

# Perspectives on Probabilistic Graphical Models

Dong Liu

*Information Science and Engineering  
KTH - Royal Institute of Technology*



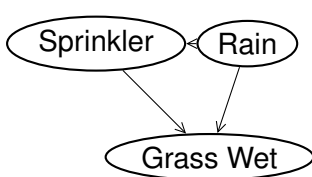
Profile page: <https://firsthandscientist.github.io/>

Slide is available at: <https://github.com/FirstHandScientist/phdthesis>

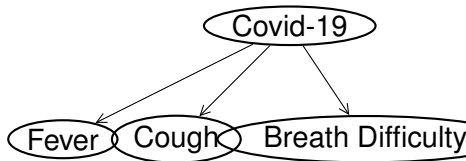
Why are Probabilistic Graphical Models interested?

# PROBABILISTIC GRAPHICAL MODELS

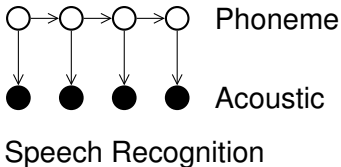
Directed: Bayesian Networks



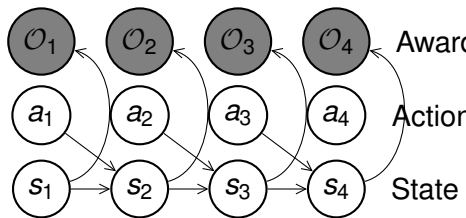
Is the sprinkler working?



Is the person get contagious by COVID?



Speech Recognition

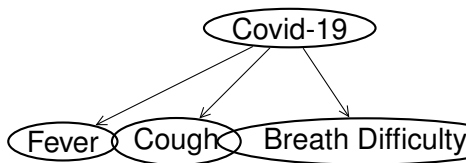
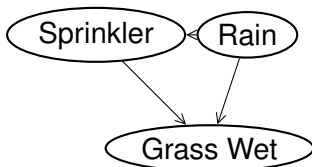


Communication system

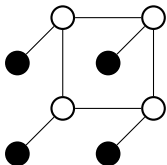
Two key aspects to encode in a graphical model:

# PROBABILISTIC GRAPHICAL MODELS

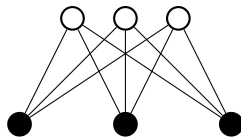
Directed: Bayesian Networks



Undirected: Markov Random Field



Computer vision



Control, reinforcement learning

Two key aspects to encode in a graphical model:

- attributes of our interests in a system → variable nodes
- relationship of these factors (dependencies or independencies) → structures of a graph

# MARKOV RANDOM FIELD

## MARKOV RANDOM FIELD (MRF) AND FUNDAMENTAL INFERENCE PROBLEMS

Let us walk through via MRF

- An MRF can be represented by a graph  $\mathcal{G}(\mathcal{V}, \mathcal{E})$  with each node  $i \in \mathcal{V}$  is associated with a random variable  $X_i$
- The probability distribution (Gibbs distribution) is

$$p(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_a \psi_a(\mathbf{x}_a; \boldsymbol{\theta}_a),$$

where  $\mathbf{x} = \{x_1, x_2, \dots, x_N\}$ ,  $a$  indexes potential functions  $\mathcal{I} = \{\psi_A, \psi_B, \dots, \psi_M\}$  and  $\boldsymbol{\theta}$  is set of potential function parameters.  $Z(\boldsymbol{\theta}) = \sum_{\mathbf{x}} \prod_a \psi_a(\mathbf{x}_a; \boldsymbol{\theta}_a)$ .

What do we do with graphical models?

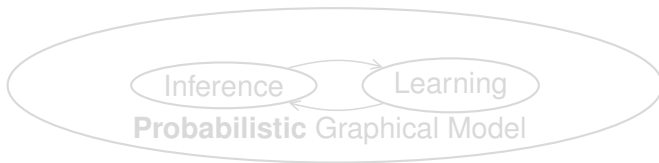
# USAGE OF GRAPHICAL MODELS

In general:

- Representation
  - In place of real systems
  - Abstraction of complex problems or systems (with subjective bias)
- Answer queries

Evidence (observation)  $\rightarrow$  ??  $\rightarrow$  Answers

Two components interacting with each other:



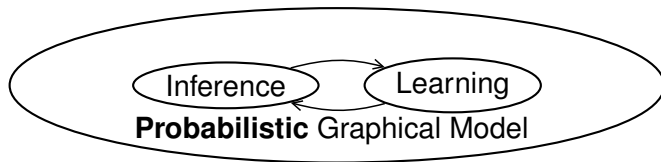
# USAGE OF GRAPHICAL MODELS

In general:

- Representation
  - In place of real systems
  - Abstraction of complex problems or systems (with subjective bias)
- Answer queries

Evidence (observation)  $\rightarrow$  ??  $\rightarrow$  Answers

Two components interacting with each other:





# USAGE OF GRAPHICAL MODELS

Why impact in two direction?

- Learning to Inference:



A graphical model

- built by expert knowledge, or
  - built by extracting information from evidence (empirical data).
- Inference to Learning:



Model learning: an error trial process that compares inferred 'fact' and actual fact (evidence).

Model learning usually needs inference as a subroutine, which sometimes are replaced by sampling in particle based methods.

# USAGE OF GRAPHICAL MODELS

Why impact in two direction?

- Learning to Inference:



A graphical model

- built by expert knowledge, or
- built by extracting information from evidence (empirical data).

- Inference to Learning:



Model learning: an error trial process that compares inferred 'fact' and actual fact (evidence).

Model learning usually needs inference as a subroutine, which sometimes are replaced by sampling in particle based methods.

# COMMON INFERENCE PROBLEMS

The common inference problems in a MRF  $\mathcal{G}(\mathcal{V}, \mathcal{E})$ :

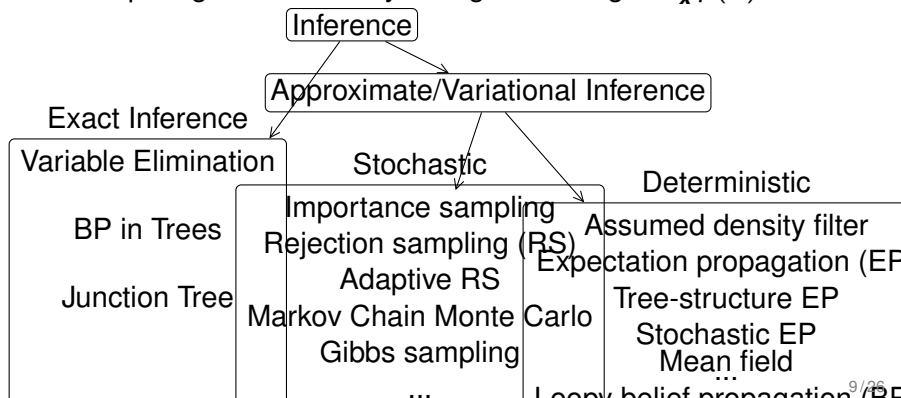
- Computing the likelihood of observed data.
- Computing the marginals distribution  $p(\mathbf{x}_A)$  over particular subset  $A \subset \mathcal{V}$  of nodes
- Computing the conditional distribution  $p(\mathbf{x}_A | \mathbf{x}_B)$ ,
- Computing the most likely configuration  $\text{argmax}_{\mathbf{x}} p(\mathbf{x})$



## COMMON INFERENCE PROBLEMS

The common inference problems in a MRF  $\mathcal{G}(\mathcal{V}, \mathcal{E})$ :

- Computing the likelihood of observed data.
- Computing the marginals distribution  $p(\mathbf{x}_A)$  over particular subset  $A \subset \mathcal{V}$  of nodes
- Computing the conditional distribution  $p(\mathbf{x}_A | \mathbf{x}_B)$ ,
- Computing the most likely configuration  $\operatorname{argmax}_{\mathbf{x}} p(\mathbf{x})$



# INFERENCE ROUTINE IN LEARNING

What is  $\theta$  in  $p(\mathbf{x}; \theta) = \frac{1}{Z(\theta)} \prod_a \psi_a(\mathbf{x}_a; \theta_a)$ ?

A direct view:

$$\max_{\theta} \log p(\mathbf{x}; \theta) = \max_{\theta} \sum_a \log \psi_a(\mathbf{x}_a; \theta_a) \quad \underbrace{- \log Z(\theta)}_{\text{Helmholtz free energy}},$$

An alternative view:

$$\frac{\partial \log p(\mathbf{x}; \theta)}{\partial \theta_a} = \frac{\partial \log \varphi_a(\mathbf{x}_a; \theta_a)}{\partial \theta_a} - \mathbb{E}_{p(\mathbf{x}_a; \theta)} \left[ \frac{\partial \log \varphi_a(\mathbf{x}_a; \theta_a)}{\partial \theta_a} \right].$$

Remark:

- This essentially requires estimation of Helmholtz free energy or marginal probabilities.
- Stationary points translate into moment matching.

## INFERENCE ROUTINE IN LEARNING

What is  $\theta$  in  $p(\mathbf{x}; \theta) = \frac{1}{Z(\theta)} \prod_a \psi_a(\mathbf{x}_a; \theta_a)$ ?

A direct view:

$$\max_{\theta} \log p(\mathbf{x}; \theta) = \max_{\theta} \sum_a \log \psi_a(\mathbf{x}_a; \theta_a) \quad \underbrace{- \log Z(\theta)}_{\text{Helmholtz free energy}},$$

An alternative view:

$$\frac{\partial \log p(\mathbf{x}; \theta)}{\partial \theta_a} = \frac{\partial \log \varphi_a(\mathbf{x}_a; \theta_a)}{\partial \theta_a} - \mathbb{E}_{p(\mathbf{x}_a; \theta)} \left[ \frac{\partial \log \varphi_a(\mathbf{x}_a; \theta_a)}{\partial \theta_a} \right].$$

Remark:

- This essentially requires estimation of Helmholtz free energy or marginal probabilities.
- Stationary points translate into moment matching.

Play with **Gibbs (variational) free energy**

$$F_V(b) = \text{KL}(b(\mathbf{x})||p(\mathbf{x}; \theta)) - \log Z(\theta)$$

with trial  $b(\mathbf{x})$ .

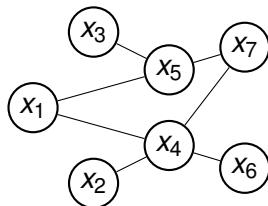
# WHAT IS THE STATE OF $x$ ?

## A TOY EXAMPLE

Assume that we are interested into the state of node  $i$  in an MRF, it can be answered by

- the probability  $p(x_i)$ , or
- an empirical version, a collection of samples  $\{x_i^n\}_{n=1}^N$

It is similar for the case when  $\mathbf{x}$  is of interests, instead of  $x_i$ .



what is the state of  $x_4$

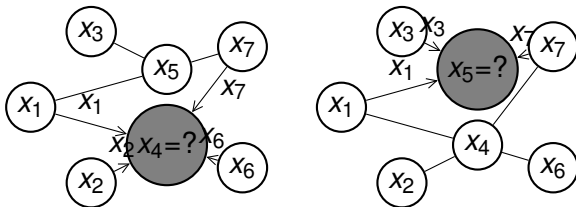


# WHAT IS THE STATE OF $x$ ?

GIBBS SAMPLING: LET US GUESS BY SAMPLING

We can approximately sample iteratively:

$$x_i \sim p(x_i | \mathbf{x}_{-i}) \sim p(x_i, \mathbf{x}_{-i})$$



This coordinate-wise sampling algorithm is called Gibbs sampling, which answers queries by collected samples  $\{\mathbf{x}^n\}_1^N$ .

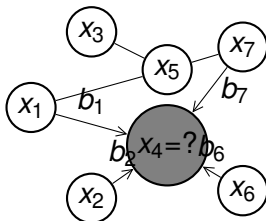
Gibbs sampling is named after the physicist Josiah Willard Gibbs, which was described by brothers Stuart and Donald Geman in 1984, some eight decades after the death of Gibbs.

$$x_i \sim \exp\left\{\sum_{a \in \text{ne}_i} \log \varphi_a(x_i, \mathbf{x}_{a-i}; \theta_a)\right\}$$

where  $\text{ne}_i$  gives the neighboring potential factors of node  $i$ .

# WHAT IS THE STATE OF $x$ ?

Naive Mean Field: **message in form of sample values** → **message in form of belief**



Corresponding to minimization of **variational free energy**  $F_V(b)$  with trial  $b$  in **fully-factorized form** for univariate  $\{b_i\}$ .

Iterative sampling → iterative belief update via

$$\log b_i(x_i) \propto \sum_{a \in \text{ne}_i} \sum_{\mathbf{x}_a \setminus x_i} \prod_{j \in a \setminus i} b_j(x_j) \log \varphi_a(\mathbf{x}_a; \theta_a).$$

# WHAT IS THE STATE OF $x$ ?

BELIEF PROPAGATION (BP): LET US GUESS BY PROPAGATING BELIEF

Proposed by Pearl (1982) for Bayesian networks

(tree-structured graphs), which then widely used for general graphs (loopy BP).

Yedidia, et al, connected the loopy BP with stationary points of **Bethe free energy**

$$F_{\text{Bethe}}(b) = \sum_{a \in \mathcal{F}} \sum_{\mathbf{x}_a} b_a(\mathbf{x}_a) \log \frac{b_a(\mathbf{x}_a)}{\varphi_a(\mathbf{x}_a)} - \sum_{i=1}^N (|\text{ne}_i| - 1) \sum_{x_i} b_i(x_i) \log b_i(x_i)$$

Corresponding to minimization of approximated **variational free energy**  $F_v(b)$  with trial  $b$  includes  $\{b_i\}$  and  $\{b_a\}$ .

$$\text{msg : factor to variable } m_{a \rightarrow i}(x_i) \propto \sum_{\mathbf{x}_a \setminus x_i} \varphi_a(\mathbf{x}_a) \prod_{j \in a \setminus i} m_{j \rightarrow a}(x_j),$$

$$\text{msg : variable to factor } m_{j \rightarrow a}(x_j) \propto \prod_{a' \in \text{ne}_j \setminus a} m_{a' \rightarrow j}(x_j)$$

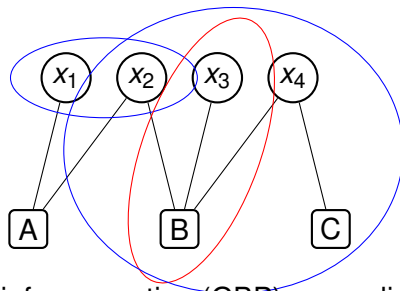
See, D. Liu, M. T. Vu, Z. Li, and Lars K. Rasmussen.  $\alpha$  belief propagation for approximate inference. 2020

D. Liu, N. N. Moghdam, L. K. Rasmussen, etc.  $\alpha$  belief propagation as fully factorized approximation. In GlobalSIP, 2019. for alternative view to loopy BP.

# WHAT IS THE STATE OF $x$ ?

YEDIDIA, FREEMAN, WEISS: A STEP TO GENERALIZATION

**Message among variables & factors** → **message among regions**



Generalized belief propagation (GBP) generalizes loopy BP

- usual better approximation than LBP
- higher complexity
- sensitive to scheduling of region messages

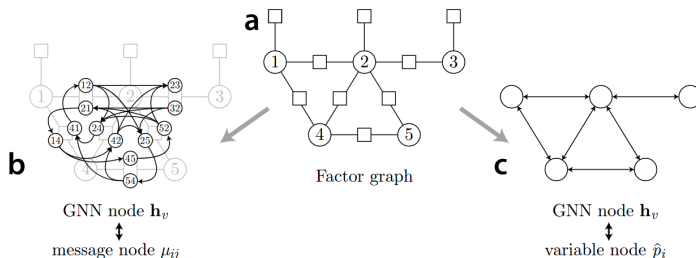
Corresponding to minimization of approximated **variational free energy**  $F_V(b)$  **with trial**  $b$  **including**  $\{b_R\}$ .

A *region*  $R$  is a set  $V_R$  of variables nodes and a set  $A_R$  of factor nodes, such that if a factor node ' $a$ ' belongs to  $A_R$ , all the variables nodes neighboring  $a$  are in  $V_R$ .

Attempts with neural networks: an imitation game of message passing, or trials under free energy?

# LEARN THE MESSAGE UPDATE RULE BY NN

An end-to-end learning process: Factor graph  $\rightarrow$  converted graph representation  $\rightarrow$  GNN  $\rightarrow$  Output



- sum-product update rule (in BP)  $\rightarrow$  NN, to learn
- pseudo probability (belief aggregation)  $\rightarrow$  NN, to learn
- end-to-end learning that requires true marginal probability, which BP, GPB and mean field do not require

For related methods, see:

Heess et al, Learning to Pass Expectation Propagation Messages

Yoon, et al. 2019, Inference in Probabilistic Graphical Models by Graph Neural Networks

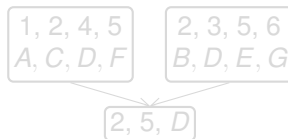
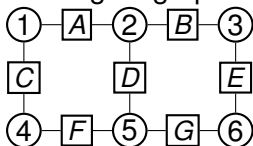
# RENN

## REGION REVISITED

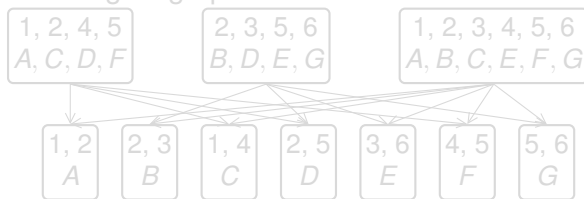
- If you cannot collect true targets ( $p(x_i)$ )
- If you are unwilling to be restricted to pre-defined inference

Factor graph representation of MRF (2-by-3 grid) with factor nodes.

MRF  $\rightarrow$  region graph:



An alternative region graph of the same MRF:



# RENN

## REGION REVISITED

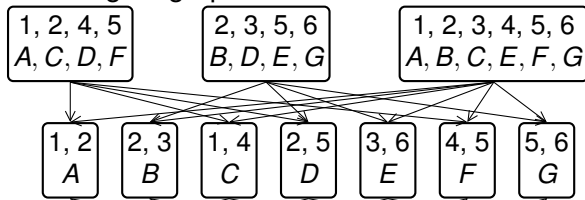
- If you cannot collect true targets ( $p(x_i)$ )
- If you are unwilling to be restricted to pre-defined inference

Factor graph representation of MRF (2-by-3 grid) with factor nodes.

MRF  $\rightarrow$  region graph:



An alternative region graph of the same MRF:



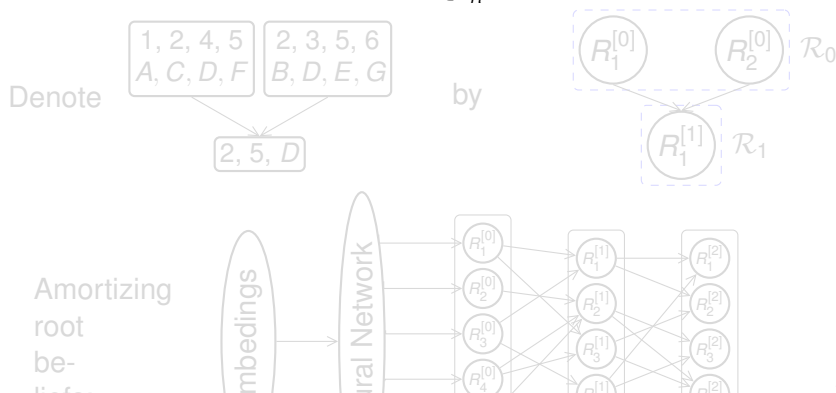


# RENN

The region-based free energy of a region graph is

$$F_R(\mathcal{B}; \theta) = \sum_{R \in \mathcal{R}} \underbrace{c_R}_{\text{counting number}} \sum_{\mathbf{x}_R} b_R(\mathbf{x}_R) \left( \underbrace{E_R(\mathbf{x}_R; \theta_R)}_{\text{region average energy}} + \ln b_R(\mathbf{x}_R) \right)$$

- counting number: balance the contribution of each region
- region average energy:  $-\sum_{a \in A_R} \ln \varphi_a(\mathbf{x}_a; \theta_a)$

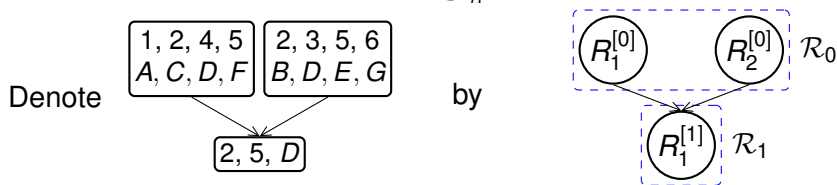


# RENN

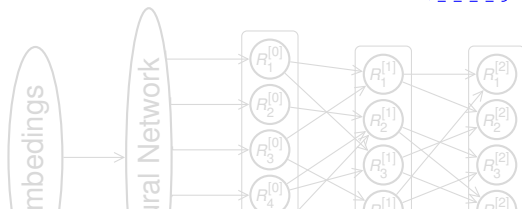
The region-based free energy of a region graph is

$$F_R(\mathcal{B}; \theta) = \sum_{R \in \mathcal{R}} \underbrace{c_R}_{\text{counting number}} \sum_{\mathbf{x}_R} b_R(\mathbf{x}_R) \left( \underbrace{E_R(\mathbf{x}_R; \theta_R)}_{\text{region average energy}} + \ln b_R(\mathbf{x}_R) \right)$$

- counting number: balance the contribution of each region
- region average energy:  $-\sum_{a \in A_R} \ln \varphi_a(\mathbf{x}_a; \theta_a)$



Amortizing  
root  
be-  
lie-  
f

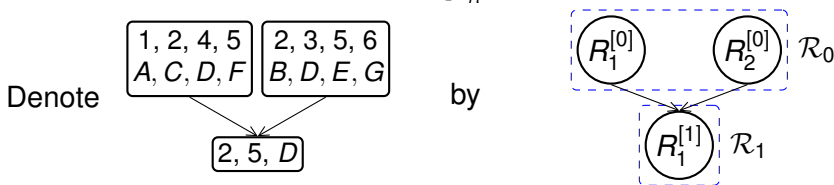


# RENN

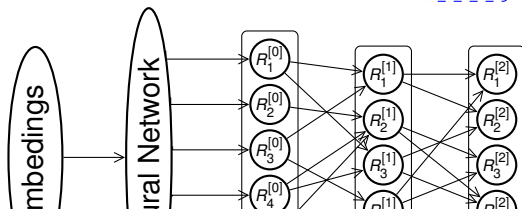
The region-based free energy of a region graph is

$$F_R(\mathcal{B}; \theta) = \sum_{R \in \mathcal{R}} \underbrace{c_R}_{\text{counting number}} \sum_{\mathbf{x}_R} b_R(\mathbf{x}_R) \left( \underbrace{E_R(\mathbf{x}_R; \theta_R)}_{\text{region average energy}} + \ln b_R(\mathbf{x}_R) \right)$$

- counting number: balance the contribution of each region
- region average energy:  $-\sum_{a \in A_R} \ln \varphi_a(\mathbf{x}_a; \theta_a)$



Amortizing  
root  
be-  
liefs



# RENN

Objective of RENN<sup>1</sup>:

min **region-based free energy**( $F_R$ ) + **panelty on belief consistency**  
*along region graph struture*

Inference only



RENN

Learning alternatives of MRFs

learn with customized optm.

$$\frac{\partial \log p(\mathbf{x}; \theta)}{\partial \theta_a} = \frac{\partial \log \varphi_a(\mathbf{x}_a; \theta_a)}{\partial \theta_a}$$

$$- \mathbb{E}_{p(\mathbf{x}; \theta)} \left[ \frac{\partial \log \varphi_a(\mathbf{x}_a; \theta_a)}{\partial \theta_a} \right]$$

learn with auto-grads

$$\max_{\theta} \log p(\mathbf{x}; \theta) = \max_{\theta} \sum_a \log \psi$$

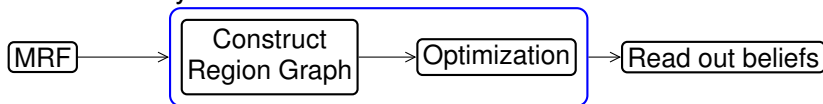
$$- \log Z$$

# RENN

Objective of RENN<sup>1</sup>:

min **region-based free energy**( $F_R$ ) + **panelty on belief consistency**  
*along region graph struture*

Inference only



RENN

Learning alternatives of MRFs

learn with customized optm.

learn with auto-grads

$$\frac{\partial \log p(\mathbf{x}; \theta)}{\partial \theta_a} = \frac{\partial \log \varphi_a(\mathbf{x}_a; \theta_a)}{\partial \theta_a}$$

$$\max_{\theta} \log p(\mathbf{x}; \theta) = \max_{\theta} \sum_a \log \psi$$

$$- \mathbb{E}_{p(\mathbf{x}; \theta)} \left[ \frac{\partial \log \varphi_a(\mathbf{x}_a; \theta_a)}{\partial \theta_a} \right]$$

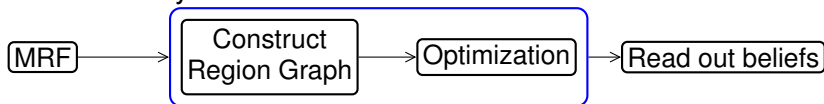
$$- \log Z$$

# RENN

Objective of RENN<sup>1</sup>:

min **region-based free energy** ( $F_R$ ) + **penalty on belief consistency**  
*along region graph structure*

Inference only



RENN

Learning alternatives of MRFs

learn with customized optm.

learn with auto-grads

$$\frac{\partial \log p(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}_a} = \frac{\partial \log \varphi_a(\mathbf{x}_a; \boldsymbol{\theta}_a)}{\partial \boldsymbol{\theta}_a} - \mathbb{E}_{p(\mathbf{x}; \boldsymbol{\theta})} \left[ \frac{\partial \log \varphi_a(\mathbf{x}_a; \boldsymbol{\theta}_a)}{\partial \boldsymbol{\theta}_a} \right]$$

$$\max_{\boldsymbol{\theta}} \log p(\mathbf{x}; \boldsymbol{\theta}) = \max_{\boldsymbol{\theta}} \sum_a \log \psi_a(\mathbf{x}_a; \boldsymbol{\theta}_a)$$

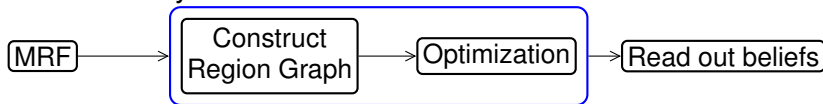
$$- \log Z$$

# RENN

Objective of RENN<sup>1</sup>:

min **region-based free energy** ( $F_R$ ) + **panelty on belief consistency**  
*along region graph struture*

Inference only



RENN

Learning alternatives of MRFs

learn with customized optm.

$$\frac{\partial \log p(\mathbf{x}; \theta)}{\partial \theta_a} = \frac{\partial \log \varphi_a(\mathbf{x}_a; \theta_a)}{\partial \theta_a} - \mathbb{E}_{p(\mathbf{x}; \theta)} \left[ \frac{\partial \log \varphi_a(\mathbf{x}_a; \theta_a)}{\partial \theta_a} \right]$$

learn with auto-grads

$$\max_{\theta} \log p(\mathbf{x}; \theta) = \max_{\theta} \sum_a \log \psi_a(\mathbf{x}_a; \theta_a) - \log Z$$

# INFERENCE RESULTS

Ising model:  $p(\mathbf{x}; \theta) = \frac{1}{Z(\theta)} \exp(\sum_{(i,j) \in \mathcal{E}_F} J_{ij} x_i x_j + \sum_{i \in \mathcal{V}} h_i x_i)$ ,

$\mathbf{x} \in \{-1, 1\}^N$ ,

- $J_{ij}$  is the pairwise log-potential between node  $i$  and  $j$ ,  
 $J_{ij} \sim \mathcal{N}(0, 1)$
- $h_i$  is the node log-potential for node  $i$ ,  $h_i \sim \mathcal{N}(0, \gamma^2)$

Inference on grid graph ( $\gamma = 0.1$ ).

Metric	$n$	Mean Field	Loopy BP	Damped BP	GBP	Inference Net	RENN
$\ell_1$ error	25	$0.271 \pm 0.051$	$0.086 \pm 0.078$	$0.084 \pm 0.076$	$0.057 \pm 0.024$	$0.111 \pm 0.072$	<b><math>0.049 \pm 0.078</math></b>
	100	$0.283 \pm 0.024$	$0.085 \pm 0.041$	$0.062 \pm 0.024$	$0.064 \pm 0.019$	$0.074 \pm 0.034$	<b><math>0.025 \pm 0.011</math></b>
	225	$0.284 \pm 0.019$	$0.100 \pm 0.025$	$0.076 \pm 0.025$	$0.073 \pm 0.013$	$0.073 \pm 0.012$	<b><math>0.046 \pm 0.011</math></b>
	400	$0.279 \pm 0.014$	$0.110 \pm 0.016$	$0.090 \pm 0.016$	$0.079 \pm 0.009$	$0.083 \pm 0.009$	<b><math>0.061 \pm 0.009</math></b>
Correlation $\rho$	25	$0.633 \pm 0.197$	$0.903 \pm 0.114$	$0.905 \pm 0.113$	$0.923 \pm 0.045$	$0.866 \pm 0.117$	<b><math>0.951 \pm 0.112</math></b>
	100	$0.582 \pm 0.112$	$0.827 \pm 0.134$	$0.902 \pm 0.059$	$0.899 \pm 0.043$	$0.903 \pm 0.049$	<b><math>0.983 \pm 0.012</math></b>
	225	$0.580 \pm 0.080$	$0.801 \pm 0.078$	$0.863 \pm 0.088$	$0.869 \pm 0.037$	$0.873 \pm 0.037$	<b><math>0.949 \pm 0.022</math></b>
	400	$0.596 \pm 0.054$	$0.779 \pm 0.059$	$0.822 \pm 0.047$	$0.852 \pm 0.024$	$0.841 \pm 0.028$	<b><math>0.912 \pm 0.025</math></b>
log $Z$ error	25	$2.512 \pm 1.060$	$0.549 \pm 0.373$	$0.557 \pm 0.369$	<b><math>0.169 \pm 0.142</math></b>	$0.762 \pm 0.439$	$0.240 \pm 0.140$
	100	$13.09 \pm 2.156$	$1.650 \pm 1.414$	$1.457 \pm 1.365$	<b><math>0.524 \pm 0.313</math></b>	$2.836 \pm 2.158$	$1.899 \pm 0.495$
	225	$29.93 \pm 4.679$	$3.348 \pm 1.954$	$3.423 \pm 2.157$	<b><math>1.008 \pm 0.653</math></b>	$3.249 \pm 2.058$	$4.344 \pm 0.813$
	400	$51.81 \pm 4.706$	$5.738 \pm 2.107$	$5.873 \pm 2.211$	<b><math>1.750 \pm 0.869</math></b>	$3.953 \pm 2.558$	$7.598 \pm 1.146$

- $\ell_1$  error of beliefs v.s. true
- correlation  $\rho$  between true and approximate marginals,
- log  $Z$  error, true v.s. free energy approximation



# LEARNING MRFs

What is  $\theta$  in  $p(\mathbf{x}; \theta)$ ?

Table of negative log-likelihood of learned MRFs

$n$	True	Exact	Mean Field	Loopy BP	Damped BP	GBP	Inference Net	RENN
Grid Graph								
25	9.000	9.004	9.811	9.139	9.196	10.56	9.252	<b>9.048</b>
100	19.34	19.38	23.48	19.92	20.02	28.61	20.29	<b>19.76</b>
225	63.90	63.97	69.01	66.44	66.25	92.62	68.15	<b>64.79</b>
Complete Graph								
9	3.276	3.286	9.558	5.201	5.880	10.06	5.262	<b>3.414</b>
16	4.883	4.934	28.74	13.64	18.95	24.45	13.77	<b>5.178</b>

learning

# SUMMARY

- Brief on probabilistic graphic models
- Overview of inference methods
- A focus on the message-passing
- Transition to inference methods with NN

Thank you for your attention.  
Q&A.