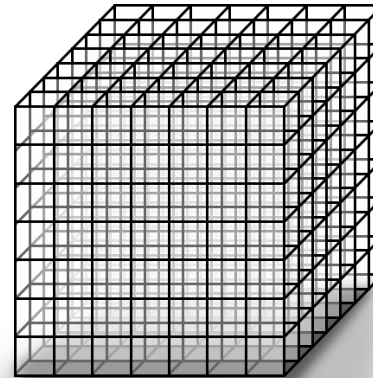
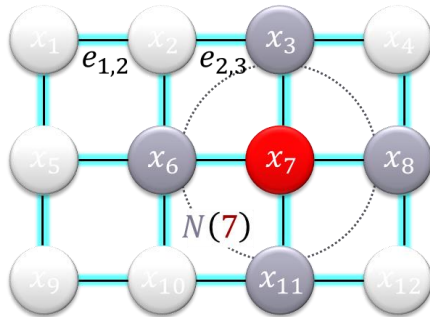


Machine Learning in Scientific Computing

CECAM/CSM/IRTG SCHOOL 2018



Lecture 3.1.2

Markov Random Fields

Markov Random Fields and Graphical Models

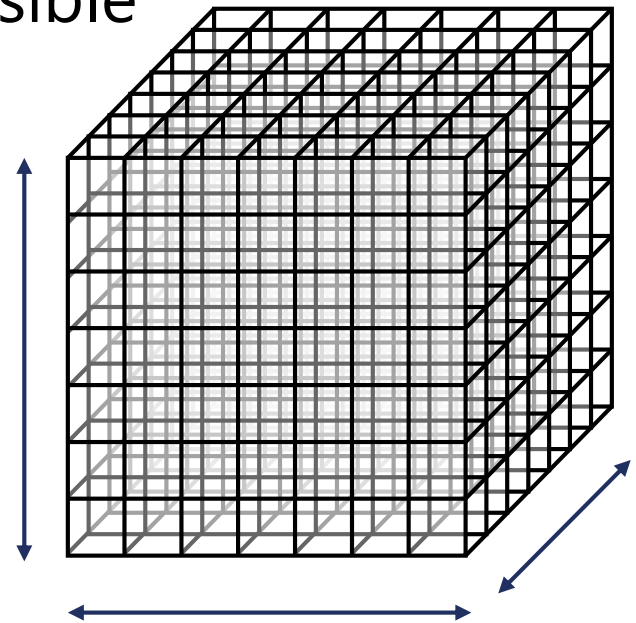
Reducing dependencies

Problem:

- $p(x_1, x_2, \dots, x_n)$ is too high-dimensional
- k States, n variables: $O(k^n)$ density entries
- General dependencies often infeasible

Reduction of Dependencies

- Parametric models (Gauss etc.)
- Markov Random Fields (MRFs)
- Deep Networks
- ...



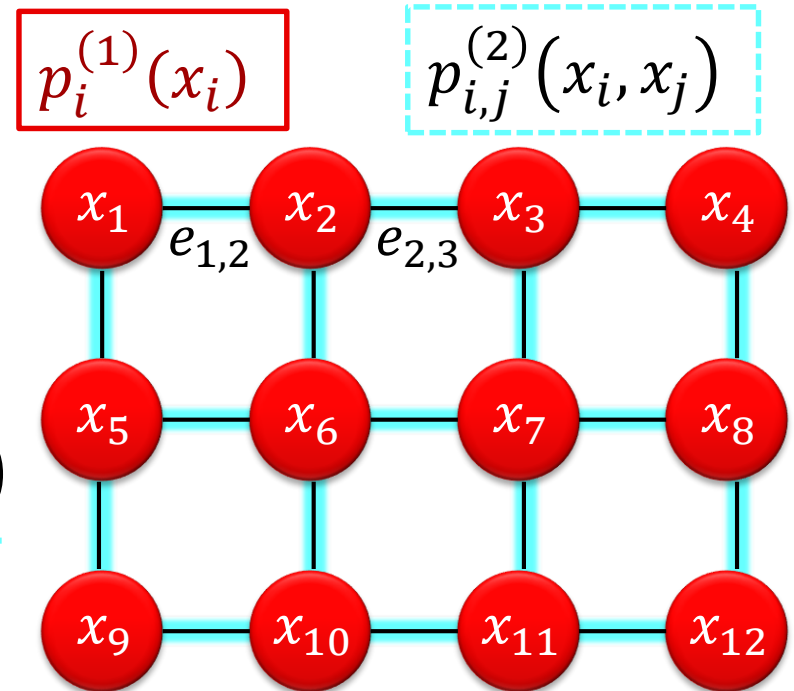
Graphical Models

Factorize Models

- Pairwise models:

$$p(x_1, \dots, x_n) = \frac{1}{Z} \prod_{i=1}^n \underbrace{p_i^{(1)}(x_i)} \prod_{i,j \in E} \underbrace{p_{i,j}^{(2)}(x_i, x_j)}$$

- Model complexity:
 - $O(nk^2)$ parameters
- Higher order models:
 - Triplets, quadruples as factors
 - Local neighborhoods



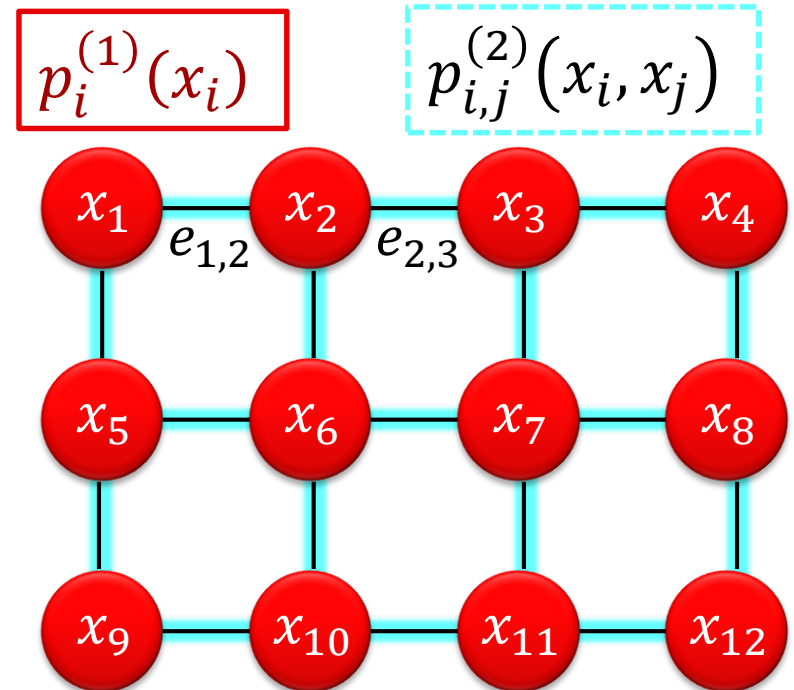
Graphical Models

Markov Random fields

- Factorize density in local "cliques"

Graphical model

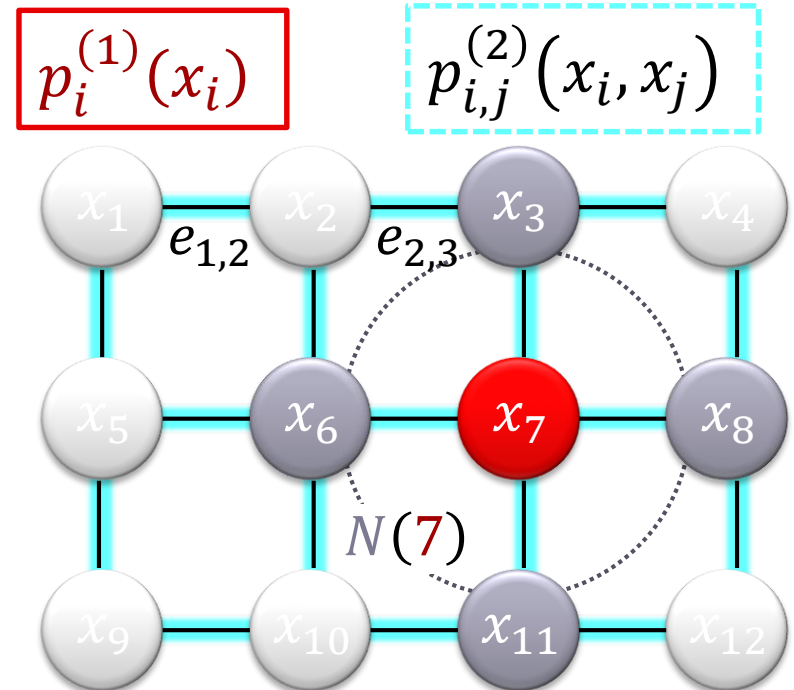
- Connect variables that are directly dependent
- Formal model:
Conditional independence



Graphical Models

Conditional Independence

- A node is conditionally independent of all others given the values of its direct neighbors
- I.e. set these values to constants, x_7 is independent of all others



Formally

- $p(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) = p(x_i | \{x_j | j \in N(i)\})$

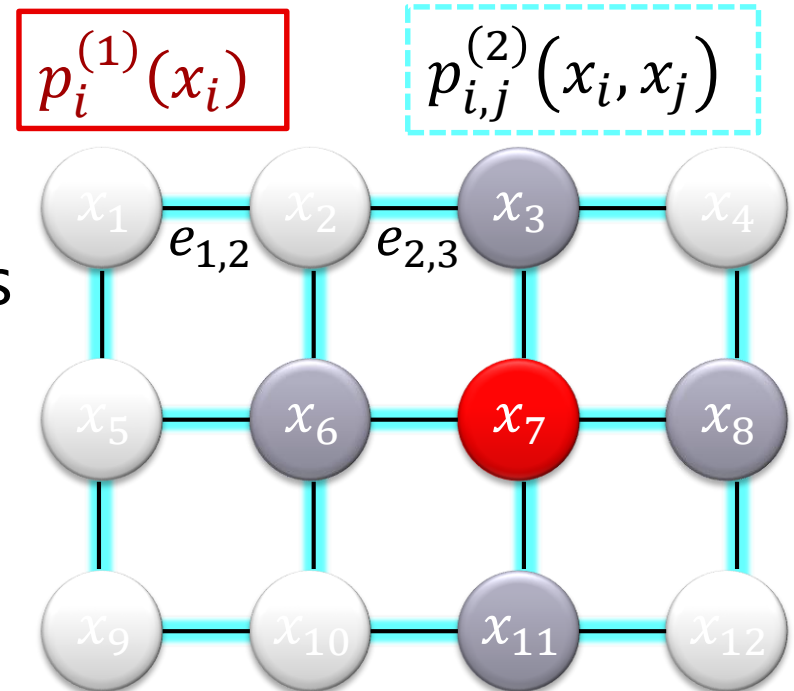
Graphical Models

Theorem (Hammersley–Clifford):

- Assuming positive densities $p(x_i) > 0$

The theorem

- Given conditional independence as graph, density factors over cliques in the graph.
- And vice versa.

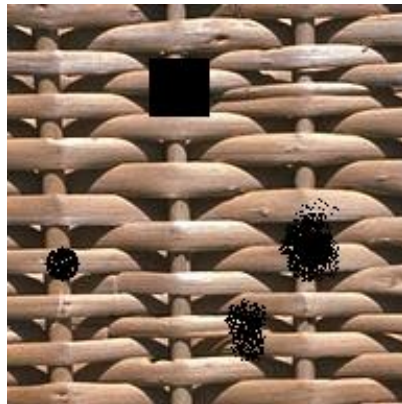


$$p(\mathbf{x}) = \frac{1}{Z} \prod_{i=1}^n \underbrace{p_i^{(1)}(x_i)}_{\text{red}} \prod_{i,j \in E} \underbrace{p_{i,j}^{(2)}(x_i, x_j)}_{\text{cyan}}$$

Example: Texture Synthesis



original



region selected

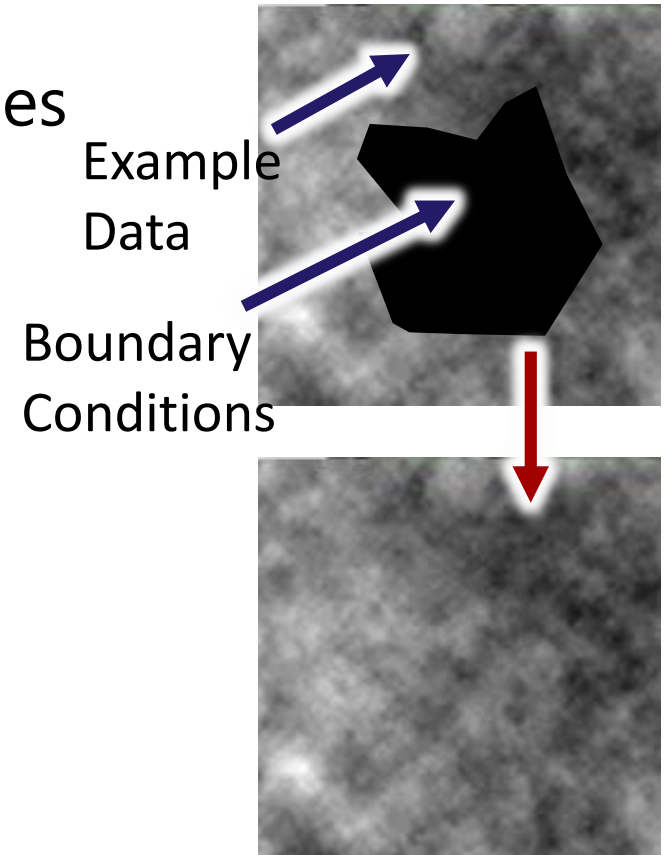


completion

Texture Synthesis

Idea

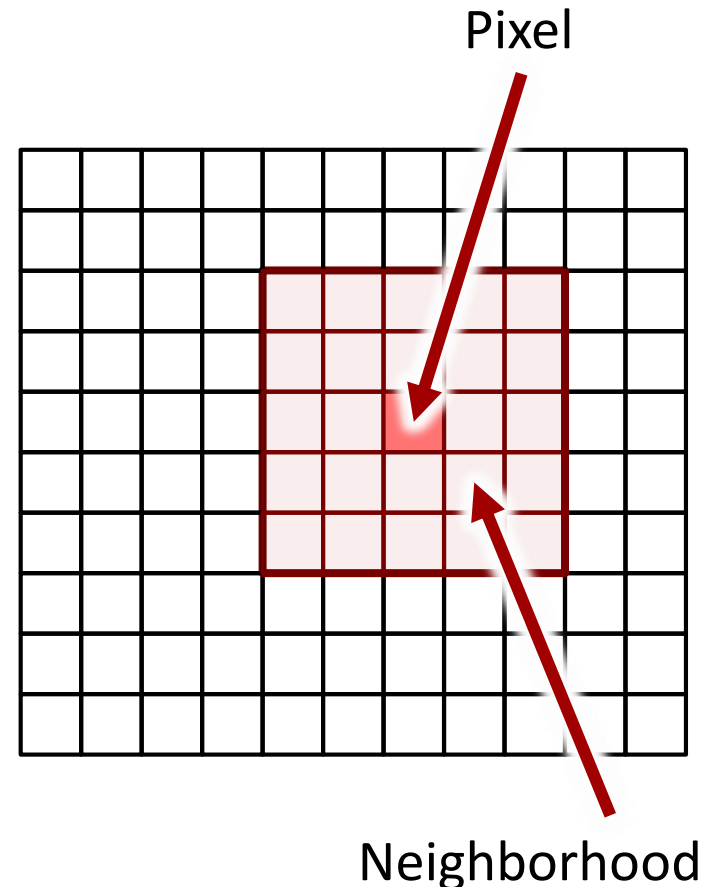
- One or more images as examples
- Learn image statistics
- Use knowledge:
 - Specify boundary conditions
 - Fill in texture



The Basic Idea

Markov Random Field Model

- Image statistics
- How pixels are colored depends on local neighborhood only (Markov Random Field)
- Predict color from neighborhood

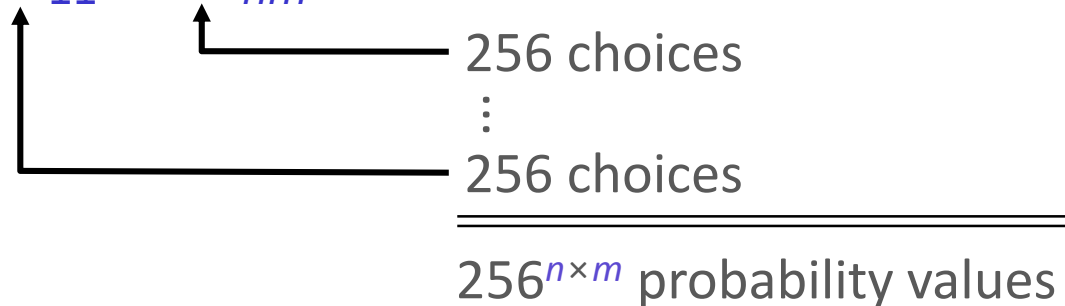


A Little Bit of Theory...

Image statistics:

- An image of $n \times m$ pixels
- Random variable: $\mathbf{x} = [x_{11}, \dots, x_{nm}] \in [0, 1, \dots, 255]^{n \times m}$
- Probability distribution:

$$p(\mathbf{x}) = p(x_{11}, \dots, x_{nm})$$



Impossible to learn full images from examples!

Simplification

Problem:

- Statistical dependencies
- Simple model can express dependencies on all kinds of combinations

Markov Random Field:

- Each pixel is *conditionally independent* of the rest of the image given a small neighborhood
- In English: likelihood only depends on neighborhood, not rest of the image

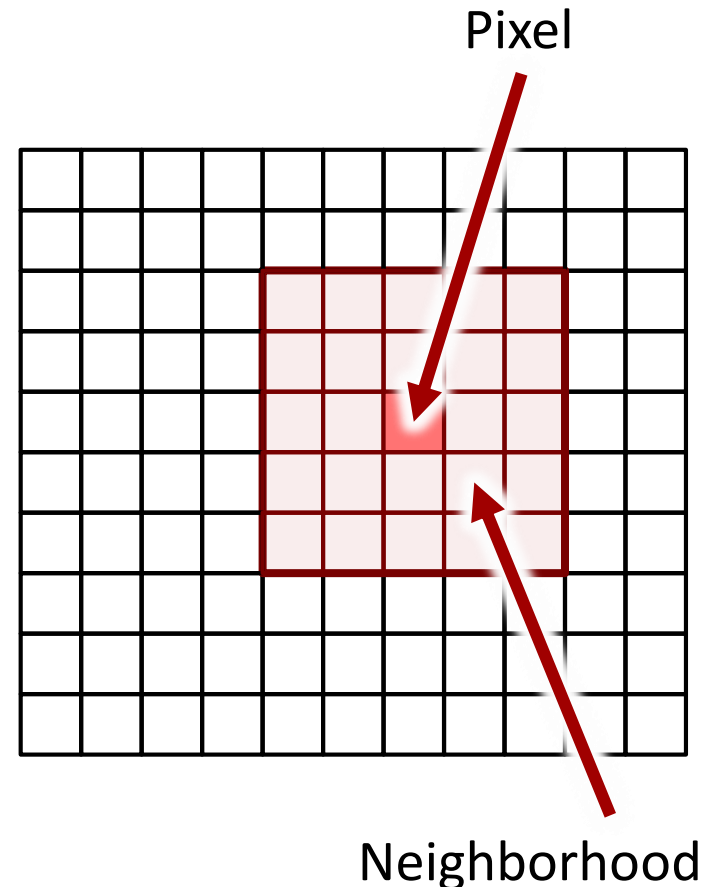
Markov Random Field

Example:

- Red pixel depends on light red region
- Not on black region
- If region is known, probability is fixed and independent of the rest

However:

- Regions overlap
- Indirect global dependency

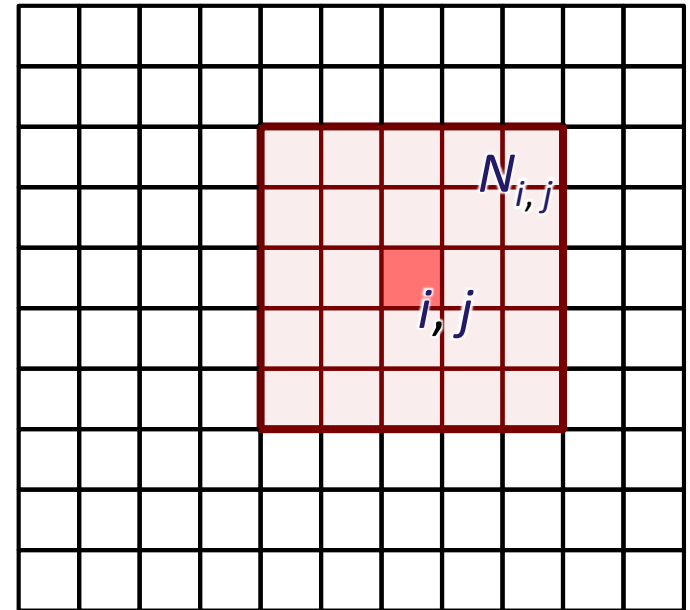


Texture Synthesis

Use for Texture Synthesis

$$p(\mathbf{x}) = \frac{1}{Z} \prod_{i=1}^n \prod_{j=1}^m p_{i,j}(N_{i,j})$$

$$\begin{aligned} p_{i,j} &= p_{i,j}(N_{i,j}) \\ &= p_{i,j}(x_{i-k,j-k} \dots, x_{i+k,j+k}) \end{aligned}$$



Inference

Inference Problem

- Computing $p(\mathbf{x})$ is trivial for known \mathbf{x} .
- Finding the \mathbf{x} that maximizes $p(\mathbf{x})$ is very complicated.
- In general: NP-hard
- No efficient solution known (not even for images)

In practice

- Different approximation strategies
("heuristics", strict approximation is also NP-hard)

Simple Practical Algorithm

Here is the short story:

- Unknown pixels:
consider known neighborhood
- Match to all of the known data
- Copy the pixel with the best matching neighborhood
- Region growing, outside in

Approximation only

- Can run into bad local minima

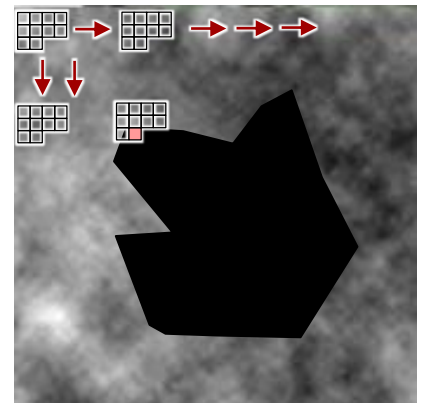
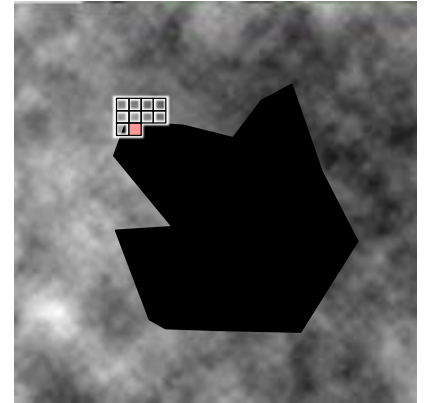


Image Analogies (Siggraph 2001)

CRF Image Segmentation

Reference (Image)

X. He, R.S. Zemel, M.A. Carreira-Perpinan:
Multiscale Conditional Random Fields for Image Labeling,
IEEE CVPR 2004.

Example: Weak Formulations of Differential Equations

Differential Equations

Example equation

$$\frac{d}{dt} f(t) = F(f(t), t)$$

Discretization

$$\frac{y_i - y_{i-1}}{h} = F(y_i, t_i)$$

Weak formulation (variational approach)

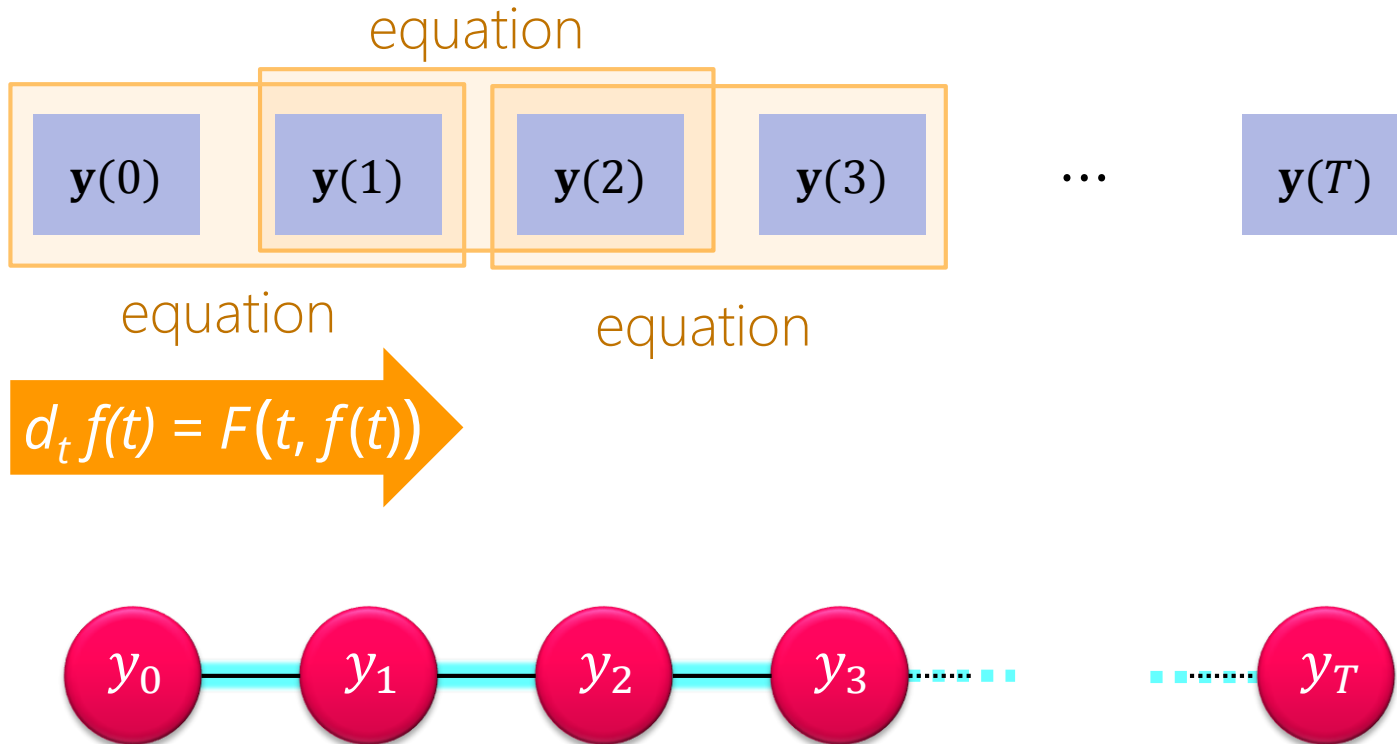
$$\left(\frac{y_i - y_{i-1}}{h} - F(y_i, t_i) \right)^2 \rightarrow \min$$

$$\arg \max_{\mathbf{y}} \frac{1}{Z} \prod_{i=1}^{n-1} \exp \left(\frac{y_i - y_{i-1}}{h} - F(y_i, t_i) \right)^2$$

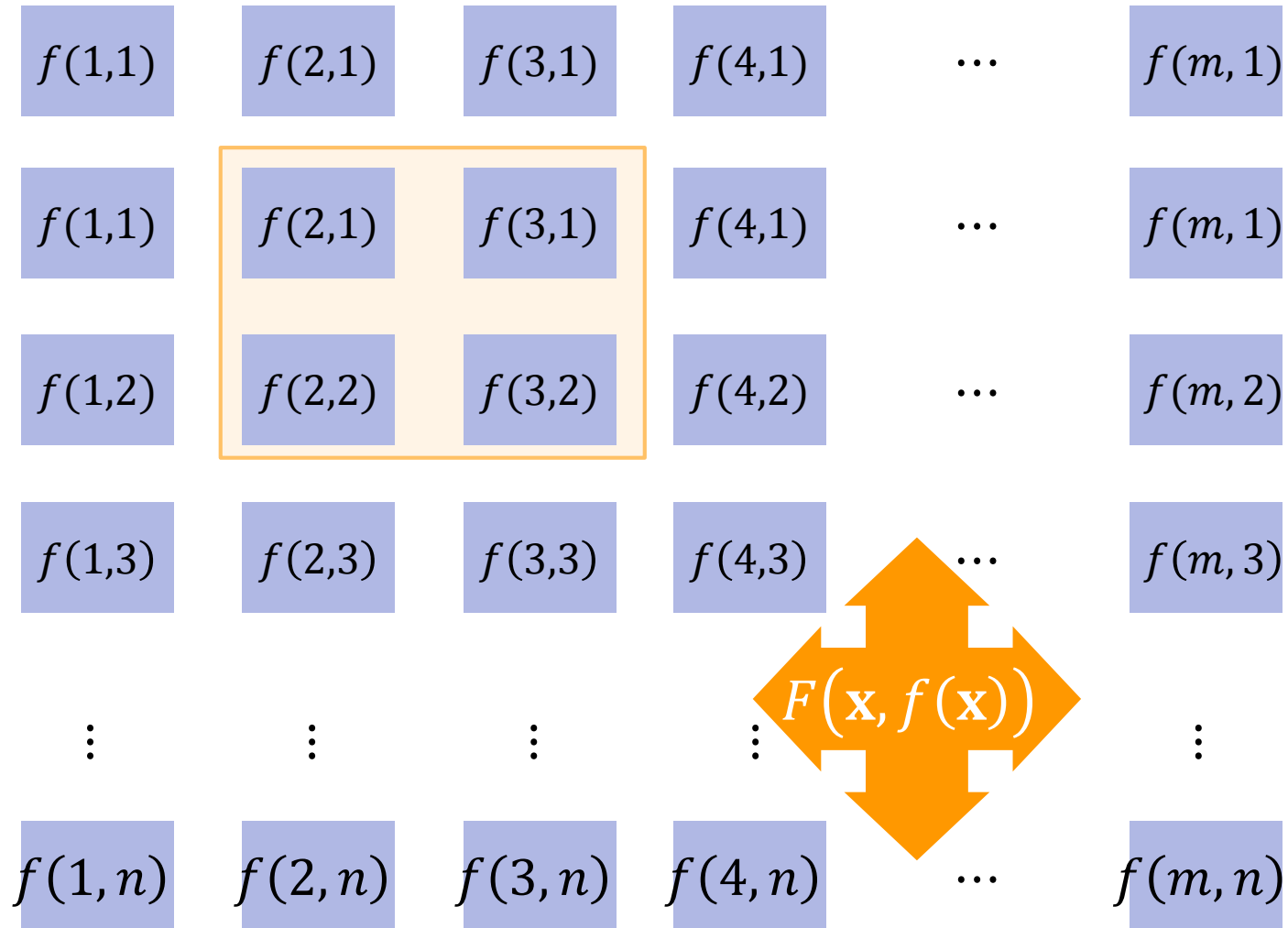
Gaussian MRF

Ordinary Differential Equations

Causal Chain



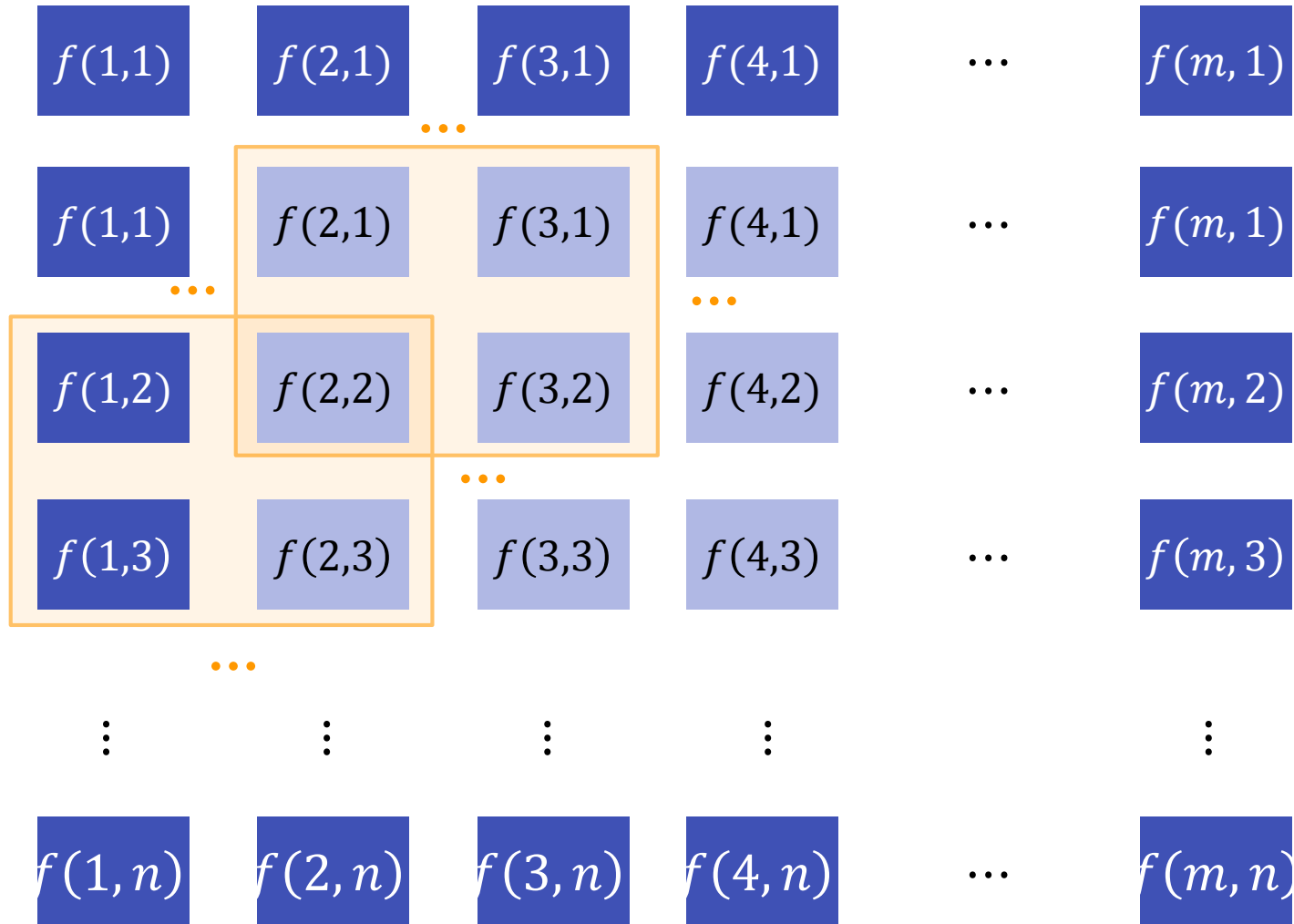
Structure of PDE



Structure of PDE

$f(1,1)$	$f(2,1)$	$f(3,1)$	$f(4,1)$...	$f(m, 1)$
$f(1,1)$	$f(2,1)$	$f(3,1)$	$f(4,1)$...	$f(m, 1)$
$f(1,2)$	$f(2,2)$	$f(3,2)$	$f(4,2)$...	$f(m, 2)$
$f(1,3)$	$f(2,3)$	$f(3,3)$	$f(4,3)$...	$f(m, 3)$
⋮	⋮	⋮	⋮		⋮
$f(1,n)$	$f(2,n)$	$f(3,n)$	$f(4,n)$...	$f(m,n)$

Boundary Value Problem



Inference in MRFs

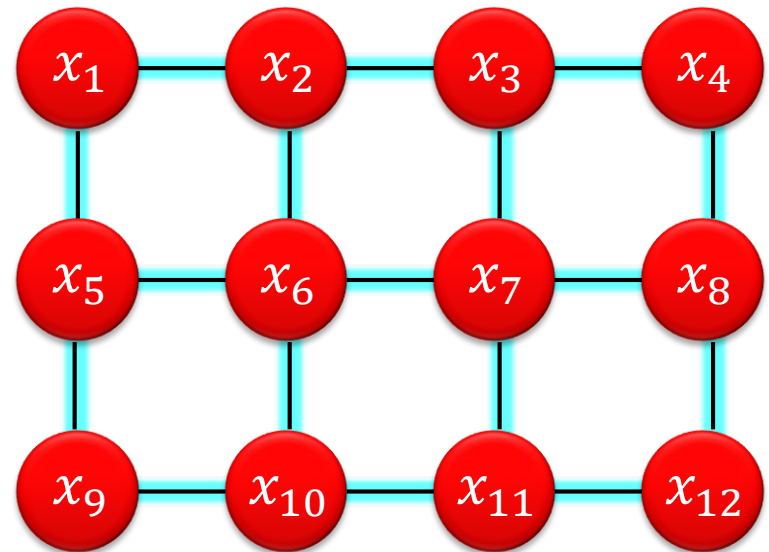
Inference

Model

$$p(x_1, \dots, x_n) = \frac{1}{Z} \prod_{i=1}^n \underline{p_i^{(1)}(x_i)} \prod_{i,j \in E} \underline{p_{i,j}^{(2)}(x_i, x_j)}$$

Inference

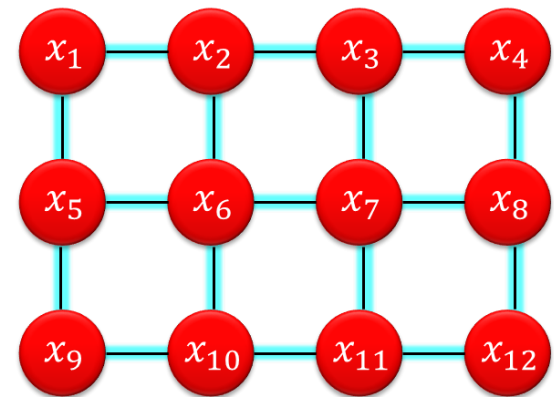
- Try all combinations of (x_1, \dots, x_n)
- Compute probability $p(x_1, \dots, x_n)$
- Determine maximum (or marginalize)



Inference

Complexity

- Brute-force-search / -integration:
 - Infeasible for large dimensions
 - Only toy-models
- Analytic maximization / integration
 - Special models only
 - Example: Gaussians
- Numerical maximization
 - All convex log-likelihoods
 - Gaussians (L_2 -error)
 - L_1 -errors
 - etc.



$$p(x_1, \dots, x_n) = \frac{1}{Z} \prod_{i=1}^n \underline{p_i^{(1)}(x_i)} \prod_{i,j \in E} \underline{p_{i,j}^{(2)}(x_i, x_j)}$$

Example:
Image Reconstruction

Image Reconstruction Model

Problem statement

- Measured 2D pixel image
- Distorted by noise
- Want to remove noise

Bayesian problem modeling

- Model of measurement process
- Prior distribution on images (this is Bayesian)

Inference: Maximum-a-posteriori

Model

Image

- $m_{i,j}$ with $i = 1 \dots w, j = 1, \dots, h$
- (continuous analogue: $f: [1, w] \times [1, h] \rightarrow \mathbb{R}$)

Probability space

- $\Omega = \mathbb{R}^{w \times h}$

Model

Bayes rule

$$P(\textcolor{blue}{M}|\textcolor{red}{D}) \sim P(\textcolor{red}{D}|\textcolor{blue}{M}) \cdot P(\textcolor{blue}{M})$$

Likelihood

- $P(\textcolor{red}{D}|\textcolor{blue}{M}) = \prod_{i=1}^w \prod_{j=1}^h P(\textcolor{red}{d}_i|\textcolor{blue}{m}_i)$ (i.i.d. noise)
 $= \prod_{i=1}^w \prod_{j=1}^h \mathcal{N}_{\textcolor{red}{d}_i, \sigma_D}(\textcolor{blue}{m}_i)$ (Gaussian noise)
 $= \prod_{i=1}^w \prod_{j=1}^h \left[\frac{1}{\sigma_D \sqrt{2\pi}} e^{-\frac{(\textcolor{blue}{m}_i - \textcolor{red}{d}_i)^2}{2\sigma_D^2}} \right]$
(Gaussian distribution)

Model

Likelihood

$$\blacksquare P(\textcolor{red}{D}|\textcolor{blue}{M}) = \prod_{i=1}^w \prod_{j=1}^h \left[\frac{1}{\sigma_D \sqrt{2\pi}} e^{-\frac{(\textcolor{blue}{m}_i - \textcolor{red}{d}_i)^2}{2\sigma_D^2}} \right]$$

Neg-log-likelihood

$$E(\textcolor{red}{D}|\textcolor{blue}{M}) := -\ln P(\textcolor{red}{D}|\textcolor{blue}{M}) = \sum_{i=1}^w \sum_{j=1}^h \frac{(\textcolor{blue}{m}_i - \textcolor{red}{d}_i)^2}{2\sigma_D^2} + \frac{\cancel{wh}}{\cancel{\sigma_D} \sqrt{2\pi}}$$

independent of $\textcolor{blue}{m}_i$

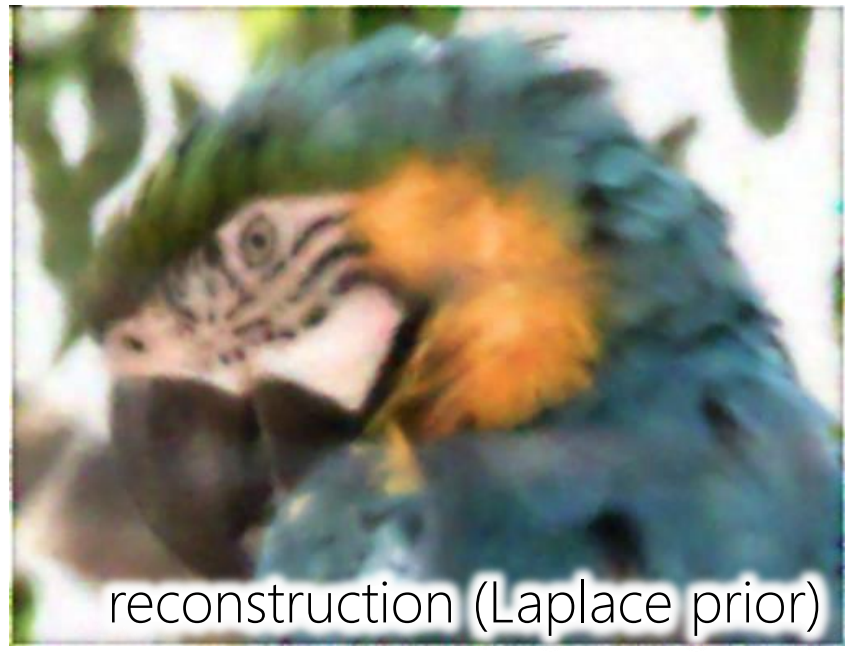
Model

Prior

- Assumption: Large image gradients are unlikely
- Gaussian distribution on Gradients
- Neg-log-likelihood: $\frac{1}{2\sigma^2} \|\nabla f\|^2$
- Discreet:

$$E(M) := -\ln P(M) = \sum_{i=1}^{w-1} \sum_{j=1}^{h-1} \frac{(m_{i+1,j} - m_{i,j})^2 + (m_{i,j+1} - m_{i,j})^2}{2\sigma_M^2} + \frac{wh}{\sigma_X \sqrt{2\pi}}$$

independent of m_i



Minimization Problem

Minimize

$$\begin{aligned} & E(\textcolor{red}{D}|\textcolor{blue}{M}) + E(\textcolor{blue}{M}) \\ &= \sum_{i=1}^w \sum_{j=1}^h \frac{(\textcolor{blue}{m}_{i,j} - \textcolor{red}{d}_{i,j})^2}{2\sigma_D^2} + \sum_{i=1}^{w-1} \sum_{j=1}^{h-1} \frac{(\textcolor{blue}{m}_{i+1,j} - \textcolor{blue}{m}_{i,j})^2 + (\textcolor{blue}{m}_{i,j+1} - \textcolor{blue}{m}_{i,j})^2}{2\sigma_M^2} \end{aligned}$$

Equivalent minimization objective

$$\sum_{i=1}^w \sum_{j=1}^h (\textcolor{blue}{m}_{i,j} - \textcolor{red}{d}_{i,j})^2 + \frac{\sigma_D^2}{\sigma_M^2} \sum_{i=1}^{w-1} \sum_{j=1}^{h-1} (\textcolor{blue}{m}_{i+1,j} - \textcolor{blue}{m}_{i,j})^2 + (\textcolor{blue}{m}_{i,j+1} - \textcolor{blue}{m}_{i,j})^2$$

Continuous

$$\int_{\Omega} (\textcolor{blue}{m}(\mathbf{x}) - \textcolor{red}{d}(\mathbf{x}))^2 d\mathbf{x} + \frac{\sigma_M^2}{\sigma_D^2} \int_{\Omega} \|\nabla \textcolor{blue}{m}(\mathbf{x})\|^2 d\mathbf{x}$$

Numerical Solution

Looks familiar?

- Solution via linear system

Variant

- Penalize l_1 norm instead of l_2 norm of gradients

$$\int_{\Omega} (m(\mathbf{x}) - d(\mathbf{x}))^2 d\mathbf{x} + \frac{\sigma_D^2}{\sigma_M^2} \int_{\Omega} \|\nabla m(\mathbf{x})\|^1 d\mathbf{x}$$

- Laplace distribution (single exponential)
- Yields sharper images (natural image statistics)

Technical Remark

Image prior

$$-\ln P(\mathbf{M}) = \sum_{i=1}^{w-1} \sum_{j=1}^{h-1} \frac{(\mathbf{m}_{i+1,j} - \mathbf{m}_{i,j})^2 + (\mathbf{m}_{i,j+1} - \mathbf{m}_{i,j})^2}{2\sigma_M^2} + \frac{wh}{\sigma_M \sqrt{2\pi}}$$

- This is an “improper prior”
 - Does not integrate to one!
 - Infinite subspaces without penalty
- Formal fix
 - Assume broader prior on function value itself: $f \sim N_{0, \sigma_{\text{very large}}}$
- For MAP estimation, this does not matter
 - We just find a point of maximum density
 - Integration not required

Fancy Inference Schemes

Belief Propagation

Model

$$p(\mathbf{x}) = \frac{1}{Z} \prod_{i=1}^n \underline{p_i^{(1)}(x_i)} \prod_{i,j \in E} \underline{p_{i,j}^{(2)}(x_i, x_j)}$$

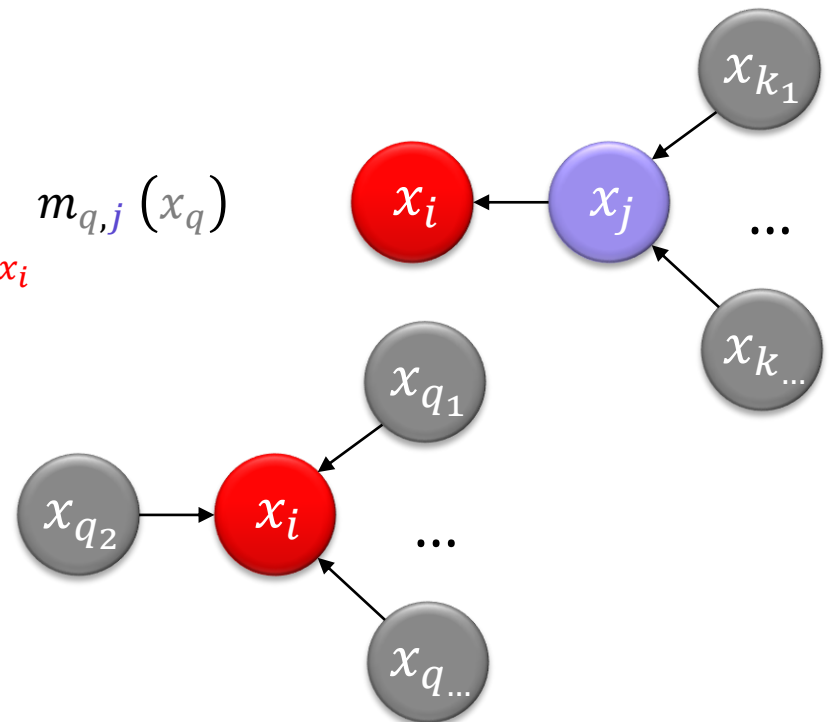
Messages

- Maximum a posteriori inference?
- Use “max-marginal”:

$$m_{j,i}(x_i) = \max_{x_j=1\dots k} p_i^{(1)}(x_j) p_{i,j}^{(2)}(x_i, x_j) \prod_{q \in N(j) \setminus x_i} m_{q,j}(x_q)$$

$$b_i(x_i) = \frac{1}{Z_i} p_i^{(1)}(x_i) \prod_{q \in N(i)} m_{q,i}(x_q)$$

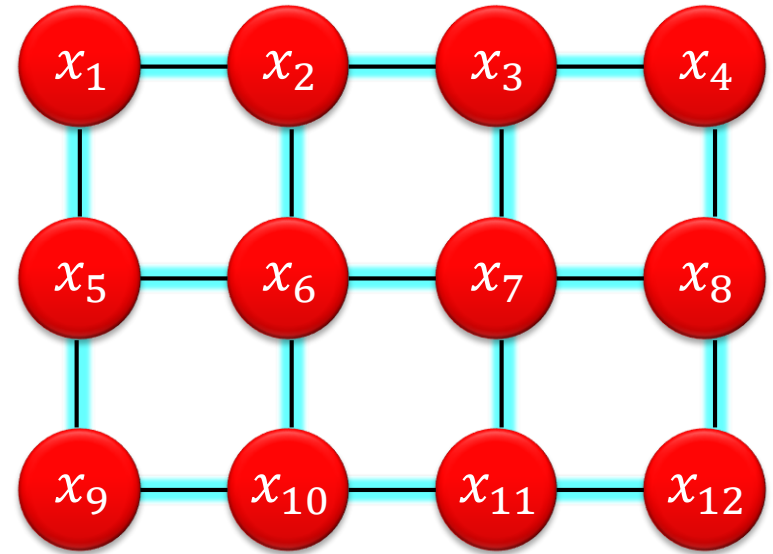
most likely state = $\arg \max_{x_i} b_i(x_i)$



Loopy BP

Loopy BP

- Loops? Which loops?
- Just run BP on loopy graph
- Arbitrary order

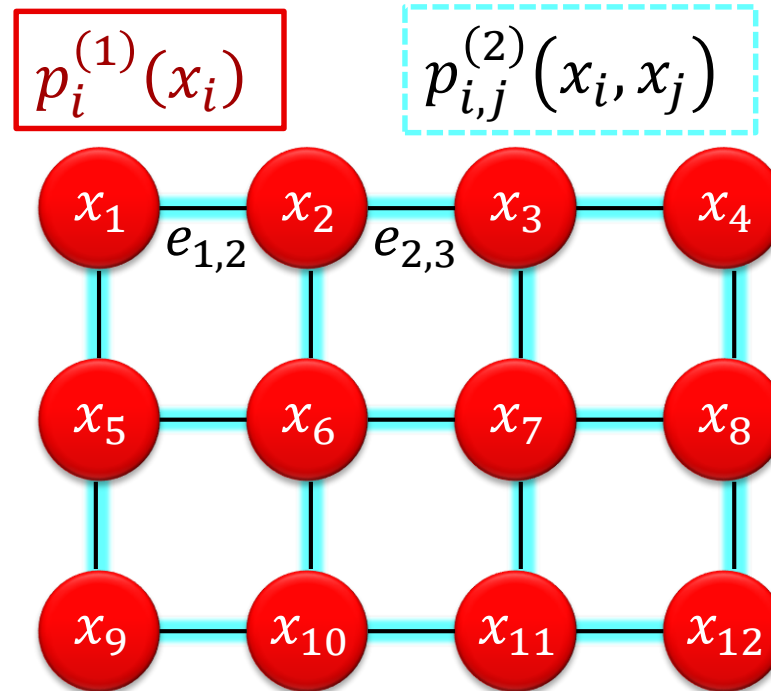


Result

- No guarantees (results can be wrong)
- Most often, results will be wrong
- Frequently still a reasonable approximation
 - Problem dependent, but worth a try. Popular 10-20 years ago.

Graph Cut

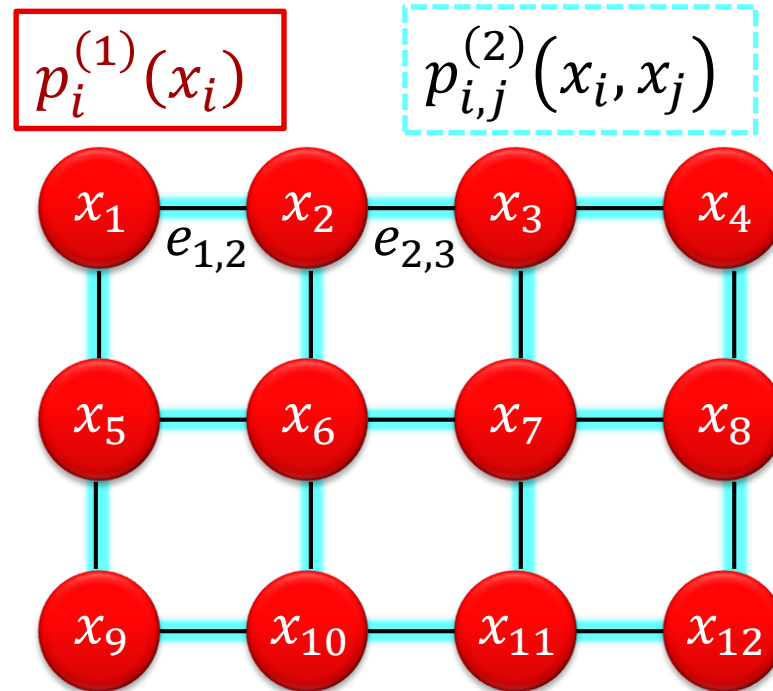
Pairwise MRF



Pairwise MRF

$$p(x_1, \dots, x_n) = \frac{1}{Z} \prod_{i=1}^n \underline{p_i^{(1)}(x_i)} \prod_{i,j \in E} \underline{p_{i,j}^{(2)}(x_i, x_j)}$$

Neg-Log Likelihood



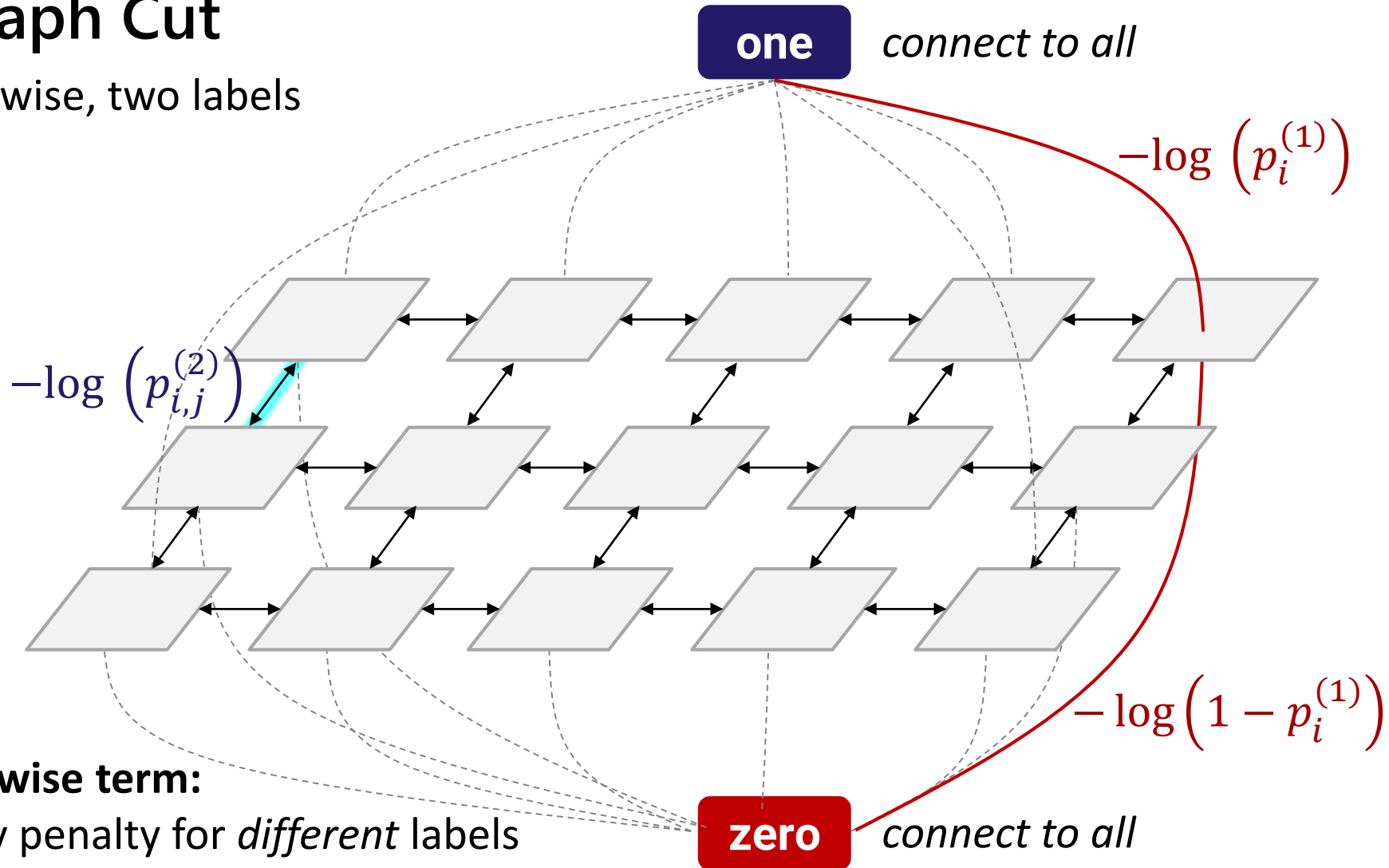
Pairwise MRF

$$-\log(p(x_1, \dots, x_n)) = \sum_{i=1}^n \underbrace{\left[-\log \left(p_i^{(1)}(x_i) \right) \right]}_{\text{red}} \sum_{i,j \in E} \underbrace{\left[-\log \left(p_{i,j}^{(2)}(x_i, x_j) \right) \right]}_{\text{cyan}}$$

Pairwise MRF Inference

Graph Cut

Pairwise, two labels



Pairwise term:

Only penalty for *different* labels possible (submodular!)

Example Application

“Grab-Cut”

- Rother, Kolmogorov, Blake
Siggraph 2004

Probabilistic Model

- Per pixel: “Gaussian Mixture”
of pixel colors
 - Foreground (red box)
and background (rest)
- Pairwise: Neighboring Pixels
 - Different Label (f/b) incur fixed cost
- Graph-Cut Inference

MCMC

More General Tools

MCMC (Markov Chain Monte Carlo)

- Gibbs Sampling
- Metropolis algorithm
- Many variants

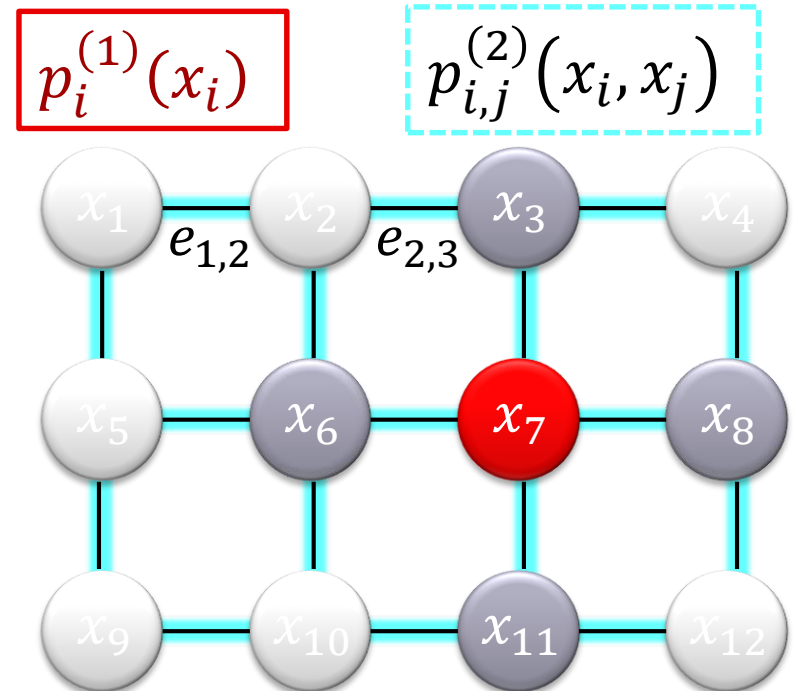
Idea

Naiive Sampling

- Infeasible (exponential)

Gibbs Sampling

- Random initialization
- Select random node
 - Fix neighbors
 - Compute local distribution
 - Sample
- Repeat



Convergence: Mixing times (hard to estimate)

Literature

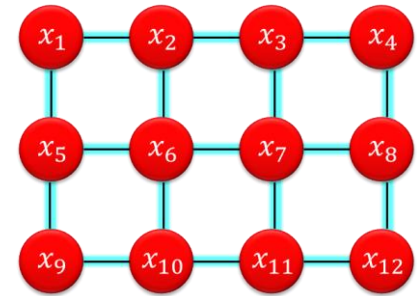
S. Geman & D Geman:

Stochastic Relaxation, Gibbs Distributions,
and the Bayesian Restoration of Images.

In: IEEE Transactions on Pattern Analysis and
Machine Intelligence (PAMI) 6(6), Nov. 1984.

Learning MRF-Models from training data

Learning



Maximum Likelihood: Maximize

$$p(x_1, \dots, x_n) = \frac{1}{Z} \prod_{i=1}^n p_i^{(1)}(x_i) \prod_{i,j \in E} p_{i,j}^{(2)}(x_i, x_j)$$

$$\begin{aligned} &= \frac{\prod_{i=1}^n p_i^{(1)}(x_i) \prod_{i,j \in E} p_{i,j}^{(2)}(x_i, x_j)}{\int_{x_1, \dots, x_n} \prod_{i=1}^n p_i^{(1)}(x_i) \prod_{i,j \in E} p_{i,j}^{(2)}(x_i, x_j) dx_1 \cdots dx_n} \\ &= \frac{P_X^{(1)}(\theta) P_X^{(2)}(\theta)}{Z(\theta)} \quad (\theta \text{ denoting set of learnable parameters,} \\ &\quad X \text{ denoting training data}) \end{aligned}$$

$Z(\theta)$ is numerically intractable for high dimensions, general distributions

In Practice

Learning MRFs in practice

- Tractable models, e.g. Gaussian
 - Often too simple
- Ignore Z
 - kind of wrong
 - but easy to implement
- Approximate Z
 - Mean-field theory; general: Variational Bayes
 - Replace integral by maximum (convex Z ; still kind of wrong)
- Direct discriminative learning (ignore Bayes rule)