

Deep Generative Models

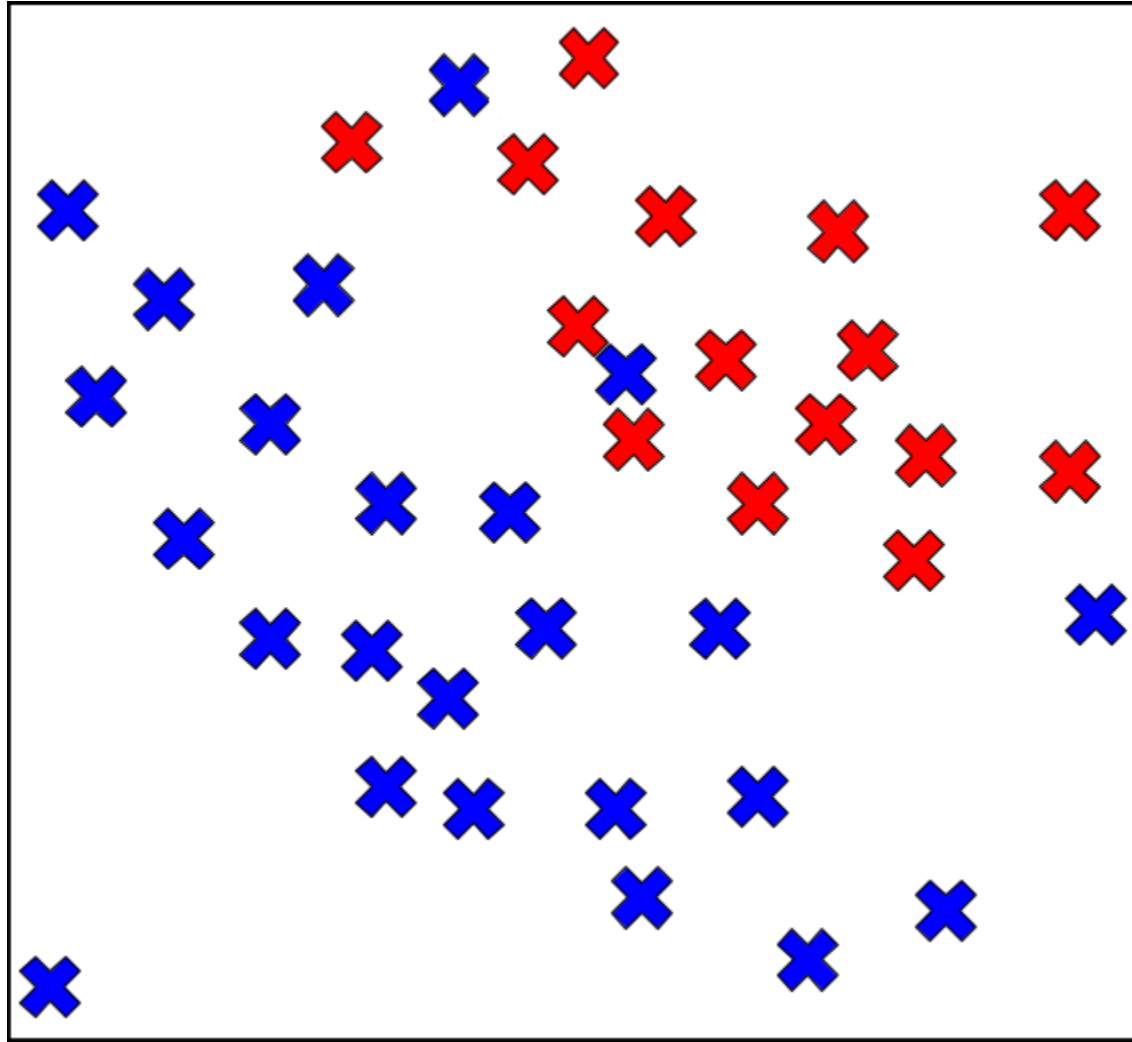
Jakub M. Tomczak
AMLAB, UvA

2017/11/29

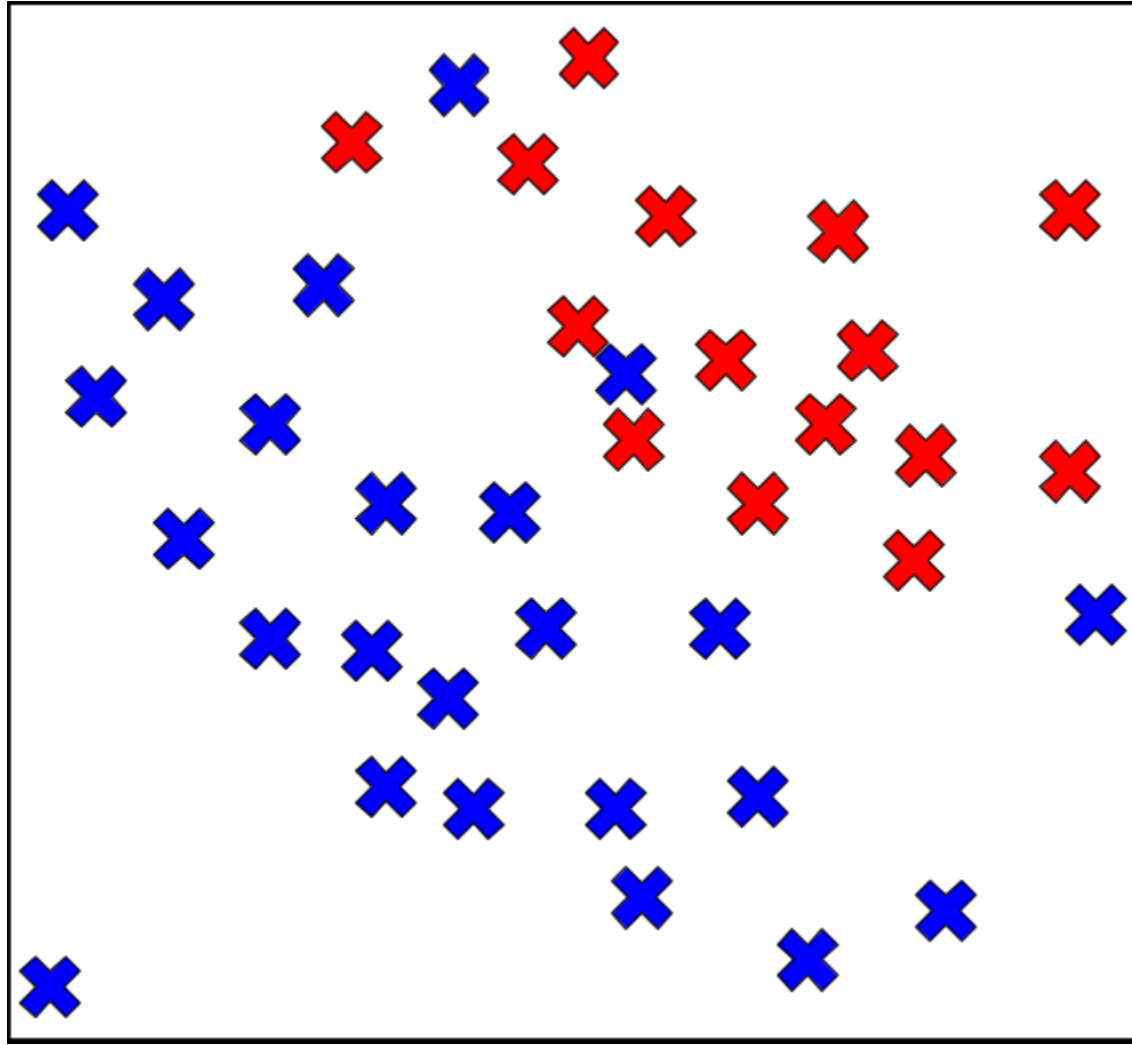


Do we need generative modeling?

Do we need generative modeling?

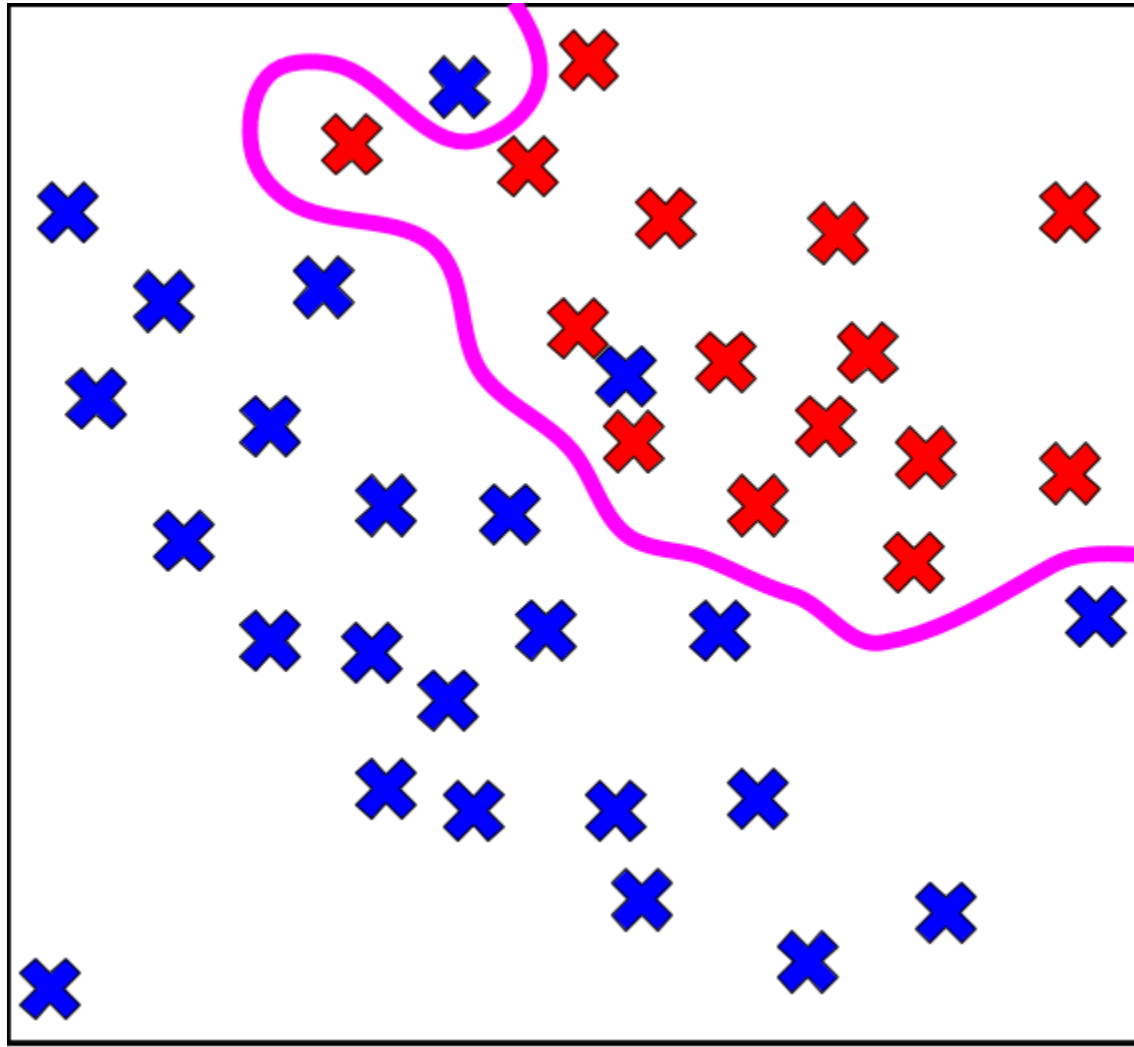


Do we need generative modeling?



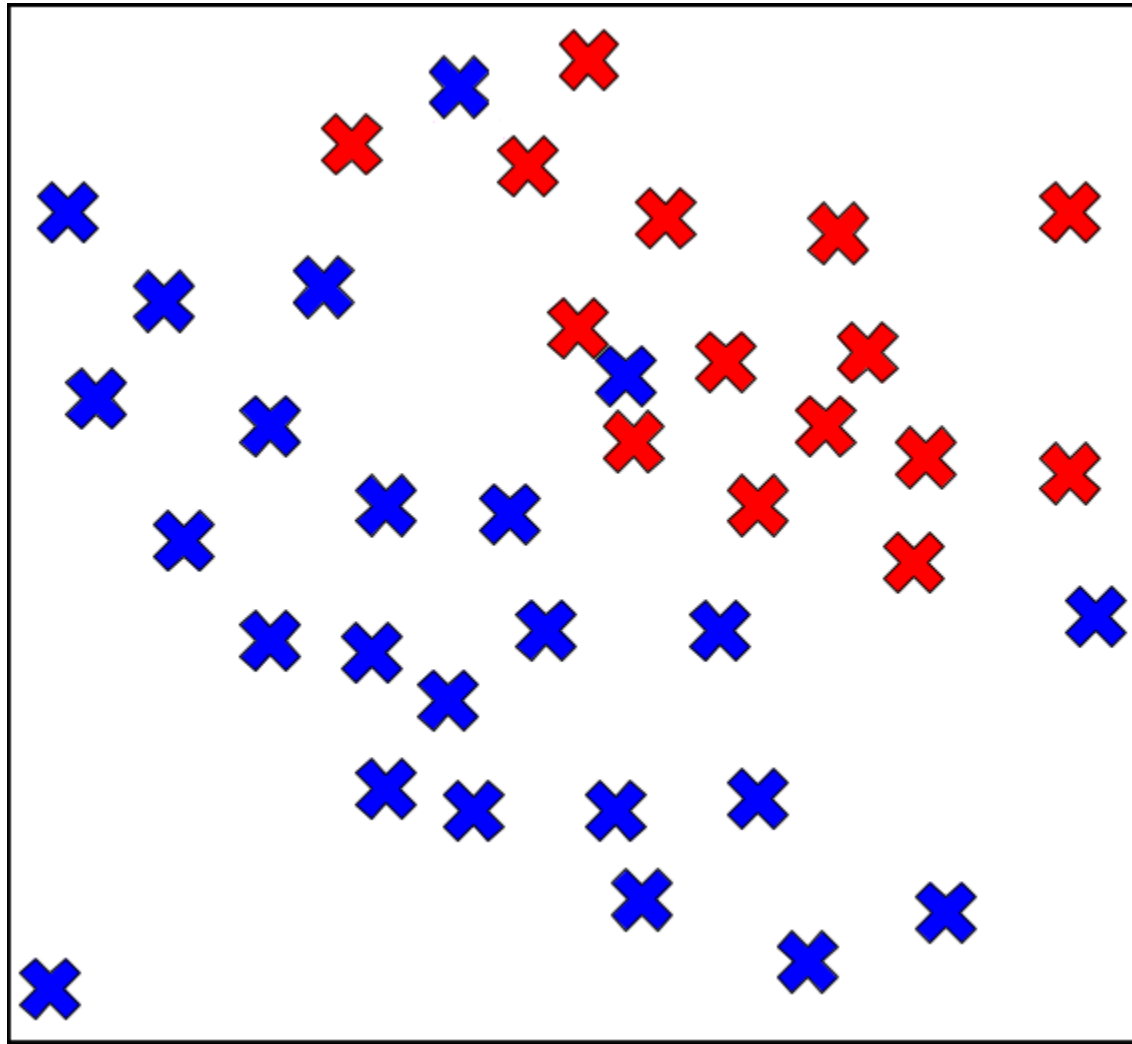
$$p_{\theta}(y|x)$$

Do we need generative modeling?



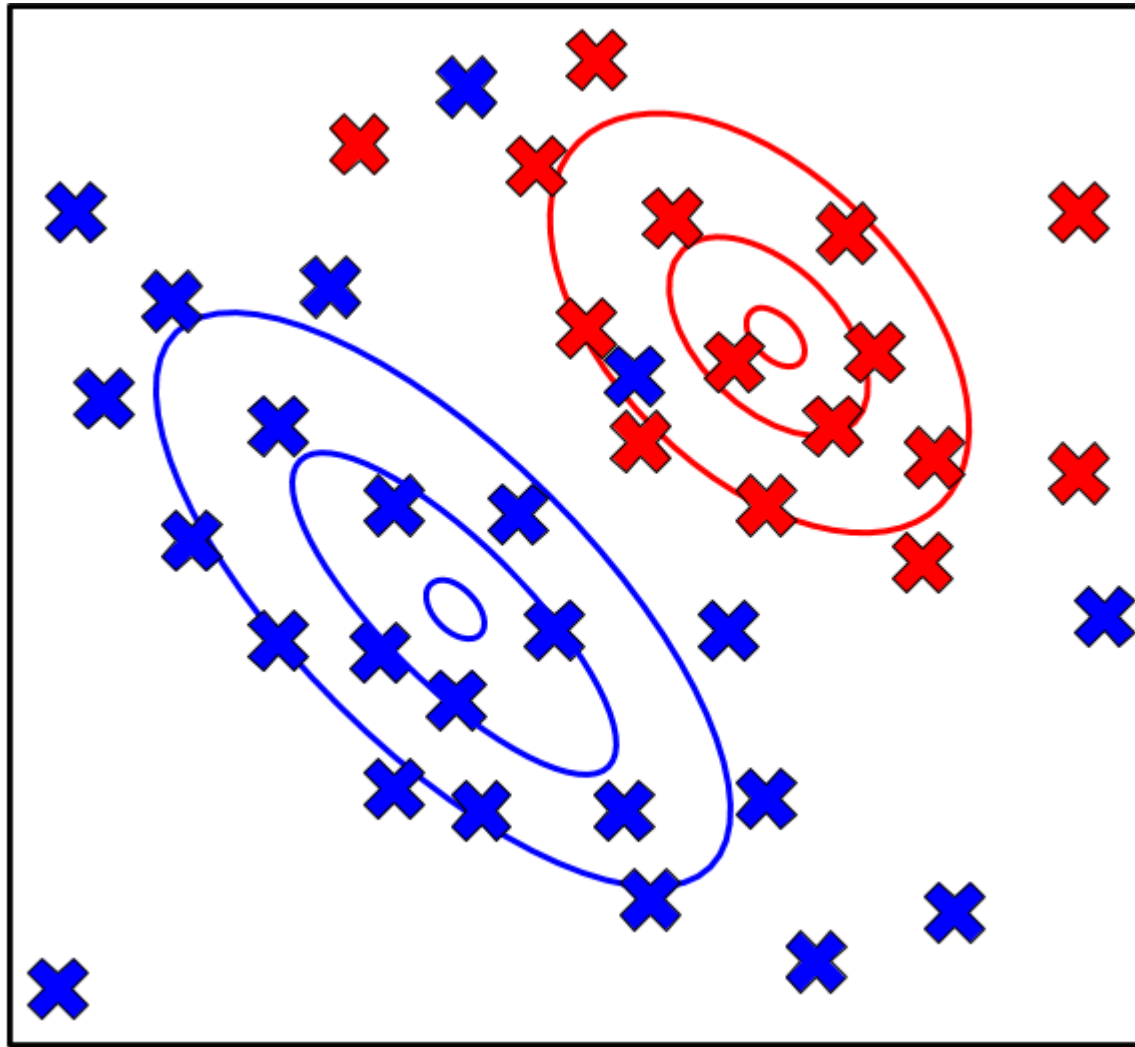
$$p_{\theta}(y|x)$$

Do we need generative modeling?



$$p_{\theta}(x, y) = p_{\theta}(y|x) p_{\theta}(x)$$

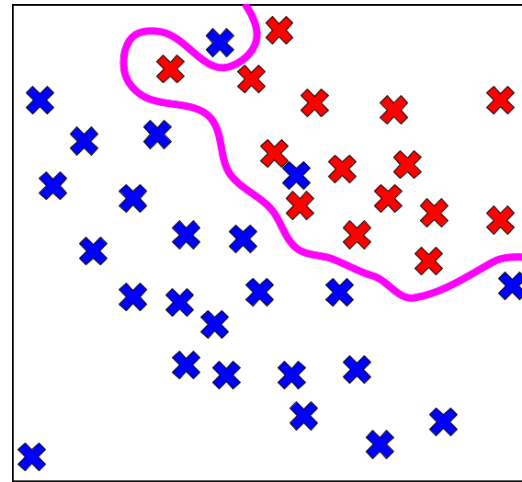
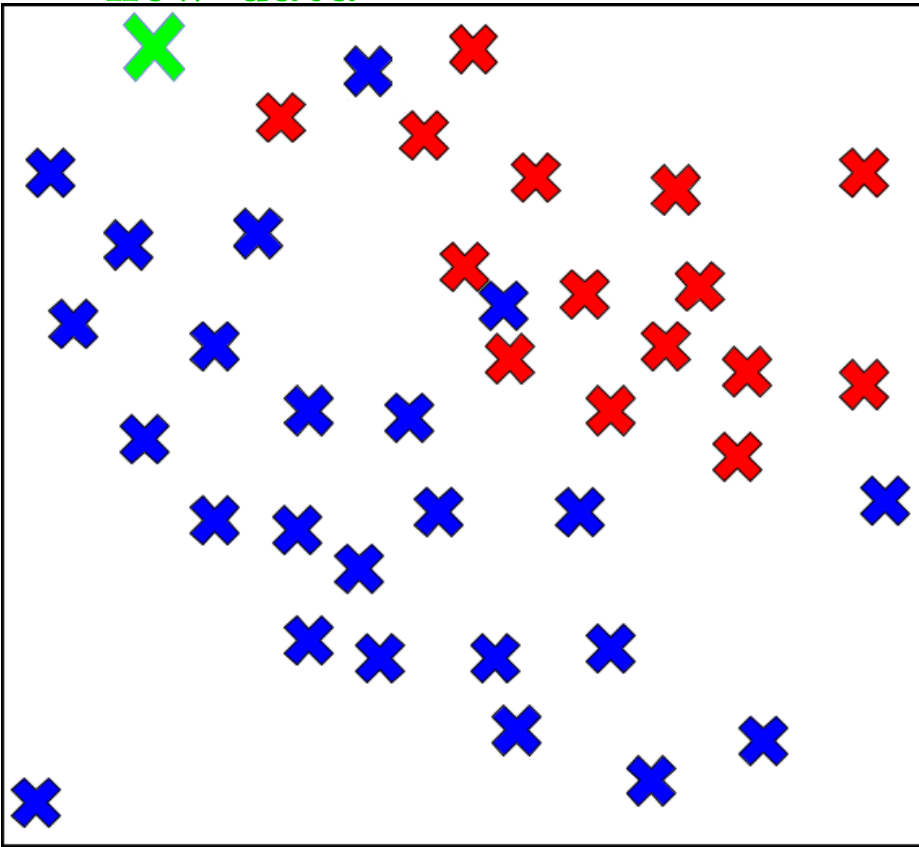
Do we need generative modeling?



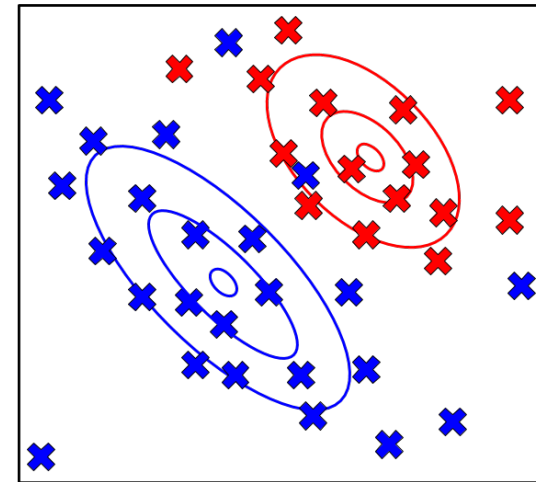
$$p_\theta(x, y) = p_\theta(y|x) p_\theta(x)$$

Do we need generative modeling?

new data



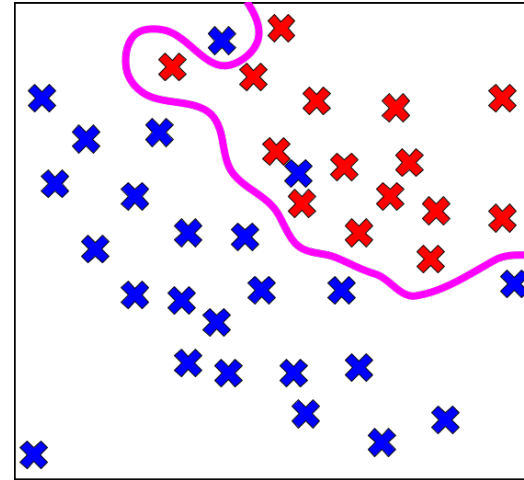
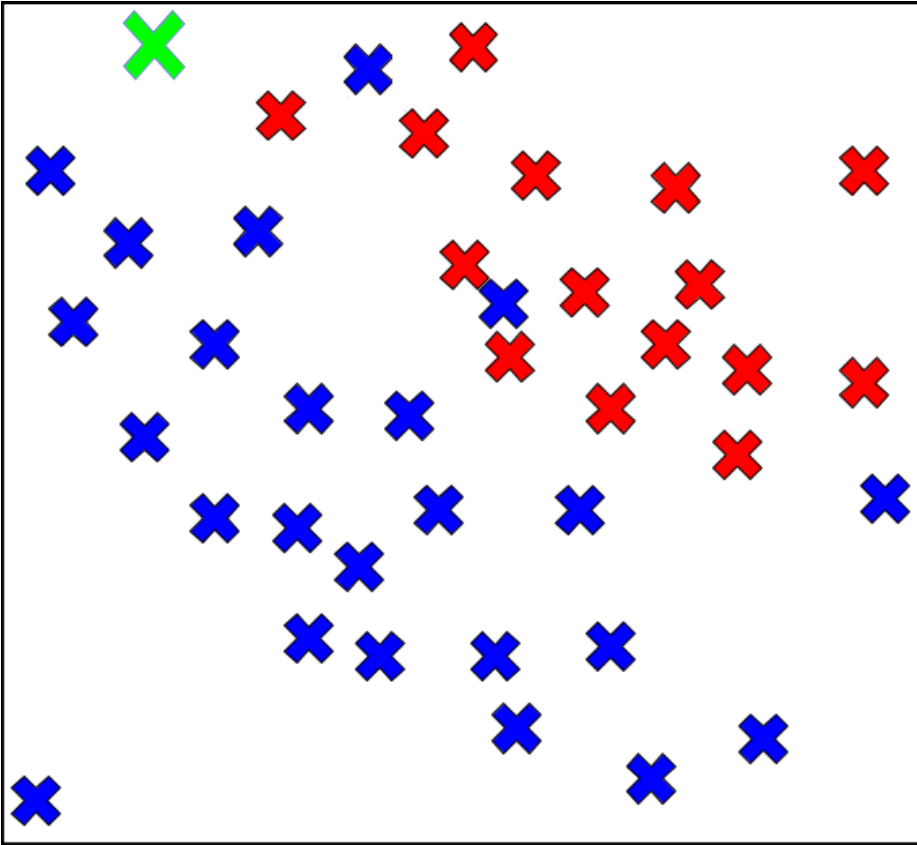
$$p_{\theta}(y|x)$$



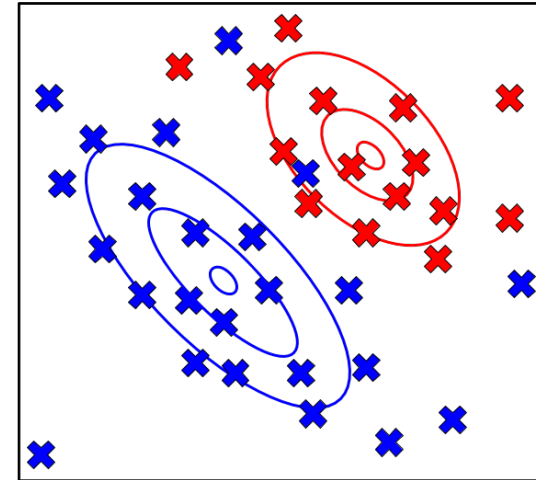
$$p_{\theta}(x, y)$$

Do we need generative modeling?

new data



$$p_{\theta}(y|x)$$



$$p_{\theta}(x, y)$$

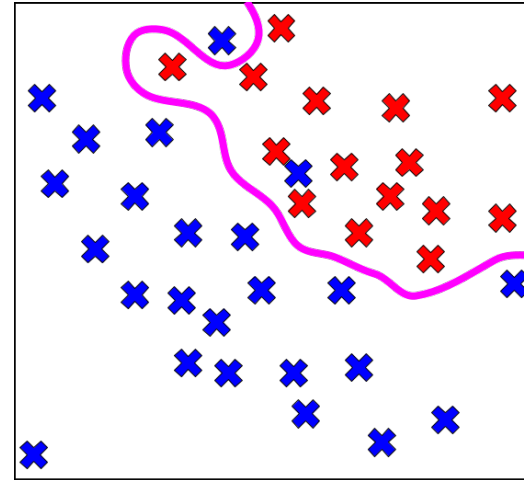
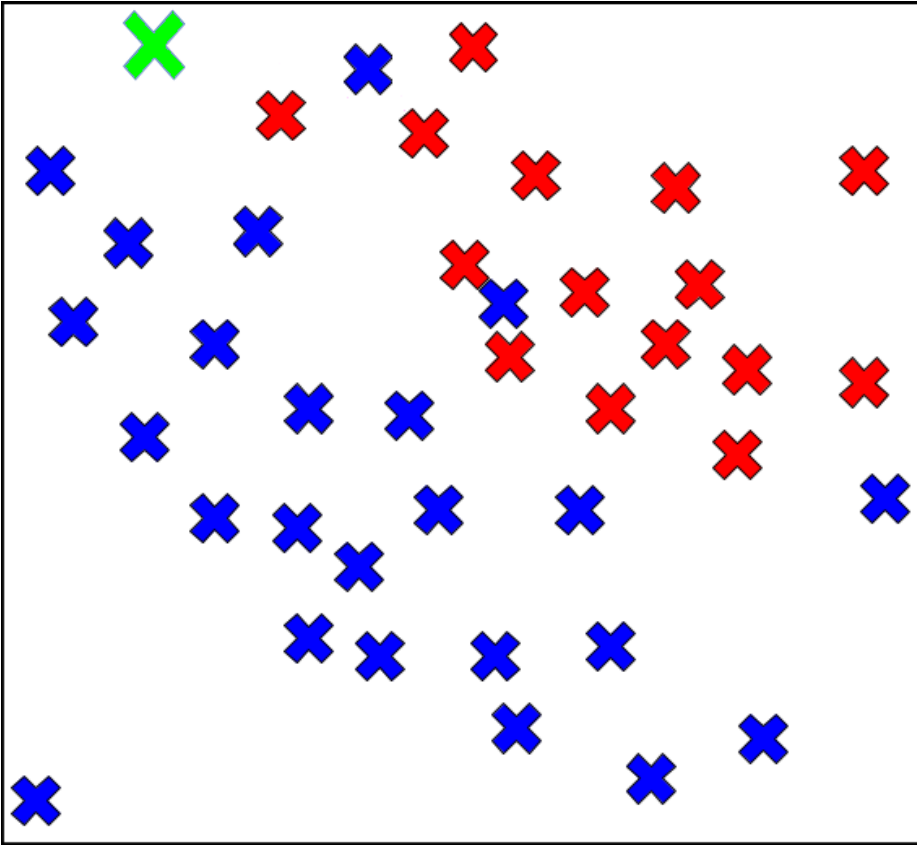
High probability
of the **blue** label.

=

Highly probable
decision!

Do we need generative modeling?

new data

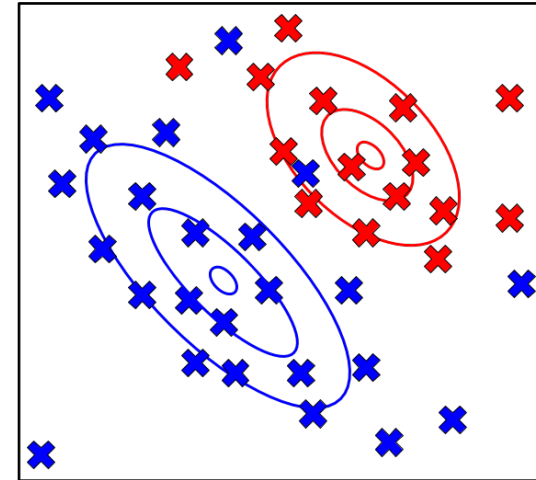


$$p_{\theta}(y|x)$$

High probability
of the **blue** label.

=

Highly probable
decision!



$$p_{\theta}(x, y)$$

High probability
of the **blue** label.

x

Low probability
of the **object**.

=

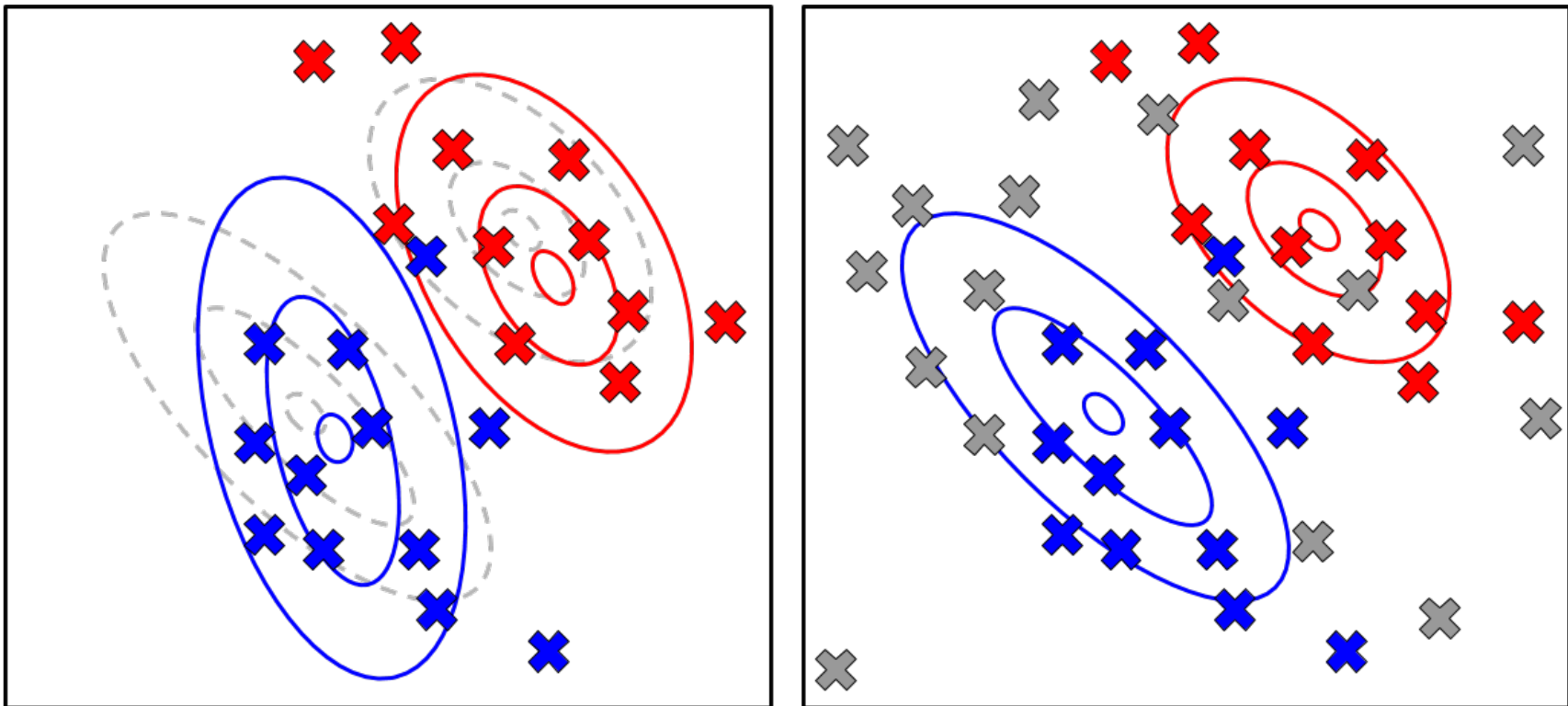
Uncertain
decision!

Generative Modeling

- Providing decision is not enough. *How to evaluate uncertainty? Distribution of y is only a part of the story.*
- Generalization problem. *Without knowing the distribution of \mathbf{x} how we can generalize to new data?*
- Understanding the problem is crucial (“**What I cannot create, I do not understand**”, Richard P. Feynman). *Properly modeling data is essential to make better decisions.*

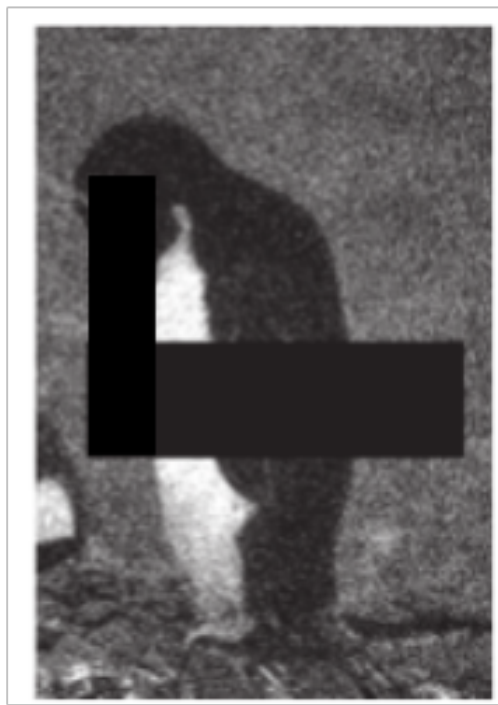
Generative Modeling

- Semi-supervised learning.
Use unlabeled data to train a better classifier.



Generative Modeling

- Handling missing or distorted data.
Reconstruct and/or denoise data.



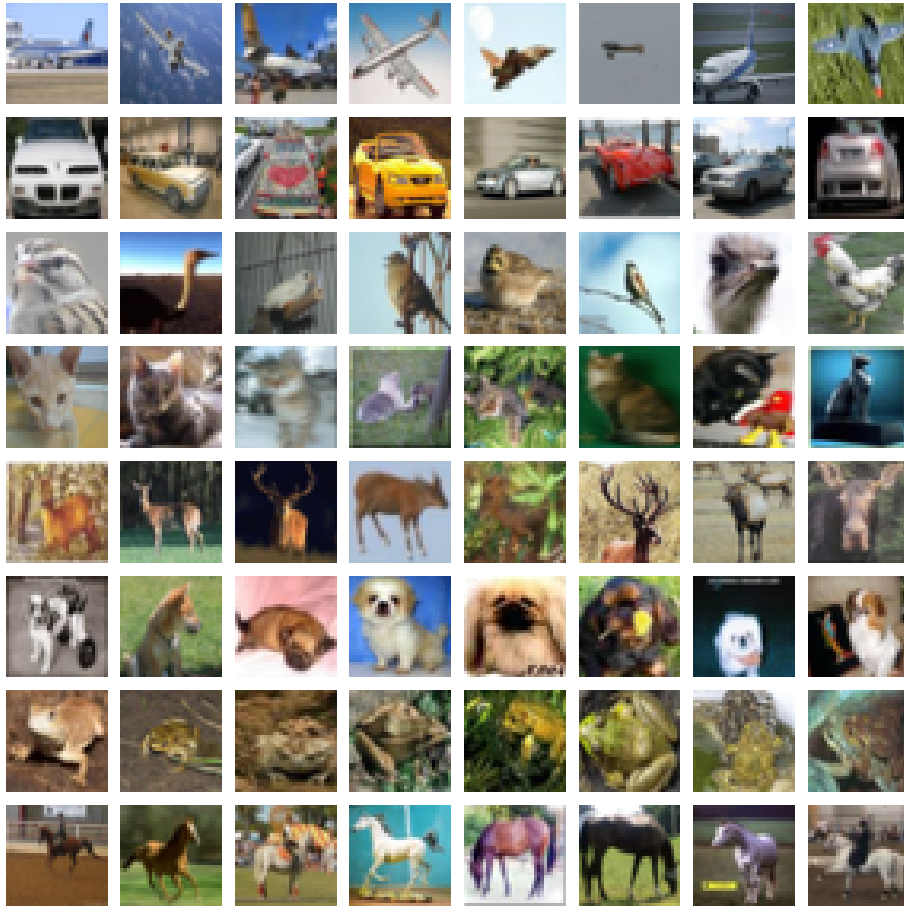
???



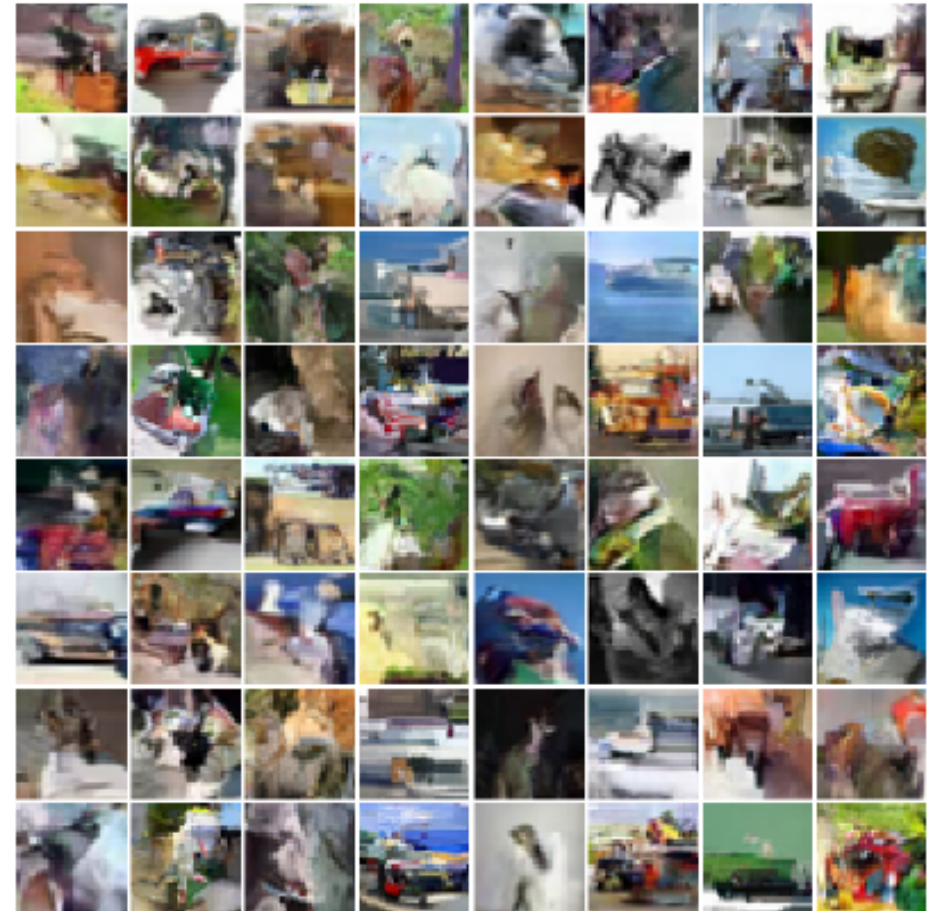
penguin!

Generative Modeling

Image generation



Real



Generated

Generative Modeling

Sequence generation

he had been unable to conceal the fact that there was a logical explanation for his inability to alter the fact that they were supposed to be on the other side of the house .

with a variety of pots strewn scattered across the vast expanse of the high ceiling , a vase of colorful flowers adorned the tops of the rose petals littered the floor and littered the floor .

atop the circular dais perched atop the gleaming marble columns began to emerge from atop the stone dais, perched atop the dais .

Generated

How to formulate a generative model?

Modeling in a high-dimensional space is difficult.



How to formulate a generative model?

Modeling in a high-dimensional space is difficult.



How to formulate a generative model?

Modeling in a high-dimensional space is difficult.

→ modeling all dependencies among pixels.

$$p(\mathbf{x}) = \prod_{d=1}^c \psi_d(\mathbf{x}_d)$$

How to formulate a generative model?

Modeling in a high-dimensional space is difficult.

→ modeling all dependencies among pixels.

$$p(x) = \prod_{d=1}^c \psi_c(x_c)$$

very inefficient!

How to formulate a generative model?

Modeling in a high-dimensional space is difficult.

→ modeling all dependencies among pixels.

$$p(x) = \prod_{d=1}^c \psi_c(x_c)$$

very inefficient!

A possible **solution**? → **Models with latent variables**

Boltzmann distribution

A joint distribution of random variables:

$$p(\mathbf{x}) = \frac{1}{Z} \exp\{-E(\mathbf{x})\}$$

Boltzmann distribution

A joint distribution of random variables:

$$p(\mathbf{x}) = \frac{1}{Z} \exp\{-E(\mathbf{x})\}$$

Energy function

Boltzmann distribution

A joint distribution of random variables:

$$p(\mathbf{x}) = \frac{1}{Z} \exp\{-E(\mathbf{x})\}$$

Partition function

Energy function

$$Z = \sum_{\mathbf{x}} \exp\{-E(\mathbf{x})\}$$

Boltzmann distribution

A joint distribution of random variables:

$$p(\mathbf{x}) = \frac{1}{Z} \exp\{-E(\mathbf{x})\}$$

It is called the **Boltzmann** (or **Gibbs**) **distribution**.

Boltzmann distribution

A joint distribution of random variables:

$$p(\mathbf{x}) = \frac{1}{Z} \exp\{-E(\mathbf{x})\}$$

Advantages:

- Joint distribution is described by $E(\mathbf{x})$.
- Well-studied distribution in statistical physics.

Boltzmann distribution

A joint distribution of random variables:

$$p(\mathbf{x}) = \frac{1}{Z} \exp\{-E(\mathbf{x})\}$$

Advantages:

- Joint distribution is described by $E(\mathbf{x})$.
- Well-studied distribution in statistical physics.

Drawbacks:

- Calculation of partition function requires enumeration of all \mathbf{x} , e.g., if \mathbf{x} 's are binary $\rightarrow 2^n$ combinations.

Boltzmann Machines

Let us consider the following energy function:

$$E(\mathbf{x}) = -\mathbf{x}^\top \mathbf{W} \mathbf{x} - \mathbf{b}^\top \mathbf{x}$$

where \mathbf{x} are binary.

Boltzmann Machines

Let us consider the following energy function:

$$E(\mathbf{x}) = -\mathbf{x}^\top \mathbf{W} \mathbf{x} - \mathbf{b}^\top \mathbf{x}$$

where \mathbf{x} are binary.

first-order term
(models “*prior*”)

Boltzmann Machines

Let us consider the following energy function:

$$E(\mathbf{x}) = -\mathbf{x}^\top \mathbf{W} \mathbf{x} - \mathbf{b}^\top \mathbf{x}$$

where \mathbf{x} are binary.

second-order term
(models correlations)

first-order term
(models “*prior*”)

Boltzmann Machines

Let us consider the following energy function:

$$E(\mathbf{x}) = -\mathbf{x}^\top \mathbf{W} \mathbf{x} - \mathbf{b}^\top \mathbf{x}$$

where \mathbf{x} are binary.

second-order term
(models correlations)

first-order term
(models “*prior*”)

Very complicated...

Can we do better?

Boltzmann Machines

Let us consider the following energy function:

$$E(\mathbf{x}) = -\mathbf{x}^\top \mathbf{W} \mathbf{x} - \mathbf{b}^\top \mathbf{x}$$

where \mathbf{x} are binary.

second-order term
(models correlations)

first-order term
(models “*prior*”)

Very complicated...

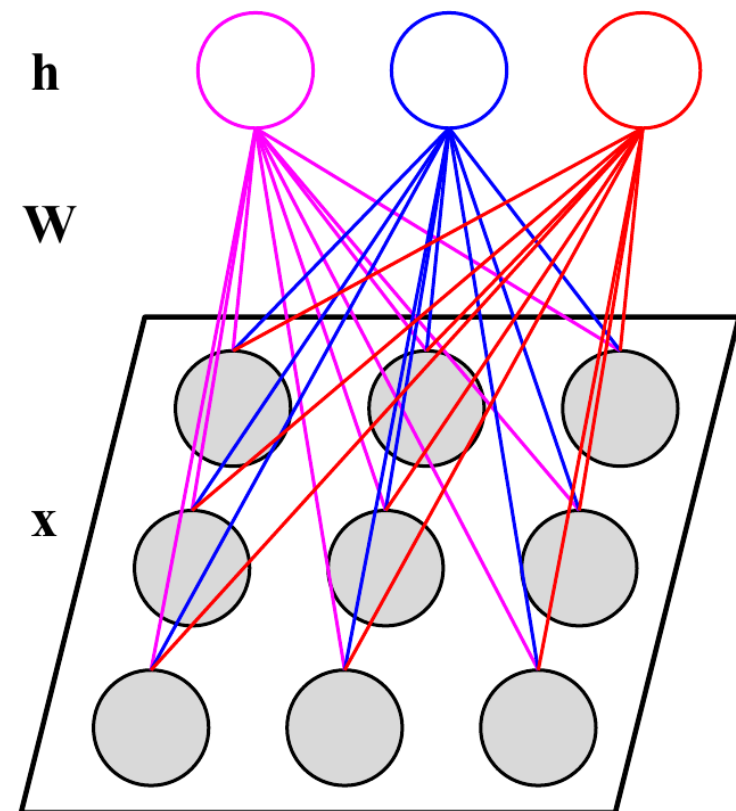
Can we do better? **Latent variables!**

Restricted Boltzmann Machines

Let us consider the following energy function:

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

where \mathbf{x} and \mathbf{h} are binary.

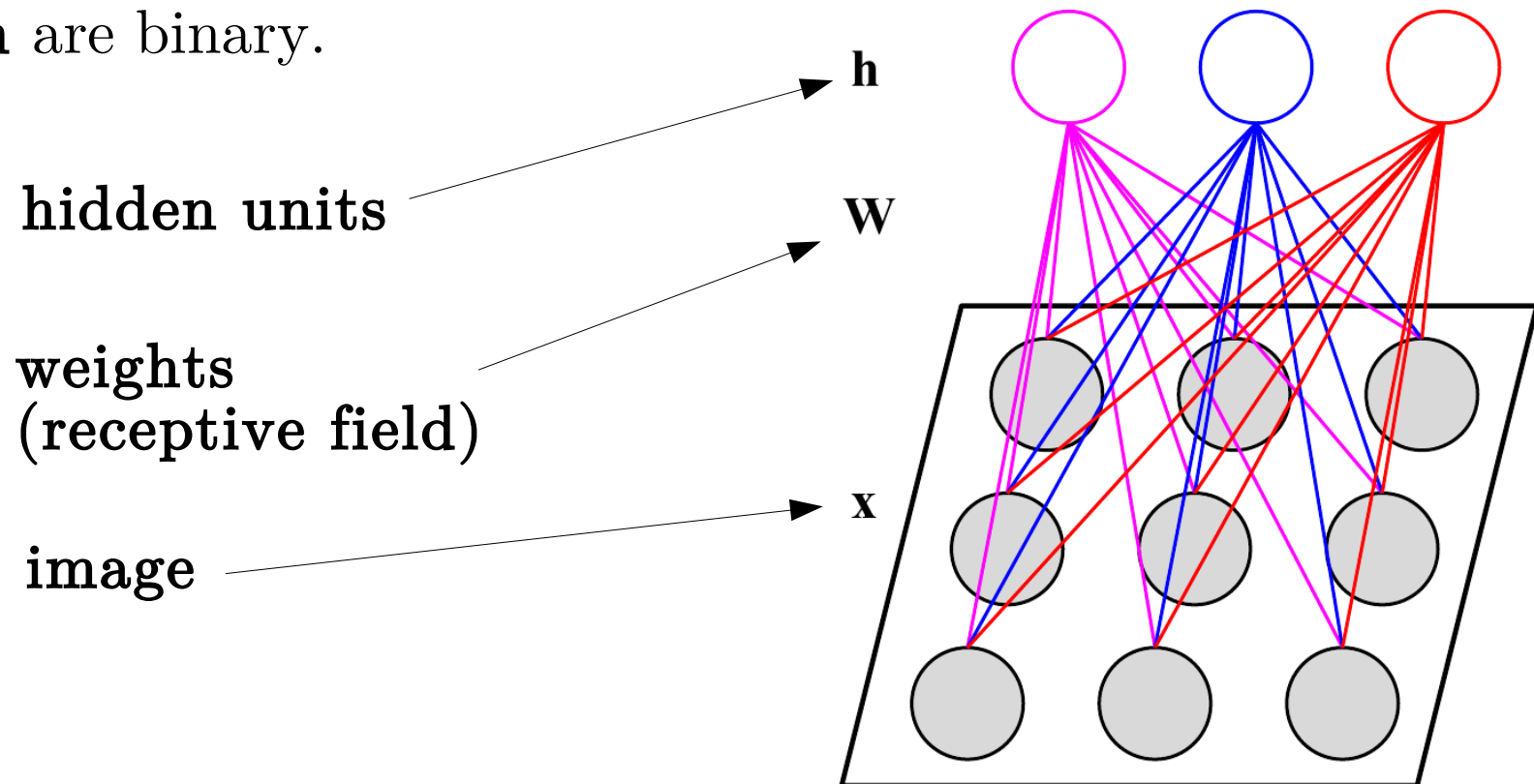


Restricted Boltzmann Machines

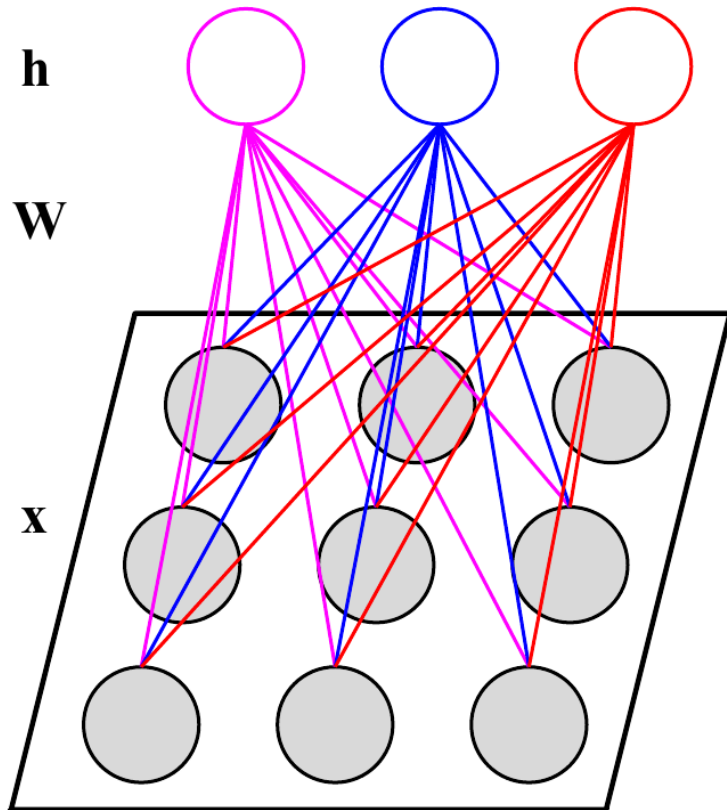
Let us consider the following energy function:

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

where \mathbf{x} and \mathbf{h} are binary.



Restricted Boltzmann Machines



Conditional dependencies:

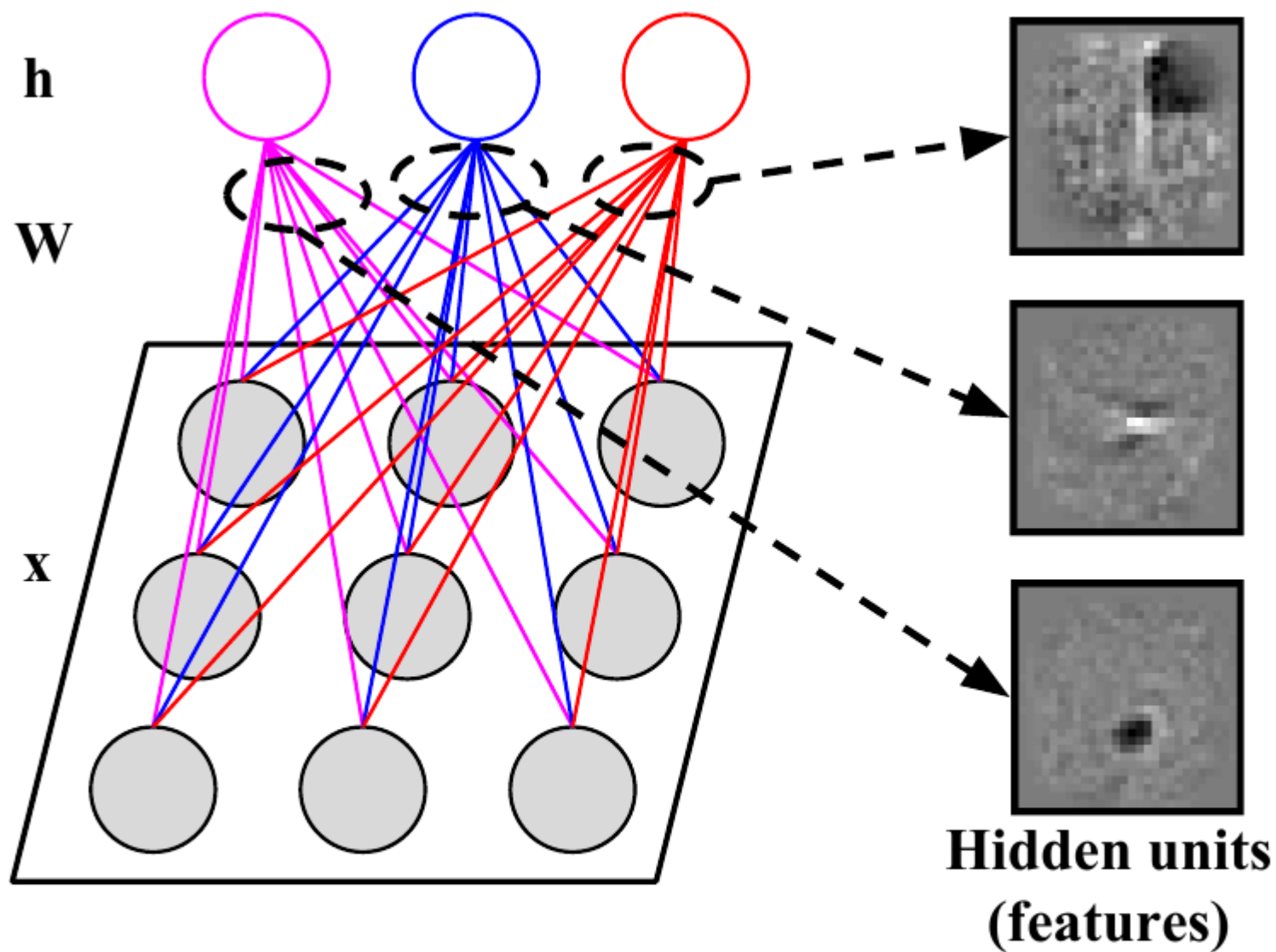
$$p(\mathbf{h}|\mathbf{x}, \theta) = \prod_j p(h_j|\mathbf{x}, \theta)$$

$$p(\mathbf{x}|\mathbf{h}, \theta) = \prod_i p(x_i|\mathbf{h}, \theta)$$

where

$$p(h_j = 1|\mathbf{x}, \theta) = \frac{1}{1 + \exp(-(\mathbf{W}_{\cdot j}^\top)\mathbf{x} - c_j)}$$

$$p(x_i = 1|\mathbf{h}, \theta) = \frac{1}{1 + \exp(-\mathbf{W}_{\cdot i}\mathbf{h} - b_j)}$$



$$P(h_j = 1 | \mathbf{x}, \theta) = \frac{1}{1 + \exp(-(\mathbf{W}_{\cdot j})^\top \mathbf{x} - c_j)}$$

RBM: Training

We train RBMs using the **log-likelihood function**:

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{n=1}^N \log p(\mathbf{x}_n | \theta)$$

where

$$\begin{aligned} p(\mathbf{x}_n | \theta) &= \frac{1}{Z} \exp(-FE(\mathbf{x}_n)) \\ &= \frac{1}{Z} \exp \left(\mathbf{b}^\top \mathbf{x}_n + \sum_j \log(1 + \exp\{(\mathbf{W}_{\cdot j})^\top \mathbf{x}_n + c_j\}) \right) \end{aligned}$$

RBM: Training

We train RBMs using the **log-likelihood function**:

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{n=1}^N \log p(\mathbf{x}_n | \theta)$$

where

$$\begin{aligned} p(\mathbf{x}_n | \theta) &= \frac{1}{Z} \exp(-FE(\mathbf{x}_n)) \\ &= \frac{1}{Z} \exp\left(\mathbf{b}^\top \mathbf{x}_n + \underbrace{\sum_j \log(1 + \exp\{(\mathbf{W}_{\cdot j})^\top \mathbf{x}_n + c_j\})}_{\text{softplus}}\right) \end{aligned}$$

Free energy

RBM: Training

Let us calculate gradient of the LL wrt parameters:

$$\begin{aligned}\frac{\partial \log p(\mathbf{x}_n|\theta)}{\partial \theta} &= -\frac{\partial}{\partial \theta} F E(\mathbf{x}_n) - \frac{\partial}{\partial \theta} \log Z \\ &= -\sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{x}, \theta) \frac{\partial}{\partial \theta} E(\mathbf{x}_n, \mathbf{h}|\theta) + \sum_{\tilde{\mathbf{x}}, \mathbf{h}} p(\tilde{\mathbf{x}}, \mathbf{h}|\theta) \frac{\partial}{\partial \theta} E(\tilde{\mathbf{x}}, \mathbf{h}|\theta)\end{aligned}$$

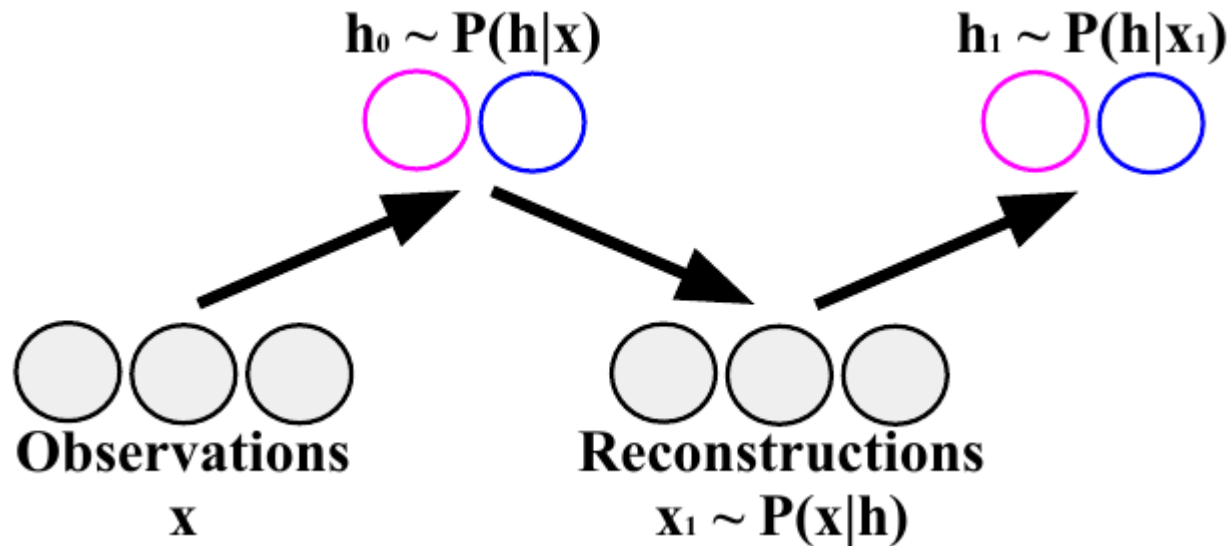
RBM: Training

Let us calculate gradient of the LL wrt parameters:

$$\begin{aligned}\frac{\partial \log p(\mathbf{x}_n|\theta)}{\partial \theta} &= -\frac{\partial}{\partial \theta} E(\mathbf{x}_n) - \frac{\partial}{\partial \theta} \log Z \\ &= -\underbrace{\sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{x}, \theta) \frac{\partial}{\partial \theta} E(\mathbf{x}_n, \mathbf{h}|\theta)}_{\text{easy}} + \underbrace{\sum_{\tilde{\mathbf{x}}, \mathbf{h}} p(\tilde{\mathbf{x}}, \mathbf{h}|\theta) \frac{\partial}{\partial \theta} E(\tilde{\mathbf{x}}, \mathbf{h}|\theta)}_{\text{difficult}}\end{aligned}$$

It requires application of an approximate inference (e.g., MCMC).

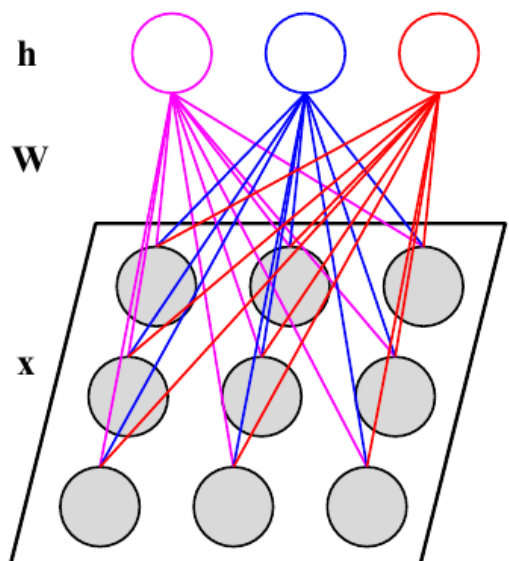
RBM: Contrastive Divergence



Apply k steps of **Gibbs sampler** and approximate the gradient as follows:

$$\frac{\partial \log p(\mathbf{x}_n | \theta)}{\partial \theta} = -\frac{\partial}{\partial \theta} E(\mathbf{x}_n, \mathbf{h}_0 | \theta) + \frac{\partial}{\partial \theta} E(\mathbf{x}_k, \mathbf{h}_k | \theta)$$

RBM with Gaussian inputs



- **Observables:** $\mathbf{x} \in \mathbb{R}^D$
- **Hiddens:** $\mathbf{h} \in \{0, 1\}^M$
- Good for modelling **natural images**.

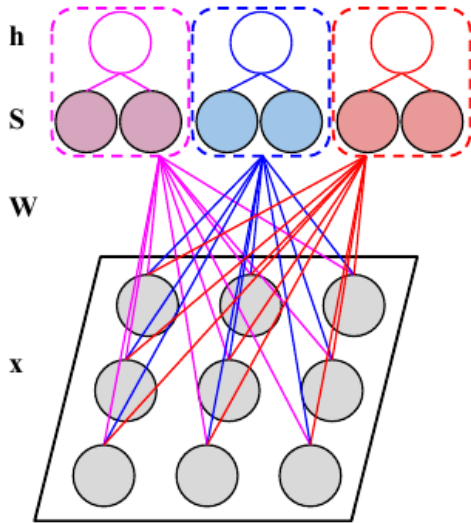
Energy for a joint configuration (\mathbf{x}, \mathbf{h}) :

$$E(\mathbf{x}, \mathbf{h}|\theta) = \sum_i \sum_j W_{ij} h_j \frac{x_i}{\sigma_i} + \sum_i \frac{(x_i - b_i)^2}{2\sigma_i^2} + \sum_j c_j h_j$$

Then:

$$P(\mathbf{x}|\mathbf{h}, \theta) = \prod_i \mathcal{N}(b_i + \sum_j W_{ij} h_j, \sigma_i^2)$$

Subspace RBM



- **Observables:** $\mathbf{x} \in \{0, 1\}^D$,
hidden variables: $\mathbf{h} \in \{0, 1\}^M$,
 $\mathbf{S} \in \{0, 1\}^{M \times K}$
- Aim: **features invariant to transformations.**

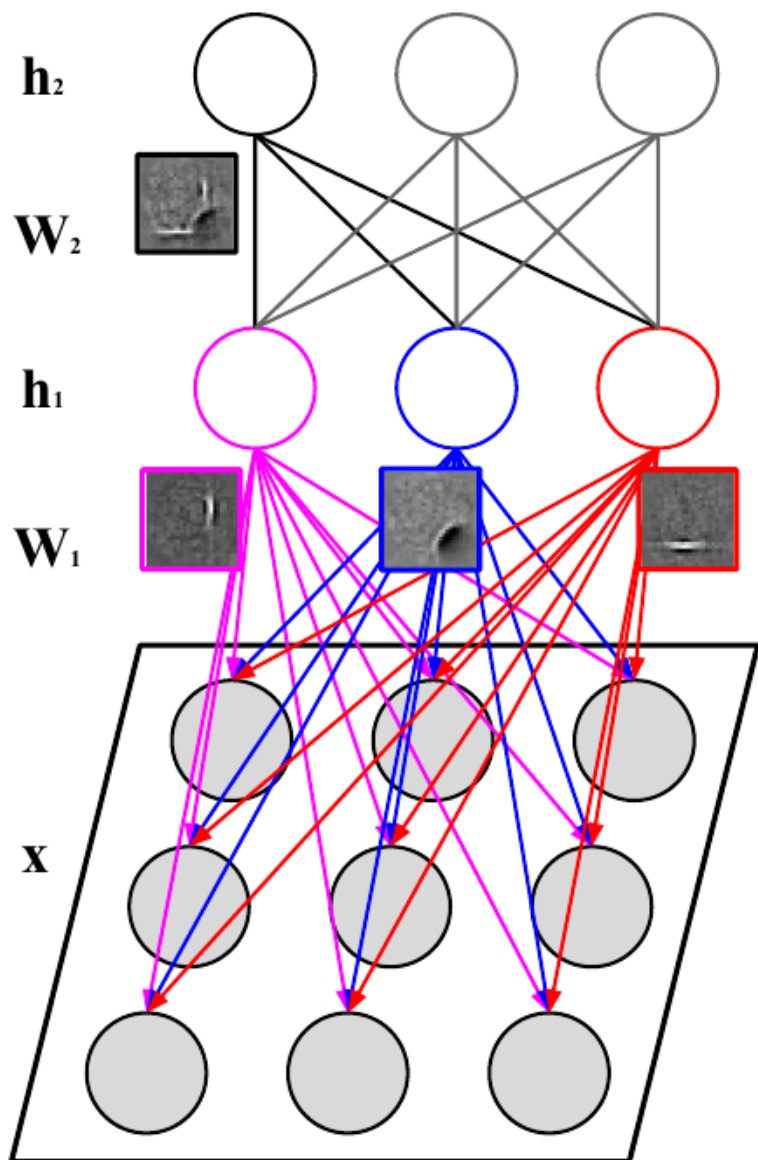
Energy for a joint configuration $(\mathbf{x}, \mathbf{h}, \mathbf{S})$:

$$E(\mathbf{x}, \mathbf{h}, \mathbf{S} | \theta) = \sum_{i,j,k} W_{ijk} x_i h_j S_{jk} + \mathbf{b}^\top \mathbf{x} + \mathbf{c}^\top \mathbf{h} + \sum_j h_j \sum_k A_{jk} S_{jk}$$

Then:

$$P(\mathbf{x}) \propto \exp(\mathbf{b}^\top \mathbf{x}) \prod_j \left[2^K + \exp(c_j) \prod_k \left(1 + \exp\left(\sum_i W_{ijk} x_i + A_{jk} \right) \right) \right]$$

Deep Belief Network



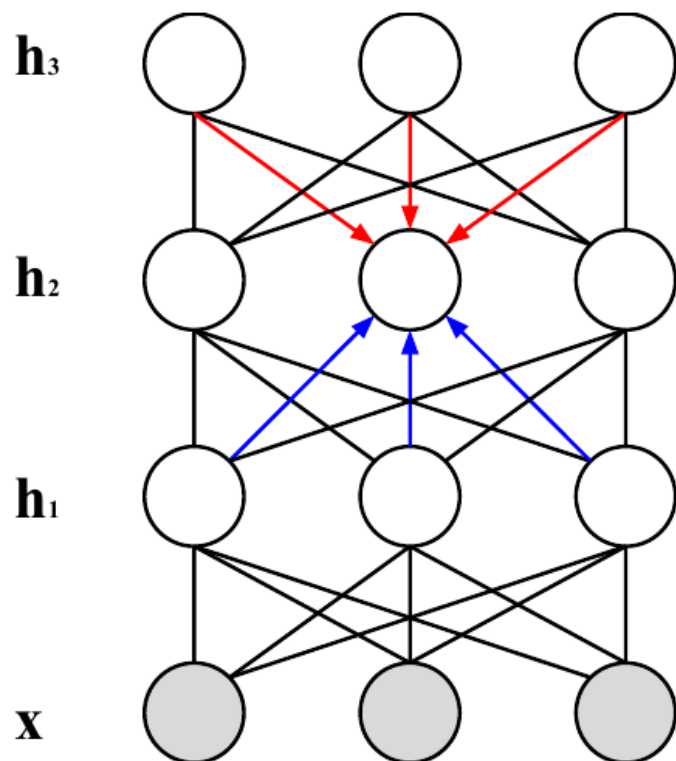
- RBM (or its extension) can be treated as a **building block**.
- Stacking RBMs leads to a **deep belief network**:

$$P(\mathbf{x}, \mathbf{h}_1, \mathbf{h}_2) = \underbrace{P(\mathbf{x}|\mathbf{h}_1)}_{\text{from RBM}} \underbrace{P(\mathbf{h}_1, \mathbf{h}_2)}_{\text{RBM}}$$

Deep Boltzmann Machines

Energy function:

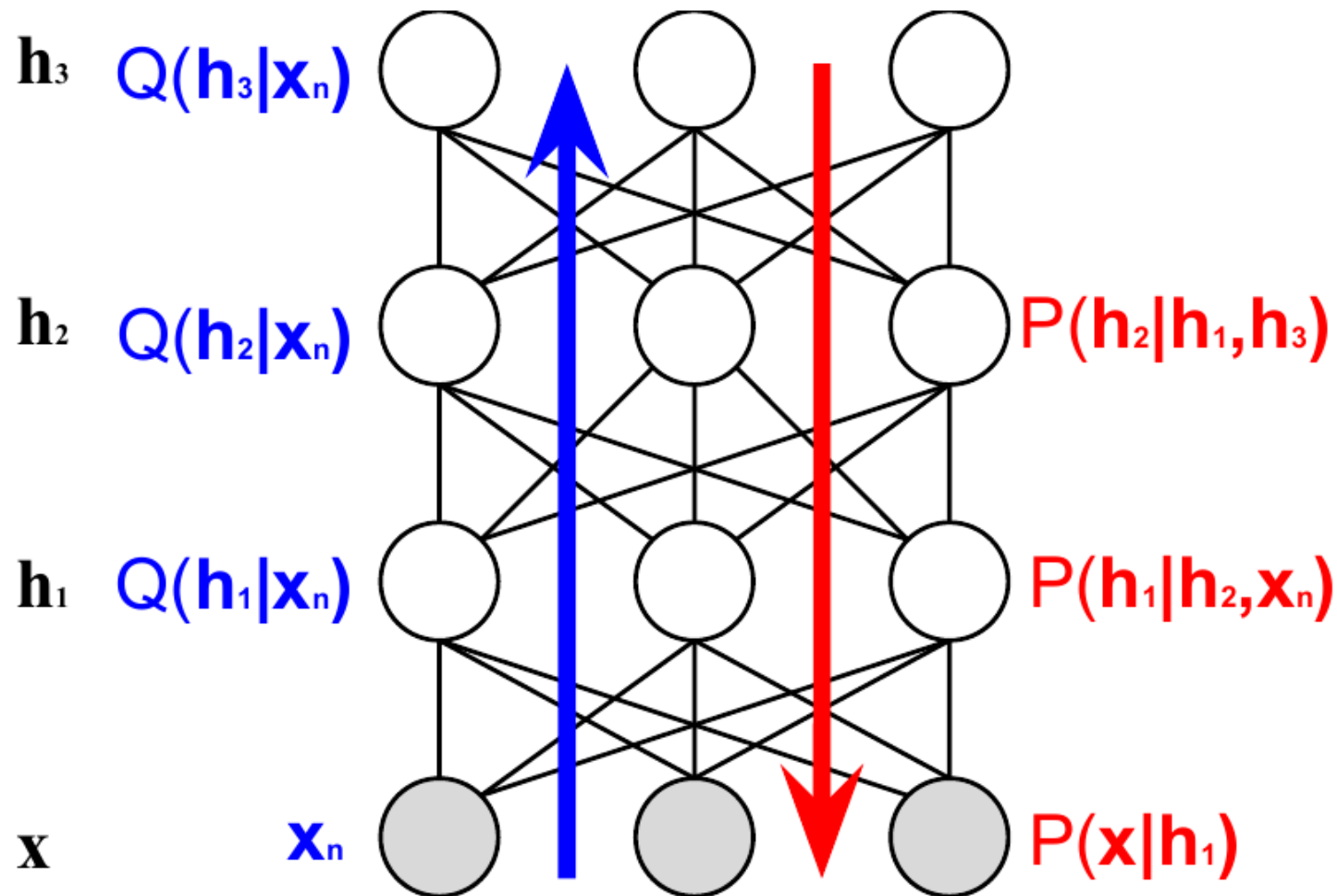
$$E(\mathbf{x}, \mathbf{h}_1, \mathbf{h}_2 | \theta) = \mathbf{x}^\top \mathbf{W}_1 \mathbf{h}_1 + \mathbf{h}_1^\top \mathbf{W}_2 \mathbf{h}_2 + \mathbf{h}_2^\top \mathbf{W}_3 \mathbf{h}_3$$



- Joint distribution – **Boltzmann (Gibbs) distribution**.
- All connections are **undirected**.
- Two types of relationships: (**top-down** and **bottom-up**):

$$P(h_2^{(k)} = 1 | \mathbf{h}_1, \mathbf{h}_3) = \text{sigm}\left(\sum_i W_1^{(jk)} h_1^{(j)} + \sum_l W_3^{(kl)} h_3^{(l)}\right)$$

Training DBM

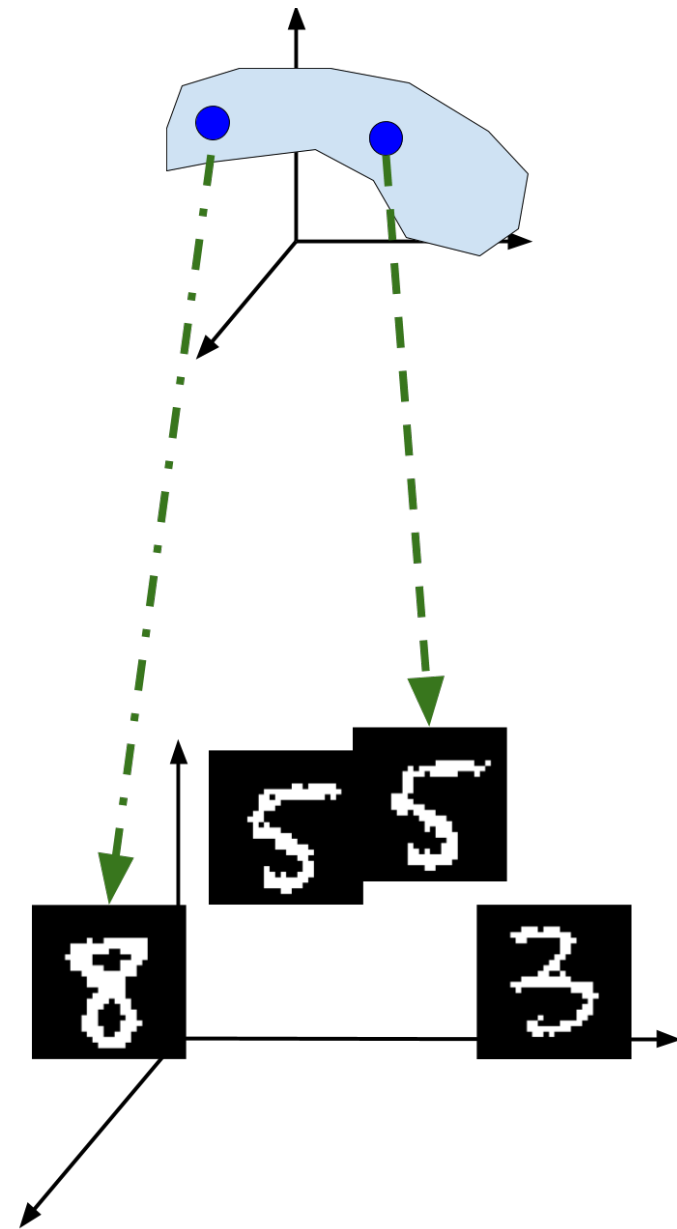


Calculating gradients is **intractable**, thus, **variational methods** (**mean-field**) and **sampling methods** are utilized.

Latent Variable Models

- Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$



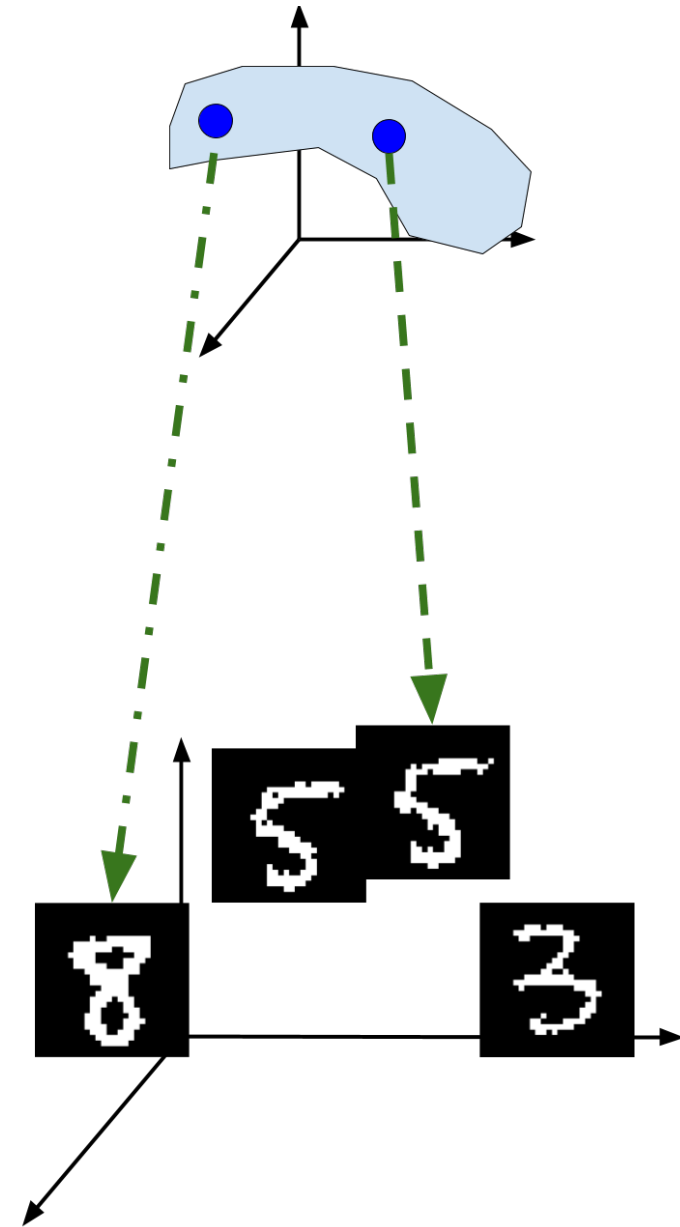
Latent Variable Models

- Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$

First sample \mathbf{z} .

Second, sample \mathbf{x} for given \mathbf{z} .



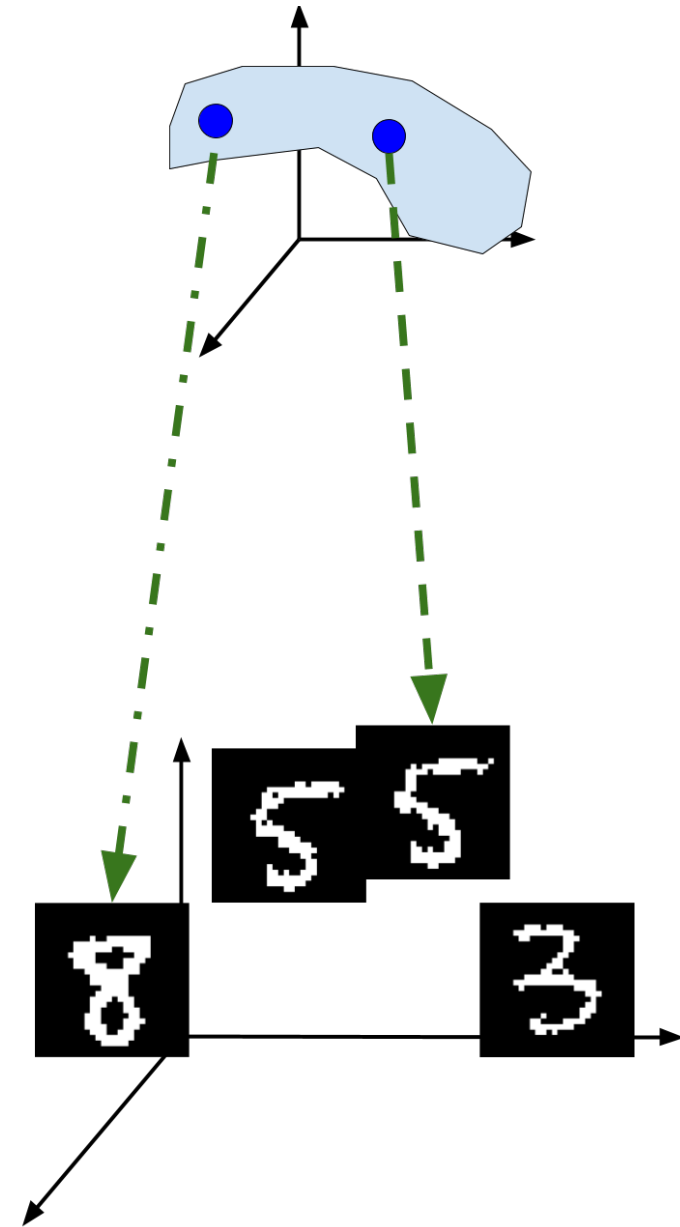
Latent Variable Models

- Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$

First sample \mathbf{z} .

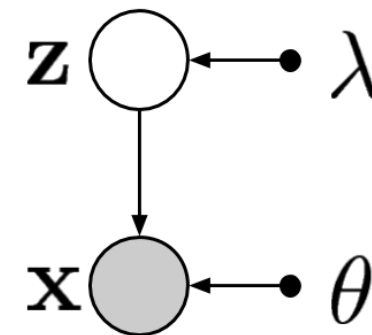
Second, sample \mathbf{x} for given \mathbf{z} .



Latent Variable Models

- Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$

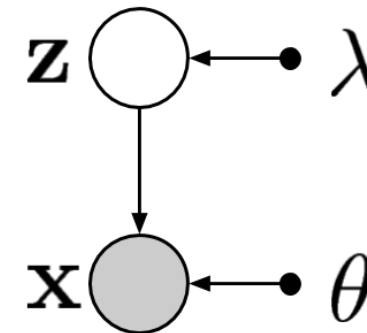


- If $p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{W}\mathbf{z} + \mathbf{b}, \Psi)$ and $p_{\lambda}(\mathbf{z}) = \mathcal{N}(\mu_0, \Sigma_0)$, then \rightarrow **Factor Analysis**.
- What if we take a non-linear transformation of \mathbf{z} ?
 \rightarrow *an infinite mixture of Gaussians*.

Latent Variable Models

- Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$



- If $p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{Wz} + \mathbf{b}, \Psi)$ and $p_{\lambda}(\mathbf{z}) = \mathcal{N}(\mu_0, \Sigma_0)$, then \rightarrow **Factor Analysis**.

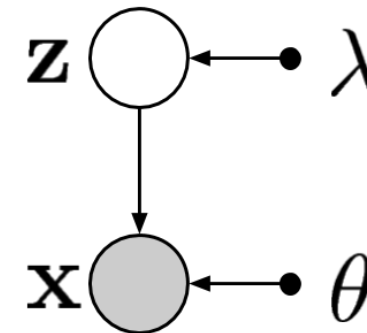
Convenient but limiting!

- What if we take a non-linear transformation of \mathbf{z} ?
 \rightarrow *an infinite mixture of Gaussians.*

Latent Variable Models

- Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$

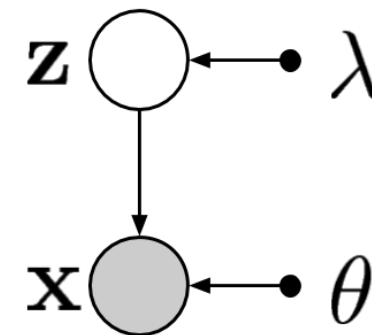


- If $p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{W}\mathbf{z} + \mathbf{b}, \Psi)$ and $p_{\lambda}(\mathbf{z}) = \mathcal{N}(\mu_0, \Sigma_0)$, then \rightarrow **Factor Analysis**.
- What if we take a **non-linear** transformation of \mathbf{z} ?
 \rightarrow *an infinite mixture of Gaussians*.

Latent Variable Models

- Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$

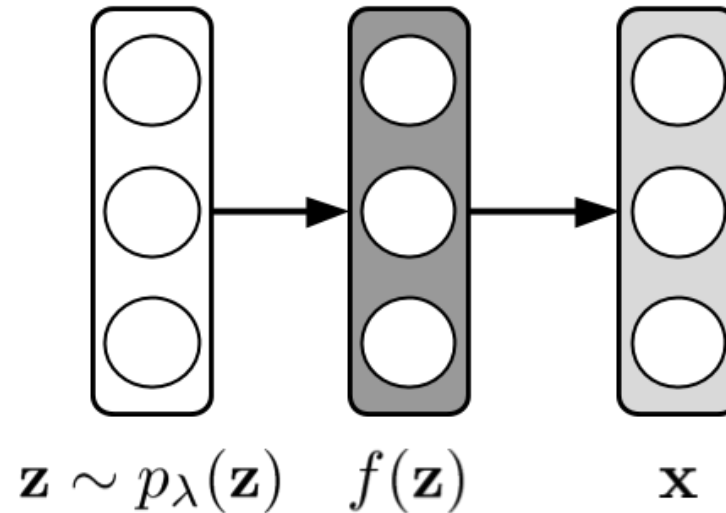


- If $p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{W}\mathbf{z} + \mathbf{b}, \Psi)$ and $p_{\lambda}(\mathbf{z}) = \mathcal{N}(\mu_0, \Sigma_0)$, then \rightarrow **Factor Analysis**.

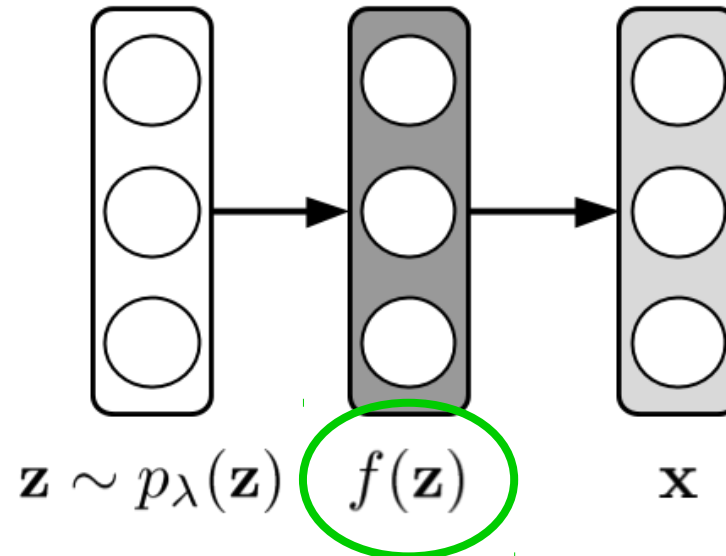
- What if we take a **non-linear transformation** of \mathbf{z} ?
 \rightarrow *an infinite mixture of Gaussians.*

Neural network

Deep Generative Models (DGM): Density Network

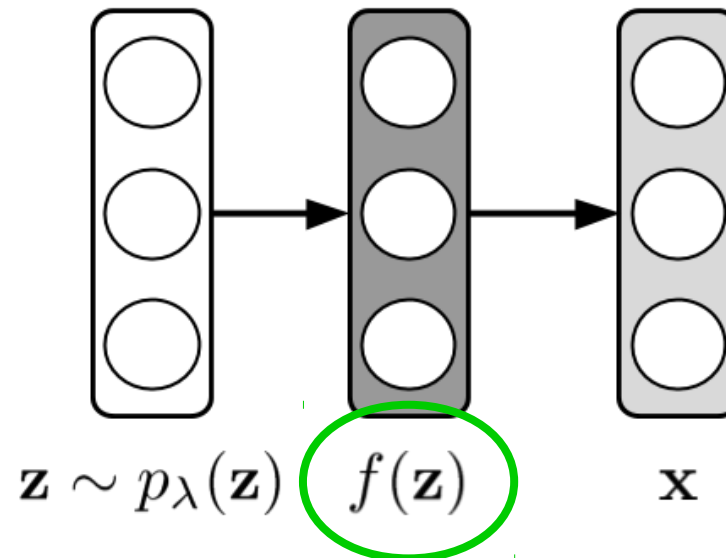


DGM: Density Network



Neural Network

DGM: Density Network



Neural Network

How to train this model?!

DGM: Density Network

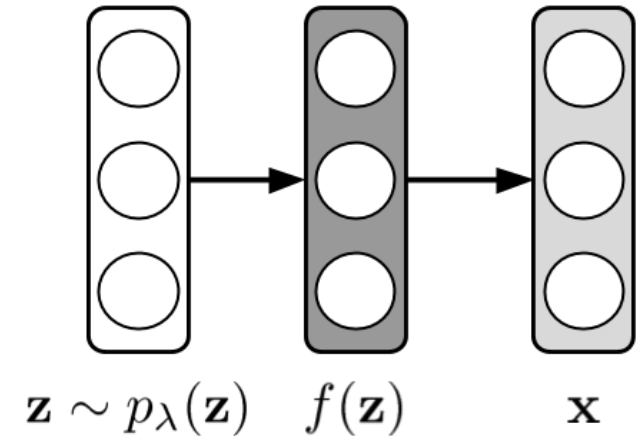
- MC approximation:

$$\log p(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$

$$\approx \log \frac{1}{S} \sum_{s=1}^S \exp \left(\log p_{\theta}(\mathbf{x}|\mathbf{z}_s) \right)$$

where:

$$\mathbf{z}_s \sim p_{\lambda}(\mathbf{z})$$



DGM: Density Network

- MC approximation:

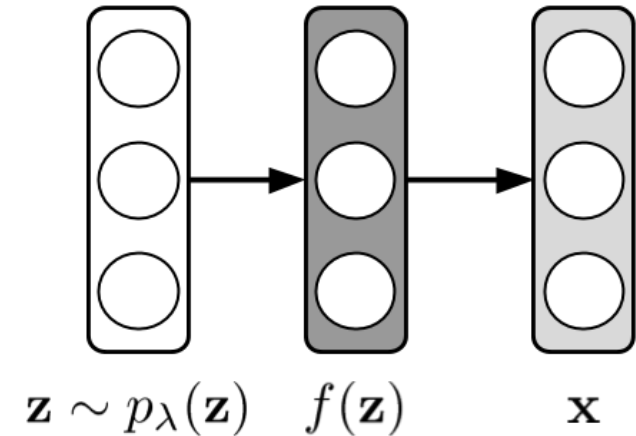
$$\log p(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$

$$\approx \log \frac{1}{S} \sum_{s=1}^S \exp \left(\log p_{\theta}(\mathbf{x}|\mathbf{z}_s) \right)$$

where:

$$\mathbf{z}_s \sim p_{\lambda}(\mathbf{z})$$

Sample \mathbf{z} many times,
apply log-sum-exp trick and
maximize log-likelihood.



DGM: Density Network

- MC approximation:

$$\log p(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$

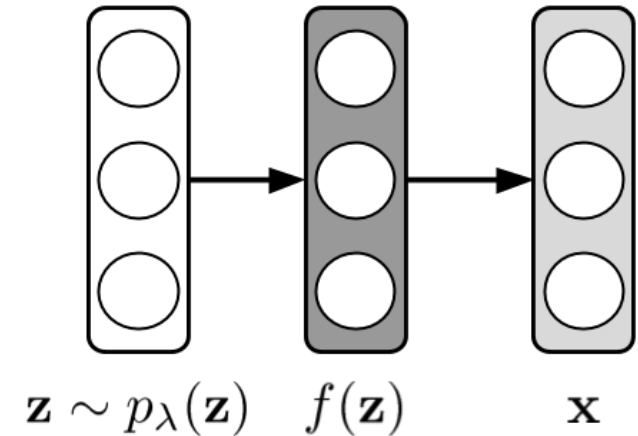
$$\approx \log \frac{1}{S} \sum_{s=1}^S \exp \left(\log p_{\theta}(\mathbf{x}|\mathbf{z}_s) \right)$$

where:

$$\mathbf{z}_s \sim p_{\lambda}(\mathbf{z})$$

Sample \mathbf{z} many times,
apply log-sum-exp trick and
maximize log-likelihood.

**It scales badly in high
dimensional cases!**



DGM: Density Network

PROS

Log-likelihood approach

Easy sampling

Training using gradient-based methods

CONS

Requires explicit models

Fails in high dim. cases

DGM: Density Network

PROS

Log-likelihood approach

Easy sampling

Training using gradient-based methods

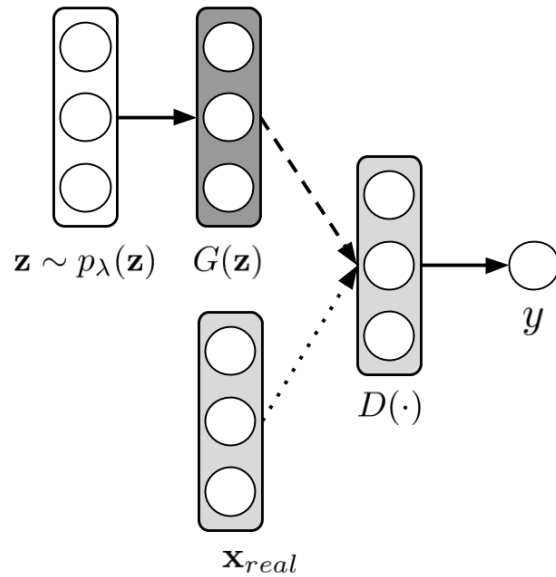
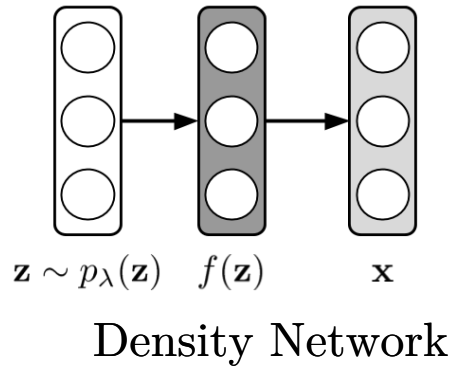
CONS

Requires explicit models

Fails in high dim. cases

Can we do better?

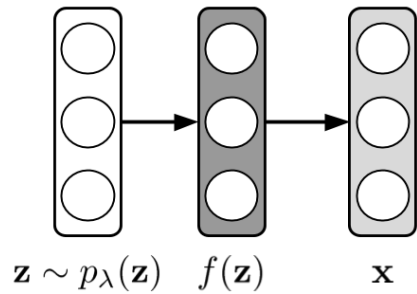
DGM: so far we have



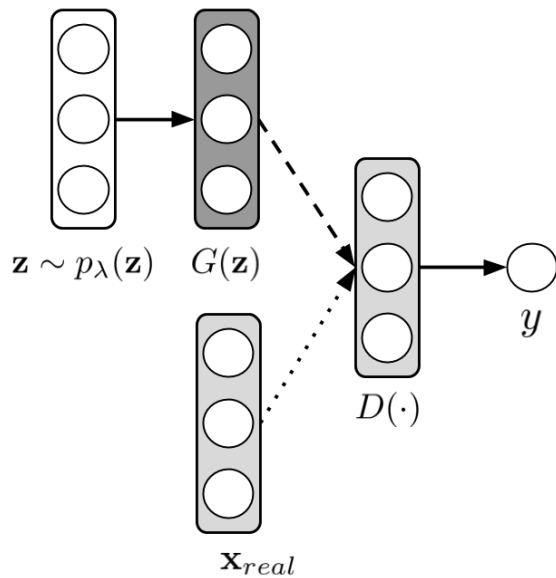
DGM: so far we have

Works only for low dim. cases...

Inefficient training...

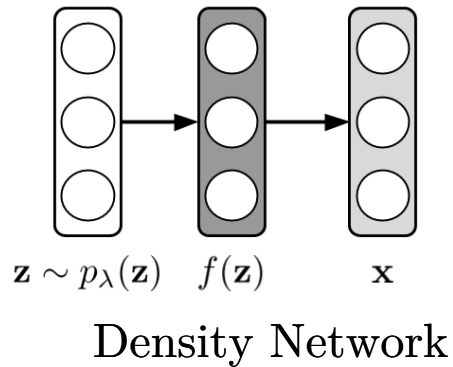


Density Network



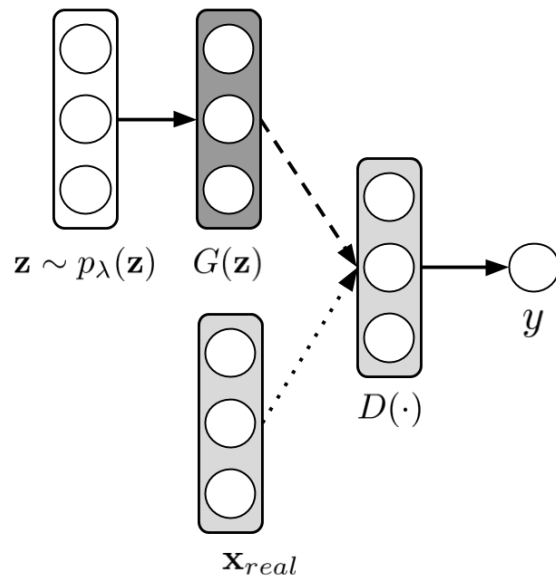
Generative Adversarial Net

DGM: so far we have



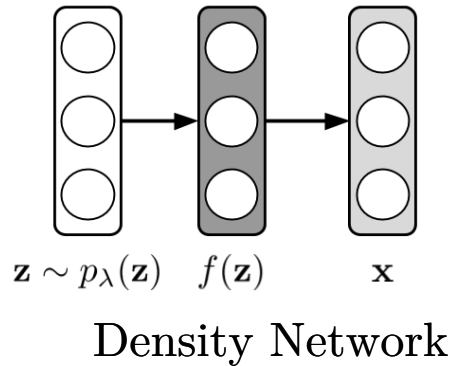
Works only for low dim. cases...

Inefficient training...



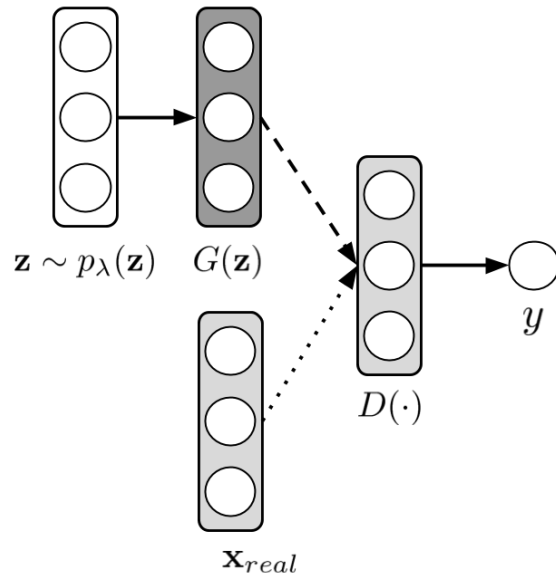
Works for high dim. cases!

DGM: so far we have



Works only for low dim. cases...

Inefficient training...

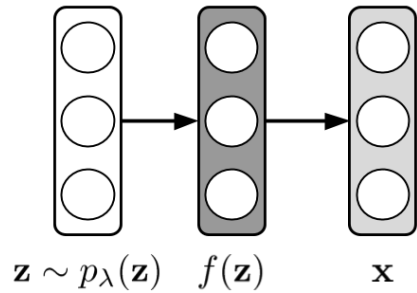


Works for high dim. cases!

Doesn't train a distribution...

Unstable training...

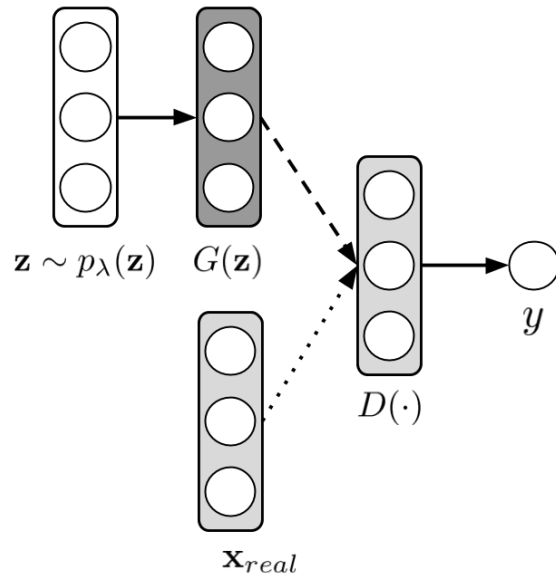
DGM: so far we have



Density Network

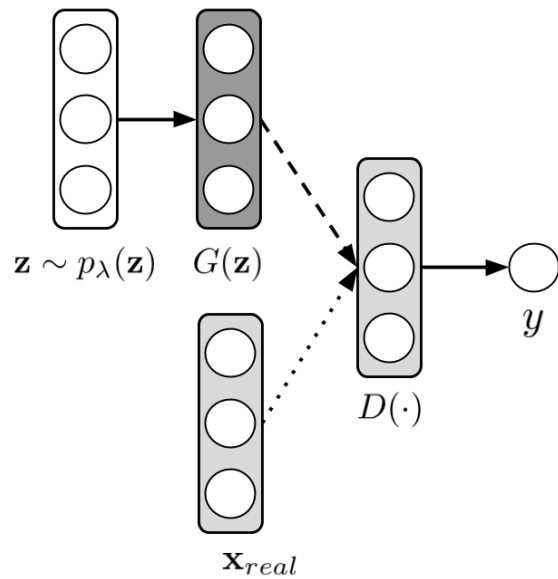
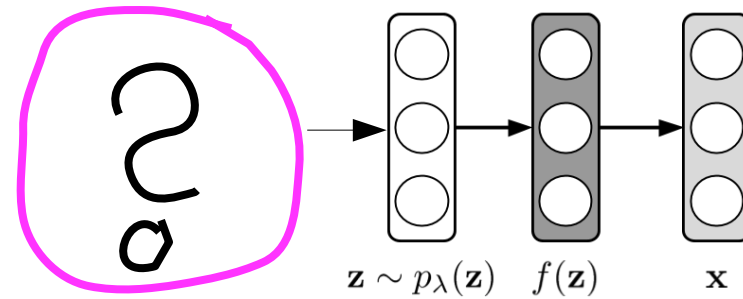
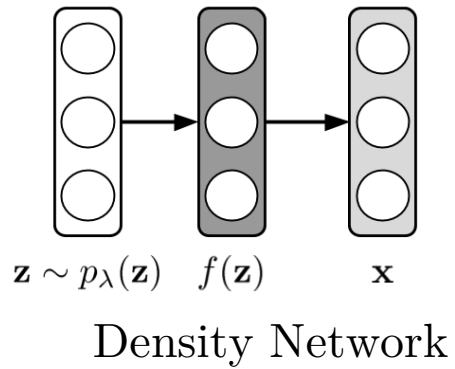
QUESTION

Can we stick to the log-likelihood approach but with a simple training procedure?



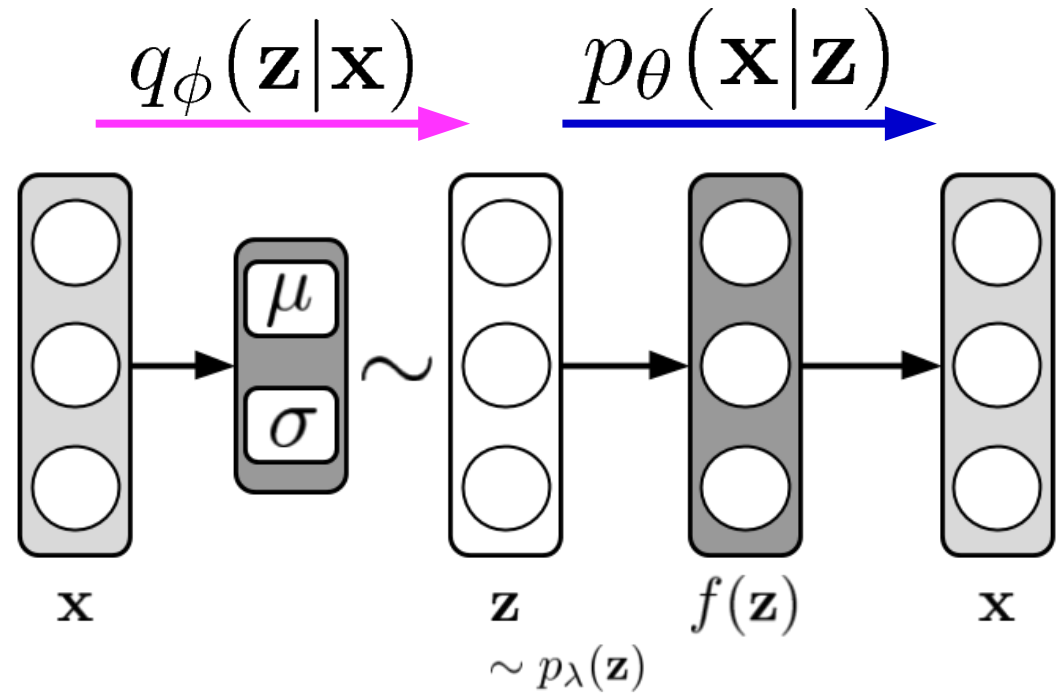
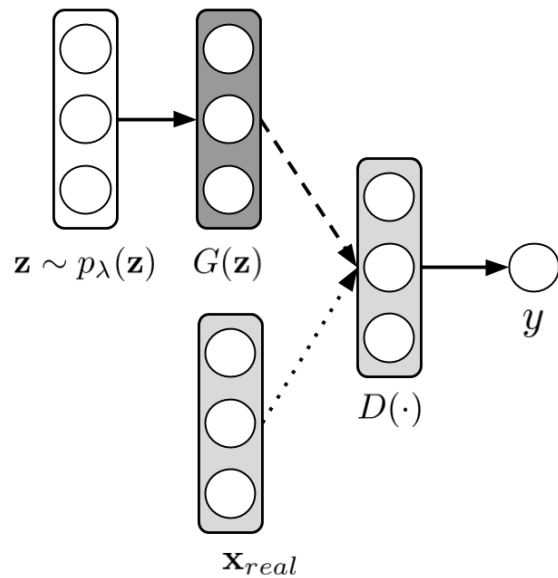
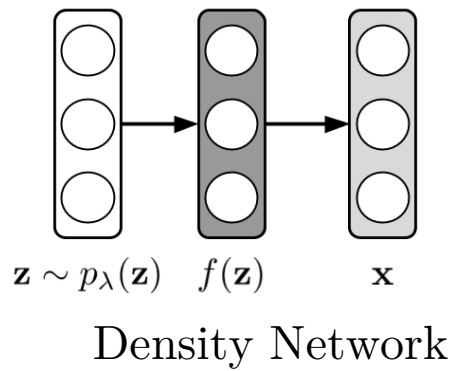
Generative Adversarial Net

DGM: so far we have



Generative Adversarial Net

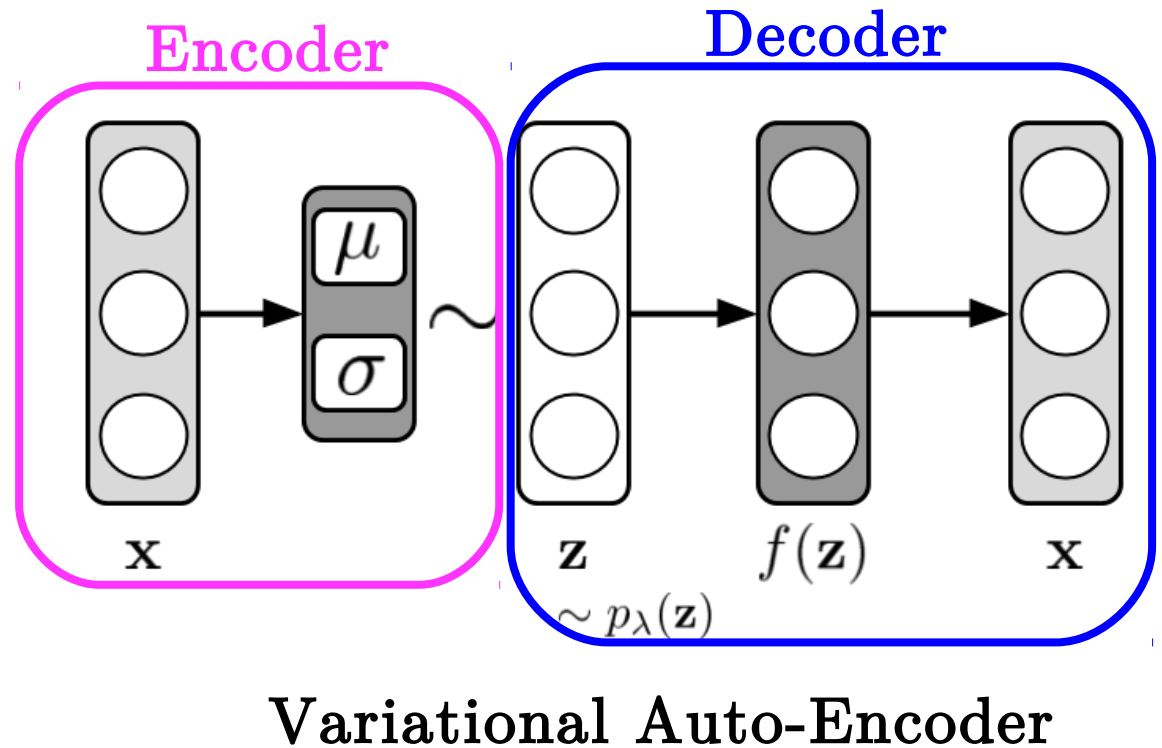
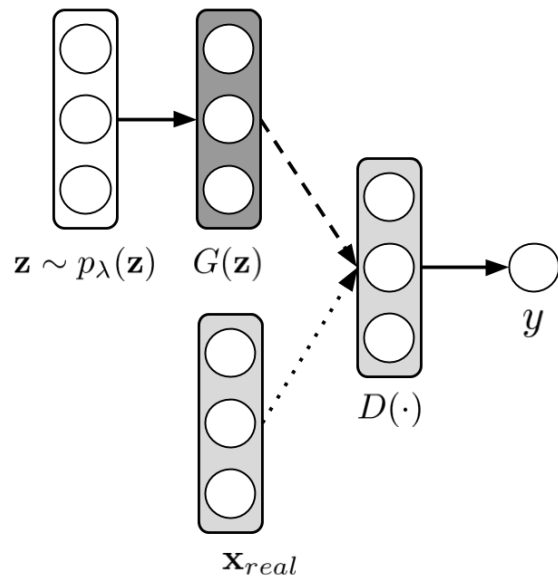
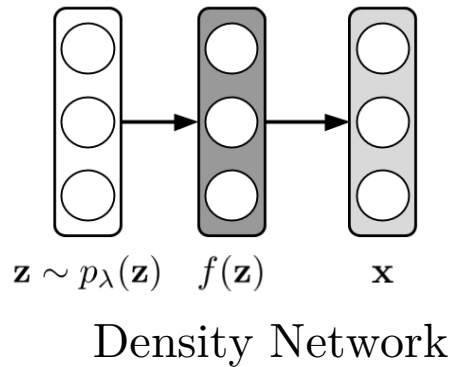
DGM: Variational Auto-Encoder



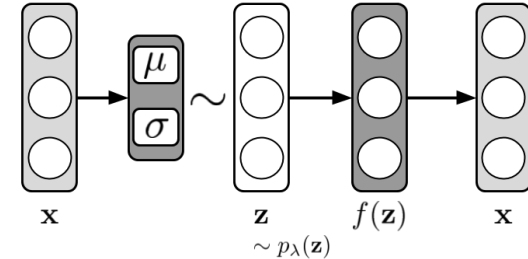
Generative Adversarial Net

Variational Auto-Encoder

DGM: Variational Auto-Encoder



DGM: Variational Auto-Encoder



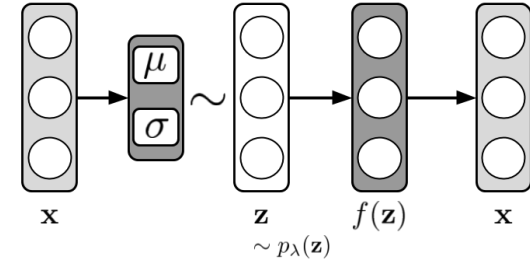
$$\log p(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$

$$= \log \int \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$

$$\geq \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} d\mathbf{z}$$

$$= \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}[q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\lambda}(\mathbf{z})]$$

DGM: Variational Auto-Encoder



$$\log p(\mathbf{x}) = \log \int p_\theta(\mathbf{x}|\mathbf{z}) p_\lambda(\mathbf{z}) d\mathbf{z}$$

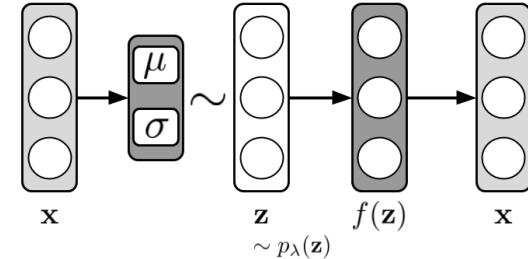
Variational posterior

$$= \log \int \frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} p_\theta(\mathbf{x}|\mathbf{z}) p_\lambda(\mathbf{z}) d\mathbf{z}$$

$$\geq \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{p_\theta(\mathbf{x}|\mathbf{z}) p_\lambda(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z}$$

$$= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \text{KL}[q_\phi(\mathbf{z}|\mathbf{x}) || p_\lambda(\mathbf{z})]$$

DGM: Variational Auto-Encoder



$$\log p(\mathbf{x}) = \log \int p_\theta(\mathbf{x}|\mathbf{z}) p_\lambda(\mathbf{z}) d\mathbf{z}$$

$$= \log \int \frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} p_\theta(\mathbf{x}|\mathbf{z}) p_\lambda(\mathbf{z}) d\mathbf{z}$$

$$\geq \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{p_\theta(\mathbf{x}|\mathbf{z}) p_\lambda(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z}$$

$$= \underbrace{\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction error}} - \underbrace{\text{KL}[q_\phi(\mathbf{z}|\mathbf{x}) || p_\lambda(\mathbf{z})]}_{\text{Regularization}}$$

DGM: Variational Auto-Encoder

Our objective is the evidence lower bound.

$$\log p(\mathbf{x}) \geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \text{KL}[q_\phi(\mathbf{z}|\mathbf{x}) || p_\lambda(\mathbf{z})]$$

We can approximate it using MC sample.

$$\mathcal{L}(\mathbf{x}) \approx \frac{1}{S} \sum_{s=1}^S [\log p_\theta(\mathbf{x}|\mathbf{z}_s) - \log q_\phi(\mathbf{z}_s|\mathbf{x}) + \log p_\lambda(\mathbf{z}_s)]$$

DGM: Variational Auto-Encoder

Our objective is the evidence lower bound.

$$\log p(\mathbf{x}) \geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \text{KL}[q_\phi(\mathbf{z}|\mathbf{x}) || p_\lambda(\mathbf{z})]$$

We can approximate it using MC sample.

$$\mathcal{L}(\mathbf{x}) \approx \frac{1}{S} \sum_{s=1}^S [\log p_\theta(\mathbf{x}|\mathbf{z}_s) - \log q_\phi(\mathbf{z}_s|\mathbf{x}) + \log p_\lambda(\mathbf{z}_s)]$$

How to properly calculate gradients (*i.e.*, train the model)?

DGM: Variational Auto-Encoder

$$\mathcal{L}(\mathbf{x}) \approx \frac{1}{S} \sum_{s=1}^S [\log p_{\theta}(\mathbf{x}|\mathbf{z}_s) - \log q_{\phi}(\mathbf{z}_s|\mathbf{x}) + \log p_{\lambda}(\mathbf{z}_s)]$$

PROBLEM: calculating gradient wrt parameters of the variational posterior (*i.e.*, sampling process).

DGM: Variational Auto-Encoder

$$\mathcal{L}(\mathbf{x}) \approx \frac{1}{S} \sum_{s=1}^S [\log p_{\theta}(\mathbf{x}|\mathbf{z}_s) - \log q_{\phi}(\mathbf{z}_s|\mathbf{x}) + \log p_{\lambda}(\mathbf{z}_s)]$$

PROBLEM: calculating gradient wrt parameters of the variational posterior (*i.e.*, sampling process).

SOLUTION: use a non-centered parameterization (a.k.a. *reparameterization trick*).

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\mu, \sigma^2)$$

$$\mathbf{z}_s = \mu + \sigma \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I})$$

DGM: Variational Auto-Encoder

$$\mathcal{L}(\mathbf{x}) \approx \frac{1}{S} \sum_{s=1}^S [\log p_{\theta}(\mathbf{x}|\mathbf{z}_s) - \log q_{\phi}(\mathbf{z}_s|\mathbf{x}) + \log p_{\lambda}(\mathbf{z}_s)]$$

PROBLEM: calculating gradient wrt parameters of the variational posterior (*i.e.*, sampling process).

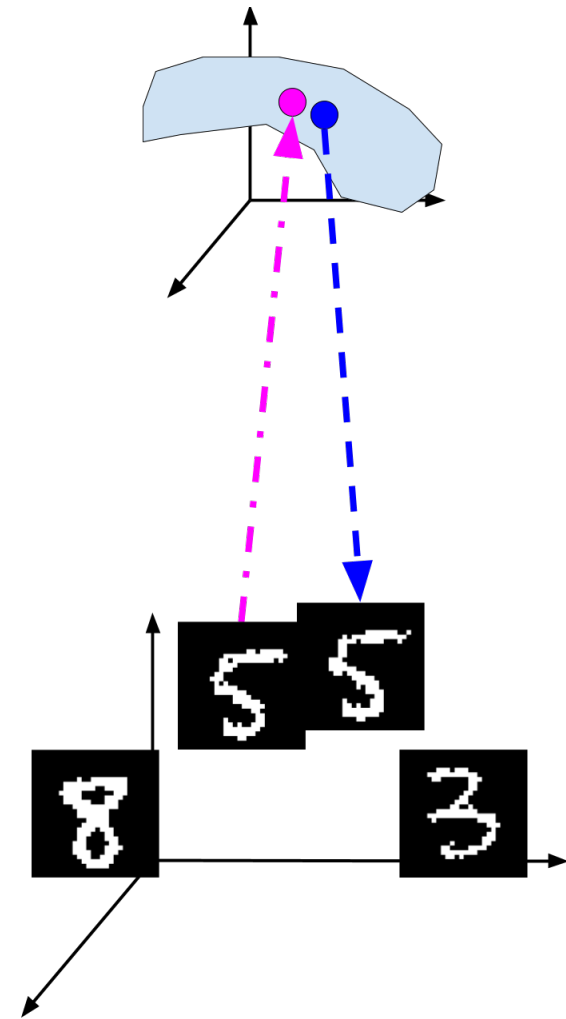
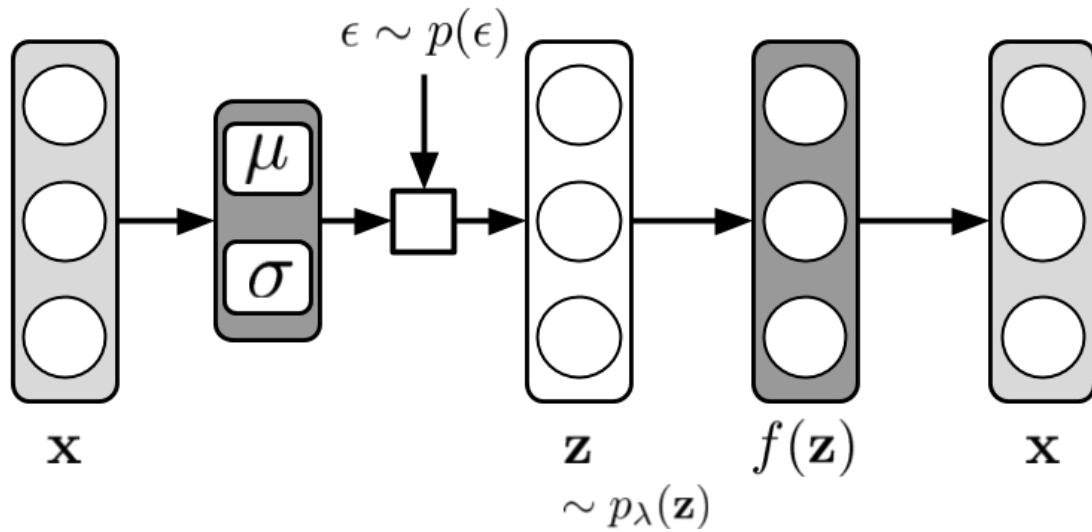
SOLUTION: use a non-centered parameterization (a.k.a. *reparameterization trick*).

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\mu, \sigma^2)$$

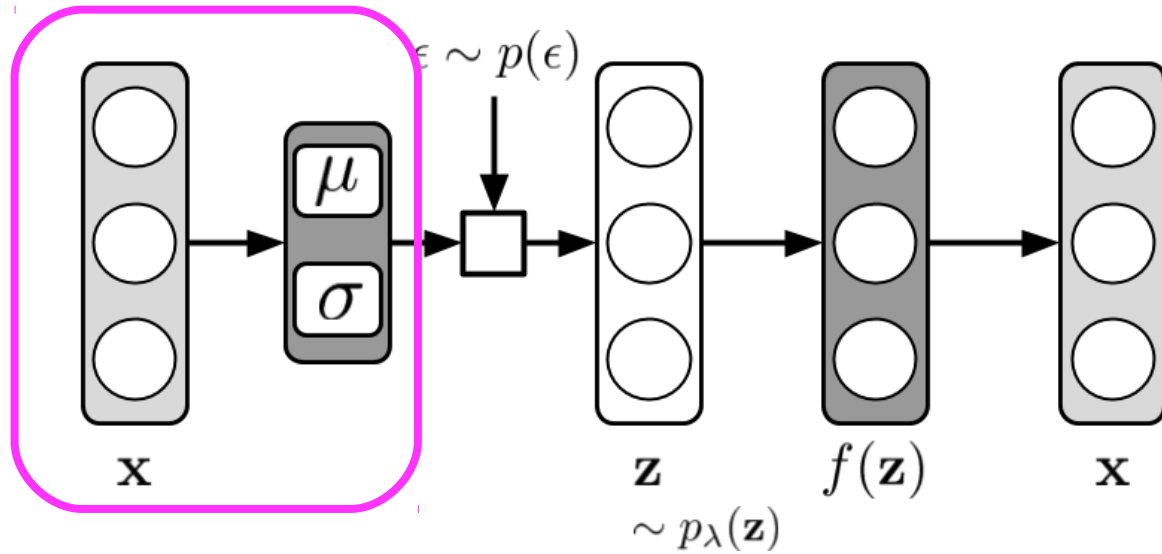
$$\mathbf{z}_s = \mu + \sigma \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I})$$

Output of a neural network

DGM: Variational Auto-Encoder

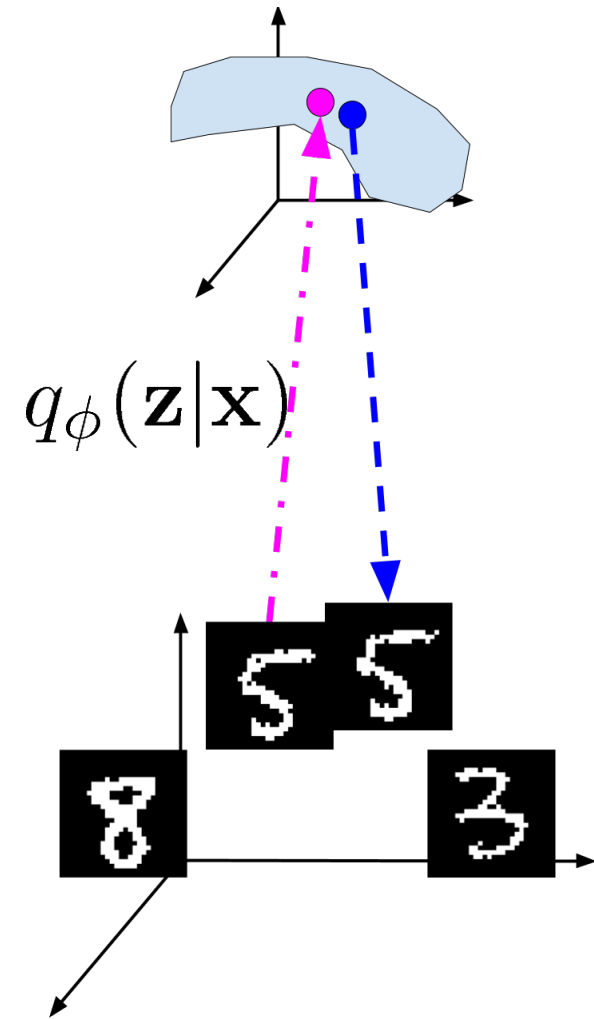


DGM: Variational Auto-Encoder

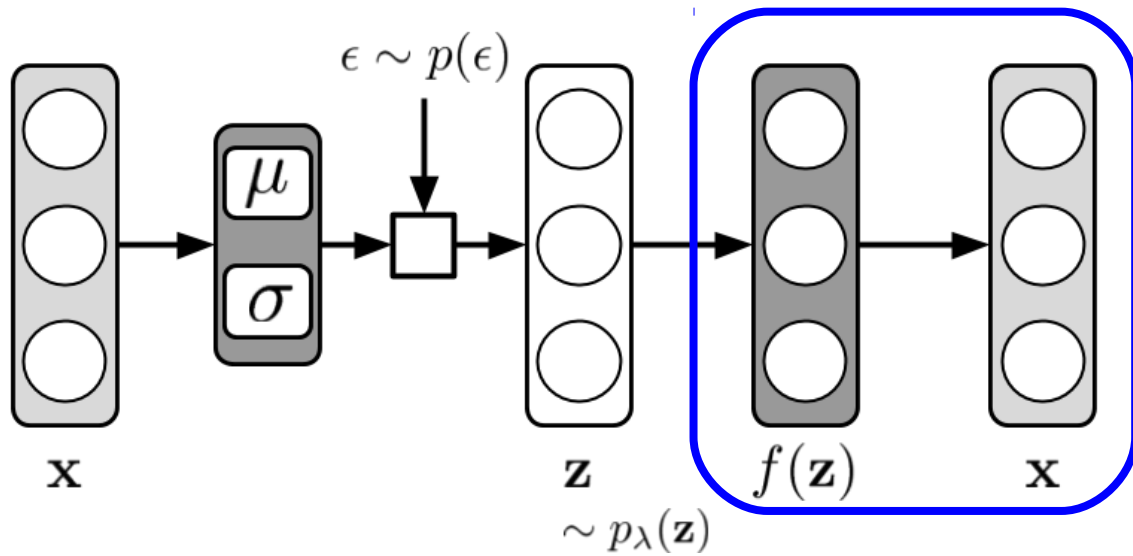


A **deep neural net** that outputs parameters of the variational posterior (**encoder**):

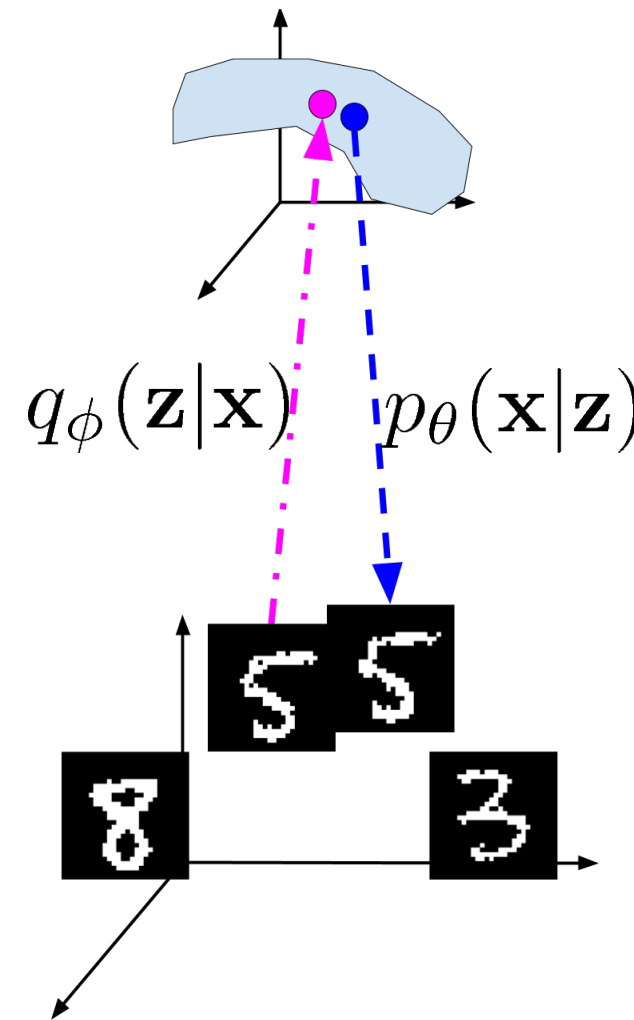
$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\underline{\mu(\mathbf{x})}, \text{diag}\{\underline{\sigma^2(\mathbf{x})}\})$$



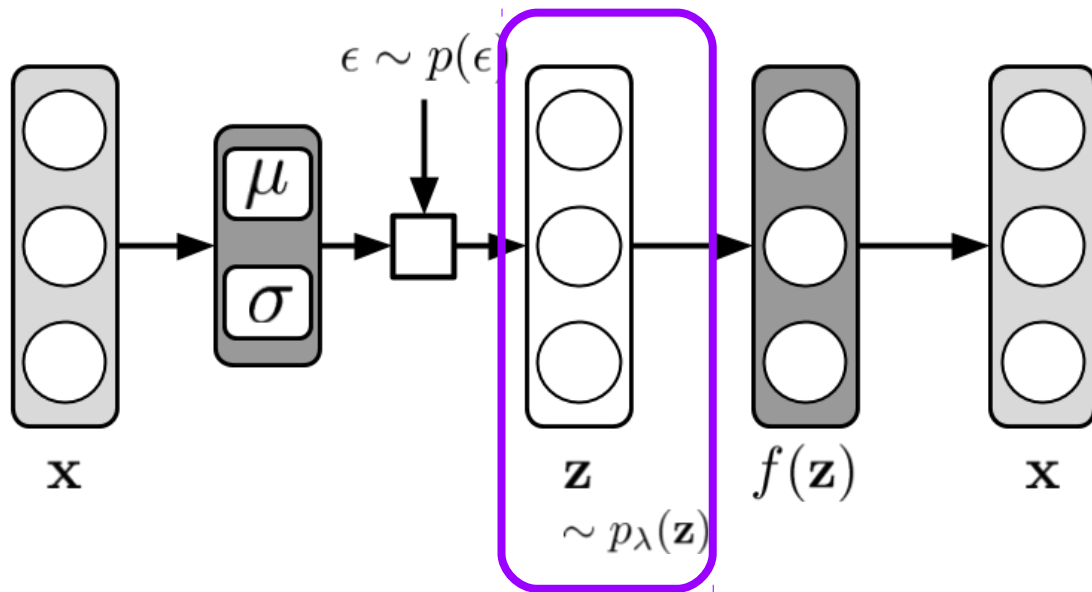
DGM: Variational Auto-Encoder



A **deep neural net** that outputs parameters of the generator (**decoder**), *e.g.*, a normal distribution or Bernoulli distribution.

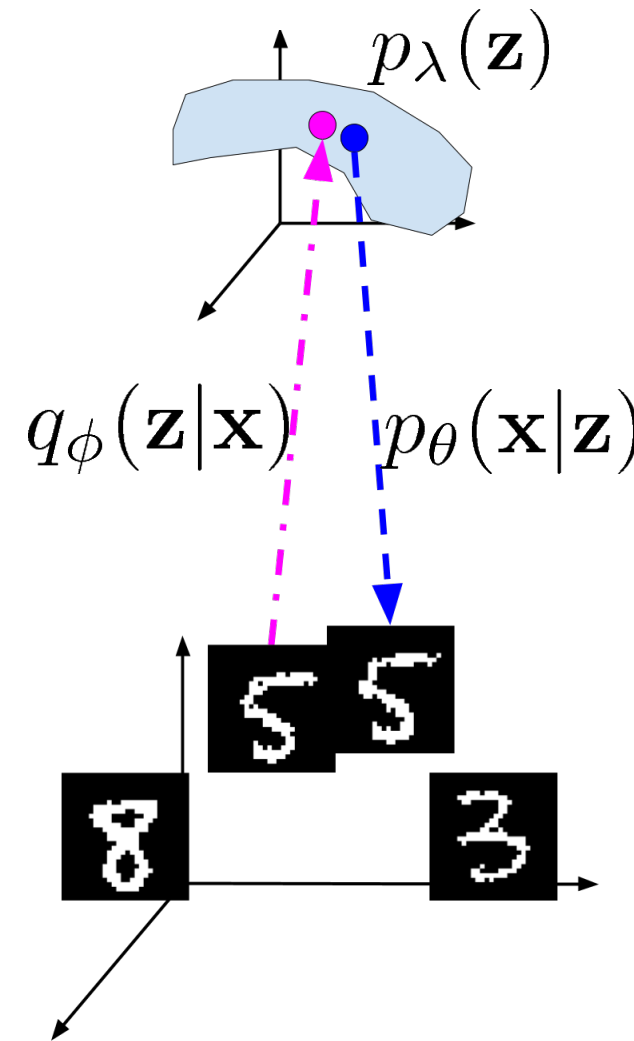


DGM: Variational Auto-Encoder



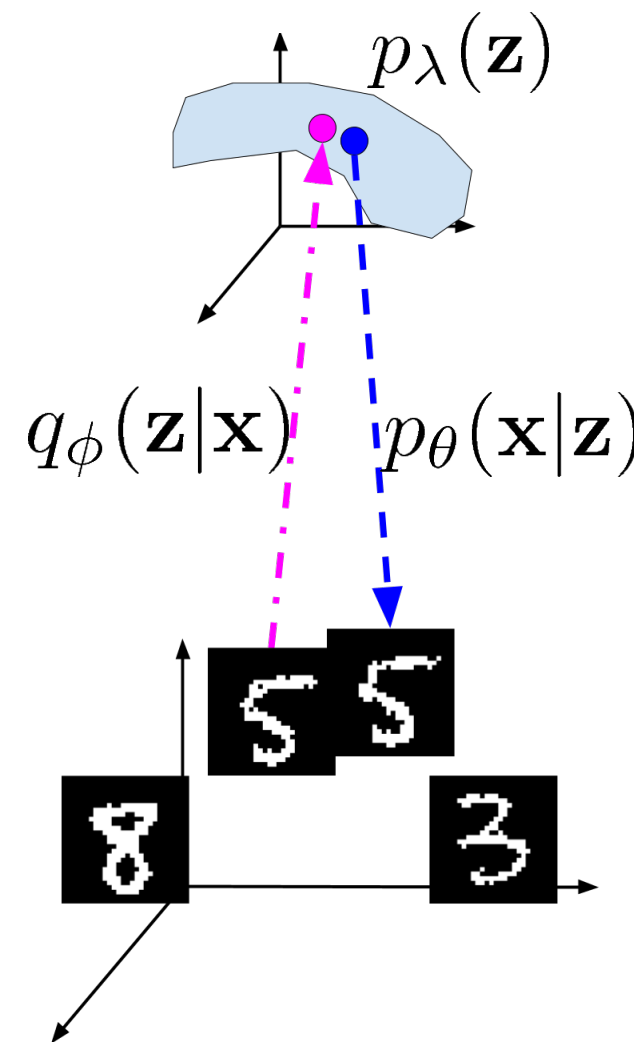
A **prior** that **regularizes** the encoder and takes part in the **generative process**.

$$p_\lambda(\mathbf{z}) = \mathcal{N}(\mathbf{z}|0, \mathbf{I})$$



DGM: Variational Auto-Encoder

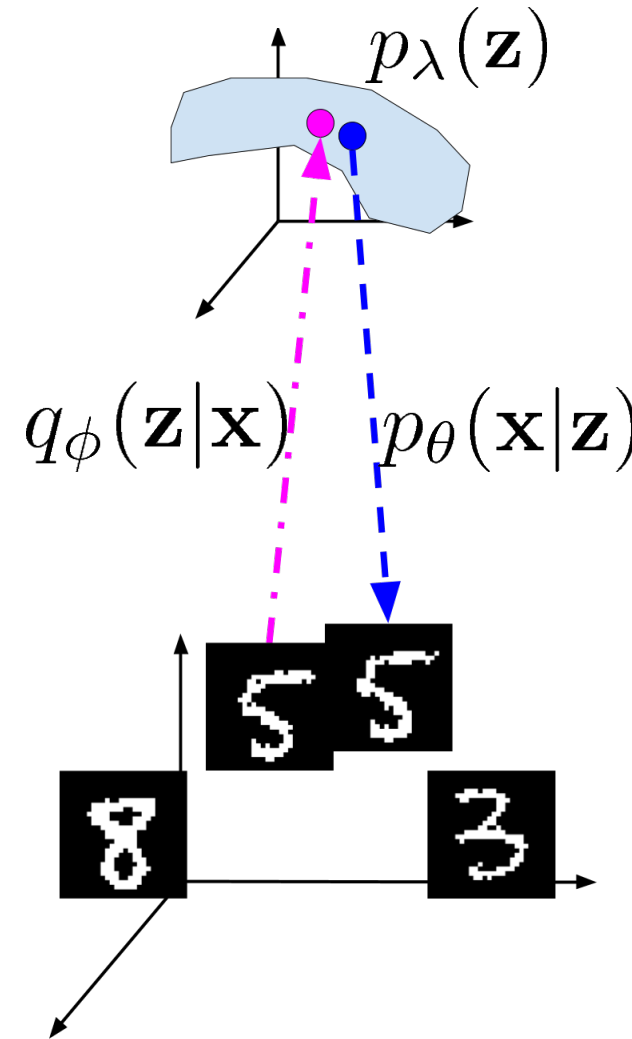
$$q_{\phi}(\mathbf{z}|\mathbf{x}) \propto p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z})$$



DGM: Variational Auto-Encoder

$$q_{\phi}(\mathbf{z}|\mathbf{x}) \propto p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z})$$

Feedforward nets
Convolutional nets
PixelCNN
Gated PixelCNN

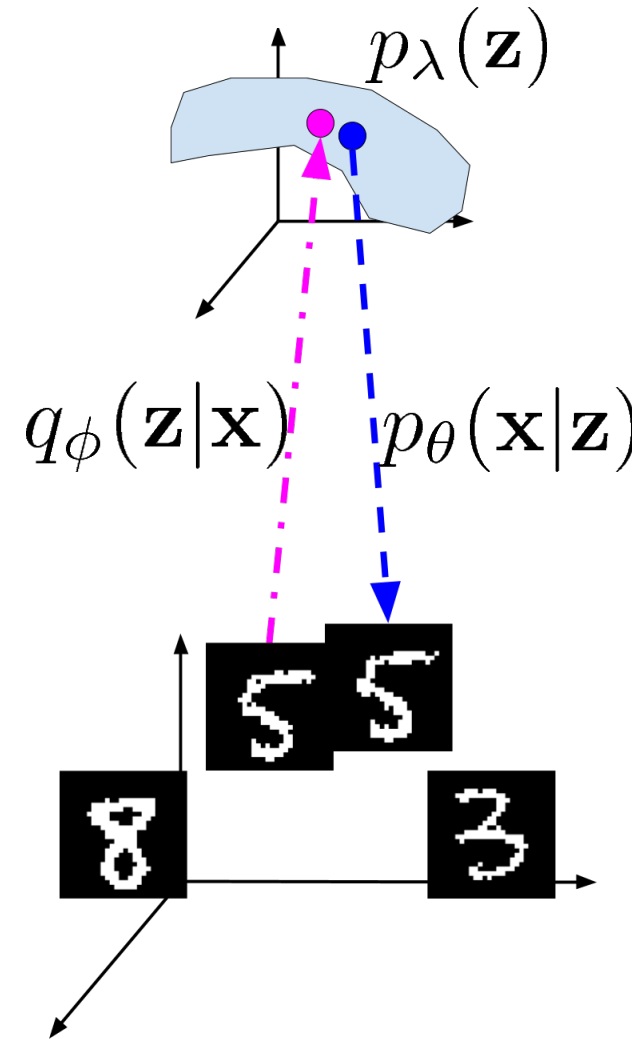


DGM: Variational Auto-Encoder

$$q_{\phi}(\mathbf{z}|\mathbf{x}) \propto p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z})$$

Normalizing flows
Volume-preserving flows
Gaussian processes
Stein Particle Descent
Operator VI

Feedforward nets
Convolutional nets
PixelCNN
Gated PixelCNN



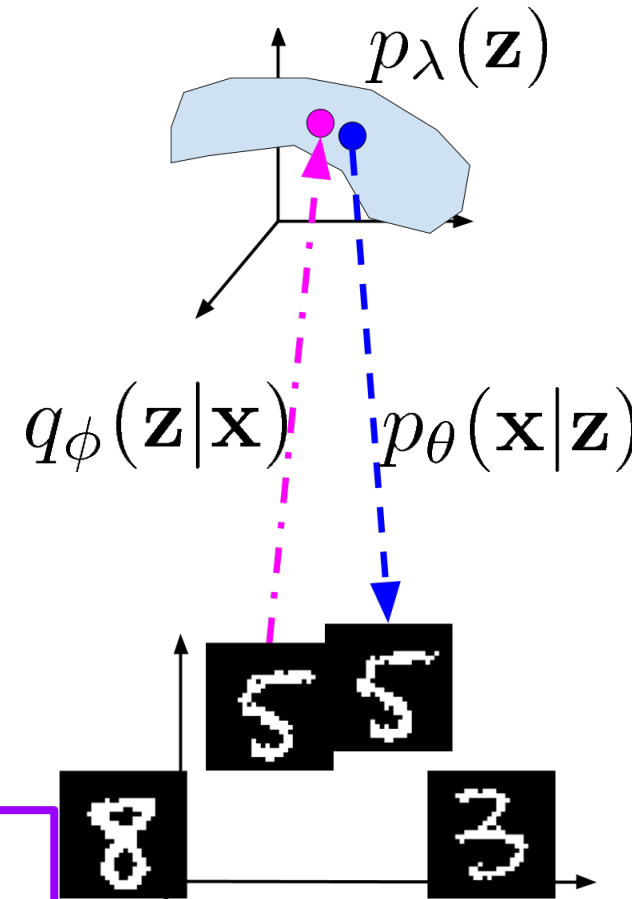
DGM: Variational Auto-Encoder

$$q_{\phi}(\mathbf{z}|\mathbf{x}) \propto p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z})$$

Normalizing flows
Volume-preserving flows
Gaussian processes
Stein Particle Descent
Operator VI

Feedforward nets
Convolutional nets
PixelCNN
Gated PixelCNN

Auto-regressive Prior
Objective Prior
Stick-Breaking Prior
VampPrior



DGM: Variational Auto-Encoder

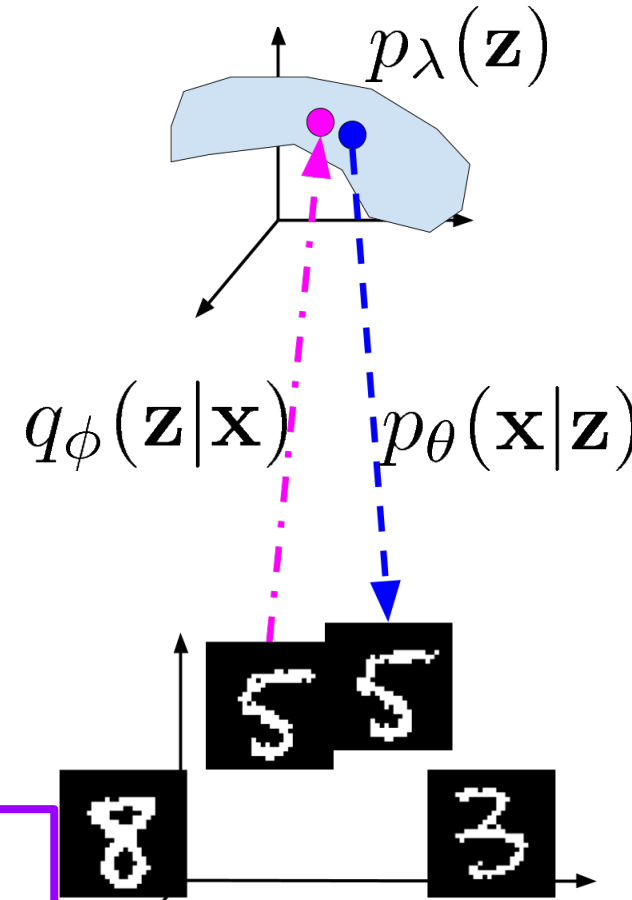
$$q_{\phi}(\mathbf{z}|\mathbf{x}) \propto p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z})$$

Normalizing flows
Volume-preserving flows
Gaussian processes
Stein Particle Descent
Operator VI

Feedforward nets
Convolutional nets
PixelCNN
Gated PixelCNN

Auto-regressive Prior
Objective Prior
Stick-Breaking Prior
VampPrior

Importance Weighted AE
Renyi Divergence
Stein Divergence



DGM: VAE

PROS

Log-likelihood framework

Easy sampling

Training using gradient-based methods

Stable training

Discovers latent representation

Could be easily combined with other probabilistic frameworks

CONS

Only **explicit** models

Produces **blurry images(?)**

In order to *make better decisions*, we need a *better understanding of reality*.

=

generative modeling

Web-page:

<https://jmtomczak.github.io>

Code on github:

<https://github.com/jmtomczak>

Contact:

J.M.Tomczak@uva.nl

jakubmkt@gmail.com

Part of the presented research was funded by the European Commission within the Marie Skłodowska-Curie Individual Fellowship (Grant No. 702666, "*Deep learning and Bayesian inference for medical imaging*").



RESEARCH & INNOVATION
Marie Skłodowska-Curie actions