# Perspectives on Probabilistic Graphical Models

## Dong Liu

*Information Science and Engineering*
*Department of Intelligent Systems*
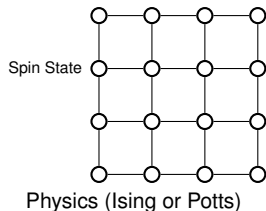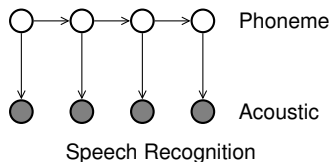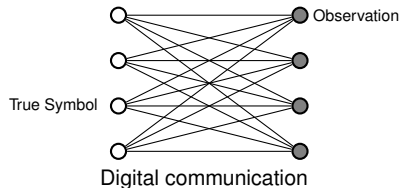*KTH - Royal Institute of Technology*



Profile page: https://firsthandscientist.github.io/
Slide is available at: https://github.com/FirstHandScientist/phdthesis

Why are probabilistic graphical models interesting?

## RICH REPRESENTATIONS



Speech Recognition

Decision-making (reinforcement learning)

Physics (Ising or Potts)

Digital communication

- Computer perception
- Error-control codes

- Computational biology
- Natural language processing
- etc.

# A GUIDE TO THIS DISSERTATION

## WHAT ARE PROBABILISTIC GRAPHICAL MODELS

Informally...
A PGM is a structured graph representation to encode

- Attributes of our interests in a system $\rightarrow$ variable nodes
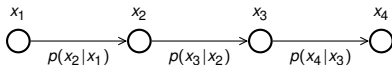- Relationship of these factors $\rightarrow$ structures of a graph

Intrinsic property: **reasoning with uncertainty**

# WHAT ARE PROBABILISTIC GRAPHICAL MODELS

### EXEMPLIFIED DEFINITIONS

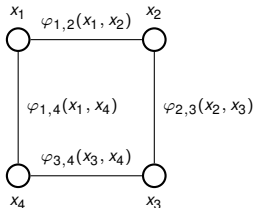A directed/undirected graph encoding dependencies/indepedencies of distribution $p(\boldsymbol{x}; \boldsymbol{\theta})$:

- A Generative model/BN is a directed graph (DAG)

$$x_1 \quad\xrightarrow{p(x_2|x_1)}\quad x_2 \quad\xrightarrow{p(x_3|x_2)}\quad x_3 \quad\xrightarrow{p(x_4|x_3)}\quad x_4$$

$$p(\boldsymbol{x}; \boldsymbol{\theta}) = \prod_{n=1}^{N} p(x_n | \underbrace{\mathcal{P}(x_n)}_{\substack{\text{parent nodes of } x_n}})$$

$$\underbrace{\phantom{p(\boldsymbol{x}; \boldsymbol{\theta})}}_{\substack{\text{the local functions} \\ \text{are proper distributions}}}$$

- An MRF denoted by an undirected graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$

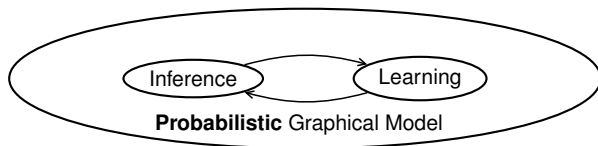$$\begin{array}{cc}
x_1 \quad \varphi_{1,2}(x_1, x_2) \quad x_2 \\
\varphi_{1,4}(x_1, x_4) \quad \varphi_{2,3}(x_2, x_3) \\
\varphi_{3,4}(x_3, x_4) \\
x_4 \qquad\qquad x_3
\end{array}$$

- The probability distribution (Gibbs distribution) is $p(\boldsymbol{x}; \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_{a \in \mathcal{I}} \underbrace{\psi_a(\boldsymbol{x}_a; \boldsymbol{\theta}_a)}_{\substack{\text{potential function is not} \\ \text{necessarily a proper distribution}}}$

- $a$ indexes potential functions $\mathcal{I} = \{\psi_A, \psi_B, \cdots, \psi_M\}$

- $Z(\boldsymbol{\theta}) = \sum_{\boldsymbol{x}} \prod_a \psi_a(\boldsymbol{x}_a; \boldsymbol{\theta}_a)$.

## WHAT TO DO WITH GRAPHICAL MODELS

- Inference

    - 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 partition function or the Helmholtz free energy (for MRFs)

- Learning

    - To model or determine $p(\mathbf{x}; \boldsymbol{\theta})$.

Two key components interacting with each other:

**Probabilistic** Graphical Model

# 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

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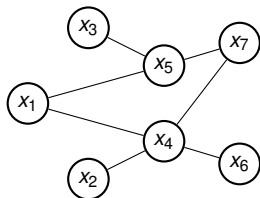


What is the state of $x_4$

# 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

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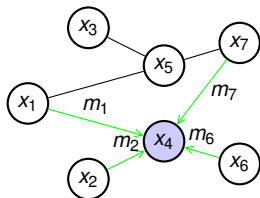- the probability $p(x_i)$



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

- Propagating beliefs iteratively
- Queries by collected beliefs $\{m_i\}$.

---

Intuition from *Gibbs (variational) free energy*

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

with trial $b(\boldsymbol{x})$. Instance: Bethe free energy.

Motivation
○

Preliminary
○○○○
○

Inference
●○○○○○
○○○○○○○○

Infer..Learn
○○○

Learning
○○○○○○○○
○○○○
○○○○○

Summary and Q&A
○
○

## ALTERNATIVE VEIW OF BP: $\alpha$-BP
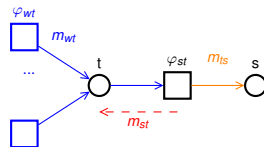
Ingredients:

- A pairwise Markov random field:
  $p(\boldsymbol{x}) \propto$
  $\prod_{s \in \mathcal{V}} \varphi_s(x_s) \prod_{(s,t) \in \mathcal{E}} \varphi_{st}(x_s, x_t)$
- A trial distribution: $q(\boldsymbol{x}) \propto$
  $\prod_{s \in \mathcal{V}} \tilde{\varphi}_s(x_s) \prod_{(s,t) \in \mathcal{E}} \tilde{\varphi}_{st}(x_s, x_t)$
  with factorization
  $\tilde{\varphi}_{s,t}(x_s, x_t) := m_{st}(x_t) m_{ts}(x_s)$
- A metric: $\alpha$-Divergence

A factor graph representation
$\mathcal{G}_F := (\mathcal{V} \cup \mathcal{F}, \mathcal{E}_F)$

---

Definition of $\alpha$-divergence $\mathcal{D}_\alpha(p \| q) = \frac{\sum_{\boldsymbol{x}} \alpha p(\boldsymbol{x}) + (1-\alpha) q(\boldsymbol{x}) - p(\boldsymbol{x})^\alpha q(\boldsymbol{x})^{1-\alpha}}{\alpha(1-\alpha)}$

Updating message via $\alpha$-BP:



$$\underbrace{m_{ts}^{\text{new}}(x_s)}_{\text{new msg. to s}} \propto \underbrace{m_{ts}(x_s)^{1-\alpha_{ts}}}_{\text{old msg. to s}} \left[ \sum_{x_t} \varphi_{ts}(x_t, x_s)^{\alpha_{ts}} \underbrace{m_{st}(x_t)^{1-\alpha_{ts}}}_{\text{old msg. to t}} \underbrace{\varphi_t(x_t) \prod_{w \in \mathcal{N}(t) \setminus s} m_{wt}(x_t)}_{\text{msg. from variable node } t \text{ to factor } \varphi_{st}} \right].$$
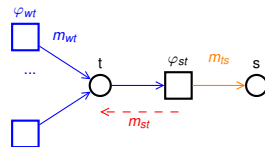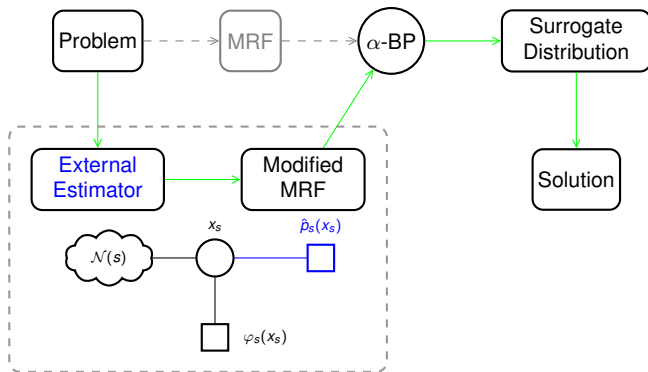
Updating message via $\alpha$-BP:



$$\underbrace{m_{ts}^{\text{new}}(x_s)}_{\text{new msg. to s}} \propto \underbrace{m_{ts}(x_s)^{1-\alpha_{ts}}}_{\text{old msg. to s}} \left[ \sum_{x_t} \varphi_{ts}(x_t, x_s)^{\alpha_{ts}} \underbrace{m_{st}(x_t)^{1-\alpha_{ts}}}_{\text{old msg. to t}} \varphi_t(x_t) \underbrace{\prod_{w \in \mathcal{N}(t) \backslash s} m_{wt}(x_t)}_{\text{msg. from variable node } t \text{ to factor } \varphi_{st}} \right].$$

# INSIGHTS OF $\alpha$-BP

### Connection to standard BP

- $\alpha \to 1$
- $\alpha$-divergence reduces to KL-divergence
- Update rule of $\alpha$-BP reduces to
  $m_{ts}^{new}(x_s) \propto$
  $\sum_{x_t} \varphi_{st}(x_s, x_t)\varphi_t(x_t) \prod_{w \in \mathcal{N}(t) \setminus s} m_{wt}(x_t)$,
  which is standard BP update rule

### Convergence

For an arbitrary pairwise Markov random field over binary variables, if the largest singular value of matrix $\boldsymbol{M}(\alpha, \theta)$ is less than one, $\alpha$-BP converges to a fixed point. The associated fixed point is unique.
See Corollary 3.1 for relaxed condition where singular value computation is avoided.

### What does that mean

- You can use $\alpha$-BP as an alternative to (loopy) BP
- You can use matrix $\boldsymbol{M}$ to check if you are guaranteed to get stable solution from $\alpha$-BP

---

Matrix $\boldsymbol{M}(\alpha, \theta)$, size $|\vec{\mathcal{E}}| \times |\vec{\mathcal{E}}|$

Each element is either 0 or a function of $\alpha$ and potentials factors

# SOME NUMERICAL RESULTS: APPLICATION CASE



Numerical results of $\alpha$-BP: symbol error of MIMO detection.

# CONTINUING: 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

Approximating **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 $V_R$.

## A TOY EXAMPLE OF REGION GRAPHS

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



- Clustering nodes
- level/layer-wise
- Hierarchical
- Msg. Scheduling
- ...
- See Section 4.1

# RENN: REGION-BASED ENERGY NEURAL NETWORK

The **region-based free energy** of a region graph is

$$F_R(\mathcal{B}; \boldsymbol{\theta}) = \sum_{R \in \mathcal{R}} \underbrace{c_R}_{\text{Counting number}} \underbrace{(\text{region average energy} - \text{region entropy})}_{\text{region free energy}},$$

# RENN: REGION-BASED ENERGY NEURAL NETWORK

The **region-based free energy** of a region graph is

$$F_R(\mathcal{B}; \boldsymbol{\theta}) = \sum_{R \in \mathcal{R}} \underbrace{c_R}_{\text{Counting number}} \underbrace{(\text{region average energy} - \text{region entropy})}_{\text{region free energy}},$$

# RENN

Non-root belief:



Marginalize to child region          Marginalize to child region

Average incoming marginalization from parents

Objective of RENN:

$$\min_{\substack{\text{parameter} \\ \text{of NN}}} \underbrace{\textbf{region-based free energy}(F_R)}_{\text{Accumulated over all regions}} + \underbrace{\textbf{penalty on belief inconsistency}}_{\text{Recursively computed via levels of region graph}}$$

# RENN

Non-root belief:



Objective of RENN:

$$\underset{\substack{\text{min} \\ \text{parameter} \\ \text{of NN}}}{}\quad \underbrace{\textbf{region-based free energy}(F_R)}_{\text{Accumulated over all regions}} + \underbrace{\textbf{penalty on belief inconsistency}}_{\text{Recursively computed via levels of region graph}}$$

RENN Inference:

## INSIGHTS OF RENN

### Generalization

Bethe free energy can be recovered from region-based free energy:

- two-level region graph representation
- constraint that each region can contain at most one factor node

Section 4.2.1

### Attributes of RENN

- RENN requires neither sampling technique nor training data (ground-truth marginal probabilities) in performing inference tasks; **on-the-fly inference**
- RENN does gradient descent w.r.t. its neural network parameter instead of iterative message-passing, and returns approximation of marginal probabilities and partition estimation **in one-shot**
- No message propagation, thus **no need of message scheduling**
- Competitive performance and efficiency

## SOME NUMERICAL COMPARISONS

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

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

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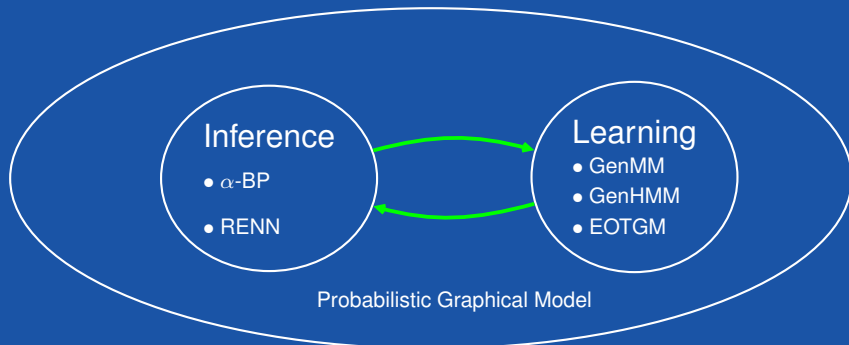- $\log Z$ error, true v.s. free energy approximation.

# SOME NUMERICAL COMPARISONS

RICHER COMPARISONS

Inference on grid and complete graphs.

| | | Metric | | Mean Field | Loopy BP | Damped BP | GBP | Inference Net | RENN |
|---|---|---|---|---|---|---|---|---|---|
| High Temp -erature | Complete graph N=16 $J_{ij} \sim \mathcal{N}(0,1)$ $h_i \sim \mathcal{N}(0,\gamma^2)$ | $\ell_1$- error | $\gamma = 1$ | $0.273 \pm 0.086$ | $0.239 \pm 0.059$ | $0.239 \pm 0.059$ | $0.260 \pm 0.086$ | $0.249 \pm 0.067$ | **0.181** $\pm 0.092$ |
| | | | $\gamma = 4$ | $0.197 \pm 0.049$ | $0.181 \pm 0.035$ | $0.180 \pm 0.034$ | $0.210 \pm 0.070$ | $0.174 \pm 0.030$ | **0.125** $\pm 0.050$ |
| | | $\log Z$ error | $\gamma = 1$ | $20.66 \pm 5.451$ | $178.7 \pm 22.18$ | $178.9 \pm 21.88$ | $153.3 \pm 25.29$ | $213.6 \pm 12.75$ | **14.41** $\pm 4.135$ |
| | | | $\gamma = 4$ | **10.74** $\pm 7.385$ | $565.7 \pm 73.33$ | $566.1 \pm 73.13$ | $106.0 \pm 54.43$ | $588.3 \pm 62.58$ | $14.72 \pm 4.155$ |
| Low Temp -erature | Grid graph N=100 $J_{ij} \sim \mathcal{U}(-u,u)$ $h_i \sim \mathcal{U}(-1,1)$ | $\ell_1$ error | 5 | $0.257 \pm 0.065$ | $0.115 \pm 0.071$ | $0.120 \pm 0.073$ | $0.250 \pm 0.024$ | $0.164 \pm 0.036$ | **0.100** $\pm 0.046$ |
| | | | 15 | $0.328 \pm 0.068$ | $0.228 \pm 0.088$ | $0.267 \pm 0.147$ | $0.303 \pm 0.026$ | $0.279 \pm 0.024$ | **0.207** $\pm 0.054$ |
| | | $\log Z$ error | 5 | $42.65 \pm 17.86$ | $7.346 \pm 7.744$ | **5.444** $\pm 4.811$ | $8.369 \pm 7.401$ | $65.60 \pm 8.786$ | $11.34 \pm 4.724$ |
| | | | 15 | $164.9 \pm 56.07$ | $58.40 \pm 41.36$ | $101.9 \pm 54.31$ | **23.10** $\pm 15.06$ | $224.3 \pm 25.52$ | $78.85 \pm 15.08$ |

Low temperature setting translates to high variance of coupling strength between nodes (larger variance of $J_{ij}$).

## INFERENCE ROUTINE IN LEARNING

What is $\boldsymbol{\theta}$ in $p(\boldsymbol{x}; \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_a \psi_a(\boldsymbol{x}_a; \boldsymbol{\theta}_a)$?
A direct view:

$$\max_{\boldsymbol{\theta}} \log p(\boldsymbol{x}; \boldsymbol{\theta}) = \max_{\boldsymbol{\theta}} \sum_a \log \psi_a(\boldsymbol{x}_a; \boldsymbol{\theta}_a) \quad \underbrace{-\log Z(\boldsymbol{\theta})}_{\text{can be est. by min } F_V} \quad ,$$

An alternative view:

$$\frac{\partial \log p(\boldsymbol{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}_a} = \frac{\partial \log \varphi_a(\boldsymbol{x}_a; \boldsymbol{\theta}_a)}{\partial \boldsymbol{\theta}_a} - \underbrace{\mathbb{E}_{p(\boldsymbol{x}_a; \boldsymbol{\theta})} \left[ \frac{\partial \log \varphi_a(\boldsymbol{x}_a; \boldsymbol{\theta}_a)}{\partial \boldsymbol{\theta}_a} \right]}_{\text{can be est. by beliefs}}.$$

Remark:

- This essentially requires estimation of partition function 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:

$$\max_{\theta} \log p(\mathbf{x}; \theta) = \max_{\theta} \sum_a \log \psi_a(\mathbf{x}_a; \theta_a) \quad \underbrace{-\log Z(\theta)}_{\text{can be est. by min } F_V} ,$$

An alternative view:

$$\frac{\partial \log p(\mathbf{x}; \theta)}{\partial \theta_a} = \frac{\partial \log \varphi_a(\mathbf{x}_a; \theta_a)}{\partial \theta_a} - \underbrace{\mathbb{E}_{p(\mathbf{x}_a; \theta)} \left[ \frac{\partial \log \varphi_a(\mathbf{x}_a; \theta_a)}{\partial \theta_a} \right]}_{\text{can be est. by beliefs}}.$$

Remark:

- This essen~~tially~~ ~~~~ ~~~~ babilities.
- Stationary ~~~~

Inference ⟶ Learning

- Two modules are not necessarily coupled
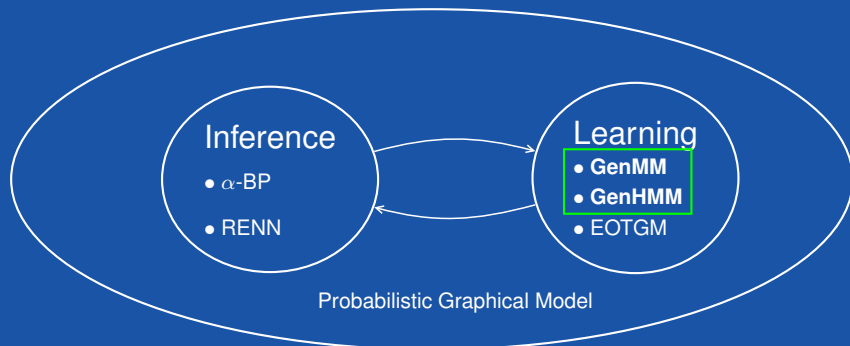- Each module may be replaced by another algorithm while the other one remains.

# LEARNING MRFS
WHAT IS $\theta$ IN $p(\boldsymbol{x}; \theta)$?

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 | **9.048** |
| 100 | 19.34 | 19.38 | 23.48 | 19.92 | 20.02 | 28.61 | 20. 29 | **19.76** |
| 225 | 63.90 | 63.97 | 69.01 | 66.44 | 66.25 | 92.62 | 68.15 | **64.79** |
| Complete Graph | | | | | | | | |
| 9 | 3.276 | 3.286 | 9.558 | 5.201 | 5.880 | 10.06 | 5.262 | **3.414** |
| 16 | 4.883 | 4.934 | 28.74 | 13.64 | 18.95 | 24.45 | 13.77 | **5.178** |

Average consumed time per epoch (unit: second) for two MRF learning cases.

| | Mean Field | Loopy BP | Damped BP | GBP | Inference Net | RENN |
|---|---|---|---|---|---|---|
| Grid $\mathcal{G}, N = 225$ | 40.09 | 335.1 | 525.1 | 12.37 | 19.49 | 16.03 |
| Complete $\mathcal{G}, N = 16$ | 2.499 | 12.40 | 5.431 | 1.387 | 0.882 | 2.262 |

# INCOMPLETE OBSERVATION

Partial observation: $\boldsymbol{x} = [\ \underbrace{\boldsymbol{x}_U}_{\text{Unobserved}}\ ,\ \underbrace{\boldsymbol{x}_O}_{\text{Observed}}\ ]$

$$l(\boldsymbol{x}_O; \boldsymbol{\theta}) = \log \sum_{\boldsymbol{x}_U} p(\boldsymbol{x}_U, \boldsymbol{x}_O; \boldsymbol{\theta}) = \underbrace{\log Z(\boldsymbol{x}_O; \boldsymbol{\theta})}_{\sum_{\boldsymbol{x}_U} \tilde{p}(\boldsymbol{x}; \boldsymbol{\theta}),\ \textit{generalize}\ \log Z(\boldsymbol{\theta})} - \underbrace{\log Z(\boldsymbol{\theta})}_{\text{0 in DAGs}}\ ,$$

$$\underbrace{\phantom{\log Z(\boldsymbol{x}_O; \boldsymbol{\theta}) - \log Z(\boldsymbol{\theta})}}_{\text{both may be est. by free energy minimization in MRFs}}$$

Connect Free Energy to Evidence Lower Bounder:

$$l(\boldsymbol{x}_O; \boldsymbol{\theta}) \geqslant -\ \underbrace{F_v(q(\boldsymbol{x}_U|\boldsymbol{x}_O))}_{\text{Variational Free Energy}}\ -\log Z(\boldsymbol{\theta})$$

$$= \mathbb{E}_{q(\boldsymbol{x}_U|\boldsymbol{x}_O)} \left[ \log \frac{p(\boldsymbol{x}_U, \boldsymbol{x}_O; \boldsymbol{\theta})}{q(\boldsymbol{x}_U|\boldsymbol{x}_O)} \right]$$

$$= \underbrace{\mathbb{E}_{q(\boldsymbol{x}_U|\boldsymbol{x}_O)} [\log p(\boldsymbol{x}_U, \boldsymbol{x}_O; \boldsymbol{\theta})] + H(q(\boldsymbol{x}_U|\boldsymbol{x}_O))}_{\text{Evidence Lower Bound } F(q, \boldsymbol{\theta})}$$

Intuition of maximizing $F(q, \boldsymbol{\theta})$

- Maximizing (incomplete) likelihood
- Minimizing free energy

This gives raise of EM as a coordinate ascent method:

$$\text{E step}: \quad q^{(t+1)} = \operatorname*{argmax}_{q} F(q, \boldsymbol{\theta}^{(t)}),$$

$$\text{M step}: \quad \boldsymbol{\theta}^{(t+1)} = \operatorname*{argmax}_{\boldsymbol{\theta}} F(q^{(t+1)}, \boldsymbol{\theta}).$$

# GENERATOR MIXED MODEL

EQUIPPING EM WITH NORMALIZING FLOWS



- Ideal case: The underline true $p^*(\boldsymbol{x})$ is in hypothesis space $\mathcal{H}$, i.e. $p^*(\boldsymbol{x}) \in \mathcal{H}$.

- Out of reach: Test $p^*(\boldsymbol{x}) \overset{?}{\in} \mathcal{H}$

- Luckily, what is at our hands is:

  $\mathcal{H}$ is large $\rightarrow$ condidate $p(\boldsymbol{x}; \boldsymbol{\theta})$ is flexible

This brings up the finite **mixture** models.

$$p(\boldsymbol{x}; \Theta) = \sum_{k=1}^{K} \pi_k p_k(\boldsymbol{x}) = \sum_{k=1}^{K} \pi_k p( \underbrace{\boldsymbol{g}(\boldsymbol{z}; \theta_k)}_{\text{Variable change via generator } \boldsymbol{g}} )$$

What to expect from GenMM:

- Flexible and expressive model, enlarging hyperspace $\mathcal{H}$
- Tractable likelihood
- Compatible with typical statistical models
- Compatible with NN tools/frameworks
- Scale to high-dimensional structured data
- Efficient in sampling (data generation)
- ...

# A HIGH-LEVEL VIEW OF GENMM: FINITE MIXTURE
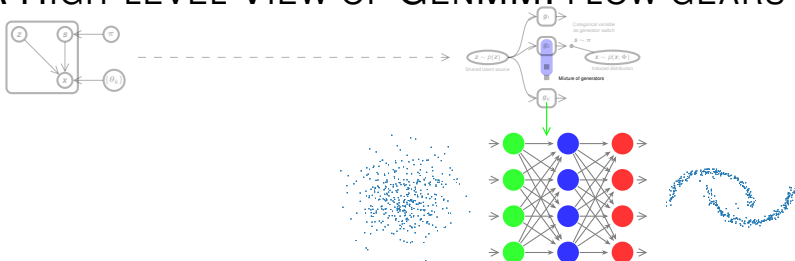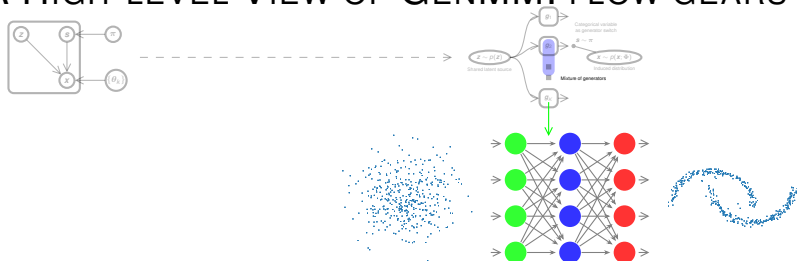
# A HIGH-LEVEL VIEW OF GENMM: FLOW GEARS



When the $k$-th generator is selected, i.e., $s_k = 1$ and $s_{k'} = 0$ for $k' \neq k$, say $\tilde{\boldsymbol{x}} = \boldsymbol{x}|_{s_k=1}$. By following the change of variable rule

$$\underbrace{p(\tilde{\boldsymbol{x}})|_{\tilde{\boldsymbol{x}}=\tilde{\boldsymbol{g}}(\boldsymbol{z})}}_{\text{Induced distribution}} = \underbrace{p(\boldsymbol{z})}_{\substack{\text{Assumed known distribution} \\ \text{easy to sample}}} \cdot \underbrace{\left| \det\left( \frac{\partial \boldsymbol{z}}{\partial \tilde{\boldsymbol{x}}} \right) \right|}_{\substack{\text{Computational load} \\ \text{depends on the mapping}}} .$$

A toy example:

Gaussian linear transform: $Z \sim \mathsf{N}(0,1) \xrightarrow{X = \sigma \cdot Z + \mu} X \sim \mathsf{N}(\mu, \sigma)$

# A HIGH-LEVEL VIEW OF GENMM: FLOW GEARS



When the $k$-th generator is selected, i.e., $s_k = 1$ and $s_{k'} = 0$ for $k' \neq k$, say $\tilde{\boldsymbol{x}} = \boldsymbol{x}|_{s_k=1}$. By following the change of variable rule

$$\underbrace{p(\tilde{\boldsymbol{x}})|_{\tilde{\boldsymbol{x}}=\tilde{\boldsymbol{g}}(\boldsymbol{z})}}_{\text{Induced distribution}} = \underbrace{p(\boldsymbol{z})}_{\substack{\text{Assumed known distribution} \\ \text{easy to sample}}} \cdot \underbrace{\left| \det\left( \frac{\partial \boldsymbol{z}}{\partial \tilde{\boldsymbol{x}}} \right) \right|}_{\substack{\text{Computational load} \\ \text{depends on the mapping}}} .$$
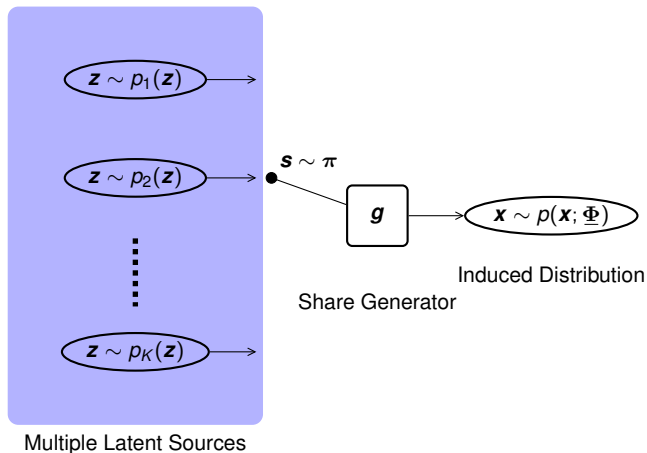
Powering it with a $L$-layer neural network implementation:

$$\boldsymbol{z} = \boldsymbol{h}_0 \; \underset{\tilde{\boldsymbol{f}}_1}{\overset{\tilde{\boldsymbol{g}}_1}{\rightleftarrows}} \; \boldsymbol{h}_1 \; \underset{\tilde{\boldsymbol{f}}_2}{\overset{\tilde{\boldsymbol{g}}_2}{\rightleftarrows}} \; \cdots\cdots \; \underset{\tilde{\boldsymbol{f}}_L}{\overset{\tilde{\boldsymbol{g}}_L}{\rightleftarrows}} \; \boldsymbol{x} = \boldsymbol{h}_L$$

# A HIGH-LEVEL VIEW OF GENMM: FLOW GEARS



When the $k$-th generator is selected, i.e., $s_k = 1$ and $s_{k'} = 0$ for $k' \neq k$, say $\tilde{\boldsymbol{x}} = \boldsymbol{x}|_{s_k=1}$. By following the change of variable rule

$$\underbrace{p(\tilde{\boldsymbol{x}})|_{\tilde{\boldsymbol{x}}=\tilde{\boldsymbol{g}}(\boldsymbol{z})}}_{\text{Induced distribution}} = \underbrace{p(\boldsymbol{z})}_{\substack{\text{Assumed known distribution} \\ \text{easy to sample} \\ \text{offloading}}} \cdot \underbrace{\left| \det\left( \frac{\partial \boldsymbol{z}}{\partial \tilde{\boldsymbol{x}}} \right) \right|}_{\substack{\text{Computational load} \\ \text{depends on the mapping}}} .$$

Layer structure:

- RealNVP
- Glow      ← - - - -
- ODE

$$\boldsymbol{z} = \boldsymbol{h}_0 \underset{\tilde{\boldsymbol{f}}_1}{\overset{\tilde{\boldsymbol{g}}_1}{\rightleftarrows}} \boldsymbol{h}_1 \underset{\tilde{\boldsymbol{f}}_2}{\overset{\tilde{\boldsymbol{g}}_2}{\rightleftarrows}} \cdots\cdots \underset{\tilde{\boldsymbol{f}}_L}{\overset{\tilde{\boldsymbol{g}}_L}{\rightleftarrows}} \boldsymbol{x} = \boldsymbol{h}_L$$

30 / 44

# LATMM: ALTERNATIVE MIXTURE



Multiple Latent Sources

Share Generator

Induced Distribution
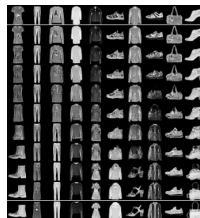
## SAMPLING EXAMPLES



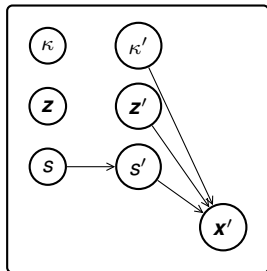Generated samples by GenMM and LatMM.



Interpolation in latent space

## APPLICATION TO CLASSIFICATION TASKS

Application to classification with maximum likelihood. Test Accuracy Table of GenMM for Classification Task

| Dataset | K=1 | K=2 | K=3 | K=4 | K=10 | K=20 | State Of Art |
|---|---|---|---|---|---|---|---|
| Letter | 0.9459 | 0.9513 | 0.9578 | 0.9581 | 0.9657 | **0.9674** | 0.9582 |
| Satimage | 0.8900 | 0.8975 | 0.9045 | 0.9085 | 0.9105 | **0.9160** | 0.9090 |
| Norb | 0.9184 | 0.9257 | 0.9406 | 0.9459 | 0.9538 | **0.9542** | 0.8920 |

## GENHMM: BRING THE CONCEPT INTO HMM



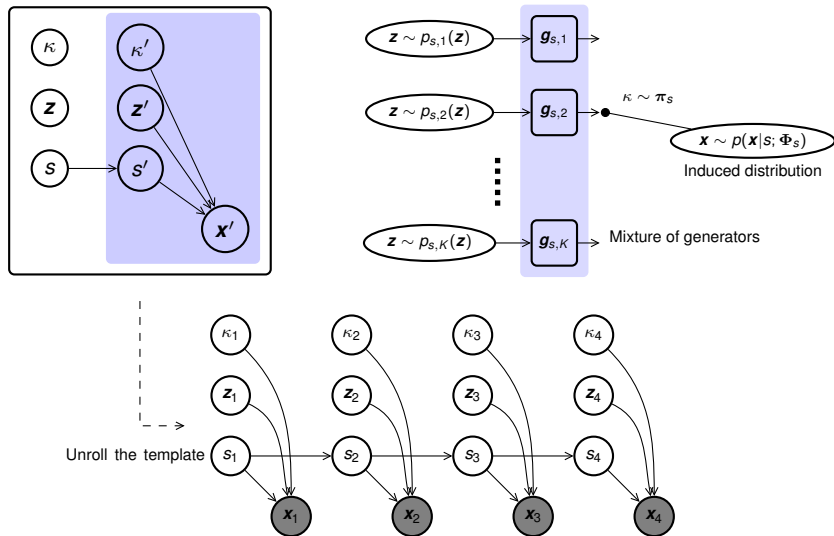At time $t$, the probabilistic model of a state $s \in \mathcal{S}$ is then given by

$$p(\boldsymbol{x}|s; \boldsymbol{\Phi}_s) = \sum_{\kappa=1}^{K} \pi_{s,\kappa} p(\boldsymbol{x}|s, \kappa; \boldsymbol{\theta}_{s,\kappa}),$$

where

- $\pi_{s,\kappa} = p(\kappa|s; \boldsymbol{H})$, naturally $\sum_{\kappa=1}^{K} \pi_{s,\kappa} = 1$
- $p(\boldsymbol{x}|s, \kappa; \boldsymbol{\theta}_{s,\kappa})$ is induced by the $k$th generator $\boldsymbol{g}_k(\boldsymbol{z}) = \boldsymbol{g}(\boldsymbol{z}; \boldsymbol{\theta}_k)$

# GENHMM: BRING THE CONCEPT INTO HMM

## ALTERNATIVE VIEW OF GENHMM



Hypothesis set of GenHMM as $\mathcal{H} := \{\boldsymbol{H} | \boldsymbol{H} = \{\mathcal{S}, \boldsymbol{q}, \boldsymbol{A}, p(\boldsymbol{x} | s; \boldsymbol{\Phi}_s)\}\}$

- $\mathcal{S}$: the set of hidden states of $\boldsymbol{H}$.
- $\boldsymbol{q} = [q_1, \cdots, q_{|\mathcal{S}|}]^{\mathsf{T}}$: the initial state distributions of $\boldsymbol{H}$. $q_i = p(s_1 = i; \boldsymbol{H})$.
- $\boldsymbol{A}$: the transition matrix of states in $\boldsymbol{H}$.
- $\boldsymbol{\Phi} = \{\boldsymbol{\Phi}_s | s \in \mathcal{S}\}$.

## ALTERNATIVE VIEW OF GENHMM



Hypothesis set of GenHMM as $\mathcal{H} := \{\boldsymbol{H} | \boldsymbol{H} = \{\mathcal{S}, \boldsymbol{q}, \boldsymbol{A}, p(\boldsymbol{x}|s; \boldsymbol{\Phi}_s)\}\}$

- $\mathcal{S}$: the set of hidden states of $\boldsymbol{H}$.
- $\boldsymbol{q} = [q_1, \cdots, q_{|\mathcal{S}|}]^\mathsf{T}$: the initial state distributions of $\boldsymbol{H}$. $q_i = p(s_1 = i; \boldsymbol{H})$.
- $\boldsymbol{A}$: the transition matrix of states in $\boldsymbol{H}$.
- $\boldsymbol{\Phi} = \{\boldsymbol{\Phi}_s | s \in \mathcal{S}\}$.

## ALTERNATIVE VIEW OF GENHMM



Hypothesis set of GenHMM as $\mathcal{H} := \{\textbf{\textit{H}}|\textbf{\textit{H}} = \{\mathcal{S}, \textbf{\textit{q}}, \textbf{\textit{A}}, p(\textbf{\textit{x}}|s; \boldsymbol{\Phi}_s)\}\}$

- $\mathcal{S}$: the set of hidden states of $\textbf{\textit{H}}$.
- $\textbf{\textit{q}} = \left[q_1, \cdots, q_{|\mathcal{S}|}\right]^{\mathsf{T}}$: the initial state distributions of $\textbf{\textit{H}}$. $q_i = p(s_1 = i; \textbf{\textit{H}})$.
- $\textbf{\textit{A}}$: the transition matrix of states in $\textbf{\textit{H}}$.
- $\boldsymbol{\Phi} = \{\boldsymbol{\Phi}_s | s \in \mathcal{S}\}$.

# APPLICATION OF GENHMM

Speech Recognition:

Phoneme classification / recognition

| Model | Criterion | K=1 | K=3 | K=5 |
|---|---|---|---|---|
| GMM-HMM | Accuracy | 62.3 | 68.0 | 68.7 |
| linear variable change | Precision | 67.9 | 72.6 | 73.0 |
| | F1 | 63.7 | 69.1 | 69.7 |
| GenHMM | Accuracy | 76.7 | **77.7** | 77.7 |
| non-linear variable change | Precision | 76.9 | **78.1** | 78.0 |
| | F1 | 76.1 | **77.1** | 77.0 |

Robustness to perturbation of noise.

| Model | Criterion | White Noise SNR | | | |
|---|---|---|---|---|---|
| | | 15dB | 20dB | 25dB | 30dB |
| GMM-HMM | Accuracy | 36.6 | 44.2 | 50.8 | 57.1 |
| | Precision | 59.2 | 64.2 | 68.4 | 70.6 |
| | F1 | 39.9 | 47.7 | 53.9 | 59.9 |
| GenHMM | Accuracy | 52.4 | 62.0 | 69.7 | **74.3** |
| | Precision | 60.0 | 65.9 | 71.7 | **74.8** |
| | F1 | 52.5 | 62.0 | 69.3 | **73.5** |

Application to sepsis detection for infants, c.f. see Section 7.5.

- generative training + discriminative training
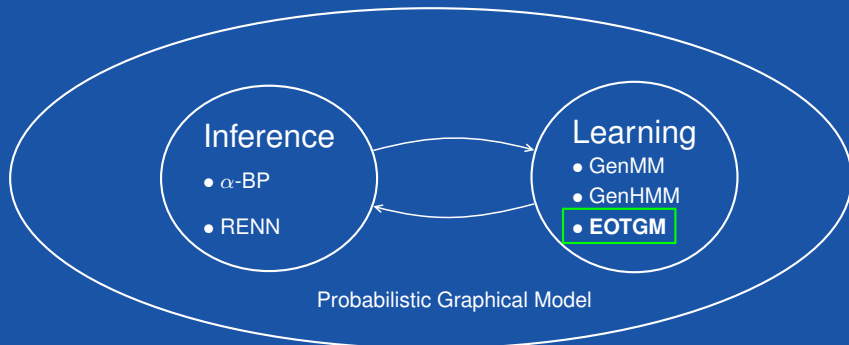- innovative feature inspired from acoustic signal feature

## REMARK ON GENMM/GENHMM
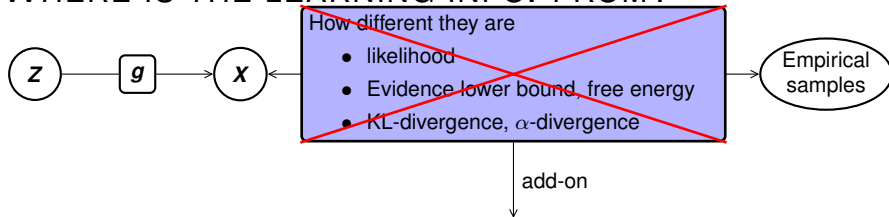
### Remarks on Learning of GenMM/GenHMM

- $F$, evidence lower bound, equivalent to a negative free energy in DAGs
  $= \mathcal{Q}$ in EM + Entropy
- E-step require inference (message-passing for posteriors)
- No optimality in M-step (NN generators, batch-size gradient descent).
- Still, guaranteed non-decreasing lklh. (c.f. Proposition 6.1, Proposition 7.1)

### Attributes

- Free dimension for flexibility: number of mixture + complexity of functional form of neural networks
- Tractable likelihood and efficient sampling
- Compatible with classic statistic methods and neural network techniques (error back-propagation, optimizer)
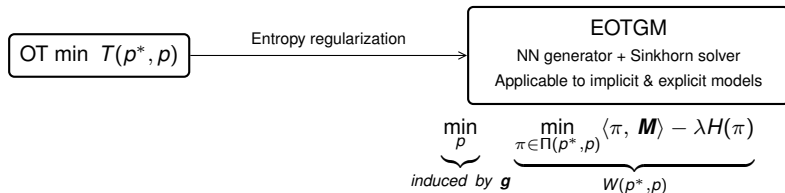
# WHERE IS THE LEARNING INFO. FROM?

$$z \longrightarrow \boxed{g} \longrightarrow x$$

How different they are
- likelihood
- Evidence lower bound, free energy
- KL-divergence, $\alpha$-divergence

Empirical samples

add-on

Optimal transport (OT): moving mass from a dist. to another

$$T(p^*, p) = \min_{\pi \in \Pi(p^*, p)} \langle \underbrace{\pi}_{\text{marginalize to } p^*, p}, \underbrace{M}_{\substack{\text{cost matrix} \\ \text{pair-wise sample difference}}} \rangle,$$

Key attributes:

- Doesnot require tractible lklh.
- Learning gradient info. from sample comparison
- High complexity, each evaluation is sovling an optimization problem

# EOTGM: EOT BASED GENERATIVE MODEL



$$\underbrace{\min_{p}}_{\substack{\text{induced by } \boldsymbol{g}}} \underbrace{\min_{\pi \in \Pi(p^*, p)} \langle \pi, \boldsymbol{M} \rangle - \lambda H(\pi)}_{W(p^*, p)}$$
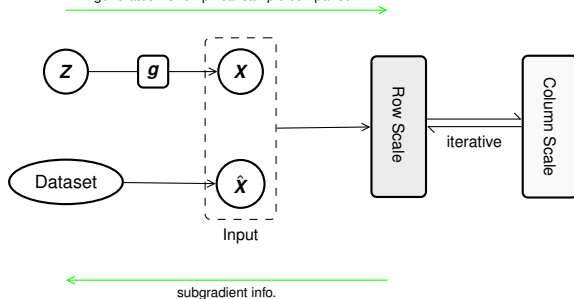
EOTGM employs:

alternatively scaling rows & columns of matrix $e^{-\boldsymbol{M}/\lambda}$ (Sinkhorn & Knopp)

to extract gradient information (sub-gradient) for generator $\boldsymbol{g}$

---

Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport

# EOTGM AND EOTGAN

EOTGM

# EOTGM AND EOTGAN

EOTGM

# EOTGM AND EOTGAN

EOTGM



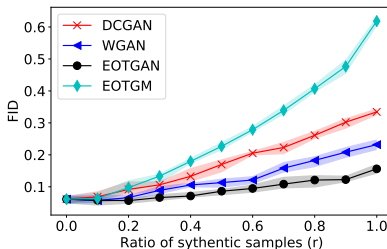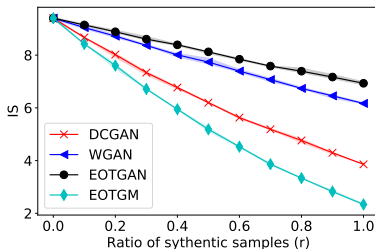generated v.s. empirical sample comparison
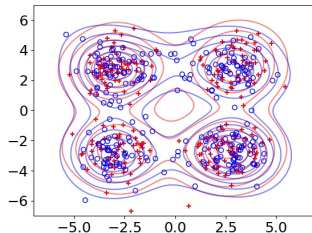
subgradient info.

EOTGAN: learn implicit distribution with feature mapping, c.f. Section 8.2.2 (Euclidean distance is not suitable for multimedia signals.)

# NUMERICAL

→ Toy distribution learning (target at 4-mixture Gaussians) using EOTGM. Real samples (red '+') and contour (red curve), versus generated samples (blue 'o') and contour (blue curve) by **g**.

↓ Comparison of semantic scores (on MNIST) versus mixing ratio r:

- IS: Inseption Score (large is good)
- FID: Frechet Inception Distance (small is good)

## SUMMARY

### Wrap-up

- Inference with message-passing and analysis
- Inference with free energy minimization by neural networks
- Inference .. Learning: their interactions
- Neural network generators in EM for more flexible modeling; A Further step into temporal models
- A bonus modeling method for likelihood-free learning

Thank you for your attention.
Q&A.