

# Lecture 13

## Bayesian Neural Networks and Variational Autoencoder

Instructor: Shibo Li

shiboli@cs.fsu.edu



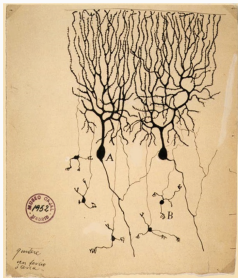
Department of Computer Science  
Florida State University

- Neural networks and Back-propagation
- Stochastic optimization
- Bayesian neural networks
- Bayes by Backprop and reparameterization trick
- Auto-encoding variational Bayes
- Generative adversarial networks



- Neural networks and Back-propagation
- Stochastic optimization
- Bayesian neural networks
- Bayes by Backprop and reparameterization trick
- Auto-encoding variational Bayes
- Generative adversarial networks

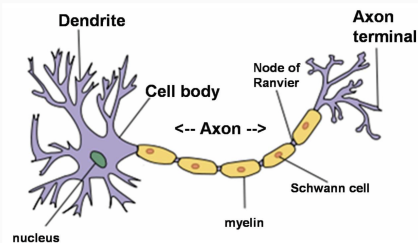
- 1943: McCullough and Pitts showed how linear threshold units can compute logical functions
- 1949: Hebb suggested a learning rule that has some physiological plausibility
- 1950s: Rosenblatt, the Peceptron algorithm for a single threshold neuron
- 1969: Minsky and Papert studied the neuron from a geometrical perspective
- 1980s: Convolutional neural networks (Fukushima, LeCun), the backpropagation algorithm (various)
- 2003-today: More compute, more data, deeper networks

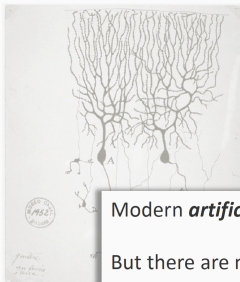


The first drawing of a brain cells by Santiago Ramón y Cajal in 1899

**Neurons:** core components of brain and the nervous system consisting of

1. Dendrites that collect information from other neurons
2. An axon that generates outgoing spikes





The first drawing of neurons by Santiago Ramón y Cajal in 1892

**Neurons:** core components of brain and the nervous system consisting of

1. Dendrites that collect information from other neurons
2. An axon that generates outgoing spikes

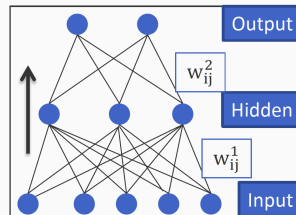
Modern **artificial** neurons are “inspired” by biological neurons

But there are many, many fundamental differences

Don't take the similarity seriously (as also claims in the news about the “emergence” of intelligent behavior)

A function that converts inputs to outputs defined by a **directed acyclic graph**

- Nodes organized in layers, correspond to neurons
- Edges carry output of one neuron to another, associated with weights



- To define a neural network, we need to specify:

- The structure of the graph
  - How many nodes, the connectivity
- The activation function on each node

Called the **architecture** of the network

Typically predefined, part of the design of the classifier

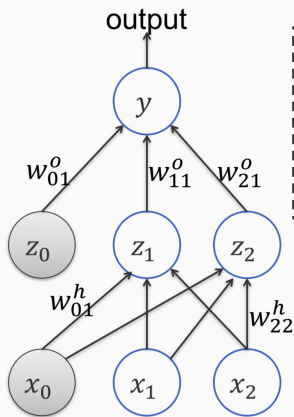
- The edge weights

Learned from data

$$\text{output} = \text{activation}(\mathbf{w}^T \mathbf{x} + b)$$

Name of the neuron	Activation function: $\text{activation}(z)$
Linear unit	$z$
Threshold/sign unit	$\text{sgn}(z)$
Sigmoid unit	$\frac{1}{1 + \exp(-z)}$
Rectified linear unit (ReLU)	$\max(0, z)$
Tanh unit	$\tanh(z)$

Many more activation functions exist (sinusoid, sinc, Gaussian, polynomial...)



Given an input  $\mathbf{x}$ , how is the output predicted

$$\text{output } y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$$

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

$$z_1 = \sigma(w_{01}^h + w_{11}^h x_1 + w_{21}^h x_2)$$

Suppose the true label for this example is a number  $y^*$

We can write the *square loss* for this example as:

$$L = \frac{1}{2} (y - y^*)^2$$

## An L-layer NN



$$\mathbf{x}_0 \rightarrow \mathbf{x}_1 \rightarrow \dots \rightarrow \mathbf{x}_{L-1} \rightarrow \mathbf{x}_L$$

Arbitrary element-wise activation function

$$\mathbf{x}_j = \sigma(\mathbf{W}_j \mathbf{x}_{j-1}) \quad (1 \leq j \leq L-1) \quad \text{middle layer}$$

$$\mathbf{f}_{\text{out}} = \mathbf{x}_L = \mathbf{W}_L \mathbf{x}_{L-1} \quad \text{output layer}$$

$$\mathbf{W}_j : d_j \times d_{j-1}$$

$d_0$  : input dim.

$d_L$  : output dim.



$$\mathbf{x}_0 \rightarrow \mathbf{x}_1 \rightarrow \dots \mathbf{x}_{L-1} \rightarrow \mathbf{x}_L$$

$$\mathbf{x}_j = \sigma(\mathbf{W}_j \mathbf{x}_{j-1}) \quad (1 \leq j \leq L-1) \quad \text{Middle layer}$$

$$\mathbf{f}_{\text{out}} = \mathbf{x}_L = \mathbf{W}_L \mathbf{x}_{L-1} \quad \text{output layer}$$

We can also recursively write

$$\mathbf{f}_{\mathcal{W}}(\mathbf{x}_0) = \mathbf{f}_{\text{out}} = \mathbf{W}_L \sigma(\mathbf{W}_{L-1} \sigma(\dots \sigma(\mathbf{W}_1 \mathbf{x}_0)))$$

$$\mathcal{W} = \{\mathbf{W}_1, \dots, \mathbf{W}_L\}$$

- To compute the output, you need to start from the bottom level and sequentially pass each layer

$$\mathbf{x}_0 \rightarrow \mathbf{x}_1 \rightarrow \dots \mathbf{x}_{L-1} \rightarrow \mathbf{x}_L$$

This is called forward pass

In general, training NN is to minimize a loss function  $\mathcal{L}(\mathcal{W}, \mathcal{D})$  where  $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$

For example, square loss:

$$\mathcal{L}(\mathcal{W}, \mathcal{D}) = \frac{1}{N} \sum_{n=1}^N [y^{(n)} - f_{\mathcal{W}}(\mathbf{x}^{(n)})]^2$$

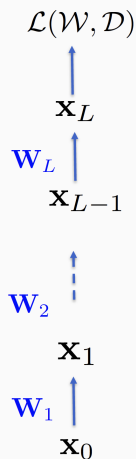
In general, training NN is to minimize a loss function  $\mathcal{L}(\mathcal{W}, \mathcal{D})$  where  $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$

e.g.,  $\mathcal{L}(\mathcal{W}, \mathcal{D}) = \frac{1}{N} \sum_{n=1}^N [y^{(n)} - f_{\mathcal{W}}(\mathbf{x}^{(n)})]^2$

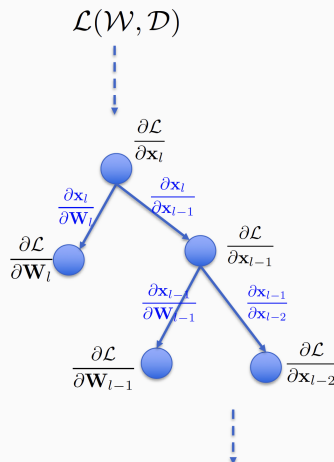
$f_{\mathcal{W}}(\mathbf{x}_0)$



How to efficiently compute gradient?  
Do it in backward!



from the root



- We will not discuss the detail because
  - It is trivial and mechanical
  - Nowadays, you never need to implement BP by yourself. TensorFlow, PyTorch, ... will do this automatically for you

- Neural networks and Back-propagation
- **Stochastic optimization**
- Bayesian neural networks
- Bayes by Backprop and reparameterization trick
- Auto-encoding variational Bayes
- General adversarial networks

- Suppose we aim to optimize an objective function that can be viewed as an expectation

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{p(u)}[g(\boldsymbol{\theta}, u)]$$

- Then we can compute a stochastic gradient for stochastic optimization

$$\nabla \mathcal{L}(\boldsymbol{\theta}) = \nabla \mathbb{E}_{p(u)}[g(\boldsymbol{\theta}, u)] = \mathbb{E}_{p(u)}[\nabla g(\boldsymbol{\theta}, u)]$$

under certainty conditions



- Suppose we aim to optimize an objective function that can be viewed as an expectation

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{p(u)}[g(\boldsymbol{\theta}, u)]$$

- Then we can compute a stochastic gradient for stochastic optimization

$$\nabla \mathcal{L}(\boldsymbol{\theta}) = \nabla \mathbb{E}_{p(u)}[g(\boldsymbol{\theta}, u)] = \mathbb{E}_{p(u)}[\nabla g(\boldsymbol{\theta}, u)]$$

under certainty conditions

- 1. Initialize  $\theta$  randomly (or 0)
- 2. For  $t = 1..T$ 
  - Sample  $u$  from  $p(u)$
  - Calculate stochastic gradient  $\nabla g(\theta, u)$
  - Update  $\theta \leftarrow \theta - \gamma_t \nabla g(\theta, u)$
- 3. Return  $\theta$

$\gamma_t$ : learning rate, many  
tweaks possible

With enough iterations, it will converge almost surely (i.e., with probability one)

Provided the step sizes are “*square summable, but not summable*”

- Step sizes  $\gamma_t$  are positive
  - Sum of squares of step sizes over  $t = 1$  to  $\infty$  is not infinite
  - Sum of step sizes over  $t = 1$  to  $\infty$  is infinity
- 
- Some examples:  $\gamma_t = \frac{\gamma_0}{1 + \frac{\gamma_0 t}{c}}$  or  $\gamma_t = \frac{\gamma_0}{1+t}$

- Learning rate is critical to convergence rate
- There are many works that develop learning rate schedules
- The main-stream is momentum-based approaches
- Most popular approaches include ADAM, Adagrad, Adadelta, etc.
- There are well developed libraries, and you do not need to implement them by yourself.

- It is the foundation of modern NN training

$$\mathcal{L}(\mathcal{W}, \mathcal{D}) = \sum_{n=1}^N \mathcal{L}(\mathcal{W}, \mathbf{x}_n, y_n)$$

- If we partition the training data into mini-batches  $\{B_1, B_2, \dots\}$  and each with size  $B$  (e.g., 100)

$$\begin{aligned}\mathcal{L}(\mathcal{W}, \mathcal{D}) &= \sum_{u=1}^{N/B} \frac{B}{N} \sum_{n \in B_u} \frac{N}{B} \mathcal{L}(\mathcal{W}, \mathbf{x}_n, y_n) \\ &= \mathbb{E}_{p(u)} \left[ \frac{N}{B} \sum_{n \in B_u} \mathcal{L}(\mathcal{W}, \mathbf{x}_n, y_n) \right]\end{aligned}$$

$$\text{Distribution: } p(u = j) = \frac{B}{N}$$

$$\text{stochastic gradient: } \sum_{n \in B_u} \nabla \mathcal{L}(\mathcal{W}, \mathbf{x}_n, y_n)$$

For each update we only need to access a small mini-batch. So it largely reduces the cost

- Neural networks and Back-propagation
- Stochastic optimization
- **Bayesian neural networks**
- Bayes by Backprop and reparameterization trick
- Auto-encoding variational Bayes
- Generative adversarial networks

- Bayesian version of NNs
- We place prior over the weights
- We use different distributions to sample the observed output

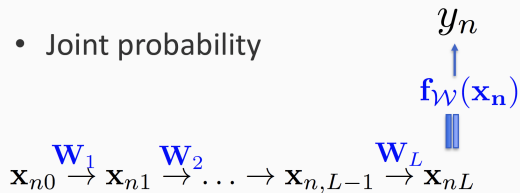
$$\mathbf{f}_{\mathcal{W}}(\mathbf{x}_0)$$



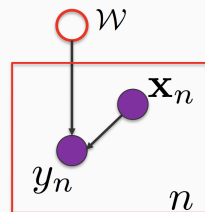
$$\mathbf{x}_0 \xrightarrow{\mathbf{W}_1} \mathbf{x}_1 \xrightarrow{\mathbf{W}_2} \dots \mathbf{x}_{L-1} \xrightarrow{\mathbf{W}_L} \mathbf{x}_L$$

$$\mathbf{f}_{\mathcal{W}}(\mathbf{x}_0) = \mathbf{f}_{\text{out}} = \mathbf{W}_L \sigma(\mathbf{W}_{L-1} \sigma(\dots \sigma(\mathbf{W}_1 \mathbf{x}_0)))$$

- Joint probability



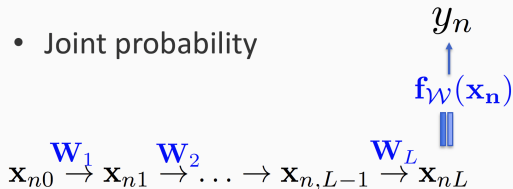
$$\mathcal{W} = \{\mathbf{W}_1, \dots, \mathbf{W}_L\}$$



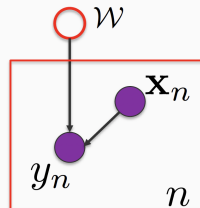
$$p(\mathcal{W}, \mathcal{D}) = p(\mathcal{W}) \prod_{n=1}^N p(y_n | f_{\mathcal{W}}(\mathbf{x}_n))$$



- Joint probability



$$\mathcal{W} = \{\mathbf{W}_1, \dots, \mathbf{W}_L\}$$



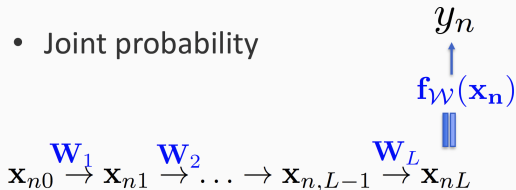
$$p(\mathcal{W}, \mathcal{D}) = p(\mathcal{W}) \prod_{n=1}^N p(y_n | f_{\mathcal{W}}(\mathbf{x}_n))$$

Example of weight priors

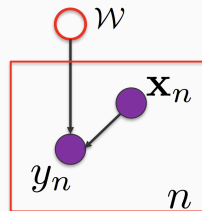
Individual Gaussian  $p(\mathcal{W}) = \prod_{w \in \mathcal{W}} \mathcal{N}(w | 0, 1)$

Spike and slab:  $p(\mathcal{W}) = \prod_{w \in \mathcal{W}} \pi \mathcal{N}(w | 0, \sigma_1^2) + (1 - \pi) \mathcal{N}(w | 0, \sigma_2^2)$  **Encourage sparsity**  
 e.g.,  $\pi = 0.5, \sigma_1^2 = 1, \sigma_2^2 = 1e-3$

- Joint probability



$$\mathcal{W} = \{\mathbf{W}_1, \dots, \mathbf{W}_L\}$$



$$p(\mathcal{W}, \mathcal{D}) = p(\mathcal{W}) \prod_{n=1}^N p(y_n | f_{\mathcal{W}}(\mathbf{x}_n))$$

Example of likelihood

Gaussian:  $p(y_n | f_{\mathcal{W}}(\mathbf{x}_n)) = \mathcal{N}(y_n | f_{\mathcal{W}}(\mathbf{x}_n), \sigma^2)$

Bernoulli:  $p(y_n | f_{\mathcal{W}}(\mathbf{x}_n)) = \text{Bern}(y_n | 1 / (1 + \exp(-f_{\mathcal{W}}(\mathbf{x}_n))))$

Categorical:  $p(y_n | \mathbf{f}_{\mathcal{W}}(\mathbf{x}_n)) = \prod_k \left( \frac{\exp([\mathbf{f}_{\mathcal{W}}(\mathbf{x}_n)]_k)}{\sum_j \exp([\mathbf{f}_{\mathcal{W}}(\mathbf{x}_n)]_j)} \right)^{\mathbb{1}(y_{nk}=1)}$  softmax

- Estimate the posterior distribution of NN weights

$$p(\mathcal{W}|\mathcal{D})$$

- Estimate the predictive distribution

$$p(y^*|\mathbf{x}^*, \mathcal{D}) = \int p(y^*|f_{\mathcal{W}}(\mathbf{x}^*))p(\mathcal{W}|\mathcal{D})d\mathcal{W}$$

- Neural networks and Back-propagation
- Stochastic optimization
- Bayesian neural networks
- Bayes by Backprop and reparameterization trick
- Auto-encoding variational Bayes
- Generative adversarial networks


- The golden-standard for BNN inference is HMC. However, it is often too slow to be practical.
- We want to use variational inference, how?

- We want to use variational inference, how?

Introduce variational posterior and construct variational evidence lower bound!

We choose fully factorized Gaussian

Estimate a free parameter

$$q(\mathcal{W}) = \prod_i q(w_i) = \prod_i \mathcal{N}(w_i | \mu_i, \log(1 + \exp(\rho_i)))$$


$$\begin{aligned} \log(p(\mathcal{D})) &\geq \mathcal{L}(\boldsymbol{\theta}) = \int q(\mathcal{W}) \log \frac{p(\mathcal{W})p(\mathcal{D}|\mathcal{W})}{q(\mathcal{W})} d\mathcal{W} & \boldsymbol{\theta} = \{(\mu_i, \rho_i)\} \\ &= \sum_i \mathbb{E}_{q(w_i)} [\log p(w_i)] + \sum_{n=1}^N \mathbb{E}_{q(\mathcal{W})} [\log p(y_n | f_{\mathcal{W}}(\mathbf{x}_n))] + \sum_i H(q(w_i)) \end{aligned}$$

$$q(\mathcal{W}) = \prod_i q(w_i) = \prod_i \mathcal{N}(w_i | \mu_i, \log(1 + \exp(\rho_i)))$$

$$\log(p(\mathcal{D})) \geq \mathcal{L}(\theta) = \int q(\mathcal{W}) \log \frac{p(\mathcal{W})p(\mathcal{D}|\mathcal{W})}{q(\mathcal{W})} d\mathcal{W}$$

$$= \sum_i \mathbb{E}_{q(w_i)} [\log p(w_i)] + \underbrace{\sum_{n=1}^N \mathbb{E}_{q(\mathcal{W})} [\log p(y_n | f_{\mathcal{W}}(\mathbf{x}_n))]}_{\text{Totally intractable, Why?}} + \sum_i H(q(w_i))$$

Analytical for  
Gaussian prior

Gaussian  
entropy

$$\log(\log(1 + \exp(\rho_i))2\pi e)$$

How to maximize  $\mathcal{L}(\theta)$  ?

- Stochastic optimization
- The key question: How to compute the stochastic gradient for each

$$\mathbb{E}_{q(\mathcal{W})}[\log p(y_n | f_{\mathcal{W}}(\mathbf{x}_n))]$$

Can we use current parameters to sample  $\mathcal{W}$ ,  
plugging into log and calculate the gradient?

$$\widehat{\mathcal{W}} \sim q(\mathcal{W} | \boldsymbol{\theta}) \quad \boldsymbol{\theta} = \{(\mu_i, \rho_i)\}$$



$$\nabla \log p(y_n | f_{\widehat{\mathcal{W}}}(\mathbf{x}_n))$$

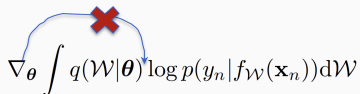


Totally wrong!



- The reason is the distribution contains unknown parameters, and so the expectation and derivative are not interchangeable!

$$\nabla_{\theta} \mathbb{E}_{q(\mathcal{W}|\theta)} [\log p(y_n | f_{\mathcal{W}}(\mathbf{x}_n))] \neq \mathbb{E}_{q(\mathcal{W}|\theta)} [\nabla_{\theta} \log p(y_n | f_{\mathcal{W}}(\mathbf{x}_n))]$$



$$\nabla_{\theta} \int q(\mathcal{W}|\theta) \log p(y_n | f_{\mathcal{W}}(\mathbf{x}_n)) d\mathcal{W}$$

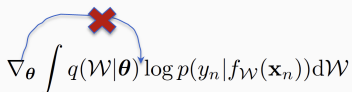


0

Why?

- The reason is the distribution contains unknown parameters, and so the expectation and derivative are not interchangeable!

$$\nabla_{\theta} \mathbb{E}_{q(\mathcal{W}|\theta)} [\log p(y_n | f_{\mathcal{W}}(\mathbf{x}_n))] \neq \mathbb{E}_{q(\mathcal{W}|\theta)} [\nabla_{\theta} \log p(y_n | f_{\mathcal{W}}(\mathbf{x}_n))]$$



$$\nabla_{\theta} \int q(\mathcal{W}|\theta) \log p(y_n | f_{\mathcal{W}}(\mathbf{x}_n)) d\mathcal{W}$$



0

Why?

Because the log likelihood itself does not include variational parameters!

- The solution is to **get rid of the unknown parameters in the distribution** under which we compute the expectation. How?

$$q(\mathcal{W}) = \prod_i q(w_i) = \prod_i \mathcal{N}(w_i | \mu_i, \log(1 + \exp(\rho_i)))$$

$$w_i = \mu_i + \epsilon_i \sqrt{\log(1 + \exp(\rho_i))} \quad \epsilon_i \sim \mathcal{N}(0, 1)$$



$$\text{vec}(\mathcal{W}) = \boldsymbol{\mu} + \text{diag}(\sqrt{\log(1 + \exp(\boldsymbol{\rho}))}) \cdot \boldsymbol{\epsilon} \quad \longrightarrow \quad \mathcal{W} = T(\boldsymbol{\theta}, \boldsymbol{\epsilon}), \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Reparameterized Gaussian sample

$$\mathbb{E}_{q(\mathcal{W}|\boldsymbol{\theta})}[\log p(y_n|f_{\mathcal{W}}(\mathbf{x}_n))] = \mathbb{E}_{p(\boldsymbol{\epsilon})}[\log p(y_n|f_{T(\boldsymbol{\theta},\boldsymbol{\epsilon})}(\mathbf{x}_n))]$$



$$\int q(\mathcal{W}|\boldsymbol{\theta}) \log p(y_n|f_{\mathcal{W}}(\mathbf{x}_n)) d\mathcal{W} = \int p(\boldsymbol{\epsilon}) \log p(y_n|f_{T(\boldsymbol{\theta},\boldsymbol{\epsilon})}(\mathbf{x}_n)) d\boldsymbol{\epsilon}$$



$$\nabla_{\boldsymbol{\theta}} \int q(\mathcal{W}|\boldsymbol{\theta}) \log p(y_n|f_{\mathcal{W}}(\mathbf{x}_n)) d\mathcal{W} = \nabla_{\boldsymbol{\theta}} \int p(\boldsymbol{\epsilon}) \log p(y_n|f_{T(\boldsymbol{\theta},\boldsymbol{\epsilon})}(\mathbf{x}_n)) d\boldsymbol{\epsilon}$$



$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{q(\mathcal{W}|\boldsymbol{\theta})}[\log p(y_n|f_{\mathcal{W}}(\mathbf{x}_n))]$$

$$= \int \nabla_{\boldsymbol{\theta}} p(\boldsymbol{\epsilon}) \log p(y_n|f_{T(\boldsymbol{\theta},\boldsymbol{\epsilon})}(\mathbf{x}_n)) d\boldsymbol{\epsilon}$$

$$= \int p(\boldsymbol{\epsilon}) \nabla_{\boldsymbol{\theta}} \log p(y_n|f_{T(\boldsymbol{\theta},\boldsymbol{\epsilon})}(\mathbf{x}_n)) d\boldsymbol{\epsilon}$$

$$= \mathbb{E}_{p(\boldsymbol{\epsilon})}[\nabla_{\boldsymbol{\theta}} \log p(y_n|f_{T(\boldsymbol{\theta},\boldsymbol{\epsilon})}(\mathbf{x}_n))]$$

$$\mathbb{E}_{q(\mathcal{W}|\boldsymbol{\theta})}[\log p(y_n|f_{\mathcal{W}}(\mathbf{x}_n))] = \mathbb{E}_{p(\boldsymbol{\epsilon})}[\log p(y_n|f_{T(\boldsymbol{\theta},\boldsymbol{\epsilon})}(\mathbf{x}_n))]$$



$$\int q(\mathcal{W}|\boldsymbol{\theta}) \log p(y_n|f_{\mathcal{W}}(\mathbf{x}_n)) d\mathcal{W} = \int p(\boldsymbol{\epsilon}) \log p(y_n|f_{T(\boldsymbol{\theta},\boldsymbol{\epsilon})}(\mathbf{x}_n)) d\boldsymbol{\epsilon}$$



$$\nabla_{\boldsymbol{\theta}} \int q(\mathcal{W}|\boldsymbol{\theta}) \log p(y_n|f_{\mathcal{W}}(\mathbf{x}_n)) d\mathcal{W} = \nabla_{\boldsymbol{\theta}} \int p(\boldsymbol{\epsilon}) \log p(y_n|f_{T(\boldsymbol{\theta},\boldsymbol{\epsilon})}(\mathbf{x}_n)) d\boldsymbol{\epsilon}$$



$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{q(\mathcal{W}|\boldsymbol{\theta})}[\log p(y_n|f_{\mathcal{W}}(\mathbf{x}_n))]$$

$$= \int \nabla_{\boldsymbol{\theta}} p(\boldsymbol{\epsilon}) \log p(y_n|f_{T(\boldsymbol{\theta},\boldsymbol{\epsilon})}(\mathbf{x}_n)) d\boldsymbol{\epsilon}$$

$$= \int p(\boldsymbol{\epsilon}) \nabla_{\boldsymbol{\theta}} \log p(y_n|f_{T(\boldsymbol{\theta},\boldsymbol{\epsilon})}(\mathbf{x}_n)) d\boldsymbol{\epsilon}$$

$$= \mathbb{E}_{p(\boldsymbol{\epsilon})} [\nabla_{\boldsymbol{\theta}} \log p(y_n|f_{T(\boldsymbol{\theta},\boldsymbol{\epsilon})}(\mathbf{x}_n))]$$

Stochastic gradient ascent!

$$\begin{aligned}
 \mathcal{L}(\theta) &= \sum_i \mathbb{E}_{q(w_i)} [\log p(w_i)] + \sum_i H(q(w_i)) \\
 &+ \underbrace{\sum_{u=1}^{N/B} \frac{B}{N} \sum_{n \in \mathcal{B}_u} \frac{N}{B} \mathbb{E}_{p(\epsilon)} [\log p(y_n | f_{T(\theta, \epsilon)}(\mathbf{x}_n))]}_{\substack{\mathbb{E}_{p(u)} \mathbb{E}_{p(\epsilon)} \sum_{n \in \mathcal{B}_u} \frac{N}{B} [\log p(y_n | f_{T(\theta, \epsilon)}(\mathbf{x}_n))]} \\
 &\quad \swarrow \\
 &\text{Constant distribution}
 \end{aligned}$$

- 1. Initialize  $\theta$  randomly
- 2. For  $t = 1..T$ 
  - Sample  $u$  from  $p(u)$ ,  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - Calculate stochastic gradient  $\nabla_{\theta} [\alpha(\theta)] + \frac{N}{B} \sum_{n \in \mathcal{B}_u} \nabla_{\theta} [\log p(y_n | f_{T(\theta, \epsilon)}(\mathbf{x}_n))]$
  - Update  $\theta \leftarrow \theta + \gamma_t \cdot \left( \nabla_{\theta} [\alpha(\theta)] + \frac{N}{B} \sum_{n \in \mathcal{B}_u} \nabla_{\theta} [\log p(y_n | f_{T(\theta, \epsilon)}(\mathbf{x}_n)) \right]$
- 3. Return  $q(\mathcal{W} | \theta) = \prod_i \mathcal{N}(w_i | \mu_i, \log(1 + \exp(\rho_i)))$

- 1. Initialize  $\theta$  randomly
- 2. For  $t = 1..T$ 
  - Sample  $u$  from  $p(u)$ ,  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - Calculate stochastic gradient  $\nabla_{\theta} [\alpha(\theta)] + \frac{N}{B} \sum_{n \in \mathcal{B}_u} \nabla_{\theta} [\log p(y_n | f_{T(\theta, \epsilon)}(\mathbf{x}_n))]$
  - Update  $\theta \leftarrow \theta + \gamma_t \cdot \left( \nabla_{\theta} [\alpha(\theta)] + \frac{N}{B} \sum_{n \in \mathcal{B}_u} \nabla_{\theta} [\log p(y_n | f_{T(\theta, \epsilon)}(\mathbf{x}_n)) \right]$
- 3. Return  $q(\mathcal{W} | \theta) = \prod_i \mathcal{N}(w_i | \mu_i, \log(1 + \exp(\rho_i)))$

output of  
the NN, so it  
needs BP!



$$\begin{aligned} p(y^*|\mathbf{x}^*, \mathcal{D}) &= \int p(y^*|f_{\mathcal{W}}(\mathbf{x}^*))p(\mathcal{W}|\mathcal{D})d\mathcal{W} \\ &\approx \int p(y^*|f_{\mathcal{W}}(\mathbf{x}^*))q(\mathcal{W}|\boldsymbol{\theta})d\mathcal{W} \end{aligned}$$

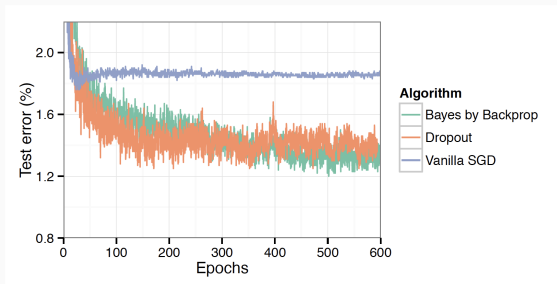
Still intractable, but we can use Monte-Carlo approximation

$$\approx \frac{1}{M} \sum_{j=1}^m p(y^*|f_{\mathcal{W}_j}(\mathbf{x}^*)) \quad \mathcal{W}_j \sim q(\mathcal{W}|\boldsymbol{\theta})$$

We can also generate samples of  $y^*$  to obtain an empirical (or histogram) distribution

Table 1. Classification Error Rates on MNIST. ★ indicates result used an ensemble of 5 networks.

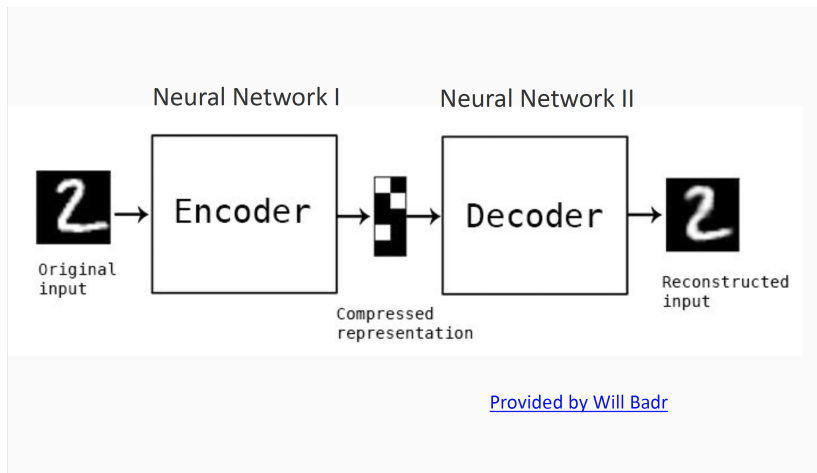
Method	# Units/Layer	# Weights	Test Error
SGD, no regularisation (Simard et al., 2003)	800	1.3m	1.6%
SGD, dropout (Hinton et al., 2012)			≈ 1.3%
SGD, dropconnect (Wan et al., 2013)	800	1.3m	<b>1.2%</b> ★
SGD	400	500k	1.83%
	800	1.3m	1.84%
	1200	2.4m	1.88%
SGD, dropout	400	500k	1.51%
	800	1.3m	1.33%
	1200	2.4m	1.36%
Bayes by Backprop, Gaussian	400	500k	1.82%
	800	1.3m	1.99%
	1200	2.4m	2.04%
Bayes by Backprop, Scale mixture	400	500k	1.36%
	800	1.3m	1.34%
	1200	2.4m	<b>1.32%</b>



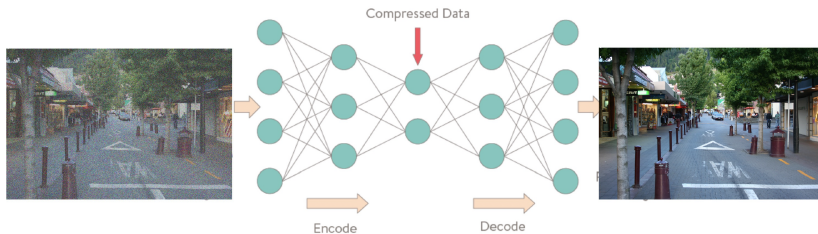
*Figure 2.* Test error on MNIST as training progresses.

- State of the art NN inference, very popular
- The same scalability to SGD, but it can estimate posteriors!
- Core idea : variational inference + reparameterization trick
- This is also the foundation of nearly all the modern Bayesian NN training.

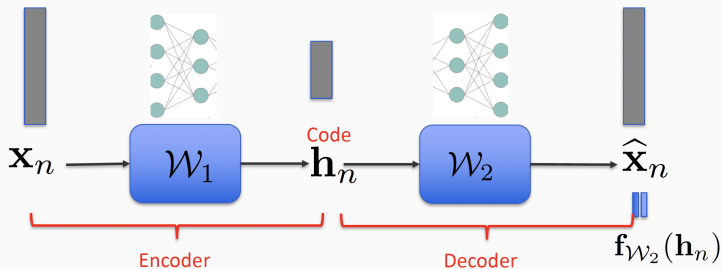
- Neural networks and Back-propagation
- Stochastic optimization
- Bayesian neural networks
- Bayes by Backprop and reparameterization trick
- **Auto-encoding variational Bayes**
- Generative adversarial networks



Dimension reduction is very important:  
compression, denoise, ...



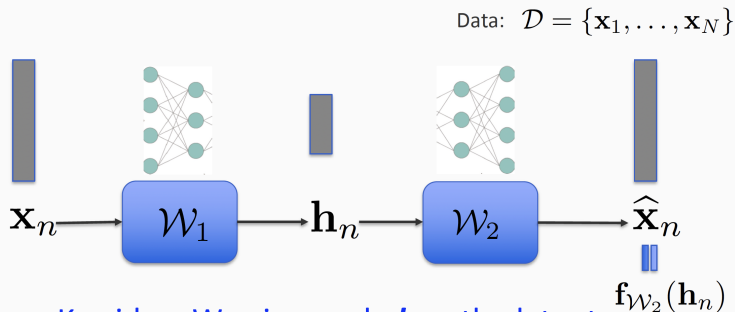
[Provided by Will Badr](#)



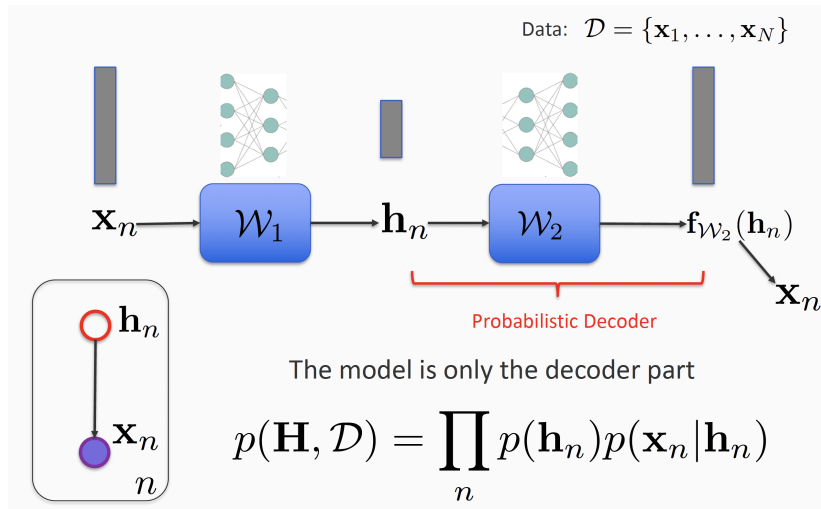
Given data  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$

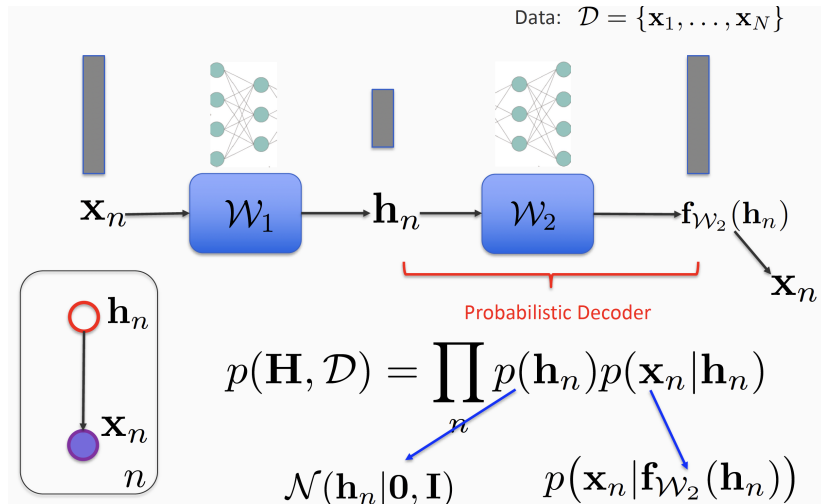
$$\text{Loss: } \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{f}_{\mathcal{W}_2}(\mathbf{h}_{\mathcal{W}_1}(\mathbf{x}_n))\|^2$$

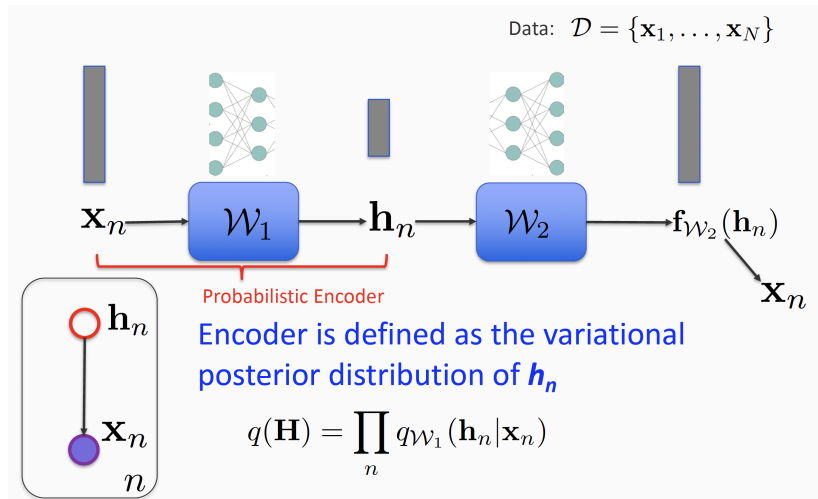


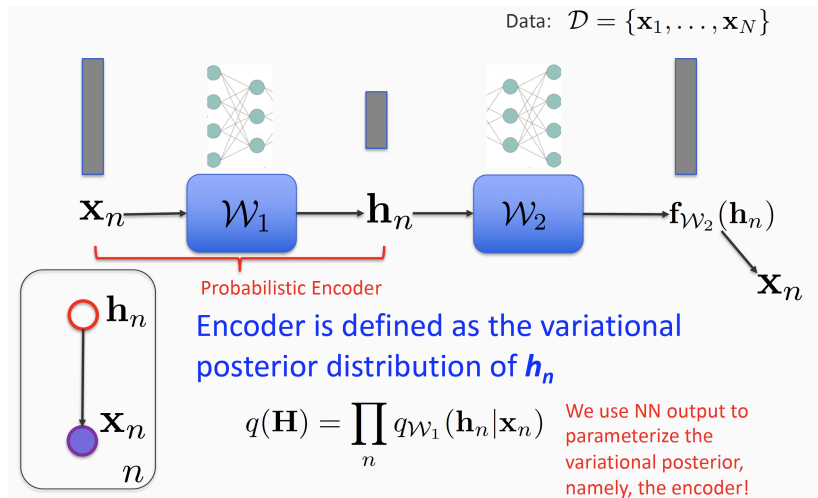


Key idea: We view code  $\mathbf{h}$  as the latent random variables. We want to estimate the posterior distribution of  $\mathbf{h}$ ; However, the NN weights are considered as hyper-parameters rather than RVs.









- Maximize the variational ELBO

$$\begin{aligned}
 \mathcal{L} &= \int q(\mathbf{H}) \log \frac{p(\mathbf{H})p(\mathbf{H}, \mathcal{D})}{q(\mathbf{H})} d\mathbf{H} \\
 &= \sum_{n=1}^N \int q_{\mathcal{W}_1}(\mathbf{h}_n | \mathbf{x}_n) \log \frac{p(\mathbf{h}_n)p(\mathbf{x}_n | \mathbf{f}_{\mathcal{W}_2}(\mathbf{h}_n))}{q_{\mathcal{W}_1}(\mathbf{h}_n | \mathbf{x}_n)} d\mathbf{h}_n \\
 &= \sum_{n=1}^N \mathbb{E}_{q_{\mathcal{W}_1}(\mathbf{h}_n | \mathbf{x}_n)} \left[ \log \frac{p(\mathbf{h}_n)p(\mathbf{x}_n | \mathbf{f}_{\mathcal{W}_2}(\mathbf{h}_n))}{q_{\mathcal{W}_1}(\mathbf{h}_n | \mathbf{x}_n)} \right]
 \end{aligned}$$

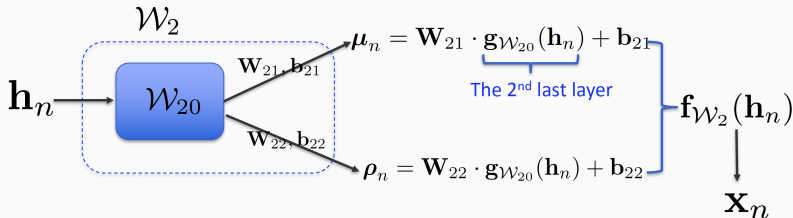
ELBO is obviously intractable, why?

Use reparameterization trick + stochastic optimization (on mini-batches)!

## Concrete example



- Likelihood for continuous output

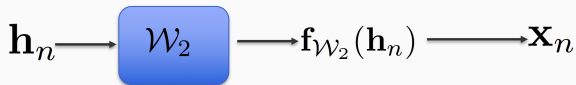
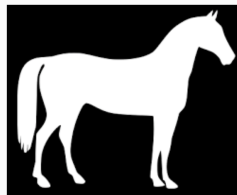


$$p(\mathbf{x}_n | \mathbf{h}_n) = p(\mathbf{x}_n | \mathbf{f}_{\mathcal{W}_2}(\mathbf{h}_n)) = \mathcal{N}(\mathbf{x}_n | \mu_n, \text{diag}(\exp(\rho_n)))$$

Gaussian with diagonal covariance

## Concrete example

- Likelihood for binary output



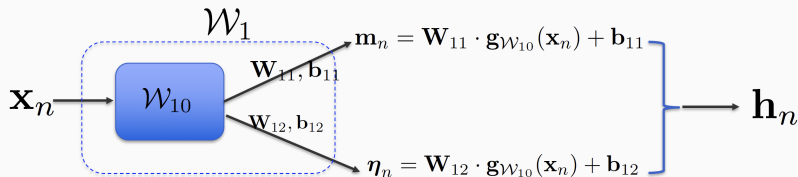
$$p(\mathbf{x}_n | \mathbf{h}_n) = p(\mathbf{x}_n | \mathbf{f}_{\mathcal{W}_2}(\mathbf{h}_n)) = \prod_j \text{Bern}([\mathbf{x}_n]_j | \alpha([\mathbf{f}_{\mathcal{W}_2}(\mathbf{h}_n)]_j))$$

Bernoulli likelihood over each element

$$\alpha(t) = 1/(1 + \exp(-t))$$



- Gaussian encoder (most commonly used)



$$q_{\mathcal{W}_1}(\mathbf{h}_n | \mathbf{x}_n) = \mathcal{N}(\mathbf{h}_n | \mathbf{m}_n, \text{diag}(\exp(\boldsymbol{\eta}_n)))$$

$$\mathcal{L} = \sum_{n=1}^N \mathbb{E}_{q_{\mathcal{W}_1}(\mathbf{h}_n | \mathbf{x}_n)} \left[ \log \frac{p(\mathbf{h}_n) p(\mathbf{x}_n | \mathbf{f}_{\mathcal{W}_2}(\mathbf{h}_n))}{q_{\mathcal{W}_1}(\mathbf{h}_n | \mathbf{x}_n)} \right]$$

Very easy to use  
reparameterization trick!

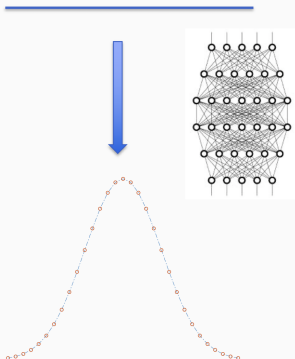
- Convert auto-encoder estimation into a probabilistic inference problem
- Trivial application of VI
- State-of-the-art
- Very popular

- Neural networks and Back-propagation
- Stochastic optimization
- Bayesian neural networks
- Bayes by Backprop and reparameterization trick
- Auto-encoding variational Bayes
- Generative adversarial networks

- Consider a uniform random variable  $\mathbf{X}$ , How can we make a transformation/mapping  $T$  such that the transformed variable follows an arbitrary distribution?
- This is classical statistical question
- Suppose the target distribution has CDF to be  $F$
- Then we should do  $T(\mathbf{X}) = F^{-1}(\mathbf{X})$

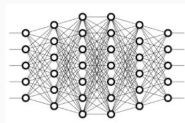
- Now let us consider an even harder problem
- Suppose I **do NOT know the CDF of the target distribution** (this is often true in practice)
- I only **have a set of samples from the target distribution** (e.g., a set of images)
- Can I **learn** such a **mapping  $T$** , such that  $T(\mathbf{X})$  follows the target distribution reflected by the given samples? (In general,  $\mathbf{X}$  can come from any convenient distribution)
- That is what GAN aims for

- We will use an NN to represent the mapping. The learning is to identify the parameters of the NN



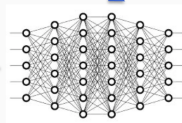
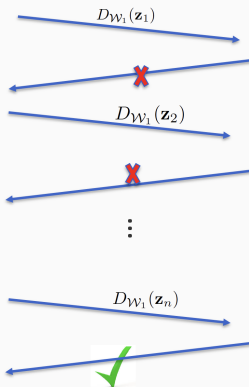
- Key idea: Adversarial Training
- How: we will introduce two NNs, one is a generative network (faker), the other is a discriminative network. (police). We want to train an excellent faker through grilling it by a stronger and stronger police.

- Key idea: Adversarial Training (Gaming)



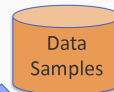
Generator (faker)  
 $G_{W_1}(\cdot)$

I want to fake  
the sample as  
good as possible



Discriminator (police)  
 $D_{W_2}(\cdot)$

I want to detect  
the faked sample  
as well as possible



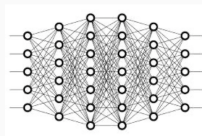


- Adversarial Training (Gaming)

Training examples



Generator (faker)



$$\mathbf{Z} \xrightarrow{G_{W_1}(\cdot)} \mathbf{X}$$

Can be generated from any easy distribution, uniform, Gaussian white noise, ...

The transformed sample, expected to follow the same distribution with the training examples

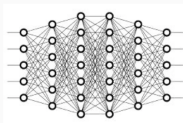
Note that they do not need to have the same dimension!

- Adversarial Training (Gaming)

Training examples



Discriminator (police)



$\mathbf{X}$   
A candidate

$$D_{W_2}(\cdot)$$

Probability of being true

The probability that the candidate can be considered as a sample from the distribution that produces the training examples

- Adversarial Training (Gaming)

Training objective: min—max problem

Training examples



$$\min_{\mathcal{W}_1} \max_{\mathcal{W}_2} \mathcal{L}(\mathcal{W}_1, \mathcal{W}_2) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_{\mathcal{W}_2}(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \in p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D_{\mathcal{W}_2}(G_{\mathcal{W}_1}(\mathbf{z})))]$$

Empirical distribution constructed  
from the training examples

So, we are searching for saddle points as solution, rather than (local) maxima and minima.

## Mini-Max Stochastic Optimization

- Randomly Initialize  $\mathcal{W}_1, \mathcal{W}_2$  and other hyper-parameters
- For  $t=1..T$ 
  - For  $k$  steps do

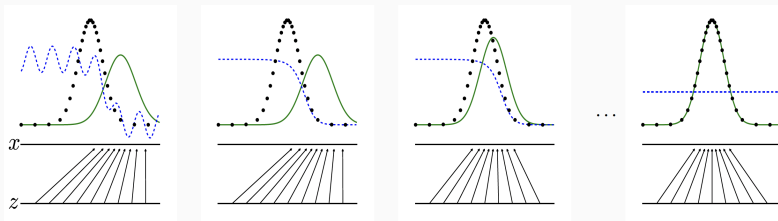
- Sample a minibatch of  $m$  samples  $\mathbf{z}_1, \dots, \mathbf{z}_m \sim p_{\mathbf{z}}(\mathbf{z})$
- Sample a minibatch of  $m$  samples  $\mathbf{x}_1, \dots, \mathbf{x}_m \sim p_{\text{data}}$
- Update Discriminator with stochastic gradient ascent

$$\mathcal{W}_2 \leftarrow \mathcal{W}_2 + \gamma_{tk} \cdot \nabla_{\mathcal{W}_2} \frac{1}{m} \sum_{i=1}^m [\log D_{\mathcal{W}_2}(\mathbf{x}_i) + \log(1 - D_{\mathcal{W}_2}(G_{\mathcal{W}_1}(\mathbf{z}_i)))]$$

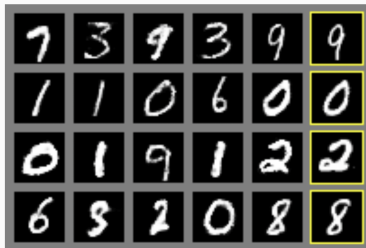
- Sample a minibatch  $m$  samples  $\mathbf{z}_1, \dots, \mathbf{z}_m \sim p_{\mathbf{z}}(\mathbf{z})$
- Update Generator with stochastic gradient descent

$$\mathcal{W}_1 \leftarrow \mathcal{W}_1 - \eta_t \cdot \nabla_{\mathcal{W}_1} \frac{1}{m} \sum_{i=1}^m \log(1 - D_{\mathcal{W}_2}(G_{\mathcal{W}_1}(\mathbf{z}_i)))$$

- Return  $\mathcal{W}_1, \mathcal{W}_2$



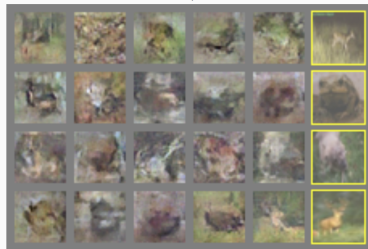
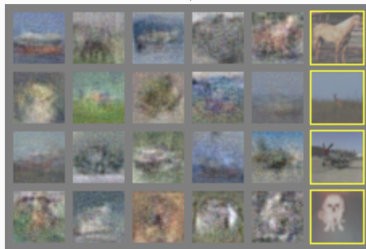
Ian Goodfellow, et. al. 2014



a)



b)





Many funny examples online....



- Deepfake
- Style transfer
- Composition
- ...

- What are Bayesian NNs?
- What are the key idea of BP and stochastic optimization?
- How to conduct variational inference for BNNs?
- What is the reparameterization trick?
- The key idea of Bayes by Backprop, variational auto-encoder and GANs
- You should be able to implement them (with TensorFlow or pyTorch) now!