# Deep Probabilistic Models (2)

Piyush Rai
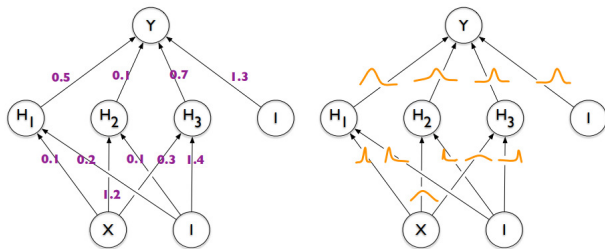
Probabilistic Machine Learning (CS772A)

Nov 2, 2017

# Recap: Bayesian Neural Networks

- Responses modeled via a suitable prob. distribution whose params are outputs of an NN, e.g.,

$$y_n \quad \sim \quad \mathcal{N}(\text{NN}(\boldsymbol{x}_n; \mathbf{W}), \sigma^2) \qquad \text{(for real-valued responses)}$$
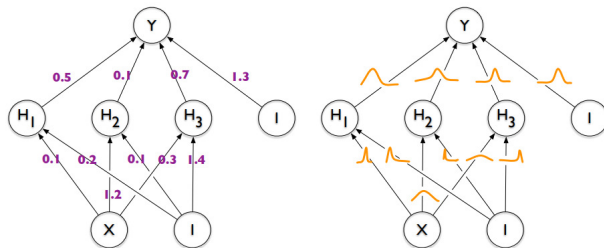


- $\text{NN}(\boldsymbol{x}_n; \mathbf{W})$ is a neural network with features $\boldsymbol{x}_n$ as its inputs and parameters $\mathbf{W}$

[†] "Weight Uncertainty in Neural Networks" (Blundell et al, 2015)

# Recap: Bayesian Neural Networks

- Responses modeled via a suitable prob. distribution whose params are outputs of an NN, e.g.,

$$y_n \quad \sim \quad \mathcal{N}(\text{NN}(\boldsymbol{x}_n; \mathbf{W}), \sigma^2) \qquad \text{(for real-valued responses)}$$
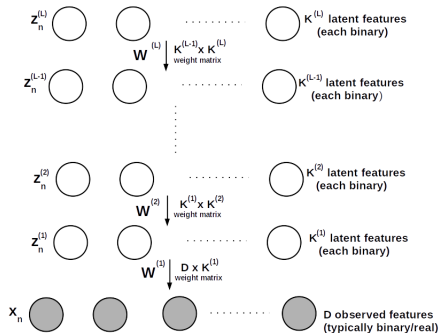


- $\text{NN}(\boldsymbol{x}_n; \mathbf{W})$ is a neural network with features $\boldsymbol{x}_n$ as its inputs and parameters $\mathbf{W}$

- Unlike standard neural networks, we learn the posterior over the unknowns

  - Non-conjugate model. MCMC or VB with Monte Carlo approximations typically used

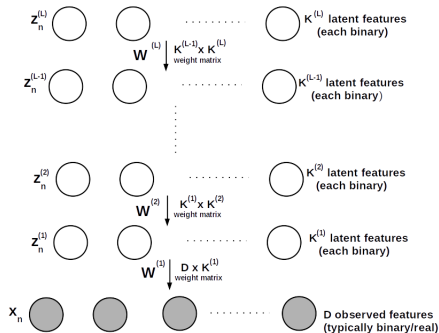[†] "Weight Uncertainty in Neural Networks" (Blundell et al, 2015)

# Recap: Sigmoid Belief Network

- An unsupervised generative model for the inputs $x_1, \ldots, x_N$
- Assumes data generated by successive nonlinear transformations of multiple layers of latent features
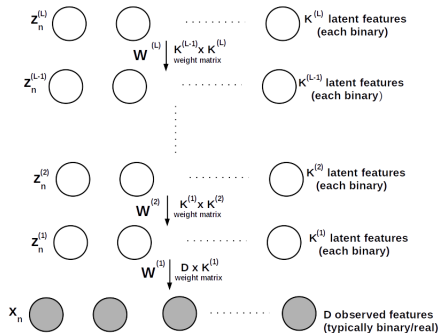
# Recap: Sigmoid Belief Network

- An unsupervised generative model for the inputs $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N$

- Assumes data generated by successive nonlinear transformations of multiple layers of latent features



- In SBN, the latent features in each hidden layer are assumed to be binary-valued
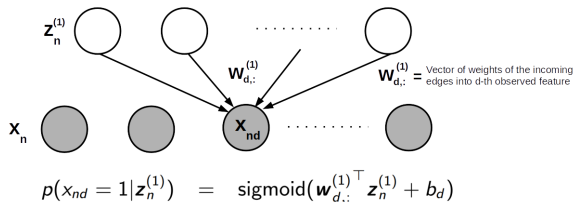
# Recap: Sigmoid Belief Network

- An unsupervised generative model for the inputs $x_1, \ldots, x_N$

- Assumes data generated by successive nonlinear transformations of multiple layers of latent features



- In SBN, the latent features in each hidden layer are assumed to be binary-valued

- The goal is to infer $\mathbf{Z}^{(1)}, \ldots, \mathbf{Z}^{(L)}$ and the other parameters of the network (MCMC or VB with Monte Carlo approximations is needed since the model is non-conjugate)

# Recap: Sigmoid Belief Network (A Zoomed-in Look)

- Layer 1 hidden nodes (latent features) generate each observed feature (assuming binary) as



$$p(x_{nd} = 1|\boldsymbol{z}_n^{(1)}) \quad = \quad \text{sigmoid}({\boldsymbol{w}_{d,:}^{(1)}}^{\top} \boldsymbol{z}_n^{(1)} + b_d)$$
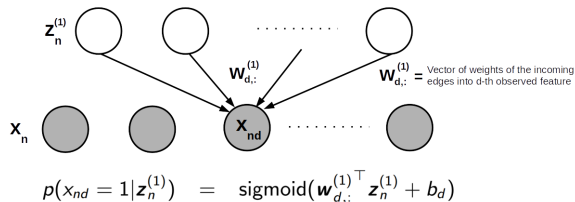
# Recap: Sigmoid Belief Network (A Zoomed-in Look)

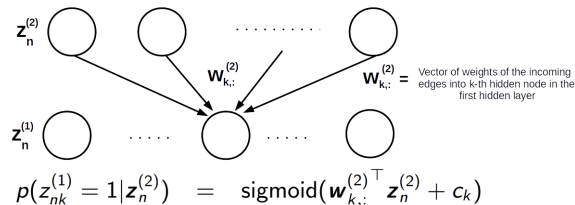- Layer 1 hidden nodes (latent features) generate each observed feature (assuming binary) as



$$p(x_{nd} = 1 | \mathbf{z}_n^{(1)}) = \text{sigmoid}(\mathbf{w}_{d,:}^{(1)\top} \mathbf{z}_n^{(1)} + b_d)$$

- Each hidden layer generates the nodes of hidden layer below it as (e.g., L2 → L1 in fig. below)



$$p(z_{nk}^{(1)} = 1 | \mathbf{z}_n^{(2)}) = \text{sigmoid}(\mathbf{w}_{k,:}^{(2)\top} \mathbf{z}_n^{(2)} + c_k)$$
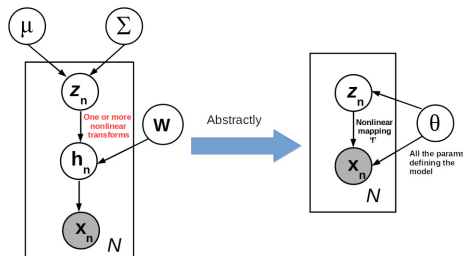
# Recap: Deep Latent Gaussian Models

- Assume that a Gaussian latent variable $z_n$ is pushed through a (deep) NN to generate $x_n$

# Recap: Deep Latent Gaussian Models

- Assume that a Gaussian latent variable $z_n$ is pushed through a (deep) NN to generate $x_n$



- A simple example with $z_n \sim \mathcal{N}(0, I_K)$ and assuming $x_n \in \mathbb{R}^D$ with Gaussian likelihood

$$x_n \sim \mathcal{N}(h_n, \sigma^2 I_D)$$

where $h_n$ is a deterministic nonlinear transform of $z_n$, e.g., $h_n = W \underbrace{\sigma(V z_n)}_{\text{sigmoid}}$ or $h_n = W \underbrace{\max\{0, V z_n\}}_{\text{ReLU}}$

# Recap: Deep Latent Gaussian Models

- Assume that a Gaussian latent variable $z_n$ is pushed through a (deep) NN to generate $x_n$
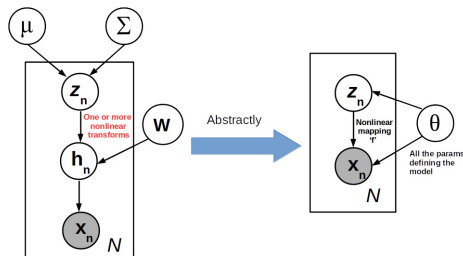


- A simple example with $z_n \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$ and assuming $x_n \in \mathbb{R}^D$ with Gaussian likelihood

$$x_n \sim \mathcal{N}(h_n, \sigma^2 \mathbf{I}_D)$$

where $h_n$ is a deterministic nonlinear transform of $z_n$, e.g., $h_n = \underbrace{\mathbf{W}\,\sigma(\mathbf{V}z_n)}_{\text{sigmoid}}$ or $h_n = \mathbf{W}\underbrace{\max\{\mathbf{0}, \mathbf{V}z_n\}}_{\text{ReLU}}$

- The goal is to infer the latent variables $z_n$'s and the other model params $\theta$ (MCMC or VB with Monte Carlo approximations is needed since the model is non-conjugate)

## Inference for Deep Latent Gaussian Models



- The usual approach for inference in such models (as in most probabilistic models) is iterative

## Inference for Deep Latent Gaussian Models



- The usual approach for inference in such models (as in most probabilistic models) is iterative, e.g.,
- Initialize $\theta$. Then iterate until converence
    - For $n = 1, \ldots, N$
        - Infer $p(z_n|x_n)$ using MCMC. If doing VB, update variational parameters $\phi_n$ of $q(z_n|\phi_n)$

# Inference for Deep Latent Gaussian Models



- The usual approach for inference in such models (as in most probabilistic models) is iterative, e.g.,
- Initialize $\theta$. Then iterate until converence
    - For $n = 1, \ldots, N$
        - Infer $p(z_n|x_n)$ using MCMC. If doing VB, update variational parameters $\phi_n$ of $q(z_n|\phi_n)$
    - Infer $\theta$ (its full posterior using MCMC or VB, or a point estimate)

# Inference for Deep Latent Gaussian Models

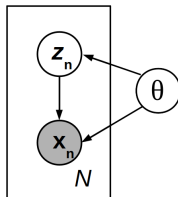

- The usual approach for inference in such models (as in most probabilistic models) is iterative, e.g.,
- Initialize $\theta$. Then iterate until converence
  - For $n = 1, \ldots, N$
    - Infer $p(z_n|x_n)$ using MCMC. If doing VB, update variational parameters $\phi_n$ of $q(z_n|\phi_n)$
  - Infer $\theta$ (its full posterior using MCMC or VB, or a point estimate)

- This iterative approach can be slow for large $N$
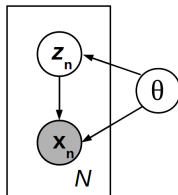
# Inference for Deep Latent Gaussian Models



- The usual approach for inference in such models (as in most probabilistic models) is iterative, e.g.,
- Initialize $\theta$. Then iterate until converence
    - For $n = 1, \ldots, N$
        - Infer $p(z_n | x_n)$ using MCMC. If doing VB, update variational parameters $\phi_n$ of $q(z_n | \phi_n)$
    - Infer $\theta$ (its full posterior using MCMC or VB, or a point estimate)
- This iterative approach can be slow for large $N$
- Also, inferring $z$ for new data point(s) $x$ would require using the same iterative procedure

# Variational Auto-encoder (VAE)

- Esseentially a DLGM, i.e., the $z$ to $x$ mapping $p(x|z)$ is defined by a neural net. Proposed almost simultaneously by Kingma & Welling (2013), and Rezende *et al* (2014)

# Variational Auto-encoder (VAE)

- Esseentially a DLGM, i.e., the $z$ to $x$ mapping $p(x|z)$ is defined by a neural net. Proposed almost simultaneously by Kingma & Welling (2013), and Rezende *et al* (2014)



- VAE uses VB for inference but has a fast, non-iterative way of computing $z_n$ for a data point $x_n$

# Variational Auto-encoder (VAE)

- Esseentially a DLGM, i.e., the $z$ to $x$ mapping $p(x|z)$ is defined by a neural net. Proposed almost simultaneously by Kingma & Welling (2013), and Rezende *et al* (2014)



- VAE uses VB for inference but has a fast, non-iterative way of computing $z_n$ for a data point $x_n$
- Key idea: For each point $x_n$, instead of learning a separate $q(z_n|\phi_n)$ with local params $\phi_n$, assume

$$q(z_n|\phi_n) = q(z_n|\text{NN}(x_n; \phi))$$

## Variational Auto-encoder (VAE)

- Esseentially a DLGM, i.e., the $z$ to $x$ mapping $p(x|z)$ is defined by a neural net. Proposed almost simultaneously by Kingma & Welling (2013), and Rezende *et al* (2014)
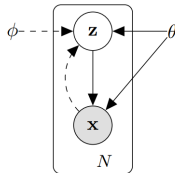


- VAE uses VB for inference but has a fast, non-iterative way of computing $z_n$ for a data point $x_n$
- Key idea: For each point $x_n$, instead of learning a separate $q(z_n|\phi_n)$ with local params $\phi_n$, assume

$$q(z_n|\phi_n) = q(z_n|\text{NN}(x_n; \phi))$$

so, basically, each $\phi_n$ is computed by a neural net with global parameters $\phi$ and input $x_n$

## Variational Auto-encoder (VAE)

- Esseentially a DLGM, i.e., the $z$ to $x$ mapping $p(x|z)$ is defined by a neural net. Proposed almost simultaneously by Kingma & Welling (2013), and Rezende *et al* (2014)



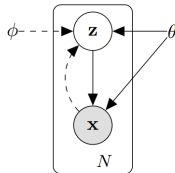- VAE uses VB for inference but has a fast, non-iterative way of computing $z_n$ for a data point $x_n$
- Key idea: For each point $x_n$, instead of learning a separate $q(z_n|\phi_n)$ with local params $\phi_n$, assume

$$q(z_n|\phi_n) = q(z_n|\text{NN}(x_n; \phi))$$

so, basically, each $\phi_n$ is computed by a neural net with global parameters $\phi$ and input $x_n$

- Once $\phi$ is learned, we can get $q(z_*|x_*) = q(z_*|\phi_*)$ for any $x_*$ by just using $\phi_* = \text{NN}(x_*; \phi)$

## Variational Auto-encoder (VAE)

- Esseentially a DLGM, i.e., the $z$ to $x$ mapping $p(x|z)$ is defined by a neural net. Proposed almost simultaneously by Kingma & Welling (2013), and Rezende *et al* (2014)



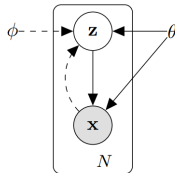- VAE uses VB for inference but has a fast, non-iterative way of computing $z_n$ for a data point $x_n$
- Key idea: For each point $x_n$, instead of learning a separate $q(z_n|\phi_n)$ with local params $\phi_n$, assume
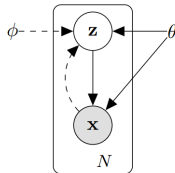
$$q(z_n|\phi_n) = q(z_n|\text{NN}(x_n; \phi))$$

so, basically, each $\phi_n$ is computed by a neural net with global parameters $\phi$ and input $x_n$

- Once $\phi$ is learned, we can get $q(z_*|x_*) = q(z_*|\phi_*)$ for any $x_*$ by just using $\phi_* = \text{NN}(x_*; \phi)$
- $p(x|z)$ is known as decoder and $q(z|x)$ is known as encoder

# Standard Auto-encoder vs Variational Auto-encoder

- A standard auto-encoder learns to (nonlinearly) compress and uncompress an input



- Model is trained to minimize the reconstruction error (difference b/w $x$ and $\hat{x}$)

## Standard Auto-encoder vs Variational Auto-encoder

- A standard auto-encoder learns to (nonlinearly) compress and uncompress an input



- Model is trained to minimize the reconstruction error (difference b/w $x$ and $\hat{x}$)
- However, it can't "generate" a "realistic" input from a random $z$ (the model isn't <u>trained</u> for that)

## Standard Auto-encoder vs Variational Auto-encoder

- A standard auto-encoder learns to (nonlinearly) compress and uncompress an input



- Model is trained to minimize the reconstruction error (difference b/w $x$ and $\hat{x}$)
- However, it can't "generate" a "realistic" input from a random $z$ (the model isn't <u>trained</u> for that)
- VAE allows this by assuming a distribution (e.g., Gaussian) over $z$ and learning to generate $x$ from random $z$'s drawn from that distribution (so the model is <u>trained</u> to do this!)

## Standard Auto-encoder vs Variational Auto-encoder
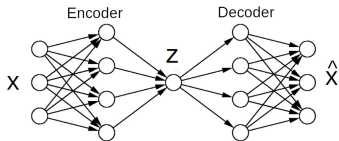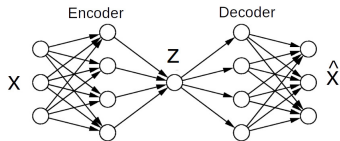
- A standard auto-encoder learns to (nonlinearly) compress and uncompress an input



- Model is trained to minimize the reconstruction error (difference b/w $x$ and $\hat{x}$)
- However, it can't "generate" a "realistic" input from a random $z$ (the model isn't <u>trained</u> for that)
- VAE allows this by assuming a distribution (e.g., Gaussian) over $z$ and learning to generate $x$ from random $z$'s drawn from that distribution (so the model is <u>trained</u> to do this!)



- Note: Simple generative models like PPCA or factor analysis also have this ability to generate data from random $z$ but the linear map from $z$ to $x$ limits the type of data that can be generated well

# VAE: The Encoder

- Role of encoder: Take $x$ as input and generate an encoding $z$



Encoder

# VAE: The Encoder

- Role of encoder: Take $x$ as input and generate an encoding $z$



Encoder

- Unlike standard autoencoders, for each $x$, VAE gives us a distribution $q(z|x)$ over its encoding

# VAE: The Encoder

- Role of encoder: Take $\boldsymbol{x}$ as input and generate an encoding $\boldsymbol{z}$



$$q_\phi(\mathbf{z}|\mathbf{x})$$

Encoder

- Unlike standard autoencoders, for each $\boldsymbol{x}$, VAE gives us a distribution $q(\boldsymbol{z}|\boldsymbol{x})$ over its encoding
- Assume $q(\boldsymbol{z}|\boldsymbol{x})$ to be Gaussian whose mean/var are computed by a NN with global params $\phi$

$$\mu_z = \text{NN}(\boldsymbol{x}; \phi) \qquad \sigma_z^2 = \text{NN}(\boldsymbol{x}; \phi)$$

# VAE: The Encoder

- Role of encoder: Take $x$ as input and generate an encoding $z$



$$q_\phi(z|x)$$

Encoder

- Unlike standard autoencoders, for each $x$, VAE gives us a distribution $q(z|x)$ over its encoding
- Assume $q(z|x)$ to be Gaussian whose mean/var are computed by a NN with global params $\phi$

$$\mu_z = \text{NN}(x; \phi) \qquad \sigma_z^2 = \text{NN}(x; \phi)$$

- Since $\mu_z, \sigma_z$ are outputs of neural networks, the $x$ to $z$ mapping is nonlinear

# VAE: The Decoder

- Role of decoder: Generate $\boldsymbol{x}$ given $\boldsymbol{z}$. Defined by the likelihood model $p_\theta(\boldsymbol{x}|\boldsymbol{z})$

$$p_\theta(\boldsymbol{x}|\boldsymbol{z})$$



Decoder

# VAE: The Decoder

- Role of decoder: Generate $\boldsymbol{x}$ given $\boldsymbol{z}$. Defined by the likelihood model $p_\theta(\boldsymbol{x}|\boldsymbol{z})$

$$p_\theta(\boldsymbol{x}|\boldsymbol{z})$$



Decoder

- Unlike PPCA (and like GPLVM), the $\boldsymbol{z}$ to $\boldsymbol{x}$ mapping is nonlinear (modeled by a neural network)

# VAE: The Decoder

- Role of decoder: Generate $x$ given $z$. Defined by the likelihood model $p_\theta(x|z)$

$$p_\theta(x|z)$$



Decoder

- Unlike PPCA (and like GPLVM), the $z$ to $x$ mapping is nonlinear (modeled by a neural network)
- Assume $p(x|z)$ to be Gaussian whose mean/var are computed by a NN with global params $\theta$

$$\mu_x = \text{NN}(z; \theta) \qquad \sigma_x^2 = \text{NN}(z; \theta)$$

# VAE: The Decoder

- Role of decoder: Generate $x$ given $z$. Defined by the likelihood model $p_\theta(x|z)$

$$p_\theta(x|z)$$



Decoder

- Unlike PPCA (and like GPLVM), the $z$ to $x$ mapping is nonlinear (modeled by a neural network)
- Assume $p(x|z)$ to be Gaussian whose mean/var are computed by a NN with global params $\theta$

$$\mu_x = \text{NN}(z; \theta) \qquad \sigma_x^2 = \text{NN}(z; \theta)$$

- Thus in the VAE, both $x$ to $z$ (encoder) and $z$ to $x$ (decoder) mappings are nonlinear

# Inference for VAE

- VAE uses variational inference (hence the name!) to learn the model parameters $\theta$ and $\phi$



Encoder

Decoder

# Inference for VAE

- VAE uses variational inference (hence the name!) to learn the model parameters $\theta$ and $\phi$



<div align="center">
$q_\phi(\mathbf{z}|\mathbf{x})$       $p_\theta(\mathbf{x}|\mathbf{z})$
</div>

Simply means that these two parts are actually "connected"

Encoder         Decoder

- Typically a prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$ is assumed on $\mathbf{z}$. The ELBO for a single $\mathbf{x}_n$ will be

$$\text{ELBO} = \mathbb{E}_{q_{\phi}}[\log p(\mathbf{x}_n, \mathbf{z}_n | \theta) - \log q(\mathbf{z}_n | \mathbf{x}_n)] \qquad (\text{note: } q_\phi \text{ and } q(\mathbf{z}_n | \mathbf{x}_n) \text{ mean the same})$$

# Inference for VAE

- VAE uses variational inference (hence the name!) to learn the model parameters $\theta$ and $\phi$



- Typically a prior $p(\boldsymbol{z}) = \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_K)$ is assumed on $\boldsymbol{z}$. The ELBO for a single $\boldsymbol{x}_n$ will be

$$
\begin{aligned}
\text{ELBO} &= \mathbb{E}_{q_\phi}[\log p(\boldsymbol{x}_n, \boldsymbol{z}_n|\theta) - \log q(\boldsymbol{z}_n|\boldsymbol{x}_n)] \qquad \text{(note: } q_\phi \text{ and } q(\boldsymbol{z}_n|\boldsymbol{x}_n) \text{ mean the same)} \\
&= \mathbb{E}_{q_\phi}[\log p(\boldsymbol{x}_n|\boldsymbol{z}_n)] - \text{KL}(q(\boldsymbol{z}_n|\boldsymbol{x}_n)||p(\boldsymbol{z}_n))
\end{aligned}
$$

# Inference for VAE

- VAE uses variational inference (hence the name!) to learn the model parameters $\theta$ and $\phi$



Encoder        Decoder

- Typically a prior $p(z) = \mathcal{N}(0, I_K)$ is assumed on $z$. The ELBO for a single $x_n$ will be

$$
\begin{aligned}
\text{ELBO} &= \mathbb{E}_{q_\phi}[\log p(x_n, z_n | \theta) - \log q(z_n | x_n)] \quad \text{(note: } q_\phi \text{ and } q(z_n | x_n) \text{ mean the same)} \\
&= \mathbb{E}_{q_\phi}[\log p(x_n | z_n)] - \text{KL}(q(z_n | x_n) || p(z_n))
\end{aligned}
$$

- ELBO intuition: Maximizing it will learn latent code $z_n$ that will give good reconstruction for $x_n$ (in expectation) and will keep $q(z_n | x_n)$ to be close to the prior $p(z_n)$ (i.e., regularization)

# Inference for VAE

- The ELBO is maximized w.r.t. the variational parameters $\phi$

$$
\begin{aligned}
\text{ELBO} &= \mathbb{E}_{q_\phi}[\log p(\boldsymbol{x}_n, \boldsymbol{z}_n | \theta) - \log q(\boldsymbol{z}_n | \boldsymbol{x}_n)] \\
&= \mathbb{E}_{q_\phi}[\log p(\boldsymbol{x}_n | \boldsymbol{z}_n)] - \text{KL}(q(\boldsymbol{z}_n | \boldsymbol{x}_n) || p(\boldsymbol{z}_n))
\end{aligned}
$$

# Inference for VAE

- The ELBO is maximized w.r.t. the variational parameters $\phi$

$$
\begin{aligned}
\text{ELBO} &= \mathbb{E}_{q_\phi}[\log p(\boldsymbol{x}_n, \boldsymbol{z}_n | \theta) - \log q(\boldsymbol{z}_n | \boldsymbol{x}_n)] \\
&= \mathbb{E}_{q_\phi}[\log p(\boldsymbol{x}_n | \boldsymbol{z}_n)] - \text{KL}(q(\boldsymbol{z}_n | \boldsymbol{x}_n) || p(\boldsymbol{z}_n))
\end{aligned}
$$

- ELBO is intractable because it requires computing expectations w.r.t. $q_\phi(\boldsymbol{z}_n | \boldsymbol{x}_n)$

# Inference for VAE

- The ELBO is maximized w.r.t. the variational parameters $\phi$

$$
\begin{aligned}
\text{ELBO} &= \mathbb{E}_{q_\phi}[\log p(\boldsymbol{x}_n, \boldsymbol{z}_n|\theta) - \log q(\boldsymbol{z}_n|\boldsymbol{x}_n)] \\
&= \mathbb{E}_{q_\phi}[\log p(\boldsymbol{x}_n|\boldsymbol{z}_n)] - \text{KL}(q(\boldsymbol{z}_n|\boldsymbol{x}_n)||p(\boldsymbol{z}_n))
\end{aligned}
$$

- ELBO is intractable because it requires computing expectations w.r.t. $q_\phi(\boldsymbol{z}_n|\boldsymbol{x}_n)$

- Note: If using the second form of the ELBO expression (the one with the KL term), the KL term has an analytic expression if $q(\boldsymbol{z}|\boldsymbol{x})$ and $p(\boldsymbol{z})$ are Gaussians, and need not be approximated

# Inference for VAE

- The ELBO is maximized w.r.t. the variational parameters $\phi$

$$
\begin{aligned}
\text{ELBO} &= \mathbb{E}_{q_\phi}[\log p(\boldsymbol{x}_n, \boldsymbol{z}_n | \theta) - \log q(\boldsymbol{z}_n | \boldsymbol{x}_n)] \\
&= \mathbb{E}_{q_\phi}[\log p(\boldsymbol{x}_n | \boldsymbol{z}_n)] - \text{KL}(q(\boldsymbol{z}_n | \boldsymbol{x}_n) || p(\boldsymbol{z}_n))
\end{aligned}
$$

- ELBO is intractable because it requires computing expectations w.r.t. $q_\phi(\boldsymbol{z}_n | \boldsymbol{x}_n)$

- Note: If using the second form of the ELBO expression (the one with the KL term), the KL term has an analytic expression if $q(\boldsymbol{z}|\boldsymbol{x})$ and $p(\boldsymbol{z})$ are Gaussians, and need not be approximated

  - E.g., if $q(\boldsymbol{z}_n | \boldsymbol{x}_n) = \mathcal{N}(\boldsymbol{\mu}_n, \text{diag}(\sigma_{n1}^2, \ldots, \sigma_{nK}^2))$ and $p(\boldsymbol{z}_n) = \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_K)$ then

$$
-\text{KL}(q(\boldsymbol{z}_n | \boldsymbol{x}_n) || p(\boldsymbol{z}_n)) = \frac{1}{2} \sum_{k=1}^{K} (1 + \log(\sigma_{nk}^2) - \mu_{nk}^2 - \sigma_{nk}^2)
$$

# Inference for VAE

- The ELBO is maximized w.r.t. the variational parameters $\phi$

$$\begin{aligned} \text{ELBO} &= \mathbb{E}_{q_\phi}[\log p(\boldsymbol{x}_n, \boldsymbol{z}_n | \theta) - \log q(\boldsymbol{z}_n | \boldsymbol{x}_n)] \\ &= \mathbb{E}_{q_\phi}[\log p(\boldsymbol{x}_n | \boldsymbol{z}_n)] - \text{KL}(q(\boldsymbol{z}_n | \boldsymbol{x}_n) || p(\boldsymbol{z}_n)) \end{aligned}$$

- ELBO is intractable because it requires computing expectations w.r.t. $q_\phi(\boldsymbol{z}_n | \boldsymbol{x}_n)$

- Note: If using the second form of the ELBO expression (the one with the KL term), the KL term has an analytic expression if $q(\boldsymbol{z}|\boldsymbol{x})$ and $p(\boldsymbol{z})$ are Gaussians, and need not be approximated

    - E.g., if $q(\boldsymbol{z}_n | \boldsymbol{x}_n) = \mathcal{N}(\boldsymbol{\mu}_n, \text{diag}(\sigma_{n1}^2, \ldots, \sigma_{nK}^2))$ and $p(\boldsymbol{z}_n) = \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_K)$ then
    $$-\text{KL}(q(\boldsymbol{z}_n | \boldsymbol{x}_n) || p(\boldsymbol{z}_n)) = \frac{1}{2} \sum_{k=1}^{K} (1 + \log(\sigma_{nk}^2) - \mu_{nk}^2 - \sigma_{nk}^2)$$

- Since we basically want derivative of the ELBO, one option is to use Monte Carlo samples $\boldsymbol{z}_n^{(1)}, \ldots, \boldsymbol{z}_n^{(L)}$ from $q(\boldsymbol{z}_n | \phi_n)$ and approximate the ELBO derivative as (recall BBVI discussion)

# Inference for VAE

- The ELBO is maximized w.r.t. the variational parameters $\phi$

$$
\begin{aligned}
\text{ELBO} &= \mathbb{E}_{q_\phi}[\log p(\boldsymbol{x}_n, \boldsymbol{z}_n | \theta) - \log q(\boldsymbol{z}_n | \boldsymbol{x}_n)] \\
&= \mathbb{E}_{q_\phi}[\log p(\boldsymbol{x}_n | \boldsymbol{z}_n)] - \text{KL}(q(\boldsymbol{z}_n | \boldsymbol{x}_n) || p(\boldsymbol{z}_n))
\end{aligned}
$$

- ELBO is intractable because it requires computing expectations w.r.t. $q_\phi(\boldsymbol{z}_n | \boldsymbol{x}_n)$

- Note: If using the second form of the ELBO expression (the one with the KL term), the KL term has an analytic expression if $q(\boldsymbol{z} | \boldsymbol{x})$ and $p(\boldsymbol{z})$ are Gaussians, and need not be approximated

    - E.g., if $q(\boldsymbol{z}_n | \boldsymbol{x}_n) = \mathcal{N}(\boldsymbol{\mu}_n, \text{diag}(\sigma_{n1}^2, \ldots, \sigma_{nK}^2))$ and $p(\boldsymbol{z}_n) = \mathcal{N}(\boldsymbol{0}, \mathbf{I}_K)$ then
    $$
    -\text{KL}(q(\boldsymbol{z}_n | \boldsymbol{x}_n) || p(\boldsymbol{z}_n)) = \frac{1}{2} \sum_{k=1}^{K} (1 + \log(\sigma_{nk}^2) - \mu_{nk}^2 - \sigma_{nk}^2)
    $$

- Since we basically want derivative of the ELBO, one option is to use Monte Carlo samples $\boldsymbol{z}_n^{(1)}, \ldots, \boldsymbol{z}_n^{(L)}$ from $q(\boldsymbol{z}_n | \phi_n)$ and approximate the ELBO derivative as (recall BBVI discussion)

$$
\nabla_\phi \mathbb{E}_q[f(\boldsymbol{z}_n)] = \mathbb{E}_q[f(\boldsymbol{z}_n) \nabla_\phi \log q(\boldsymbol{z}_n)] \approx \frac{1}{L} \sum_{\ell=1}^{L} f(\boldsymbol{z}_n^{(\ell)}) \nabla_\phi \log q(\boldsymbol{z}_n^{(\ell)})
$$

where $f(\boldsymbol{z}_n)$ denotes all the difficult terms in the ELBO expression

## Reparametrization Trick

- The Monte Carlo estimator of ELBO's gradient (previous slide) suffers from high variance

- Reason: The drawn samples $z_n^{(\ell)} \sim q_\phi(z_n|x_n)$ depend on $\phi$ and ELBO itself depends on $\phi$

$$\text{ELBO} \;\;=\;\; \mathbb{E}_{q_\phi}[\log p(x_n, z_n|\theta) - \log q_\phi(z_n|x_n)]$$

## Reparametrization Trick

- The Monte Carlo estimator of ELBO's gradient (previous slide) suffers from high variance

- Reason: The drawn samples $z_n^{(\ell)} \sim q_\phi(z_n|x_n)$ depend on $\phi$ and ELBO itself depends on $\phi$

$$\text{ELBO} = \mathbb{E}_{q_\phi}[\log p(x_n, z_n|\theta) - \log q_\phi(z_n|x_n)]$$

- Reparametrization Trick (Kingma & Welling, 2014): Make the randomness in $z$ independent of $\phi$

## Reparametrization Trick

- The Monte Carlo estimator of ELBO's gradient (previous slide) suffers from high variance

- Reason: The drawn samples $z_n^{(\ell)} \sim q_\phi(z_n|x_n)$ depend on $\phi$ and ELBO itself depends on $\phi$

$$\text{ELBO} = \mathbb{E}_{q_\phi}[\log p(x_n, z_n|\theta) - \log q_\phi(z_n|x_n)]$$

- Reparametrization Trick (Kingma & Welling, 2014): Make the randomness in $z$ independent of $\phi$

- Suppose our variational distribution $q_\phi(z|x) = \mathcal{N}(z|\mu, \sigma^2)$, where $\mu$, $\sigma$ are functions of $x$

## Reparametrization Trick

- The Monte Carlo estimator of ELBO's gradient (previous slide) suffers from high variance
- Reason: The drawn samples $z_n^{(\ell)} \sim q_\phi(z_n|x_n)$ depend on $\phi$ and ELBO itself depends on $\phi$

$$\text{ELBO} = \mathbb{E}_{q_\phi}[\log p(x_n, z_n|\theta) - \log q_\phi(z_n|x_n)]$$

- Reparametrization Trick (Kingma & Welling, 2014): Make the randomness in $z$ independent of $\phi$
- Suppose our variational distribution $q_\phi(z|x) = \mathcal{N}(z|\mu, \sigma^2)$, where $\mu, \sigma$ are functions of $x$
- We can reparametrize $z$ as $z = g(\epsilon, \phi, x) = \mu + \sigma\epsilon$, where $\epsilon \sim \mathcal{N}(\epsilon|0, 1)$

## Reparametrization Trick

- The Monte Carlo estimator of ELBO's gradient (previous slide) suffers from high variance
- Reason: The drawn samples $z_n^{(\ell)} \sim q_\phi(z_n|x_n)$ depend on $\phi$ and ELBO itself depends on $\phi$

$$\text{ELBO} = \mathbb{E}_{q_\phi}[\log p(x_n, z_n|\theta) - \log q_\phi(z_n|x_n)]$$

- Reparametrization Trick (Kingma & Welling, 2014): Make the randomness in $z$ independent of $\phi$
- Suppose our variational distribution $q_\phi(z|x) = \mathcal{N}(z|\mu, \sigma^2)$, where $\mu, \sigma$ are functions of $x$
- We can reparametrize $z$ as $z = g(\epsilon, \phi, x) = \mu + \sigma\epsilon$, where $\epsilon \sim \mathcal{N}(\epsilon|0, 1)$
- Now the above ELBO can be replaced by

$$\text{ELBO} = \mathbb{E}_{p(\epsilon)}[\log p(x_n, g(\epsilon, \phi, x_n)|\theta) - \log q_\phi(g(\epsilon, \phi, x_n))|x_n)]$$

## Reparametrization Trick

- The Monte Carlo estimator of ELBO's gradient (previous slide) suffers from high variance
- Reason: The drawn samples $z_n^{(\ell)} \sim q_\phi(z_n|x_n)$ depend on $\phi$ and ELBO itself depends on $\phi$

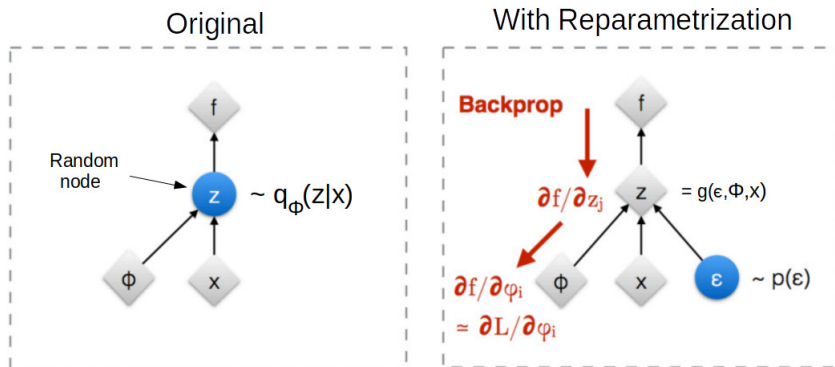$$\text{ELBO} = \mathbb{E}_{q_\phi}[\log p(x_n, z_n|\theta) - \log q_\phi(z_n|x_n)]$$

- Reparametrization Trick (Kingma & Welling, 2014): Make the randomness in $z$ independent of $\phi$
- Suppose our variational distribution $q_\phi(z|x) = \mathcal{N}(z|\mu, \sigma^2)$, where $\mu, \sigma$ are functions of $x$
- We can reparametrize $z$ as $z = g(\epsilon, \phi, x) = \mu + \sigma\epsilon$, where $\epsilon \sim \mathcal{N}(\epsilon|0, 1)$
- Now the above ELBO can be replaced by

$$\text{ELBO} = \mathbb{E}_{p(\epsilon)}[\log p(x_n, g(\epsilon, \phi, x_n)|\theta) - \log q_\phi(g(\epsilon, \phi, x_n))|x_n)]$$

- Given samples from $p(\epsilon)$, can easily approximate the ELBO (these samples don't depend on $\phi$)

# Reparametrization Trick

- Decoupling the randomness of **z** from $\phi$ using the reparametrization $z = g(\epsilon, \phi, x)$ also helps backpropagate easily through **z** when taking derivatives
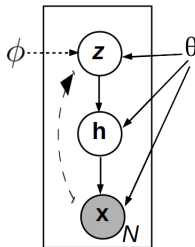
## VAE Architecture: An Example

- Assume a generative model (decoder) of the form $p_\theta(x|z)$ as specified below

$$\log p(\mathbf{x}|\mathbf{z}) = \log \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2\mathbf{I})$$
$$\text{where } \boldsymbol{\mu} = \mathbf{W}_4\mathbf{h} + \mathbf{b}_4$$
$$\log \boldsymbol{\sigma}^2 = \mathbf{W}_5\mathbf{h} + \mathbf{b}_5$$
$$\mathbf{h} = \tanh(\mathbf{W}_3\mathbf{z} + \mathbf{b}_3)$$



where $\theta$ consists of all the weights and bias terms

## VAE Architecture: An Example

- Assume a generative model (decoder) of the form $p_\theta(x|z)$ as specified below

$$\log p(\mathbf{x}|\mathbf{z}) = \log \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2\mathbf{I})$$
$$\text{where } \boldsymbol{\mu} = \mathbf{W}_4\mathbf{h} + \mathbf{b}_4$$
$$\log \boldsymbol{\sigma}^2 = \mathbf{W}_5\mathbf{h} + \mathbf{b}_5$$
$$\mathbf{h} = \tanh(\mathbf{W}_3\mathbf{z} + \mathbf{b}_3)$$



where $\theta$ consists of all the weights and bias terms

- The encoder $q_\phi(z|x)$ can also be defined to have the same form (swap $z$ and $x$ above) and the $\phi$'s will be the weights and bias terms

## VAE Architecture: An Example
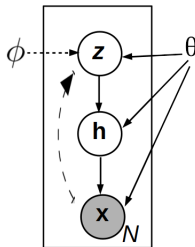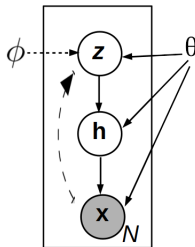
- Assume a generative model (decoder) of the form $p_\theta(x|z)$ as specified below

$$\log p(\mathbf{x}|\mathbf{z}) = \log \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2 \mathbf{I})$$
$$\text{where} \ \ \boldsymbol{\mu} = \mathbf{W}_4 \mathbf{h} + \mathbf{b}_4$$
$$\log \boldsymbol{\sigma}^2 = \mathbf{W}_5 \mathbf{h} + \mathbf{b}_5$$
$$\mathbf{h} = \tanh(\mathbf{W}_3 \mathbf{z} + \mathbf{b}_3)$$
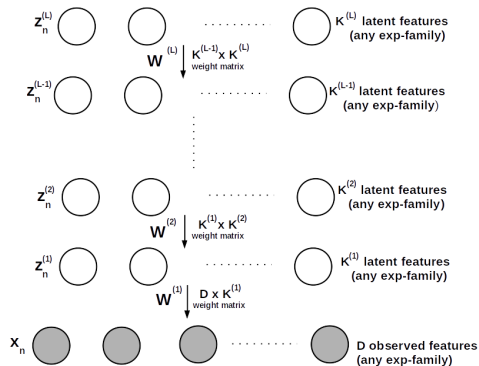
where $\theta$ consists of all the weights and bias terms

- The encoder $q_\phi(z|x)$ can also be defined to have the same form (swap $z$ and $x$ above) and the $\phi$'s will be the weights and bias terms

- Note: Recent work on VAE uses richer priors $p(z)$ as well as richer variational approx. $q(z|x)$

   - For standard VAEs, these are simple Gaussians with diagonal covariances

## Other Architectures: Deep Exponential Families

- Similar in architecture to sigmoid belief networks
- However, latent variables in every layer, as well as observations, are from exp. family distributions
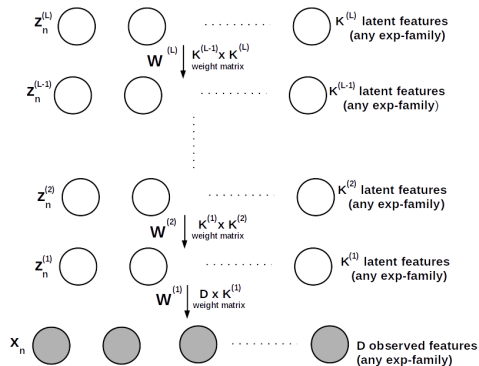
## Other Architectures: Deep Exponential Families
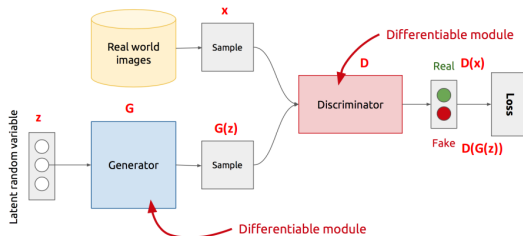
- Similar in architecture to sigmoid belief networks
- However, latent variables in every layer, as well as observations, are from exp. family distributions



- Overall model not conjugate but BBVI (Ranganath et al, 2013) or MCMC methods can be used

# Other Architectures: Generative Adversarial Networks

- Based on a game between a generator and a discriminator (Goodfellow et al, 2013)

# Other Architectures: Generative Adversarial Networks

- Based on a game between a generator and a discriminator (Goodfellow et al, 2013)



- Generator (a neural net) generates realistic looking "fake" data $x$ from random $z$

# Other Architectures: Generative Adversarial Networks

- Based on a game between a generator and a discriminator (Goodfellow et al, 2013)



- Generator (a neural net) generates realistic looking "fake" data **x** from random **z**
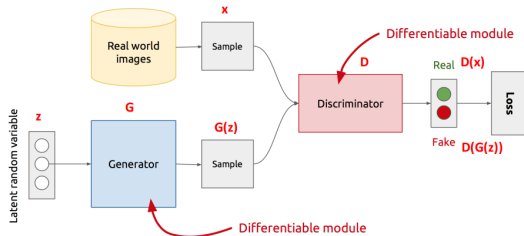- Discriminator tries to detect fake data from real data

# Other Architectures: Generative Adversarial Networks

- Based on a game between a generator and a discriminator (Goodfellow et al, 2013)



- Generator (a neural net) generates realistic looking "fake" data $x$ from random $z$
- Discriminator tries to detect fake data from real data
- At game's equilibrium, the $p_{gen} = p_{data}$ and success rate of discriminator $= 50\%$ (i.e., random)

---

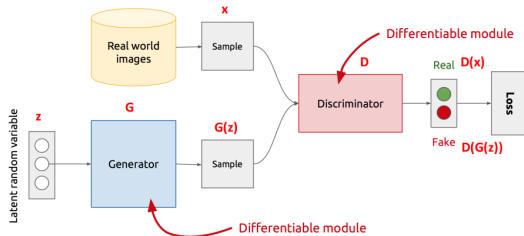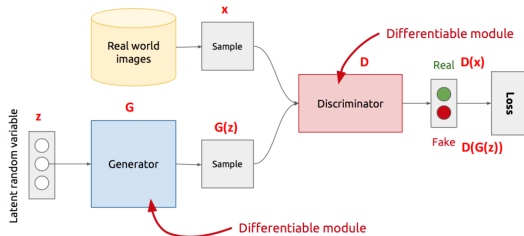# Other Architectures: Generative Adversarial Networks

- Based on a game between a generator and a discriminator (Goodfellow et al, 2013)



- Generator (a neural net) generates realistic looking "fake" data $x$ from random $z$
- Discriminator tries to detect fake data from real data
- At game's equilibrium, the $p_{gen} = p_{data}$ and success rate of discriminator $= 50\%$ (i.e., random)
- Originally designed mainly for synthetic data generation tasks but recent work extends GANs for many other problems such as latent variable inference, semi-supervised learning, etc.

# Summary

- Probabilistic modeling allows developing very flexible deep learning models

- Much of the recent progress is fuelled by advances in probabilistic modeling and inference

## Summary

- Probabilistic modeling allows developing very flexible deep learning models

- Much of the recent progress is fuelled by advances in probabilistic modeling and inference

- State-of-the-art results on a variety of tasks such as

  - Representation Learning (latent variables used as a new learned representation of data)

  - Density Estimation (i.e., $p(x)$)

  - Data generation (models like VAE and GAN can generate very realistic looking synthetic data)

  - Semi-supervised learning (models like VAE and GAN can be combined with supervised learning)

# Summary

- Probabilistic modeling allows developing very flexible deep learning models

- Much of the recent progress is fuelled by advances in probabilistic modeling and inference

- State-of-the-art results on a variety of tasks such as

  - Representation Learning (latent variables used as a new learned representation of data)

  - Density Estimation (i.e., $p(x)$)

  - Data generation (models like VAE and GAN can generate very realistic looking synthetic data)

  - Semi-supervised learning (models like VAE and GAN can be combined with supervised learning)

- Several other models such as Deep Boltzmann Machines, Neural Autoregressive Density Estimator, etc. that we didn't cover here

# Summary

- Probabilistic modeling allows developing very flexible deep learning models

- Much of the recent progress is fuelled by advances in probabilistic modeling and inference

- State-of-the-art results on a variety of tasks such as

  - Representation Learning (latent variables used as a new learned representation of data)

  - Density Estimation (i.e., $p(x)$)

  - Data generation (models like VAE and GAN can generate very realistic looking synthetic data)

  - Semi-supervised learning (models like VAE and GAN can be combined with supervised learning)

- Several other models such as Deep Boltzmann Machines, Neural Autoregressive Density Estimator, etc. that we didn't cover here

- An important distinction between explicit and implicit generative models

# Summary

- Probabilistic modeling allows developing very flexible deep learning models

- Much of the recent progress is fuelled by advances in probabilistic modeling and inference

- State-of-the-art results on a variety of tasks such as

  - Representation Learning (latent variables used as a new learned representation of data)

  - Density Estimation (i.e., $p(x)$)

  - Data generation (models like VAE and GAN can generate very realistic looking synthetic data)

  - Semi-supervised learning (models like VAE and GAN can be combined with supervised learning)

- Several other models such as Deep Boltzmann Machines, Neural Autoregressive Density Estimator, etc. that we didn't cover here

- An important distinction between explicit and implicit generative models

  - Models like PPCA, FA, DLGM, SBN, VAE, etc. have an explicit likelihood model for data

# Summary

- Probabilistic modeling allows developing very flexible deep learning models

- Much of the recent progress is fuelled by advances in probabilistic modeling and inference

- State-of-the-art results on a variety of tasks such as

    - Representation Learning (latent variables used as a new learned representation of data)

    - Density Estimation (i.e., $p(x)$)

    - Data generation (models like VAE and GAN can generate very realistic looking synthetic data)

    - Semi-supervised learning (models like VAE and GAN can be combined with supervised learning)

- Several other models such as Deep Boltzmann Machines, Neural Autoregressive Density Estimator, etc. that we didn't cover here

- An important distinction between explicit and implicit generative models

    - Models like PPCA, FA, DLGM, SBN, VAE, etc. have an explicit likelihood model for data

    - A model like GAN only defines $p(x)$ implicitly (no "likelihood" model for data)