

# Generative Models in Deep Learning

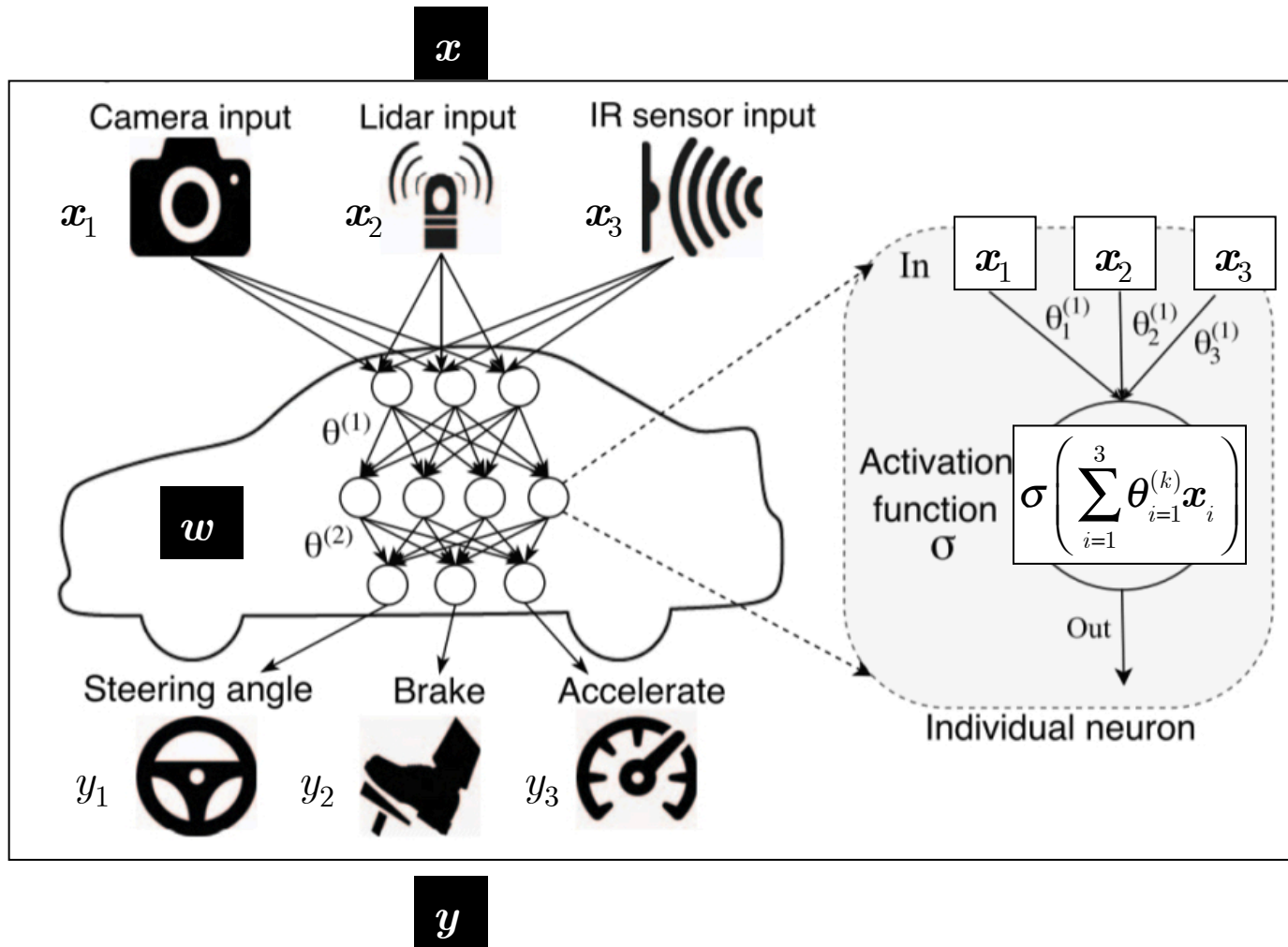
Sargur N. Srihari

[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics

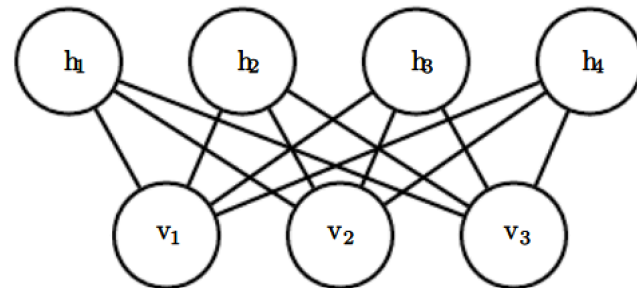
1. Need for Probabilities in Machine Learning
2. Representations
  1. Generative and Discriminative Models
  2. Directed/Undirected Representations
    - Complexity of Learning
    - Sampling and Variational Inference
3. Deep Generative Models: RBMs, DBNs, DBMs
4. Directed Generative Nets: VAEs, GANs
5. Remarks on Deep Learning as a science

# Neural Network for Autonomous Vehicle

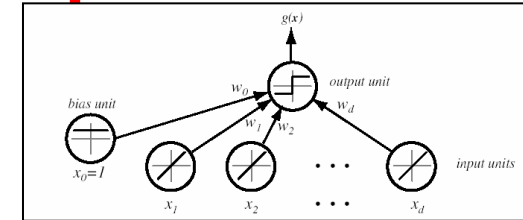


# Why are probabilities used in ML?

1. Conditional probability distributions model the output, e.g.,  $p(\mathbf{y}|\mathbf{x};\mathbf{w})$
2. Likelihood over  $p(\mathbf{y}|\mathbf{x};\mathbf{w})$  is used to define a loss function  $E(\mathbf{w})$  (instead of accuracy) which we optimize to determine weights
3. Probability distributions model the joint distribution  $p(\mathbf{v},\mathbf{h})$  from which we can infer the output  $p(\mathbf{h}|\mathbf{v})$  and  $p(\mathbf{v}|\mathbf{h})$



# Perceptron did not use probabilities

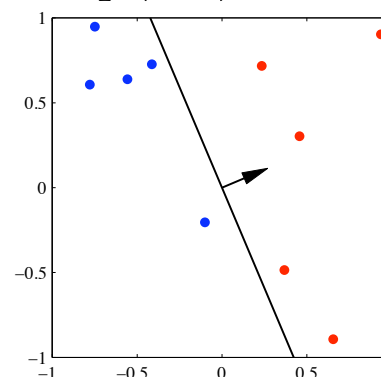
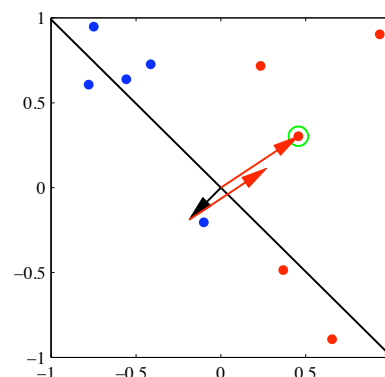
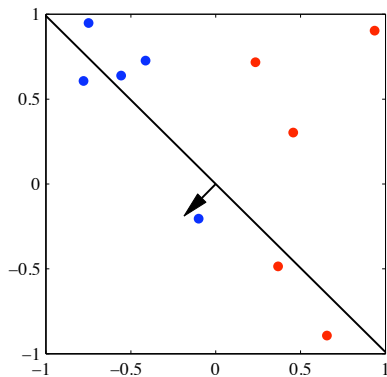


- Output is  $g(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x})$ 
  - Where  $f$  is a threshold function
- Seek  $\mathbf{w}$ : if  $\mathbf{x}_n \in C_1$ ,  $\mathbf{w}^T \mathbf{x}_n \geq 0$ , if  $\mathbf{x}_n \in C_2$ ,  $\mathbf{w}^T \mathbf{x}_n < 0$ 
  - With targets  $t_n \in \{+1, -1\}$ , we need  $\mathbf{w}^T \mathbf{x}_n t_n > 0$ 
    - Minimize no of misclassified samples  $M$

$$E_P(\mathbf{w}) = -\sum_{n \in M} \mathbf{w}^T \mathbf{x}_n t_n$$

- Implicit 0-1 loss function

– Update rule:  $\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E_P(\mathbf{w}^{\tau}) = \mathbf{w}^{\tau} + \eta \mathbf{x}_n t_n$



Cycle through  $\mathbf{x}_n$   
in turn  
If incorrectly  
classified add to  
weight vector

# Disadvantage of Accuracy as Loss

- Counting no. of misclassified samples is:

$$E_P(\mathbf{w}) = -\sum_{n \in M} \mathbf{w}^T \mathbf{x}_n t_n$$

- Which we wish to minimize

- Disadvantages

- Piecewise constant function of  $\mathbf{w}$  with discontinuities

- No closed form solution

- Intractable even for a linear classifier

- Gradient descent requires gradient direction

- Solution: define the output to be probabilistic

- Output is a distribution  $p(y|\mathbf{x};\mathbf{w})$

# Surrogate Loss Function

- When output can take one of two values
  - Assuming Bernoulli distribution,  $\mathbf{t} = \{t_n\} \quad n=1, \dots, N$ ,  $y_n = \sigma(\mathbf{w}^T \mathbf{x}_n)$  likelihood of observed data using model is

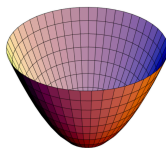
$$p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n} \quad \text{which we wish to maximize}$$

- *Negative log-likelihood* is cross-entropy, which we wish to minimize

$$E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

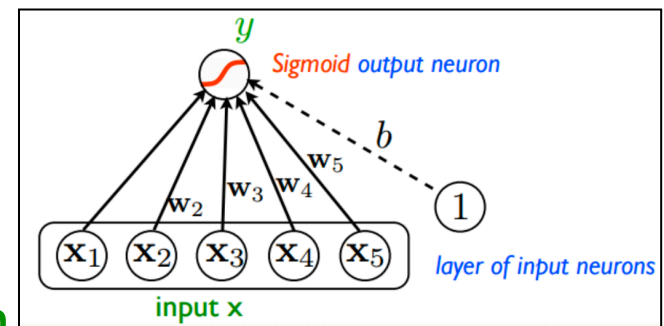
Cross entropy is concave, can use gradient descent

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



- The model is called **logistic regression**

- Using learnt  $\mathbf{w}$  we can predict  $p(y | \mathbf{x}; \mathbf{w})$



# Surrogate loss for continuous case

- Linear regression  $y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x})$ 
  - Assume output is corrupted by Gaussian noise

$$p(t | \mathbf{x}) \sim N(t; y(\mathbf{x}, \mathbf{w}), \sigma^2)$$

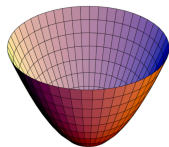
- Negative log-likelihood:

$$E(\mathbf{w}) = -\sum_{n=1}^N \log p(y_n | \mathbf{x}_n; \mathbf{w})$$

$$E(\mathbf{w}) = -N \log \sigma - \frac{N}{2} \log 2\pi - \sum_{n=1}^N \frac{\|y_n - t_n\|^2}{2\sigma^2}$$

Minimize Squared error

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



- Using learnt  $\mathbf{w}$  we can predict  $p(y | \mathbf{x}; \mathbf{w})$ 
  - 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 Conditionals are very different

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

- Then  $p(\mathbf{x}, y)$  is

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

- Conditional distribution  $p(y|\mathbf{x})$

- Then  $p(y|\mathbf{x})$  is

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

- But we can use joint to generate conditional using Bayes rule as seen next

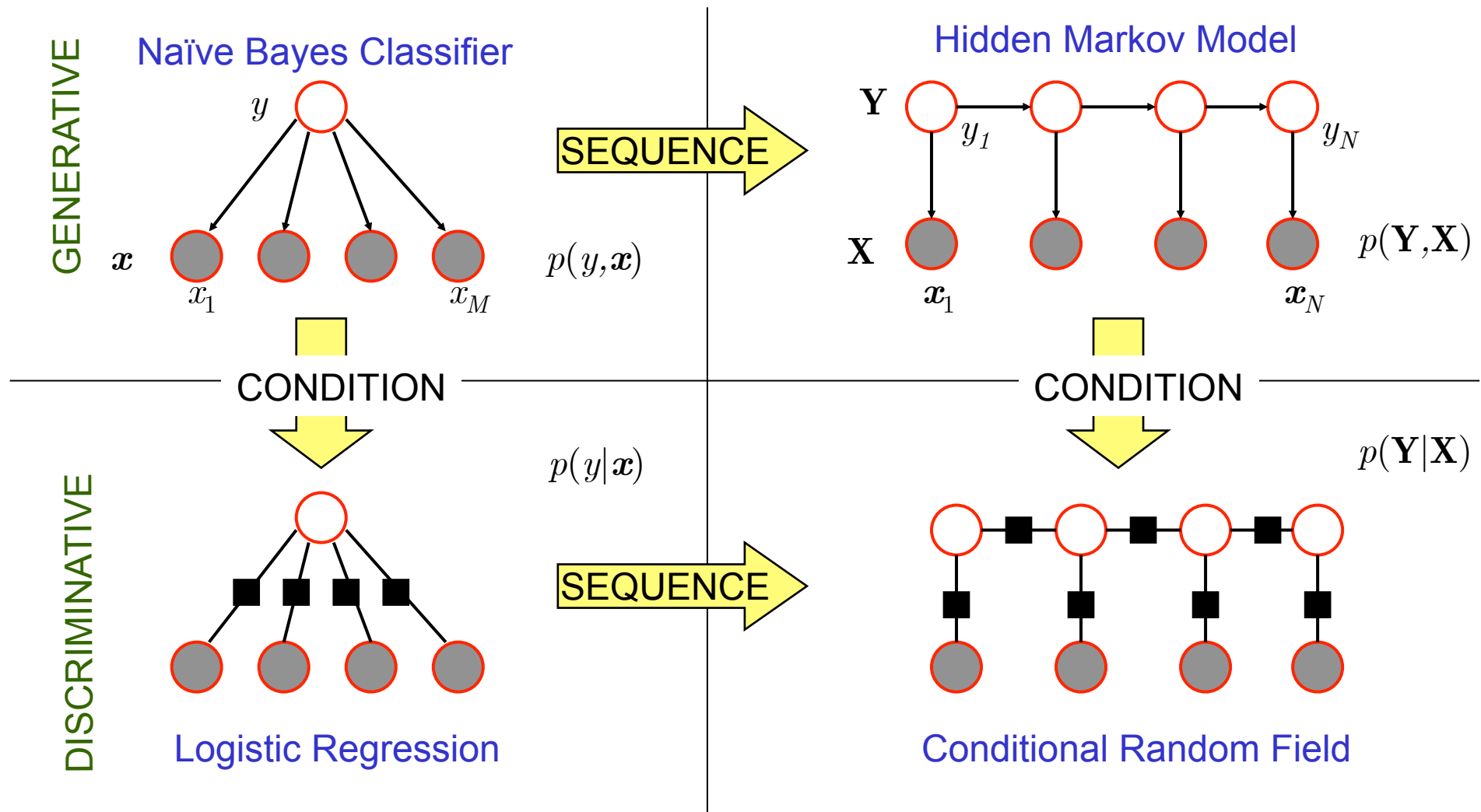
# Classification → determine conditional

- Conditional  $p(y|\mathbf{x})$  is a natural for classification
  1. Discriminative algorithms model  $p(y|\mathbf{x})$  directly
  2. 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

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

- Where  $p(\mathbf{x}, y) = p(\mathbf{x}|y)p(y)$ , i.e., class-conditional  $\mathbf{x}$  prior
- However,  $p(\mathbf{x}, y)$  is useful for other purposes.
  - E.g., you could use  $p(\mathbf{x}, y)$  to generate likely  $(\mathbf{x}, y)$  pairs
- There are examples of both types of classifiers

# Single decision vs Sequential decisions



# Generative or Discriminative?

- Generative models, which model  $p(\mathbf{x}, y)$ , appear to be more general and therefore better
- But Discriminative models, which model  $p(y|\mathbf{x})$ , outperform generative models in classification
  1. Fewer adaptive parameters
    - With  $M$  features logistic regression needs  $M$  parameters
    - Bayes classifier with  $p(\mathbf{x}|y)$  modeled as Gaussian:
      - $2M$  for means +  $M(M+1)/2$  for  $\Sigma$  + Class priors  $\rightarrow M(M+5)/2 + 1$
      - Grows quadratically with  $M$
  2. Performance depends on  $p(\mathbf{x}|y)$  approximations

# Deep Generative Models: Overview

- Generative models are built/trained using:
  - PGMs, Monte Carlo methods, Partition Functions, Approximate Inference
- Models represent probability distributions
  - Some allow explicit probability evaluation
  - Others don't but allow sampling
  - Some described by graphs and factors, others not

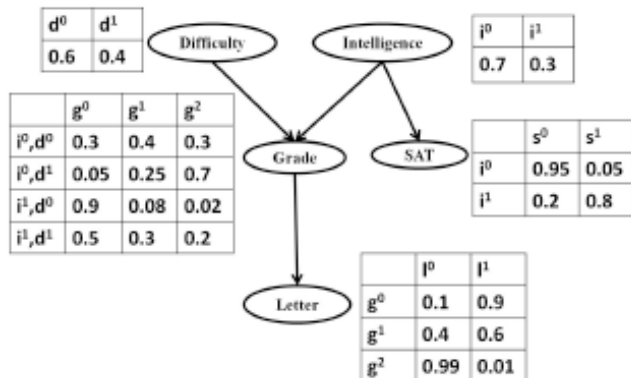
# Directed graphical models

## Representation (Bayesian Network)

- CPDs:  $p(x_i | pa(x_i))$
- Joint Distribution over  $N$  variables

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

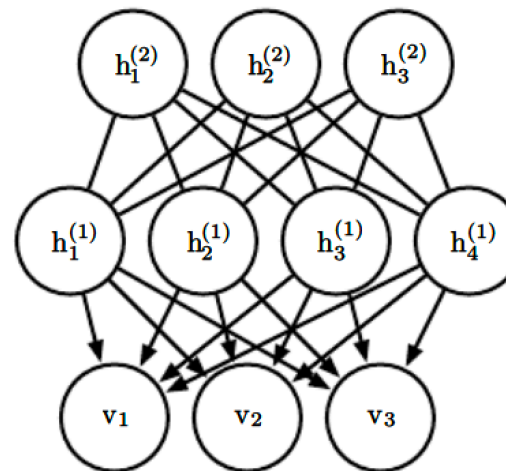
$$\mathbf{x} = \{D, I, G, S, L\}$$



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

## Deep Learning Example

### Deep Belief Network



## Intractability of Inference

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

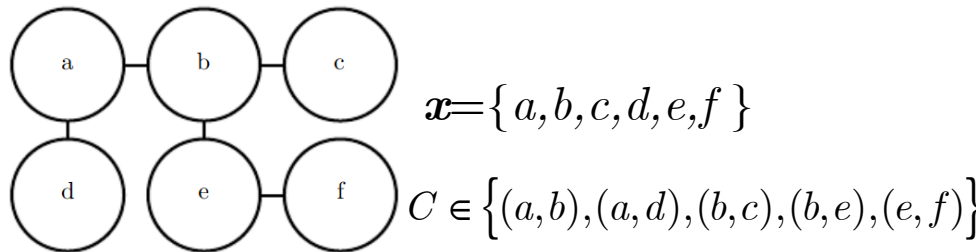
Summing out variables is #P complete

# Undirected Graphical models

## Representation (Markov Network)

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

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



$$p(\mathbf{x}) = \frac{1}{Z} \phi_{a,b}(a, b) \phi_{b,c}(b, c) \phi_{a,d}(a, d) \phi_{b,e}(b, e) \phi_{e,f}(e, f)$$

- **Energy Model**

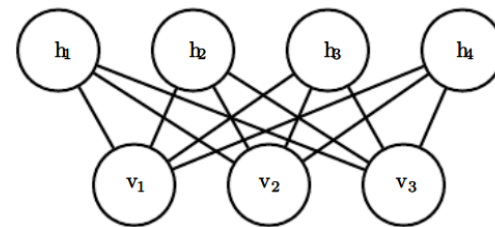
$$E(C) = -\ln \phi(C), \quad \phi(C) = \exp(-E(C))$$

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

$$p(\mathbf{x}) = \frac{1}{Z} \exp\{-E(\mathbf{x})\}$$

## Deep Learning Example

### Restricted Boltzmann machine



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

$$p(\mathbf{h} | \mathbf{v}) = \prod_i p(h_i | \mathbf{v}) \quad \text{and}$$

$$p(\mathbf{v} | \mathbf{h}) = \prod_i p(v_i | \mathbf{h})$$

## Intractability of Inference:

Partition Function:

$$Z = \sum_{\mathbf{x}} E(\mathbf{x})$$



# Generative Models more complex

- Modeling is more complex than classification
  - In classification
    - Both input and output are given
    - Optimization only needs to learn the mapping
  - In generative modeling
    - Data doesn't specify both input  $z$  and output  $x$
    - How to arrange  $z$  space usefully and how to map  $z$  to  $x$
    - Requires optimizing intractable criteria

# RBM training: Contrastive Divergence

Maximize log-likelihood

$$L(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\}; \boldsymbol{\theta}) = \sum_m \log \tilde{p}(\mathbf{x}^{(m)}; \boldsymbol{\theta}) - \sum_m \log Z(\boldsymbol{\theta})$$

For stochastic gradient ascent, take derivatives:

$$g_m = \nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}^{(m)}; \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(m)}; \boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}} \log Z(\boldsymbol{\theta})$$

## Gradient Update

**Positive phase:**

Sample a minibatch of examples  
From the training set

**Negative Phase:**

Gibbs samples from the  
model initialized with training  
data are used to

# Approximate Inference Methods

## 1. Particle-based (sampling) methods

– To evaluate  $E[f] = \int f(z)p(z)dz$  use  $L$  samples  $z^{(l)}$

$$\hat{f} = \frac{1}{L} \sum_{i=1}^L f(z^{(i)})$$

- Directed models: Ancestral (or Forward) sampling
- Undirected models: Gibbs sampling

## 2. Variational inference methods

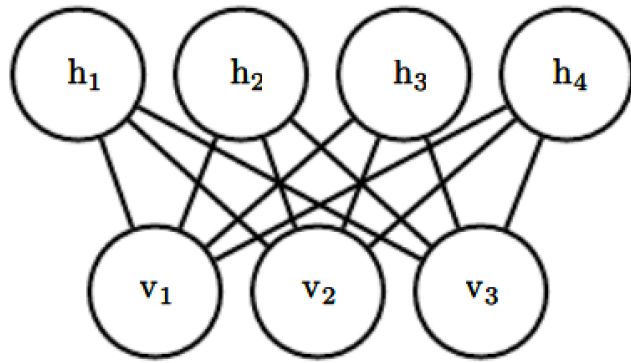
– ELBO  $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q)$  is a lower bound on  $\log p(\mathbf{h}, \mathbf{v}; \boldsymbol{\theta})$

- Inference is viewed as maximizing  $\mathcal{L}$  wrt  $h$

– Mean field approach:

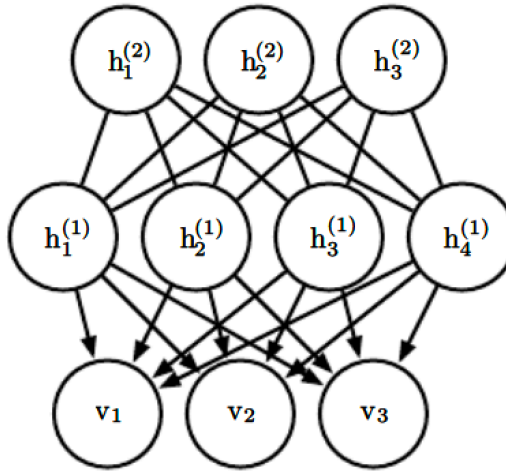
$$q(\mathbf{h} | \mathbf{v}) = \prod_i q(h_i | \mathbf{v})$$

# Deep Generative Models use both directed/undirected graphs



RBM is an undirected graphical model based on a bipartite graph

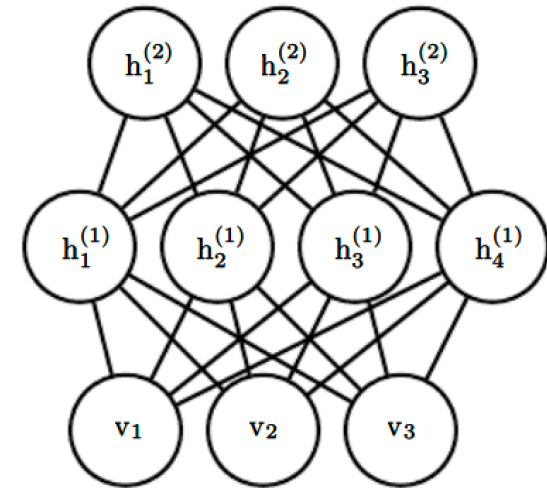
Efficient evaluation and differentiation of  $P(\mathbf{v})$   
Efficient sampling



Deep Belief Network  
Hybrid graphical model with multiple hidden layers

Sampling: Gibbs for top layers, Ancestral for lower

Trained using contrastive divergence



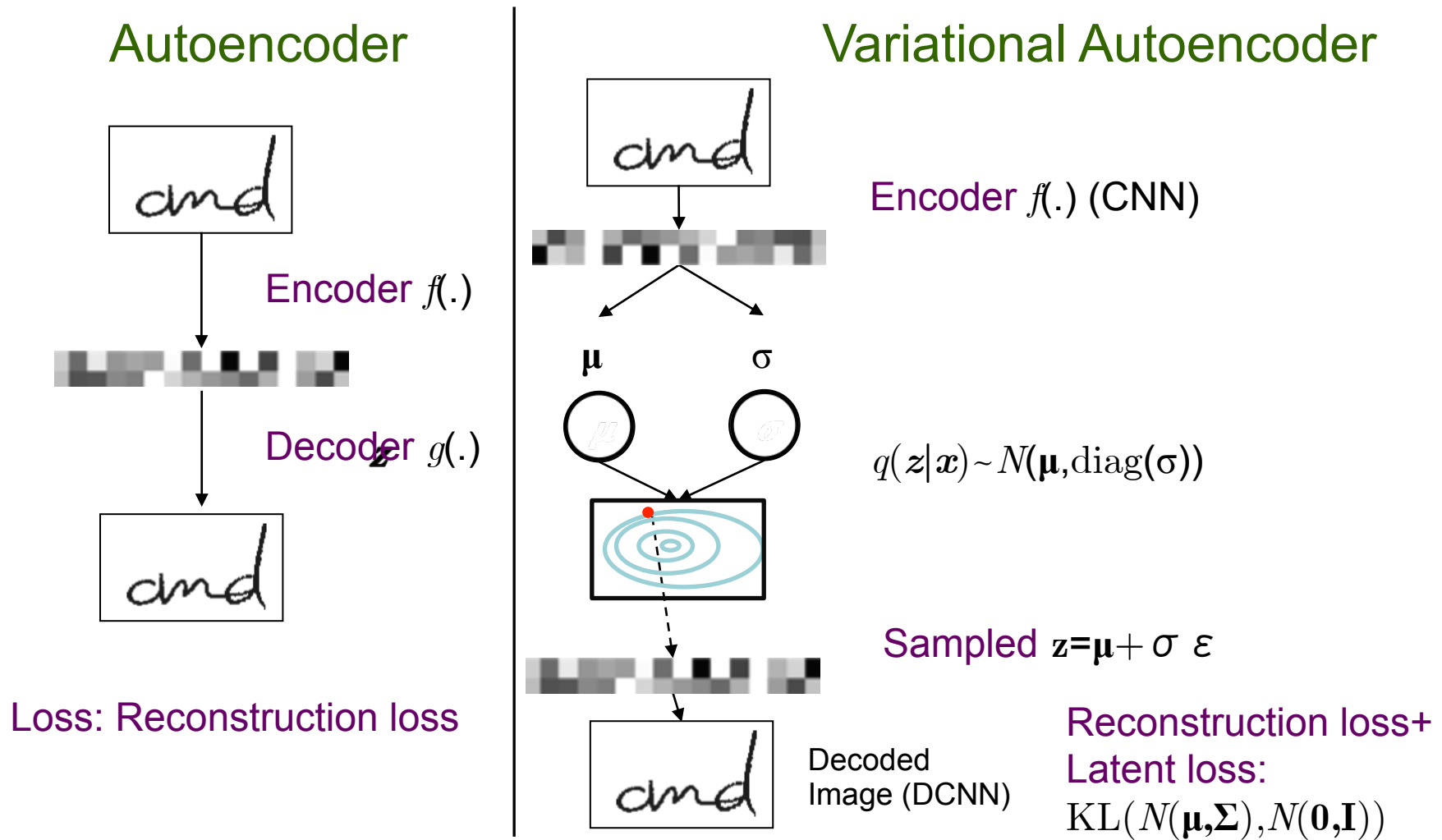
Deep Boltzmann Machine is an undirected graphical model with several layers of latent Variables

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

# Differentiable Generator Nets

- Many generative models are based on the idea of using a *differentiable generator network*
- Model transforms samples of latent variables  $\mathbf{z}$  to samples  $\mathbf{x}$  or to distributions over samples  $\mathbf{x}$  using a differentiable function  $g(\mathbf{z};\theta^{(g)})$  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 vs VAE



We have  $\mathbf{z} \sim \text{Enc}(x) = q(\mathbf{z}|x)$  and  $y \sim \text{Dec}(\mathbf{z}) = p(x|\mathbf{z})$

# Adversarial Example: Driving



(a)



(b)

Self driving car correctly decides to turn left for first image  
But turns right and crashes into guardrail for a slightly darkened version of (a)

# Adversarial Approach

- Carefully craft synthetic images
  - By adding minimal perturbations to an existing image
  - So that they get classified differently than the original picture but still look the same to the human eye.



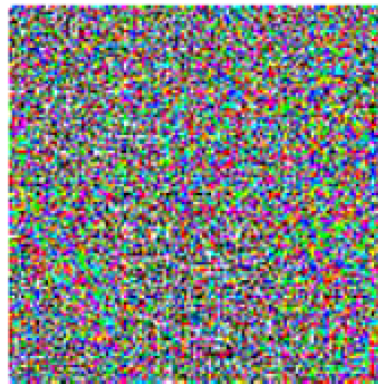
# Adversarial Example Generation

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 ×



=



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

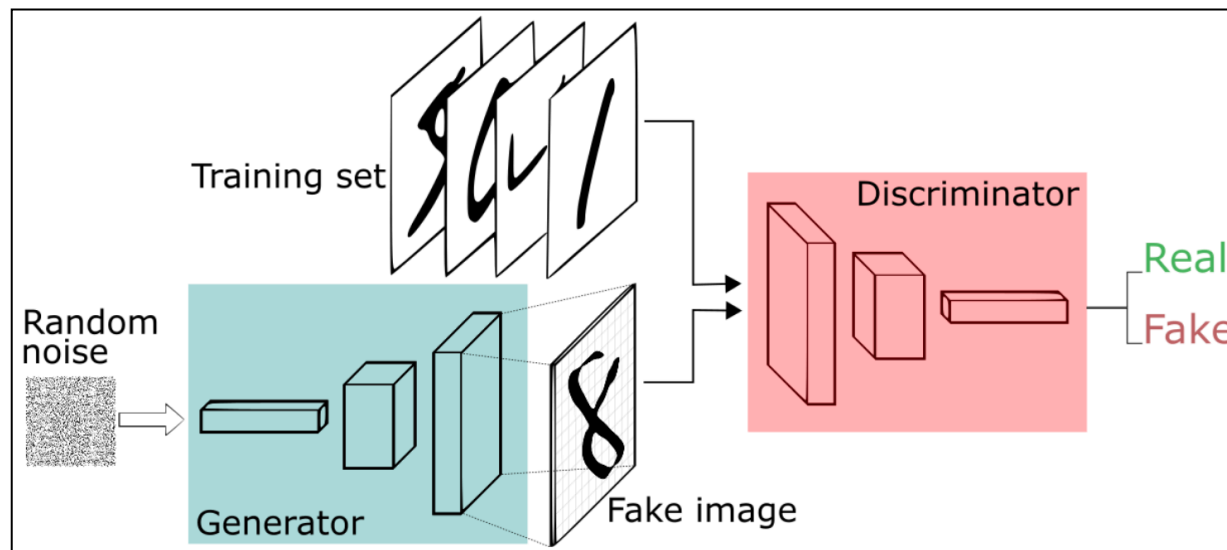
$y$  = "panda"  
with 58% confidence

$y$  = "nematode"  
With 8.2% confidence

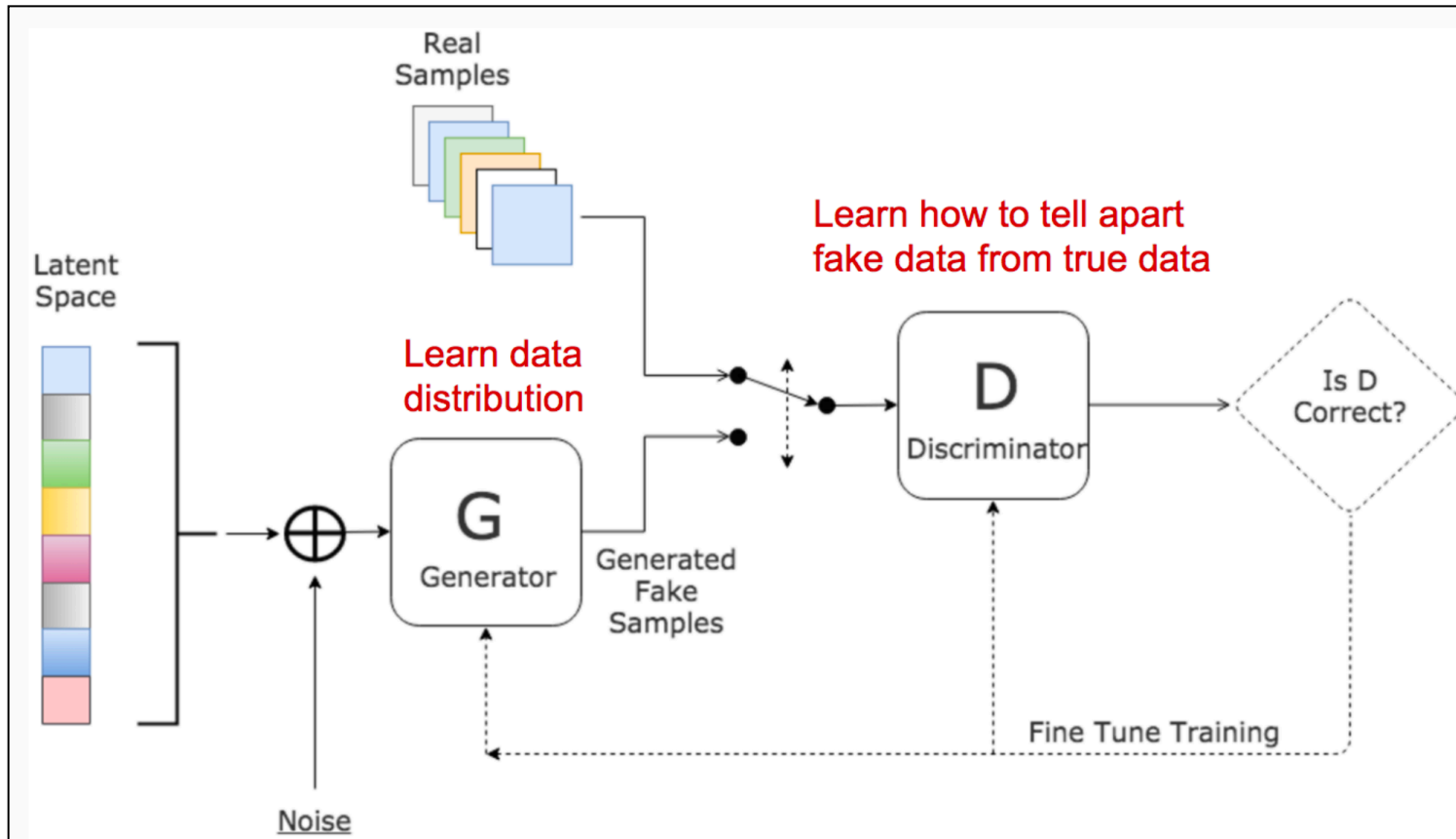
$y$  = "gibbon"  
With 99% confidence

# Generative Adversarial Network

- *A generative model* (G-network)
  - is trained to fool a feedforward classifier
    - A feedforward network that generates images from noise
- *A discriminative model* (D-network)
  - A feedforward classifier that attempts to recognize samples from G as fake and samples from training set as real



# Training a GAN



Training procedure for G is to maximize probability that D 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

# Loss function for GAN

- Both D and G are playing a minimax game in which we should optimize the following loss function

$$\begin{aligned}\min_G \max_D L(D, G) &= \mathbb{E}_{x \sim p_r(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))] \\ &= \mathbb{E}_{x \sim p_r(x)} [\log D(x)] + \mathbb{E}_{x \sim p_g(x)} [\log(1 - D(x))]\end{aligned}$$

# Problems in GANs

- Hard to achieve Nash equilibrium
  - Two models trained simultaneously may not converge
- Vanishing gradient
  - If discriminator behaves badly, generator does not have feedback
  - If discriminator does a great job, gradient drops to zero and learning is superslow
- Mode Collapse
  - During training generator may collapse to a setting where it always produces the same outputs

# Wasserstein GAN (WGAN)

- Minimizes a reasonable and efficient approximation of the Earth Mover's distance instead of KL divergence
- GANs output a probability, WGAN outputs a value
- Advantages
  - Improve the stability of learning,
  - Get rid of problems like mode collapse
  - Provide meaningful learning curves useful for debugging and hyperparameter searches.

# Importance of GANs

- GANs have been at the forefront of research on generative models in the last couple of years.
- GANs have been used for:
  - Image generation
  - image processing
  - image synthesis from captions
  - data generation for visual recognition
  - many other applications,

# Data Sets for GANs

- **LSUN (Large Scale Scene Understanding)**

• <b>Category</b>	<b>Training</b>	<b>Validation</b>
• Bedroom	3,033,042 images (43 GB)	300 images
• Bridge	818,687 images (16 GB)	300 images
• Church Outdoor	126,227 images (2.3 GB)	300 images
• Classroom	168,103 images (3.1 GB)	300 images
• Conference Room	229,069 images (3.8 GB)	300 images
• Dining Room	657,571 images (11 GB)	300 images
• Kitchen	2,212,277 images (34 GB)	300 images
• Living Room	1,315,802 images (22 GB)	300 images
• Restaurant	626,331 images (13 GB)	300 images
• Tower	708,264 images (12 GB)	300 images
• <b>Testing Set</b>	10,000 images	



# Visual Quality of Generated Images



WGAN 2017  
128x128

ICLR 2018  
256 x 256

Architecture/hyperparameters are carefully selected  
Latent variables capture important FoVs

# Images generated from CelebA



# Shortcoming of Adversarial Approach

- While it exposes some erroneous behaviors of a deep learning model, it must limit its perturbations to tiny invisible changes or require manual checks.
- Adversarial images only cover a small part (52.3%) of deep learning system logic

# Final Remarks on Deep Learning

- “Artificial Intelligence is poised to disrupt the world”
  - NITI Aayog, National Strategy for Artificial Intelligence, #Aifor all
- Philosophy of Science
  - Karl Popper
  - Thomas Kuhn