# Partial Order Structure and Information Geometry
## (順序構造と情報幾何)

Mahito Sugiyama (ISIR, Osaka University, PRESTO)
（杉山 麿人; 大阪大学産業科学研究所, JST さきがけ）

# Today's Model on Poset ($S, \leq$)

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

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

# Today's Model on Poset ($S, \leq$)



Probability

Zeta function

Coefficient of log-linear model
(Bias/weight in Boltzmann machines)
(Natural parameter of exponential family)

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

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

Möbius function

Expectation
(Frequency in pattern mining)
(Sufficient statistics in exponential family)

# Outcome

- Given a poset $(S, \leq)$ and consider distributions on $S$

  – The least element $\perp \in S$

1. KL divergence decomposition:

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

   with $Q$ s.t. $\theta_Q(x) = \theta_R(x)$ or $\eta_Q(x) = \eta_P(x)$ for all $x \in S \setminus \{\perp\}$

2. The set of probability distributions on $(S, \leq)$ is
   a dually flat manifold w.r.t. $\theta$ and $\eta$

   – $p$, $\theta$, and $\eta$ are coordinate systems
   – $\theta$ and $\eta$ are orthogonal
   – $\theta$ introduces the structure of exponential family
   – $\eta$ introduces the structure of mixture family

# Partially Ordered Sets

{x, y, z} ○  Power set

{x, y} ○  {x, z} ○  {y, z} ○

{x} ○  {y} ○  {z} ○

∅ ○

# Partially Ordered Sets

{x, y, z} ◯  Power set

{x, y} ◯  {x, z} ◯  {y, z} ◯

{x} ◯  {y} ◯  {z} ◯

∅ ◯

Positive integers

◯ 3

◯ 2

◯ 1

◯ 0

# Partially Ordered Sets

{x, y, z}  Power set

{x, y}  {x, z}  {y, z}

{x}  {y}  {z}

∅

Positive integers

3

2

1

0

Prefixes

000  001  010  011  100  101  110  111

00  01  10  11

0  1

λ

# Partially Ordered Sets

$\{x, y, z\}$    Power set

$\{x, y\}$  $\{x, z\}$  $\{y, z\}$

$\{x\}$  $\{y\}$  $\{z\}$

$\varnothing$

Directed Acyclic Graph

Positive integers

3

2

1

0

Prefixes

000  001 010  011 100  101 110  111

00  01  10  11

0  1

$\lambda$

# Posets with Probability Distribution

Probability distribution
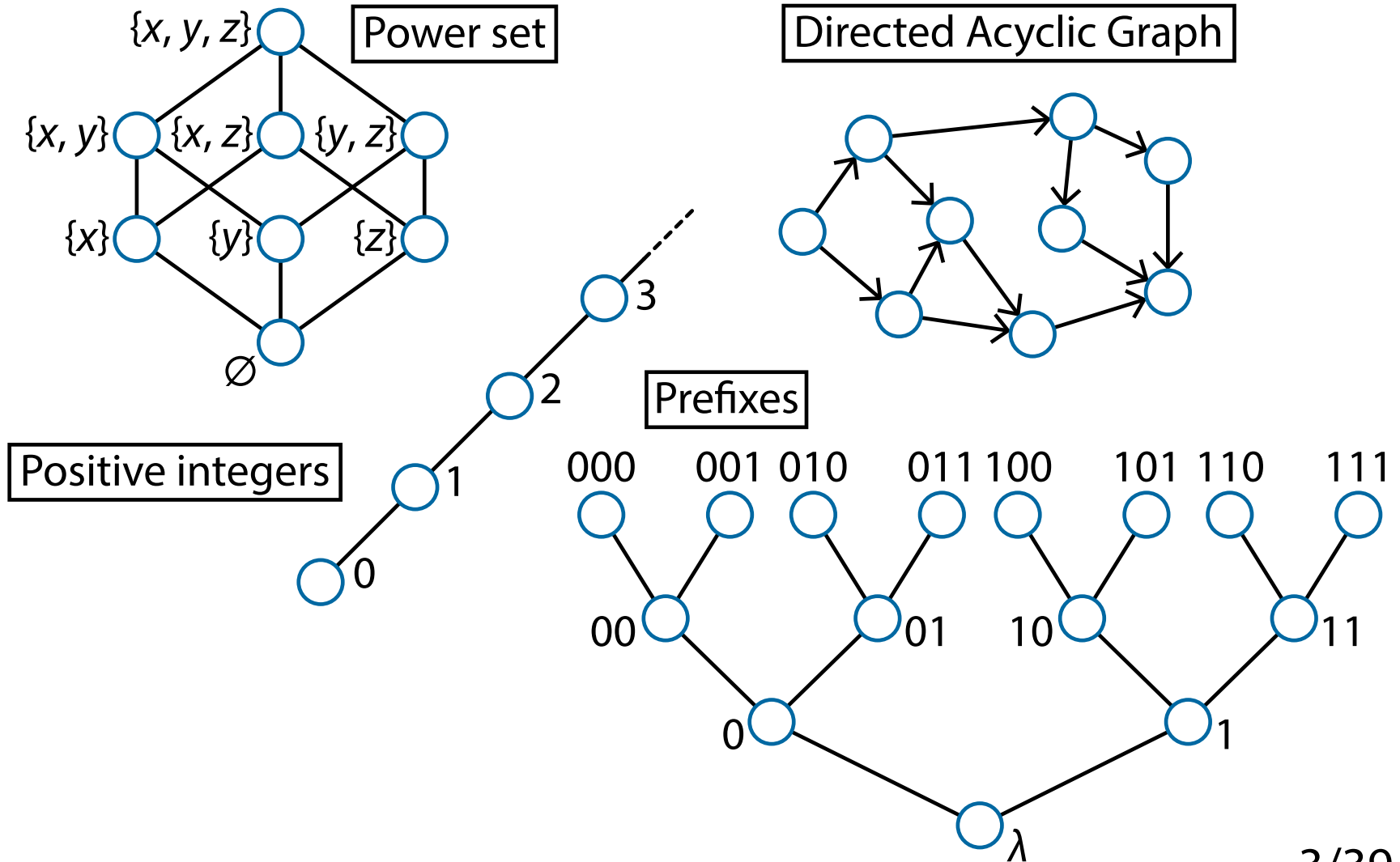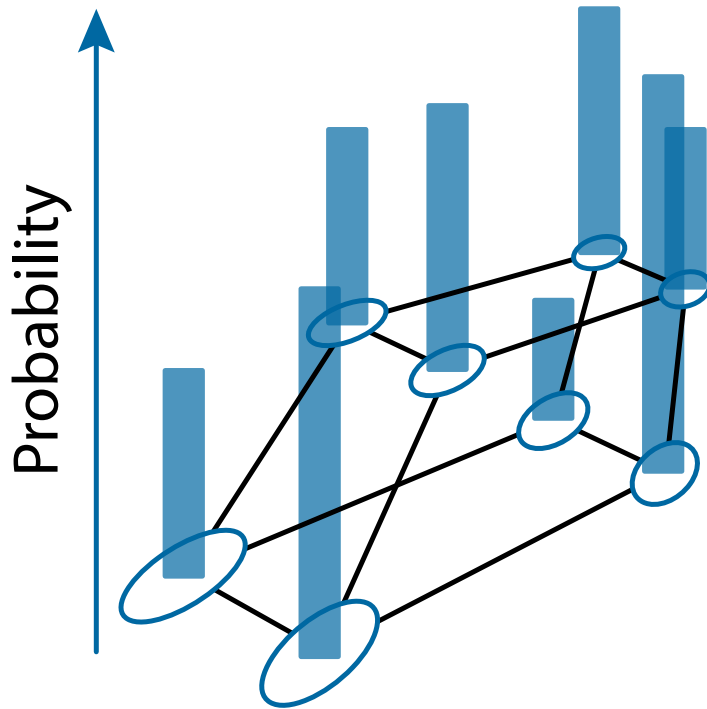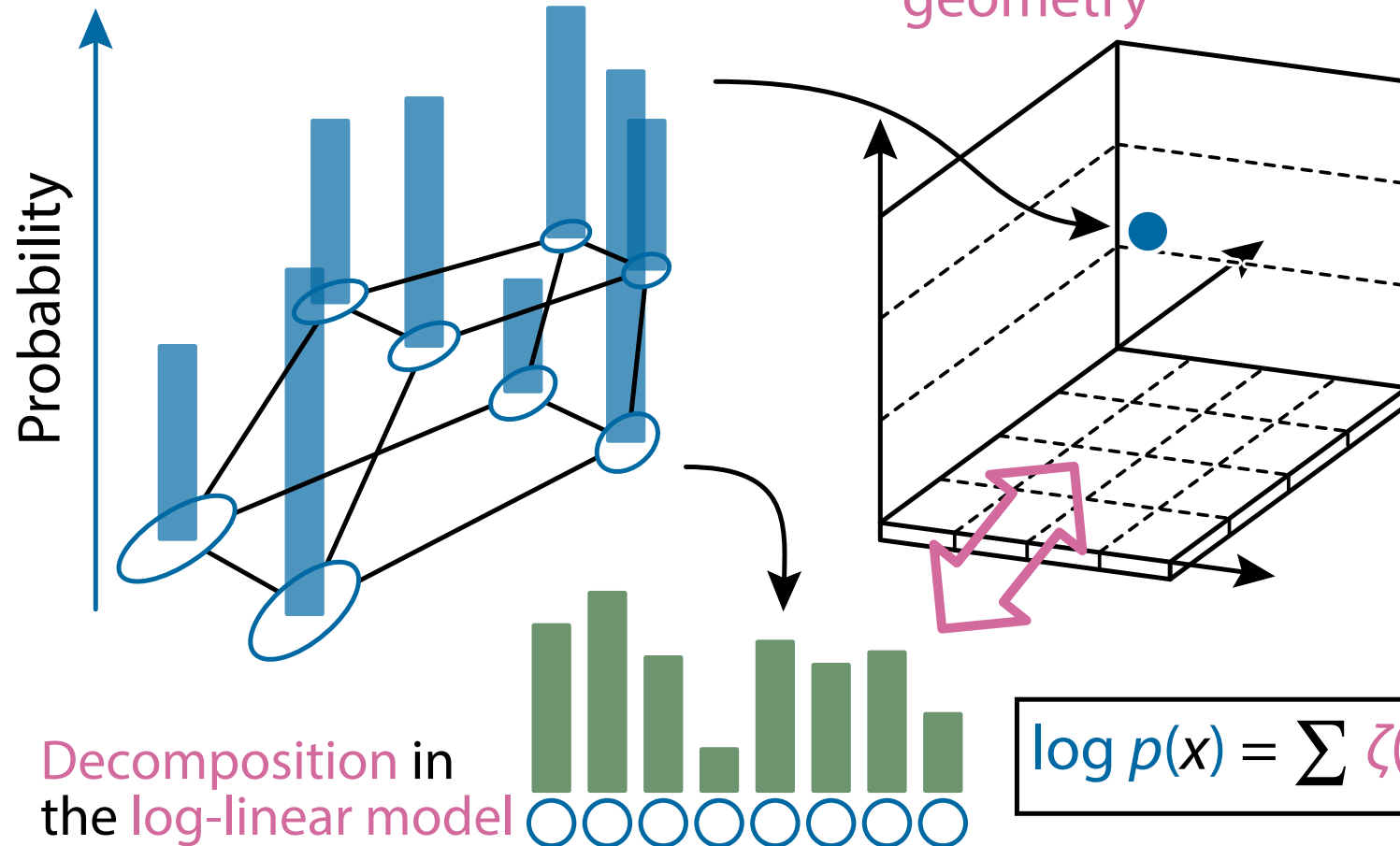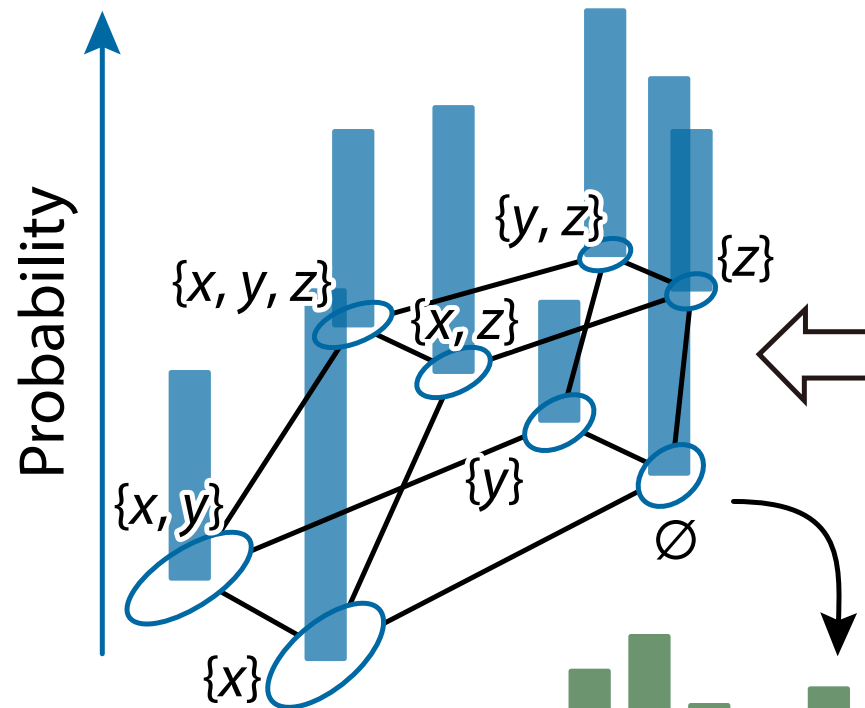on posets (partially ordered sets)

# Posets with Probability Distribution

Probability distribution
on posets (partially ordered sets)

Information geometry

Probability

Decomposition in
the log-linear model

$$\log p(x) = \sum \zeta(s, x) \theta(s)$$

# Posets with Probability Distribution

Probability distribution on posets (partially ordered sets)



Probability

$\{x, y, z\}$  $\{y, z\}$  $\{x, z\}$  $\{z\}$

$\{x, y\}$  $\{y\}$  $\varnothing$

$\{x\}$

Decomposition in the log-linear model

| $x$ | $y$ | $z$ | (e.g. Neurons, SNPs, …) |
|---|---|---|---|
| ○ | ○ | ○ | … |
| 0 | 0 | 1 | … |
| 1 | 0 | 0 | … |
| 1 | 1 | 1 | … |
| 0 | 0 | 0 | … |
| 1 | 1 | 0 | … |
| 0 | 1 | 1 | … |
| 1 | 0 | 1 | … |
| 1 | 0 | 1 | … |
| 1 | 0 | 1 | … |
| 1 | 1 | 0 | … |

Numerical score (KL divergence) and the *p*-value for higher-order intractions

$$\log p(x) = \sum \zeta(s, x) \theta(s)$$

Binary vectors (Transaction database)

|       | 🔵 | 🔴 | 🟢 |
|-------|-----|-----|-----|
| ID 1: | 1 | 1 | 0 |
| ID 2: | 1 | 1 | 1 |
| ID 3: | 1 | 1 | 0 |
| ID 4: | 1 | 1 | 1 |
| ID 5: | 1 | 1 | 0 |
| ID 6: | 1 | 0 | 1 |
| ID 7: | 1 | 0 | 1 |
| ID 8: | 1 | 1 | 1 |
| ID 9: | 1 | 0 | 0 |
| ID10: | 0 | 1 | 0 |

Poset (itemset lattice)

{🔵,🔴,🟢}

{🔵,🔴}   {🔵,🟢}   {🔴,🟢}

{🔵}   {🔴}   {🟢}

∅

Binary vectors (Transaction database)

ID 1: 1 1 0
ID 2: 1 1 1
ID 3: 1 1 0
ID 4: 1 1 1
ID 5: 1 1 0
ID 6: 1 0 1
ID 7: 1 0 1
ID 8: 1 1 1
ID 9: 1 0 0
ID10: 0 1 0

Poset (itemset lattice)

Frequency = 0.3

{🔵,🔴,🟢}

{🔵,🔴} 0.6   {🔵,🟢} 0.5   {🔴,🟢} 0.3

{🔵} 0.9   {🔴} 0.7   {🟢} 0.5

∅ 1.0

# Binary vectors (Transaction database)

| | 🔵 | 🔴 | 🟢 |
|------|---|---|---|
| ID 1: | 1 | 1 | 0 |
| ID 2: | 1 | 1 | 1 |
| ID 3: | 1 | 1 | 0 |
| ID 4: | 1 | 1 | 1 |
| ID 5: | 1 | 1 | 0 |
| ID 6: | 1 | 0 | 1 |
| ID 7: | 1 | 0 | 1 |
| ID 8: | 1 | 1 | 1 |
| ID 9: | 1 | 0 | 0 |
| ID10: | 0 | 1 | 0 |

## Poset (itemset lattice)

{🔵,🔴,🟢}  Frequency = 0.3

{🔵,🔴} 0.6   {🔵,🟢} 0.5   {🔴,🟢} 0.3

{🔵} 0.9   {🔴} 0.7   {🟢} 0.5

∅ 1.0

Binary vectors (Transaction database)

Poset (itemset lattice)

{🔵,🔴,🟢}  Frequency = 0.3
Probability = 0.3

| | 🔵 | 🔴 | 🟢 |
|---|---|---|---|
| ID 1: | 1 | 1 | 0 |
| ID 2: | 1 | 1 | 1 |
| ID 3: | 1 | 1 | 0 |
| ID 4: | 1 | 1 | 1 |
| ID 5: | 1 | 1 | 0 |
| ID 6: | 1 | 0 | 1 |
| ID 7: | 1 | 0 | 1 |
| ID 8: | 1 | 1 | 1 |
| ID 9: | 1 | 0 | 0 |
| ID10: | 0 | 1 | 0 |

{🔵,🔴}  0.6  0.3
{🔵,🟢}  0.5  0.2
{🔴,🟢}  0.3  0.0

{🔵}  0.9  0.1
{🔴}  0.7  0.1
{🟢}  0.5  0.0

∅  1.0, 0.0

*Upward =*
*Pattern mining*

Poset (itemset lattice)

{🔵,🔴,🟢} Frequency = 0.3
Probability = 0.3

{🔵,🔴} {🔵,🟢} {🔴,🟢}
0.6 0.5 0.3
0.3 0.2 0.0

{🔵} {🔴} {🟢}
0.9 0.7 0.5
0.1 0.1 0.0

$\eta$: Frequency
$p$: Probability

∅ 1.0, 0.0

$\eta(\{🔵,🔴\}) = p(\{🔵,🔴\}) + p(\{🔵,🔴,🟢\})$

*Upward = Pattern mining*
*Downward = Log-linear analysis*

Poset (itemset lattice)

Frequency = 0.3
Probability = 0.3

$\eta$: Frequency
$p$: Probability
$\theta$: Coefficient of **log-linear model**

$\eta(\{🔵,🔴\}) = p(\{🔵,🔴\}) + p(\{🔵,🔴,🟢\})$

$\log p(\{🔵,🔴\}) = \theta(\{🔵,🔴\}) + \theta(\{🔵\}) + \theta(\{🔴\}) + \theta(\varnothing)$

$$\log p(x) = \sum \zeta(s, x)\theta(s)$$

$\{ \textcolor{cyan}{\bullet}, \textcolor{red}{\bullet}, \textcolor{green}{\bullet} \}$

$\{ \textcolor{cyan}{\bullet} \ \textcolor{red}{\bullet} \}$    $\{ \textcolor{cyan}{\bullet} \ \textcolor{green}{\bullet} \}$    $\{ \textcolor{red}{\bullet}, \textcolor{green}{\bullet} \}$

$\{ \textcolor{cyan}{\bullet} \}$    $\{ \textcolor{red}{\bullet} \}$    $\{ \textcolor{green}{\bullet} \}$

$\emptyset$

e.g. Gaussian

$$\log p(x) = \sum \zeta(s, x) \theta(s)$$

Natural parameter

Exponential family:
$$p(x) = \exp\left( \sum \theta(s) F_s(x) - \psi(\theta) \right)$$

10/39

$$\eta(x) = \sum \zeta(x, s)p(s)$$

$$\eta(x) = \mathbb{E}[\, F_x(s)\, ]$$

**Sufficient statistics** of exponential family

$\{ \bullet, \bullet, \bullet \}$

$\{ \bullet \bullet \}$  $\{ \bullet \bullet \}$  $\{ \bullet, \bullet \}$

$\{ \bullet \}$  $\{ \bullet \}$  $\{ \bullet \}$

$\varnothing$

e.g. Gaussian

$$\log p(x) = \sum \zeta(s, x)\theta(s)$$

**Natural parameter**

Exponential family:  $p(x) = \exp\!\big(\sum \theta(s)F_s(x) - \psi(\theta)\big)$

# Möbius Inversion on Posets

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

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

- Möbius function $\mu : S \times S \to \mathbb{Z}$, defined as $\mu = \zeta^{-1}$:

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

- The Möbius inversion formula [Rota (1964)]:

$$g(x) = \sum_{s \in S} \zeta(s, x) f(s) \iff f(x) = \sum_{s \in S} \mu(s, x) g(s)$$

# 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 $A_1, A_2, \ldots, A_n$,

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

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

# Mathematical Formulation

- Log-linear model and its sufficient statistics:

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

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

  – Generalization of the log-linear model on binary vectors:

$$\log p(\boldsymbol{x}) = \sum_{i} \theta^{i} x^{i} + \sum_{i<j} \theta^{ij} x^{i} x^{j} + \cdots + \theta^{1\ldots n} x^{1} x^{2} \ldots x^{n},$$

- From the Möbius inversion formula,

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

Triple for each node

$p$
$\eta$
$\theta$

{🔵,🔴}

0.2
0.2
−1.79

{🔵}

{🔴}

0.3
0.5
1.10

0.4
0.6
1.39

∅

0.1
1.0
−2.30

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

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

Triple for each node

$p$
$\eta$
$\theta$

{🔵,🔴}

0.2
0.2
−1.79

{🔵}

{🔴}

0.3
0.5
1.10

0.4
0.6
1.39

∅

0.1
1.0
−2.30

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

$p(\{🔴\})$

$p(\{🔵,🔴\})$

$p(\{🔵\})$

# Triple for each node

$\begin{array}{c} p \\ \eta \\ \theta \end{array}$

{🔵,🔴}

0.2
0.2
−1.79

{🔵}

{🔴}

0.3
0.5
1.10

0.4
0.6
1.39

∅

0.1
1.0
−2.30

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

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

# Probability distribution is a "point" in 3D space



$p(\{🔴\})$

$p(\{🔵,🔴\})$

0.3

0.4

0.2

$p(\{🔵\})$

# Triple for each node

$p$
$\eta$
$\theta$

$\{$●,●$\}$

0.2
0.2
−1.79

$\{$●$\}$                    $\{$●$\}$

0.3                        0.4
0.5                        0.6
1.10                       1.39

∅

0.1
1.0
−2.30

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

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

# Probability distribution is a "point" in 3D space

$\eta(\{$●$\})$

0.5

0.6

0.2

$\eta(\{$●,●$\})$

$\eta(\{$●$\})$

# Triple for each node

$p$
$\eta$
$\theta$

$\{\bullet,\bullet\}$

$\{\bullet\}$
0.2
0.2
−1.79

$\{\bullet\}$

0.3
0.5
1.10

0.4
0.6
1.39

∅

0.1
1.0
−2.30
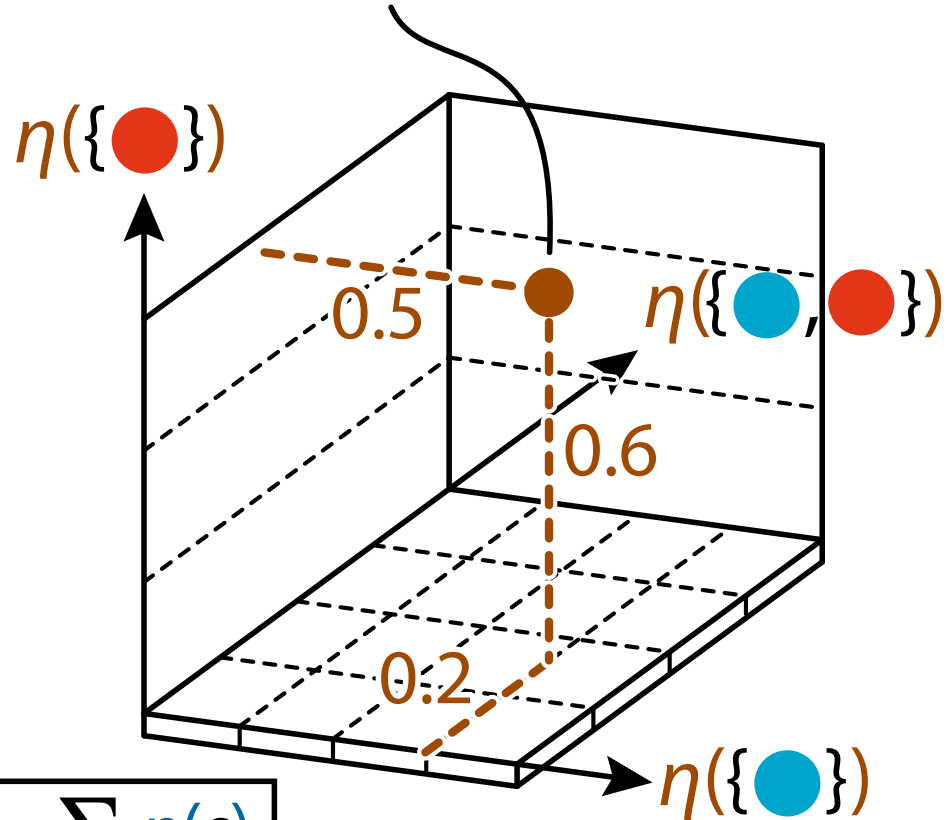
$$\eta(x) = \sum_{s \geq x} p(s)$$
$$\log p(x) = \sum_{s \leq x} \theta(s)$$

# Probability distribution is a "point" in 3D space



$\theta(\{\bullet\})$

$\theta(\{\bullet,\bullet\})$

1.10

1.39

−1.79

$\theta(\{\bullet\})$

Triple for each node

$p$
$\eta$
$\theta$

$\{\bullet, \bullet\}$

0.2
0.2
−1.79

$\{\bullet\}$

0.3
0.5
1.10

∅

$\{\bullet\}$

0.4
0.6
1.39

0.1
1.0
−2.30

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

one-to-one

$p$

$\eta \longleftrightarrow \theta$

$\{\bullet\}$

$\{\bullet, \bullet\}$

$\{\bullet\}$

# Triple for each node

$p$
$\eta$
$\theta$

$\{ \bullet , \bullet \}$

0.2
0.2
−1.79

$\{ \bullet \}$

0.3
0.5
1.10

$\{ \bullet \}$

0.4
0.6
1.39

$\varnothing$

0.1
1.0
−2.30

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

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

one-to-one

$p$

$\eta \longleftrightarrow \theta$

$\theta(\{ \bullet \})$

$\{ \bullet , \bullet \}$

$\eta(\{ \bullet \})$

*θ and η are dually orthogonal*

# Orthogonality of $\theta$ and $\eta$

- From Möbius inversion,

$$\sum_{s \in S} \zeta(x, s)\mu(s, y) = \delta_{x,y}, \qquad \delta_{x,y} = \begin{cases} 1 & \text{if } x = y, \\ 0 & \text{otherwise} \end{cases}$$

- $\theta$ and $\eta$ are dually orthogonal:

$$\mathrm{E}\left[ \frac{\partial}{\partial\theta(x)} \log p(s) \frac{\partial}{\partial\eta(y)} \log p(s) \right] = \sum_{s \in S} \zeta(x, s)\mu(s, y) = \delta_{x,y}$$

- Partial order structure leads to the same dually flat structure with the exponential family

# Existing Approach Limited To Power Set



{x, y, z}    Power set

{x, y}  {x, z}  {y, z}

{x}  {y}  {z}

∅

# Our Approach Applies To Any Posets



{x, y, z}

Subset of power set

{x, y}          {y, z}

{x}    {y}

∅

# Our Approach Applies To Any Posets

$\{x, y, z\}$

Subset of power set

Directed Acyclic Graph

$\{x, y\}$     $\{y, z\}$

$\{x\}$    $\{y\}$

$\varnothing$

3

2

Positive integers

1

0

Prefixes

000    001   010    011   100    101   110    111

00     01    10     11

0     1

$\lambda$

# KL Divergence Decomposition

- KL divergence decomposition:

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

with $Q$ s.t. $\theta_Q(x) = \theta_R(x)$ or $\eta_Q(x) = \eta_P(x)$ for all $x \in S$

- $Q$ is called the mixed distribution of $(P, R)$
- It is known as the (generalized) Pythagoras theorem in Information Geometry

- We can derive from Möbius inversion:

$$D_{KL}[P, Q] + D_{KL}[Q, R] - D_{KL}[P, R]$$

$$= \sum_{s \in S} \left( \eta_Q(s) - \eta_P(s) \right) \left( \theta_Q(s) - \theta_R(s) \right)$$

Dist. P

| | |
|---|---|
| $p$ | |
| $\eta$ | |
| $\theta$ | |

0.2
0.2
−1.79

0.3
0.5
1.10

0.4
0.6
1.39

∅

Dist. R

0.1
0.1
−1.95

0.7
0.8
1.95

0.1
0.2
0.0

∅

Dist. *P*

$p$
$\eta$
$\theta$

0.2
0.2
−1.79

0.3
0.5
1.10
∅

0.4
0.6
1.39

choose η

Mixed
distribution *Q*

???
0.2
???

???
???
1.95
∅

???
???
0.0

Dist. *R*

0.1
0.1
−1.95

0.7
0.8
1.95
∅

0.1
0.2
0.0

choose θ

24/39

Dist. *P*

*p*
*η*
*θ*

0.2
0.2
−1.79

0.3          0.4
0.5          0.6
1.10  ∅    1.39

choose *η*

Mixed
distribution *Q*

0.2
0.2
−1.13

0.62          0.089
0.5           0.6
1.95   ∅     0.0

Dist. *R*

0.1
0.1
−1.95

0.7          0.1
0.8          0.2
1.95  ∅    0.0

choose *θ*

Dist. P

p
η
θ

0.2
0.2
−1.79

0.3
0.5
1.10
∅

0.4
0.6
1.39

choose η

KL[P, R]

Dist. R

0.1
0.1
−1.95

0.7
0.8
1.95
∅

0.1
0.2
0.0

choose θ

Mixed
distribution Q

0.2
0.2
−1.13

0.62
0.5
1.95
∅

0.089
0.6
0.0

Dist. *P*

$p$
$\eta$
$\theta$

0.2
0.2
−1.

$KL[P, R] =$
$KL[P, Q] + KL[Q, R]$

Nonnegative decomposition
of the KL divergence

Dist. *R*

0.1
0.1
1.95

0.3
0.5
1.10

0.6
1.39

0.8
1.95

0.1
0.2
0.0

∅

choose $\eta$

Mixed
distribution *Q*

0.2
0.2
−1.13

choose $\theta$

0.62
0.5
1.95

0.089
0.6
0.0

∅

Dist. $P$

$p$
$\eta$
$\theta$

0.2
0.2
−1.

0.3
0.5
1.10 ∅

0.6
1.39

choose $\eta$

0.4390 =
0.3946 + 0.0444

Nonnegative decomposition
of the KL divergence

Dist. $R$

0.1
0.1
1.95

0.1
0.2
0.0

0.8
1.95 ∅

choose $\theta$

Mixed
distribution $Q$

0.2
0.2
−1.13

0.62
0.5
1.95 ∅

0.089
0.6
0.0

28/39

Dist. $P$

$p$
$\eta$
$\theta$

0.2
0.2
−1.79

0.3
0.5
1.10

0.4
0.6
1.39

∅

Uniform dist. $P_0$

0.25
0.25
0.0

0.25
0.5
0.0

0.25
0.5
0.0

∅

Dist. $P$

$p$
$\eta$
$\theta$

0.2
0.2
−1.79

0.3
0.5
1.10

0.4
0.6
1.39

∅

choose $\eta$

Mixed
distribution $Q$

∅

???
???
0.0

???
0.5
???

???
0.6
???

Uniform dist. $P_0$

0.25
0.25
0.0

0.25
0.5
0.0

0.25
0.5
0.0

∅

choose $\theta$
(KNOCK DOWN)

Log-linear model
$\log p(x) = \sum_{s \le x} \theta(s)$

Dist. $P$
Uniform dist. $P_0$

Contribution of the node
$= KL[P, Q] = 0.086$

$p$
$\eta$
$\theta$

0.2
0.2
−1.79

0.25
0.25
0.0

0.3
0.5
1.10

0.4
0.6
1.39

0.25
0.5
0.0

0.25
0.5
0.0

∅

choose $\eta$

0.3
0.3
0.0

∅

choose $\theta$
(KNOCK DOWN)

Mixed
distribution $R$

0.2
0.5
0.0

0.3
0.6
0.405

∅

Log-linear model
$\log p(x) = \sum_{s \leq x} \theta(s)$

Dist. $P$

Uniform dist. $P_0$

$p$
$\eta$
$\theta$

Contribution of the node
$= KL[P, Q] = 0.086$

0.2
0.2
−1.79

0.25

0.3
0.5
1.10

0.4
0.6
1.39

∅

choose

The statistics $\lambda$:
$\lambda = 2\cdot$[sample size]$\cdot KL[P, Q]$
follows $\chi^2$-distribution
with d.f. [#nodes − 1]
$\Rightarrow$ $p$-value can be obtained!

Mixed
distribution $R$

0.3
0.3
0.0

choose 0
(KNOCK DOWN)

0.2
0.5
0.0

0.3
0.6
0.405

∅

Log-linear model
$\log p(x) = \sum_{s \leq x} \theta(s)$

31/39

# Log-Linear Model on Subgraphs



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

Sufficient statistics
of exponential family

Log-linear model:
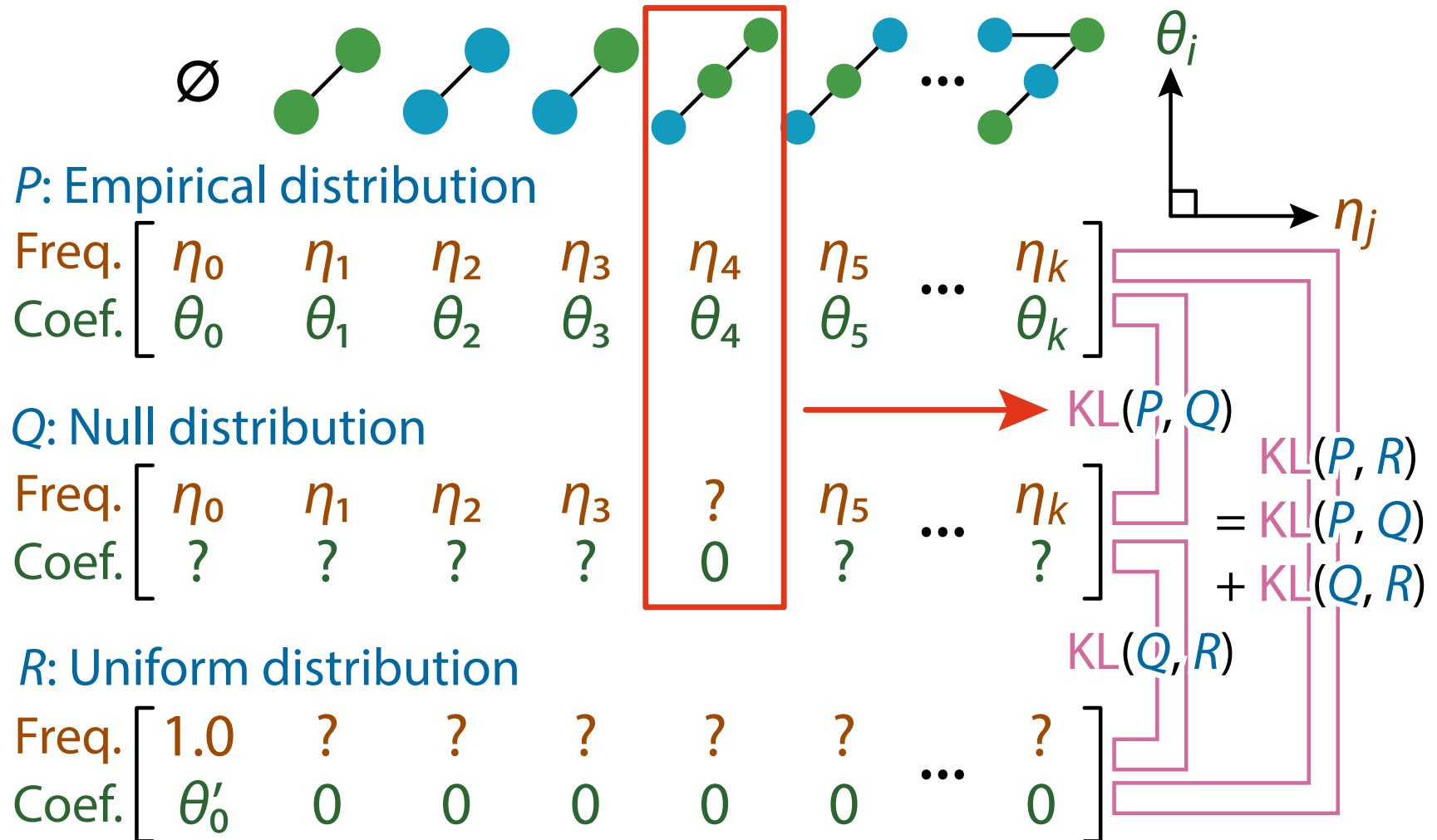$$\log p(x) = \sum_{s \sqsubseteq x} \theta(s)$$

Natural parameter
of exponential family

# Information of Each Subgraph



$P$: Empirical distribution

Freq. $\begin{bmatrix} \eta_0 & \eta_1 & \eta_2 & \eta_3 & \eta_4 & \eta_5 & \cdots & \eta_k \\ \theta_0 & \theta_1 & \theta_2 & \theta_3 & \theta_4 & \theta_5 & & \theta_k \end{bmatrix}$
Coef.

# Information of Each Subgraph

# Information of Each Subgraph



$\emptyset$

$\theta_i$

$\eta_j$

**$P$: Empirical distribution**

$$\begin{array}{l} \text{Freq.} \\ \text{Coef.} \end{array} \begin{bmatrix} \eta_0 & \eta_1 & \eta_2 & \eta_3 & \eta_4 & \eta_5 & \cdots & \eta_k \\ \theta_0 & \theta_1 & \theta_2 & \theta_3 & \theta_4 & \theta_5 & & \theta_k \end{bmatrix}$$

$\text{KL}(P, Q)$

$\text{KL}(P, R)$

**$Q$: Null distribution**

$$\begin{array}{l} \text{Freq.} \\ \text{Coef.} \end{array} \begin{bmatrix} \eta_0 & \eta_1 & \eta_2 & \eta_3 & ? & \eta_5 & \cdots & \eta_k \\ ? & ? & ? & ? & 0 & ? & & ? \end{bmatrix}$$

$= \text{KL}(P, Q)$

$+ \text{KL}(Q, R)$

**$R$: Uniform distribution**

$\text{KL}(Q, R)$

$$\begin{array}{l} \text{Freq.} \\ \text{Coef.} \end{array} \begin{bmatrix} 1.0 & ? & ? & ? & ? & ? & & ? \\ \theta_0' & 0 & 0 & 0 & 0 & 0 & \cdots & 0 \end{bmatrix}$$

# Make a Poset from Data

Dataset

🔵 🔴 🟢

ID 1:  1  1  0

ID 2:  1  1  1

ID 3:  1  1  0

ID 4:  1  1  1

ID 5:  1  1  0

ID 6:  1  0  1

ID 7:  1  0  1

ID 8:  1  1  1

ID 9:  1  0  0

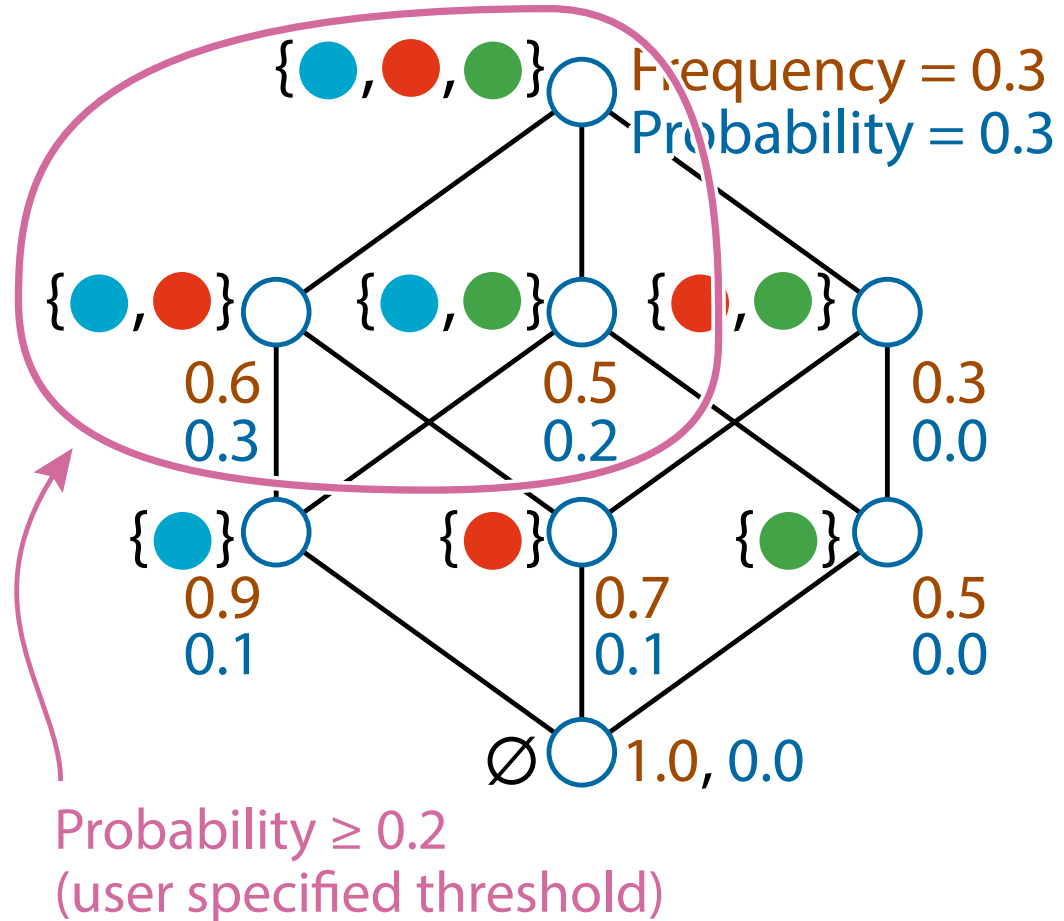ID10:  0  1  0

{🔵,🔴,🟢}   Frequency = 0.3
             Probability = 0.3

{🔵,🔴}   {🔵,🟢}   {🔴,🟢}
  0.6       0.5       0.3
  0.3       0.2       0.0

{🔵}       {🔴}       {🟢}
  0.9       0.7       0.5
  0.1       0.1       0.0

∅   1.0, 0.0

Number of nodes = $2^{\#features}$
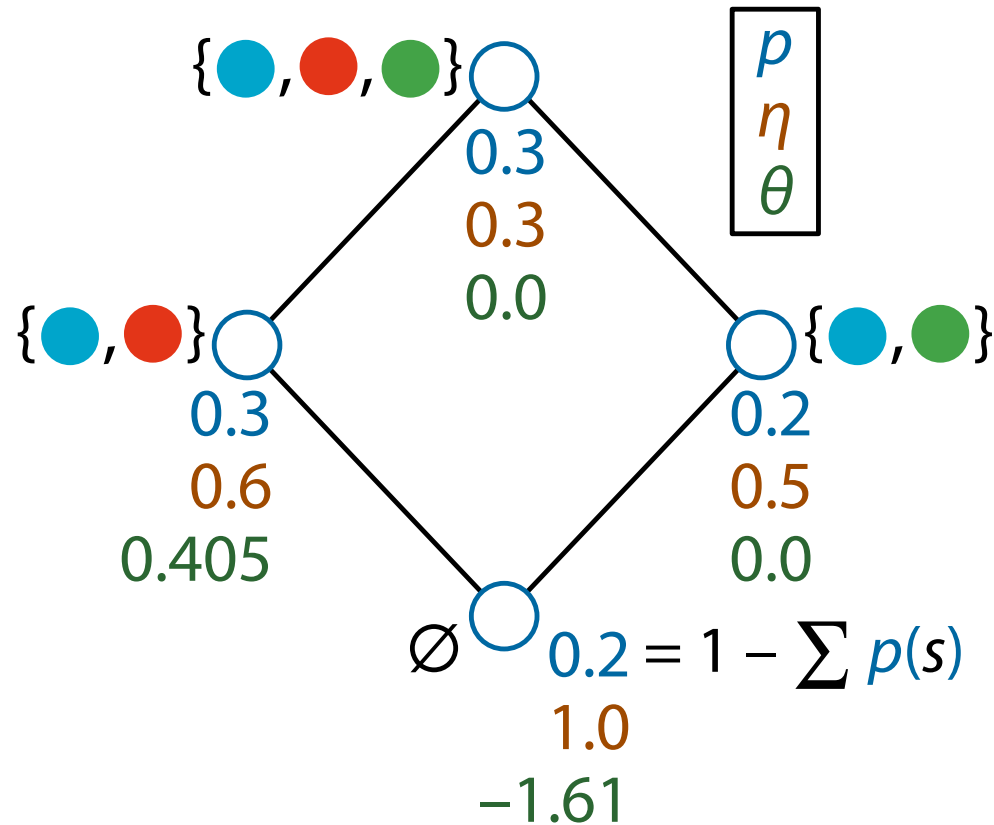⇒ combinatorial explosion!

# Make a Poset from Data



Dataset

ID 1: 1 1 0
ID 2: 1 1 1
ID 3: 1 1 0
ID 4: 1 1 1
ID 5: 1 1 0
ID 6: 1 0 1
ID 7: 1 0 1
ID 8: 1 1 1
ID 9: 1 0 0
ID10: 0 1 0

{●,●,●} Frequency = 0.3
Probability = 0.3

{●,●} 0.6 0.3
{●,●} 0.5 0.2
{●,●} 0.3 0.0

{●} 0.9 0.1
{●} 0.7 0.1
{●} 0.5 0.0

∅ 1.0, 0.0

Probability ≥ 0.2
(user specified threshold)

# Remove Nodes with Probability 0
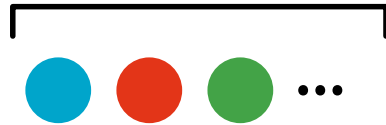
Dataset



ID 1:  1  1  0

ID 2:  1  1  1

ID 3:  1  1  0

ID 4:  1  1  1

ID 5:  1  1  0

ID 6:  1  0  1

ID 7:  1  0  1

ID 8:  1  1  1

ID 9:  1  0  0

ID10:  0  1  0

$\{\bullet,\bullet,\bullet\}$

0.3
0.3
0.0

$p$
$\eta$
$\theta$

$\{\bullet,\bullet\}$

0.3
0.6
0.405

$\{\bullet,\bullet\}$

0.2
0.5
0.0

$\varnothing$  $0.2 = 1 - \sum p(s)$

1.0
−1.61

# Example on Real Data (kosarak)

# features: 41,270

🔵 🔴 🟢 ...

| ID 1: | 1 | 1 | 0 | |
| ID 2: | 1 | 1 | 1 | |
| ID 3: | 1 | 1 | 0 | ... |
| ID 4: | 1 | 1 | 1 | |
| ID 5: | 1 | 1 | 0 | |

Sample size: 990,002

Total runtime: 4.95 seconds

# nodes: 3,253
(Threshold: $10^{-5}$)

# significant interactions: 583

Single feature: 537

Pairwise interactions: 41
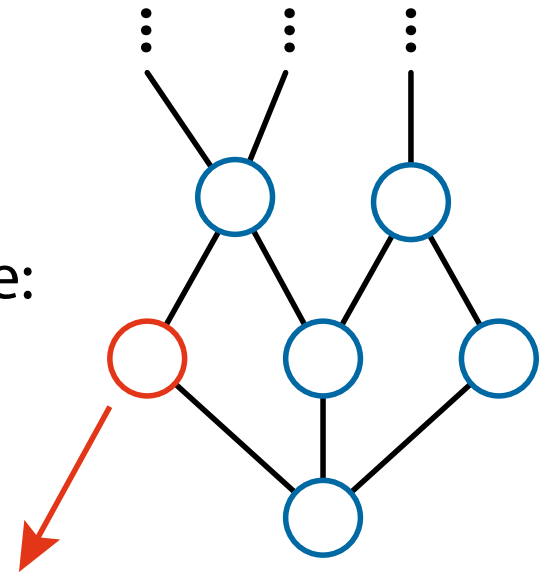
Triple interactions: 5

# Example on Real Data (accidents)

# features: 468

🔵 🔴 🟢 ···

ID 1:  1  1  0
ID 2:  1  1  1
ID 3:  1  1  0 ···
ID 4:  1  1  1
ID 5:  1  1  0
⋮       ⋮

Sample size:
340,183

Total runtime:
4.95 seconds

# nodes: 281
(Threshold: $5 \times 10^{-6}$)

# significant interactions: 280
# features in each interaction
is between 26 to 41

# Conclusion

- A close connection between the partial order structure and information geometry

  - Möbius inversion leads to the dually flat manifolds

    - M. Sugiyama, H. Nakahara, K. Tsuda, *Information Decomposition on Structured Space*, IEEE ISIT (2016)
    - S. Amari, *Information geometry on hierarchy of probability distributions*, IEEE Trans. Info. Theory (2001)
    - H. Nakahara, S. Amari, *Information-geometric measure for neural spikes*, Neural Computation (2002)

- We can decompose the KL divergence and asses the significance on any posets