

TAT GIO: March 2nd, 2017

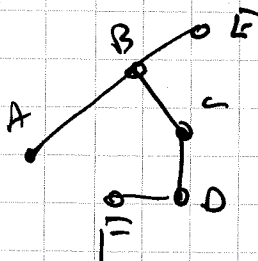
A. Rhodes

Intro To Graphical Models

- Why deploy graphical models in ML-related domains? A:
- ① They are often interpretable and give a principled, model-based framework of a phenomenon of study (compare to a "black box" NN)
 - ② We can easily encode domain/structural knowledge into a GM.
 - ③ They offer computational savings, since they compactly encode: $p(x|\theta)$.

In general, there are ② paradigms for GMs:
C.J. Pearl 1988

① MRFs (undirected GMs)



vs.

② Bayes Nets (Directed) ("Causal Models")
(Belief Nets)



NB: MRFs & Bayes Nets are not equivalent!

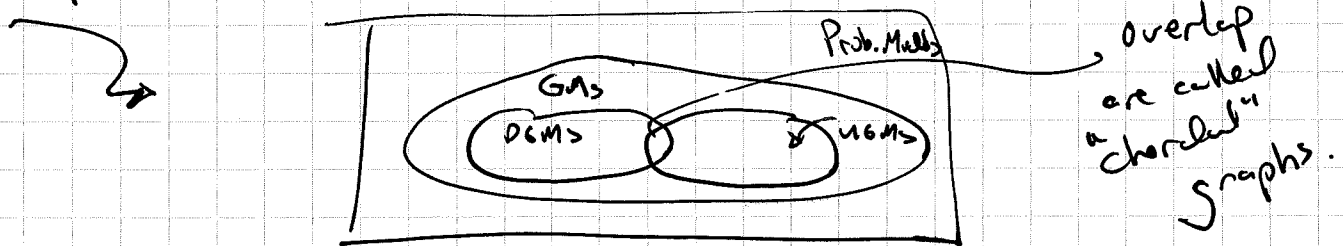
2

More generally, we may ask: What is the expressive power of an MRF vs. a Bayes Net?

We say: G (a graph) is an I-map of a distribution P if $I(G) \subseteq I(P)$, which is to say G does not make any conditional independence statements that are not true of P itself (Here: $I(P)$ stands for the collection of independence assumptions encoded by P).

If $I(G) = I(P)$ we say G is a "perfect" map of P .

It turns out that the set of all probability distributions that admit of an MRF representation is different from the set of distributions that can be represented by DAGs (i.e. Bayes Nets)



Let's further draw out the distinction b/w UGMs & DGMs for deeper understanding.

UGMs (MRFs)

Advantages of UGMs over DGMs:

- ① UGMs are "symmetric" & thus "natural" in certain domains (e.g. vision)
- ② They can be used discriminatively (so they can be fine-tuned using labeled data)

Disadvantages:

- ① Parameters can be less interpretable
- ② Parameter estimation is computationally burdensome.

Key Properties of MRFs

③ Equivalent Notions of CI (conditional indep.)

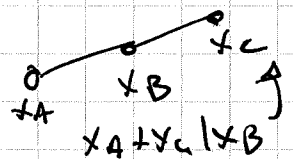
① Global Markov Property:

$X_A \perp X_B \mid X_C$ iff Nodes in C separate A & B in graph.



② Local Markov Property:

$X_A \perp X_B \mid N(X_A)$, $X_B \notin N(X_A)$



③ Pairwise Markov Property

4

$$X_A \perp X_B \mid V(G) \setminus \{X_A, X_B\} \iff G_{AB} = \emptyset$$

↗ This says X_A & X_B are conditionally independent given all other vertices in the graph iff $X_A \not\sim X_B$ in G .

Recall: ① \iff ② \iff ③

More succinctly, Define the Markov Blanket

of a node: X_A as the smallest set of nodes in G that renders X_A conditionally independent of all other nodes in the graph.

Notationally: $X_A \perp V(G) \setminus \text{closure}(X_A) \mid \text{MB}(X_A)$

where $\text{closure}(X_A) \triangleq N(X_A) \cup \{X_A\}$ &

$\text{MB}(X_A) \triangleq$ Markov Blanket of X_A .

[NB]: $\text{MB}(X_A) \triangleq N(X_A)$ for an MRF,

so (A) is equivalent to ②, the Local Markov Property.

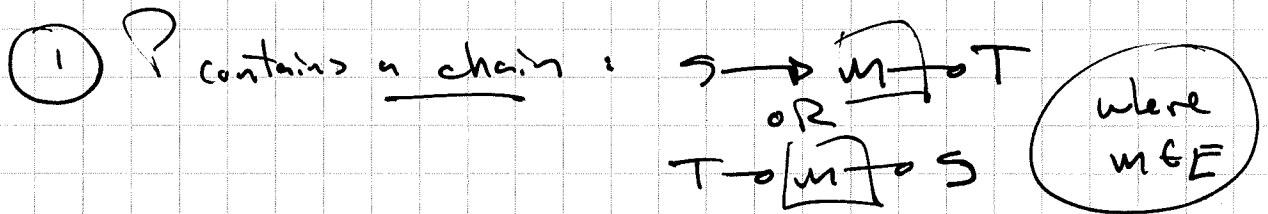
Conditional Independence Properties of DGMs (Bonus)

Note that CI properties are less "intuitive" for DGMs vs. JGMs.

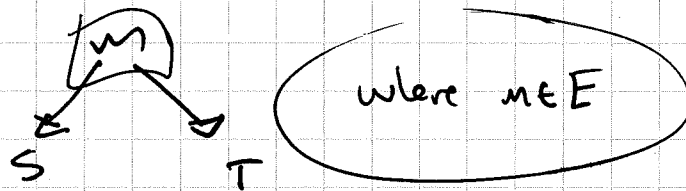
CI properties ^{for DGMs} are codified as the notion of cl-separation.

All told, there are 3 distinct cases we need to consider:

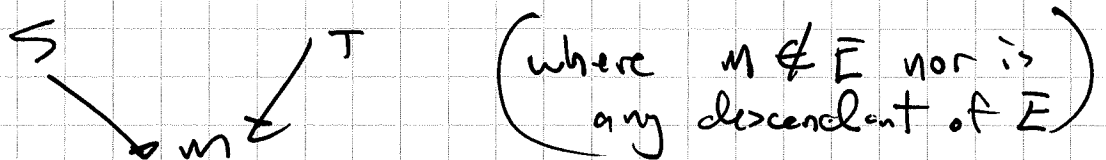
We say an undirected path P is cl-separated by a set of nodes " E " (containing evidence/observation) iff at least one of the following holds:



2) P contains a fork:



3) P contains a v-structure:



Note: ↘

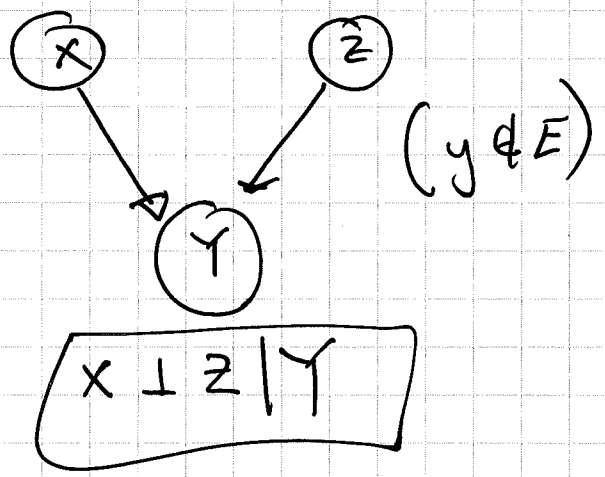
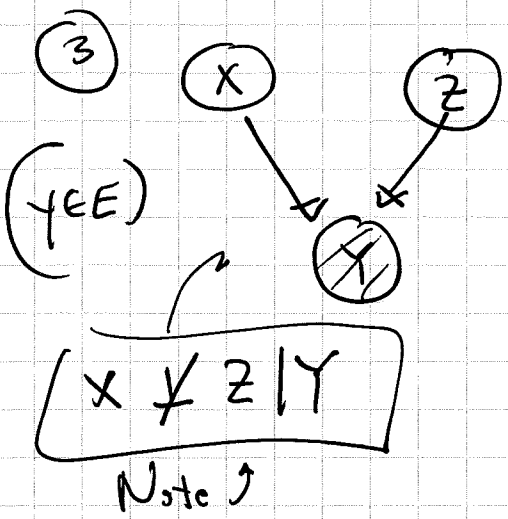
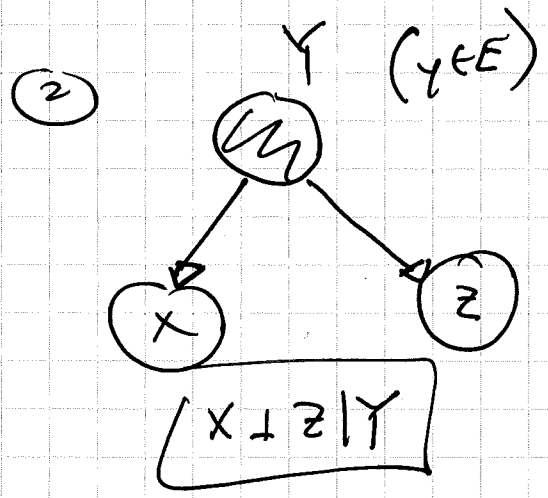
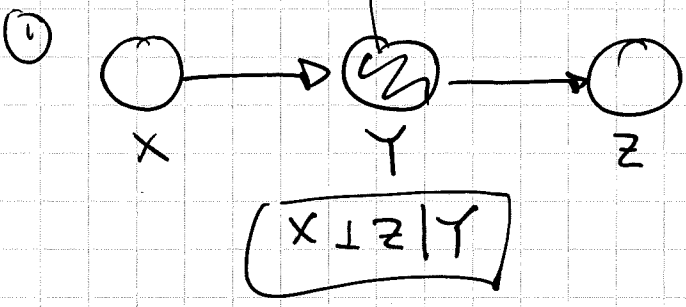
FIVE STAR. ★★★★★ FIVE STAR. ★★★★★ FIVE STAR. ★★★★★ FIVE STAR. ★★★★★

The CI properties of a DGM (Technically: an acyclic DGM) are defined thus:

$X_A \perp X_B \mid X_E \iff A \text{ is d-separated from } B \text{ given } E$

Remark: It is helpful to think of d-separation as "stopping the flow of information" from one set of nodes to another. ("Bayes Ball Theorem")

Examples



FIVE STAR. ★★★★★

FIVE STAR. ★★★★★

FIVE STAR. ★★★★★

FIVE STAR. ★★★★★

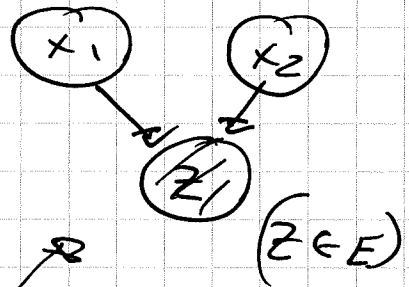
Remark: Don't be unphased by case (3).

The v-structure case is simply the effect known as "explaining away" evidence (Berkeley's Paradox).

Ex. For easy illustration: where $X_1 \perp X_2$

Suppose $X_1, X_2 \sim \text{Ber}(.5)$ $S = \{0, 1\}$

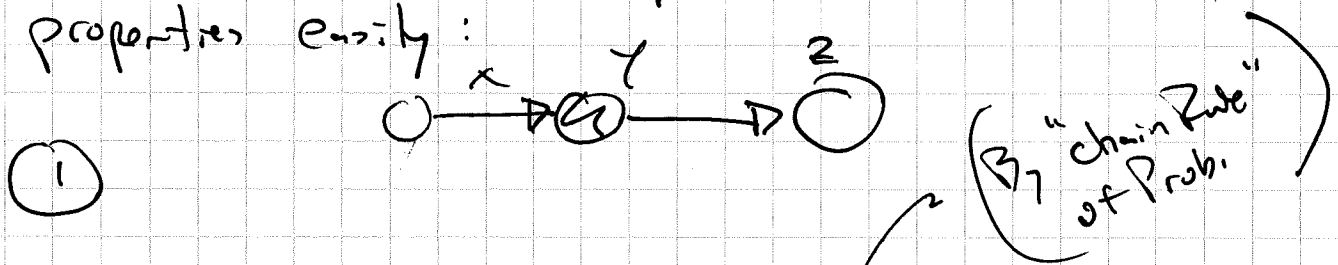
Define $Z = X_1 + X_2$



Then $P(X_1 = 0 | Z = 0) = 1$,
 $P(X_2 = 1 | Z = 0) = 0$, etc.

so, $X_1 \not\perp X_2 | Z$, as promised.

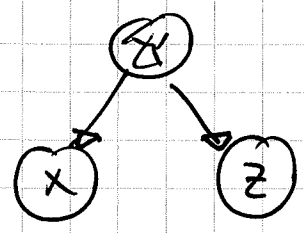
In more detail, we can prove the CI/d-separation properties easily:



$$p(x, y, z) = p(x)p(y|x)p(z|y)$$

condition on y:
$$p(x, z|y) = \frac{p(x)p(y|x)p(z|y)}{p(y)} = \frac{p(x, y)p(z|y)}{p(y)} = p(x|y)p(z|y)$$

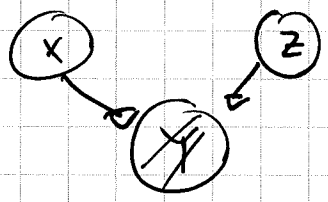
2



$$p(x, y, z) = p(y)p(x|y)p(z|y)$$

$$\rightarrow p(x, z|y) = \frac{p(y)p(x|y)p(z|y)}{p(y)} = p(x|y)p(z|y) \quad \square$$

3



$$p(x, y, z) = p(x)p(z)p(y|x, z)$$

$$\text{Suppose } y \in E \rightarrow p(x, z|y) = \frac{p(x)p(z)p(y|x, z)}{p(y)} \neq p(x|y) \cdot p(z|y)$$

clearly, $x \not\perp z | y$

However, $p(x, z) = p(x) \cdot p(z)$ for the unconditional distribution.

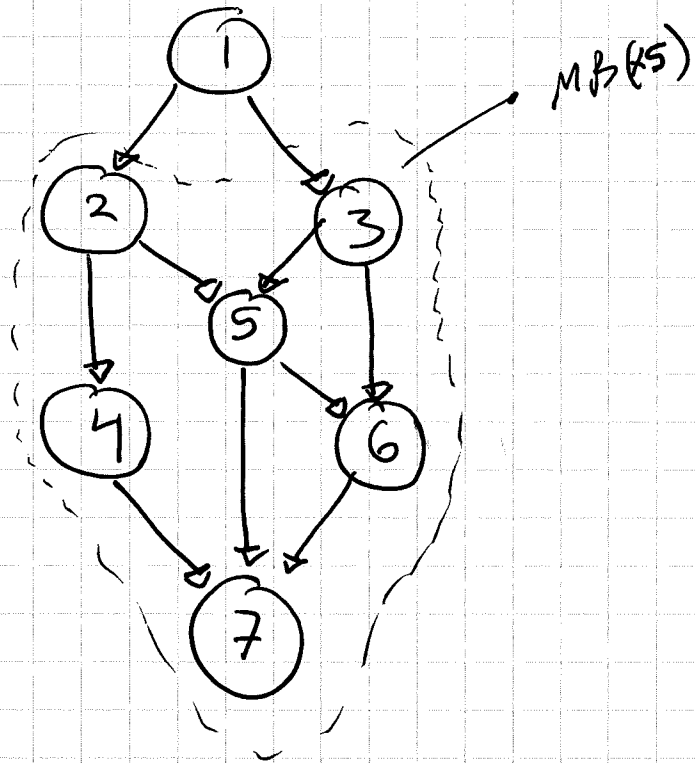
By analogy with MRTs, we can define the

Markov Blanket for DGMs.

Recall that the M.B. of a node in a graph effectively renders that node conditionally independent of all other nodes in the graph.

For a DGM, $MB(x_A) \triangleq \text{Children}(x_A) \cup \text{Parents}(x_A) \cup \text{Co-parents}(x_A)$

Example:



$MB(x_5) = \{x_2, x_3, x_4, x_5, x_6\}$

$\underbrace{x_4}_{\text{Children}}$
 $\underbrace{x_2, x_3}_{\text{Parents}}$
 $\underbrace{x_6}_{\text{Co-parent}}$

You are welcome to verify, for instance, that:

$$P(x_5 | x_{-5}) \propto P(x_5 | x_2, x_3) P(x_6 | x_3, x_5) P(x_7 | x_4, x_5, x_6)$$

Returning now to our first discussion point — remember that DGMs & UGMs are perfect maps for different sets of distributions.

(However, the latter regime is more “powerful” than the other)

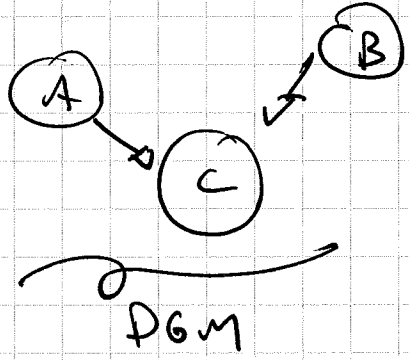
We are now equipped to see this distinction directly.

Q: Can we simply convert, say, a DGM to JGM?

A: Not quite. In practice, one can "moralize" an JGM (by connecting co-parents) but this process unfortunately loses some CI information.

(Also: converting JGM to DOM is problematic because of normalization constraints).

As an example of some CI relationships that can't be modeled perfectly by a JGM - but can be modeled by a DGM consider:



→ implicitly makes (2) CI claims:

(1) $A + B \neq C$

(2) $A \neq B + C$

You are welcome to confirm that no JGM can represent precisely (d only) these (2) CI statements!

Inference in GMs

Most often, we are interested in computing $p(x|\theta)$ for some density of interest. As we show, sampling can furnish us with a mechanism to approximate (or even exactly solve) problems of high complexity (e.g. NP-hard TSP).

In an explicit optimization setting, we can use a sampling-based procedure called simulated annealing.

Remember that GMs provide many useful advantages as a ML tool - (least of not which are their (potential) computational savings).

In the main, if $x_i \in \{1, \dots, K\}$ then the full specification of $p(x)$, for $x \in \mathbb{R}^D$ $x \in \{1, \dots, K\}^D$ is $\mathcal{O}(K^D)$ - which is infeasible in a general setting.

By encoding various CJ conditions, a GM greatly reduces this computation.

In domains for which "causality"/directionality attributions are appropriate (e.g. Medical Diagnosis), we would use a DGM. Here, however, we focus on inference for UGM, with specific applications to Gibbs sampling & image segmentation in Computer vision. (Our OMs are thus more conducive to spatial - as opposed to causal - structures).

Hammersley-Clifford Theorem

Since UGM have no inherent topological ordering, we can't use the "chain rule" to represent $p(x|\theta)$. Instead, we associate potential functions with maximal cliques in the graph; H-M asserts

that $p(x|\theta)$ factors as a product of potentials over maximal cliques. (if $I(\theta) \subseteq I(\phi)$, then we can write $p(x|\theta)$)

H-M:

$$p(x|\theta) = \frac{1}{z(\theta)} \prod_{c \in C} \psi_c(x_c|\theta_c)$$

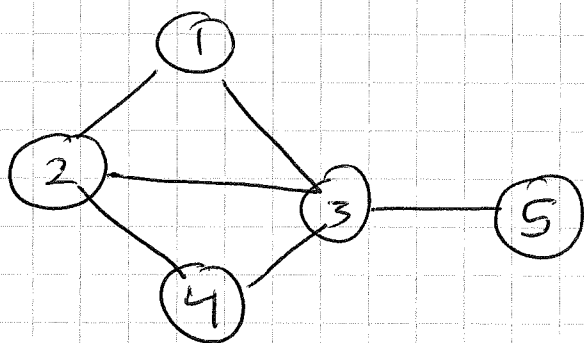
$z(\theta) \triangleq \sum_x \prod_{c \in C} \psi_c(x_c|\theta_c)$ (partition function)

↓ see Koller for proof

FIVE STAR. ★★★★★
 FIVE STAR. ★★★★★
 FIVE STAR. ★★★★★

EX.

G →



$$p(x|\theta) = \frac{1}{Z(\theta)} \psi_{123}(x_1, x_2, x_3) \psi_{234}(x_2, x_3, x_4) \psi_{35}(x_3, x_5)$$

In fact, as a consequence of $[H-C]$, one can show that $p(x|\theta)$ can, equivalently, be parameterized according to the edges (only), instead of maximal cliques. (Called a Pairwise MRF)

Returning to the example above, we have:

$$p(x|\theta) \propto \prod_{s \sim t} \psi_{st}(x_s, x_t)$$

$$= \psi_{12}(x_1, x_2) \psi_{13}(x_1, x_3) \psi_{23}(x_2, x_3)$$

$$\cdot \psi_{24}(x_2, x_4) \psi_{34}(x_3, x_4) \psi_{35}(x_3, x_5)$$

Broadly speaking, inference for GMs is commonly performed using "message-passing" (Pearl) algorithms (e.g. Belief Propagation / sum-product).

Importantly, exact inference ($\mathcal{O}(N(b))$) is possible for tree-structures (we simply propagate "messages" from the root). However, in most practical applications only approximate inference is possible

(using an iterative, "loopy belief" propagation.)

- convergence (or approximate convergence), is often not guaranteed & difficult, moreover, to ascertain.

We content ourselves to consider a classical, Gibbs model, inspired by statistical physics.

Define the potential (over a clique/edge):

$\psi_c(x_c | \theta_c) \stackrel{\Delta}{=} \exp(-E(x_c | \theta_c))$, for some energy function: $E(x_c) > 0$.

The Gibbs distribution is defined:

15

$$p(x|\theta) = \frac{1}{Z(\theta)} \exp\left(-\sum_c E(x_c|\theta_c)\right)$$

This is an "energy-based" model - whereupon high probability states correspond to low energy configurations.

In the binary case (i.e. $x_i \in \{-1, 1\}$) (e.g. Ising Model)

we can write: $\log p(x|\theta) = \delta \sum_{s,t} x_s x_t$

$$p(x|\theta) = \frac{1}{Z(\theta)} \exp\left[\delta \sum_{s,t} x_s x_t\right]$$

Oftentimes, we consider the presence of an external field using a local energy term.

Finally, this gives us:

$$p(x|\theta) = \frac{1}{Z(\theta)} \left[\sum_{s \in S} \lambda_s x_s + \delta \sum_{s,t} x_s x_t \right]$$

General Version of Gibbs Sampling

We want to sample from: $p(x|\theta)$ ($x \in \mathbb{R}^d$)

Basic Idea: Sample, sequentially, from full conditionals.

Gibbs Sampling

(1) Initialize: $\{x_i : i=1, \dots, D\}$

(2) For $t=1, \dots, T$

Sample: $x_1^{(t+1)} \sim p(x_1 | X_{(-1)}^{(t)})$
 "full conditional"

Sample: $x_2^{(t+1)} \sim p(x_2 | X_{(-2)}^{(t)})$

\vdots

Sample: $x_D^{(t+1)} \sim p(x_D | X_{(-D)}^{(t)})$

Complications: (i) how to "order" samples (can be random, but careful)
 (ii) Need full conditionals - can use an approximation or M-H procedure

Under "nice" assumptions (ergodicity), we get $x \sim p(x|\theta)$.

Sampling Algorithms More Generally

(17)

Q: How to perform posterior inference more generally?

Often times we can't simply rely on strong parametric assumptions (e.g. conjugacy/exp. family structures).

Monte Carlo Approximation / inference can get around this.

Basic Idea: ① Draw samples: $x^s \sim p(x|\theta)$

② Compute quantity of interest, e.g.

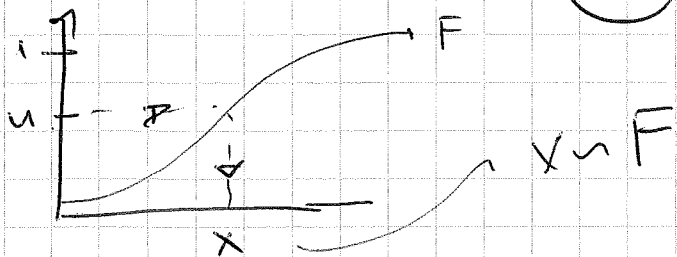
Moments: $E[p(x|\theta)] = \frac{1}{S} \sum_{i=1}^S x_{i,i}$, etc.

In general:

$$E[f(x)] = \int f(x) p(x|\theta) dx = \frac{1}{S} \sum_{i=1}^S f(x_i^s)$$

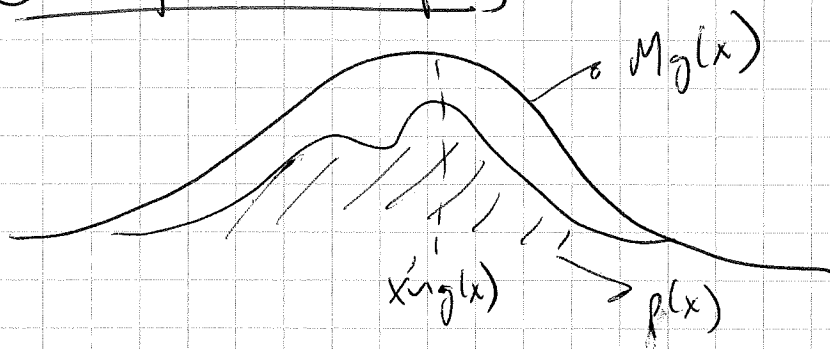
Common MC Sampling Techniques:

- ① CDF method: ① sample $U \sim U(0,1)$
② $F^{-1}(U) \sim F$



② Rejection Sampling

18



Want: $x \sim p(x)$

Idea: ① Sample $x' \sim g(x)$ where $\tilde{p} \leq Mg(x)$

② Sample $U \sim U(0,1)$

③ Reject x' if $u > \frac{\tilde{p}(x)}{Mg(x)}$, else accept.

One can show that $x' \sim p(x)$.

Issues: need "Good" $g(x)$, M value & rejection rate can grow astronomically in high dimensions!

Pros of MC sampling: samples are independent!

Cons of MC sampling: very inefficient in high dimensions!

Alternatively, we can use MCMC Methods

Pros: works in H-D

Cons: dependent samples, convergence can be problematic.

MCMC Sampling Methods

19

Idea: Build a MC (Markov Chain) - Take a random walk (for a sufficiently long time - called the "burn in" phase); Then $x^{(T)} \sim p(x|\theta)$.

The All-Time Classic MCMC Method:

Metropolis-Hastings Algorithm (Los Alamos)

Idea: Want sample: $x \sim p(x|\theta)$, where, possibly, $p(x|\theta)$ is high-dimensional.
(Initialize: x)

① Define proposal distribution: $g(x'|x)$.

The probability we transition: $x \rightarrow x'$.

Usually: $g(x'|x) = N(x'|x, \Sigma)$.

② Sample: $x' \sim g(x'|x)$

③ Compute: (acceptance probability)

$$\alpha = \frac{\tilde{p}(x') g(x|x')}{\tilde{p}(x) g(x'|x)}$$

let $r = \min(1, \alpha)$

④ Sample $u \sim U(0,1)$

⑤

$$x^{(t+1)} = \begin{cases} x' & \text{if } u < r \\ x & \text{if } u \geq r \end{cases}$$

(only need $\tilde{p} = \frac{p}{Z}$)

Remark: One can show that Gibbs sampling is a special case of M-H (where we only accept).

Gibbs Sampling for Binary Segmentation

Define:
$$P(x_T | x_{(-T)}, \theta) = \prod_{S \in N(T)} \psi_{ST}(x_S, x_T)$$

Re "full conditional" J : coupling strength

Define Re edge potentials:
$$\psi(x_S, x_T) = \exp(J x_S x_T)$$

$\forall x_T \in \{-1, 1\}$.

$$P(x_T = \oplus | x_{(-T)}, \theta) = \frac{\prod_{S \in N(T)} \psi_{ST}(x_T = \oplus, x_S)}{\prod_{S \in N(T)} \psi(x_T = \oplus, x_S) + \prod_{S \in N(T)} \psi(x_T = \ominus, x_S)}$$

$$\frac{\exp\left[J \sum_{S \in N(T)} x_S\right]}{\exp\left[J \sum_{S \in N(T)} x_S\right] + \exp\left[-J \sum_{S \in N(T)} x_S\right]}$$

$$= \frac{\exp\left[J \sum_{S \in N(T)} x_S\right]}{\exp\left[J \sum_{S \in N(T)} x_S\right] + \exp\left[-J \sum_{S \in N(T)} x_S\right]}$$

$$= \frac{\exp\left[J \eta_T\right]}{\exp\left[J \eta_T\right] + \exp\left[-J \eta_T\right]} \quad \left(\eta_T = \sum_{S \in N(T)} x_S\right)$$

$$= \frac{\exp\left[J \eta_T\right]}{\exp\left[J \eta_T\right] + \exp\left[-J \eta_T\right]} = \phi\left(\frac{2J \eta_T}{\sqrt{2J \text{mod}}}\right)$$

Remark: It is not difficult to see that:

21

$$\eta_T = x_T (a_T - d_T)$$

a_T : # agreeing neighbors

d_T : # disagreeing neighbors

Consequently, if, say $a_T = d_T$, we get a uniform conditional.

For image segmentation/denoising, we include a local evidence term: ψ_T . (e.g. $\psi_T(x_T) = N(x_T | \mu, \sigma^2)$)

With local evidence added, then, we have:

$$\underbrace{P(x_T = \pm 1 | x_{(-T)}, y, \theta)}_{\substack{\text{local} \\ \text{evidence}}} = \frac{\exp[\eta_T] \psi_T(+1)}{\exp[\eta_T] \psi_T(+1) + \exp[-\eta_T] \psi_T(-1)}$$

$$= \frac{1}{2} \left(1 + \tanh \left(\eta_T - \log \frac{\psi_T(+1)}{\psi_T(-1)} \right) \right)$$

Now the ~~probability~~ probability of x_T entering each state is determined by both compatibility with its neighbors & compatibility with the data (local likelihood term).

(*) For segmentation, we would execute the above Gibbs sampling technique for many iterations.

Lastly, as an addendum: one can also apply the technique of Simulated annealing to the preceding procedure.

Simulated Annealing

Idea: A stochastic algorithm used in optimization of a black-box function.

We assume: $p(x|\theta) \propto \exp\left(-\frac{f(x)}{T}\right)$,

where $f(x)$ is the "energy" of the system & T is the Temperature.

① Given $x^{(t)}$, sample $x' \sim q(x'|x^{(t)})$ ↙ proposal

② Set $\alpha = \exp\left[\frac{f(x) - f(x')}{T}\right]$

③ Accept new state x' with probability: $\min(1, \alpha)$.

(*) In practice we impose a cooling schedule: $T \rightarrow 0$.

(*) NB: This algorithm allows "down-hill" moves, not less frequently as $T \rightarrow 0$. If we cool sufficiently slowly, algorithm finds global optimum.

FIVE STAR. ★★★★★

FIVE STAR. ★★★★★

FIVE STAR. ★★★★★

FIVE STAR. ★★★★★