

Probabilistic Modeling meets Deep Learning

Piyush Rai

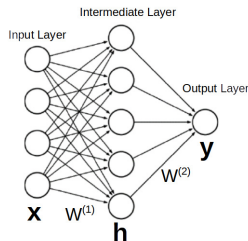
Topics in Probabilistic Modeling and Inference (CS698X)

April 3, 2019



Neural Networks

- A simple neural network with one intermediate (also called “hidden”) layer and a single output

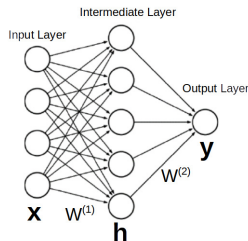


- Each intermediate layer computes a **nonlinear transformation** of its previous layer's nodes



Neural Networks

- A simple neural network with one intermediate (also called “hidden”) layer and a single output

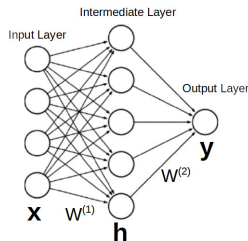


- Each intermediate layer computes a **nonlinear transformation** of its previous layer's nodes
- In traditional neural nets, \mathbf{h} is a **Linear transform** (e.g., $\mathbf{W}^{(1)}\mathbf{x}$ in the above picture) **followed by a nonlinearity** (e.g., sigmoid, ReLU, tanh, etc)



Neural Networks

- A simple neural network with one intermediate (also called “hidden”) layer and a single output

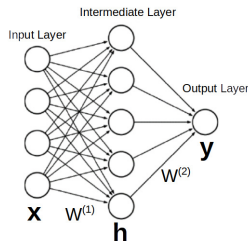


- Each intermediate layer computes a **nonlinear transformation** of its previous layer's nodes
- In traditional neural nets, \mathbf{h} is a **Linear transform** (e.g., $\mathbf{W}^{(1)}\mathbf{x}$ in the above picture) **followed by a nonlinearity** (e.g., sigmoid, ReLU, tanh, etc)
- Neural nets are awesome but brittle in many ways



Neural Networks

- A simple neural network with one intermediate (also called “hidden”) layer and a single output

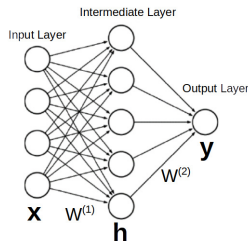


- Each intermediate layer computes a **nonlinear transformation** of its previous layer's nodes
- In traditional neural nets, \mathbf{h} is a **Linear transform** (e.g., $\mathbf{W}^{(1)}\mathbf{x}$ in the above picture) **followed by a nonlinearity** (e.g., sigmoid, ReLU, tanh, etc)
- Neural nets are awesome but brittle in many ways
 - Lots of parameters, difficult to train, need lots of data to train



Neural Networks

- A simple neural network with one intermediate (also called “hidden”) layer and a single output

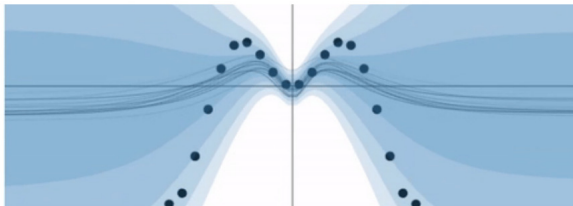


- Each intermediate layer computes a **nonlinear transformation** of its previous layer's nodes
- In traditional neural nets, \mathbf{h} is a **Linear transform** (e.g., $\mathbf{W}^{(1)}\mathbf{x}$ in the above picture) **followed by a nonlinearity** (e.g., sigmoid, ReLU, tanh, etc)
- Neural nets are awesome but brittle in many ways
 - Lots of parameters, difficult to train, need lots of data to train
 - Do not provide uncertainty estimates



What We Want..

- Neural networks with additional benefits of probabilistic/Bayesian modeling

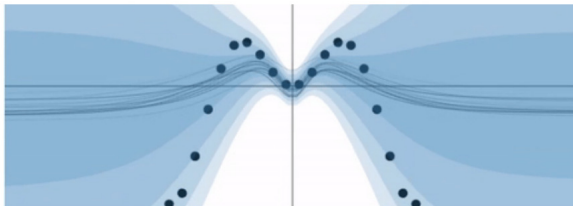


- Basically, nonlinear models with estimates of uncertainty in the model/its predictions



What We Want..

- Neural networks with additional benefits of probabilistic/Bayesian modeling

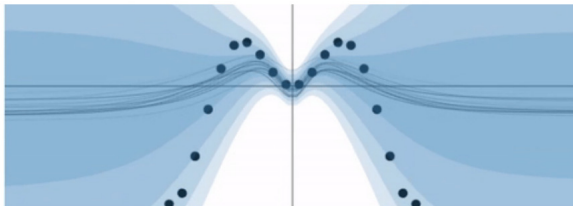


- Basically, nonlinear models with estimates of uncertainty in the model/its predictions
- Note: We already have seen something that accomplishes this - [Gaussian Processes](#)



What We Want..

- Neural networks with additional benefits of probabilistic/Bayesian modeling



- Basically, nonlinear models with estimates of uncertainty in the model/its predictions
- Note: We already have seen something that accomplishes this - [Gaussian Processes](#)
- Probabilistic/Bayesian neural nets are another alternative to this



Probabilistic/Bayesian Neural Networks

- A probabilistic model for neural network for supervised learning

$$y_n \sim \mathcal{N}(\text{NN}(\mathbf{x}_n; \mathbf{W}), \beta^{-1}) \quad (\text{for real-valued responses})$$

$$y_n \sim \text{Bernoulli}(\sigma(\text{NN}(\mathbf{x}_n; \mathbf{W}))) \quad (\text{for binary responses})$$

$$y_n \sim \text{ExpFam}(\text{NN}(\mathbf{x}_n; \mathbf{W})) \quad (\text{for general types of responses modeled by exp-family})$$

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



Probabilistic/Bayesian Neural Networks

- A probabilistic model for neural network for supervised learning

$$y_n \sim \mathcal{N}(\text{NN}(\mathbf{x}_n; \mathbf{W}), \beta^{-1}) \quad (\text{for real-valued responses})$$

$$y_n \sim \text{Bernoulli}(\sigma(\text{NN}(\mathbf{x}_n; \mathbf{W}))) \quad (\text{for binary responses})$$

$$y_n \sim \text{ExpFam}(\text{NN}(\mathbf{x}_n; \mathbf{W})) \quad (\text{for general types of responses modeled by exp-family})$$

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

- This enables learning probabilistic nonlinear input-to-output mappings



Probabilistic/Bayesian Neural Networks

- A probabilistic model for neural network for supervised learning

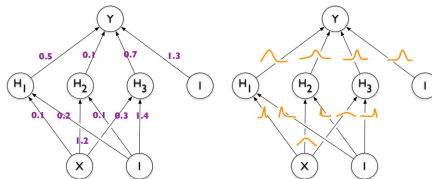
$$y_n \sim \mathcal{N}(\text{NN}(\mathbf{x}_n; \mathbf{W}), \beta^{-1}) \quad (\text{for real-valued responses})$$

$$y_n \sim \text{Bernoulli}(\sigma(\text{NN}(\mathbf{x}_n; \mathbf{W}))) \quad (\text{for binary responses})$$

$$y_n \sim \text{ExpFam}(\text{NN}(\mathbf{x}_n; \mathbf{W})) \quad (\text{for general types of responses modeled by exp-family})$$

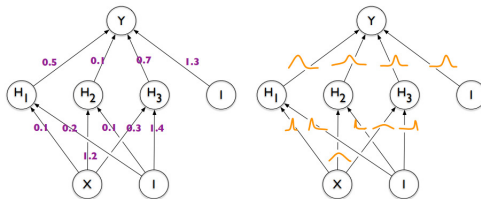
where $\text{NN}(\mathbf{x}_n; \mathbf{W})$ is a neural network with features \mathbf{x}_n as inputs and parameters \mathbf{W}

- This enables learning probabilistic nonlinear input-to-output mappings
- We can perform point estimation or fully Bayesian inference for such probabilistic neural networks



Left: Standard NN or NN with point estimation, Right: Bayesian Neural Network

Learning Bayesian Neural Networks

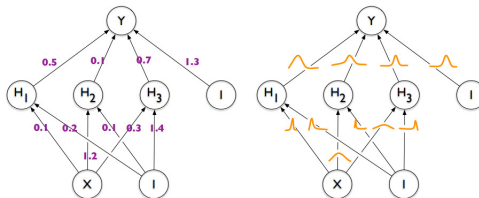


- Computing posterior over \mathbf{W} or computing posterior predictive is intractable in general

[†] “Weight Uncertainty in Neural Networks” (Blundell et al, 2015), [‡] “A Simple Baseline for Bayesian Uncertainty in Deep Learning” (Maddox et al, 2019), * “A Scalable Laplace Approximation for Neural Networks” (Ritter et al, 2018)



Learning Bayesian Neural Networks

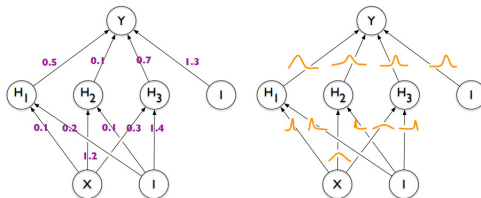


- Computing posterior over \mathbf{W} or computing posterior predictive is intractable in general
- A variety of methods exist to perform inference in such models

[†] “Weight Uncertainty in Neural Networks” (Blundell et al, 2015), [‡] “A Simple Baseline for Bayesian Uncertainty in Deep Learning” (Maddox et al, 2019), * “A Scalable Laplace Approximation for Neural Networks” (Ritter et al, 2018)



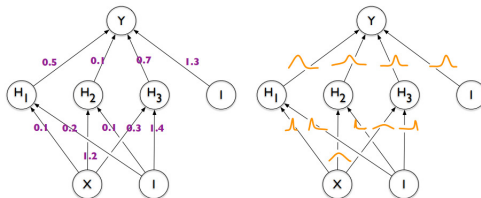
Learning Bayesian Neural Networks



- Computing posterior over \mathbf{W} or computing posterior predictive is intractable in general
- A variety of methods exist to perform inference in such models, e.g.,
 - VI based methods (e.g., BBVI, reparametrization trick, Bayes-by-Backprop[†], etc)

[†] “Weight Uncertainty in Neural Networks” (Blundell et al, 2015), [‡] “A Simple Baseline for Bayesian Uncertainty in Deep Learning” (Maddox et al, 2019), * “A Scalable Laplace Approximation for Neural Networks” (Ritter et al, 2018)

Learning Bayesian Neural Networks

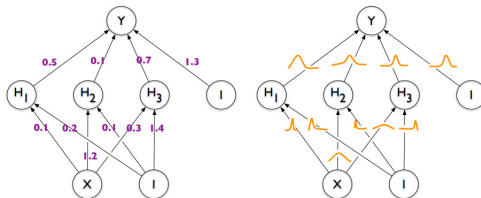


- Computing posterior over \mathbf{W} or computing posterior predictive is intractable in general
- A variety of methods exist to perform inference in such models, e.g.,
 - VI based methods (e.g., BBVI, reparametrization trick, Bayes-by-Backprop[†], etc)
 - Classic MCMC or Hamiltonian Monte Carlo (HMC), etc.

[†] "Weight Uncertainty in Neural Networks" (Blundell et al, 2015), [‡] "A Simple Baseline for Bayesian Uncertainty in Deep Learning" (Maddox et al, 2019), * "A Scalable Laplace Approximation for Neural Networks" (Ritter et al, 2018)



Learning Bayesian Neural Networks

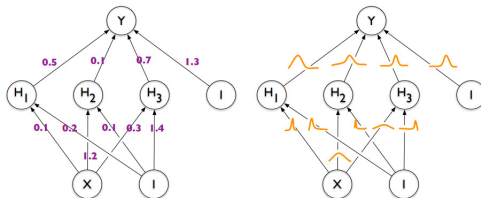


- Computing posterior over \mathbf{W} or computing posterior predictive is intractable in general
- A variety of methods exist to perform inference in such models, e.g.,
 - VI based methods (e.g., BBVI, reparametrization trick, Bayes-by-Backprop[†], etc)
 - Classic MCMC or Hamiltonian Monte Carlo (HMC), etc.
 - SGD-inspired methods (e.g., SGLD, online HMC, SWAG[‡])

[†] “Weight Uncertainty in Neural Networks” (Blundell et al, 2015), [‡] “A Simple Baseline for Bayesian Uncertainty in Deep Learning” (Maddox et al, 2019), * “A Scalable Laplace Approximation for Neural Networks” (Ritter et al, 2018)



Learning Bayesian Neural Networks



- Computing posterior over \mathbf{W} or computing posterior predictive is intractable in general
- A variety of methods exist to perform inference in such models, e.g.,
 - VI based methods (e.g., BBVI, reparametrization trick, Bayes-by-Backprop[†], etc)
 - Classic MCMC or Hamiltonian Monte Carlo (HMC), etc.
 - SGD-inspired methods (e.g., SGLD, online HMC, SWAG[‡])
 - Scalable versions of classic approximations, e.g., Laplace with diag/block-diag approx. of Hessian^{*}

[†] "Weight Uncertainty in Neural Networks" (Blundell et al, 2015), [‡] "A Simple Baseline for Bayesian Uncertainty in Deep Learning" (Maddox et al, 2019), ^{*} "A Scalable Laplace Approximation for Neural Networks" (Ritter et al, 2018)

Regularization in Deep Neural Nets: A Bayesian View

- Dropout is a popular way to regularize deep neural nets



Regularization in Deep Neural Nets: A Bayesian View

- Dropout is a popular way to regularize deep neural nets
- Basic idea: Perturb each NN weight or activations using multiplicative noise, e.g.,
 - Gaussian dropout: $\hat{w}_{ij} = w_{ij}\epsilon_{ij}$ where $\epsilon_{ij} \sim \mathcal{N}(1, \alpha)$
 - Bernoulli dropout: $\hat{w}_{ij} = w_{ij}\epsilon_{ij}$ where $\epsilon_{ij} \sim \text{Bernoulli}(\alpha)$



Regularization in Deep Neural Nets: A Bayesian View

- Dropout is a popular way to regularize deep neural nets
- Basic idea: Perturb each NN weight or activations using multiplicative noise, e.g.,
 - Gaussian dropout: $\hat{w}_{ij} = w_{ij}\epsilon_{ij}$ where $\epsilon_{ij} \sim \mathcal{N}(1, \alpha)$
 - Bernoulli dropout: $\hat{w}_{ij} = w_{ij}\epsilon_{ij}$ where $\epsilon_{ij} \sim \text{Bernoulli}(\alpha)$
- Assuming Gaussian dropout, we can estimate \mathbf{W} using stoch. grad., e.g.

$$\begin{aligned}\nabla_{\mathbf{W}} \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) &= \nabla_{\mathbf{W}} \log p(\mathbf{y}|\mathbf{X}, \mathbf{W}\epsilon) \\ \epsilon &\sim \mathcal{N}(\mathbf{1}, \alpha \mathbf{I})\end{aligned}$$



Regularization in Deep Neural Nets: A Bayesian View

- Dropout is a popular way to regularize deep neural nets
- Basic idea: Perturb each NN weight or activations using multiplicative noise, e.g.,
 - Gaussian dropout: $\hat{w}_{ij} = w_{ij}\epsilon_{ij}$ where $\epsilon_{ij} \sim \mathcal{N}(1, \alpha)$
 - Bernoulli dropout: $\hat{w}_{ij} = w_{ij}\epsilon_{ij}$ where $\epsilon_{ij} \sim \text{Bernoulli}(\alpha)$
- Assuming Gaussian dropout, we can estimate \mathbf{W} using stoch. grad., e.g.

$$\begin{aligned}\nabla_{\mathbf{W}} \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) &= \nabla_{\mathbf{W}} \log p(\mathbf{y}|\mathbf{X}, \mathbf{W}\epsilon) \\ \epsilon &\sim \mathcal{N}(\mathbf{1}, \alpha \mathbf{I})\end{aligned}$$

- The above stoch. grad. is the same as the stoch. grad of (verify using reparam trick)

$$\nabla_{\mathbf{W}} \int \mathcal{N}(\hat{\mathbf{W}}|\mathbf{W}, \alpha \mathbf{W}^2) \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) d\hat{\mathbf{W}} = \nabla_{\mathbf{W}} \mathbb{E}[\log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}})]$$

$$\text{where } \mathcal{N}(\hat{\mathbf{W}}|\mathbf{W}, \alpha \mathbf{W}^2) = \prod_{ij} \mathcal{N}(\hat{w}_{ij}|w_{ij}, \alpha w_{ij}^2)$$



Regularization in Deep Neural Nets: A Bayesian View

- So dropout is equivalent to optimizing the following w.r.t. \mathbf{W}

$$\int q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) d\hat{\mathbf{W}}$$

where $q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) = \prod_{ij} \mathcal{N}(\hat{w}_{ij}|w_{ij}, \alpha w_{ij}^2)$



Regularization in Deep Neural Nets: A Bayesian View

- So dropout is equivalent to optimizing the following w.r.t. \mathbf{W}

$$\int q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) d\hat{\mathbf{W}}$$

where $q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) = \prod_{ij} \mathcal{N}(\hat{w}_{ij}|w_{ij}, \alpha w_{ij}^2)$

- The above is basically ELBO with variational approx $q(\hat{\mathbf{W}}|\mathbf{W}, \alpha)$ but without the KL term

$$\int q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) d\hat{\mathbf{W}} - KL(q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) || p(\hat{\mathbf{W}}))$$



Regularization in Deep Neural Nets: A Bayesian View

- So dropout is equivalent to optimizing the following w.r.t. \mathbf{W}

$$\int q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) d\hat{\mathbf{W}}$$

where $q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) = \prod_{ij} \mathcal{N}(\hat{w}_{ij}|w_{ij}, \alpha w_{ij}^2)$

- The above is basically ELBO with variational approx $q(\hat{\mathbf{W}}|\mathbf{W}, \alpha)$ but without the KL term

$$\int q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) d\hat{\mathbf{W}} - KL(q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) || p(\hat{\mathbf{W}}))$$

- Gaussian dropout will be exactly equivalent to VI based on ELBO maximization if



Regularization in Deep Neural Nets: A Bayesian View

- So dropout is equivalent to optimizing the following w.r.t. \mathbf{W}

$$\int q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) d\hat{\mathbf{W}}$$

where $q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) = \prod_{ij} \mathcal{N}(\hat{w}_{ij}|w_{ij}, \alpha w_{ij}^2)$

- The above is basically ELBO with variational approx $q(\hat{\mathbf{W}}|\mathbf{W}, \alpha)$ but without the KL term

$$\int q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) d\hat{\mathbf{W}} - KL(q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) || p(\hat{\mathbf{W}}))$$

- Gaussian dropout will be exactly equivalent to VI based on ELBO maximization if
 - The prior $p(\hat{\mathbf{W}})$ is such that KL depends only on α



Regularization in Deep Neural Nets: A Bayesian View

- So dropout is equivalent to optimizing the following w.r.t. \mathbf{W}

$$\int q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) d\hat{\mathbf{W}}$$

where $q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) = \prod_{ij} \mathcal{N}(\hat{w}_{ij}|w_{ij}, \alpha w_{ij}^2)$

- The above is basically ELBO with variational approx $q(\hat{\mathbf{W}}|\mathbf{W}, \alpha)$ but without the KL term

$$\int q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) d\hat{\mathbf{W}} - KL(q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) || p(\hat{\mathbf{W}}))$$

- Gaussian dropout will be exactly equivalent to VI based on ELBO maximization if
 - The prior $p(\hat{\mathbf{W}})$ is such that KL depends only on α
 - α is some fixed constant



Regularization in Deep Neural Nets: A Bayesian View

- So dropout is equivalent to optimizing the following w.r.t. \mathbf{W}

$$\int q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) d\hat{\mathbf{W}}$$

where $q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) = \prod_{ij} \mathcal{N}(\hat{w}_{ij}|w_{ij}, \alpha w_{ij}^2)$

- The above is basically ELBO with variational approx $q(\hat{\mathbf{W}}|\mathbf{W}, \alpha)$ but without the KL term

$$\int q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) d\hat{\mathbf{W}} - KL(q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) || p(\hat{\mathbf{W}}))$$

- Gaussian dropout will be exactly equivalent to VI based on ELBO maximization if
 - The prior $p(\hat{\mathbf{W}})$ is such that KL depends only on α
 - α is some fixed constant
- Priors such as log-uniform prior $p(\hat{\mathbf{W}}) = \prod_{ij} p(\hat{w}_{ij}), p(\hat{w}_{ij}) \propto \frac{1}{|\hat{w}_{ij}|}$ have such a property



Regularization in Deep Neural Nets: A Bayesian View

- So dropout is equivalent to optimizing the following w.r.t. \mathbf{W}

$$\int q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) d\hat{\mathbf{W}}$$

where $q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) = \prod_{ij} \mathcal{N}(\hat{w}_{ij}|w_{ij}, \alpha w_{ij}^2)$

- The above is basically ELBO with variational approx $q(\hat{\mathbf{W}}|\mathbf{W}, \alpha)$ but without the KL term

$$\int q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{W}}) d\hat{\mathbf{W}} - KL(q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) || p(\hat{\mathbf{W}}))$$

- Gaussian dropout will be exactly equivalent to VI based on ELBO maximization if
 - The prior $p(\hat{\mathbf{W}})$ is such that KL depends only on α
 - α is some fixed constant
- Priors such as log-uniform prior $p(\hat{\mathbf{W}}) = \prod_{ij} p(\hat{w}_{ij}), p(\hat{w}_{ij}) \propto \frac{1}{|\hat{w}_{ij}|}$ have such a property
- Sparse variational dropout: If we use $q(\hat{\mathbf{W}}|\mathbf{W}, \alpha) = \prod_{ij} \mathcal{N}(\hat{w}_{ij}|w_{ij}, \alpha_{ij} w_{ij})$



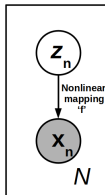
Constructing Generative Models using Neural Nets

- Probabilistic view enables neural net based **latent variable models** (“deep generative models”)



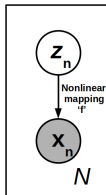
Constructing Generative Models using Neural Nets

- Probabilistic view enables neural net based **latent variable models** (“deep generative models”)
- Useful for **unsupervised learning**. Nonlinear latent variable to data mapping f modeled by NN



Constructing Generative Models using Neural Nets

- Probabilistic view enables neural net based **latent variable models** (“deep generative models”)
- Useful for **unsupervised learning**. Nonlinear latent variable to data mapping f modeled by NN



- Example: A probabilistic neural network for **latent variable modeling** (e.g., PPCA)

$$\mathbf{x}_n \sim \mathcal{N}(\text{NN}(\mathbf{z}_n; \mathbf{W}), \sigma^2 \mathbf{I}_D) \quad (\text{for real-valued features})$$

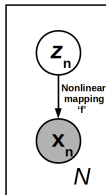
$$\mathbf{x}_n \sim \text{ExpFam}(\text{NN}(\mathbf{z}_n; \mathbf{W})) \quad (\text{for general types of features modeled by exp-family})$$

where $\text{NN}(\mathbf{z}_n; \mathbf{W})$ is a **neural network** with latent variables \mathbf{z}_n as inputs and parameters \mathbf{W}



Constructing Generative Models using Neural Nets

- Probabilistic view enables neural net based **latent variable models** (“deep generative models”)
- Useful for **unsupervised learning**. Nonlinear latent variable to data mapping f modeled by NN

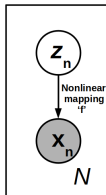


- Example: A probabilistic neural network for **latent variable modeling** (e.g., PPCA)
 - $\mathbf{x}_n \sim \mathcal{N}(\text{NN}(\mathbf{z}_n; \mathbf{W}), \sigma^2 \mathbf{I}_D)$ (for real-valued features)
 - $\mathbf{x}_n \sim \text{ExpFam}(\text{NN}(\mathbf{z}_n; \mathbf{W}))$ (for general types of features modeled by exp-family)
- where $\text{NN}(\mathbf{z}_n; \mathbf{W})$ is a **neural network** with latent variables \mathbf{z}_n as inputs and parameters \mathbf{W}
- The NN enables learning a nonlinear latent variable to data mapping f



Constructing Generative Models using Neural Nets

- Probabilistic view enables neural net based **latent variable models** (“deep generative models”)
- Useful for **unsupervised learning**. Nonlinear latent variable to data mapping f modeled by NN



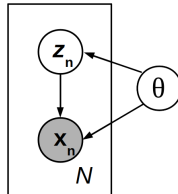
- Example: A probabilistic neural network for **latent variable modeling** (e.g., PPCA)
$$\mathbf{x}_n \sim \mathcal{N}(\text{NN}(\mathbf{z}_n; \mathbf{W}), \sigma^2 \mathbf{I}_D) \quad (\text{for real-valued features})$$
$$\mathbf{x}_n \sim \text{ExpFam}(\text{NN}(\mathbf{z}_n; \mathbf{W})) \quad (\text{for general types of features modeled by exp-family})$$

where $\text{NN}(\mathbf{z}_n; \mathbf{W})$ is a **neural network** with latent variables \mathbf{z}_n as inputs and parameters \mathbf{W}
- The NN enables learning a nonlinear latent variable to data mapping f
- If \mathbf{z}_n has a Gaussian prior, such models are called “**Deep Latent Gaussian Models**” (DLGM)



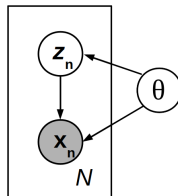
Inference for Deep Latent Gaussian Models

- Assume θ to be the global parameters of the model (params defining $p(\mathbf{z})$, $p(\mathbf{x}|\mathbf{z})$, etc.)



Inference for Deep Latent Gaussian Models

- Assume θ to be the global parameters of the model (params defining $p(\mathbf{z})$, $p(\mathbf{x}|\mathbf{z})$, etc.)

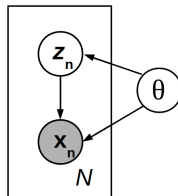


- The usual approach for inference in such models (as in most Bayesian models) is iterative



Inference for Deep Latent Gaussian Models

- Assume θ to be the global parameters of the model (params defining $p(\mathbf{z})$, $p(\mathbf{x}|\mathbf{z})$, etc.)

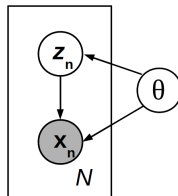


- The usual approach for inference in such models (as in most Bayesian models) is iterative, e.g.,
- Initialize θ . Then iterate until convergence



Inference for Deep Latent Gaussian Models

- Assume θ to be the global parameters of the model (params defining $p(\mathbf{z})$, $p(\mathbf{x}|\mathbf{z})$, etc.)

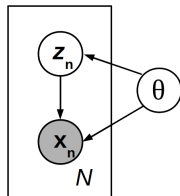


- The usual approach for inference in such models (as in most Bayesian models) is iterative, e.g.,
- Initialize θ . Then iterate until convergence
 - For $n = 1, \dots, N$
 - Infer $p(\mathbf{z}_n|\mathbf{x}_n)$ using MCMC. If doing VB, update variational parameters ϕ_n of $q(\mathbf{z}_n|\phi_n)$



Inference for Deep Latent Gaussian Models

- Assume θ to be the global parameters of the model (params defining $p(\mathbf{z})$, $p(\mathbf{x}|\mathbf{z})$, etc.)

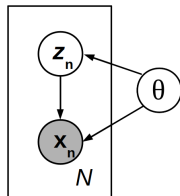


- The usual approach for inference in such models (as in most Bayesian models) is iterative, e.g.,
- Initialize θ . Then iterate until convergence
 - For $n = 1, \dots, N$
 - Infer $p(z_n|x_n)$ using MCMC. If doing VB, update variational parameters ϕ_n of $q(z_n|\phi_n)$
 - Infer θ (its full posterior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



Inference for Deep Latent Gaussian Models

- Assume θ to be the global parameters of the model (params defining $p(\mathbf{z})$, $p(\mathbf{x}|\mathbf{z})$, etc.)

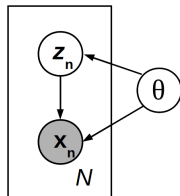


- The usual approach for inference in such models (as in most Bayesian models) is iterative, e.g.,
- Initialize θ . Then iterate until convergence
 - For $n = 1, \dots, N$
 - Infer $p(z_n|x_n)$ using MCMC. If doing VB, update variational parameters ϕ_n of $q(z_n|\phi_n)$
 - Infer θ (its full posterior using MCMC or VB, or a point estimate)
- This iterative approach can be slow for large N



Inference for Deep Latent Gaussian Models

- Assume θ to be the global parameters of the model (params defining $p(\mathbf{z})$, $p(\mathbf{x}|\mathbf{z})$, etc.)

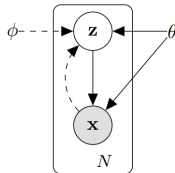


- The usual approach for inference in such models (as in most Bayesian models) is iterative, e.g.,
- Initialize θ . Then iterate until convergence
 - For $n = 1, \dots, N$
 - Infer $p(\mathbf{z}_n|\mathbf{x}_n)$ using MCMC. If doing VB, update variational parameters ϕ_n of $q(\mathbf{z}_n|\phi_n)$
 -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



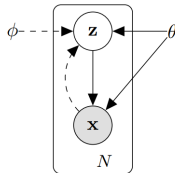
Variational Auto-encoder (VAE)

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

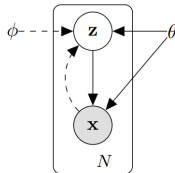


- VAE uses VB for inference but has a **fast, non-iterative** way of computing \mathbf{z}_n for a data point \mathbf{x}_n



Variational Auto-encoder (VAE)

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



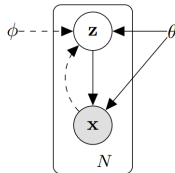
- VAE uses VB for inference but has a **fast, non-iterative** way of computing \mathbf{z}_n for a data point \mathbf{x}_n
- Key idea: For each point \mathbf{x}_n , instead of learning a separate $q(\mathbf{z}_n|\phi_n)$ with local params ϕ_n , assume

$$q(\mathbf{z}_n|\phi_n) = q(\mathbf{z}_n|\text{NN}(\mathbf{x}_n; \phi))$$



Variational Auto-encoder (VAE)

- Esseentially a DLGM, i.e., the \mathbf{z} to \mathbf{x} mapping $p(\mathbf{x}|\mathbf{z})$ 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 \mathbf{z}_n for a data point \mathbf{x}_n
- Key idea: For each point \mathbf{x}_n , instead of learning a separate $q(\mathbf{z}_n|\phi_n)$ with local params ϕ_n , assume

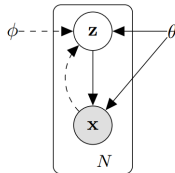
$$q(\mathbf{z}_n|\phi_n) = q(\mathbf{z}_n|\text{NN}(\mathbf{x}_n; \phi))$$

so, basically, each ϕ_n is computed by a neural net with global parameters ϕ and input \mathbf{x}_n



Variational Auto-encoder (VAE)

- Esseentially a DLGM, i.e., the \mathbf{z} to \mathbf{x} mapping $p(\mathbf{x}|\mathbf{z})$ 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 \mathbf{z}_n for a data point \mathbf{x}_n
- Key idea: For each point \mathbf{x}_n , instead of learning a separate $q(\mathbf{z}_n|\phi_n)$ with local params ϕ_n , assume

$$q(\mathbf{z}_n|\phi_n) = q(\mathbf{z}_n|\text{NN}(\mathbf{x}_n; \phi))$$

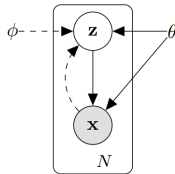
so, basically, each ϕ_n is computed by a neural net with global parameters ϕ and input \mathbf{x}_n

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



Variational Auto-encoder (VAE)

- Esseentially a DLGM, i.e., the \mathbf{z} to \mathbf{x} mapping $p(\mathbf{x}|\mathbf{z})$ 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 \mathbf{z}_n for a data point \mathbf{x}_n
- Key idea: For each point \mathbf{x}_n , instead of learning a separate $q(\mathbf{z}_n|\phi_n)$ with local params ϕ_n , assume

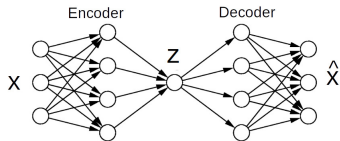
$$q(\mathbf{z}_n|\phi_n) = q(\mathbf{z}_n|\text{NN}(\mathbf{x}_n; \phi))$$

so, basically, each ϕ_n is computed by a neural net with global parameters ϕ and input \mathbf{x}_n

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

Standard Auto-encoder vs Variational Auto-encoder

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

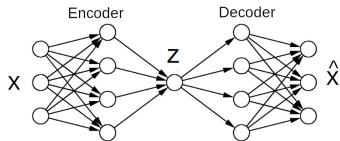


- Model is trained to minimize the **reconstruction error** (difference b/w x and \hat{x})



Standard Auto-encoder vs Variational Auto-encoder

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

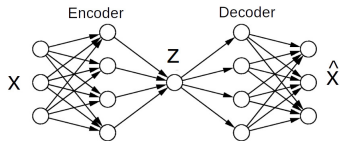


- Model is trained to minimize the **reconstruction error** (difference b/w x and \hat{x})
- However, it can't "generate" a "realistic" input from a random z (the model isn't trained for that)

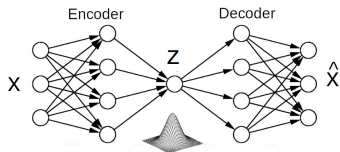


Standard Auto-encoder vs Variational Auto-encoder

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

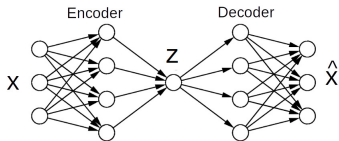


- Model is trained to minimize the **reconstruction error** (difference b/w x and \hat{x})
- However, it can't "generate" a "realistic" input from a random z (the model isn't trained for that)
- VAE allows this by assuming a distribution (e.g., Gaussian) over z and learning to generate x from random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

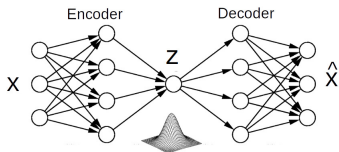


Standard Auto-encoder vs Variational Auto-encoder

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



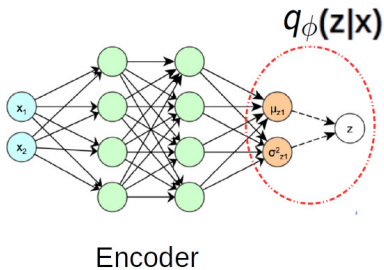
- Model is trained to minimize the **reconstruction error** (difference b/w x and \hat{x})
- However, it can't "generate" a "realistic" input from a random z (the model isn't trained for that)
- VAE allows this by assuming a distribution (e.g., Gaussian) over z and learning to generate x from random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

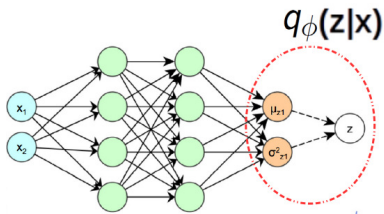
VAE: The Encoder

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



VAE: The Encoder

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



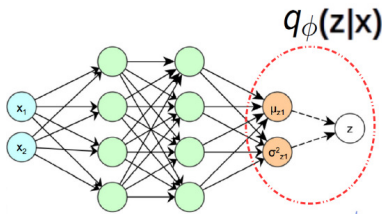
Encoder

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



VAE: The Encoder

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



Encoder

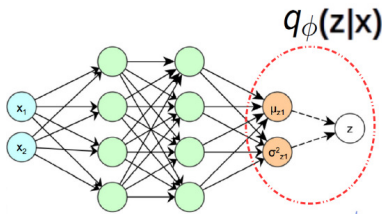
- Unlike standard autoencoders, for each \mathbf{x} , VAE gives us a distribution $q(\mathbf{z}|\mathbf{x})$ over its encoding
- Assume $q(\mathbf{z}|\mathbf{x})$ to be Gaussian whose mean/var are computed by a NN with global params ϕ

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



VAE: The Encoder

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



Encoder

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

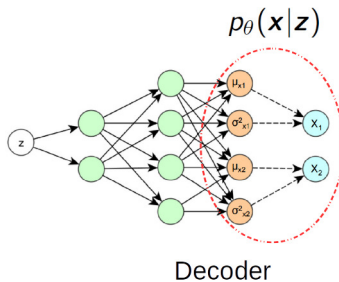
$$\mu_z = \text{NN}(\mathbf{x}; \phi) \quad \sigma_z^2 = \text{NN}(\mathbf{x}; \phi)$$

- Since μ_z, σ_z are outputs of neural networks, the \mathbf{x} to \mathbf{z} mapping is nonlinear



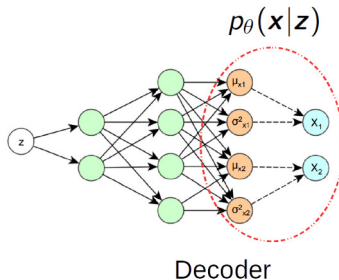
VAE: The Decoder

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



VAE: The Decoder

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

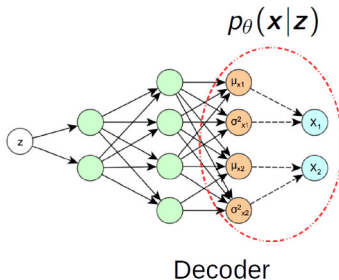


- Unlike PPCA, the \mathbf{z} to \mathbf{x} mapping is nonlinear (modeled by a neural network)



VAE: The Decoder

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



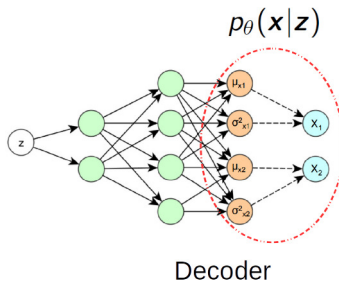
- Unlike PPCA, the \mathbf{z} to \mathbf{x} mapping is nonlinear (modeled by a neural network)
- Assume $p(\mathbf{x}|\mathbf{z})$ to be Gaussian whose mean/var are computed by a NN with global params θ

$$\mu_{\mathbf{x}} = \text{NN}(\mathbf{z}; \theta) \quad \sigma_{\mathbf{x}}^2 = \text{NN}(\mathbf{z}; \theta)$$



VAE: The Decoder

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



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

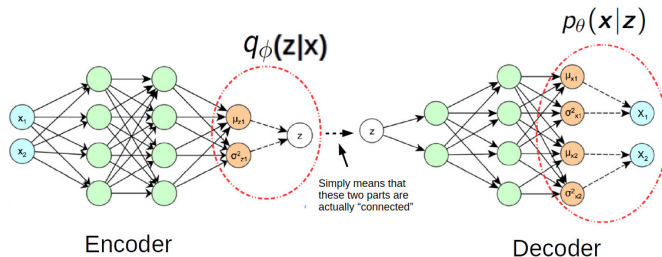
$$\mu_{\mathbf{x}} = \text{NN}(\mathbf{z}; \theta) \quad \sigma_{\mathbf{x}}^2 = \text{NN}(\mathbf{z}; \theta)$$

- Thus in the VAE, both \mathbf{x} to \mathbf{z} (encoder) and \mathbf{z} to \mathbf{x} (decoder) mappings are nonlinear



Inference for VAE

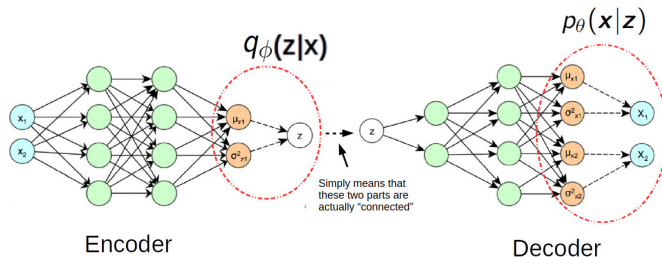
- VAE uses variational inference (hence the name!) to learn the model parameters θ and ϕ



[†] "Auto-encoding Variational Bayes" (Kingma and Welling, 2013)

Inference for VAE

- VAE uses variational inference (hence the name!) to learn the model parameters θ and ϕ



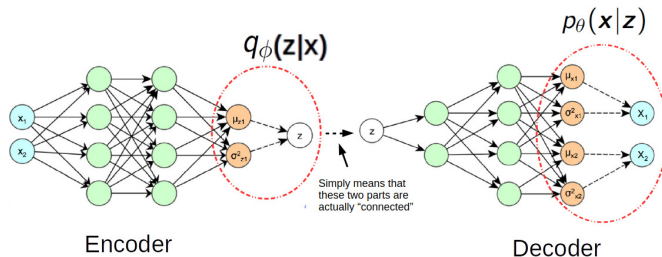
- Typically a prior $p(z) = \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$ is assumed on z . The ELBO for a single \mathbf{x}_n will be

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

[†] "Auto-encoding Variational Bayes" (Kingma and Welling, 2013)

Inference for VAE

- VAE uses variational inference (hence the name!) to learn the model parameters θ and ϕ



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

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

- Variational inference uses the reparametrization trick[†] for computing ELBO derivatives

[†] "Auto-encoding Variational Bayes" (Kingma and Welling, 2013)

Some Other Architectures: Deep Exponential Families

- Standard VAE has only one layer of latent z and a neural net to transform z into x



Some Other Architectures: Deep Exponential Families

- Standard VAE has only one layer of latent \mathbf{z} and a neural net to transform \mathbf{z} into \mathbf{x}
- Many other deep architectures have multiple layers of latent variables



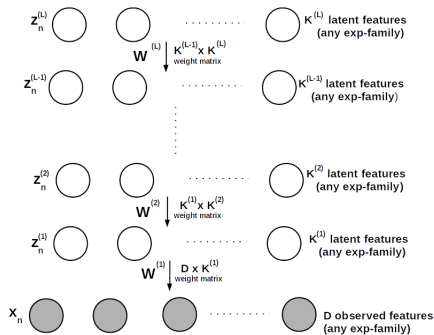
Some Other Architectures: Deep Exponential Families

- Standard VAE has only one layer of latent z and a neural net to transform z into x
- Many other deep architectures have multiple layers of latent variables
- Deep Exponential Family (DEF) is one such recently proposed popular model



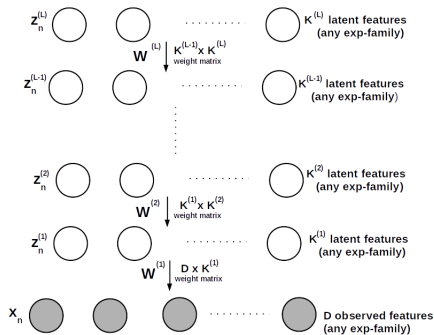
Some Other Architectures: Deep Exponential Families

- Standard VAE has only one layer of latent z and a neural net to transform z into x
- Many other deep architectures have multiple layers of latent variables
- Deep Exponential Family (DEF) is one such recently proposed popular model
- in DEF, latent variables in every layer, as well as observations, are from exp. family distributions



Some Other Architectures: Deep Exponential Families

- Standard VAE has only one layer of latent z and a neural net to transform z into x
- Many other deep architectures have multiple layers of latent variables
- Deep Exponential Family (DEF) is one such recently proposed popular model
- in DEF, 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

Learning the right size of a deep neural network

- How to decide the number of layers and width of each layer?
- Nonparametric Bayesian methods can help here



Learning the right size of a deep neural network

- How to decide the number of layers and width of each layer?
- Nonparametric Bayesian methods can help here



- A [cascaded Indian Buffet Prior](#) can model the relationships between nodes in adjacent layers



Learning the right size of a deep neural network

- How to decide the number of layers and width of each layer?
- Nonparametric Bayesian methods can help here



- A **cascaded Indian Buffet Prior** can model the relationships between nodes in adjacent layers
 - The bottom-most layer is the data layer (fixed/known size)
 - Width of each intermediate layer and active connections can be inferred by the IBP prior



Learning the right size of a deep neural network

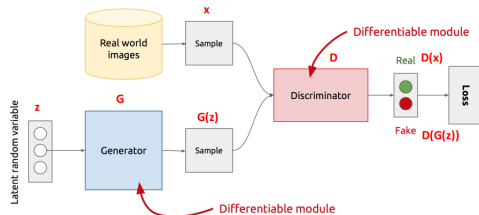
- How to decide the number of layers and width of each layer?
- Nonparametric Bayesian methods can help here



- A **cascaded Indian Buffet Prior** can model the relationships between nodes in adjacent layers
 - The bottom-most layer is the data layer (fixed/known size)
 - Width of each intermediate layer and active connections can be inferred by the IBP prior
- Another option is to use sparsity inducing priors on the connection weights

Generative Adversarial Networks (GAN)

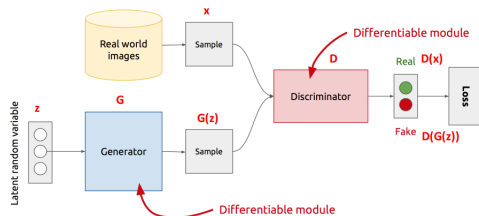
- Based on a game between a generator and a discriminator (Goodfellow et al, 2013)[†]



[†] Generative Adversarial Nets (Goodfellow et al, 2013), Figure: <https://www.slideshare.net/xavigiro/deep-learning-for-computer-vision-generative-models-and-adversarial-training-upc-2016>

Generative Adversarial Networks (GAN)

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



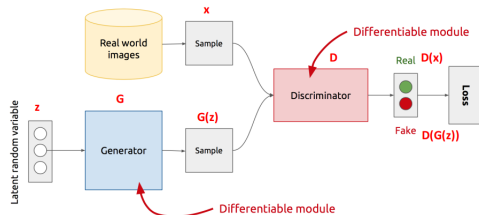
- Can be thought of as a two-player minimax game

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

[†] Generative Adversarial Nets (Goodfellow et al, 2013), Figure: <https://www.slideshare.net/xavigiro/deep-learning-for-computer-vision-generative-models-and-adversarial-training-upc-2016>

Generative Adversarial Networks (GAN)

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



- Can be thought of as a two-player minimax game

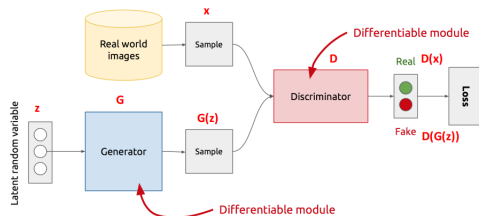
$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

- With the generator G fixed, the optimal discriminator $D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$

[†] Generative Adversarial Nets (Goodfellow et al, 2013), Figure: <https://www.slideshare.net/xavigiro/deep-learning-for-computer-vision-generative-models-and-adversarial-training-upc-2016>

Generative Adversarial Networks (GAN)

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



- Can be thought of as a two-player minimax game

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

- With the generator G fixed, the optimal discriminator $D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$
- At the global minimum of the objective, $p_g = p_{data}$

[†] Generative Adversarial Nets (Goodfellow et al, 2013), Figure: <https://www.slideshare.net/xavigiro/deep-learning-for-computer-vision-generative-models-and-adversarial-training-upc-2016>

Summary

- Probabilistic modeling allows developing very flexible deep learning 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



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)

