

Structured Probabilistic Models (chapter 16, part 2)

Monte Carlo Methods (chapter 17)

University of Padova, A.A. 2022/23

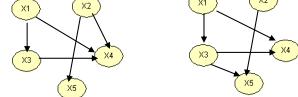
Learning about Dependencies

Learning Graph Structure

- Try out several graphs
- See which graph does best job of some criterion
 - Fitting training set with small model complexity
 - Fitting validation set
- Iterative search, propose new graphs similar to best graph so far
(remove edge / add edge / flip edge)

Unknown structure - learn graph and parameters

➤ Complete data:
optimization (search in space of graphs)



$$S = \arg \max_S \text{Score}(S)$$

1. Log-likelihood-based scores

- MLC - Maximum Likelihood criterion

2. Penalized log-likelihood scores

- BIC - Bayesian Information Criterion (Schwarz, 1978)
- MDL - Minimum Description Length (Rissanen, 1978)
- AIC - Akaike's Information Criterion

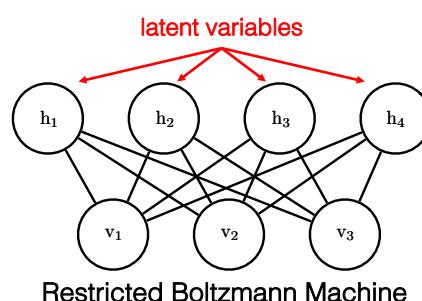
3. Bayesian scores: Bayes, BDe, K2

4. Predictive scores

- Cross-validation (k-Fold-CV)
- Prequential (Preq)

Use Latent Variables

- Use one graph structure
- Many latent variables
- Dense connections of latent variables to observed variables
- Parameters learn that each latent variable interacts strongly with only a small subset of observed variables
- Trainable just with gradient descent; no discrete search over graphs

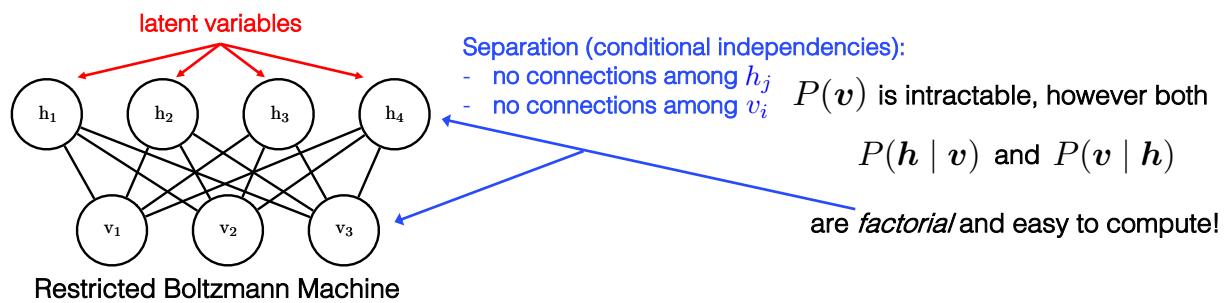


Inference and Approximate Inference

- Inferring marginal distribution over some nodes or conditional distribution of some nodes given other nodes is #P hard
 - NP-hardness describes decision problems
 - #P-hardness describes counting problems, e.g., how many solutions are there to a problem where finding one solution is NP-hard
- Deep Learning usually rely on approximate inference

A.A. 2022/23: Deep Learning

Example: (shallow) Restricted Boltzmann Network



Energy-based Undirected Model

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

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

hidden units:
 not observable variables

visible units:
 observable variables

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

$$\begin{aligned} P(\mathbf{h} | \mathbf{v}) &= \frac{P(\mathbf{h}, \mathbf{v})}{P(\mathbf{v})} \\ &\stackrel{\mathbf{v} \text{ observed: constant}}{=} \frac{1}{P(\mathbf{v})} \frac{1}{Z} \exp \left\{ \mathbf{b}^\top \mathbf{v} + \mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h} \right\} \\ &= \frac{1}{Z'} \exp \left\{ \mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h} \right\} \quad h_j \text{ separation} \\ &= \frac{1}{Z'} \exp \left\{ \sum_{j=1}^{n_h} c_j h_j + \sum_{j=1}^{n_h} \mathbf{v}^\top \mathbf{W}_{:,j} h_j \right\} \\ &= \frac{1}{Z'} \prod_{j=1}^{n_h} \exp \left\{ c_j h_j + \mathbf{v}^\top \mathbf{W}_{:,j} h_j \right\} \end{aligned}$$

A.A. 2022/23: Deep Learning

Example: (shallow) Restricted Boltzmann Network

From unnormalized to normalized distributions

$$\begin{aligned} P(h_j = 1 \mid \mathbf{v}) &= \frac{\tilde{P}(h_j = 1 \mid \mathbf{v})}{\tilde{P}(h_j = 0 \mid \mathbf{v}) + \tilde{P}(h_j = 1 \mid \mathbf{v})} \\ &= \frac{\exp \{c_j + \mathbf{v}^\top \mathbf{W}_{:,j}\}}{\exp \{0\} + \exp \{c_j + \mathbf{v}^\top \mathbf{W}_{:,j}\}} \\ &= \sigma(c_j + \mathbf{v}^\top \mathbf{W}_{:,j}) . \end{aligned}$$

Recall that:
 $\sigma(x) = \frac{\exp(x)}{\exp(x) + \exp(0)}$

Finally, for the hidden layer:

$$P(\mathbf{h} \mid \mathbf{v}) = \prod_{j=1}^{n_h} \sigma((2\mathbf{h} - 1) \odot (\mathbf{c} + \mathbf{W}^\top \mathbf{v}))_j$$

$$1 - \sigma(x) = \sigma(-x)$$

...and with a similar derivation, for the visible layer:

$$P(\mathbf{v} \mid \mathbf{h}) = \prod_{i=1}^{n_v} \sigma((2\mathbf{v} - 1) \odot (\mathbf{b} + \mathbf{W}\mathbf{h}))_i$$

Example: (shallow) Restricted Boltzmann Network

Why do we care about efficient computation of

$$P(\mathbf{h}|\mathbf{v}) \text{ and } P(\mathbf{v}|\mathbf{h})?$$

Because of learning:

- Given training data \mathbf{x} (to be clamped to the visible units) we need to maximize $p(\mathbf{x}; \boldsymbol{\theta})$
- We know it is easier to maximize the Log-Likelihood $\log p(\mathbf{x}; \boldsymbol{\theta})$, however its gradient with respect to $\boldsymbol{\theta}$ is

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

- Computing the term

$$\nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}; \boldsymbol{\theta})$$

is not a problem for RBMs, however the term

$$\nabla_{\boldsymbol{\theta}} \log Z(\boldsymbol{\theta})$$

is more problematic!

Monte Carlo Methods

Under some conditions satisfied by RBMs (as well as most ML models), it can be shown that

$$\nabla_{\theta} \log Z = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \nabla_{\theta} \log \tilde{p}(\mathbf{x})$$

which is not possible to compute exactly, however we can use

Monte Carlo methods

to compute a (good) approximation of it!

Randomized Algorithms	Las Vegas	Monte Carlo
Type of Answer	Exact	Random amount of error
Runtime	Random (until answer found)	Chosen by user (longer runtime gives less error)

A.A. 2022/23: Deep Learning

Monte Carlo Methods: basic idea

$$s = \sum_{\mathbf{x}} p(\mathbf{x}) f(\mathbf{x}) = E_p[f(\mathbf{x})]$$

$$s = \int p(\mathbf{x}) f(\mathbf{x}) d\mathbf{x} = E_p[f(\mathbf{x})]$$

Estimating sums/integrals with samples

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

$\mathbf{x}^{(i)} \sim p$

A.A. 2022/23: Deep Learning

Monte Carlo Methods: basic idea

Unbiased:

- The expected value for finite n is equal to the correct value

$$\mathbb{E}[\hat{s}_n] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[f(\mathbf{x}^{(i)})] = \frac{1}{n} \sum_{i=1}^n s = s.$$

- The value for any specific n samples will have random error, but the errors for different sample sets cancel out

Low variance:

- Variance is $O(\frac{1}{n})$
- For very large n , the error converges “almost surely” to 0

However sampling from $p(\mathbf{x})$ (i.e., $\mathbf{x}^{(i)} \sim p$) is not always possible!

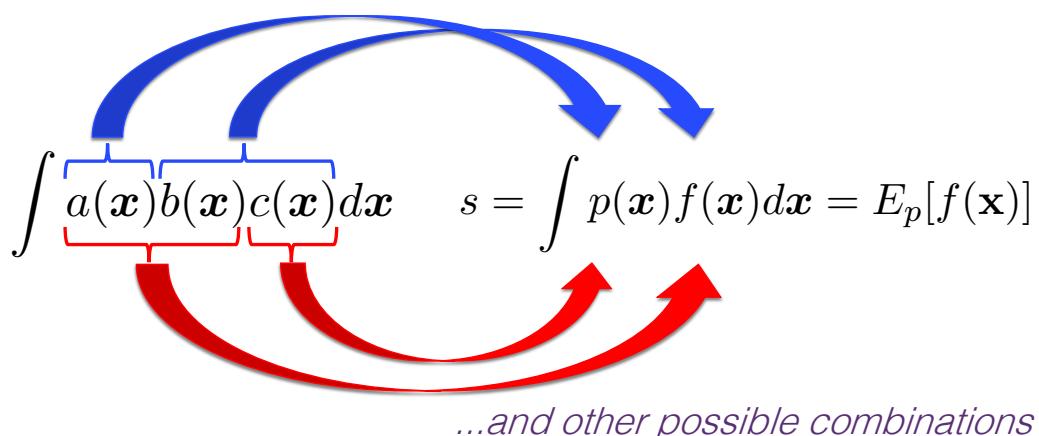
A.A. 2022/23: Deep Learning

Importance Sampling

Suppose we want to compute (an approximation of)

$$\int a(\mathbf{x})b(\mathbf{x})c(\mathbf{x})d\mathbf{x}$$

? Non-unique decomposition!



A.A. 2022/23: Deep Learning

Importance Sampling

$$p(\mathbf{x})f(\mathbf{x}) = q(\mathbf{x}) \frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}$$

This is our new p , meaning it is the distribution we will draw samples from

This ratio is our new f , meaning we will evaluate it at each sample

Why is this useful ?

- Maybe it is feasible to sample from q but not from p
- A good q can reduce the variance of the estimate
- Importance sampling is still unbiased for every q

A.A. 2022/23: Deep Learning

Optimal q (optimal importance sampling)

importance sampling estimator

$$\hat{s}_q = \frac{1}{n} \sum_{i=1, \mathbf{x}^{(i)} \sim q}^n \frac{p(\mathbf{x}^{(i)})f(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})}$$

Unbiased!

$$\text{Var}[\hat{s}_q] = \text{Var}\left[\frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}\right]/n$$

$$\mathbb{E}_q[\hat{s}_q] = \mathbb{E}_q[\hat{s}_p] = s$$

We can choose q so to reduce variance!!

The minimum variance occurs when q is

$$q^*(\mathbf{x}) = \frac{p(\mathbf{x})|f(\mathbf{x})|}{Z}$$

- Determining the optimal q requires solving the original integral, so not useful in practice
- Useful to understand intuition behind importance sampling
- Places more mass on points where the weighted function is larger

A.A. 2022/23: Deep Learning

Biased Importance Sampling

If we need to use only the unnormalized distribution

$$\begin{aligned}\hat{s}_{BIS} &= \frac{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})}} \\ &= \frac{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})}} \quad \text{Asymptotically unbiased} \\ &= \frac{\sum_{i=1}^n \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})}} \quad \text{poor choice of } q \text{ can make the efficiency much worse}\end{aligned}$$

Markov Chain

So far we have assumed we can sample from p or q easily

- This is true when p or q has a directed graphical model representation
 - Use ancestral sampling
 - Sample each node given its parents, moving from roots to leaves
- Sampling from undirected models is more difficult
 - Cannot get a fair sample in one pass
 - Use a Monte Carlo algorithm that incrementally updates samples, comes closer to sampling from the right distribution at each step
 - This is called a Markov Chain

Markov Chain

A Markov chain is defined by

- a random state x
- a transition distribution $T(x'|x)$ specifying the probability that a random update will go to state x' if it starts in state x

Running the Markov chain means repeatedly updating the state x to a value x' sampled from $T(x'|x)$

Single Chain (can have more chains running in parallel!)

- $q^{(t)}(x)$: distribution at iteration t from which x is drawn
- Aim: $q^{(t)}(x) \xrightarrow{t} p(x)$

A.A. 2022/23: Deep Learning

Markov Chain

- positive integer i to represent a state (countably many states)
- can describe $q(x)$ as a vector v where $v_i = q(x = i)$
- probability of a single state landing in state x' is given by

$$q^{(t+1)}(x') = \sum_x q^{(t)}(x)T(x'|x)$$

- can represent the effect of the transition operator T using a *stochastic matrix* A (each of its columns represents a probability distribution)

$$A_{i,j} = T(x' = i|x = j)$$

Perron-Frobenius theorem \Rightarrow the largest eigenvalue of A is real and equal to 1

- then $v^{(t)} = Av^{(t-1)} \Rightarrow v^{(t)} = A^t v^{(0)} \Rightarrow v^{(t)} = V \text{diag}(\lambda)^t V^{-1} v^{(0)}$
- largest eigenvalue is 1 \Rightarrow converges to a stationary distribution (equilibrium distribution)

$$v' = Av = v$$

converges to p for the right choice of T

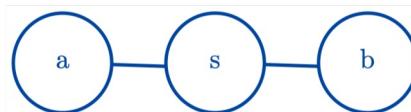
A.A. 2022/23: Deep Learning

Simple Markov Chain: Gibbs sampling

- Repeatedly cycle through all variables
 - For each variable, randomly sample that variable given its Markov blanket
 - For an undirected model, the Markov blanket is just the neighbors in the graph
- *Block Gibbs* trick: conditionally independent variables may be sampled simultaneously (used in RBMs!)

Example:

- Initialize a , s , and b
- For n repetitions
 - Sample a from $P(a|s)$ and b , from $P(b|s)$
 - Sample s from $P(s|a, b)$



Block Gibbs trick lets us sample a and b in parallel

Problems...

Generally infeasible to...

- ... know ahead of time how long *mixing* will take
- ... know how far a chain is from equilibrium
- ... know whether a chain is at equilibrium

Usually in deep learning we just run for n steps, for some n that we think will be big enough, and hope for the best

In practice

- Mixing can take an infeasibly long time
- This is especially true for
 - High-dimensional distributions
 - Distributions with strong correlations between variables
 - Distributions with multiple highly separated modes

