

PROBABILISTIC GRAPHICAL MODELS

Lecture 1 - BELIEF NETWORKS AND MARKOV

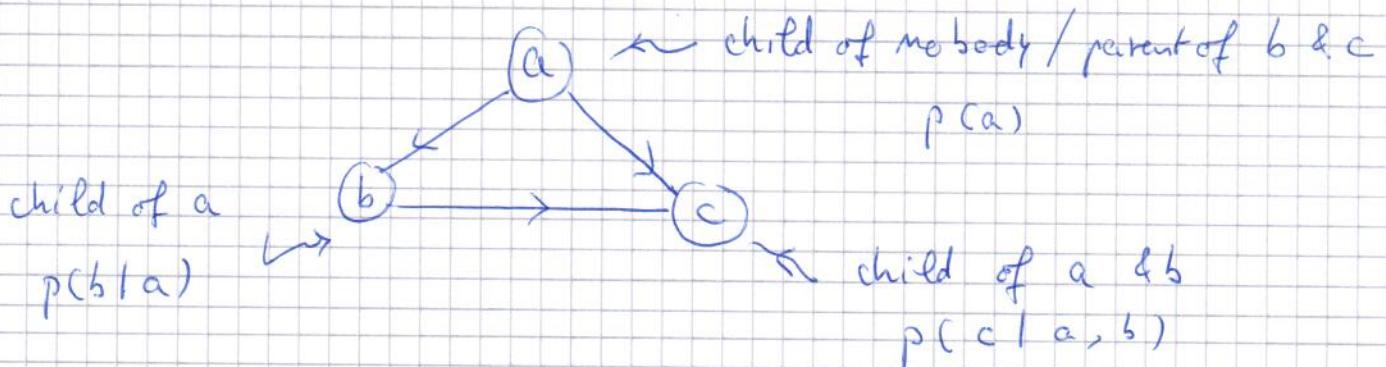
RANDOM FIELDS.

I. BELIEF NETWORKS (also called Bayesian Networks).

Let a, b, c three random variables (r.v) with joint distribution $p(a, b, c)$. By Bayes rule we can write

$$p(a, b, c) = p(c | a, b) p(a, b) = p(c | a, b) p(b | a) p(a)$$

The right hand side can be represented graphically as



This representation is called a Directed Acyclic Graph

DAG : No cycle is found when you follow the arrows.

The representation is highly non-unique in general.

It is always possible to write a joint distri as a DAG in many different ways. The representation becomes useful when the prob distri has a structure and

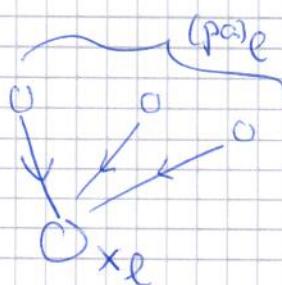
(2)

Some edges are missing corresponding to the fact that some conditional distributions vanish.

Definition: A Belief Network (BN) is a

Directed Acyclic Graph (DAG) which defines a joint p.d.f of the form

$$p(x_1, \dots, x_K) = \prod_{l=1}^K p(x_l | \underbrace{\text{pa}_l}_{\text{parents of } l.})$$



Each conditional distribution $p(x_l | (\text{pa})_l)$ is "attached" to node x_l .

Exercise: Let $x_i \in \mathcal{A}$ some discrete alphabet.

$\underline{x} = (x_1, \dots, x_K) \in \mathcal{A}^K$. Check that for a BN

we automatically have a normalization ;

$$\sum_{\underline{x} \in \mathcal{A}^K} p(\underline{x}) = 1.$$

hint: Start summing later generations (children) and go upwards in generations.

Examples :

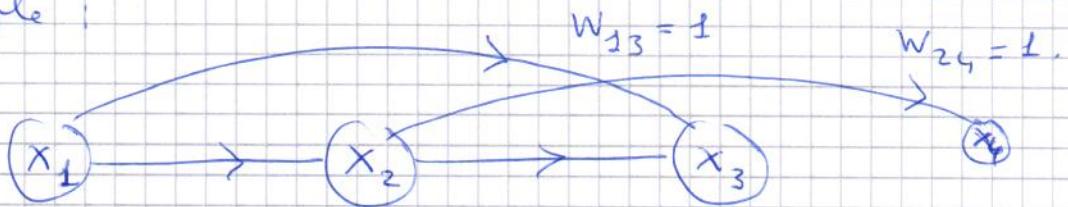
1. Gaussian Belief Network

$$p(\underline{x}) = \prod_{i=1}^K \mathcal{N}(x_i | \sum_{j \in (\text{par}_i)} w_{ij} x_j + b_i; \sigma_i^2)$$

where

$$\mathcal{N}(x | \mu; \sigma^2) = \frac{e^{-\frac{(x-\mu)^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}}$$

w_{ij} is the adjacency matrix of the DAG : For example :



$$w_{12} = 1 \quad w_{23} = 1$$

In the exercises ; compute $E(x_i)$ & $\text{Cov}(x_i, x_j)$ by recursive marginalization. You could then learn w_{ij} , b_i , σ_i^2 by matching these expressions to empirical moments.

2. Markov Chain

Let $\{X^{(t)}, t=0 \dots T\}$ a stochastic process with state space $X^{(t)} \in S$ (the "alphabet").

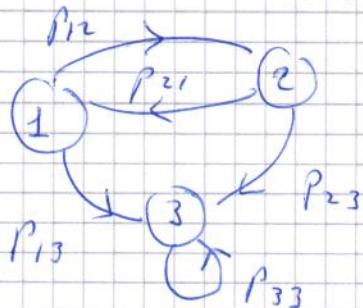
The process is a "Markov chain" if it satisfies

$$p(X^{(t)} | X^{(t-1)} \dots X^{(0)}) = p(X^{(t)} | X^{(t-1)}).$$

Furthermore we say it is "homogeneous" if the transition probability satisfies:

$$p(X^{(t)} = j | X^{(t-1)} = i) = p_{i \rightarrow j} \text{ index of time.}$$

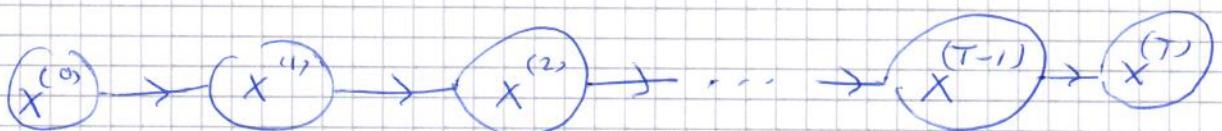
Recall the state graph representation of the Markov chain



← This is NOT a DAG.

Here we illustrate another representation, The DAG

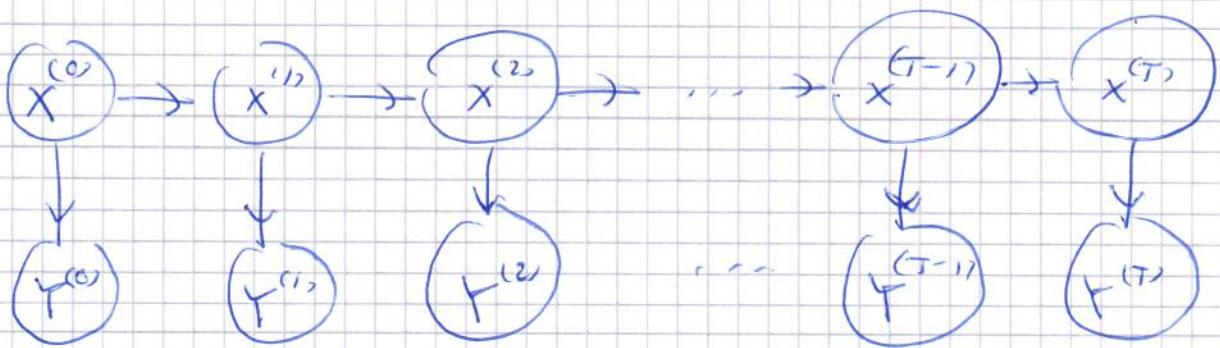
for the joint probability of the MC:



$$p(X^{(0)} \dots X^{(T)}) = p(X^{(0)}) p(X^{(1)} | X^{(0)}) \dots p(X^{(T)} | X^{(T-1)}).$$

(5)

3. Hidden Markov Model.



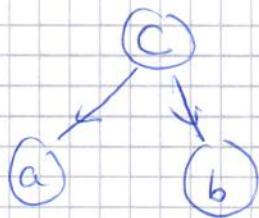
$$\underbrace{P(X^{(0)}, X^{(1)}, \dots, X^{(T)}, Y^{(0)}, \dots, Y^{(T)})}_{\text{full joint prob}} = p(X^{(0)}) p(X^{(1)}|X^{(0)}) \dots p(X^{(T)}|X^{(T-1)}) \\ \cdot p(Y^{(0)}|X^{(0)}) p(Y^{(1)}|X^{(1)}) \dots p(Y^{(T)}|X^{(T)})$$

Here usually this models a situation where a time series $Y^{(0)}, \dots, Y^{(T)}$ is visible (observed) and we want to learn the transition probability matrix P_{ij} of a Hidden Markov Chain $(X^{(0)}, \dots, X^{(T)})$.

(6)

II. CONDITIONAL INDEPENDENCE EXPRESSED BY DAG's.

Example :



$$p(a, b | c) = p(a | c) p(b | c) p(c)$$

Claim : $a \perp\!\!\!\perp b | c$ "a and b independent given c"

Proof : $p(a, b | c) = \frac{p(a, b, c)}{p(c)} = \frac{p(a | c) p(b | c) p(c)}{p(c)}$

$$= p(a | c) p(b | c). \quad \blacksquare$$

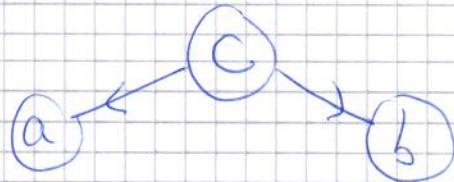
Remark - $p(a, b) \neq p(a) p(b)$ so a and b are not independent. They are independent only when conditioned on c -

Three Main situations that can occur in DAG's.

a) Tail to Tail :

$$p(a, b | c) = p(a | c) p(b | c)$$

$a \perp\!\!\!\perp b | c$. "a & b indep given c".



b) Head to Tail :



$a \perp\!\!\!\perp b | c$ "a & b indep given c".

7

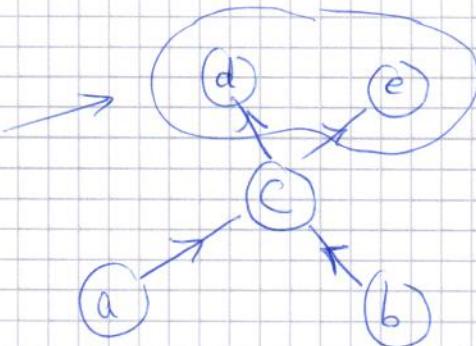
proof : $p(a, b | c) = \frac{p(a, b, c)}{p(c)} = \frac{p(b|c) p(c|a) p(a)}{p(c)}$

 $= p(b|c) \frac{p(c|a)}{p(c)} = p(b|c) p(a|c)$

■

c) Head to Head.

descendents of (c).



Here $a \perp\!\!\!\perp b | c$

\nwarrow a & b are indep if not conditioned
on c neither its descendants

Proof : $p(a, b) = \sum_{c, d, e} p(a, b, c, d, e)$

 $= \sum_{c, d, e} p(c|ab) p(d|c) p(e|c) p(a) \cdot p(b)$
 $= p(a) p(b) \underbrace{\sum_c p(c|ab)}_1 \underbrace{\sum_d p(d|c)}_2 \underbrace{\sum_e p(e|c)}_2$
 $= p(a) \cdot p(b)$

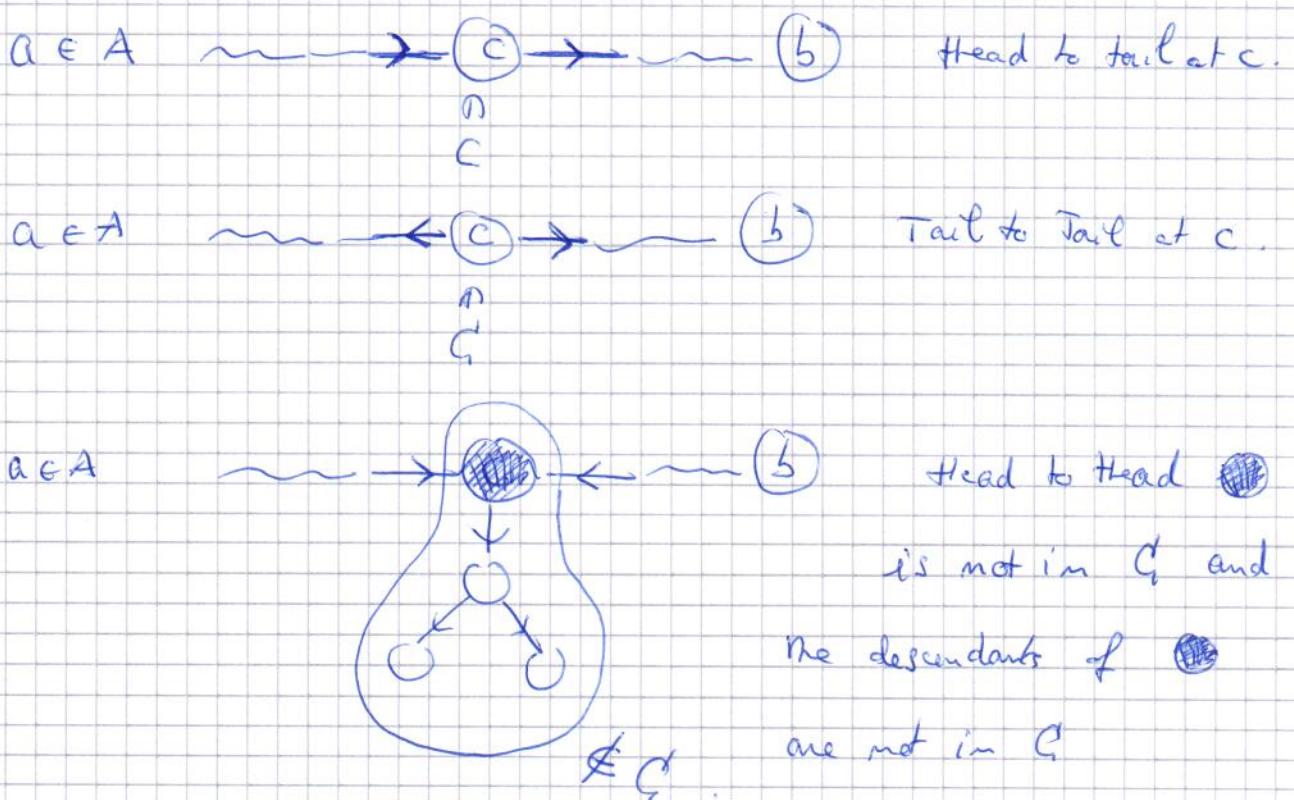
■

General Conditional Independence Principle - Also called
sometimes "d-separation principle".

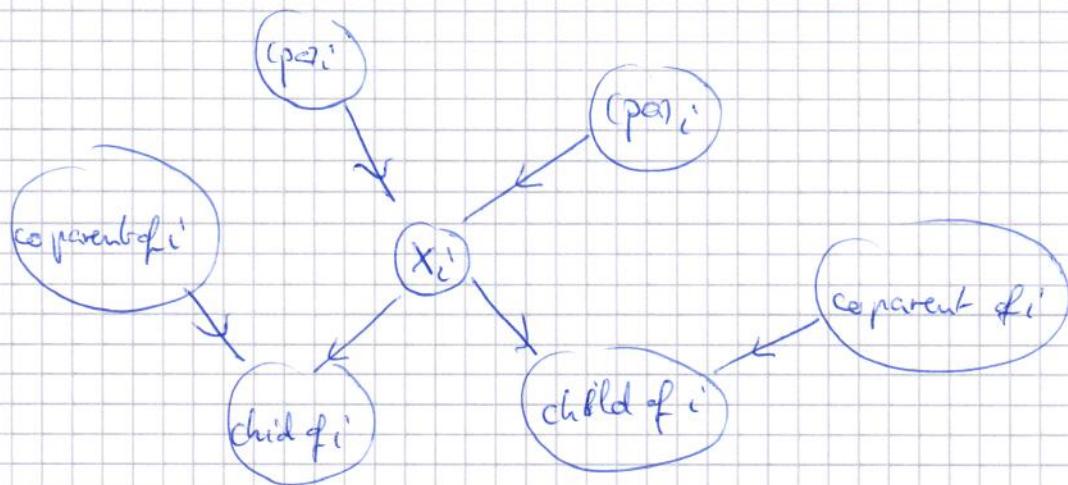
Let A, B, C three disjoint subsets of r.v. If
all undirected paths from A to B are "blocked"
by C , Then $A \perp\!\!\!\perp B | C$

[the converse is not true in general as we can arrange for
numerical coincidences].

Definition of "Blocked": A path from (a) to (b)
is blocked by C if either of the three happens.



Notion of Marker Blanket of a node in a BN.



Marker Blanket of x_i = $MB(i) = \{ (\text{par})_i, (\text{child})_i, \text{copar}(i) \}$.

Lemma : exercise .

$$p(x_i | \{x_j\}_{j \neq i}) = p(x_i | MB(i)) .$$

$$p(x_S | x_{V \setminus S}) = p(x_S | (MB)_{(S)})$$

where S is a set of nodes, $V \setminus S$ the complementary set and $MB(S)$ the Marker Blanket of the set S .

III. MARKOV RANDOM FIELDS (and Gibbs dist.)

Markov random fields are undirected graphical models that again allow to express easily independence conditional statements. These are closely related to Gibbs distributions of statistical physics. They are a sort of generalization of Markov Stochastic Processes.

Consider $G = (V, E)$ undirected with n, v attached to vertices $\{x_1, x_2, \dots, x_N\}$, $n = |V|$ and E a set of edges.

Definition of a MRF:

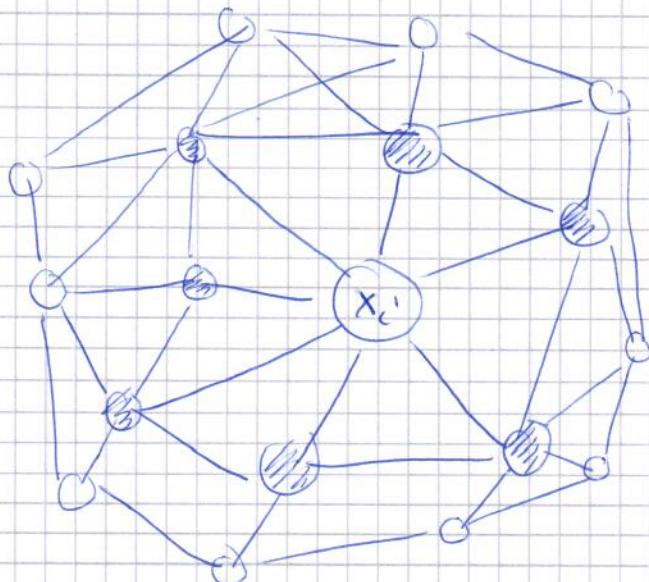
Given a graph $G = (V, E)$ we say that $p(x)$ is a MRF if for any set $S \subset V$ with "external boundary" defined by

$$B_S = \{j \in V \mid (i, j) \in E \text{ with } i \in S, j \in V \setminus S\}$$

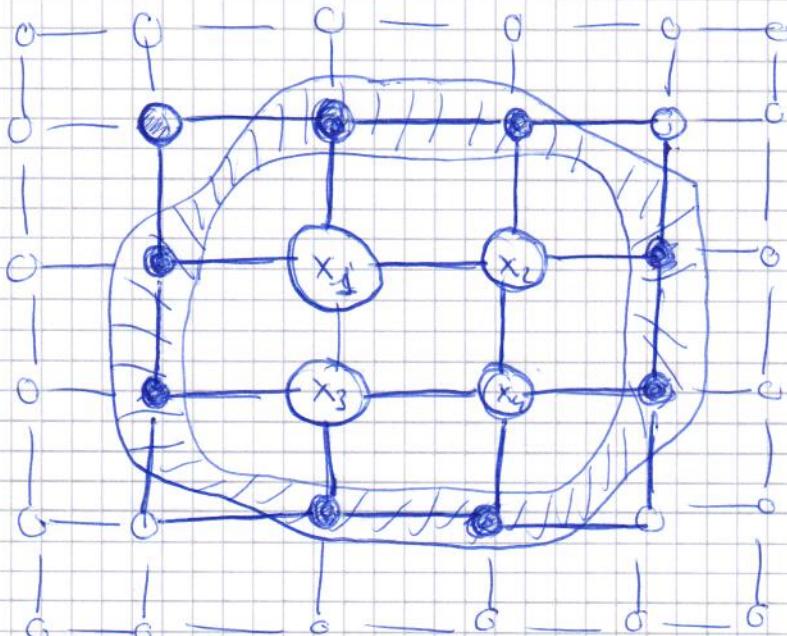
we have

$$P(\{x_i\}_{i \in S} \mid \{x_j\}_{j \in V \setminus S}) = P(\{x_i\}_{i \in S} \mid \{x_j\}_{j \in B_S}).$$

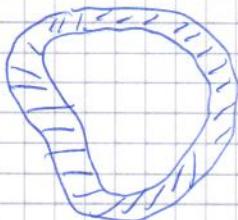
(11.)



- $G = (V, E)$.
 - here $S = i$.
 - \mathcal{B}_S = shaded nodes.
 - \mathcal{B}_S is also called the Markov Blanket of S .
- $MB(S)$.



- $S = \{x_1, x_2, x_3, x_4\}$
- $MB(S) = \mathcal{B}_S$
= ● nodes inside



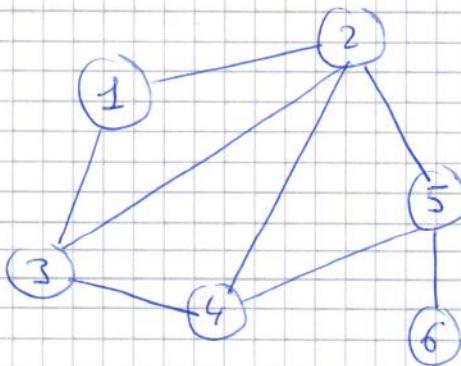
Notation: We adopt the notation $X_S = \{x_i\}_{i \in S}$
for a set of r.v.

Notion of Clique: A clique is a subset $C \subset V$

such that the restriction of $G = (V, E)$ to this subset is a complete graph. All nodes in a clique are pairwise connected.

Maximal Clique: A clique C is maximal if when we add an extra node the set $C \cup \{\text{extra node}\}$ is not a clique.

Example



(1,2), (1,3), (2,3), (2,4), (2,5), (4,5), (5,6) are cliques.

(5,6) is also a maximal clique (the others are not maximal).

(1,2,3) maximal clique

(2,4,5) maximal clique

(2,3,4) maximal clique

Basic Thm for MRF : Hammersley-Clifford ~ 1970's.

Fix $G = (V, E)$, Let $p(\underline{x}) > 0$ (strictly positive otherwise careful!) be a MRF according to the previous definition.

(i) Then $p(\underline{x})$ must be of the form:

$$p(\underline{x}) = \frac{1}{Z} \prod_{\substack{\text{Maximal} \\ \text{cliques} \\ C}} \psi_C(x_C) \quad (\star)$$

where $\psi_C(x_C)$ is a syst of > 0 weights (un-normalized), associated to maximal cliques. Also, Z is called the normalization factor and:

$$Z = \sum_{\underline{x} \in \Omega^V} \prod_{\substack{\text{Max cliques} \\ C}} \psi_C(x_C).$$

(ii) The converse is also true i.e., $(\star) \Rightarrow$ previous conditional independence of MRF.

In other words (\star) implies that for

any $S \subset V$:

$$p(x_S | x_{V \setminus S}) = p(x_S | \underbrace{(MB)(S)}_{\text{or } B_S}).$$

Remarks:

- a) The $\psi_c(x_c)$ are not prob distr. They are positive but not normalized.
- b) The normalization factor Z is also called the partition function.
- c) The sum over $x \in \Omega^N$ to compute Z contains an exponential number of terms and is intractable when the number variables grows large.
- d) Since $\psi_c(x_c) > 0$ we can always write them as
- $$\psi_c(x_c) = \exp(-E_c(x_c))$$
- where $E_c(x_c) \in \mathbb{R}$ are called "energy functions" or "potentials". This representation is used in physics and called the Boltzmann weight repr.
- e) $p(x)$ is called a Gibbs distribution if \exists energy functions $E_S(x_S)$ st S is a clique and
- $$p(x) = \frac{1}{Z} \exp\left\{-\sum_{S \text{ cliques}} E_S(x_S)\right\},$$
- usually the sum does not carry necessarily over maximal cliques. For finite lattices MRF \Leftrightarrow Gibbs. For infinite lattices the theory of Gibbs distr is much deeper and goes back to Dobrushin - Leontov - Ruelle (DLR).

f) MRF are a special case of exponential family.

But note that not all exp families are MRF. To see a MRF as an exponential family note that:

$$p(\underline{x}) = \frac{1}{Z} \prod_{c \in \text{max clique}} \varphi_c(x_c) = \exp\left\{-\sum_{c \in \text{max clique}} \psi_c(x_c) - \log Z\right\}$$

Take $\delta_i \rightarrow \delta_{c, z_c} \equiv -E_c(z_c)$ for $C = \text{max clique}$
 $z_c \in \partial \mathbb{H}^4$.

$$T_i(\underline{x}) \rightarrow T_{c, z_c}(\underline{x}) \equiv \delta_{x_c, z_c}.$$

$$\therefore \rightarrow (c, z_c) \in \{\text{max clique}\} \times \partial \mathbb{H}^4.$$

Then we see that:

$$p(\underline{x}) = h(\underline{x}) \exp\left\{\sum_i \delta_i T_i(\underline{x}) - \psi(\underline{\delta})\right\} \quad (**)$$

with $h(\underline{x}) = 1$ and $\psi(\underline{\delta}) = \log Z$.

{(**)} is the definition of an exponential family with
 prior $h(\underline{x})$, sufficient statistic $T_i(\underline{x})$ and
 parameters δ_i . The $\exp(-\psi(\underline{\delta}))$ is a
 normalization factor.

Proof of (iii) in Hammersley-Clifford Thm.

The proof of (i) is more trivial and here we only prove the converse (ii). (See Grimmett chap 7 for (i).).

Moreover for simplicity we prove the claim for $S = \{i\}$ a singleton only. So we want to show that the factorization property (*) page 13 implies:

$$p(x_i | \{x_j\}_{j \neq i}) = p(x_i | MB(i)).$$

Note that $MB(i) = \text{set of max cliques containing node } i$.

$$\text{l.h.s} = p(x_i | \{x_j\}_{j \neq i}) = \frac{p(x)}{\int dx_i p(x)}$$

$$= \frac{\prod_c \psi_c(x_c)}{\int dx_i \prod_c \psi_c(x_c)}$$

$$\prod_c \psi_c(x_c) = \left\{ \prod_{c: x_i \in c} \psi_c(x_c) \right\} \left\{ \prod_{c: x_i \notin c} \psi_c(x_c) \right\}$$

$$\int dx_i \prod_c \psi_c(x_c) = \left\{ \int dx_i \prod_{c: x_i \in c} \psi_c(x_c) \right\} \left\{ \prod_{c: x_i \notin c} \psi_c(x_c) \right\}$$

$$\Rightarrow \text{l.h.s} = \frac{\prod_{c: x_i \in c} \psi_c(x_c)}{\int dx_i \prod_{c: x_i \in c} \psi_c(x_c)}.$$

(17)

$$n.h.s = p(x_i | MB(i)) = \frac{p(x_i, MB(i))}{\int dx_i p(x_i, MB(i))}$$

$$= \frac{\int dx_{\sim i, MB(i)} p(x)}{\int dx_i \int dx_{\sim i, MB(i)} p(x)}$$

where $\int dx_{\sim i, MB(i)}$ means integral with respect to all variables except $i, MB(i)$.

$$\text{Numerator} = \int dx_{\sim i, MB(i)} p(x) = \int dx_{\sim i, MB(i)} \underbrace{\prod_{c: x_i \in c} \psi_c(x_c) \prod_{c: x_i \notin c} \psi_c(x_c)}$$

contains i & $MB(i)$

only so we can pull it out of integral.

$$= \prod_{c: x_i \in c} \psi_c(x_c) \left\{ \int dx_{\sim i, MB(i)} \prod_{c: x_i \notin c} \psi_c(x_c) \right\}$$

$$\text{Denominator} = \left\{ \int dx_i \prod_{c: x_i \in c} \psi_c(x_c) \right\} \left\{ \int dx_{\sim i, MB(i)} \prod_{c: x_i \notin c} \psi_c(x_c) \right\}$$

$$\Rightarrow n.h.s = \frac{\prod_{c: x_i \in c} \psi_c(x_c)}{\int dx_i \prod_{c: x_i \in c} \psi_c(x_c)}$$

We see that l.h.s = r.h.s

QED

Examples of MRF .

1) Ising Model. (1920's !)

$G = (V, E)$. $s_i, i \in V$ with $s_i \in \{-1, +1\}$
are called "spins".

$$p(\underline{s}) = \frac{1}{Z} \exp \left\{ - \sum_{(i,j) \in E} J_{ij} s_i s_j - \sum_{i \in V} h_i s_i \right\}.$$

$J_{ij}, h_i \in \mathbb{R}$.

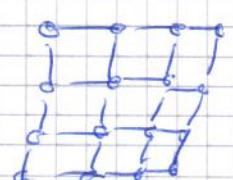
- It can be shown that this is a MRF (on a Gibbs distn).
- Usually in statistical physics J_{ij}, h_i are given and we want to compute marginals of $p(\underline{s})$. In learning theory we are given samples $\underline{s}^{(1)}, \dots, \underline{s}^{(N)} \sim p(\underline{s})$ and we want to learn J_{ij}, h_i as well as G . This is also called the "Inverse Ising Problem".

- $G = \text{complete graph}$, $J_{ij} = \frac{\beta}{|V|}$, $\beta > 0 \rightarrow$ Curie-Weiss model

$\exists \underbrace{\text{One clique}}_{\max}!$

- $G = \text{Cube} \subset \mathbb{Z}^d$, $J_{ij} = \begin{cases} \beta & \text{for } |i-j|=1 \\ 0 & \text{otherwise} \end{cases}$

\rightarrow Traditional Ising model on square grid.



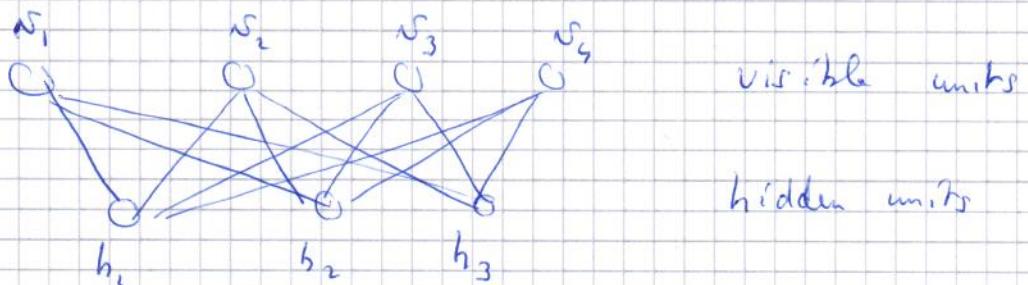
\rightarrow Application in Image Processing for example.
magnetism etc--.

2) Restricted Boltzmann Machine (RBMs)

In the learning theory community the Ising Model is also (obviously) called a "Boltzmann machine". There also exist a variant that is popular with many applications (e.g. as a generative model for pdf's) called the Restricted Boltzmann Machine (RBMs).

$$G = \left(\underbrace{\text{Visible} \cup \text{Hidden}}_{\text{vertices}}, \underbrace{E}_{\substack{\text{edges between} \\ \text{visible \& hidden only}}} \right)$$

a "complete bipartite" graph:



$$P(\underline{v}, \underline{h}) = \frac{1}{Z} \exp \left\{ - \sum_{\substack{i \in \text{Vis} \\ j \in \text{Hid}}} w_{ij} v_i h_j - \sum_{i \in \text{Vis}} b_i v_i - \sum_{j \in \text{Hid}} c_j h_j \right\}$$

Here $\underbrace{\text{cliques}}_{\max}$ are pairs (i, j) .

$$P_{\text{target}}(\underline{v}) = \sum_{\underline{h} \in \text{Hidden}} P(\underline{v}, \underline{h}), \text{ learn } w_{ij}, b_i, c_j \text{ from }$$

samples $(\underline{v}^{(1)}, \dots, \underline{v}^{(n)})$ of visible variables $\sim P_{\text{target}}(\underline{v})$

N Conversion of BN into MRF.

This conversion is useful often. Here we illustrate with two examples and give the general algorithm later on.

a) Simple example of a chain graph (e.g a Markov Process).



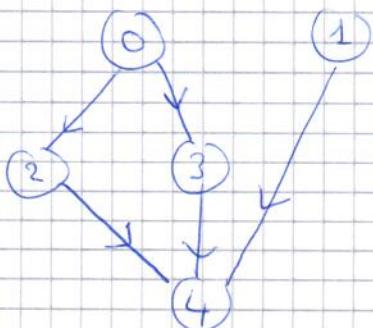
$$\begin{aligned}
 p(\underline{x}) &= p(x_1) \underbrace{p(x_2 | x_1)}_{\varphi_{12}(x_1, x_2)} \cdots \underbrace{p(x_m | x_{m-1})}_{\varphi_{23}(x_2, x_3)} \\
 &= \frac{1}{Z} \underbrace{\varphi_{12}(x_1, x_2)}_{\varphi_{23}(x_2, x_3)} \cdots \underbrace{\varphi(x_m, x_{m-1})}_{\varphi_{m-1,m}(x_{m-1}, x_m)}
 \end{aligned}$$

with $Z = 1$ here (recall BN are automatically normalized!).

Cliques are $(12), (23), \dots, (m-1, m)$. The MRF is the same without arrows (it is undirected).



b) Generic example



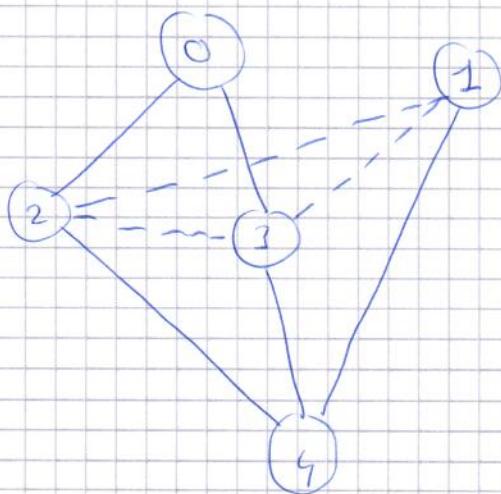
$$\begin{aligned}
 p(\underline{x}) &= p(x_0) p(x_1) p(x_2 | x_0) p(x_3 | x_0) \\
 &\quad \cdot p(x_4 | x_2 x_3 x_1).
 \end{aligned}$$

(21).

$$P(x) = \underbrace{p(x_0) p(x_1 | x_0) p(x_2 | x_0)}_{\psi_{023}(x_0 x_1 x_2)} \underbrace{p(x_3 | x_0 x_1 x_2) p(x_4 | x_1 x_2 x_3 x_4)}_{\psi_{1234}(x_1 x_2 x_3 x_4)}$$

with $Z = 1$.

MRF graph is :



- added edges between parents of child
- This process of "adding edges between parents of a child" is called the "moralization step".
- The MRF graph obtained from the BN DAG is called the "Moralization graph".

c) Conversion Algorithm $\text{BN} \rightarrow \text{MRF}$

- Take a DAG.
- For each Node connect all its parents and remove arrows. We get the Moralization graph.
- Take each Max Clique in Moralization graph and define $\psi_c(x_c)$ by the product of cond prob dist appearing in DAG of each node in the clique. (See example b) above)