

# Need for Sampling in Machine Learning

Sargur Srihari

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

# Topics in “Need for Sampling”

- Computing Expected Value
- Computing Answers to queries
  - From summation to expected value
- Justification of approximation
- Sampling in Bayesian prediction
- Examples of samples
- Samples in a discrete case

# Common Task: Computing Expectation

- Find expectation of some function  $f(\mathbf{x})$ 
  - with respect to a probability distribution  $p(\mathbf{x})$
  - Components of  $\mathbf{x}$  can be discrete, continuous or combination
- In continuous case

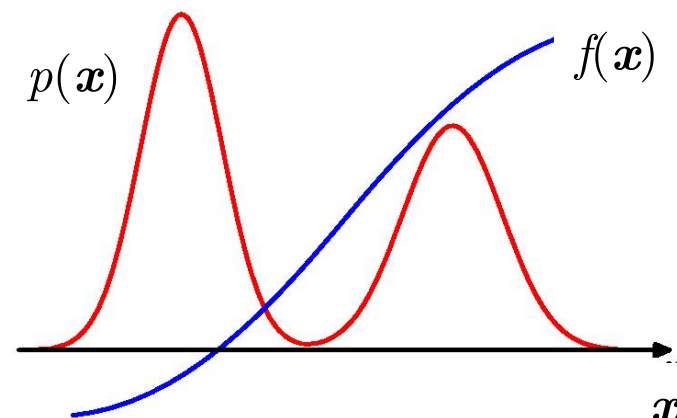
$$E[f(x)] = \int f(x)p(x)dx$$

- In discrete case

$$x \in \{1,2,3,4,5,6\}$$

$$E[x] = \sum_{i=1}^6 x_i p_i = 3.5$$

Illustration



In Bayesian regression,  
From sum and product rules

$$p(t|x) = \int p(t|x,w)p(w)dw$$

where  $p(t|x,w) \sim N(y(x,w), \beta^{-1})$

In general such expectations  
are too complex to be evaluated  
analytically

# Models and queries

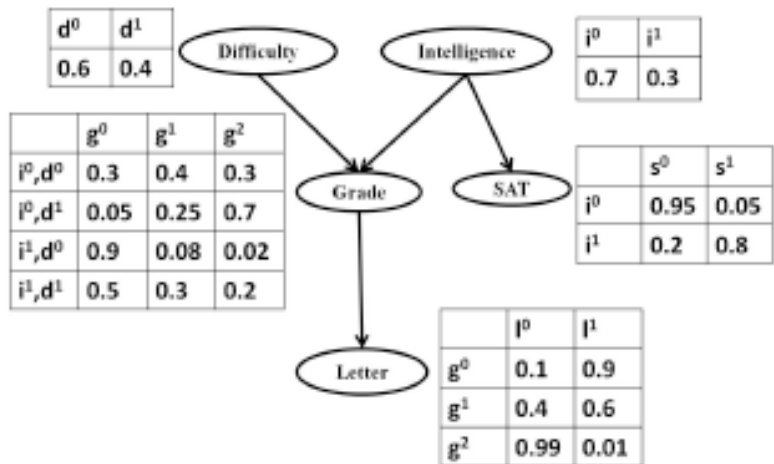
1. ML methods model data with probability distributions. E.g.,  $p(\mathbf{x}, y; \theta)$ 
  - input  $\mathbf{x}$  (*image*), output  $y$  (*steering*), parameters  $\theta$
2. Models are used to answer queries (inference)
  - E.g.,  $p(y|\mathbf{x}; \theta)$ , distribution of  $y$  given  $\mathbf{x}$
3. Sometimes inference is intractable

$$p(y|x) = \frac{p(x,y)}{\sum_y p(x,y)} \quad \text{E.g., } y \text{ is continuous}$$

4. Sampling provides an approximate answer
  - Generative models (Bayesian Networks, Markov Networks) can be readily used to generate samples

# Distributions as graphical models

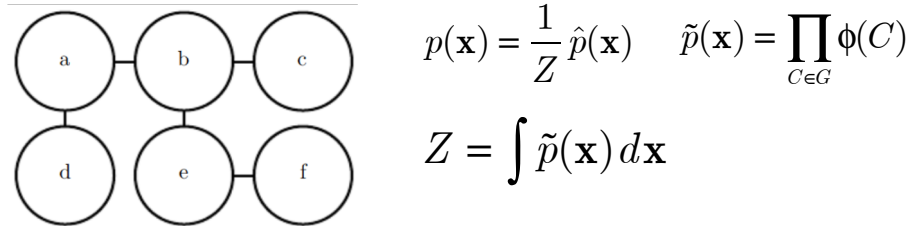
## Bayesian Network



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

$$P(D, I, G, S, L) = P(D)P(I)P(G | D, I)P(S | I)P(L | G)$$

## Markov Network



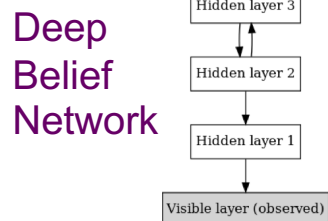
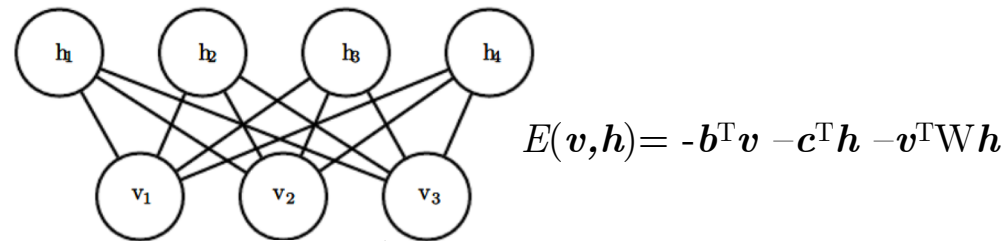
$$p(a, b, c, d, e, f) = \frac{1}{Z} \phi_{a,b}(a, b) \phi_{b,c}(b, c) \phi_{a,d}(a, d) \phi_{b,e}(b, e) \phi_{e,f}(e, f)$$

**Energy model**  $E(a, b, c, d, e, f) = E_{a,b}(a, b) + E_{b,c}(b, c) + E_{a,d}(a, d) + E_{b,e}(b, e) + E_{e,f}(e, f)$

$$\phi_{a,b}(a, b) = \exp(-E(a, b))$$

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

## Restricted Boltzmann machine



# Probability Queries

- Most common type of query is a probability query
- Query has two parts
  - *Evidence*: a subset  $E$  of variables and their instantiation  $e$
  - *Query Variables*: a subset  $Y$  of random variables in network
- **Inference Task:**  $P(Y|E=e)$ 
  - *Posterior probability distribution* over values  $y$  of  $Y$
  - *Conditioned* on the fact  $E=e$
  - Can be viewed as Marginal over  $Y$  in distribution we obtain by conditioning on  $e$

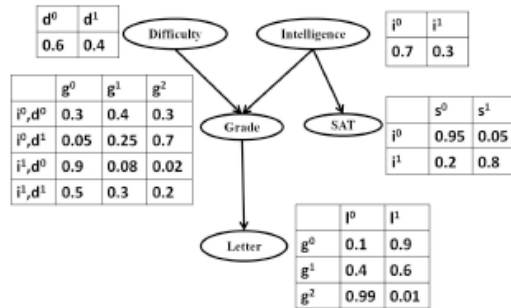
$$P(Y | E = e) = \frac{P(Y, E = e)}{P(E = e)} \quad P(E = e) = \sum_{X \setminus E} \prod_{i=1}^n P(X_i | pa(X_i)) |_{E=e}$$

- **An intractable problem**

- #P complete

# PGMs and queries

## Bayesian Network



$$P(D, I, G, S, L) =$$

$$P(D)P(I)P(G | D, I)P(S | I)P(L | G)$$

## Probabilistic Query

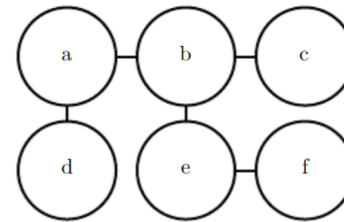
$$P(Y = y^i | E = e) = \frac{P(Y = y^i, E = e)}{P(E = e)}$$

$$P(I = i^1 | L = l^0, S = s^1) = ?$$

$$P(I = i^1, L = l^0, S = s^1) = \sum_{D, I=i^1, G, L=l^0, S=s^1} P(D, I, G, S, L)$$

$$P(L = l^0, S = s^1) = \sum_{D, I, G, L=l^0, S=s^1} P(D, I, G, S, L)$$

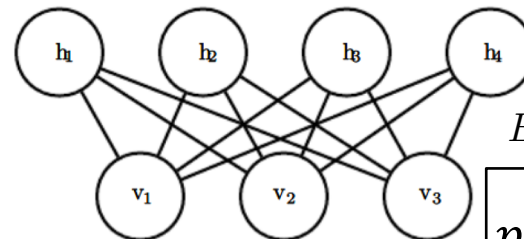
## Markov Network



**Energy model**  $E(a, b, c, d, e, f) =$   
 $E_{a,b}(a, b) + E_{b,c}(b, c) + E_{a,d}(a, d) + E_{b,e}(b, e) + E_{e,f}(e, f)$

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

## Restricted Boltzmann machine



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

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

## Probabilistic Query

$$p(h|\mathbf{v}) = \frac{p(\mathbf{v}, h)}{\sum_{\mathbf{h}} p(\mathbf{v}, h)}$$

Answering queries requires summing over all values of variables to be marginalized

# Basics of MC sampling

- When a sum or integral is intractable
  - E.g., has exponential no of terms and no exact simplification is known
  - It can often be approximated using MC sampling
- The idea is to view the sum or integral as an expectation under some distribution and to approximate the expectation by a corresponding average



# Monte Carlo Methods

- We cannot obtain precise answers to many problems in machine learning
- We can use either deterministic approximate algorithms or Monte Carlo approximations
  - Both approaches are common in machine learning
  - Monte Carlo algorithms return answers to queries with random amount of error
    - For a fixed computational budget, a Monte Carlo algorithm can provide an approximate answer

# Summation $\rightarrow$ Expectation $\rightarrow$ Average

- Sum or integral to approximate is

$$s = \sum_{\mathbf{x}} p(\mathbf{x})f(\mathbf{x}) \quad \text{or} \quad s = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x}$$

- Rewriting expression as an expectation

$$s = E_p[f(\mathbf{x})] \quad \text{or} \quad s = E_p[f(\mathbf{x})]$$

– with the constraint that  $p$  is a probability distribution (for the sum) or a pdf (for the integral)

- Approximate  $s$  by drawing  $n$  samples  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$  from  $p$  and then forming the empirical average

$$\hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}^{(i)})$$

# Justification of Approximation

- The sample average approximation is justified by a few different properties

1. The estimator  $\hat{s}$  is unbiased, since

$$\begin{aligned}
 E[\hat{s}_n] &= \frac{1}{n} \sum_{i=1}^n E[f(\mathbf{x}^{(i)})] && \text{since } \hat{s}_n \text{ is the sample average } \hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}^{(i)}) \\
 &= \frac{1}{n} \sum_{i=1}^n s && \text{since } s = E_p[f(\mathbf{x})] \\
 &= s
 \end{aligned}$$

2. The law of large numbers states that if the samples  $\mathbf{x}^{(i)}$  are i.i.d. then the average converges almost surely to the expected value

$$\lim_{n \rightarrow \infty} \hat{s}_n = s$$

Provided the variance of the individual terms  $\text{Var}[f(\mathbf{x}^{(i)})]$  is bounded

# Variance of estimate

- Consider variance of  $\hat{s}_n$  as  $n$  increases

–  $\text{Var}[\hat{s}_n]$  converges & decreases to 0 if  $\text{Var}[f(\mathbf{x}^{(i)})] < \infty$ :

$$\begin{aligned} \text{Var}[\hat{s}_n] &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}[f(\mathbf{x}^{(i)})] \\ &= \frac{\text{Var}[f(\mathbf{x})]}{n} \end{aligned}$$

- Result tells how to estimate error in MC average

– Equivalently expected error of the approximation

– Compute empirical average of the  $f(\mathbf{x}^{(i)})$  and their empirical variance to determine estimator of  $\text{Var}[\hat{s}_n]$

- Central Limit Theorem tells us that distribution of the average  $\hat{s}_n$  has a normal distribution with mean  $s$  and variance  $\text{Var}[\hat{s}_n]$ . This allows us to estimate confidence intervals around the estimate  $\hat{s}_n$  using the Normal cdf

# Sampling from base distribution $p(\mathbf{x})$

- Sampling from  $p(\mathbf{x})$  is not always possible
- In such a case use importance sampling
- A more general approach is Monte Carlo Markov chains
  - To form a sequence of estimators that converge towards the distribution of interest

# Query Types

## 1. Probability Queries

- Given  $x$  give distribution of  $y$

## 2. MAP (Maximum a posteriori probability)

- What is the most likely setting of  $y$

## 3. Marginal MAP Queries

1. When some variables are known

# Need for Sampling in Bayesian Prediction

- Given training data  $\mathbf{x}$  and  $\mathbf{t}$  and new test point  $x$ , goal is to predict value of  $t$ 
  - *i.e*, wish to evaluate *predictive distribution*  $p(t|x, \mathbf{x}, \mathbf{t})$
- **Predictive distribution** (where parameter has prior/ posteriors)

$$p(t | x, \mathbf{x}, \mathbf{t}) = \int p(t|x, \mathbf{w}) p(\mathbf{w} | \mathbf{x}, \mathbf{t}) d\mathbf{w}$$

With Gaussian noise:  $p(t | x, \mathbf{w}) = N(t | y(x, \mathbf{w}), \beta^{-1})$

Convolution of two Gaussians is Gaussian giving a closed form solution.

- When distributions are complex the integration can be replaced by sampling

# What is a sample?

- Given a set of variables  $\mathbf{x} = \{x_1, \dots, x_d\}$
- A sample is an instantiation of an assignment to all variables

$\mathbf{x}^t = \{x_1^t, \dots, x_d^t\}$ , where  $t$  indicates sample index

- Each variable in a sample takes one possible allowable value in its domain according to a probability distribution defined over  $\mathbf{x}$



# An example of samples

- Scalar variable  $x$  which can take one of  $K$  values  $\{x^0, \dots, x^{K-1}\}$
- Probability distribution  $P(x)$

– With  $K=4$

	$P(x=x^j)$
$x^0$	0.3
$x^1$	0.25
$x^2$	0.27
$x^3$	0.18

- Examples of samples  $x_1, x_2, x_3, x_4 \dots$

$$x_1 = x^2, \quad x_2 = x^0, \quad x_3 = x^1, \quad x_4 = x^0, \dots$$

–  $x^0$  repeats more often than others

- since it is more probable

# Samples generated by GANs

- Deep Convolutional Generative Adversarial Networks



- Laplacian Pyramid GAN

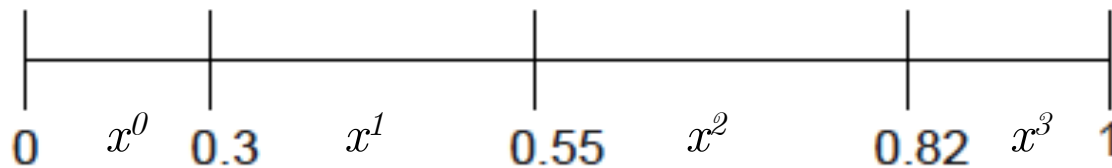


# Algorithm to generate univariate samples

- Domain of  $x$  is  $\{x^0, \dots, x^{K-1}\}$
- Probability distribution is discrete  $P(x)$

1. Divide a real line  $[0,1]$  into  $K$  intervals

such that width of interval is proportional to  $P(x=x^j)$



	$P(x=x^j)$
$x^0$	0.3
$x^1$	0.25
$x^2$	0.27
$x^3$	0.18

1. Draw a random number  $r \in [0,1]$

2. Determine region  $j$  in which  $r$  lies

3. Output  $x^j$

–Random number  $r = 0.2929$ ,  $x = ?$

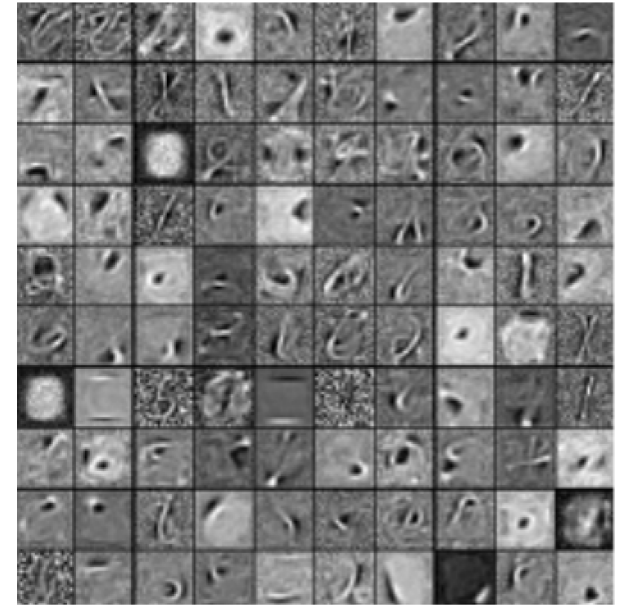
–Random number  $r = 0.5209$ ,  $x = ?$

# Ancestral Sampling for BNs

- Start with lowest numbered node
- Draw a sample from the distribution  $p(x_1)$  which we call  $\hat{x}_1$
- Work through each of the nodes in order
  - For node  $n$  we draw a sample from conditional distribution  $p(x_n | pa_n)$
  - Where parent variables are set to their sampled values
- Once final variable  $x_K$  is sampled
  - Achieved objective of obtaining a single sample from joint distribution
- To sample from marginal distribution
  - Sample from full distribution and discard unnecessary values
  - E.g., to draw from distribution  $p(x_2, x_4)$  simply sample from full distribution, retain values  $x_2^{\wedge}, x_4^{\wedge}$  and discard remaining values  $\{\hat{x}_{j \neq 2,4}\}$

# Samples from a trained RBM

- Gibbs sampling (model trained on MNIST data)



Corresponding weight vectors

- Each column is a separate Gibbs process
- Each row represents the output of another 1000 steps of Gibbs sampling
  - Successive samples are highly correlated