So far we've discussed 1) pdfs $p(x|y,z)$ (but not multivariate,
e.g. $p(\underline{x}|Y,Z)$

2) directed graphic models (Bayes nets)

Are there other ways to represent complex distributions
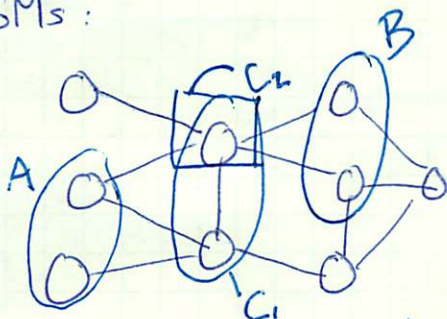(which is to say <u>knowledge</u>)

Bayes Nets (or DAGs) are one way to factor
(and therefore simplify a complex joint prob distr.

$$P(x_{1:N}) = \prod_{i=1}^{N} P(x_i | pa(x_i))$$

$$p(x_i | pa(x_i)) = p(x_i) \quad \text{if } pa(x_i) = \emptyset.$$

Are there other ways? Yes.

Undirected GMs:



nodes still represent
variables, but now there
is no causality implied,
i.e. no arrows.

what are the independence properties of this graph like?

$$A \perp\!\!\!\perp B \,|\, C \,?$$

Much simpler than in DAGs: If all paths
from A to B pass through C, the C "blocks"
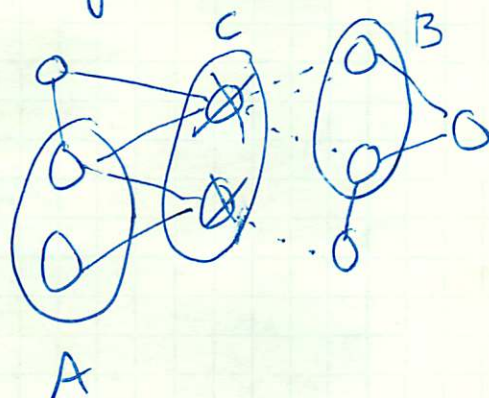A & B and $A \perp\!\!\!\perp B | C_1$.

But if there is a path that isn't blocked the
independ. cond. no longer holds $A \not\!\perp\!\!\!\perp B | C_2$

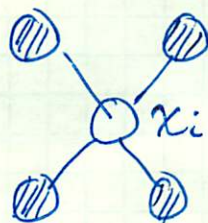There is no "explaining away" phenom.

Another way to think about it.



If we remove nodes in C from graph, and there is no path from A to B then
$$A \perp\!\!\!\perp B \mid C.$$

- Markov blanket is simpler too.



~~X's only depends on its~~
~~neighboring~~

- $X_i$ is cond. indep of all other nodes given only its neighbors.

- What factorization model do we use?

  What does a ~~direct~~ link mean? (or absense of one?)

  

  If ~~that~~ the is no connection between nodes $X_i$ & $X_j$ then ~~then~~ $X_i \perp\!\!\!\perp X_j \mid$ rest of graph.

  because all other paths are blocked.

$$P(X_i, X_j \mid X_{k \neq i,j})$$
$$\text{or } X_{k \setminus \{i,j\}}$$
$$= P(X_i \mid X_{k \neq i,j}) P(X_j \mid X_{k \neq i,j})$$
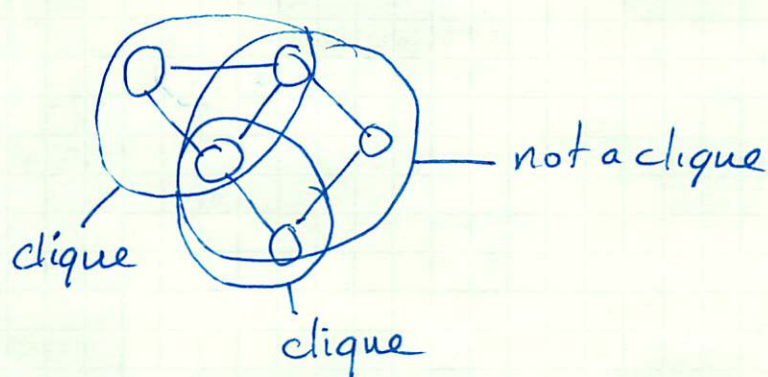
Before we define joint prob. we need to introduce concept of cliques
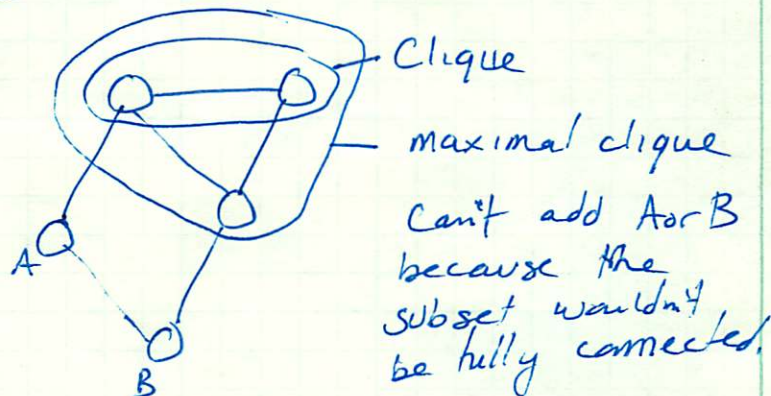
<u>Cliques</u> of a graph :

    a subset of nodes s.t. ∃ a link between all
    pairs of nodes in the subset.  ※ The set of
                                      nodes in a
                                      clique is

                                      <u>fully connected.</u>

                     — not a clique

clique

            clique

<u>Maximal cliques :</u>

    A clique where it is not possible to add
    any more nodes without breaking the clique,
    i.e. by including any more nodes the subset will
    no longer be fully connected and cease to be
    a clique :
                                  — Clique

                              — maximal clique

                                 Can't add A or B
              A                        because the
                                  subset wouldn't
                                  be fully connected.
               B

Define joint distribution to be functions of
maximal cliques — using any subset
of a maximal clique would be redundant.

Let C ~~blahblah~~ denote a clique

$$x_c = \text{set of vars in } C.$$

$$\rightarrow = \underline{x}$$

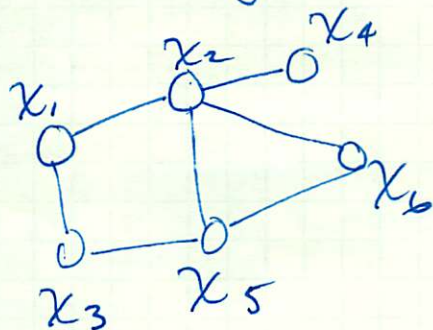$$P(x_{1:N}) = \frac{1}{Z} \prod_C \psi_c(\underline{x}_c) \quad , \quad \psi_c(\underline{x}_c) \geq 0$$

This factorizes the joint prob. distri into small joint pdfs defined on cliques.

$Z$ = normalization const to make $P(x_{1:N})$ a valid pdf. Also called <u>partition function</u>.

$$Z = \sum_{\underline{x}} \prod_C \psi_c(\underline{x}_c)$$

( assume discrete vars, but could use continuous vars and integration

$$Z = \int \prod_C \psi_c(\underline{x}_c) \, d\underline{x}_c$$



$$P(\underline{x}) = \frac{1}{Z} \psi(x_1, x_2) \psi(x_1, x_3) \cancel{\psi\psi\psi\psi\psi} \psi(x_{\frac{2}{2}}, x_{\frac{4}{4}})$$

$$\cdot \psi(x_{\frac{3}{3}}, x_{\frac{5}{5}}) \cdot \psi(x_2, x_5, x_6)$$

This is also valid, don't have to restrict pdf to maximal cliques:

$$P(\underline{x}) = \frac{1}{Z} \psi(x_1, x_2) \psi_{13} \cdot \psi_{25} \psi_{24} \underset{25}{\cancel{\psi_{\text{BA}}}} \psi_{26} \psi_{56} \qquad \text{Why?}$$

Equivalent to assuming $\psi_{256} = \psi_{25} \psi_{26} \psi_{56}$, i.e. it factorizes.

What are "potential functions"? Are they pdfs?

No (or not necessarily)

Only need $\Psi(\underline{x}) \geq 0$ so that $p(\underline{x}) \geq 0$

$\Psi(x)$ does not have a specific interpretation like conditional or marginal, unlike DAGs.

| Problem: | does not produce ~~a~~ a normalized joint pdf.

Need $Z$, which is often hard to compute.

⤷ can be a major limitation of UDGs

How bad is it? How big is the sum?

$$Z = \sum_{\underline{x}} \prod_{c} \Psi_c(\underline{x}_c)$$

⤷ this sums over all values of $\underline{x}$.

- Typically M nodes, each with K states
  $\Rightarrow$ $K^M$ different states!

- Need $Z$ for parameter learning:
  $$\frac{\partial \Psi_c(x_c | \theta_c) \frac{1}{Z}}{\partial \theta_c}$$

- But not for local cond. probs:
  $$P(\underline{x}_a | \underline{x}_b) = \frac{\frac{1}{Z} P(\underline{x}_a, \underline{x}_b)}{\frac{1}{Z} P(\underline{x}_b)} \qquad Z\text{'s cancel}$$

- There are some techniques which we will come across later.

Hammersley-Clifford theorem : Relating factorization to cond. indep.

#1) Undirected GM $G$ is an MRF if two nodes are conditionally indep whenever they are separated by evidence nodes.

$$P(X_i | X_{G \setminus i}) = P(X_i | X_{N_i})$$

$G \setminus i$ = all nodes in graph $G$ except $i$

$N_i$ = neighboring nodes of $X_i$

#2) a pdf $P(X)$ on an undirected GM is is a <u>Gibbs distribution</u> if it can be factored into positive functions defined on cliques that cover all nodes and edges of $G$.

$$P(X) = \frac{1}{Z} \prod_{c \in C_G} \psi_c(X_c)$$

$C_G$ = all (maximal) cliques

H-C theorem says ~~~~
$$\text{Def \#1} \Longleftrightarrow \text{Def \#2}$$

$\psi_c(\underline{x}_c) > 0$  so  it's  convenient  to  write

$$\psi_c(\underline{x}_c) = \exp\left[-E(\underline{x}_c)\right]$$

$$\boxed{E(\underline{x}_c) = \text{"energy function"}}$$

as  in  "energy-based" models

Exponential representation =

$$\boxed{\exp\left[-E(\underline{x}_c)\right] = \text{"Boltzman distribution"}}$$

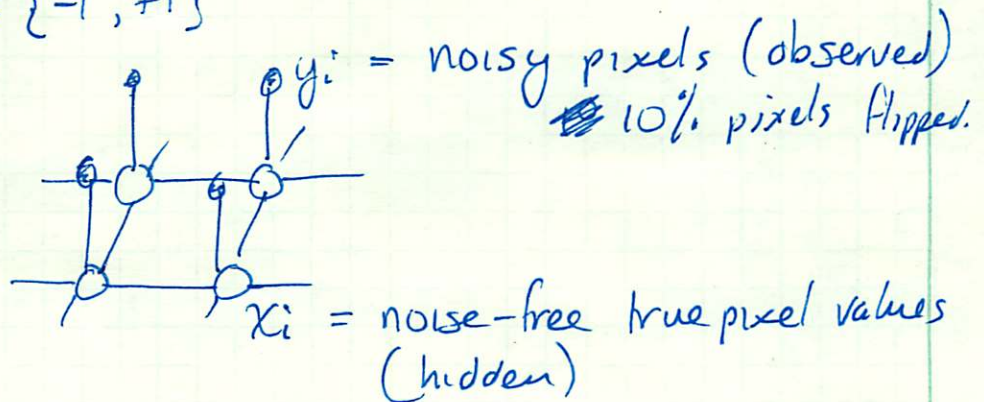Total energy is the sum of energies in each clique.

- potentials (and energy functions) do not have a specific probabilistic interpretation.
  So, what are they? How do we choose them?

- Potential fns express "good" configurations of local vars.

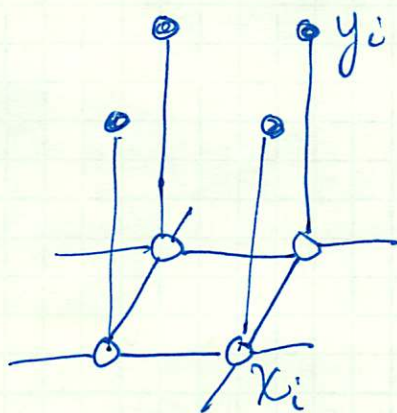Example: image denoising with binary pixels
$$y_i \in \{-1, +1\}$$

$y_i$ = noisy pixels (observed)
10% pixels flipped.

$x_i$ = noise-free true pixel values
(hidden)

How do we express knowledge in the network?
Eg. we know pixels are correlated.
→ links specify degree of correlation.

How do ~~uses~~ we set the energy functions?



$x_i$ & $y_i$ should be correlated

~~can use~~ can use $- \eta \, x_i y_i$

same $\Rightarrow$ low energy (good)

diff $\Rightarrow$ high energy

$\eta = const > 0$

The rest of the graph:

$-\beta x_i x_j$     same $\Rightarrow$ ~~high~~ lower energy

$\beta = 0 \Rightarrow$ no links     diff $\Rightarrow$ higher energy     $\beta = const > 0$

Can also model tendency for pixels to be on or off:

$$h \, x_i \qquad h = const > 0 \qquad \text{"bias"} \qquad h = 0 \Rightarrow \begin{matrix} equal \\ prob \end{matrix} +1, -1$$

Still valid because only cond. is for energy fn to be $> 0$.

$$E(\underline{x}, y) = h \sum_i x_i - \beta \sum_{i,j} x_i x_j - \cancel{\sharp} \eta \sum_i x_i y_i$$

Just add energies of all cliques $\Rightarrow$

$$p(\underline{x}, y) = \frac{1}{Z} \exp\left[ - E(\underline{x}, y) \right]$$

We are given observed pixels, so we have

$$p(\underline{x} \mid y) \quad \text{defined implicitly}$$

How do we obtain $x_i$?

Start at some sol'n, ~~to~~ change ~~if~~ $x_i$ to maximize energy.

iterated conditional modes (ICM)

coord-wise gradient ascent.

1) $x_i = y_i \; \forall i$

2) for each $x_i$, (or at random)

calc $E \mid x_i = +1$    everything else is
$E \mid x_i = -1$    held fixed.

3) change $x_i$ to state with lower energy.

can do efficiently because only
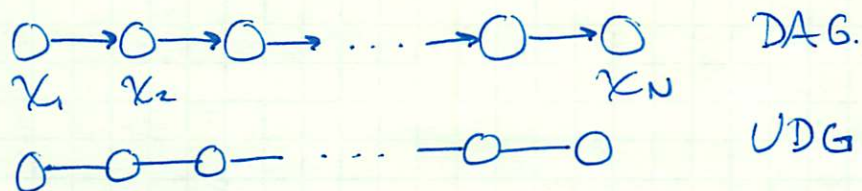one term in $E(\underline{x})$ changes.

4) repeat until stable (or run "long enough")

will converge to local maximum of $E(\underline{x})$ (not global)

[see slides for example]

How do these related to directed graphs?
Can we convert one to another?
Are there adv./disadv. to each?

Consider:

$$O \longrightarrow O \longrightarrow O \longrightarrow \cdots \longrightarrow O \longrightarrow O \qquad \text{DAG.}$$
$$X_1 \quad X_2 \qquad\qquad\qquad X_N$$

$$O\!\!-\!\!O\!\!-\!\!O\!\!-\!\! \cdots -O-O \qquad \text{UDG}$$

For DAG   $p(x) = \cancel{p(x_1)p(x_2)p(x_3)p(x_N)}$

$$p(x_1)\, p(x_2|x_1)\, p(x_3|x_2) \cdots p(x_N|x_{N-1})$$

For UDG:

$$p(x) = \frac{1}{z}\, \psi_{1,2}(x_1, x_2)\, \psi_{2,3}(x_2, x_3) \cdots \psi_{N-1,N}(\,)$$

$\Rightarrow$ can set  $\psi_{1,2}(x_1, x_2) = p(x_1)p(x_2|x_1)$

etc.
$\Rightarrow z = 1$   since it's already normalized.
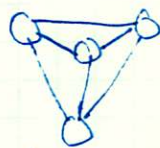
Can we do this in general?



$$p(x) = p(x_1)p(x_2)p(x_3)\, p(x_4|x_{1:3})$$

$X_1 \quad X_3$
$X_2$
$X_4$

Can we use the same trick?

$$p_{UG}(x) = p(x_1)p(x_2)p(x_3)\, \underline{p(x_4|x_1, x_2, x_3)}$$

but this is
fully connected $\Rightarrow$
no cond indep properties, no adv.

This involves all 4.
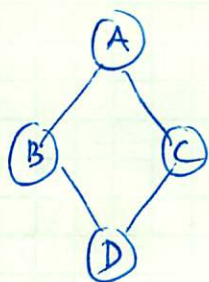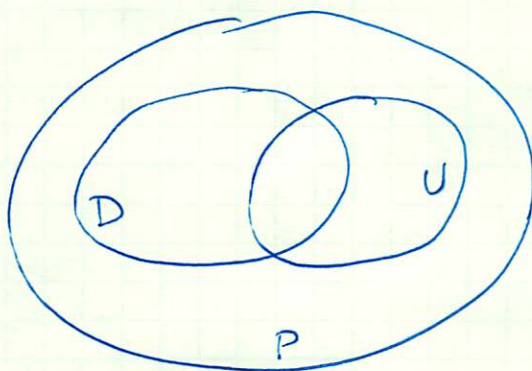$\Rightarrow$ ~~they are~~ all belong
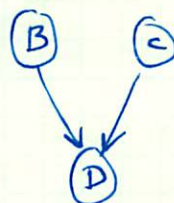to a single clique.

closuring

UGs vs DGs



No DG can represent only

$A \perp D \mid B,C$

$B \perp C \mid A,D$

No UG can rep. only

$B \perp C \mid \emptyset$