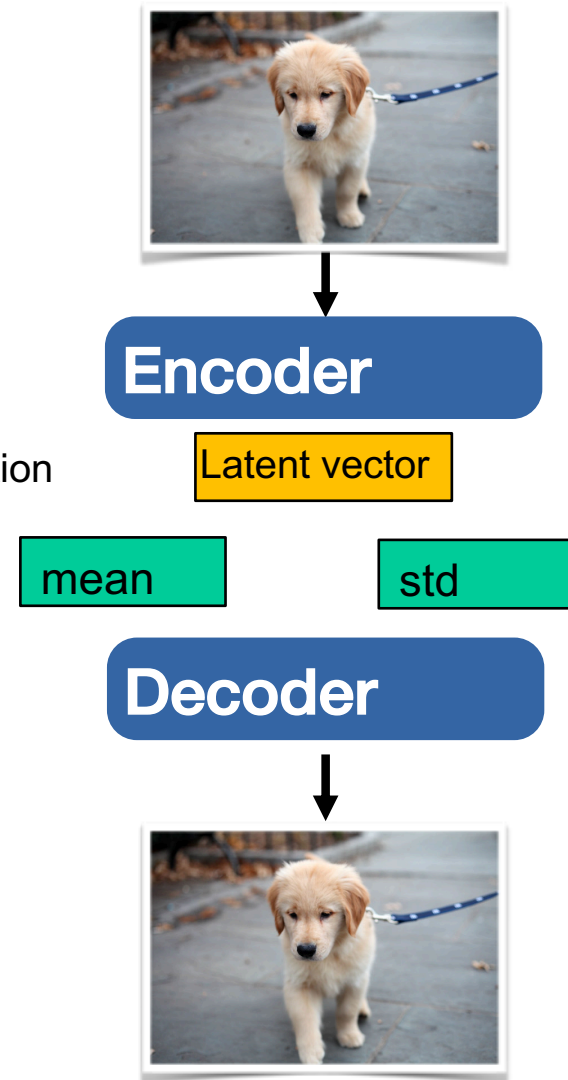
- Sampling $\mathbf{f} \sim q(\mathbf{z} | \mathbf{x})$

- Decoder

- $P(\mathbf{x} | \mathbf{f})$ Sample latent vector from distribution

- Approximately learns $P(\mathbf{x})$

- Variational lower bound



Variational Autoencoder

- Encoder
$$P(\mathbf{x}) = \int P(\mathbf{x}|\mathbf{z})P(\mathbf{z})d\mathbf{z}$$

- $q(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu_{\theta}(\mathbf{x}), \sigma_{\theta}^2(\mathbf{x})\mathbf{I})$

- Sampling $\mathbf{f} \sim q(\mathbf{z} | \mathbf{x})$

- Decoder

- $P(\mathbf{x} | \mathbf{f})$



Encoder

Latent vector

mean

std

Loss function

$$l_i(\theta, \phi) = -\mathbb{E}_{z \sim q_{\theta}(z|x_i)}[\log p_{\phi}(x_i | z)] + \text{KL}(q_{\theta}(z | x_i) || p(z))$$

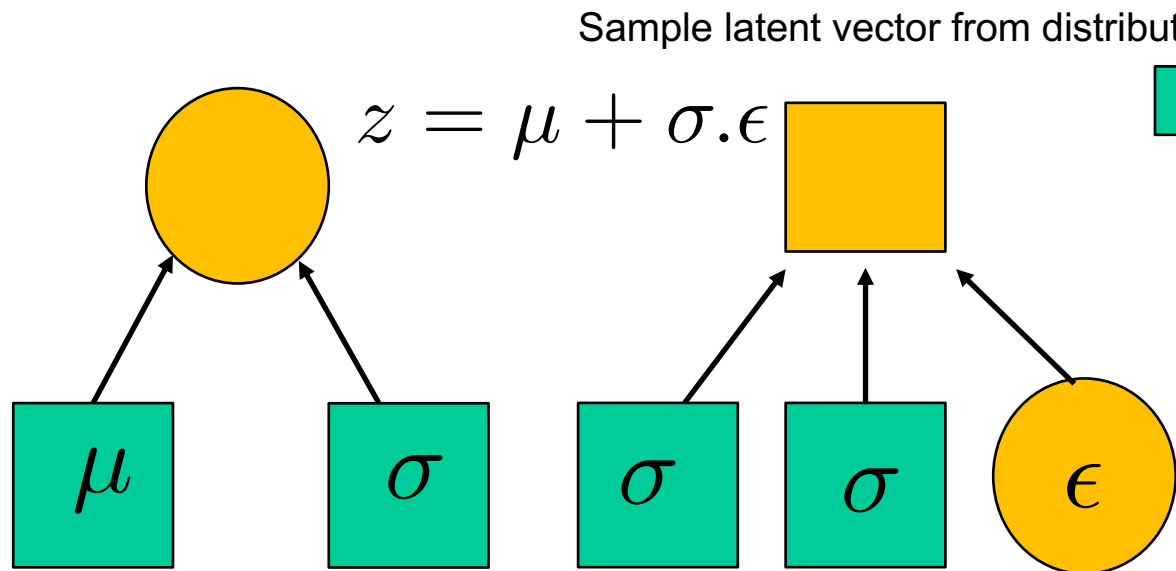
Reconstruction loss and KL divergence
 $p(z)$ is $N(0,1)$



Variational Autoencoder

Instead of sampling, from the distribution
To get the latent vector
use reparametrization trick

Make the latent vector mean
And add std and then corrupt std
with $N(0,1)$ gaussian



Encoder

Latent vector

mean

std

Decoder



VAE faces

[Alec Radford Face Manifold from variational autoencoder](#)

