

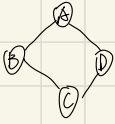


Markov Network

teach made a misconception that may propagate in class

A and C never speak to each other, B and D never speak to each other.

A → C
A ← B
B → C
C → D
C ← D



instead of CPD, the prob of one variable given another.

we need a more symmetric parameterization

e.g. affinity, how likely 2 people are to agree with each other

factor is a function $\phi: \text{val}(D) \rightarrow \mathbb{R}$ where D is the domain of ϕ

$\phi_1(A, B)$		
a^0	b^0	30
a^0	b^1	5
a^1	b^0	1
a^1	b^1	10

$\phi_2(B, C)$		
b^0	c^0	100
b^0	c^1	1
b^1	c^0	1
b^1	c^1	10

$\phi_3(C, D)$		
c^0	d^0	1
c^0	d^1	100
c^1	d^0	100
c^1	d^1	1

$\phi_4(D, A)$		
d^0	a^0	100
d^0	a^1	1
d^1	a^0	1
d^1	a^1	100

Assignment	Unnormalized	Normalized
$a^0 b^0 c^0 d^0$	300,000	0.04
$a^0 b^0 c^1 d^0$	300,000	0.04
$a^0 b^0 c^0 d^1$	30	$4.1 \cdot 10^{-6}$
$a^0 b^0 c^1 d^1$	500	$6.9 \cdot 10^{-5}$
$a^0 b^1 c^0 d^0$	500	$6.9 \cdot 10^{-5}$
$a^0 b^1 c^1 d^0$	5,000,000	0.69
$a^0 b^1 c^0 d^1$	500	$6.9 \cdot 10^{-5}$
$a^0 b^1 c^1 d^1$	100	$1.4 \cdot 10^{-5}$
$a^1 b^0 c^0 d^0$	1,000,000	0.14
$a^1 b^0 c^1 d^0$	100	$1.4 \cdot 10^{-5}$
$a^1 b^0 c^0 d^1$	100	$1.4 \cdot 10^{-5}$
$a^1 b^0 c^1 d^1$	10	$1.4 \cdot 10^{-6}$
$a^1 b^1 c^0 d^0$	100,000	0.014
$a^1 b^1 c^1 d^0$	100,000	0.014
$a^1 b^1 c^0 d^1$	100,000	0.014

local interactions between directly related variables

$$p(a, b, c, d) = \frac{1}{Z} \phi_1(a, b) \cdot \phi_2(b, c) \cdot \phi_3(c, d) \cdot \phi_4(d, a)$$

$Z = \sum_{a,b,c,d} \phi_1(a, b) \cdot \phi_2(b, c) \cdot \phi_3(c, d) \cdot \phi_4(d, a)$ is the partition function

$$(X \perp Y | Z) \iff p(X, Y, Z) = p(X, Z) \cdot p(Y, Z)$$

$$\text{right to left: } p(X, Z) = \sum_Y p(X, Y, Z) = p(X, Z) \cdot \sum_Y p(Y, Z)$$

$$p(Y, Z) = \sum_X p(X, Y, Z) = p_Y(Y, Z) \cdot \sum_X p(X, Z)$$

$$p(Z) = \sum_a \sum_y p(a, y, z) = \sum_a \sum_y \phi_1(a, z) \cdot \phi_2(y, z)$$

$$p(X|Z) = \frac{p(X, Z) \cdot p(Y, Z)}{p(Z)} = \frac{\phi_1(X, Z) \cdot \phi_2(Y, Z) \cdot (\sum_a \phi_1(a, Z)) \cdot (\sum_y \phi_2(y, Z))}{p(Z)}$$

$$= \frac{\phi_1(X, Z) \cdot \phi_2(Y, Z)}{p(Z)} = \frac{p(X, Y, Z)}{p(Z)} = p(X, Y | Z)$$

$$\therefore p(X|Z) p(Y|Z) = p(X, Y | Z)$$

$$\therefore (X \perp Y | Z)$$

$$\begin{aligned} \text{left to right: } p(X, Y, Z) &= p(Z) \cdot p(X|Z) \cdot p(Y|Z) \\ &= \underbrace{p(Z) \cdot p(X|Z)}_{\phi_1(X, Z)} \cdot \underbrace{p(Y|Z)}_{\phi_2(Y, Z)} \end{aligned}$$

$$(X \perp Y | Z) \text{ iff } p(\dots) = \phi_1(X, Z) \cdot \phi_2(Y, Z)$$

$$p(A, B, C, D) = \left[\frac{1}{Z} \phi_1(A, B) \cdot \phi_2(B, C) \right] \phi_3(C, D) \cdot \phi_4(A, D)$$

$$= f_1(B, A, C) \cdot f_2(D, A, C)$$

$$\therefore (B \perp D | A, C)$$

$$\text{similarly } (A \perp C | B, D)$$

Gibbs Distribution and Markov Network

A distribution P_{θ} is a Gibbs distribution parameterized by a set of factors

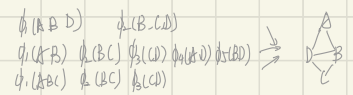
$$\theta = \{\phi_1(D_1), \dots, \phi_k(D_k)\}$$

$$P_{\theta}(x_1, \dots, x_n) = \frac{1}{Z} P_{\theta}(x_1, \dots, x_n)$$

$$= \frac{1}{\sum_{x_1, \dots, x_n} \phi_1(D_1) \times \dots \times \phi_k(D_k)} \phi_1(D_1) \times \dots \times \phi_k(D_k)$$

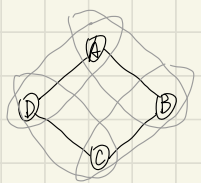
A factor is not a joint distribution, it's only one contribution to the overall joint distribution
The distribution as a whole has to take into consideration contributions from all of the factors involved

A distribution P_{θ} with $\theta = \{\phi_1(D_1) \dots \phi_k(D_k)\}$ factorizes over a Markov network \mathcal{H} if each D_k is a complete subgraph of \mathcal{H}



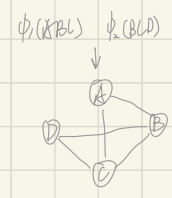
Cannot read the factorisation from graph

the factors that parameterize a Markov network are often called clique potentials



4 cliques and the clique potentials for each of them

$\{A, B\}$	$\begin{array}{c c c} A & B & \phi(A,B) \\ \hline a^1 & b^1 & 0.5 \\ a^1 & b^2 & 0.8 \\ a^2 & b^1 & 0.1 \\ a^2 & b^2 & 0 \end{array}$	$\{C, D\}$	$\begin{array}{c c c} C & D & \phi(C,D) \\ \hline c^1 & d^1 & 0.5 \\ c^1 & d^2 & 0.7 \\ c^2 & d^1 & 0.1 \\ c^2 & d^2 & 0 \end{array}$
$\{B, C\}$	$\begin{array}{c c c} B & C & \phi(B,C) \\ \hline b^1 & c^1 & 0.5 \\ b^1 & c^2 & 0.7 \\ b^2 & c^1 & 0.1 \\ b^2 & c^2 & 0 \end{array}$	$\{A, D\}$	$\begin{array}{c c c} A & D & \phi(A,D) \\ \hline a^1 & d^1 & 0.5 \\ a^1 & d^2 & 0.7 \\ a^2 & d^1 & 0.1 \\ a^2 & d^2 & 0 \end{array}$



Reduced Markov Networks

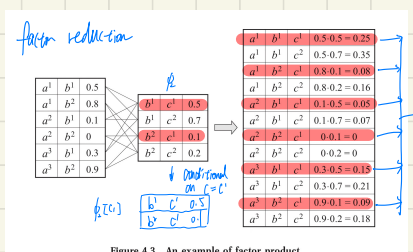


Figure 4.3 An example of factor product

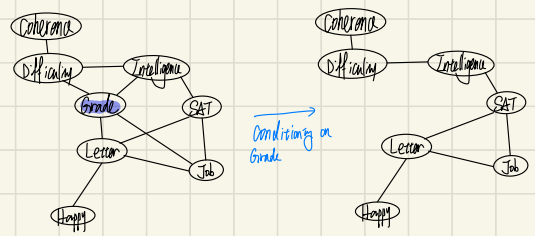
\mathcal{H} is a Markov network over X and $U = u$ is a context.

the reduced Markov network $\mathcal{H}(u)$ is a Markov network over nodes $X \setminus U$

(edges not associated with U stay the same)

In a Markov network, conditioning on a variable eliminates the node and edges.

In a Bayesian network, conditioning on a variable can activate v-structures, creating new dependencies



Martov Independencies

In Markov networks, probabilistic influence flows along undirected paths in the graph, but is blocked if we condition on the intervening nodes

A set of nodes Z separates X and Y in \mathcal{H} , denoted $\text{Sep}_{\mathcal{H}}(X; Y | Z)$

if there is no active path between nodes $X \in X$ and $Y \in Y$ given Z

global independencies in \mathcal{H} $I(\mathcal{H}) = \{ (X \perp Y | Z) : \text{Sep}_{\mathcal{H}}(X; Y | Z) \}$

P is a Gibbs distribution that factorizes over $\mathcal{H} \implies \mathcal{H}$ is an I-map for P

let X, Y, Z be any 3 disjoint subsets of all variables s.t. $\text{Sep}_{\mathcal{H}}(X; Y | Z)$

Z_X is the set of cliques that are contained in $X \cup Z$

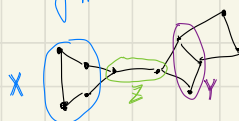
Z_Y is the set of cliques that are contained in remaining nodes

$$P(x_1, \dots, x_n) = \frac{1}{Z} \prod_{C \in Z_X} \psi_C(x_C) \cdot \prod_{C \in Z_Y} \psi_C(x_C)$$

$$= \frac{1}{Z} f(x; Z) \cdot g(Y; Z)$$

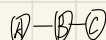
$$\therefore (X \perp Y | Z)$$

any conditional independence revealed in \mathcal{H} holds true in P



positive P factorizes over \mathcal{H}

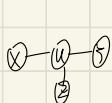
\mathcal{H} is an I-map of P



If \mathcal{H} is an I-map for positive P , then P is a Gibbs distribution that factorizes over \mathcal{H}

If X and Y are not separated given Z in \mathcal{H} ,

then X and Y are dependent given Z in some distribution that factorizes over \mathcal{H}



$$\phi_1(x, u, y) = \begin{array}{c|cc} & x & u & y \\ \hline 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \end{array}$$

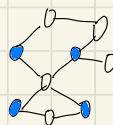
$$\phi_2(u, y) = \begin{array}{c|cc} & u & y & \phi \\ \hline 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \end{array}$$

pairwise independencies: $I_p(\mathcal{H}) = \{ (X \perp Y | X - \{X, Y\}) : X, Y \text{ are not directly connected} \}$

For a graph \mathcal{H} , the Markov blanket of X in \mathcal{H} $MB_{\mathcal{H}}(X)$ is the neighbors of X

local independencies: $I_L(\mathcal{H}) = \{ (X \perp X - \{X\} - MB_{\mathcal{H}}(X) | MB_{\mathcal{H}}(X)) : X \in \mathcal{X} \}$

a variable is independent of the rest given its immediate neighbors



given variables block the flow of influence

$$P \models I_L(\mathcal{H}) \iff P \models I_p(\mathcal{H}) \iff P \models I(\mathcal{H}) \quad \text{for a positive distribution}$$

From distribution to Graphs

2 approaches to encode independencies of a distribution with a graph structure

1° pairwise independencies

add edges between all pairs of X, Y such that $P \not\models (X \perp Y \mid \mathcal{C} - \{X, Y\})$

2° local independencies

A set U is a Markov Blanket of X in a distribution P if

- $X \notin U$
- U is a minimal set of nodes such that $(X \perp \mathcal{C} - \{X\} - U \mid U) \in \text{IP}(P)$

define a graph \mathcal{H} by introducing an edge $\{X, Y\}$ for all X and all $Y \in \text{MB}_P(X)$

both of these 2 approaches produce the unique minimal I-map of P for a positive distribution

ex. a non positive distribution P over 4 binary variables A, B, C, D ,
 $P(A, B, C, D) \neq 0$ only when $A=B=C=D$

	A	B	C	D
0°	0°	0°	0°	0°
1°	0°	1°	0°	0°
2°	0°	0°	1°	0°
3°	0°	0°	0°	1°

1° local independency construction

$P \models (A \perp C, D \mid B) \rightarrow$ edge $A-B$

$\therefore B$ is a valid choice for $\text{MB}_P(A)$

$P \models (C \perp A, B \mid D) \rightarrow$ edge $C-D$

$P \models (D \perp A, B \mid C) \rightarrow$ edge $D-C$

$P \models (B \perp C, D \mid A) \rightarrow$ edge $B-A$

but \mathcal{H} is not a I-map since $(A \perp C) \in \text{I}(\mathcal{H})$ $(A \perp C) \notin \text{IP}(P)$



2° pairwise independencies construction

none of A, B, C, D separates $(X, Y \mid \mathcal{C} - \{X, Y\})$

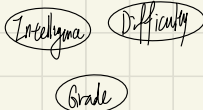
\therefore gives an empty graph

not an I-map for P



deterministic relations between variables can lead to failure in construction based on local and pairwise independence

ex



$(D \perp G) \mid (I \perp G)$
 $(I \perp D) \mid G$

} minimal I-map

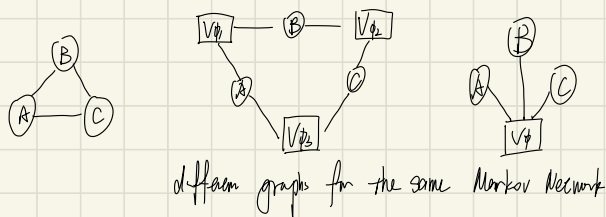


does not capture $(I \perp D)$

not every distribution P has a perfect map

Factor Graphs

A factor graph \mathcal{F} is an undirected graph containing 2 types of nodes: variable nodes or factor nodes.
 each edge connects two nodes of different types
 each factor node V_f has one factor ϕ whose scope is its neighbors



Bayesian Network to Markov Network

Given a Bayesian network B , how to represent the distribution P_B as a parameterized Markov network
 finding a minimal \mathcal{I} -map for distribution P_B

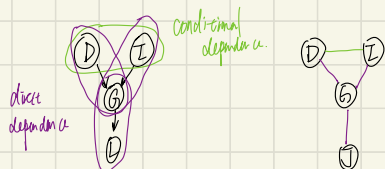
Given a graph G , how to represent the independencies in G using a undirected graph \mathcal{H}
 finding a minimal \mathcal{I} -map for the independencies $\mathcal{I}(G)$

let B be a Bayesian Network over \mathcal{X} and $E=e$ an observation. let $W=\mathcal{X}-E$
 $P_B(W|E=e)$ is a Gibbs distribution defined by $\mathcal{E}=\{\phi_{x_i}\}$
 $\phi_{x_i} = P_B(x_i | \text{parents}(x_i)) [E=e]$

The partition function for this Gibbs distribution is $P(e)$

the moral graph $M[G]$ of a Bayesian network structure G over \mathcal{X}
 is the undirected graph over \mathcal{X} that contains an undirected edge between x and y if

- there's a directed edge between x and y
- or
- x and y are both parents of the same node



Let G be a Bayesian network structure, for any distribution P such that

B is a parametrization of G , MIG is an I-map for P

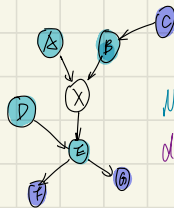
MIG captures local dependencies (direct independencies and v-structure conditional dependencies)

Let G be any Bayesian network graph, the moralized graph MIG is a minimal I-map

for any $X \in G$, select the smallest set U such that $(X \perp\!\!\!\perp U - \{X\} | U)$

the Markov blanket of X in a Bayesian Network G $MB_G(X)$ is

- $\{$
 - X 's parents
 - X 's children
 - other parents of X 's children



$MB_G(X)$

d-separated nodes

$d\text{-sep}(X; C | B)$

$d\text{-sep}(X; \{E, G\} | E)$

For a Bayesian Network graph G , $MB_G(X)$ d-separates X from all other variables

no subset of $MB_G(X)$ does

The addition of the moralizing edges to the Markov network π leads to the loss of independency information

But independency information is not always lost



MIG lose the information that X and Y are marginally independent



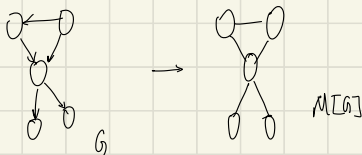
independency information is not lost

A Bayesian Network G is moral if it contains no immoralities

(all v-structures have a covering edge)

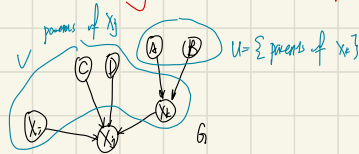


if a directed graph G is moral, then its moralized graph MIG is a perfect map of G



Markov Network to Bayesian Network

Let \mathcal{H} be a Markov network structure, G is any BN minimal I-map for \mathcal{H} .
 G can have no immoralities



Suppose there is immorality in G

$$\therefore G \models (X_1 \perp X_2 \mid X_3) \quad (X_4 \perp X_5 \mid X_6)$$

$$\therefore \mathcal{H} \models (X_1 \perp X_2 \mid X_3) \quad (X_4 \perp X_5 \mid X_6)$$

$\therefore \mathcal{H}$ contains one or more paths between X_1 and X_2 not cut by X_3
 between X_4 and X_5 not cut by X_6



\therefore there is one or more paths in \mathcal{H} between X_1 and X_2 via X_3

$$\therefore G \models \{X_1 \perp X_2 \mid U\} \quad \therefore \mathcal{H} \models \{X_1 \perp X_2 \mid U\}$$

$\therefore U$ block either paths between X_1 and X_2 or paths between X_3 and X_6 in \mathcal{H}

$$\therefore G \models (X_3 \perp U \mid V) \quad \therefore \mathcal{H} \models (X_3 \perp U \mid V)$$

$\therefore U$ cannot be in the path between X_3 and X_6

Contradiction

A Markov network may not have a BN perfect map
 of BN perfect map of MN



$$(A \perp C \mid B, D) \\ (B \perp D \mid A, C)$$



$$(B \perp D \mid A) \\ \text{not a perfect map}$$



$$(B \perp D \mid A, C) \\ (A \perp C) \\ \text{not a perfect map}$$

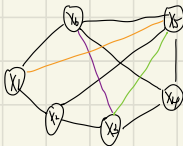


$$(A \perp C \mid B, D) \\ (D \perp B \mid A, C) \\ \text{not a perfect map}$$

Let $X_1 - X_2 \dots X_k$ be a loop in the graph; a chord in the loop is an edge connecting

X_i and X_j for two nonconsecutive nodes X_i and X_j .

A undirected graph \mathcal{H} is said to be chordal if any loop $X_1 - X_2 \dots X_k - X_1$ for $k \geq 4$ has a chord



Chord for loop $X_1 - X_2 - X_3 - X_4$

Chord for loop $X_2 - X_3 - X_4 - X_5$

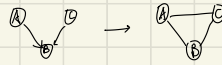
Chord for loop $X_3 - X_4 - X_5 - X_6$

Let \mathcal{H} be a nonchordal Markov network, then there is no Bayesian network G which is a perfect map for \mathcal{H} ($I(G) = I(\mathcal{H})$)

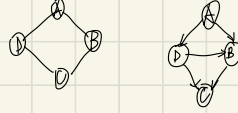
Let \mathcal{H} be a chordal Markov network, there is a Bayesian network G such that $I(\mathcal{H}) = I(G)$

Converting Bayesian Network \leftrightarrow Markov Network loses independencies

BN \rightarrow MN: lose independencies in V-structure



MN \rightarrow BN: add triangulating edges to loops



Conditional Random Field

A conditional Random Field is an undirected graph \mathcal{H} whose nodes correspond to $X \cup Y$, the network is annotated with a set of factors $\phi(D_1), \dots, \phi(D_m)$ such that each $D_i \subseteq X$ the network encodes a conditional distribution as

$$P(Y|X) = \frac{1}{Z(X)} \tilde{P}(Y, X)$$

$$= \frac{1}{\sum_{Y'} \prod_{i=1}^m \phi(D_i)} \prod_{i=1}^m \phi(D_i)$$

Two variables in \mathcal{H} are connected by an undirected edge whenever they appear together in the same scope of a factor. Instead of modeling a Joint Distribution $P(X, Y)$, CRF models a conditional distribution $P(Y|X)$ (we caring the distribution of X)

eg. logistic regression

X_i	Y	ϕ_i
0	0	$\exp(w_i)$
0	1	$\exp(0) = 1$
1	0	$\exp(0) = 1$
1	1	$\exp(w_i)$

$$\phi_i(X_i, Y) = \exp(w_i \cdot \mathbb{I}\{X_i=1, Y=1\})$$

$$\phi_i(X_i, Y) = \begin{cases} \exp(w_i) & \text{if } X_i=1 \\ \exp(0)=1 & \text{if } X_i=0 \end{cases} \Rightarrow \exp(w_i \cdot X_i) \quad (\text{how much } X_i \text{ contribute to } Y=1)$$

$$\tilde{P}_\theta(X, Y=1) = \prod \exp(w_i \cdot X_i) = \exp(\sum w_i \cdot X_i)$$

$$\tilde{P}_\theta(X, Y=0) = 1$$

$$P_\theta(Y=1|X) = \frac{\tilde{P}_\theta(X, Y=1)}{\tilde{P}_\theta(X, Y=1) + \tilde{P}_\theta(X, Y=0)} = \frac{\exp(\sum w_i X_i)}{1 + \exp(\sum w_i X_i)}$$

the sigmoid function

Log-Linear Models

rewrite a factor $\phi(D)$ as

$$\phi(D) = \exp(-E(D))$$

$$E(D) = -\ln(\phi(D))$$

$E(D) = -\ln \phi(D)$ is called an energy function,

in statistical physics, the probability of a physical state (eg. configuration of set of electrons) depends inversely on its energy.

any Markov network parameterized w/ positive factors can be converted to a log-linear representation

$$p(x_1, \dots, x_n) \propto \exp\left[-\sum_i E_i(D_i)\right]$$

eg.

$\phi_1(A, B)$	$\phi_2(B, C)$	$\phi_3(C, D)$	$\phi_4(D, A)$
$a^0 \ b^0 \ 30$	$b^0 \ c^0 \ 100$	$c^0 \ d^0 \ 1$	$d^0 \ a^0 \ 100$
$a^0 \ b^1 \ 5$	$b^0 \ c^1 \ 1$	$c^0 \ d^1 \ 100$	$d^0 \ a^1 \ 1$
$a^1 \ b^0 \ 1$	$b^1 \ c^0 \ 1$	$c^1 \ d^0 \ 100$	$d^1 \ a^0 \ 1$
$a^1 \ b^1 \ 10$	$b^1 \ c^1 \ 100$	$c^1 \ d^1 \ 1$	$d^1 \ a^1 \ 100$
(a)	(b)	(c)	(d)

$e_1(A, B)$	$e_2(B, C)$	$e_3(C, D)$	$e_4(D, A)$
$a^0 \ b^0 \ -3.4$	$b^0 \ c^0 \ -4.61$	$c^0 \ d^0 \ 0$	$d^0 \ a^0 \ -4.61$
$a^0 \ b^1 \ -1.61$	$b^0 \ c^1 \ 0$	$c^0 \ d^1 \ -4.61$	$d^0 \ a^1 \ 0$
$a^1 \ b^0 \ 0$	$b^1 \ c^0 \ 0$	$c^1 \ d^0 \ -4.61$	$d^1 \ a^0 \ 0$
$a^1 \ b^1 \ -2.3$	$b^1 \ c^1 \ -4.61$	$c^1 \ d^1 \ 0$	$d^1 \ a^1 \ -4.61$
(a)	(b)	(c)	(d)

-4.61 when 2 var agree and 0 otherwise

let D be a set of variables. we define a feature $f(D)$ to be a function $f: \text{val}(D) \mapsto \mathbb{R}$

$$\bar{P}_{\theta} = \prod_i \phi_i(D_i)$$

$$= \prod_j \exp(-w_j f_j(D))$$

f_j is a feature

eg.

x_1	x_2	f_i
0	0	a_0
0	1	a_{-1}
1	0	a_{10}
1	1	a_{11}

$$f^0(x_1, x_2) = \mathbb{1}\{x_1=0, x_2=0\}$$

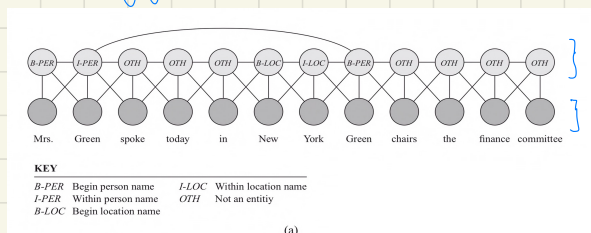
$$f^1(x_1, x_2) = \mathbb{1}\{x_1=0, x_2=1\}$$

$$\vdots$$

$$\phi_{\theta}(x_1, x_2) = \exp\left(-\sum_k w_k \cdot f^k(x_1, x_2)\right)$$

$$w_k = -\log a_k$$

eg. Features of Language.



γ : Category

x : word

$$f_i(x_i, \gamma_i) = \mathbb{1}\{\gamma_i = \text{person name}, x_i \text{ capitalized}\}$$

how important Capitalization is in recognizing a person

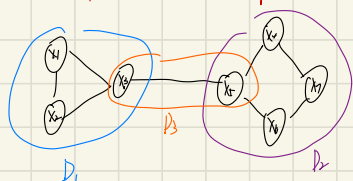
$$f_j(x_i, \gamma_i) = \mathbb{1}\{\gamma_i = \text{location name}, x_i \text{ appears in atlas}\}$$

A distribution P is a log-linear model over a Markov network \mathcal{H} if P is associated with

- a set of features $f = \{f_1(\mathbf{a}), \dots, f_k(\mathbf{a})\}$ where each D_i is a complete subgraph in \mathcal{H}
- a set of weights w_1, \dots, w_k

such that

$$p(x_1, \dots, x_n) = \frac{1}{Z} \exp \left[- \sum_{i=1}^k w_i f_i(\mathbf{a}) \right]$$



$$p(x_1, \dots, x_n) = \frac{1}{Z} \psi(D_1) \times \psi(D_2) \times \psi(D_3)$$

$$= \frac{1}{Z} \exp(-\epsilon_1(D_1)) \cdot \exp(-\epsilon_2(D_2)) \cdot \exp(-\epsilon_3(D_3))$$

$$= \frac{1}{Z} \exp \left[- \sum_{i=1}^k \epsilon_i(D_i) \right]$$

$$= \frac{1}{Z} \exp \left[- \sum_{i=1}^k w_i f_i(\mathbf{a}) \right]$$

• Metric MRFs

All x_i takes value in same label space V

Distance function $u: V \times V \mapsto \mathbb{R}^+$

reflexivity $u(x, x) = 0 \quad \forall x \in V$

symmetry $u(x_1, x_2) = u(x_2, x_1)$

triangle inequality: $u(x_1, x_3) \leq u(x_1, x_2) + u(x_2, x_3)$

$$f_{ij}(x_i, x_j) = u(x_i, x_j)$$

$\exp(-w_{ij} f_{ij}(x_i, x_j))$ where $w_{ij} > 0$ low distance \rightarrow high probability

$(x_i) - (x_j)$ want x_i and x_j to take "similar" values

eg

$$u(x_i, x_j) = |x_i - x_j|$$

$$u(x_i, x_j) = \begin{cases} 1 & \text{if } x_i \neq x_j \\ 0 & \text{if } x_i = x_j \end{cases}$$

eg Image segmentation

no penalty when adjacent super pixels take same class

some penalty otherwise