

Probabilistic Graphical Models

Florence Forbes

Equipe Mistis

INRIA Grenoble RHONE-ALPES

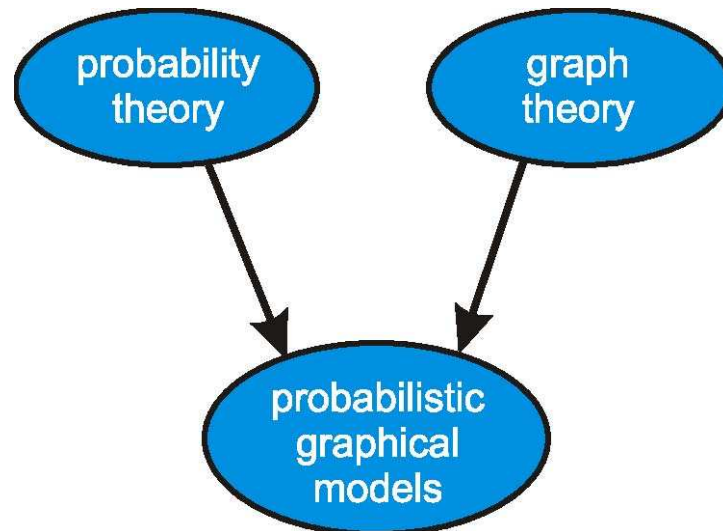
3 juillet 2014

Probabilistic graphical models

- Graphical models are used in various domains:
 - Machine learning and artificial intelligence
 - Computational biology
 - Statistical signal and image processing
 - Communication and information theory
 - Statistical physics.....
- Based on correspondences between graph theory and probability theory
- Important but difficult problems:
 - Computing **likelihoods**, **marginal distributions**, **modes**
 - Estimating model **parameters** and **structure** from noisy data

Probabilistic Graphical Models

- **Role of the graphs:**
 - graphical representations of probability distributions
 - Visualize the structure of a model
 - Design and motivate new models
 - Design graph based algorithms for inference



Probability Theory

- Sum rule

$$p(\mathbf{x}) = \sum_{\mathbf{y}} p(\mathbf{x}, \mathbf{y})$$

- Product rule

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}|\mathbf{y})p(\mathbf{y})$$

- From these we have Bayes' theorem

$$p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{y})p(\mathbf{y})}{p(\mathbf{x})}$$

– with normalization

$$p(\mathbf{x}) = \sum_{\mathbf{y}} p(\mathbf{x}|\mathbf{y})p(\mathbf{y})$$

Outline of the talk

- Directed graphs: Bayesian Networks
- Conditional independence and Markov properties
- Undirected graphs: Markov Random Fields
- Inference and learning
- Some illustrations

Directed graphs

Bayesian Networks

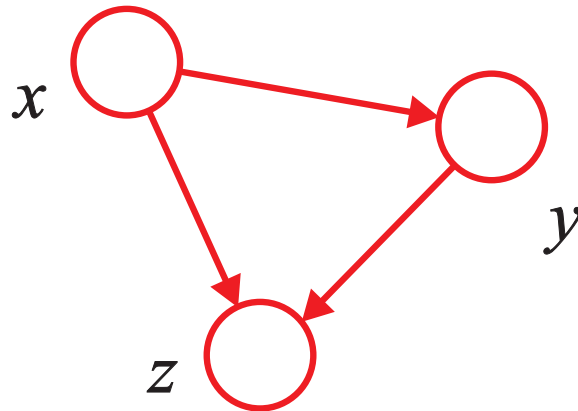
Directed Graphs: Decomposition

- Consider an arbitrary joint distribution

$$p(x, y, z)$$

- By successive application of the product rule

$$\begin{aligned} p(x, y, z) &= p(x)p(y, z|x) \\ &= p(x)p(y|x)p(z|x, y) \end{aligned}$$



General Case

- **Arbitrary** joint distribution,

$$P(x_1, \dots, x_n)$$

- Successive application of the product rule

$$P(x_1, \dots, x_n) = P(x_1)P(x_2|x_1) \dots P(x_n|x_1 \dots x_{n-1})$$

- Can be represented by a **fully connected graph** (links to all lower-numbered nodes)

Information is in the absence of links

General relationship

- Factorization property

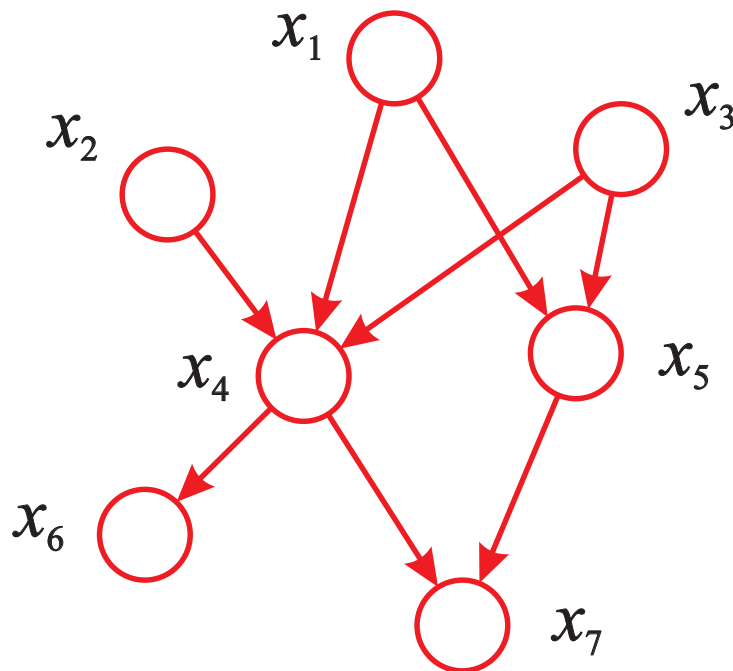
$$P(x_1, \dots, x_n) = \prod_{k=1}^n P(x_k | pa_k)$$

Where pa_k denotes the parents of x_k

- Missing link imply conditional independencies

Directed Acyclic Graphs: Bayesian Networks

- The graph can be used to impose constraints on the random vector (x_1, \dots, x_7) (ie. on the distribution P):

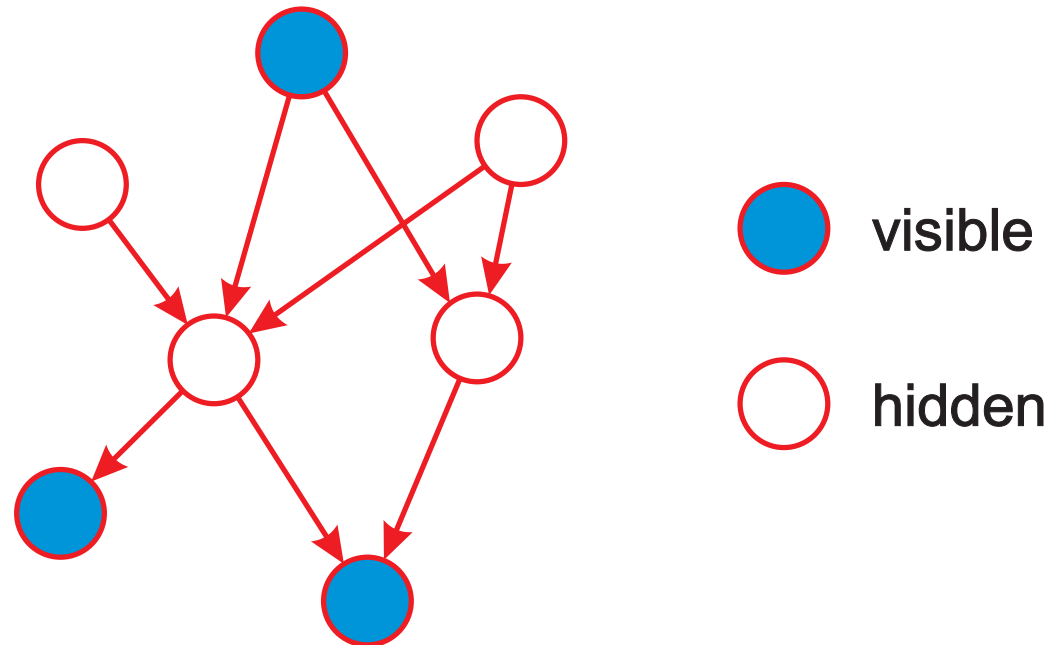


$$P(x_1)P(x_2)P(x_3)$$
$$P(x_4|x_1, x_2, x_3)$$
$$P(x_5|x_1, x_3)$$
$$P(x_6|x_4)$$
$$P(x_7|x_4, x_5)$$

No directed cycles

Hidden variables

- Variables may be hidden (latent) or visible (observed)



- Latent variables may have a specific interpretation, or may be introduced to permit a richer class of distribution

Example 1: Mixtures of Gaussians

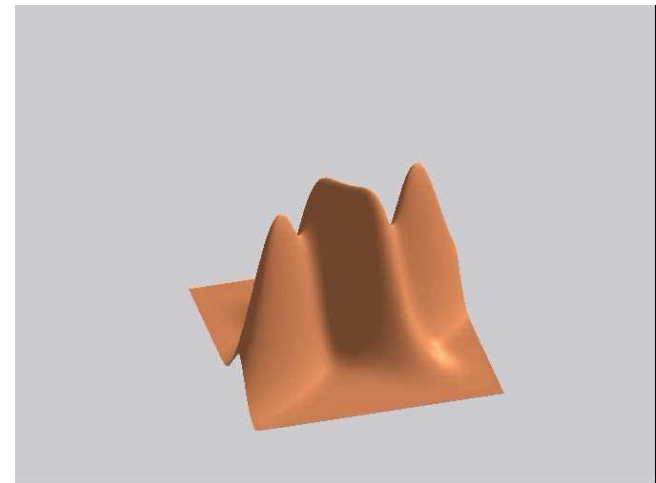
- Linear super-position of K Gaussians

$$P(y) = \sum_{k=1}^K \pi_k \mathcal{N}(y | \mu_k, \sigma_k^2)$$

- Normalization and positivity require

$$\sum_{k=1}^K \pi_k = 1 \quad 0 \leq \pi_k \leq 1$$

- illustration: mixture of 3 Gaussians



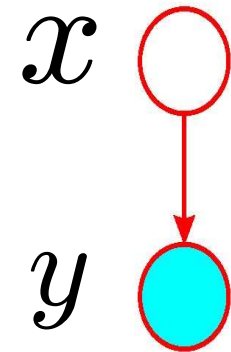
Latent Variable Viewpoint

- Discrete latent variable $x \in \{1, \dots, K\}$ describing which component generated data point y
- Conditional distribution of observed variable

$$P(y|X = k) = \mathcal{N}(y|\mu_k, \sigma_k^2)$$

- Prior distribution of latent variable

$$P(X = k) = \pi_k$$



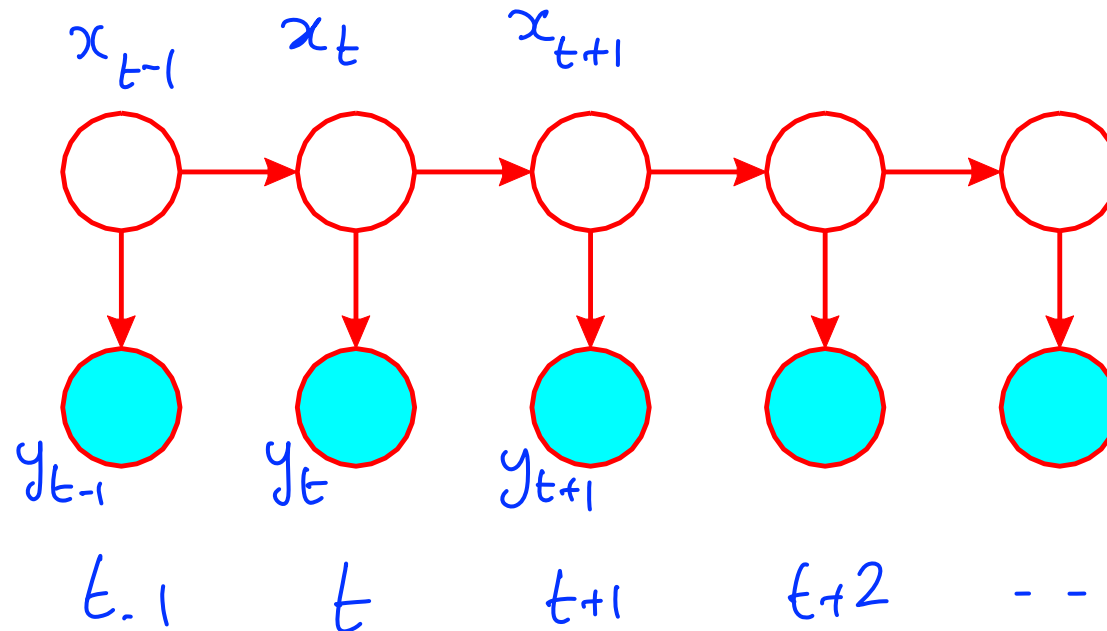
- Marginalizing over the latent variable we obtain

$$P(y) = \sum_{k=1}^K \pi_k \mathcal{N}(y|\mu_k, \sigma_k^2)$$

Example 2: State Space Models

- Hidden Markov chain
- Kalman filter

$$\dots p(x_t | x_{t-1}) p(y_t | x_t) p(x_{t+1} | x_t) \dots$$

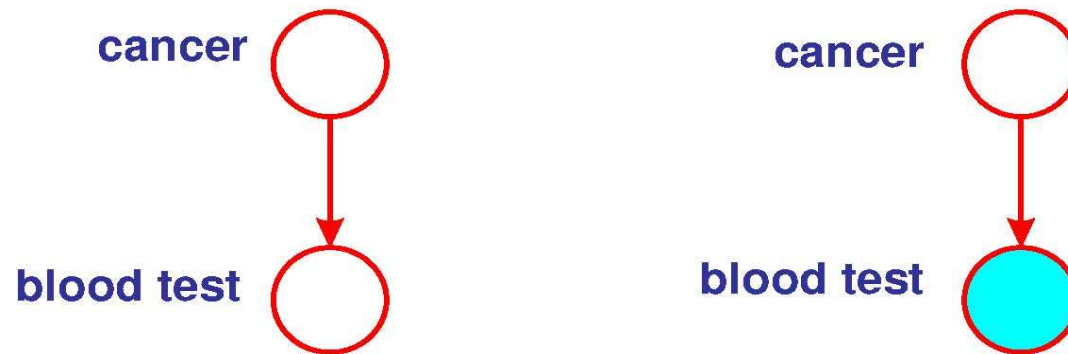


- Frequently wish to solve the problem of computing

$$p(x_t \mid y_1, \dots, y_n)$$

Causality

- **Directed graphs** can express **causal** relationships
- Often we **observe child** variables and wish to **infer** the posterior distribution of **parent** variables
- Example:



- Note: inferring causal structure from data is subtle

Conditional independence and Markov properties

Conditional independence

- **X independent of Y given Z** if for all values of z ,

$$P(x|y, z) = P(x|z)$$

- Notation:

$$X \perp Y | Z$$

- Equivalently

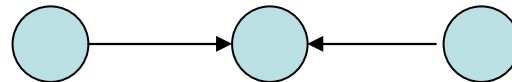
$$\begin{aligned} P(x, y|z) &= P(x|y, z)P(y|z) \\ &= P(x|z)P(y|z) \end{aligned}$$

- Conditional independence crucial in practical applications since we **can rarely work with a general joint distribution**

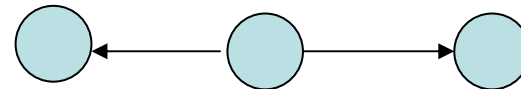
Markov properties

- Can we determine the conditional independence properties of a distribution directly from its graph?
- YES: “**d-separation**”, one subtleties due to the presence of head-to-head nodes, *explaining away effect*

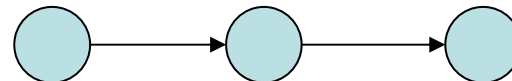
Head-to-head node



Tail-to-tail



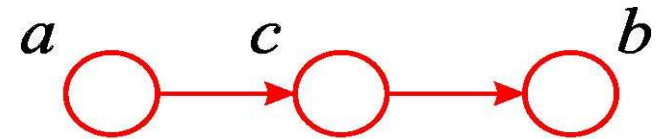
Head-to-tail



Example 1: Tail-to-head node

- Joint distribution

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



$a \not\perp b$ (c not observed)

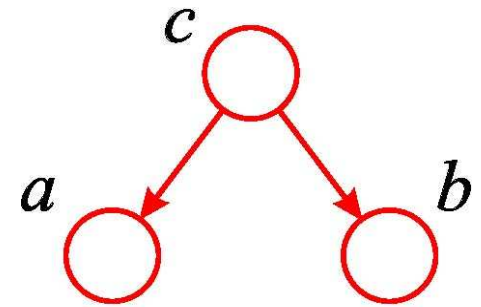
$$P(a, b|c) = P(a|c)P(b|c) \implies a \perp b|c \quad (c \text{ observed})$$

- An **observed c blocks the path** from a to b

Example 2: Tail-to-tail node

- Joint distribution

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



$a \not\perp b$ (c not observed)

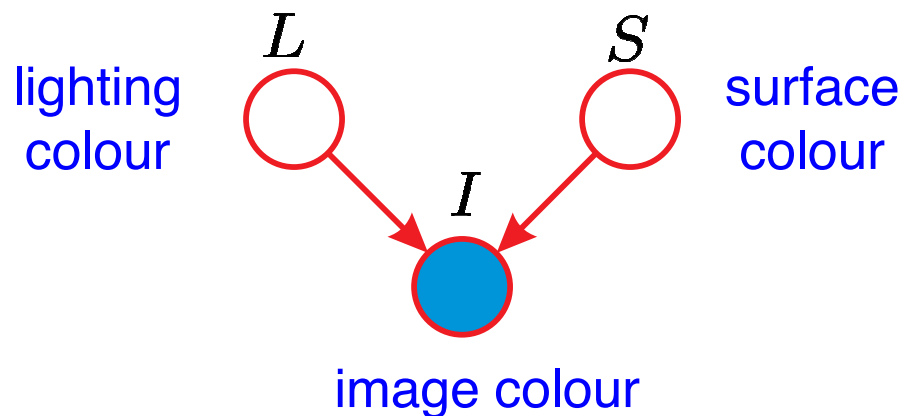
$$P(a, b|c) = P(a|c)P(b|c) \implies a \perp b|c \quad (c \text{ observed})$$

- An **observed c blocks the path** from a to b

Example 3: “Explaining Away”

Illustration: pixel colour in an image

$$p(I, L, S) = p(I|L, S)p(L)p(S)$$



$$p(L, S) = p(L)p(S)$$

$$p(L, S|I) \neq p(L|I)p(S|I)$$

An **observed I** *unblocks* the path from S to L

d-separation

- Consider 3 groups of nodes A, B, C
- To determine whether $A \perp B | C$ is true, consider all possible paths from any node in A to any node in B
- Any such **path is blocked** if there is a node X which is head-to-tail or tail-to-tail with respect to the path and X is in C

Or

if the node is head-to-head and neither the node nor any of its descendants is in C

Undirected graphs

Markov Random Fields

Undirected graphical models

- The second major class of graphical models
- Graphs specify **factorizations** of distributions and sets of conditional independence relations (**Markov properties**)
- **Markov Random Fields** or Markov network

Undirected Graphs: Factorization

- Provided $p(\mathbf{x}) > 0$ then joint distribution is product of non-negative functions over the *cliques* of the graph

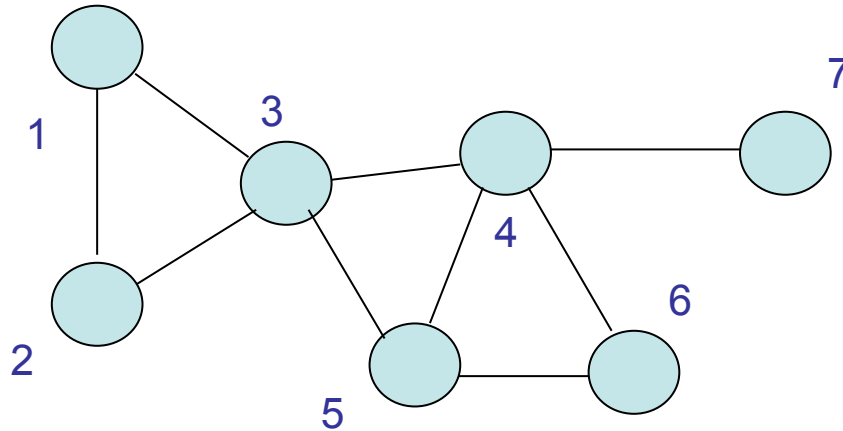
$$p(\mathbf{x}) = \frac{1}{Z} \prod_C \psi_C(\mathbf{x}_C)$$

- where $\psi_C(\mathbf{x}_C)$ are the *clique potentials*, and Z is a *normalization constant*

$$X = \{X_i, i \in V\} \quad X_C = \{X_i, i \in C\}$$

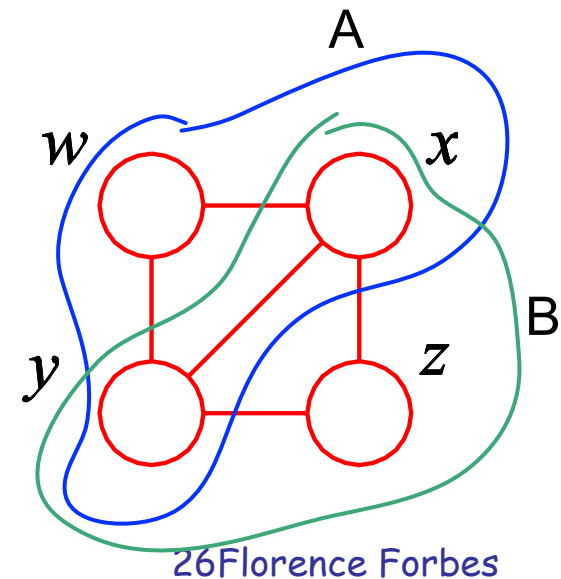
Cliques and maximal cliques

- A clique C is a subset of vertices all joined by edges



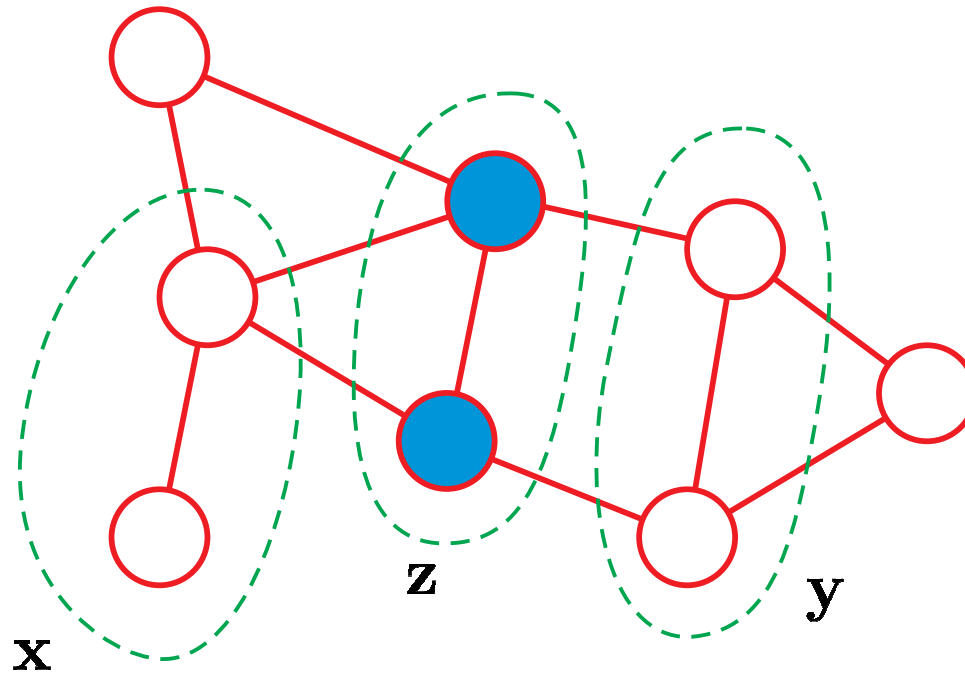
- Cliques: (1), (2),(12), (23).....
- Maximal cliques: (123), (345), (456), (47)

$$p(w, x, y, z) = \frac{1}{Z} \psi_A(w, x, y) \psi_B(x, y, z)$$



Undirected graphs: conditional independencies

- Conditional independence given by graph **separation**
x independent of y given z



Conditional independencies: Markov properties

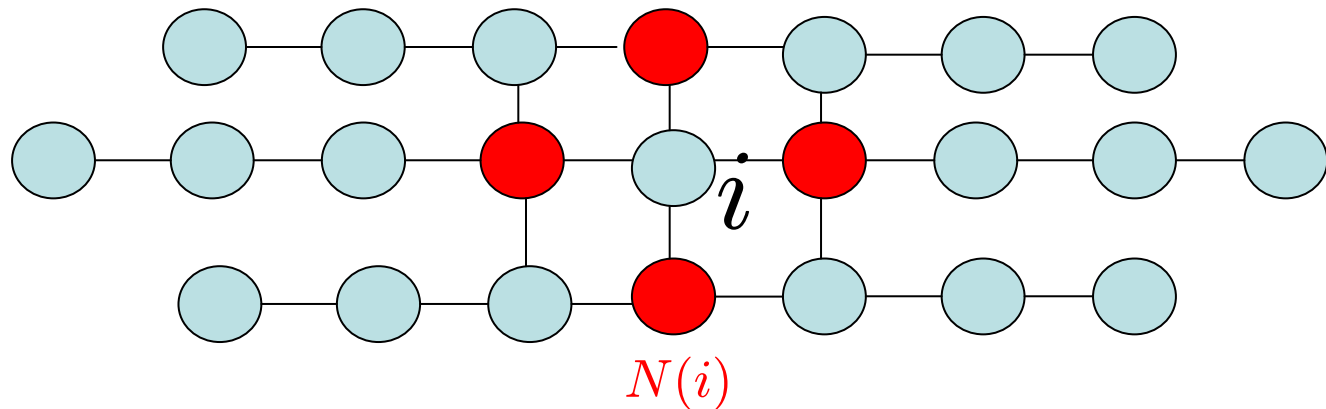
Markov blanket or Markov Boundary

of a node x_i is the set of nodes $N(i)$ such that

$$P(x_i | x_{-i}) = P(x_i | x_{N(i)})$$

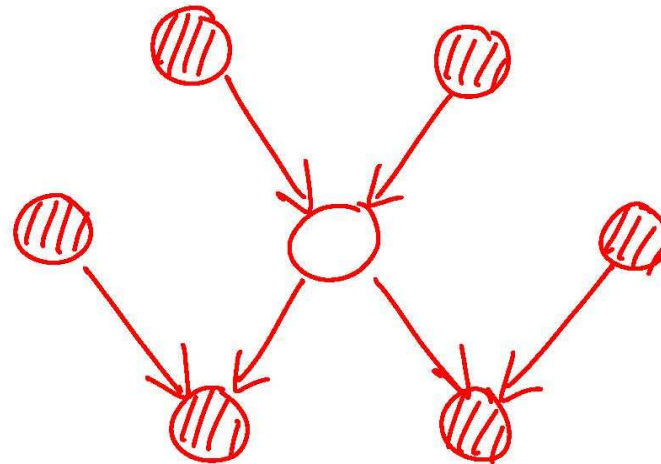
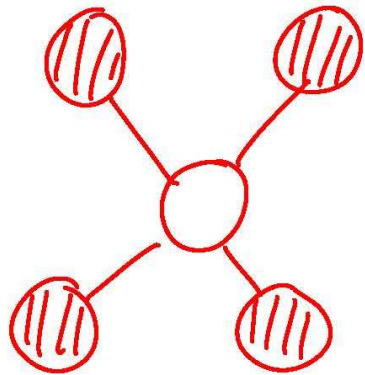
or equivalently

$$X_i \perp X_{-i \cup N(i)} | X_{N(i)}$$



Markov blankets

- Directed case: Parents, Children, Co-parents
- Undirected case: Neighbors



Markov property

- Graph $G=(V,E)$
- $X = \{X_i, i \in V\}$ random vector
- $X_A = \{X_i, i \in A\}$
- X is Markov wrt G
 - if X_A and X_B are *conditionally independent* given X_C
whenever C *separates* A and B
- Specifying conditional independencies using the neighborhood $N(i)$ is enough (V finite)

Hammersley-Clifford theorem

Makes the connection between conditional independencies (Markov properties) and factorization property

- Boltzmann-Gibbs representation

$$\Psi_c(x_c) = \exp(-E(x_c))$$

- P is a **positive MRF** (satisfies Markov properties) is equivalent to **P is a Gibbs distribution**

$$P(x) = \frac{1}{Z} \exp(-E(x))$$

- Energy function

$$E(x) = \sum_c E_c(x_c)$$

Example: pairwise Markov Random Fields

- Cliques: pairs, singletons

$$E(x) = \sum_i \{ \Psi_i(x_i) + \frac{1}{2} \sum_{j \in N(i)} \Psi_{ij}(x_i, x_j) \}$$

- Famous ones:
 - **Ising** model: **binary** variables on a graph G with pairwise interactions

$$P(x; \theta) = \frac{1}{Z} \exp\left(\sum_i \theta_i x_i + \sum_{i \sim j} \theta_{ij} x_i x_j\right)$$

- **Potts** model: **K-ary** variables

Interaction parameters+ external field parameters

Example: graph representation of a Pairwise MRF

- Typical application: image region labelling

$$p(\mathbf{x}, \mathbf{y}) = \frac{1}{Z} \prod_i \phi_i(x_i, y_i)$$

$$\prod_{i,j} \psi_{ij}(x_i, x_j)$$

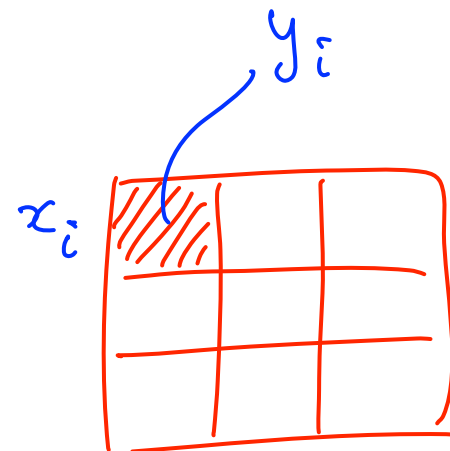
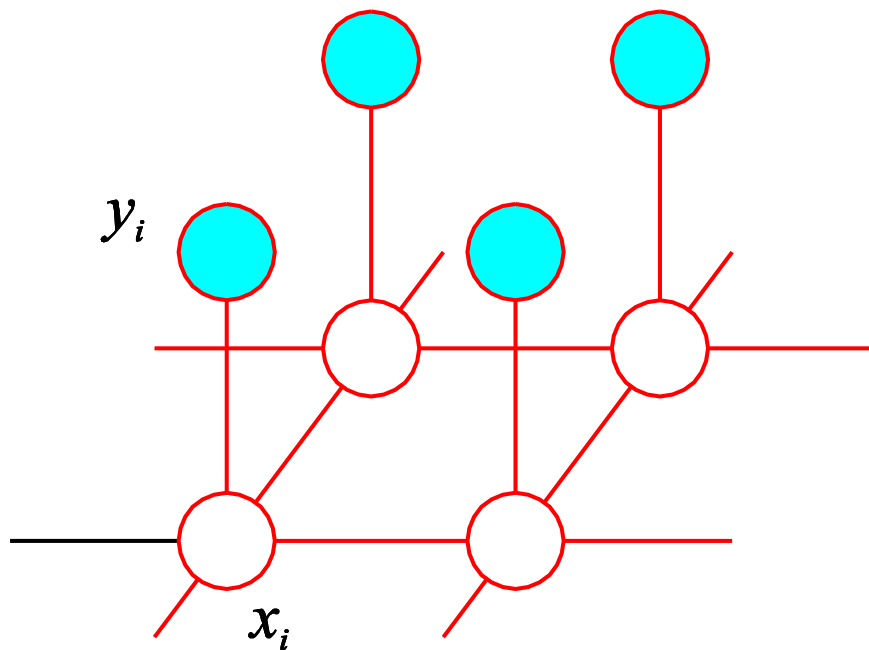
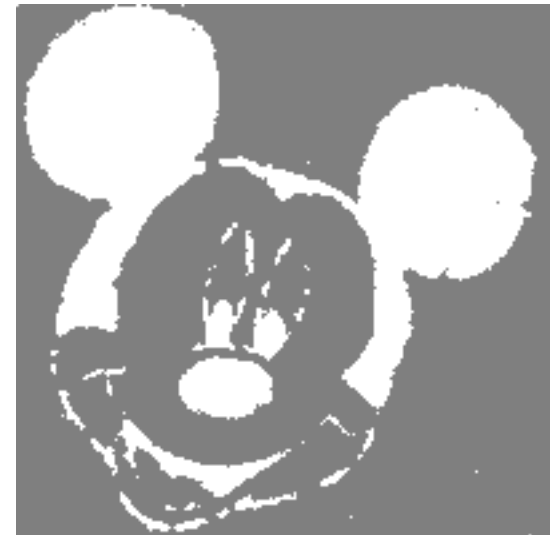
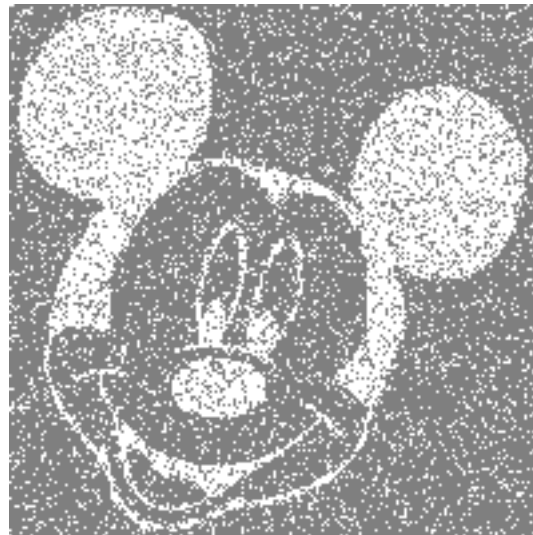
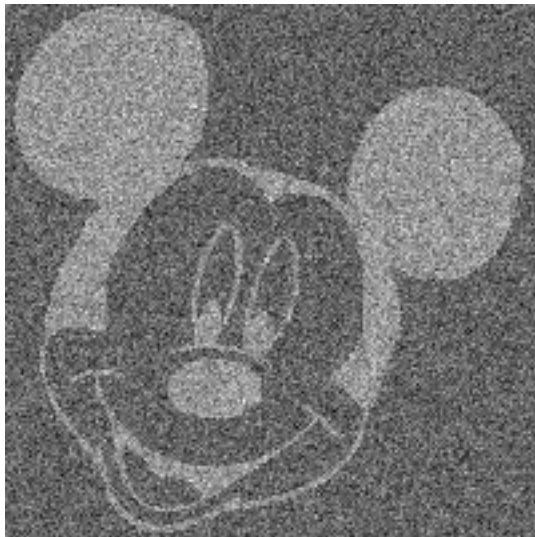
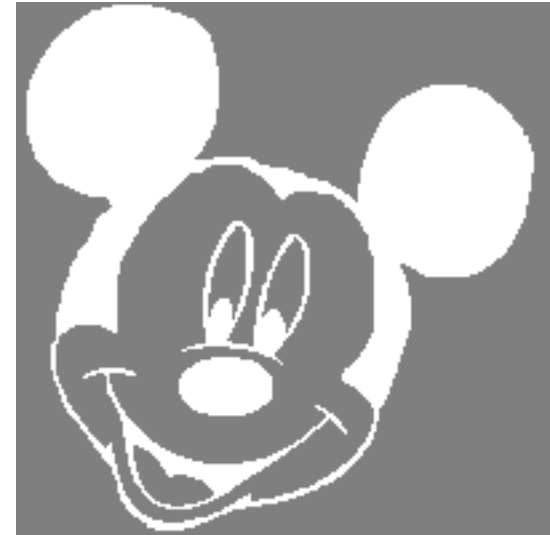


Illustration: image segmentation

site/vertex i : pixel,

Y_i : observed grey level,

X_i : label/0 or 1/ binary variable



Challenging computational problems

- Frequently, it is of interest to compute various quantities associated with an undirected graphical model:
 - The log normalization constant $\log Z$
 - Local marginal distributions ($p(x_i)$) or other local statistics
 - Modes and most probable configurations
- Often grow rapidly with graph size and max clique size
- Example: Computing the normalization constant for binary random variables

$$Z = \sum_{x \in \{0,1\}^n} \prod_{c \in C} \psi_c(x_c)$$

Complexity scales exponentially as 2^n

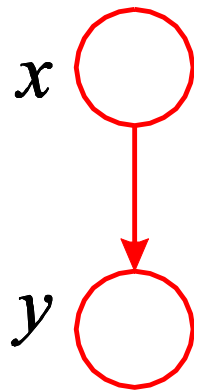
Inference and learning

Inference in Graphical models

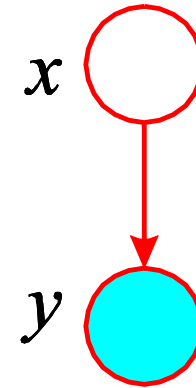
- Exploit the graphical structure to find efficient algorithm for inference and to make the structure of these algorithms clear (eg propagation of local messages around the graph)
- Exact inference
- Approximate inference

Inference

- Simple example: Bayes' theorem



$$p(x, y) = p(x)p(y|x)$$

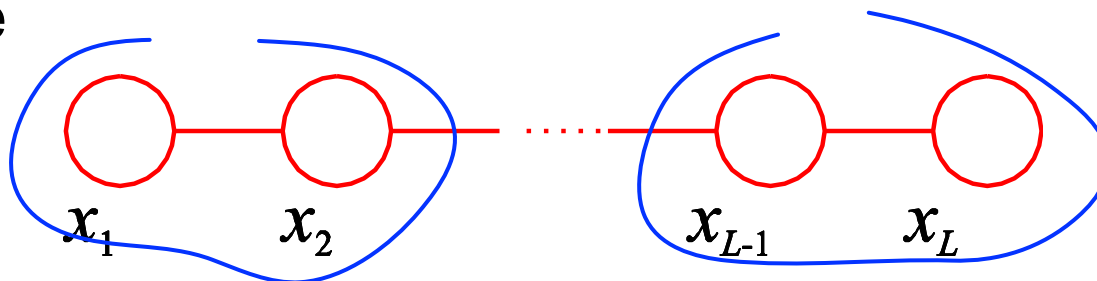


$$p(x|y) = \frac{p(x)p(y|x)}{p(y)}$$

$$\sum_x p(x) p(y|x)$$

Message Passing: compute marginals

- Example



- Find marginal for a particular node

$$p(x_i) = \sum_{x_1} \cdots \sum_{x_{i-1}} \sum_{x_{i+1}} \cdots \sum_{x_L} p(x_1, \dots, x_L)$$

- for M-state nodes, cost is $O(M^L)$
- exponential in length of chain
- but, we can exploit the graphical structure (conditional independences)

Message Passing

- Joint distribution

$$p(x_1, \dots, x_L) = \frac{1}{Z} \psi(x_1, x_2) \dots \psi(x_{L-1}, x_L)$$

- Exchange sums and products: $ab + ac = a(b+c)$

$$p(x_i) = \frac{1}{Z} \cdots \sum_{x_2} \psi(x_2, x_3) \left[\sum_{x_1} \psi(x_1, x_2) \right] \cdots \sum_{x_{L-1}} \psi(x_{L-2}, x_{L-1}) \left[\sum_{x_L} \psi(x_{L-1}, x_L) \right]$$

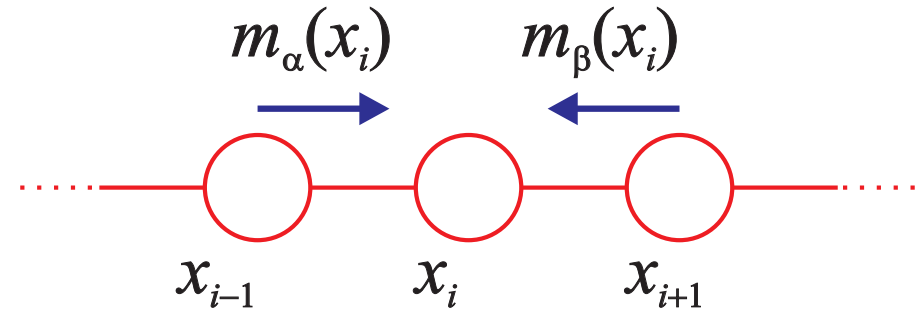
$m_\alpha(x_i)$ before x_i

$m_\beta(x_i)$ after x_i

Message Passing

- Express as product of messages

$$p(x_i) = \frac{1}{Z} m_\alpha(x_i) m_\beta(x_i)$$



- Recursive evaluation of messages: Linear in L

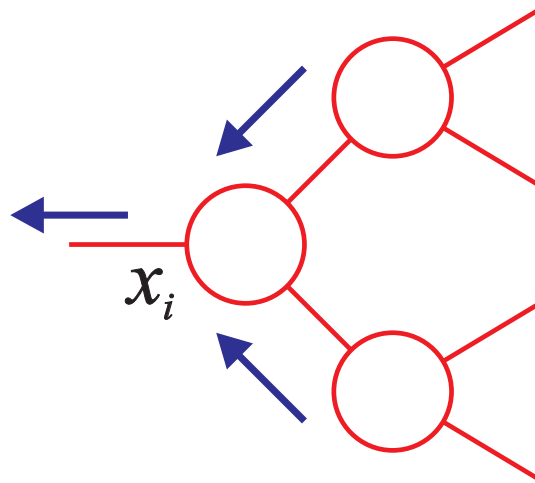
$$m_\alpha(x_i) = \sum_{x_{i-1}} \psi(x_{i-1}, x_i) m_\alpha(x_{i-1})$$

$$m_\beta(x_i) = \sum_{x_{i+1}} \psi(x_i, x_{i+1}) m_\beta(x_{i+1})$$

- Find Z by normalizing $p(x_i)$

Belief Propagation

- Extension to general tree-structured graphs
- At each node:
 - form product of *incoming* messages and local evidence
 - marginalize to give *outgoing* message
 - one message in each direction across every link



- Fails if there are loops

Junction Tree Algorithm

- An efficient exact algorithm for a general graph
 - applies to both directed and undirected graphs
 - compile original graph into a tree of cliques
 - then perform message passing on this tree
- Problem:
 - cost is exponential in size of largest clique
 - many vision models have intractably large cliques

Loopy Belief Propagation

- Apply belief propagation directly to general graph
 - possible because message passing rules are local
 - need to keep iterating
 - might not converge
- State-of-the-art performance in some applications

Max-product Algorithm: most probable x

- Goal: find

$$\mathbf{x}^{\text{MAP}} = \arg \max_{\mathbf{x}} p(\mathbf{x})$$

- define

$$\phi(x_i) = \max_{x_1} \cdots \max_{x_{i-1}} \max_{x_{i+1}} \cdots \max_{x_L} p(x_1, \dots, x_L)$$

- then

$$x_i^{\text{MAP}} = \arg \max_{x_i} \phi(x_i)$$

- Message passing algorithm with “sum” replaced by “max”
- Example:
 - Viterbi algorithm for HMMs

Inference and learning

In general: Hidden or latent X (underlying scene) and Observed Y (image)

- Inference: computing $P(x|y)$ (“posterior”)
- Learning: computing $P(y)$ (likelihood) usually $P_{\theta}(y)$
(θ : parameter estimation based on ML)

Likelihood of the data y $L(\theta) = P_{\theta}(y)$

Maximum (log) likelihood

$$\theta_{ML} = \arg \max_{\theta} \log L(\theta)$$

Example: classification with context

- The labeling problem

- ★ n objects/individuals ($i \in V = \{1, \dots, n\}$)
- ★ K labels ($k \in \mathcal{L} = \{1, \dots, K\}$)
- ★ $n * \dots$ observations ($y = (y_1, y_2, \dots)$)

assign a label to each object consistently with y :

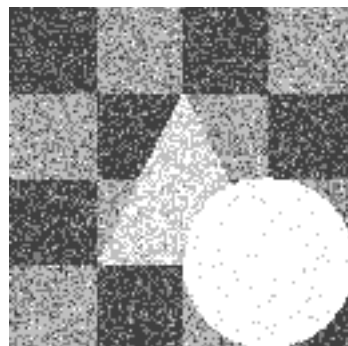
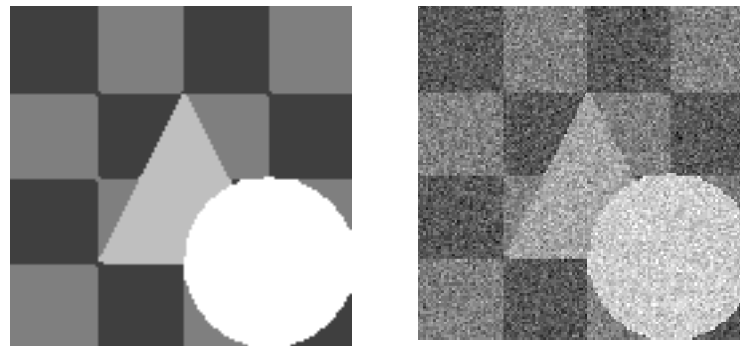
$$\mathbf{x} : V \rightarrow \mathcal{L}$$

$$x = (x_1, \dots, x_n \in \mathcal{L}^n)$$

(assignement, colouring (graph), configuration (random fields))

Contextual constraints: distance, similarity, compatibility, etc.

- Image analysis, segmentation, etc.
- Biometrics: spatially related observations
- Documents analysis: hyperlinks between documents



No context



Too much context



Good compromise

Assignment criterion: $x : V \longrightarrow \mathcal{L}$

★ assignment cost

$c(i, k)$ [likelihood of k at site i] or $c_y(i, k)$ [data term]

★ Neighborhood cost:

i and j nearby $\Rightarrow x_i$ and x_j similar/compatible

\rightarrow graph $G = (V, E)$: if $(i, j) \in E$

\rightarrow cost $w_{ij} \times d_{ij}(x_i, x_j)$ $[\Psi_{ij}(x_i, x_j)]$

Total cost:
$$E(x) = \sum_{i \in S} c(i, x_i) + \sum_{(i, j) \in E} w_{ij} d_{ij}(x_i, x_j)$$

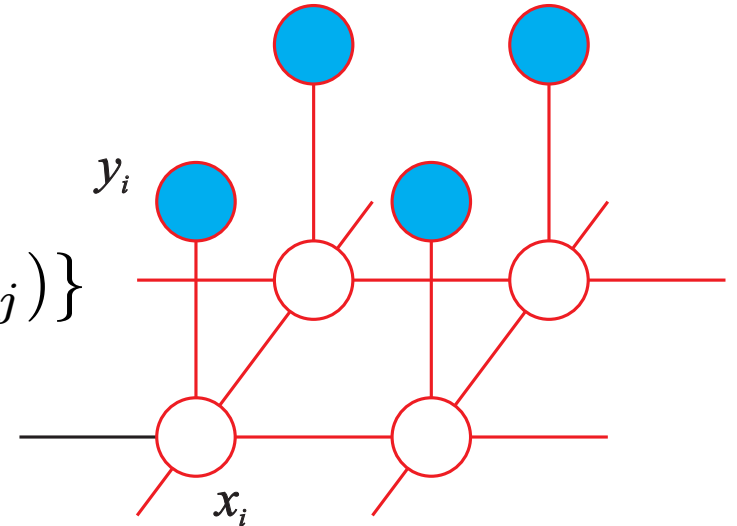
- **Goal: find x that maximizes E**
- Discrete optimization, **NP-hard, find approximations, satisfying assignments**

Optimal configuration for Pairwise MRF with energy E

Markovian approach and MAP rule

- Corresponding graphical model: Pairwise MRF

$$E(x) = \sum_i \{ \Psi_i(x_i) + \frac{1}{2} \sum_{j \in N(i)} \Psi_{ij}(x_i, x_j) \}$$



- Maximum A Posteriori (MAP) principle:

$$\hat{x} = \arg \max_{x \in \mathcal{L}^n} P(x|y)$$

Hidden MRF: accounting for observations

- Observations, eg. Measures $Y = \{Y_i, i \in S\}$
- Hidden data, eg. Labels, X discrete MRF

$$P(x) = \frac{1}{Z} \exp(-E(x))$$

- **Data term,**

$$P(y|x) = \exp(-E(y|x))$$

Conditional MRF (posterior):

$$P(x|y) = \frac{1}{Z_y} \exp(-E_y(x))$$

$$E_y(x) = E(x) + E(y|x)$$

E(x): Regularizing term (prior, context)

E(y | x): Data term

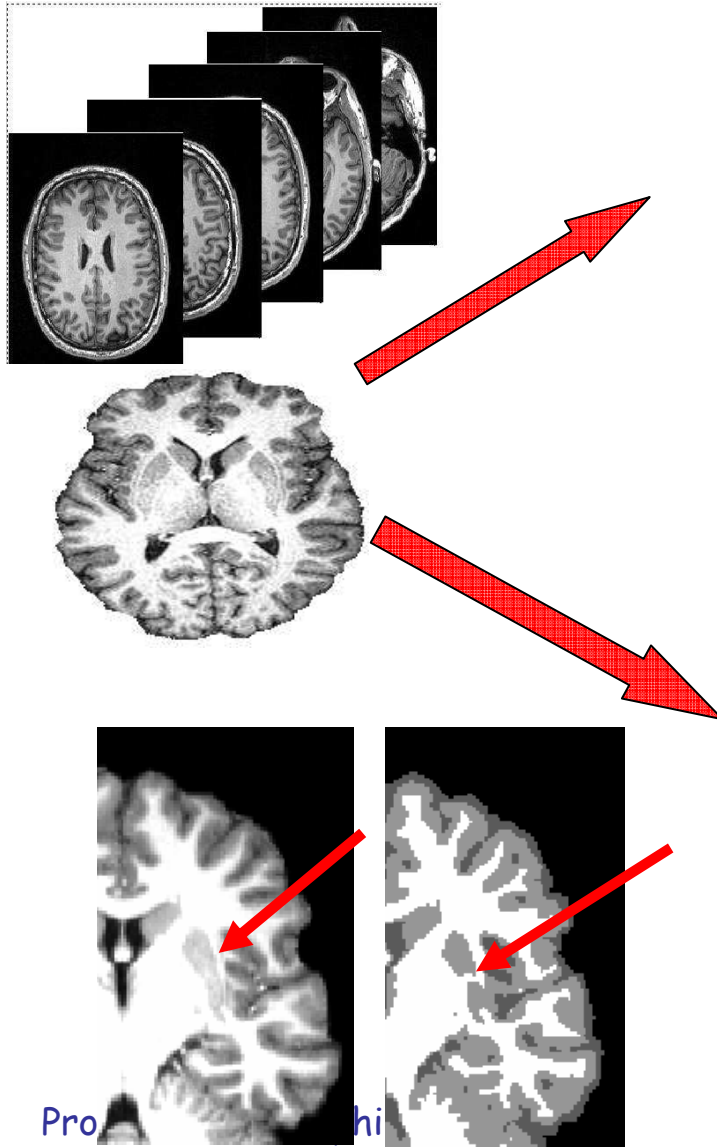
$$\text{MAP solution } \hat{x} = \arg \min_{x \in \mathcal{L}^n} E_y(x)$$

Approximate solutions

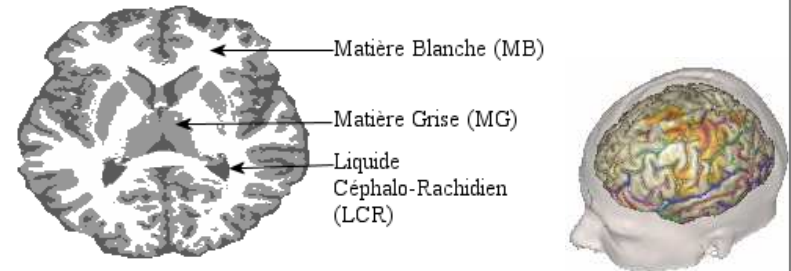
- Deterministic approaches: relaxation, variational methods (mean field, etc.)
- Stochastic approaches: Gibbs sampling, simulation methods (MC)
- Classification approaches: hard clustering, ICM, K-means
- Parameter estimation approaches: soft clustering, EM

Example 1: MRI Brain scan segmentation

Assign each voxel to a class (label) (among K classes)



Tissue segmentation (WM, GM, CSF)



→ Cortex 3D reconstruction

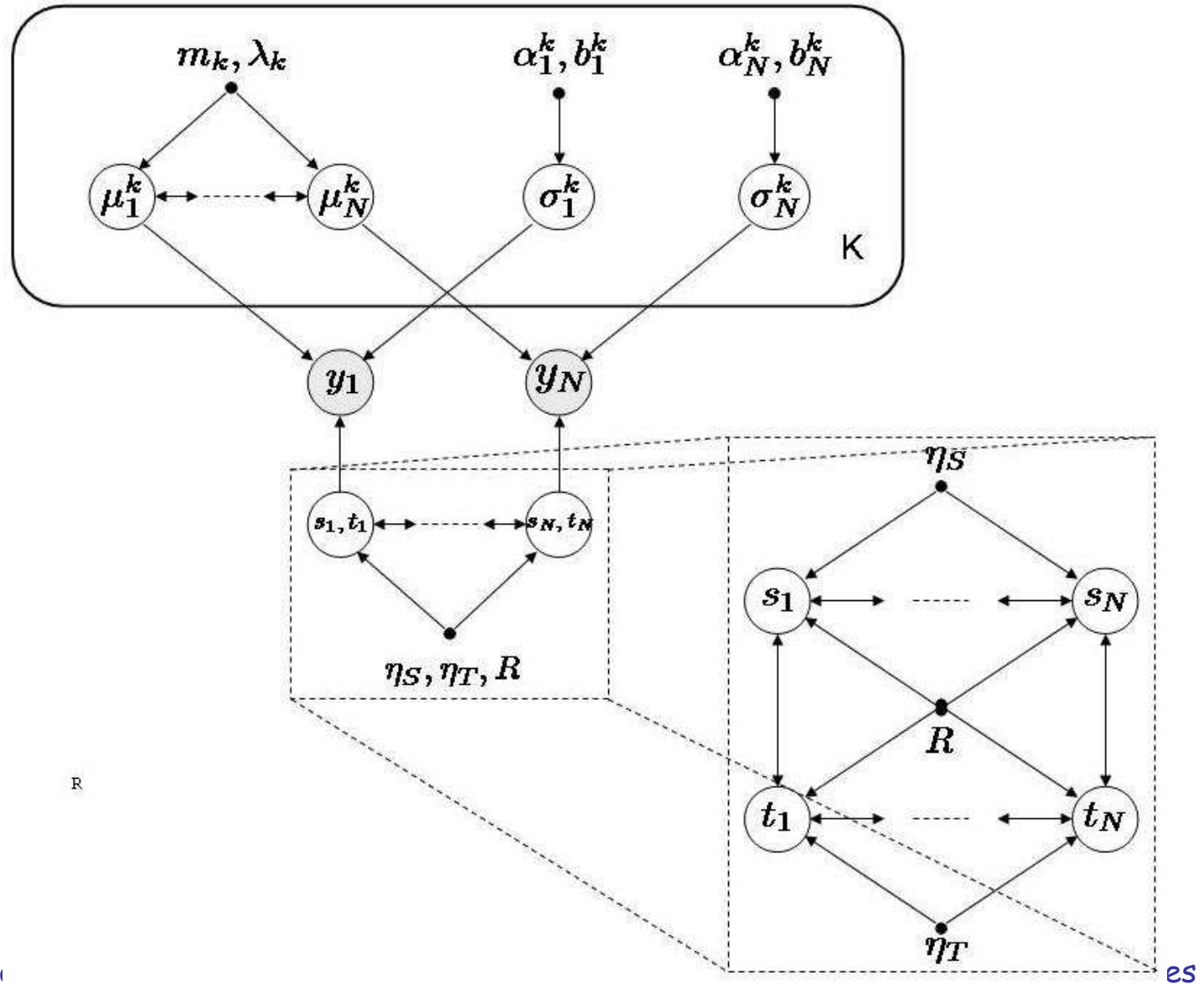
Structure segmentation



→ Useful for :

- Distinguishing Cortex GM from Nuclei GM
- volumetric studies
- ...

Graphical model representation

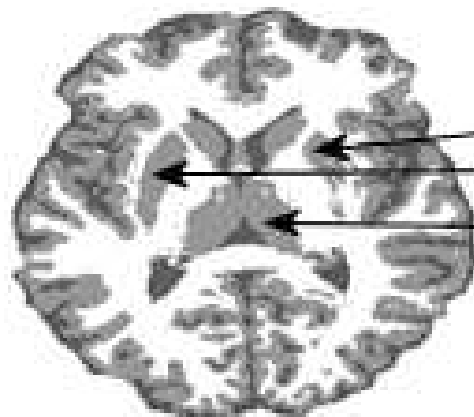


Cooperative segmentation of tissues and structures

observations



No anatomical information

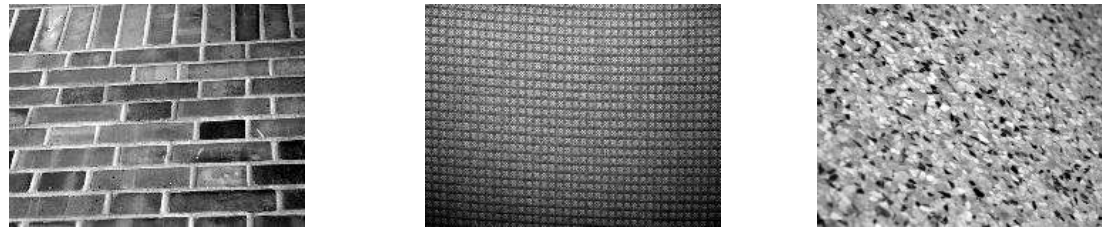


Meilleure segmentation
incontestable des putamens
et des thalamus

Cooperative method

Example 2: texture recognition

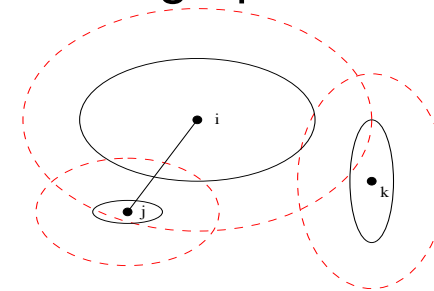
- Learning step: model estimation



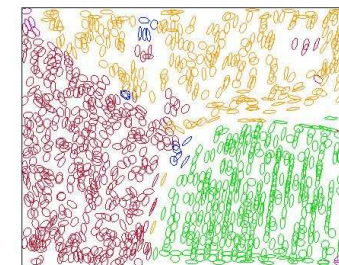
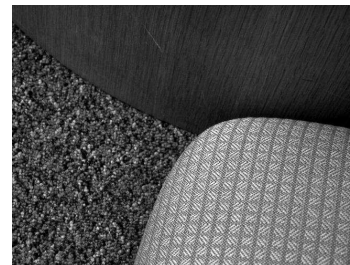
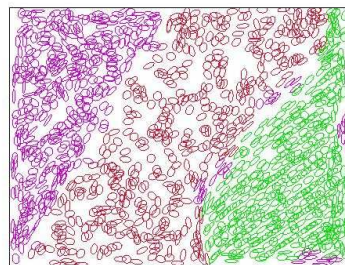
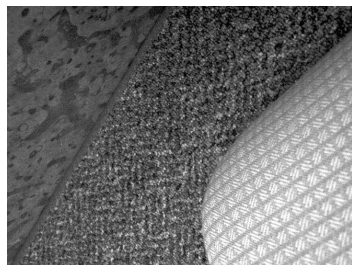
- Interest points



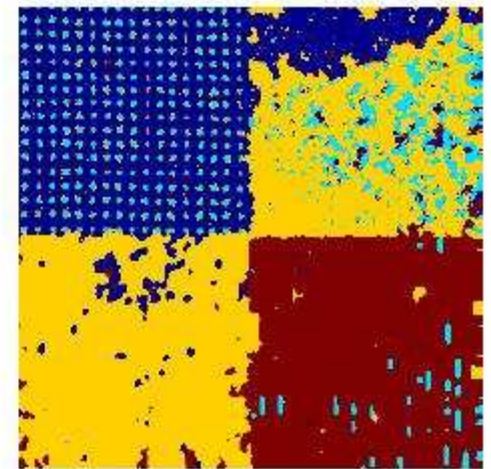
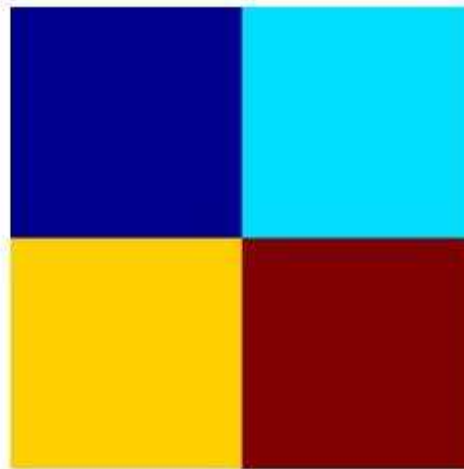
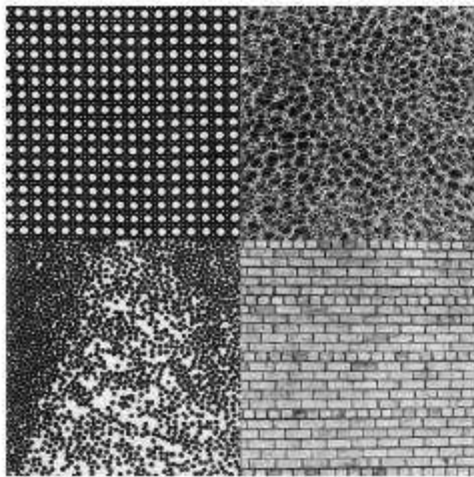
neighborhood graph



- Test step: classification



Example 2: texture recognition



The End

- Thank you for your attention