July 18, 2019

Inter-University Research Institute Corporation /
Research Organization of Information and Systems
**National Institute of Informatics**

*PRESTO*
SAKIGAKE

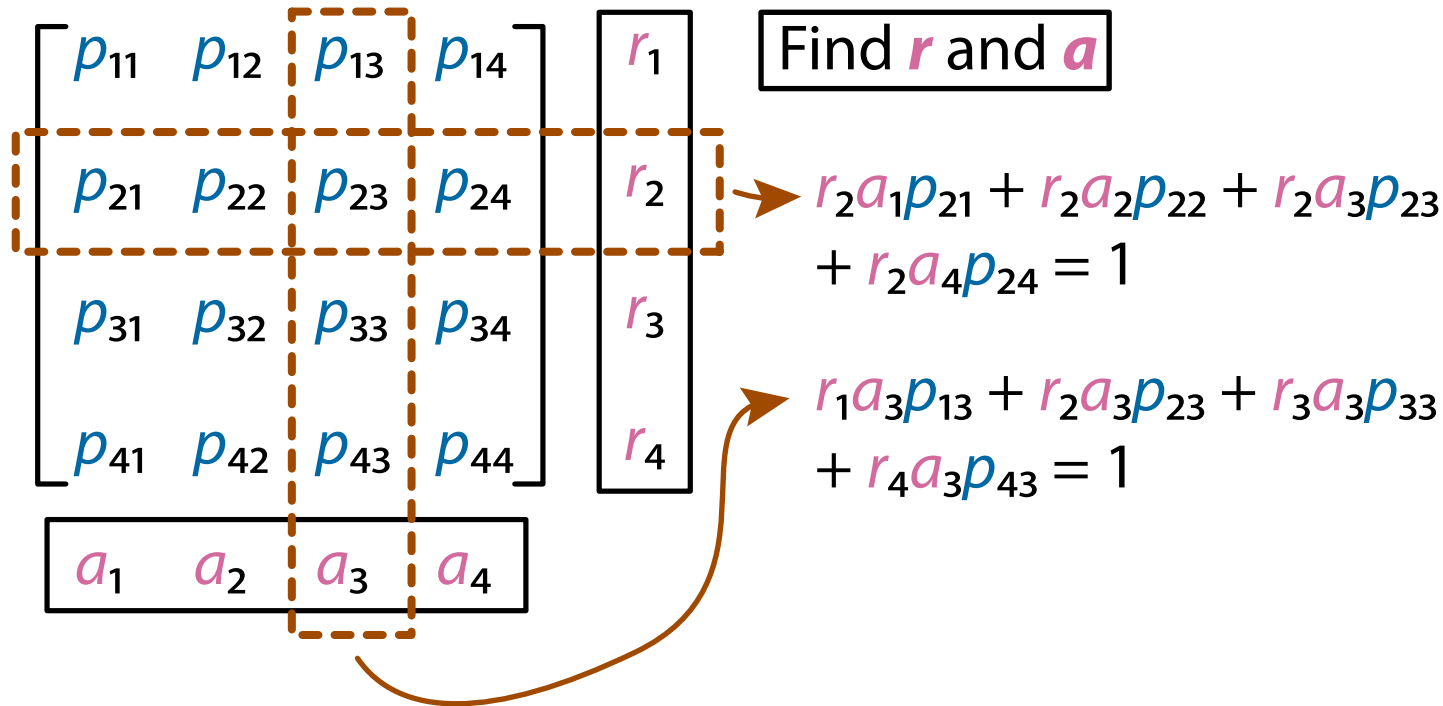# Machine Learning with Information Geometry

Mahito Sugiyama

# Summary

- We study machine learning with its applications

  - We focus on the relationship between computational processes, discrete structures, and machine learning models

- Ongoing topics

  - Machine Learning with Information Geometry
  - Machine Learning with Discrete Structure
  - Significant Pattern Mining
  - Machine Learning Applications
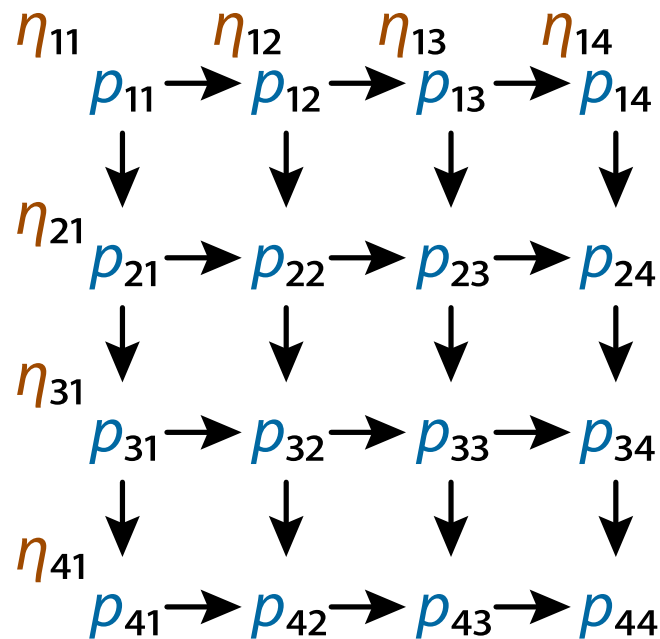
# Tensor (Matrix) Balancing

$$
\begin{bmatrix}
p_{11} & p_{12} & p_{13} & p_{14} \\
p_{21} & p_{22} & p_{23} & p_{24} \\
p_{31} & p_{32} & p_{33} & p_{34} \\
p_{41} & p_{42} & p_{43} & p_{44}
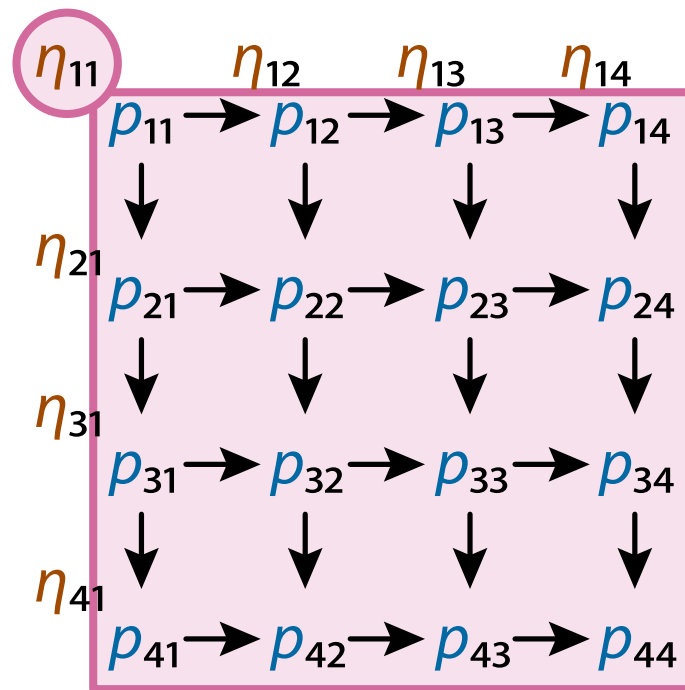\end{bmatrix}
$$

# Tensor (Matrix) Balancing



Find $r$ and $a$

$r_2 a_1 p_{21} + r_2 a_2 p_{22} + r_2 a_3 p_{23} + r_2 a_4 p_{24} = 1$

$r_1 a_3 p_{13} + r_2 a_3 p_{23} + r_3 a_3 p_{33} + r_4 a_3 p_{43} = 1$

# Introduce $\eta$

$$\begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \\ p_{41} & p_{42} & p_{43} & p_{44} \end{bmatrix} \longrightarrow$$

$\eta_{11}$ $\quad$ $\eta_{12}$ $\quad$ $\eta_{13}$ $\quad$ $\eta_{14}$

$p_{11} \rightarrow p_{12} \rightarrow p_{13} \rightarrow p_{14}$

$\eta_{21}$

$p_{21} \rightarrow p_{22} \rightarrow p_{23} \rightarrow p_{24}$

$\eta_{31}$

$p_{31} \rightarrow p_{32} \rightarrow p_{33} \rightarrow p_{34}$

$\eta_{41}$

$p_{41} \rightarrow p_{42} \rightarrow p_{43} \rightarrow p_{44}$

$$\begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \\ p_{41} & p_{42} & p_{43} & p_{44} \end{bmatrix} \longrightarrow$$
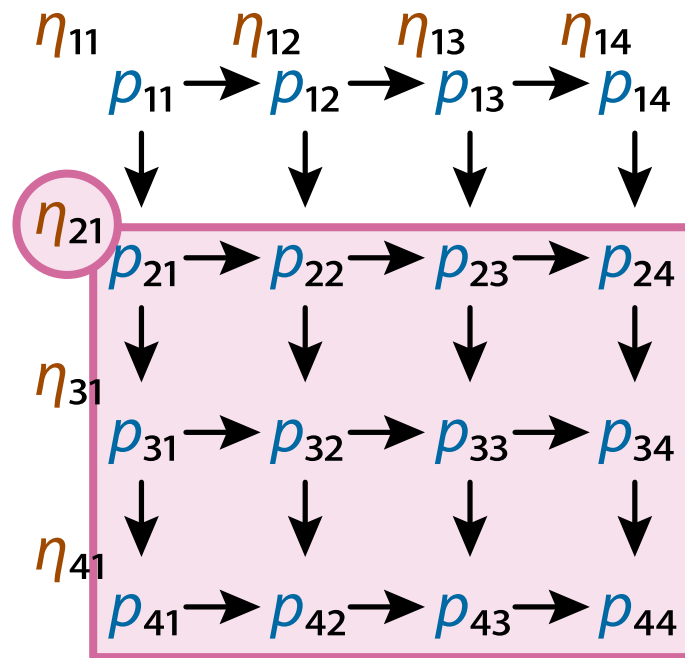
# Introduce $\eta$

$$\begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \\ p_{41} & p_{42} & p_{43} & p_{44} \end{bmatrix}$$
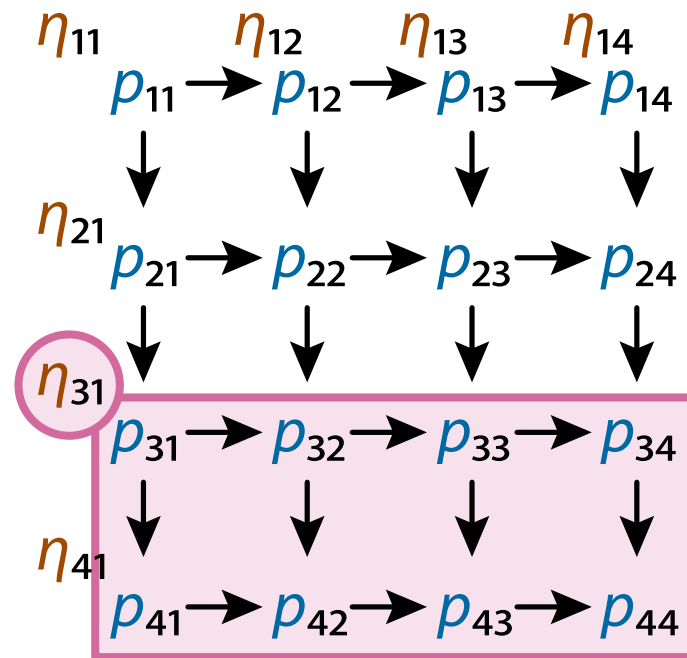
# Introduce $\eta$

$$\begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \\ p_{41} & p_{42} & p_{43} & p_{44} \end{bmatrix} \longrightarrow$$
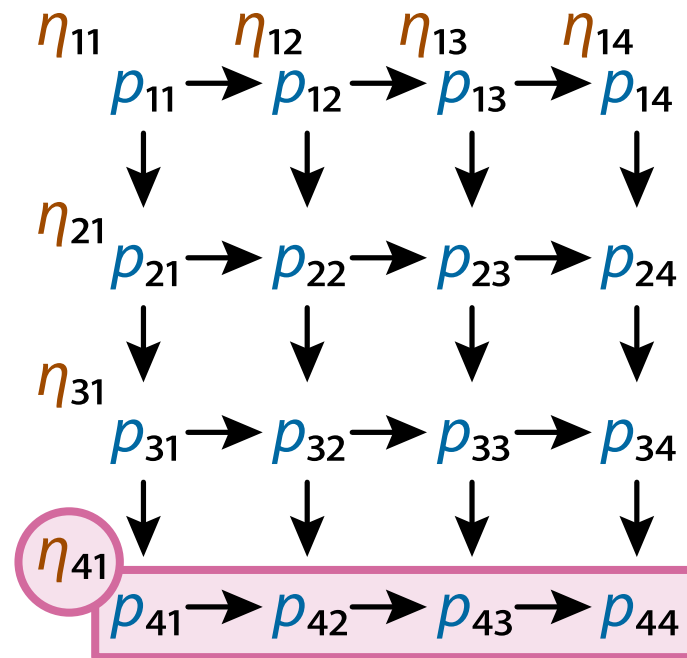
# Introduce $\eta$

$$\begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \\ p_{41} & p_{42} & p_{43} & p_{44} \end{bmatrix} \longrightarrow$$
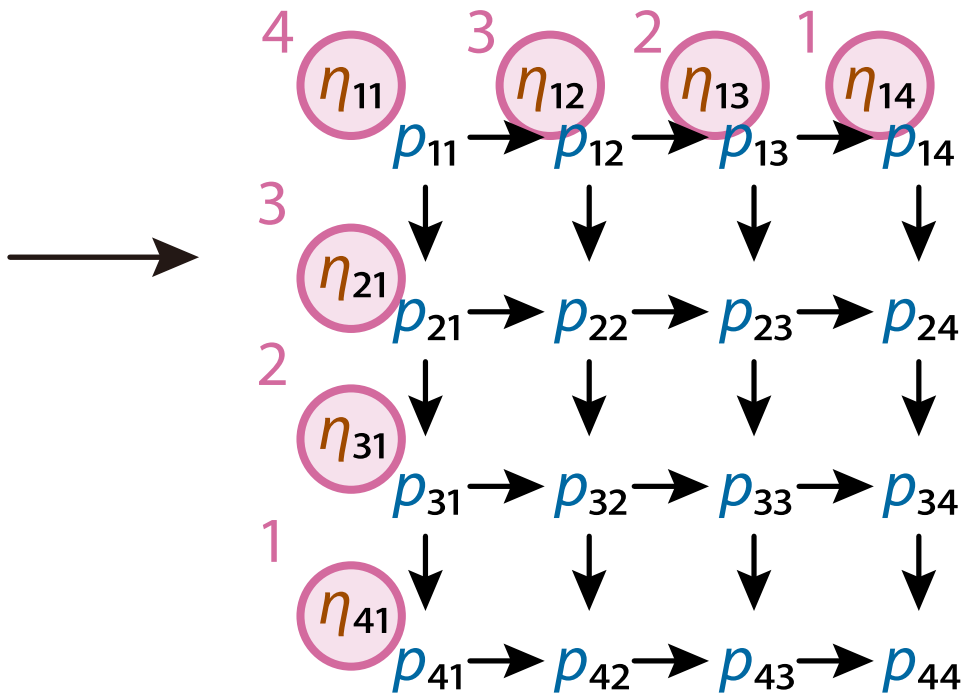
$\eta_{11}$  $\eta_{12}$  $\eta_{13}$  $\eta_{14}$

$p_{11} \rightarrow p_{12} \rightarrow p_{13} \rightarrow p_{14}$

$\eta_{21}$

$p_{21} \rightarrow p_{22} \rightarrow p_{23} \rightarrow p_{24}$

$\eta_{31}$

$p_{31} \rightarrow p_{32} \rightarrow p_{33} \rightarrow p_{34}$

$\eta_{41}$

$p_{41} \rightarrow p_{42} \rightarrow p_{43} \rightarrow p_{44}$

# Constraints on $\eta$

# Introduce $\theta$

$$
\begin{bmatrix}
p_{11} & p_{12} & p_{13} & p_{14} \\
p_{21} & p_{22} & p_{23} & p_{24} \\
p_{31} & p_{32} & p_{33} & p_{34} \\
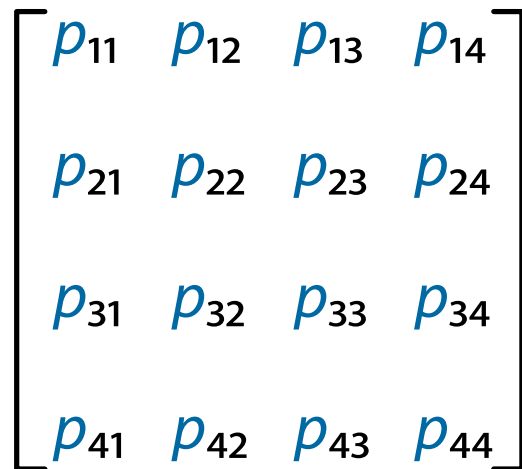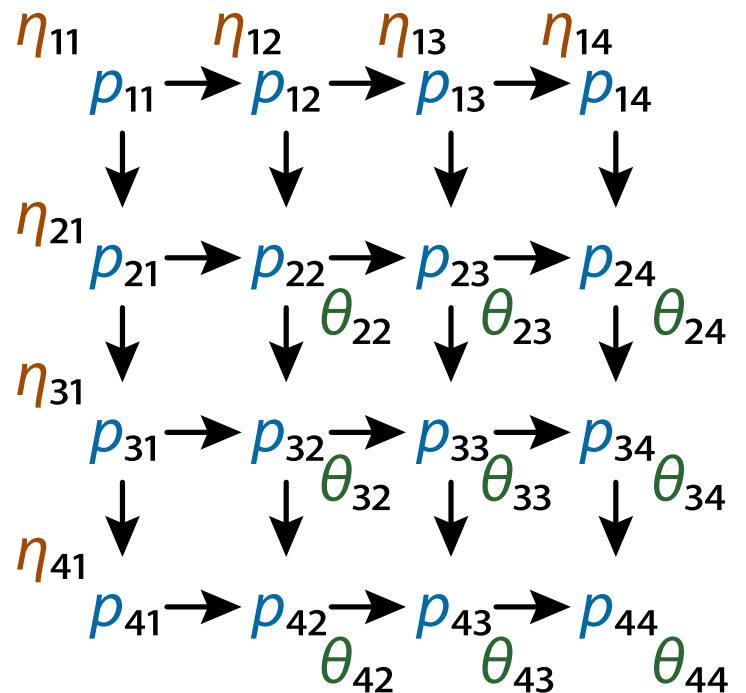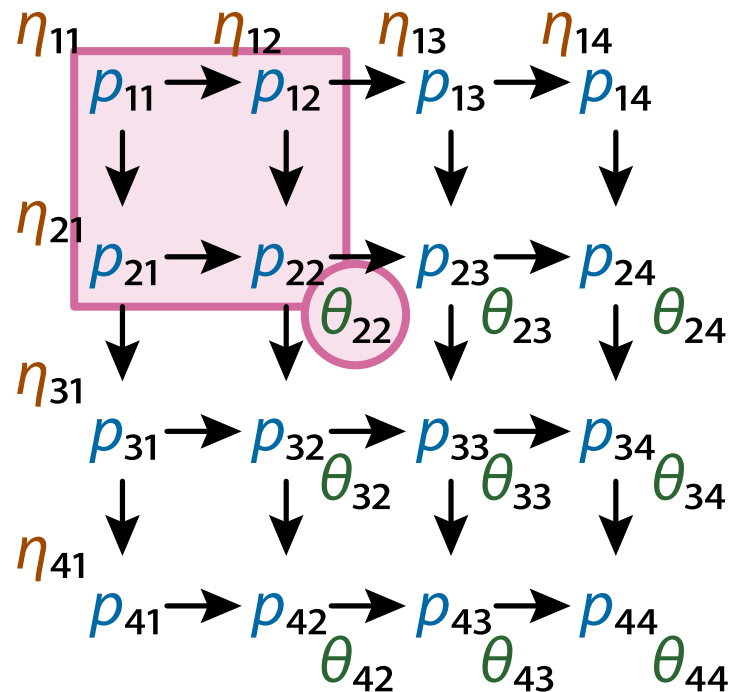p_{41} & p_{42} & p_{43} & p_{44}
\end{bmatrix}
\longrightarrow
$$

$\theta_{ij} = \log p_{ij} - \log p_{i-1j} - \log p_{ij-1} + \log p_{i-1j-1}$

$\eta_{11}$ $\quad$ $\eta_{12}$ $\quad$ $\eta_{13}$ $\quad$ $\eta_{14}$

$p_{11} \rightarrow p_{12} \rightarrow p_{13} \rightarrow p_{14}$

$\eta_{21}$

$p_{21} \rightarrow p_{22} \rightarrow p_{23} \rightarrow p_{24}$

$\theta_{22}$ $\quad$ $\theta_{23}$ $\quad$ $\theta_{24}$

$\eta_{31}$

$p_{31} \rightarrow p_{32} \rightarrow p_{33} \rightarrow p_{34}$

$\theta_{32}$ $\quad$ $\theta_{33}$ $\quad$ $\theta_{34}$

$\eta_{41}$

$p_{41} \rightarrow p_{42} \rightarrow p_{43} \rightarrow p_{44}$

$\theta_{42}$ $\quad$ $\theta_{43}$ $\quad$ $\theta_{44}$
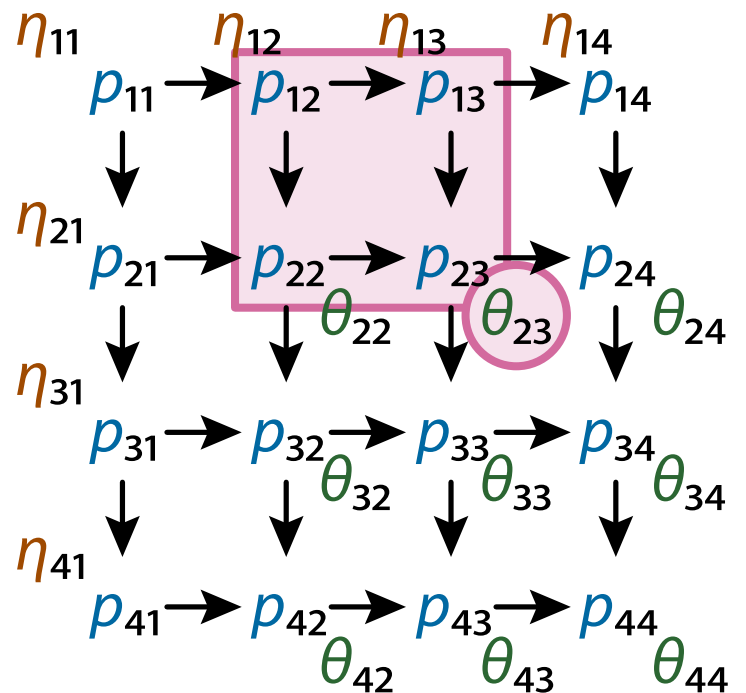
# Introduce $\theta$

$$\begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \\ p_{41} & p_{42} & p_{43} & p_{44} \end{bmatrix} \longrightarrow$$

$$\theta_{ij} = \log p_{ij} - \log p_{i-1j} - \log p_{ij-1} + \log p_{i-1j-1}$$

# Introduce $\theta$

$$\begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \\ p_{41} & p_{42} & p_{43} & p_{44} \end{bmatrix} \longrightarrow$$
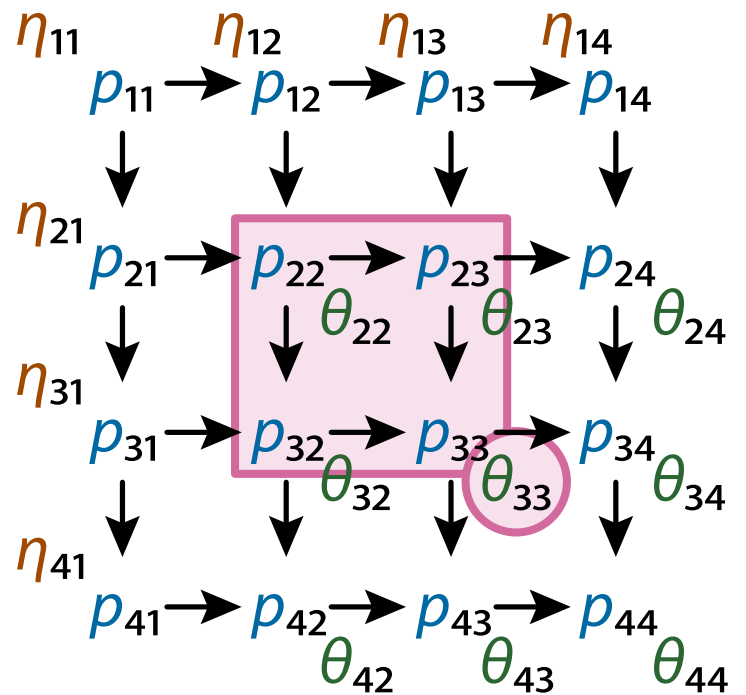
$$\theta_{ij} = \log p_{ij} - \log p_{i-1j} - \log p_{ij-1} + \log p_{i-1j-1}$$

# Introduce $\theta$



$$\begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \\ p_{41} & p_{42} & p_{43} & p_{44} \end{bmatrix}$$
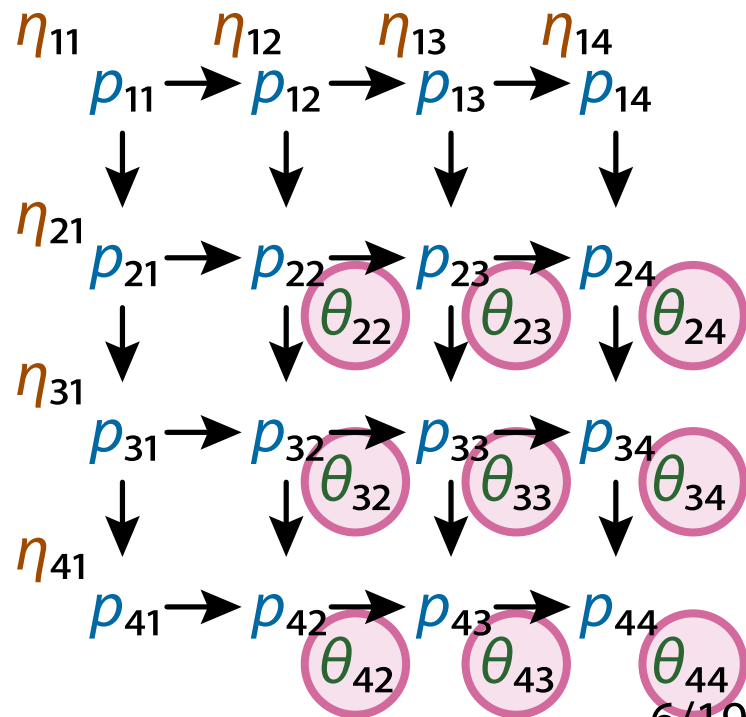
$$\theta_{ij} = \log p_{ij} - \log p_{i-1j} - \log p_{ij-1} + \log p_{i-1j-1}$$
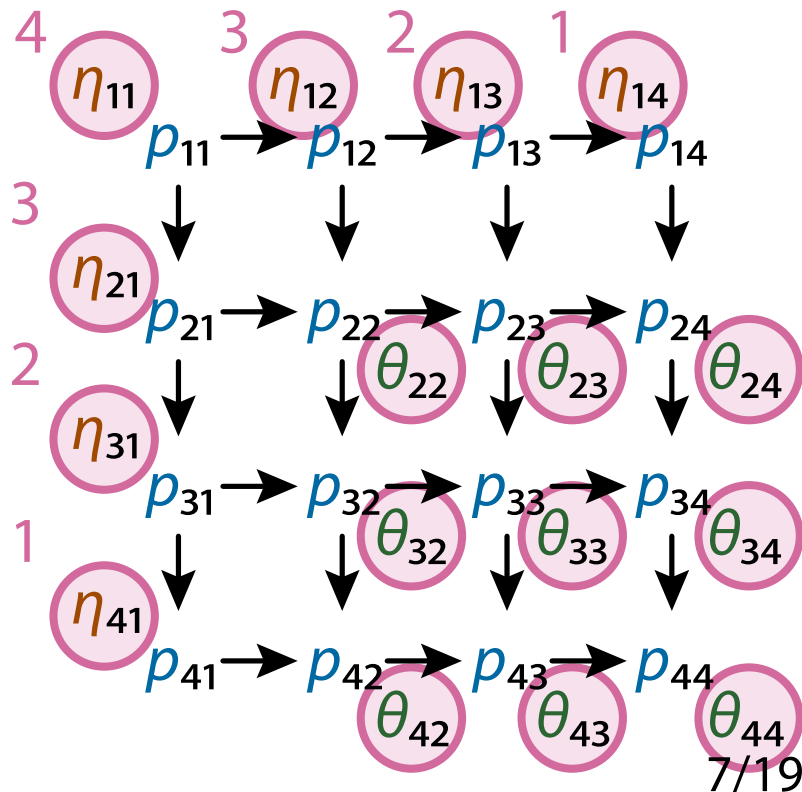
# Introduce $\theta$

$$\begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \\ p_{41} & p_{42} & p_{43} & p_{44} \end{bmatrix} \longrightarrow$$

$$\theta_{ij} = \log p_{ij} - \log p_{i-1j} - \log p_{ij-1} + \log p_{i-1j-1}$$

# Balancing as Constraints on $\eta$ and $\theta$

$$\begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \\ p_{41} & p_{42} & p_{43} & p_{44} \end{bmatrix} \longrightarrow$$

Matrix balancing is achieved iff:
$\eta_{11} = 4$, $\eta_{21} = 3$, ..., $\eta_{41} = 1$
without changing any $\theta_{ij}$

# Information Geometric View



All matrices (= Riemannian manifold)

Gradient is obtained
as Riemannian metric

Solution!

Given matrix

Balanced matrices
($\eta$ satisfied)

Achievable matrices
($\theta$ satisfied)

# Empirical Performance

# From Matrix to Poset (DAG)

$$\begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \\ p_{41} & p_{42} & p_{43} & p_{44} \end{bmatrix} \longrightarrow$$
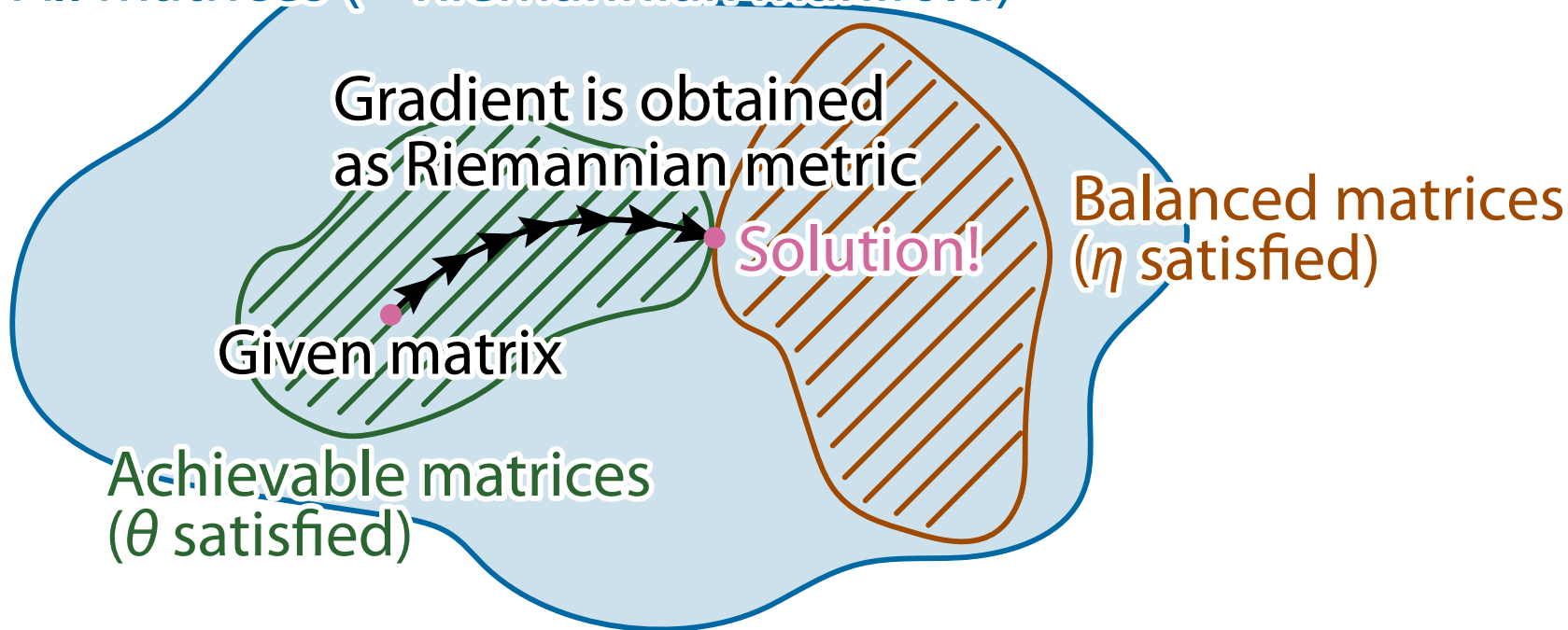
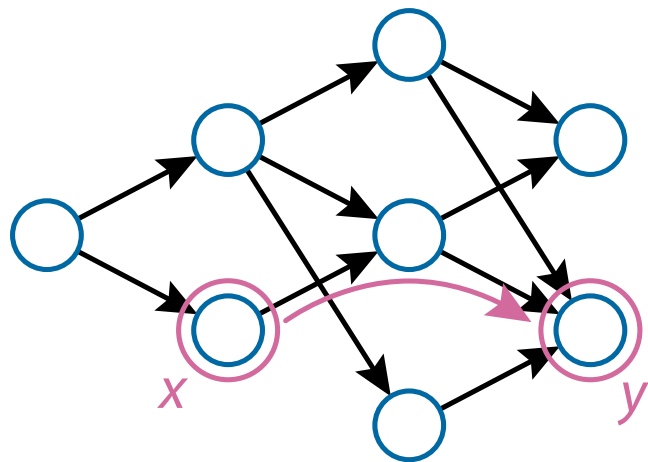# Partially Ordered Set



- Partially ordered set (poset) $(S, \leq)$
    - (i)   $x \leq x$ (reflexivity)
    - (ii)  $x \leq y, y \leq x \Rightarrow x = y$ (antisymmetry)
    - (iii) $x \leq y, y \leq z \Rightarrow x \leq z$ (transitivity)
    - – We assume that $S$ is finite and includes the least element (bottom) $\bot \in S$

- Equivalent to a DAG
    - – Each $x \in S$ is a node
    - – $x \leq y \iff y$ is reachable from $x$

# Log-Linear Model on Poset

Each $x \in S$ has a triple:
$(p(x), \theta(x), \eta(x))$



- A probability vector $p : S \rightarrow (0, 1)$
  s.t. $\sum_{x \in S} p(x) = 1$
  - (Normalized) weight for each node

- We introduce $\theta : S \rightarrow \mathbb{R}$ and $\eta : S \rightarrow \mathbb{R}$ as

$$\log p(x) = \sum_{s \leq x} \theta(s),$$

$$\eta(x) = \sum_{s \geq x} p(s)$$

# Log-Linear Model on Poset

Each $x \in S$ has a triple:
($p(x)$, $\theta(x)$, $\eta(x)$)



- A probability vector $p: S \rightarrow (0, 1)$
  s.t. $\sum_{x \in S} p(x) = 1$
  - (Normalized) weight for each node

- We introduce $\theta: S \rightarrow \mathbb{R}$ and $\eta: S \rightarrow \mathbb{R}$ as

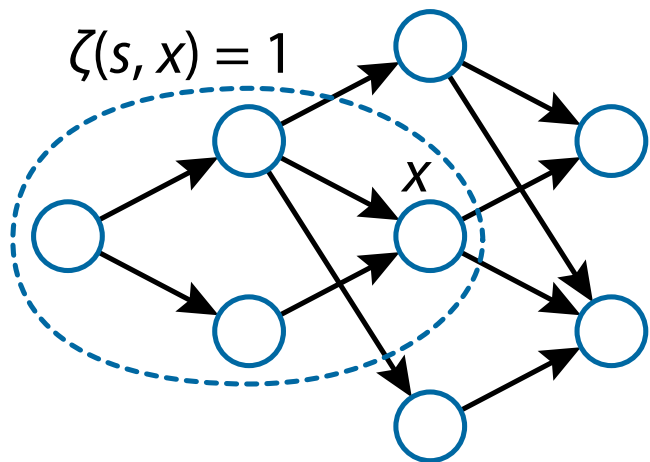$$\log p(x) = \sum_{s \leq x} \theta(s), \ \theta(x) = \sum_{s \in S} \mu(s, x) \log p(s)$$

$$\eta(x) = \sum_{s \geq x} p(s), \ p(x) = \sum_{s \in S} \mu(x, s) \eta(s)$$

# Möbius Function



- Zeta function $\zeta{:}S \times S \to \{0, 1\}$

$$\zeta(s, x) = \begin{cases} 1 & \text{if } s \leq x, \\ 0 & \text{otherwise.} \end{cases}$$

- Möbius function $\mu{:}S \times S \to \mathbb{Z}$

$$\mu(x, y) = \begin{cases} 1 & \text{if } x = y, \\ -\sum_{x \leq s < y} \mu(x, s) & \text{if } x < y, \\ 0 & \text{otherwise.} \end{cases}$$

  - We have $\zeta\mu = I$, that is;
    $\sum_{s \in S} \zeta(s, y)\mu(x, s) = \sum_{x \leq s \leq y} \mu(x, s) = \delta_{xy}$

# Möbius Function Is Generalization of Inclusion-Exclusion Principle

- For sets $A, B, C$,

$$|A \cup B \cup C| = |A| + |B| + |C| - |A \cap B| - |B \cap C| - |A \cap C| + |A \cap B \cap C|$$

- In general, for $A_1, A_2, \ldots, A_n$,

$$\left| \bigcup_i A_i \right| = \sum_{J \subseteq \{1,\ldots,n\}, J \neq \varnothing} (-1)^{|J|-1} \left| \bigcap_{j \in J} A_j \right|$$

- The Möbius function $\mu$ is the generalization of "$(-1)^{|J|-1}$"

# Riemannian Manifold with Info. Geometry



$$D_{\mathrm{KL}}[P, R] = D_{\mathrm{KL}}[P, Q] + D_{\mathrm{KL}}[Q, R]$$

Solution $Q$

$R$   $\eta$   $e$-flat submanifold

*e*-projection = Maximum Likelihood Estimation

*m*-flat submanifold

$P$

$\theta$

Riemannian metric:

$$\frac{\partial \eta(x)}{\partial \theta(y)} = \sum_{s \in S} \zeta(x, s) \zeta(y, s) p(s) - \eta(x) \eta(y)$$

$$\frac{\partial \theta(x)}{\partial \eta(y)} = \sum_{s \in S} \mu(s, x) \mu(s, y) p(s)^{-1}$$

# Example of Learning Prob. Dist. (1/2)

Dataset

| | Bread | Milk | Apple |
|------|:-----:|:----:|:-----:|
| ID 1 | ○ | × | × |
| ID 2 | ○ | ○ | × |
| ID 3 | ○ | × | × |
| ID 4 | × | ○ | ○ |
| ID 5 | × | ○ | ○ |
| ID 6 | ○ | ○ | × |
| ID 7 | ○ | × | × |
| ID 8 | ○ | ○ | × |

Sample space
(Lattice)

Bread,Milk,
Apple
Freq: 0/8 = 0

Bread,Milk
Freq: 3/8 = 0.375

Bread,Apple
Freq: 0/8 = 0

Milk,Apple
Freq: 2/8 = 0.25

Bread
Freq: 6/8 = 0.75

Milk
Freq: 5/8 = 0.625

Apple
Freq: 2/8 = 0.25

Freq = $\eta$

∅

# Example of Learning Prob. Dist. (2/2)

$$\log(\text{prob.}) = -10.41 + 9.43[\text{Bread}] + 8.52[\text{Milk}] - 9.84[\text{Apple}]$$
$$- 9.03[\text{Bread,Milk}] + 9.43[\text{Milk,Apple}]$$

MLE

Parameter $= \theta$

Boltzmann
machine



| Bread | Milk | Apple | Prob. from data | Learned prob. |
|:---:|:---:|:---:|:---:|:---:|
| × | × | × | ? | 0.0000300109 |
| ○ | × | × | 0.375 | 0.3749599867 |
| × | ○ | × | ? | 0.1499903954 |
| × | × | ○ | ? | 0.0000000016 |
| ○ | ○ | × | 0.375 | 0.2250096042 |
| ○ | × | ○ | ? | 0.0000200043 |
| × | ○ | ○ | 0.25 | 0.0999895960 |
| ○ | ○ | ○ | ? | 0.1500004008 |

# Summary of Our Approach



Distribution + Partial order structure (DAG) = Manifold in Information Geometry

Parameter $\theta$

Projection

Parameter $\eta$

Data

Projection: MLE, balancing

# Conclusion

- We have established information geometric formulation for partial order structures

  - Learning process can be achieved as a projection in the parameter space (dually flat manifold)

- We have studied several applications

  - Sugiyama, M., Nakahara, H., Tsuda, K., **Tensor Balancing on Statistical Manifold**, ICML2017
  - Sugiyama, M., Nakahara, H., Tsuda, K., **Legendre Decomposition for Tensors**, NeurIPS2018
  - **Truncated Boltzmann machines** (submitted)