# Proabilistic/Bayesian Models for Deep Learning
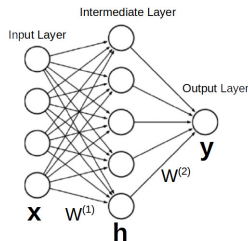
Piyush Rai

Topics in Probabilistic Modeling and Inference (CS698X)

April 12, 2018

# Neural Networks
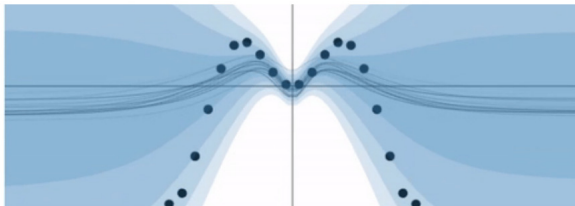
- A simple neural network with one intermediate (also called "hidden") layer and a single output



- Each intermediate layer computes a nonlinear transformation of its previous layer's nodes
- In traditional neural nets, $\boldsymbol{h}$ is a Linear transform (e.g., $\mathbf{W}^{(1)}\boldsymbol{x}$ in the above picture) followed by a nonlinearlity (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
- Note: We already have seen something that accomplishes this - Gaussian Processes
- Probabilistic/Bayesian neural nets are another alternative to this

## Neural Networks as Probabilistic Models

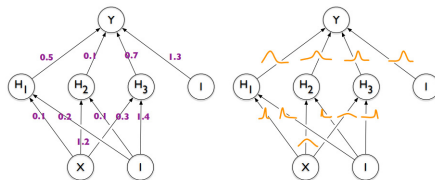- A probabilistic model for neural network for supervised learning

$$y_n \sim \mathcal{N}(\text{NN}(x_n; W), \beta^{-1}) \quad \text{(for real-valued responses)}$$
$$y_n \sim \text{Bernoulli}(\sigma(\text{NN}(x_n; W))) \quad \text{(for binary responses)}$$
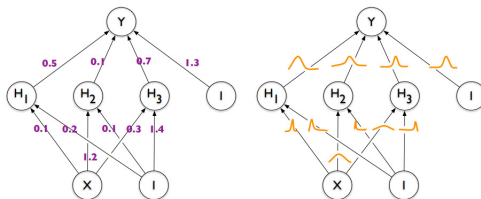$$y_n \sim \text{ExpFam}(\text{NN}(x_n; W)) \quad \text{(for general types of responses modeled by exp-family)}$$

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

- 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

"Weight Uncertainty in Neural Networks" (Blundell et al, 2015)
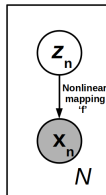
# Learning Bayesian Neural Networks



- Even if the prior $p(\mathbf{W})$ on NN weights is Gaussian, the model is not conjugate
- MCMC methods can be used to learn the posterior $p(\mathbf{W}|\mathcal{D})$ but can be slow
  - However, methods such as SGLD[†] allow efficient MCMC inference for such models (recall that SGLD only requires gradient expressions of the log-joint probability $\log p(\mathcal{D}, \mathbf{W})$ of the model)
- Variational inference is another popular alternative to MCMC for such models
  - But ELBO is intractable (due to non-conjugacy); has to be approximated via Monte-Carlo/BBVI
- Note: Hybrid architectures also possible (only last layer modeled in a fully Bayesian way)

[†] "Preconditioned Stochastic Gradient Langevin Dynamics for Deep Neural Networks" (Li et al, 2016)

# Constructing Generative Models using Neural Nets

- Probabilistic view enables neural net based latent variable models ("deep generative models")
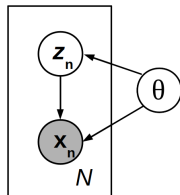- Useful for unsupervised learning. <u>Nonlinear</u> latent variable to data mapping $f$ modeled by NN



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

$$\boldsymbol{x}_n \sim \mathcal{N}(\text{NN}(\boldsymbol{z}_n; \mathbf{W}), \sigma^2 \mathbf{I}_D) \qquad \text{(for real-valued features)}$$
$$\boldsymbol{x}_n \sim \text{ExpFam}(\text{NN}(\boldsymbol{z}_n; \mathbf{W})) \qquad \text{(for general types of features modeled by exp-family)}$$

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

- The NN enables learning a nonlinear latent variable to data mapping $f$
- If $\boldsymbol{z}_n$ has a Gaussian prior, such models are called "Deep Latent Gaussian Models" (DLGM)
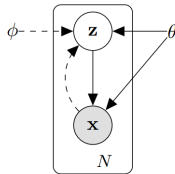
## Inference for Deep Latent Gaussian Models

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



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



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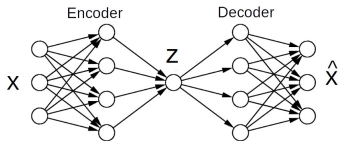- Key idea: For each point $x_n$, instead of learning a separate $q(z_n|\phi_n)$ with local params $\phi_n$, assume

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

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

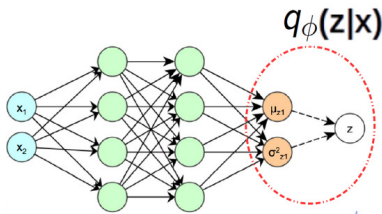- A standard auto-encoder learns to (nonlinearly) compress and uncompress an input



- Model is trained to minimize the reconstruction error (difference b/w $x$ and $\hat{x}$)
- 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
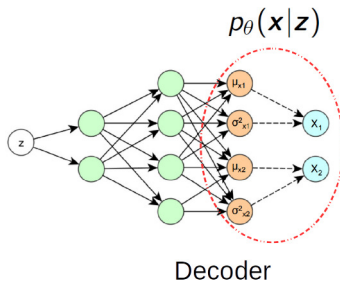
- Unlike standard autoencoders, for each $x$, VAE gives us a distribution $q(z|x)$ over its encoding
- Assume $q(z|x)$ to be Gaussian whose mean/var are computed by a NN with global params $\phi$

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

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

## VAE: The Decoder

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

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



Decoder
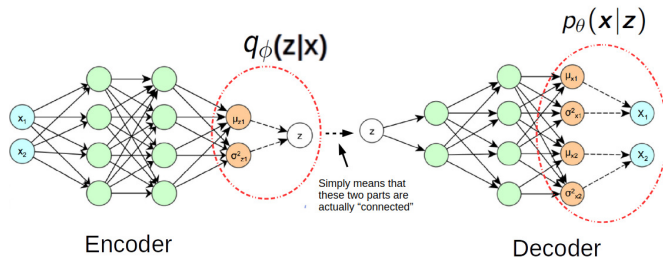
- Unlike PPCA, the $z$ to $x$ mapping is nonlinear (modeled by a neural network)

- 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 $\phi$



$q_\phi(z|x)$

$p_\theta(x|z)$

Simply means that these two parts are actually "connected"

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

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

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

---

[†] "Auto-encoding Variational Bayes" (Kingman 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 $x$
- 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

# Some 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.

https://www.slideshare.net/xavigiro/deep-learning-for-computer-vision-generative-models-and-adversarial-training-upc-2016

## 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

"Learning the Structure of Deep Sparse Graphical Models" (Adams et al, 2009), "Model Selection in Bayesian Neural Networks via Horseshoe Priors", Ghosh and Doshi-Velez (2017)

## Neural Nets vs Gaussian Processes

- Both can be used learn nonlinear input to output mappings, e.g.,

$$
\begin{aligned}
y_n &\sim \mathcal{N}(\mathsf{NN}(\boldsymbol{x}_n; \mathbf{W}), \beta^{-1}) \\
y_n &\sim \mathcal{N}(f(\boldsymbol{x}_n), \beta^{-1}) \qquad \text{where } f \sim \mathsf{GP}
\end{aligned}
$$

- Both have their pros and cons

- NN pros: Fast to train (e.g., using SGD methods) and also fast at test time

- NN cons: Difficult to train, also not Bayesian (but can be made Bayesian)

- GP pros: Simple formulation, especially for regression settings; natively Bayesian in formulation

- GP cons: Slow to train and also slow at test time

- Nowadays Bayesian NN and GPs are competitive in many applications

# 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)