

Generative Models in Deep Learning

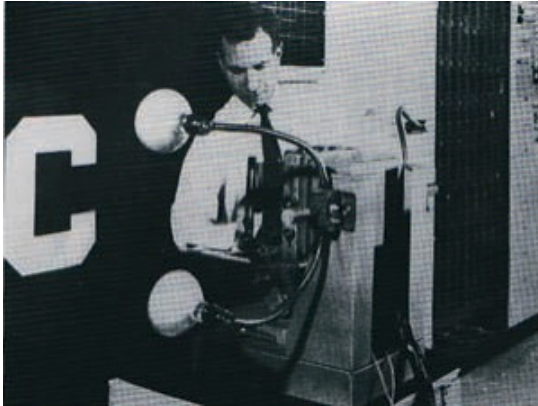
Sargur N. Srihari

srihari@cedar.buffalo.edu

Plan of Presentation

1. Early neural networks: Optimization Criteria
2. Need for Probabilities
3. Representations
 1. Generative and Discriminative Models
 2. Directed/Undirected Representations
 - Complexity of Learning
 - Sampling and Variational Inference
4. Deep Generative Models: RBMs, DBNs, DBMs
5. Directed Generative Nets: VAEs, GANs

Perceptron Hardware (Analog)



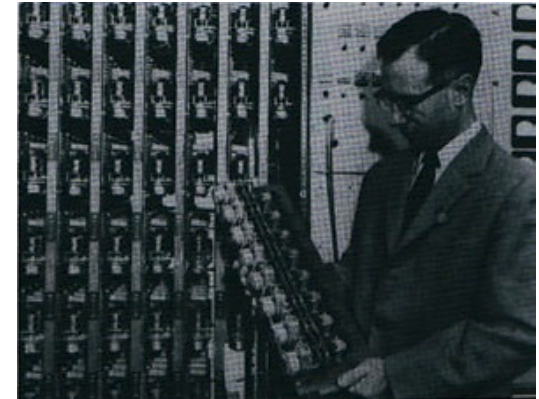
Learning to discriminate
shapes of characters

20x20 cell

Image of character



Patch-board
to allow
different
configurations of
input features ϕ



Racks of
Adaptive Weights
Implemented
as potentiometers

Built in Buffalo, NY. Known as Mark 1 Perceptron. It is now in the Smithsonian.

Perceptron Criterion: Classifier Accuracy

- Seek \mathbf{w} such that $\mathbf{x}_n \in C_1$ will have $\mathbf{w}^T \mathbf{x}_n \geq 0$
whereas $\mathbf{x}_n \in C_2$ will have $\mathbf{x} = \mathbf{w}^T \mathbf{x}_n < 0$
- With targets $t_n \in \{+1, -1\}$, we need $\mathbf{w}^T \mathbf{x}_n t_n > 0$
 - Minimize $E_P(\mathbf{w}) = -\sum_{n \in M} \mathbf{w}^T \mathbf{x}_n t_n$
 M : set of all misclassified patterns
 - Implicit 0-1 loss function
- Update rule: $\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E_P(\mathbf{w}^{\tau}) = \mathbf{w}^{\tau} + \eta \mathbf{x}_n t_n$
 η is learning rate, τ indexes the steps
Cycle through \mathbf{x}_n in turn
If incorrectly classified add to weight vector

Surrogate for Accuracy

- Disadvantage of (0-1) loss
 - Piecewise constant function of w with discontinuities
 - No closed form solution, Intractable even for a linear classifier, Gradient descent requires gradient direction
- Surrogate loss function (Probabilistic)
 - *Negative log-likelihood* of correct output/class

$$E(w) = -\sum_{i=1}^N \log p(y^{(i)} | x^{(i)}; w)$$

$$\text{Linear Regression:}$$
$$p(y | x) \sim N(y; \hat{y}(x, w), \sigma^2)$$

$$E(w) = -m \log \sigma - \frac{m}{2} \log 2\pi - \sum_{i=1}^N \frac{\|\hat{y}^{(i)} - y^{(i)}\|^2}{2\sigma^2}$$

Minimizes Squared error

$$\text{Logistic Regression}$$
$$y = \sigma(w^T x)$$

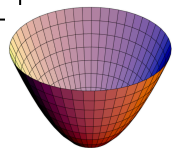
$$E(w) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}, \quad y_n = \sigma(w^T x_n)$$

$$\nabla E(w) = \sum_{n=1}^N (y_n - t_n) x_n$$

Cross entropy is concave, can use gradient descent

– Using learnt w we can predict $p(y | x; w)$

– Can pick class that yields least classification error in expectation



Surrogate may learn more

- Using log-likelihood surrogate,
 - Test set 0-1 loss continues to decrease for a long time after the training set 0-1 loss has reached zero when training
 - Because one can improve classifier robustness by further pushing the classes apart
 - Results in a more confident and robust classifier
 - Thus extracting more information from the training data than with minimizing 0-1 loss
- This sets stage for why we need to determine the conditional distribution $p(y|\mathbf{x})$

Joint and Conditional Distributions

- Input x , output y
- Joint distribution $p(x, y)$
 - Given 4 data points: (x, y) : $(1, 0)$, $(1, 0)$, $(2, 0)$, $(2, 1)$

- Then $p(x, y)$ is

	<u>$y=0$</u>	<u>$y=1$</u>
$x=1$	$\frac{1}{2}$	0
$x=2$	$\frac{1}{4}$	$\frac{1}{4}$

- Conditional distribution $p(y|x)$

- Then $p(y|x)$ is

	<u>$y=0$</u>	<u>$y=1$</u>
$x=1$	1	0
$x=2$	$\frac{1}{2}$	$\frac{1}{2}$

Classification with either model

- The conditional distribution $p(y|\mathbf{x})$ is natural for classifying a given example \mathbf{x} into a class y
 - which is why algorithms that model this directly are called discriminative algorithms

- Generative algorithms model $p(\mathbf{x}, y)$, which can be transformed into $p(y|\mathbf{x})$ by applying Bayes rule and then used for classification

$$p(y | \mathbf{x}) = \frac{p(\mathbf{x}, y)}{p(\mathbf{x})} = \frac{p(\mathbf{x} | y)p(y)}{\sum_y p(\mathbf{x} | y)p(y)}$$

- However, the distribution $p(\mathbf{x}, y)$ can also be used for other purposes.
 - E.g., you could use $p(\mathbf{x}, y)$ to *generate* likely (\mathbf{x}, y) pairs.

Generative Models preserve FoVs

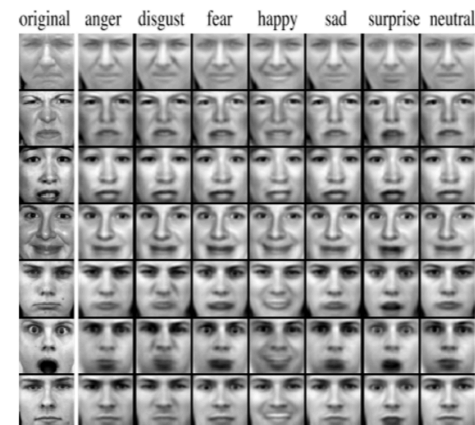
- In designing features we seek features invariance to Factors of Variation

- Face:

- pose, expression, morphology

- Speech:

- age, sex, accent, words spoken



- Discriminative models: *disentangle* FoVs and discard those we do not care about
- Generative models: *preserve* informative FoVs

Which is better: Generative vs Discriminative?

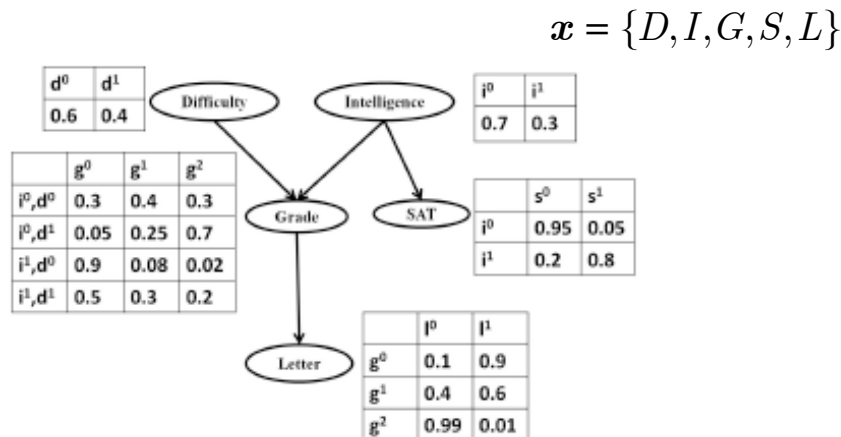
- Generative models $p(\mathbf{x}, y)$ may appear to be more useful, but it's not as simple as that
- Discriminative models $p(y|\mathbf{x})$ generally outperform generative models in classification tasks
 - Fewer adaptive parameters
 - Improved predictive performance when $p(\mathbf{x}|y=C_k)$ assumptions are poor approximations
- How do we represent the distributions?

PGM representations reduce parameters

Bayesian Networks

- CPDs: $p(x_i \mid pa(x_i))$
- Joint Distribution

$$p(\mathbf{x}) = \prod_{i=1}^M p(x_i \mid pa(x_i))$$

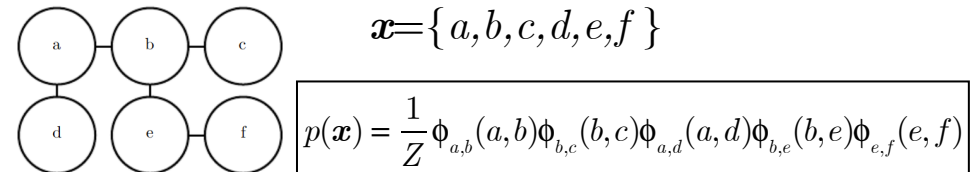


$$p(\mathbf{x}) = p(D)p(I)p(G \mid D, I)p(S \mid I)p(L \mid G)$$

Markov Networks

- Factors: $\phi(C)$
- Joint Distribution

$$p(\mathbf{x}) = \frac{1}{Z} \tilde{p}(\mathbf{x}) \quad Z = \int \tilde{p}(\mathbf{x}) d\mathbf{x} \quad \tilde{p}(\mathbf{x}) = \prod_{C \in \mathcal{G}} \phi(C)$$



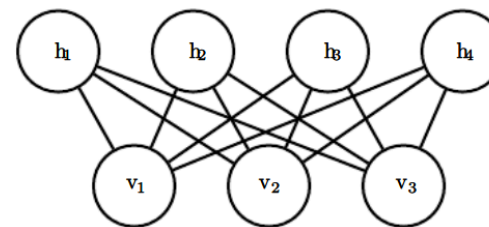
- Energy Model: Boltzmann machine

$$E(C) = -\ln \phi(C), \quad \phi(C) = \exp(-E(C))$$

$$E(\mathbf{x}) = E_{a,b}(a, b) + E_{b,c}(b, c) + E_{a,d}(a, d) + E_{b,e}(b, e) + E_{e,f}(e, f)$$

$$\tilde{p}(\mathbf{x}) = \exp(-E(\mathbf{x}))$$

Restricted Boltzmann machine



$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^T \mathbf{v} - \mathbf{c}^T \mathbf{h} - \mathbf{v}^T \mathbf{W} \mathbf{h}$$

$$P(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} \exp(-E(\mathbf{v}, \mathbf{h}))$$

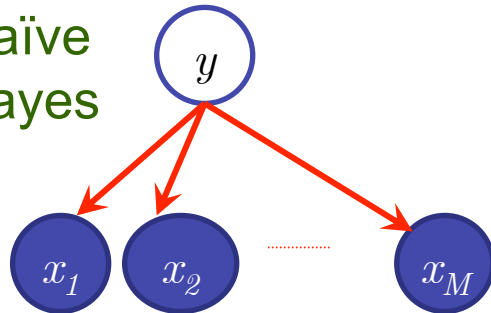
$$Z = \sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp\{-E(\mathbf{v}, \mathbf{h})\}$$

Complexity of Training with BNs, MNs

Variables $\mathbf{x} = \{x_1, \dots, x_M\}$ and classifier target y

1. Generative: estimate parameters of variables independently

Naïve
Bayes



For classification:
Determine joint:

$$p(y, \mathbf{x}) = p(y) \prod_{i=1}^M p(x_i | y)$$

From joint get required
conditional $p(y | \mathbf{x})$

Simple estimation

independently estimate M sets of parameters

But independence is usually false

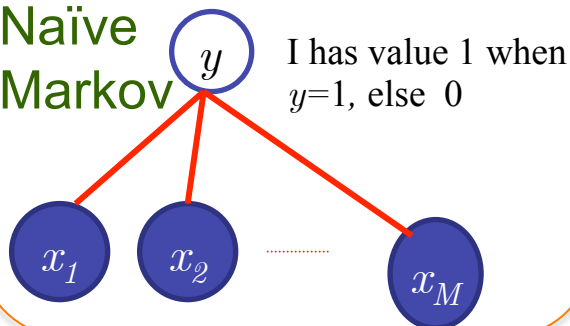
We can estimate $M(M+1)/2$ covariance matrix

2. Discriminative: estimate joint parameters w_i

Potential Functions (log-linear)

$$\phi_i(x_i, y) = \exp\{w_i x_i \mathbb{I}\{y=1\}\},$$
$$\phi_o(y) = \exp\{w_o \mathbb{I}\{y=1\}\}$$

Naïve
Markov



For classification:

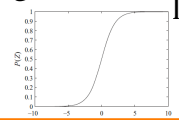
Unnormalized

$$\tilde{P}(y=1 | \mathbf{x}) = \exp\left\{w_0 + \sum_{i=1}^M w_i x_i\right\} \quad \tilde{P}(y=0 | \mathbf{x}) = \exp\{0\} = 1$$

Normalized

$$P(y=1 | \mathbf{x}) = \text{sigmoid}\left\{w_0 + \sum_{i=1}^M w_i x_i\right\} \quad \text{where } \text{sigmoid}(z) = \frac{e^z}{1+e^z}$$

Logistic Regression



Jointly optimize M parameters

More complex estimation but correlations
accounted for

Can use much richer features:

Edges, image patches sharing same pixels

multiclass

$$p(y_i | \phi) = y_i(\phi) = \frac{\exp(a_i)}{\sum_j \exp(a_j)}$$

where $a_j = \mathbf{w}_j^T \phi$

Complexity of Inference with BNs, MNs

- Bayesian Networks

$$p(E = e) = \sum_{x \setminus E} \prod_{i=1}^D p(x_i | pa(x_i))|_{E=e}$$

- An intractable problem

- #P complete

- Markov Networks

$$\begin{aligned} E(\mathbf{v}, \mathbf{h}) &= -\mathbf{b}^T \mathbf{v} - \mathbf{c}^T \mathbf{h} - \mathbf{v}^T W \mathbf{h} \\ P(\mathbf{v}, \mathbf{h}) &= \frac{1}{Z} \exp(-E(\mathbf{v}, \mathbf{h})) \\ Z &= \sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp\{-E(\mathbf{v}, \mathbf{h})\} \end{aligned}$$

- Computing partition function Z by exhaustively summing over all states is intractable

Approximate Inference: Sampling

- We wish to approximate $s = \sum_{\mathbf{x}} p(\mathbf{x})f(\mathbf{x})$ or $s = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x}$
- Rewrite as an expectation $s = E_p[f(\mathbf{x})]$ or $s = E_p[f(\mathbf{x})]$
- Draw n samples $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$ from p and then form average $\hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}^{(i)})$
 - Bayesian Networks: Ancestral sampling
 - Markov networks: Gibbs sampling
- Sampling useful for more than inference
 - Optimization much faster
 - if allowed to rapidly compute approximate estimates rather than slowly compute exact gradient
 - Training set may be redundant

Approximate Inference: Variational

- Given observed \mathbf{v} and latent \mathbf{h}
 - Too expensive to marginalize \mathbf{h}
- To infer log probability of observed $\log p(\mathbf{v} ; \theta)$
 - We compute its evidence lower bound (ELBO)

$$L(\mathbf{v}, \theta, q) = \log p(\mathbf{v}; \theta) - D_{KL}(q(\mathbf{h} | \mathbf{v}) || p(\mathbf{h} | \mathbf{v}))$$

- Thus inference is to find q that minimizes D_{KL}

- Mean field approach

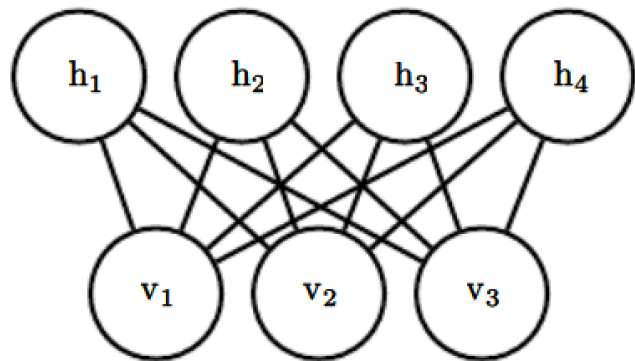
- $q(\mathbf{h} | \mathbf{v}) = \prod_i q(h_i | \mathbf{v})$

- Mean field parameters $\hat{\mathbf{h}}$ are determined by fixed point equations

Deep Generative Models

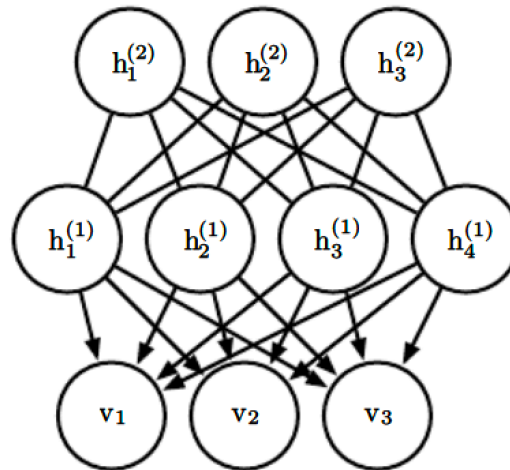
- Represent probability distributions in some way
- Some allow probability distribution functions to be evaluated explicitly
- Some do not allow that evaluation but support sampling from the distribution
- Some can be described as PGMs, others not as factors

Deep Generative Models: RBMs



RBM is an undirected graphical model based on a bipartite graph

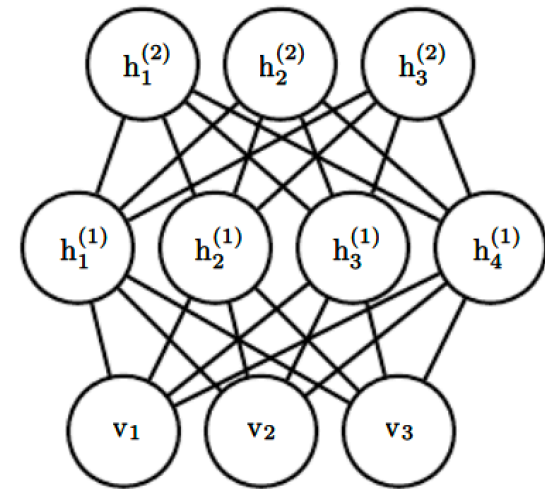
Efficient evaluation and differentiation of $P(\mathbf{v})$
Efficient sampling



Deep Belief Network
Hybrid graphical model with multiple hidden layers

Sampling: Gibbs for top layers, Ancestral for lower

Trained using contrastive divergence



Deep Boltzmann Machine is an undirected graphical model with several layers of latent Variables

Can rewrite as a bipartite Structure and use same Equations as for RBM

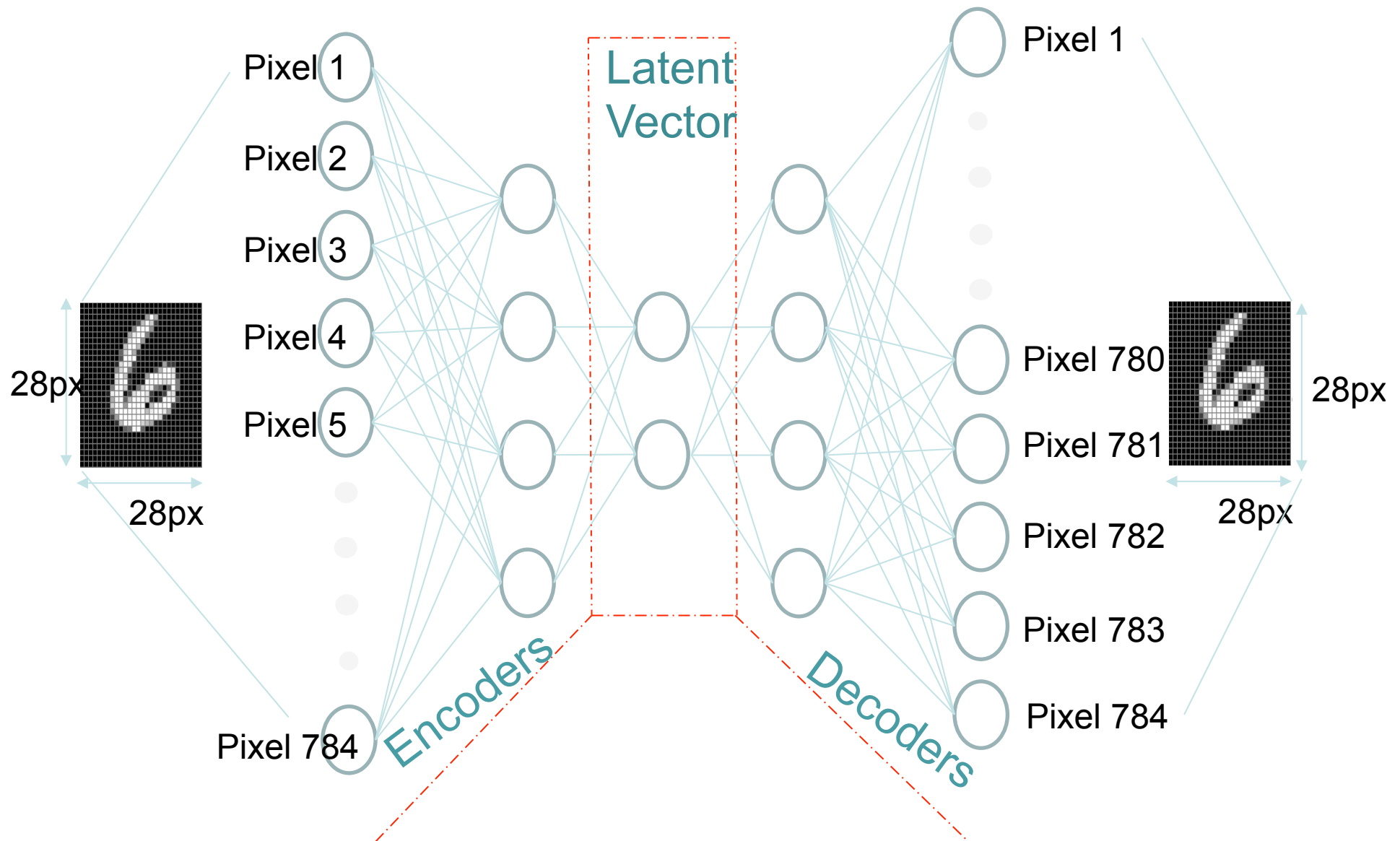
Directed Generative Nets

- *Directed* graphical models have been very popular with the general machine learning community
- But not in deep learning community until 2013
 - Overshadowed by undirected models, e.g., RBMs
- There are now a variety of deep fully directed models
- Many generative models are based on the idea of using a *differentiable generator network*

Differentiable Generator Nets

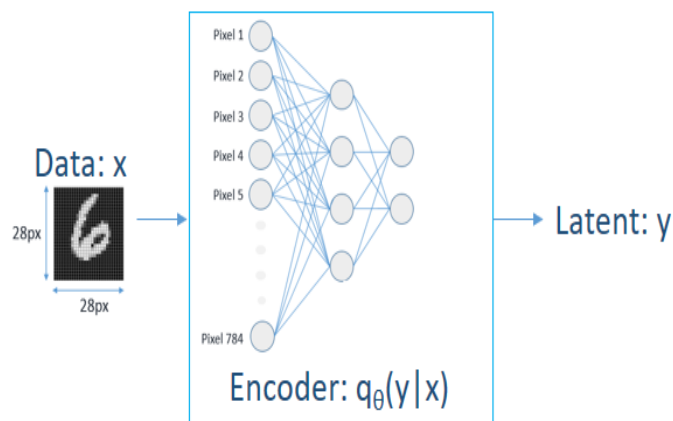
- Model transforms samples of latent variables z to samples x or to distributions over samples x using a differentiable function $g(z; \theta^{(g)})$ represented by a neural network
- Model class includes:
 - Variational autoencoders (VAE)
 - Which pair the generator net with an inference net
 - Generative adversarial networks (GAN)
 - Which pair the generator net with a discriminator net

Autoencoder



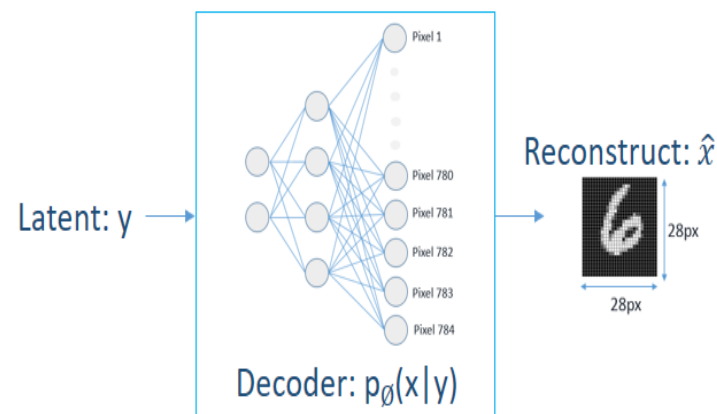
Variational Autoencoder

Encoder (generator net)



- Input x is 28x28=784 image,
- Latent output y is obtained by sampling from Gaussian $q_{\theta}(y|x)$
- Weights and biases of encoder: θ
- Compressed lower dimensional space is stochastic and imitates unit Gaussian $q_{\theta}(y|x)$.

Decoder (inference net)



- Input latent: y
- Outputs parameters to distribution of the data: x
- Weights and biases of decoder: ϕ
- Decoder $p(x|y)$ gets as input the latent vector and outputs real nos. -- probability distribution of 784 pixels [Bernoulli]
- Losses due to reconstruction measured using log-likelihood of $p_{\phi}(x|y)$

The VAE model

- Method for modeling a data distribution using a collection of independent latent variables
 - **Generative model:** $p(x, z) = p(x|z)p(z)$
 - x is a r.v. representing observed data
 - z is a collection of latent variables
 - $p(x|z)$ is parameterized by a deep neural network (decoder)
 - Components of z are independent Bernoulli or Gaussian
 - Learned approx inference trained using gradient descent
 - $q(z|x) = N(\mu, \sigma^2 I)$ whose parameters are given by another deep network (encoder)
 - Thus we have $z \sim \text{Enc}(x) = q(z|x)$ and $y \sim \text{Dec}(z) = p(x|z)$

Generative Adversarial Network (GAN)

- GANs are a generative modeling approach based on differentiable generator networks
 - Based on game-theoretic scenario in which the generator network competes against an adversary

Adversarial Example Generation

Add to x an imperceptibly small vector

Its elements are equal to the sign of the elements of the gradient of the cost function wrt the input. It changes Googlenet's classification of the image



x

$+ .007 \times$



$=$



$$x + \epsilon \operatorname{sign}(\nabla_x J(\theta, x, y))$$

$y = \text{"panda"}$
with 58% confidence

$y = \text{"nematode"}$
With 8.2% confidence

$y = \text{"gibbon"}$
With 99% confidence

Virtual Adversarial Examples

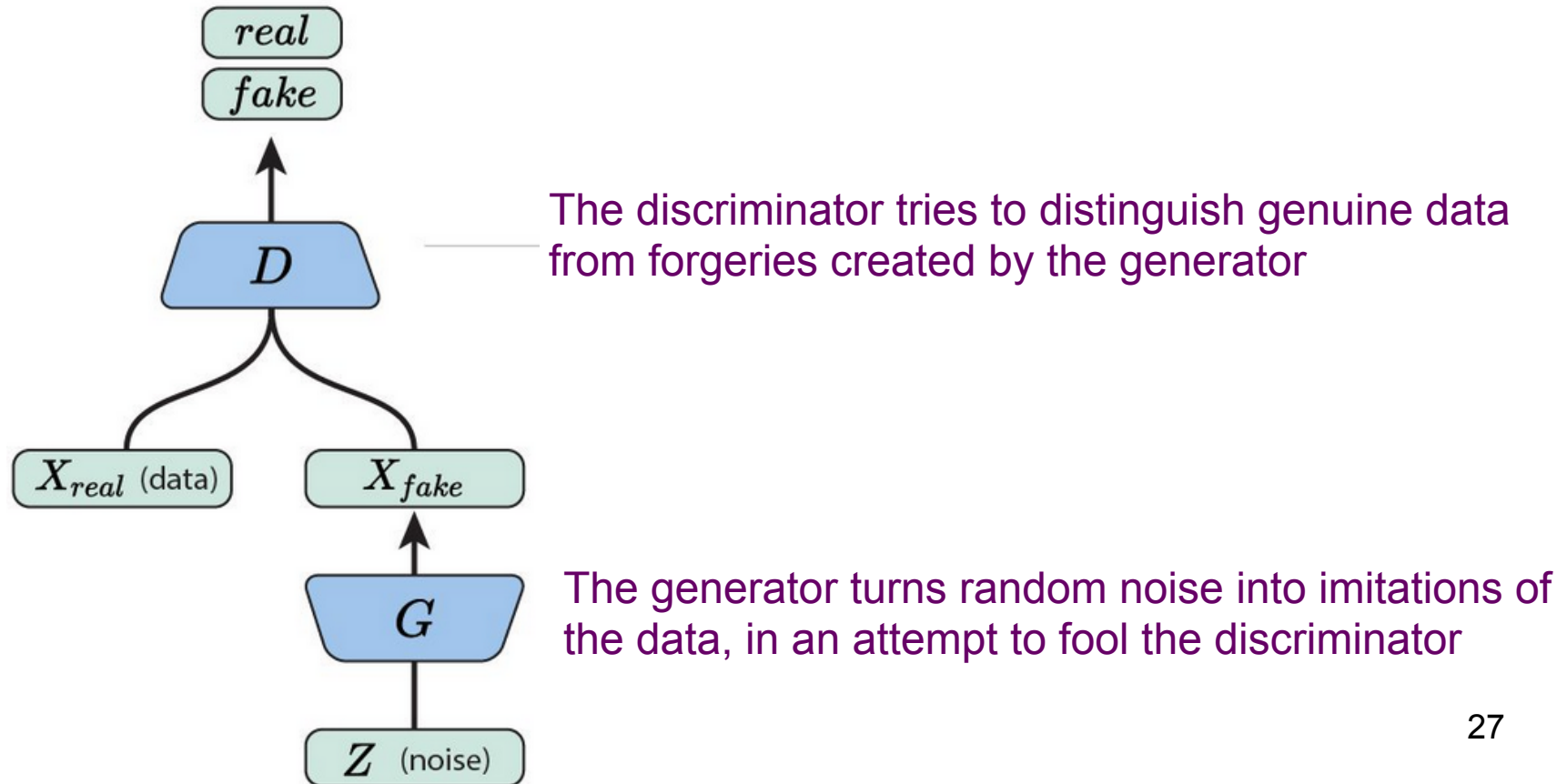
- Adversarial examples generated with using not the true label but a label provided by a trained model are called *Virtual Adversarial Examples*
 - The classifier may then be trained to assign the same label to x and x'
 - This encourages the classifier to learn a function that is robust to small changes anywhere along the manifold where the unlabeled data lie
- Assumption motivating this approach
 - different classes lie on disconnected manifolds
 - A small perturbation should not be able to jump from one class manifold to another class manifold

Generative Adversarial Net (GAN)

- Simultaneously train two models
 - A generative model G that captures the data distribution
 - A discriminative model D that estimates a sample came from the training data rather than G
 - Training procedure for G is to maximize probability that G makes a mistake
- Framework corresponds to a minimax two-player game
 - For arbitrary G and D , a unique solution exists
 - If G and D are MLPs can be trained using backprop

Generative Adversarial Network

- GANs are a way to make a generative model by having two neural networks compete with each other



Generator and Discriminator Networks

1. G produces samples $x = G(z; \theta^{(G)})$
 - Train G to minimize $\log(1 - D(G(z)))$
 2. D , its adversary, attempts to distinguish samples drawn from training data and samples drawn from generator
- D emits a probability value
 - $D(x; \theta^{(D)})$ indicating probability that x is a real training example rather than a fake sample drawn from the model

Deep Convolutional GAN (DCGAN)

- GAN architecture performs well when model architecture and hyperparameters are carefully selected
- A crafted DCGAN performs very well for image synthesis tasks
 - Latent representation space captures important factors of variation

Conditional GAN

- GAN learning problem can be simplified by breaking the generation process into many levels of detail
- Can train conditional GANs that learn to sample from a distribution $p(\mathbf{x}|y)$ rather than simply sampling from the marginal distribution $p(\mathbf{x})$

The LAPGAN Model

- A series of conditional GANs can be trained to first generate a very low resolution version of an image
- Then incrementally add details to the image
- Called LAPGAN due to the use of a Laplacian pyramid to generate varying levels of detail

GAN Image Synthesis

- Architecture/hyperparams are carefully selected
 - Latent variables capture important FoVs

Images of bedrooms generated by Deep Convolutional GAN



Images of churches generated by Laplacian Pyramid of GANs



Disentangling FoVs

- During training, only supervision is class labels
- Specified FoVs
 - Images captured from different viewpoints
 - Strong supervision: pairs of images with two different objects at same viewing angle
- Unspecified FoVs
 - Labels unavailable
- A disentanglement method
 - Combine *variational autoencoder* with *adversarial training*

Summary

- Representing probability distributions is central to deep learning
- PGMs: Undirected, directed and hybrid graphical representations are all useful
- While discriminative models have long been associated with neural nets, generative models are now playing an important role
- Computational considerations make approximate inference methods such as variational and sampling methods play a role