



Capturing Independence Graphically; Undirected Graphs

COMPSCI 276, Spring 2011
Set 2: Rina Dechter

(Reading: Pearl chapters 3, Darwiche chapter 4)



The Qualitative Notion of Dependence

motivations and issues

- The traditional definition of independence uses equality of numerical quantities as in $P(x,y)=P(x)P(y)$
- People can easily and confidently detect dependencies, but not provide numbers
- The notion of relevance and dependence are far more basic to human reasoning than the numerical
- Assertions about dependency relationships should be expressed first.



Dependency graphs

- The nodes represent propositional variables and the arcs represent local dependencies among conceptually related propositions.
- Graph concepts are entrenched in our language (e.g., “thread of thoughts”, “lines of reasoning”, “connected ideas”). One wonders if people can reason any other way except by tracing links and arrows and paths in some mental representation of concepts and relations.
- What types of (in)dependencies are deducible from graphs?
- For a given probability distribution P and any three variables $\mathbf{X, Y, Z}$, it is straightforward to verify whether knowing Z renders X independent of Y , but P does not dictate which variables should be regarded as neighbors.
- Some useful properties of dependencies and relevancies cannot be represented graphically.

Conditional Independence

4

DEFINITION: Let $U = \{\alpha, \beta, \dots\}$ be a finite set of variables with discrete values. Let $P(\cdot)$ be a joint probability function over the variables in U , and let X , Y , and Z stand for any three subsets of variables in U . X and Y are said to be *conditionally independent given Z* if

$$P(x|y, z) = P(x|z) \text{ whenever } P(y, z) > 0.$$

- For any ^{assignment} configuration x of the variables in the set X and for any configurations y and z of the variables in Y and Z satisfying $P(Y = y, Z = z) > 0$, we have

$$P(X = x|Y = y, Z = z) = P(X = x|Z = z).$$

$$\underline{I(X, Z, Y)}_P \text{ iff } P(x|y, z) = P(x|z)$$

for all values x , y , and z such that $P(y, z) > 0$.

- Marginal independence will be denoted by $I(X, \emptyset, Y)$, i.e.,
 $I(X, \emptyset, Y)$ iff $P(x|y) = P(x)$ whenever $P(y) > 0$.

Also $\underline{I(x, y|z)} \approx \underline{I(x, z, y)}$

Law of
roots

Implied independencies

5

- Partial list of (equivalent) properties satisfied by the conditional independence relation $I(X, Z, Y)$ [Lauritzen 1982]:

$$I(X, Z, Y) \Leftrightarrow P(x, y|z) = P(x|z) P(y|z), \quad \bullet$$

$$I(X, Z, Y) \Leftrightarrow \exists f, g : P(x, y, z) = f(x, z) g(y, z),$$

$$I(X, Z, Y) \Leftrightarrow P(x, y, z) = P(x|z) P(y, z).$$



Properties of Conditional Independance

$I_{pr}(X, Y, Z)$

6

THEOREM 1: Let X , Y , and Z be three disjoint subsets of variables from U . If $I(X, Z, Y)$ stands for the relation “ X is independent of Y , given Z ” in some probabilistic model P , then I must satisfy the following four independent conditions:

- Symmetry:

$$I(X, Z, Y) \Leftrightarrow I(Y, Z, X) \quad (1.a)$$

- Decomposition:

$$I(X, Z, Y \cup W) \Rightarrow I(X, Z, Y) \ \& \ I(X, Z, W) \quad (1.b)$$

- Weak Union:

$$I(X, Z, Y \cup W) \Rightarrow I(X, Z \cup W, Y) \quad (1.c)$$

- Contraction:

$$I(X, Z, Y) \ \& \ I(X, Z \cup Y, W) \Rightarrow I(X, Z, Y \cup W) \quad (1.d)$$

- If P is strictly positive, then a fifth condition holds:

- Intersection:

$$I(X, Z \cup W, Y) \ \& \ I(X, Z \cup Y, W) \Rightarrow I(X, Z, Y \cup W) \quad (1.e)$$

Properties of independence

7

INTUITIVE INTERPRETATION OF THE AXIOMS

- **Symmetry:** In any state of knowledge Z , if Y tells us nothing new about X , then X tells us nothing new about Y .
- **Decomposition:** If two combined items of information are judged irrelevant to X , then each separate item is irrelevant as well.
- **Weak union:** Learning irrelevant information W cannot help the irrelevant information Y become relevant to X .
- **Contraction:** If we judge W irrelevant to X after learning some irrelevant information Y , then W must have been irrelevant before we learned Y .
- **Together:** Irrelevant information should not alter the relevance of other propositions in the system; what was relevant remains relevant, and what was irrelevant remains irrelevant.
- **Intersection:** Unless Y affects X when W is held constant or W affects X when Y is held constant, neither W nor Y nor their combination can affect X .

Symmetry:

- $I(X,Z,Y) \rightarrow I(Y,Z,X)$

Decomposition:

- $I(X,Z,YW) \rightarrow I(X,Z,Y)$ and $I(X,Z,W)$

Weak union:

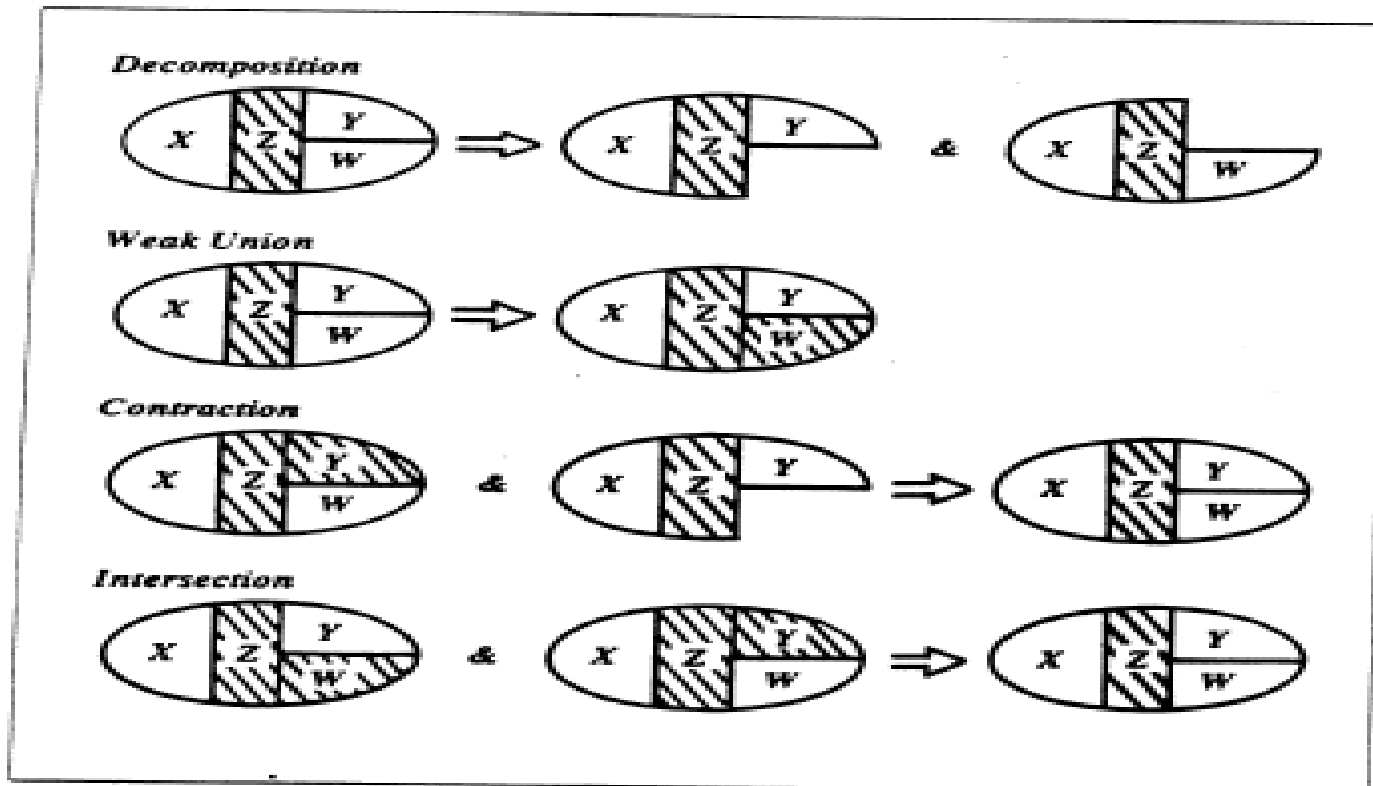
- $I(X,Z,YW) \rightarrow I(X,Z,W,Y)$

Contraction:

- $I(X,Z,Y)$ and $I(X,ZY,W) \rightarrow I(X,Z,YW)$

Intersection:

- $I(X,ZY,W)$ and $I(X,ZW,Y) \rightarrow I(X,Z,YW)$



Graphical interpretation of the axioms governing conditional independence.

Decomposition

If some information is irrelevant, then any part of it is also irrelevant.

$I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{Y} \cup \mathbf{W})$ only if $I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{Y})$ and $I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{W})$.

If learning $\mathbf{y}\mathbf{w}$ does not influence our belief in \mathbf{x} , then learning \mathbf{y} alone, or learning \mathbf{w} alone, will not influence our belief in \mathbf{x} either.

Decomposition

The opposite of Decomposition, called **Composition**:

$$I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{Y}) \text{ and } I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{W}) \text{ only if } I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{Y} \cup \mathbf{W})$$

does not hold in general.

Two pieces of information may each be irrelevant on their own, yet their combination may be relevant.

Example: Two coins and a bell

Contraction

$$I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{Y}) \text{ and } I_{Pr}(\mathbf{X}, \mathbf{Z} \cup \mathbf{Y}, \mathbf{W}) \text{ only if } I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{Y} \cup \mathbf{W})$$

If after learning the irrelevant information \mathbf{y} , the information \mathbf{w} is found to be irrelevant to our belief in \mathbf{x} , then the combined information \mathbf{yw} must have been irrelevant from the beginning.

Compare Contraction with Composition:

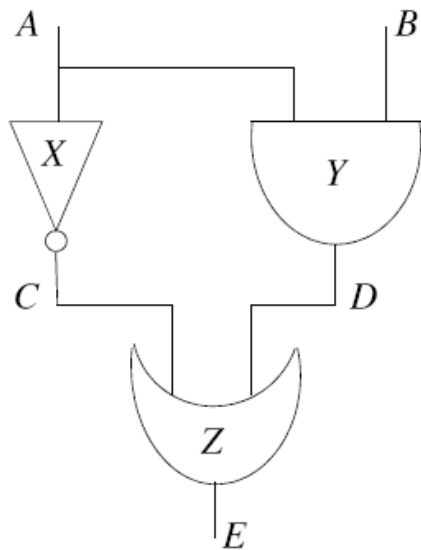
$$I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{Y}) \text{ and } I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{W}) \text{ only if } I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{Y} \cup \mathbf{W})$$

One can view Contraction as a weaker version of Composition. Recall that Composition does not hold for probability distributions.

Strictly Positive Distributions

Definition

A strictly positive distribution assign a non-zero probability to every consistent event.



Example

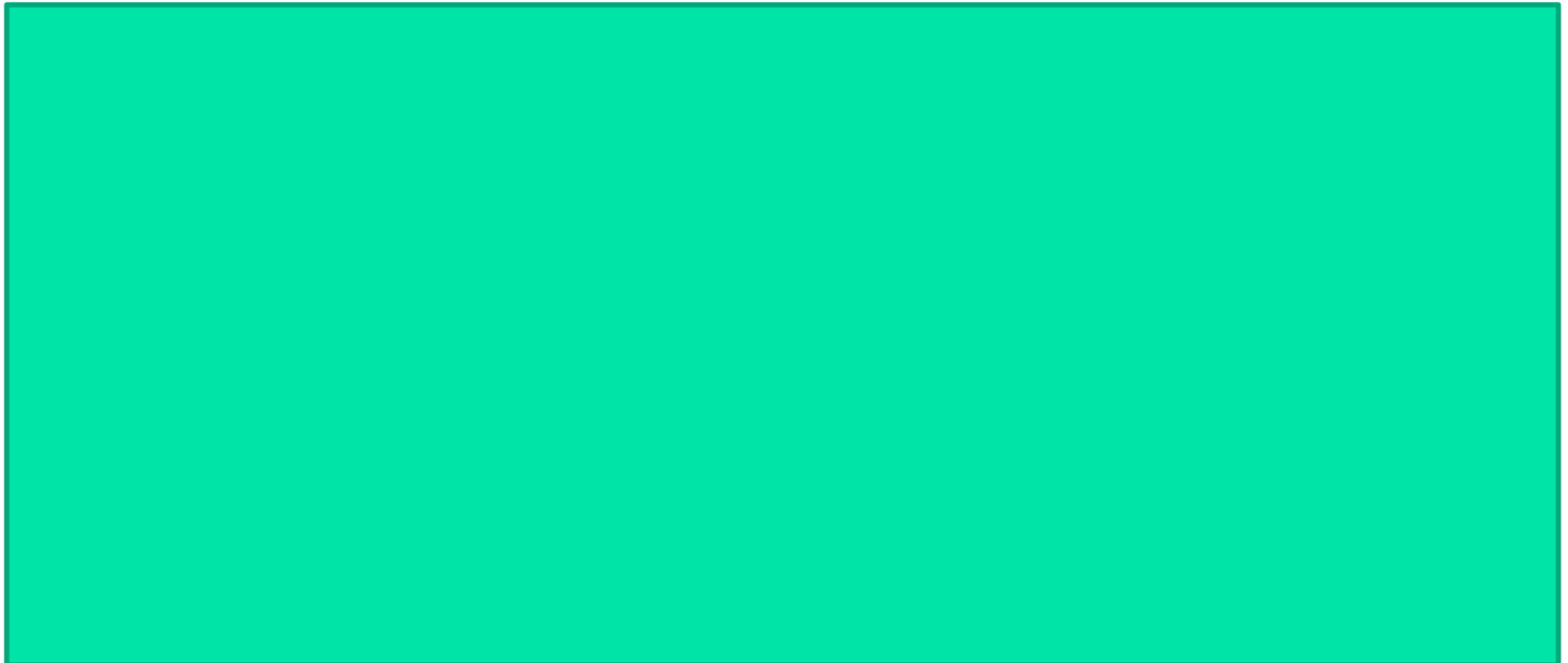
A strictly positive distribution cannot represent the behavior of Inverter X as it will have to assign the probability zero to the event $A = \text{true}, C = \text{true}$.

A strictly positive distribution cannot capture logical constraints.

Intersection

Holds only for strictly positive distributions

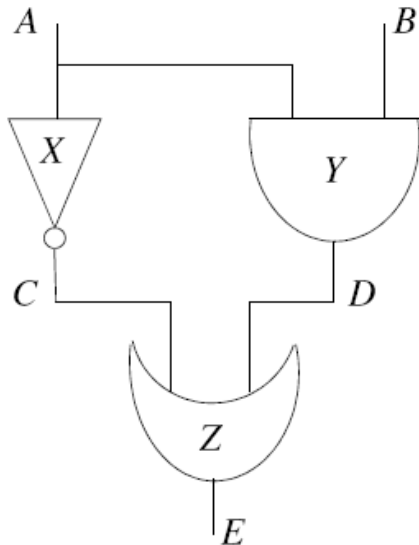
$I_{Pr}(\mathbf{X}, \mathbf{Z} \cup \mathbf{W}, \mathbf{Y})$ and $I_{Pr}(\mathbf{X}, \mathbf{Z} \cup \mathbf{Y}, \mathbf{W})$ only if $I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{Y} \cup \mathbf{W})$
If information \mathbf{w} is irrelevant given \mathbf{y} , and \mathbf{y} is irrelevant given \mathbf{w} ,
then combined information \mathbf{yw} is irrelevant to start with.



Intersection

Holds only for strictly positive distributions

$I_{Pr}(\mathbf{X}, \mathbf{Z} \cup \mathbf{W}, \mathbf{Y})$ and $I_{Pr}(\mathbf{X}, \mathbf{Z} \cup \mathbf{Y}, \mathbf{W})$ only if $I_{Pr}(\mathbf{X}, \mathbf{Z}, \mathbf{Y} \cup \mathbf{W})$
If information \mathbf{w} is irrelevant given \mathbf{y} , and \mathbf{y} is irrelevant given \mathbf{w} , then combined information \mathbf{yw} is irrelevant to start with.



- If we know the input A of inverter X , its output C becomes irrelevant to our belief in the circuit output E .
- If we know the output C of inverter X , its input A becomes irrelevant to this belief.
- Yet, variables A and C are not irrelevant to our belief in the circuit output E .

Graphs vs Graphoids



- Symmetry:
 - $I(X,Z,Y) \rightarrow I(Y,Z,X)$
- Decomposition:
 - $I(X,Z,YW) \rightarrow I(X,Z,Y)$ and $I(X,Z,W)$
- Weak union:
 - $I(X,Z,YW) \rightarrow I(X,ZW,Y)$
- Contraction:
 - $I(X,Z,Y)$ and $I(X,ZY,W) \rightarrow I(X,Z,YW)$
- Intersection:
 - $I(X,ZY,W)$ and $I(X,ZW,Y) \rightarrow I(X,Z,YW)$
- **Graphoid**: satisfy all 5 axioms
- **Semi-graphoid**: satisfies the first 4.
- Decomposition is only one way while in graphs it is iff.
- Weak union states that w should be chosen from a set that, like Y should already be separated from X by Z



Why axiomatic characterization?

- Allow deriving conjectures about independencies that are clearer
- Axioms serve as inference rules
- Can capture the principal differences between various notions of relevance or independence

DEPENDENCY MODELS AND DEPENDENCY MAPS

- A *dependency model* M is a rule that assigns truth values to the three-place predicate $I(X, Z, Y)_M$
- M determines a subset I of triplets (X, Z, Y) for which the assertion " X is independent of Y given Z " is true.
- Any probability distribution P is a dependency model.
- An undirected graph $G = (V, E)$ is a *graphical representation* of a dependency model M , if there is a correspondence between the elements in U (of M) and the set of vertices in V (of G), such that the topology of G reflects some properties of M .
- If a subset Z of nodes in a graph G intercepts all paths between the nodes of X and those of Y we write $\langle X \mid Z \mid Y \rangle_G$.



I-map and D-maps

3-13

DEFINITION: An undirected graph G is a *dependency map* (or *D-map*) of M if there is a one-to-one correspondence between the elements of U and the nodes V of G , such that for all disjoint subsets X, Y, Z of elements we have

$$I(X, Z, Y)_M \Rightarrow \langle X \mid Z \mid Y \rangle_G.$$

Similarly, G is an *independency map* (or *I-map*) of M if

$$I(X, Z, Y)_M \Leftarrow \langle X \mid Z \mid Y \rangle_G.$$

G is said to be a *perfect map* of M if it is both a *D-map* and an *I-map*.

- A *D-map* guarantees that vertices found to be connected are indeed dependent in M .
 - An *I-map*, guarantees that vertices found to be separated correspond to independent variables.
 - Empty graphs are trivial *D-maps*, while complete graphs are trivial *I-maps*.
- A model with induced dependencies cannot be i-map and d-map
 - Example: two coins and a bell... try it
 - **How we then represent two causes leading to a common consequence?**



Axiomatic characterization of Graphs

- **Definition:** A model M is graph-isomorph if there exists a graph which is a perfect map of M .
- **Theorem (Pearl and Paz 1985):** A necessary and sufficient condition for a dependency model to be graph-isomorph is that it satisfies
 - Symmetry: $I(X,Z,Y) \rightarrow I(Y,Z,X)$
 - Decomposition: $I(X,Z,YW) \rightarrow I(X,Z,Y)$ and $I(X,Z,Y)$
 - Intersection: $I(X,ZW,Y)$ and $I(X,ZY,W) \rightarrow I(X,Z,YW)$
 - Strong union: $I(X,Z,Y) \rightarrow I(X,ZW, Y)$
 - Transitivity: $I(X,Z,Y) \rightarrow$ exists t s.t. $I(X,Z,t)$ or $(I(t,Z,Y)$
- This properties are satisfied by graph separation



Markov Networks

- **Graphs and probabilities:**
 - Given P , can we construct a graph I-map with minimal edges?
 - Given (G, P) can we test if G is an I-map? a perfect map?
- **Markov Network:** A graph G which is a minimal I-map of a dependency model P , namely deleting any edge destroys its i-mapness, is called a **Markov network** of P .



Markov Networks

- **Theorem (Pearl and Paz 1985):** A dependency model satisfying symmetry decomposition and intersection has a unique minimal graph as an i-map, produced by deleting every edge (a,b) for which $I(a, \mathbf{U}-a-b, b)$ is true.
- The theorem defines an edge-deletion method for constructing G_0
- **Markov blanket** of a is a set S for which $I(a, S, \mathbf{U}-S-a)$.
- **Markov Boundary:** a minimal Markov blanket.

- **Theorem (Pearl and Paz 1985):** if symmetry, decomposition, weak union and intersection are satisfied by P , the Markov boundary is unique and it is the neighborhood in the Markov network of P



Markov Networks

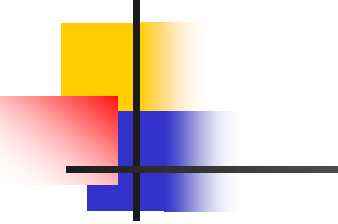
- **Corollary:** the Markov network G of any strictly positive distribution P can be obtained by connecting every node to its Markov boundary.
- The following 2 interpretations of direct neighbors are identical:
 - Neighbors as blanket that shields a variable from the influence of all others
 - Neighborhood as a tight influence between variables that cannot be weakened by other elements in the system
- So, given P (positive) how can we construct G ?
- Given (G,P) how do we test the G is an I-map of P ?
- Given G , can we construct P which is a perfect i-map? (Geiger and Pearl 1988)



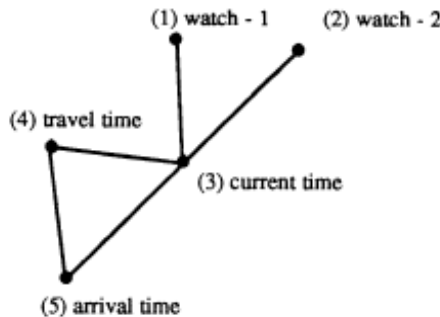
Testing I-mapness

- Theorem: Given a positive P and a graph G the following are equivalent:
 - G is an I-map of P iff G is a super-graph of the Markov network of P
 - G is locally Markov w.r.t. P (the neighbors of a in G is a Markov blanket.) iff G is a super-graph of the Markov network of P
- There appear to be no test for i-mapness of undirected graph that works for extreme distributions without testing every cutset in G (ex: $x=y=z=t$)
- Representations of probabilistic independence using undirected graphs rest heavily on the intersection and weak union axioms.
- In contrast, we will see that directed graph representations rely on the contraction and weak union axiom, with intersection playing a minor role.

CONCEPTUAL DEPENDENCIES AND THEIR MARKOV NETWORKS

- 
- An agent identifies the following variables as having influence on the main question of being late to a meeting:
 1. The time shown on the watch of Passerby 1.
 2. The time shown on the watch of Passerby 2.
 3. The correct time.
 4. The time it takes to travel to the meeting place.
 5. The arrival time at the meeting place.
 - The construction of G_0 can proceed by one of two methods:
 - The *edge-deletion* method.
 - The *Markov boundary* method.
 - The first method requires that for every pair of variables (α, β) we determine whether fixing the values of all other variables in the system will render our belief in α sensitive to β .
 - For example, the reading on Passerby 1's watch (1) will vary with the actual time (3) even if all other variables are known, so connect node 1 to node 3

- The Markov boundary method requires that for every variable α in the system, we identify a minimal set of variables sufficient to render the belief in α insensitive to all other variables in the system.
- For instance, once we know the current time (3), no other variable can affect what we expect to read on passerby 1's watch (1).



The unusual edge (3,4) reflects the reasoning that if we fix the arrival time (5) the travel time (4) must depend on current time (3).

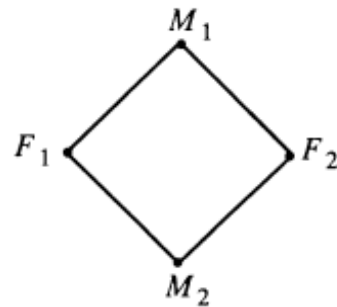
Figure 3.6. The Markov network representing the prediction of A's arrival time.

- G_0 can be used as an inference instrument.
- For example, knowing the current time (3) renders the time on Passerby 1's watch (1) irrelevant for estimating the travel time (4) (i.e., $I(1,3,4)$); 3 is a cutset in G_0 , separating 1 from 4.

Summary

- The essential qualities of conditional independence are captured by five logical axioms: (a) symmetry, (b) decomposition, (c) weak union, (d) contraction and (e) intersection.
- Intersection holds only for strictly positive distributions (i.e., reflecting no functional or definitional constraints) and is essential to the construction of undirected graphs.
- Symmetry, decomposition, and intersection enable us to construct a minimal graph G_0 (Markov network), in which every cutset corresponds to a genuine independence condition.
- The weak union axiom is needed to guarantee that the set of neighbors that G_0 assigns to each variable α is the smallest set required to shield α from the effects of all other variables.
- If we identify the Markov boundaries associated with each proposition in the system and treat them as neighborhood relations defining a graph G_0 , then we can correctly identify independence relationships by testing whether the set of known propositions constitutes a cutset in G_0 .
- Not all probabilistic dependencies can be captured by undirected graphs because graph separation is strictly normal and transitive.

MARKOV NETWORK AS A KNOWLEDGE BASE



How can we construct a probability Distribution that will have all these independencies?

Figure 3.2. An undirected graph representing interactions among four individuals.

QUANTIFYING THE LINKS

- If couple (M_1, F_2) meet less frequently than the couple (M_1, F_1) , then the first link should be weaker than the second
- The model must be consistent, complete and a Markov field of G .
- Arbitrary specification of $P(M_1, F_1)$, $P(F_1, M_2)$, $P(M_2, F_2)$, and $P(F_2, M_1)$ might lead to inconsistencies.
- If we specify the pairwise probabilities of only three pairs, incompleteness will result.

- A safe method (called *Gibbs' potential*) for constructing a complete and consistent quantitative model while preserving the dependency structure of an arbitrary graph G .

1. Identify the cliques[†] of G , namely, the largest subgraphs whose nodes are all adjacent to each other.
2. For each clique C_i , assign a nonnegative compatibility function $g_i(c_i)$, which measures the relative degree of compatibility associated with the value assignment c_i to the variables included in C_i .
3. Form the product $\prod_i g_i(c_i)$ of the compatibility functions over all the cliques.
4. Normalize the product over all possible value combinations of the variables in the system

$$P(x_1, \dots, x_n) = K \prod_i g_i(c_i), \quad (3.13)$$

where

$$K = \left[\sum_{x_1, \dots, x_n} \prod_i g_i(c_i) \right]^{-1}.$$

[†] We use the term *clique* for the more common term *maximal clique*.

Example: The dependency graph has four cliques, corresponding to the four edges

$$C_1 = \{M_1, F_1\}, C_2 = \{M_1, F_2\},$$

$$C_3 = \{M_2, F_1\}, \text{ and } C_4 = \{M_2, F_2\},$$

the compatibility functions g_i are assessed to be

$$g_i(x_{i_1}, x_{i_2}) = \begin{cases} \alpha_i & \text{if } x_{i_1} = x_{i_2} \\ \beta_i & \text{if } x_{i_1} \neq x_{i_2}, \end{cases} \quad (3.14)$$

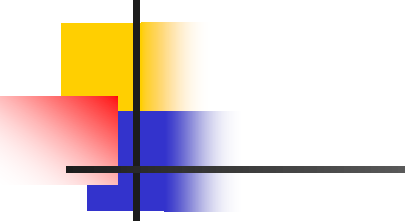
where x_{i_1} and x_{i_2} are the states of disease associated with the male and female, respectively, of couple C_i .

- The overall probability distribution function is given by the normalized product

$$\begin{aligned} P(M_1, M_2, F_1, F_2) &= K g_1(M_1, F_1) g_2(M_1, F_2) g_3(M_2, F_1) g_4(M_2, F_2) \\ &= K \prod_i \beta_i^{|x_{i_1} - x_{i_2}|} \alpha_i^{1 - |x_{i_1} - x_{i_2}|}, \end{aligned} \quad (3.15)$$

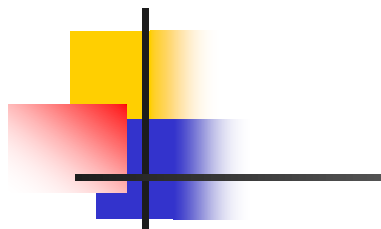
where K is a constant that makes P sum to unity over all states of the system, i.e.,

$$K^{-1} = \prod_i (\alpha_i + \beta_i) + \prod_i \alpha_i \sum_j \frac{\beta_j}{\alpha_j} + \prod_i \beta_i \sum_j \frac{\alpha_j}{\beta_j}. \quad (3.16)$$

- 
- For example, the state in which only the males carry the disease, $(m_1, \neg f_1, m_2, \neg f_2)$, will have a probability measure $K \beta_1 \beta_2 \beta_3 \beta_4$
 - The state $(m_1, f_1, \neg m_2, \neg f_2)$, on the other hand, has the probability $K \alpha_1 \beta_2 \beta_3 \alpha_4$
 - P is a Markov field of G because

$$P = f(M_1, F_1, F_2) g(F_1, F_2, M_2) = f'(F_1, M_1, M_2) g'(M_1, M_2, F_2).$$

Thus, $I(M_1, F_1 \cup F_2, M_2)_P$ and $I(F_1, M_1 \cup M_2, F_2)_P$.



G is locally markov
 If neighbors make every
 Variable independent
 From the rest.

THEOREM 6 [Hammersley and Clifford 1971]: A probability function P formed by a normalized product of nonnegative functions on the cliques of G is a **Markov field** relative to G , i.e., G is an I -map of P .

Proof: G is guaranteed to be an I -map if P is locally Markov relative to G (Theorem 5). It is sufficient, therefore, to show that the neighbors in G of each variable α constitute a Markov blanket of α relative to P , i.e., that $I(\alpha, \mathbf{B}_G(\alpha), U - \alpha - \mathbf{B}_G(\alpha))$ or (using Eq. (3.5c)) that

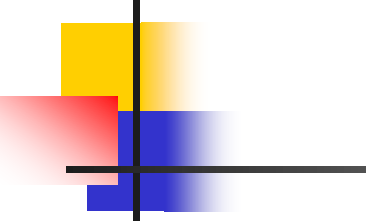
$$P(\alpha, \mathbf{B}_G(\alpha), U - \alpha - \mathbf{B}_G(\alpha)) = f_1(\alpha, \mathbf{B}_G(\alpha)) f_2(U - \alpha). \quad (3.17)$$

- Let J_α stand for the set of indices marking all cliques in G that include α , $J_\alpha = \{j : \alpha \in C_j\}$. Since P is in product form, we can write

$$P(\alpha, \beta, \dots) = K \prod_j g_j(\mathbf{c}_j) = K \prod_{j \in J_\alpha} g_j(\mathbf{c}_j) \prod_{j \notin J_\alpha} g_j(\mathbf{c}_j). \quad (3.18)$$

- The first product in Eq. (3.18) contains only variables that are adjacent to α in G ; otherwise, C_j would not be a clique. According to the definition of J_α , the second product does not involve α . Thus, Eq. (3.17) is established. Q.E.D.

INTERPRETING THE LINK PARAMETERS

- 
- It is difficult to assign meanings to the parameters of the compatibility functions.
 - Given the joint probability $P(M_1, F_1, F_2, M_2)$, to infer the compatibility functions g_i we must solve a set of simultaneous nonlinear equations for g_i
 - The solution for g_i will not be applicable to new situations.
 - For a parameter to be meaningful, it must be an abstraction of some invariant property of one's experience.
 - The quantities $P(f_1|m_1, \neg m_2)$ and $P(f_1|\neg m_1, \neg m_2)$ and their relations to the frequency of interaction of couple $\{M_1, F_1\}$ are perceived as invariant characteristics of the disease.
 - The Markov network formulation does not allow the direct specification of such judgmental input.
 - Judgments about low-order conditional probabilities (e.g., $P(m_1|f_1, \neg m_2)$) can be taken only as constraints that the joint probability distribution (Eq. (3.13)) must satisfy; from them, we might be able to calculate the actual values of the compatibility parameters.