

Variational Methods
Convex Relations
Monte Carlo Methods

Admin: tutorials this week!
Midterm solutions: tonight
Marked A7; due Friday
A8: out tonight, due Dec. 17th
Coding project: due Dec. 17th
Final Project: marking scheme out this week
report due Dec. 17th

Loopy Belief Propagation

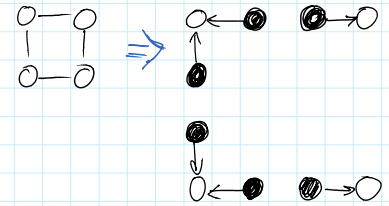
- does not require $q(x)$ to be valid distribution, only "locally consistent" (adjacent nodes agree on marginals)
- relation to "turbo codes" in information theory
- algorithm: apply sum-product updates repeatedly (exact for trees)
- convergence issues, but now exist convex/convergent versions.
- LBP is only for Gaussian/Discrete, generalization is "expectation propagation"

⊗ Pseudo-likelihood

- For learning parameters of UGM, common and consistent and convex approximation is "pseudo-likelihood"

$$p(x) \propto \prod_c \Psi_c(x) \approx \prod_j p(x_j | x_{r_j})$$

- training is now easy as in DAGs
convex, fast



Convex Relaxations to Decoding

multisite

$$\max_{x \in \{1, 2, \dots, S\}^d} \prod_{j=1}^d \phi_j(x_j) \prod_{(j,k) \in E} \phi_{jk}(x_j, x_k)$$

$$z_j = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \leftarrow \text{the one we choose}$$

$$z_{jk} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

binary

$$\max_{\substack{z_j^s \in \{0, 1\} \\ z_{jk}^{ss'} \in \{0, 1\}}}$$

$$\sum_{j=1}^d \sum_{s=1}^S z_j^s \log \phi_j(s) + \sum_{(j,k) \in E} z_{jk}^{ss'} \log \phi_{jk}(s, s') \Rightarrow \text{"Integer linear program"} \rightarrow \text{NP-Hard}$$



subject to $\sum_{s=1}^S z_j^s = 1$ (can only pick one state)

$$\sum_{s,s'} z_{jk}^{ss'} = z_j^s \quad (\text{edge labellings agree with node labellings})$$

$$\sum_{s,s'} z_{jk}^{ss'} = z_k^{s'}$$

$$z_j^s \in [0, 1] \\ z_{jk}^{ss'} \in [0, 1] \quad \text{interval, not binary}$$

"Linear Programming Relaxation"

↳ polynomial time

if get 0, 1 only: equivalent solution to integer (but unlikely to get)

Whole values: part of "some optimal solution"

Monte Carlo Methods

Methods for intractable inference problems

variational: approximate $p(x)$ with simple "simple" $q(x)$

Monte Carlo: - approximate $p(x)$ with S samples $x^s \sim p(x)$

- use x^1, x^2, \dots, x^S to approximate

good for:

- $p(x_s)$
- $p(\hat{\theta} | D)$
- $\max_x p(x)$

Ex: $\theta \sim \text{Dirichlet}$

θ_1	θ_2	θ_3
0.1	0.9	0
0.2	0.8	0
0.2	0.7	0.1
0.1	0.7	0.2

$E[\theta_1]$?

$$= (0.1 + 0.2 + 0.2 + 0.1) / 4 = 0.15$$

But the hard part is selecting good samples!

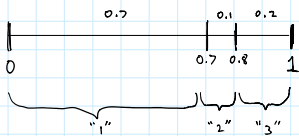
Variational: fast but biased

Monte Carlo: slow but consistent (add more samples - slower but better)

Sampling from Standard Distribution

(⊕ assume you can generate $u \in U(0,1)$)

If $x \sim \text{Multinomial}(0.7, 0.1, 0.2)$



Generate $u \sim U(0,1)$, if it falls in "2", $x^S = "2"$

For univariate continuous, use inverse CDF

For R.V. X , CDF is $F(x) = p(X \leq x)$

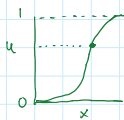
Inverse CDF: $F^{-1}(p) = \{x^i \mid F(x^i) = p\}$

"quantile function"
ex: $F^{-1}(0.5) = \text{median}$

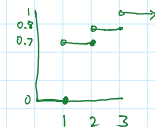
Ex: Gaussian



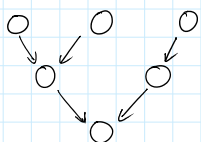
- generate u , put in F^{-1} , get value in x



(discrete represented Equivalency:
in this form:)



Sampling From DAG



← sample first

← sample conditioned on parents

← sample last

"Ancestral Sampling"

Special case:

x is a multivariate normal

$$x \sim N(\mu, \Sigma)$$

- Generate $x \sim N(0, I)$

- set: $y = \mu + Lx$, where $LL^T = \Sigma$

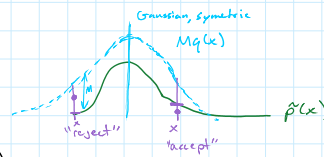
Rejection Sampling

- don't have inverse CDF or too expensive

- sample from proposal $q(x)$, where

$$Mq(x) \geq \tilde{p}(x) \quad (M \text{ should be small})$$

Sample x from $q(x)$, $u \in U(0,1)$
 "accept sample if $u \geq \frac{\tilde{p}(x)}{Mq(x)}$ "



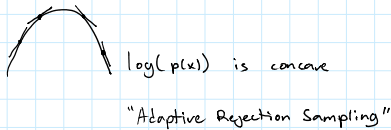
Eg: Sample from posterior: $p(\theta|D) = p(D|\theta)p(\theta)/p(D)$

Target: $\tilde{p}(\theta) = p(D|\theta)p(\theta)$

Proposal: $q(\theta) = p(\theta)$

We have: $M = p(D|\theta_{MLE})$ (MLE: satisfies inequality)

Eg 2



Importance Sampling

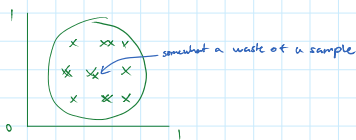
Want to compute $E[h(x)] = \int h(x)p(x) dx$

- represent as: $\int h(x) \frac{p(x)}{q(x)} q(x) dx$ q non zero if p non zero

- sample from $q(x)$, weigh each sample by $\frac{p(x)}{q(x)}$

$$E[h(x)] \approx \sum_{s=1}^S w_s h(x^s)$$

Something to keep in mind:



Sequential Importance Sampling (SIS, Particle filter)

- generate from sequence of proposed distributions $q_n(x)$
- used for tracking, non-Gaussian/non-decrete time series modelling.

MCMC (Markov Chain - Monte Carlo)

High-dimensional: integration

key idea: construct Markov chain with $p(x)$ as stationary distribution

Eg we will do a random walk through X , where asymptotically we spend $p(x)$ time at x .
 - use these dependent samples for Monte Carlo integration



Gibbs Sampling

(sampling version of coordinate descent)

- sample $X_1 \sim p(X_1 | X_2, X_3)$
- sample $X_2 \sim p(X_2 | X_1, X_3)$
- sample $X_3 \sim p(X_3 | X_1, X_2)$
- repeat

- Reduce high dimensional sampling to 1D sampling

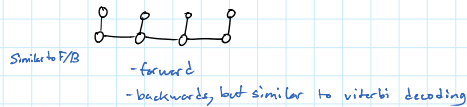
- May require "burn in" (if you start in area of low probability)
(don't want to be biased by start point)

- For graphical models, only depend on neighbours
(and co-parents for DAGs)



- "Block Gibbs Sampling"

- "Rao-Blackwellized"



Metropolis-Hastings Algorithm

- Generalization that allows arbitrary "proposal"

Metropolis algorithm for symmetric proposal:

- at some x , propose x'
- accept with probability $\min(1, \frac{p(x')}{p(x)})$
- sample from arbitrary density $p(x)$
- adapted to asymmetric

Simulated Annealing
 $\min(1, \exp(\frac{-f(x) - f(x')}{T}))$

T large, can go anywhere
 T small "only uphill"