

MATH 829: Introduction to Data Mining and Analysis

Graphical Models I

Dominique Guillot

Departments of Mathematical Sciences
University of Delaware

May 2, 2016

3/32

Example (cont.)

We compute the correlation between the grades of the students:

mec h	1.0				
vec t	0.55	1.0			
a l g	0.55	0.61	1.0		
a n a l	0.41	0.49	0.71	1.0	
s t a t	0.39	0.44	0.66	0.61	1.0
	mec h	vec t	a l g	a n a l	s t a t

- We observe that the grades are positively correlated between subjects (performance of a student across subjects (good or bad) is consistent).

We now examine the inverse correlation matrix:

mec h	1.60				
vec t	-0.56	1.80			
a l g	-0.51	-0.66	3.04		
a n a l	0.00	-0.15	-1.11	2.18	
s t a t	-0.04	-0.04	-0.86	-0.52	1.92
	mec h	vec t	a l g	a n a l	s t a t

3/32

Independence and conditional independence: motivation

We begin with a classical example (Whittaker, 1990):

- We study the examination marks of 88 students in five subjects: mechanics, vectors, algebra, analysis, statistics (Mardia, Kent, and Bibby, 1979).
- Mechanics and vectors were closed books.
- All the remaining exams were open books.

We can examine the results using a stem and leaf plot.

alg e b r a	m e c h a n i c s
0:	0
10:	46 770000
20:	012 46 5 690
30:	11 66077000
40:	01 13333305 566667777000090000
50:	000001 112 133333046651 666770000
60:	0001 11 1124051 70
70:	0
80:	0
90:	0

Note: Data appears to be normally distributed.

3/32

Example (cont.)

Interpreting the inverse correlation matrix:

- Diagonal entries = $1/(1 - R^2)$ are related to the proportion of variance explained by regressing the variable on the other variables.
- Off-diagonal entries: proportional to the correlation of pairs of variables, given the rest of the variables.

For example, $R_{mec h}^2 = (1.60 - 1)/1.60 = 37.5\%$.

For the off-diagonal entries, we first scale the inverse correlation matrix

$\Omega = (\omega_{ij})$ by computing $\frac{\omega_{ij}}{\sqrt{\omega_{ii}\omega_{jj}}}$:

mec h	1.0				
vec t	-0.33	1.0			
a l g	-0.23	-0.28	1.0		
a n a l	0.00	-0.08	-0.43	1.0	
s t a t	-0.02	-0.02	-0.36	-0.25	1.0
	mec h	vec t	a l g	a n a l	s t a t

The off-diagonal entries of the scaled inverse correlation matrix are the **negative of the conditional correlation coefficients** (i.e., the correlation coefficients after conditioning on the rest of the variables).

3/32

Example (cont.)

Notation:

- We denote the fact that two random variables X and Y are independent by $X \perp\!\!\!\perp Y$.
- We write $X \perp\!\!\!\perp Y | \{Z_1, \dots, Z_n\}$ when X and Y are independent given Z_1, \dots, Z_n .
- When the context is clear (i.e. when working with a fixed collection of random variables $\{X_1, \dots, X_n\}$), we write

$$X_i \perp\!\!\!\perp X_j \mid \text{rest}$$

instead of $X_i \perp\!\!\!\perp X_j | \{X_k : 1 \leq k \leq n, k \neq i, j\}$.

Important: In general, uncorrelated variables are not independent. This is true however for the multivariate Gaussian distribution.

5/32

Example (cont.)

We noted before that our data appears to be Gaussian. Therefore it appears that:

- $\text{anal} \perp\!\!\!\perp \text{mech} \mid \text{rest}$.
- $\text{anal} \perp\!\!\!\perp \text{vect} \mid \text{rest}$.
- $\text{stat} \perp\!\!\!\perp \text{mech} \mid \text{rest}$.
- $\text{stat} \perp\!\!\!\perp \text{vect} \mid \text{rest}$.

We represent these relations using a **graph**:



We put **no edge** between two variables iff they are conditionally independent (given the rest of the variables).

6/32

Independence and factorizations

Graphical models (a.k.a Markov random fields) are multivariate probability models whose independence structure is characterized by a graph.

Recall: Independence of random vectors is characterized by a factorization of their joint density:

- **Independent variables:** For two random vectors X, Y :

$$X \perp\!\!\!\perp Y \Leftrightarrow f_{X,Y}(x,y) = g(x)h(y) \quad \forall x,y$$

for some functions g, h .

- **Conditionally independent variables:** Similarly, for three random vectors X, Y, Z :

$$X \perp\!\!\!\perp Y | Z \Leftrightarrow f_{X,Y,Z}(x,y,z) = g(x,z)h(y,z)$$

for all x, y and all z for which $f_Z(z) > 0$.

7/32

Independence graphs

Let $X = (X_1, \dots, X_p)$ be a random vector.

- The **conditional independence graph** of X is the graph $G = (V, E)$ where $V = \{1, \dots, p\}$ and

$$(i, j) \notin E \Leftrightarrow X_i \perp\!\!\!\perp X_j \mid \text{rest}.$$

- A subset $S \subset V$ is said to separate $A \subset V$ from $B \subset V$ if every path from A to B contains a vertex in S .

Notation: If $X = (X_1, \dots, X_p)$ and $A \subset \{1, \dots, p\}$, then $X_A := (X_i : i \in A)$.

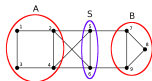
Theorem: (the separation theorem) Suppose the density of X is positive and continuous. Let $V = A \cup S \cup B$ be a partition of V such that S separates A from B . Then

$$X_A \perp\!\!\!\perp X_B \mid X_S.$$

8/32

Independence graphs (cont.)

Example: $X = (X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8, X_9)$:



Then

$$(X_1, X_2, X_3, X_4) \perp\!\!\!\perp (X_7, X_8, X_9) \mid (X_5, X_6).$$

9/32

Markov properties

Let $X = (X_1, \dots, X_p)$ be a random vector and let G be a graph on $\{1, \dots, p\}$. The vector is said to satisfy:

- The **pairwise Markov property** if $X_i \perp\!\!\!\perp X_j \mid \text{rest}$ whenever $(i, j) \notin E$.
- The **local Markov property** if for every vertex $i \in V$,

$$X_i \perp\!\!\!\perp X_{V \setminus \text{cl}(i)} \mid X_{\text{ne}(i)},$$

where $\text{ne}(i) = \{j \in V : (i, j) \in E, j \neq i\}$ and $\text{cl}(i) = \{i\} \cup \text{ne}(i)$.

- The **global Markov property** if for every disjoint subsets $A, S, B \subset V$ such that S separates A from B in G , we have

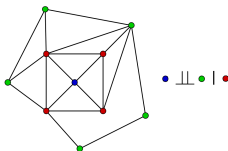
$$X_A \perp\!\!\!\perp X_B \mid X_S.$$

- Clearly, **global** \Rightarrow **local** \Rightarrow **pairwise**.
- When X has a positive and continuous density, by the separation theorem, **pairwise** \Rightarrow **global** and so all three properties are equivalent.

10/32

Example: the local Markov property

Illustration of the local Markov property:



11/32

The Hammersley–Clifford theorem

- An **undirected graphical model** (a.k.a. Markov random field) is a set of random variables satisfying a Markov property.
- Independence and conditional independence correspond to a factorization of the joint density.
- It is natural to try to characterize Markov properties via a factorization of the joint density.
- The Hammersley–Clifford theorem provides a necessary and sufficient condition for a random vector to have a Markov random field structure.

Theorem: (Hammersley–Clifford) Let X be a random vector with a positive and continuous density f . Then X satisfies the *pairwise Markov property* with respect to a graph G if and only if

$$f(x) = \prod_{C \in \mathcal{C}} \psi_C(x_C),$$

where \mathcal{C} is the set of (maximal) cliques (complete subgraphs) of G .

12/32