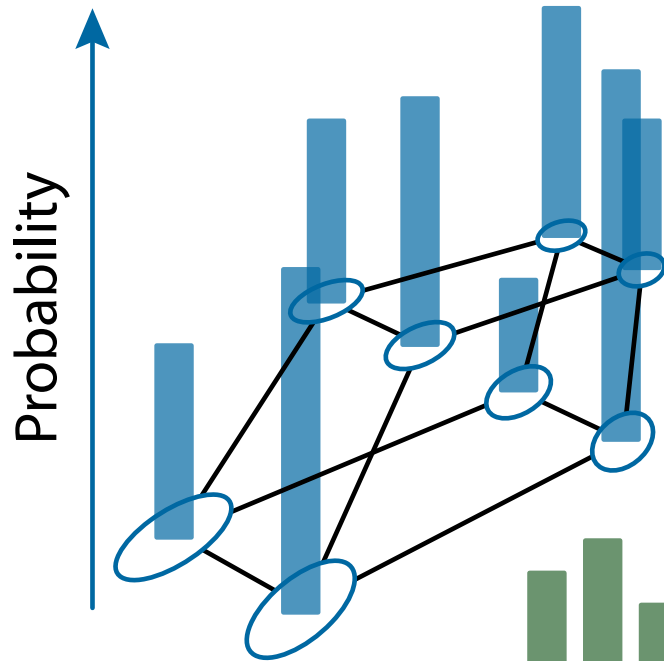# Information Decomposition on Structured Space

Mahito Sugiyama (Osaka Univ.)

Hiroyuki Nakahara (RIKEN), Koji Tsuda (UTokyo)

# Contributions

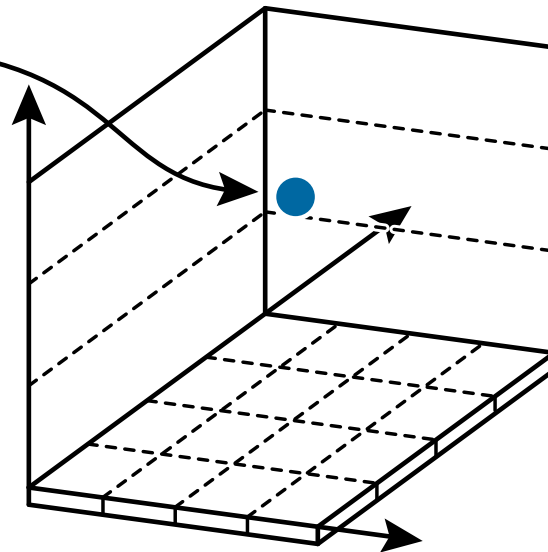- We build information geometry for posets (partially ordered sets)

  – Decomposition of KL divergence

- Key observations:

  – $\theta$-coordinate → principal ideals (lower sets) → $p$-coordinate

    ○ $\theta$-coordinate: coefficients of a log-linear model
    ○ $p$-coordinate: probabilities

  – $p$-coordinate → principal filters (upper sets) → $\eta$-coordinate

    ○ $\eta$-coordinate: frequencies (sufficient statistics)
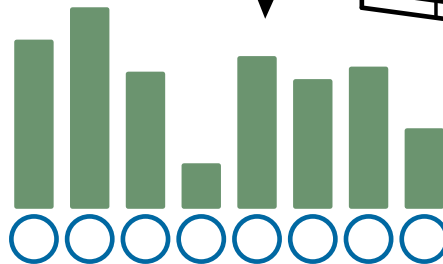
- Code: `https://git.io/decomp`

# Summary

Probability distribution
on posets (partially ordered sets)

Information
geometry

Probability
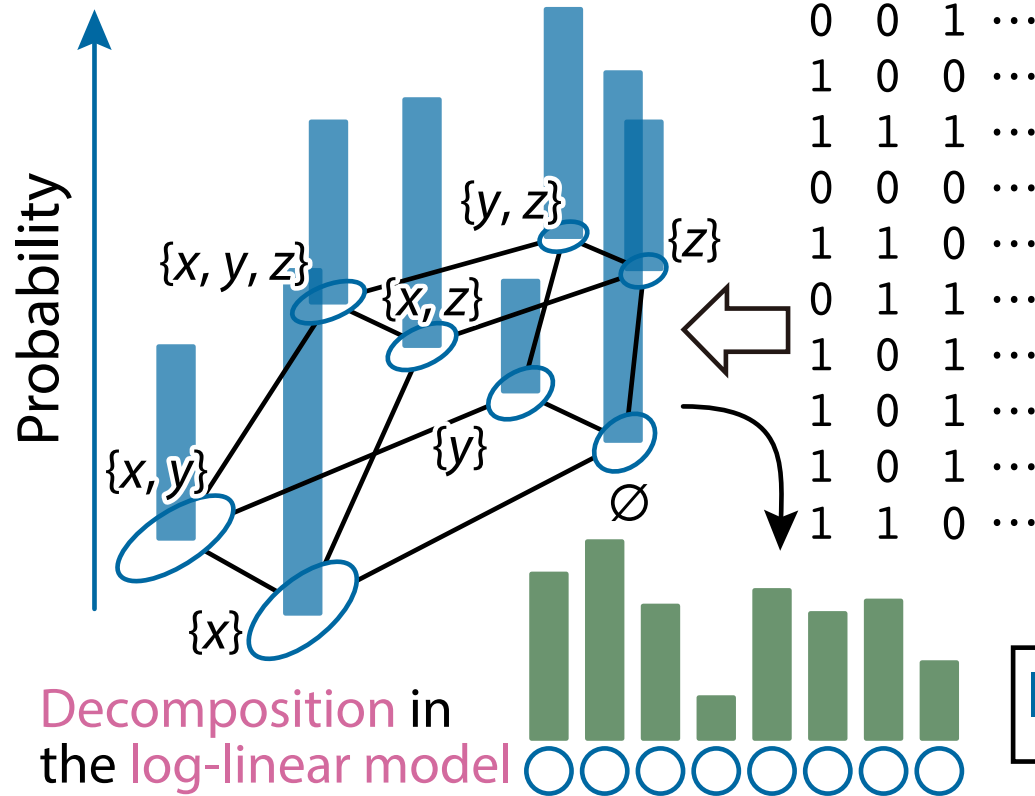
Decomposition in
the log-linear model

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

# Summary



Probability distribution on posets (partially ordered sets)

Probability

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

Decomposition in the log-linear model

x  y  z (e.g. Neurons, SNPs, ...)

| x | y | z |   |
|---|---|---|---|
| 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
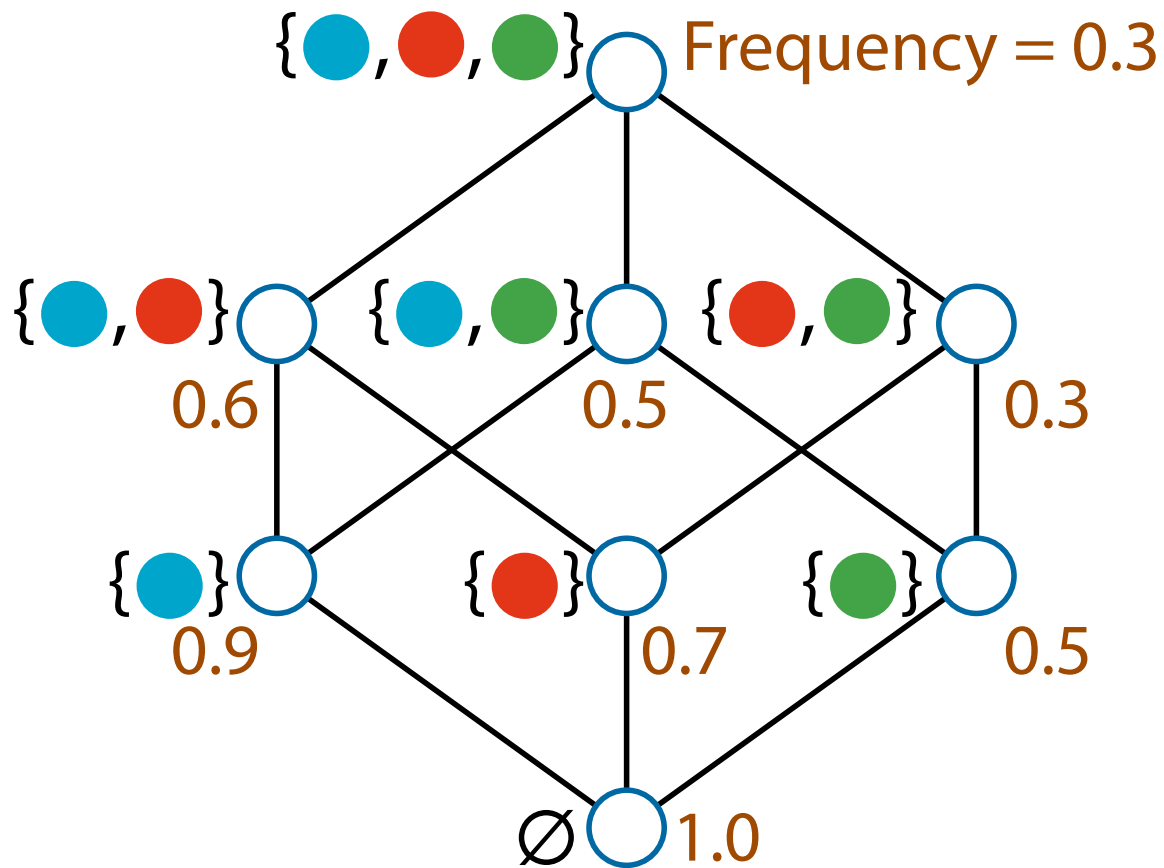
$$\log p(x) = \sum \theta(s)$$

# 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 |
| ID 10: | 0 | 1 | 0 |

Itemset lattice



3/15

Transaction database

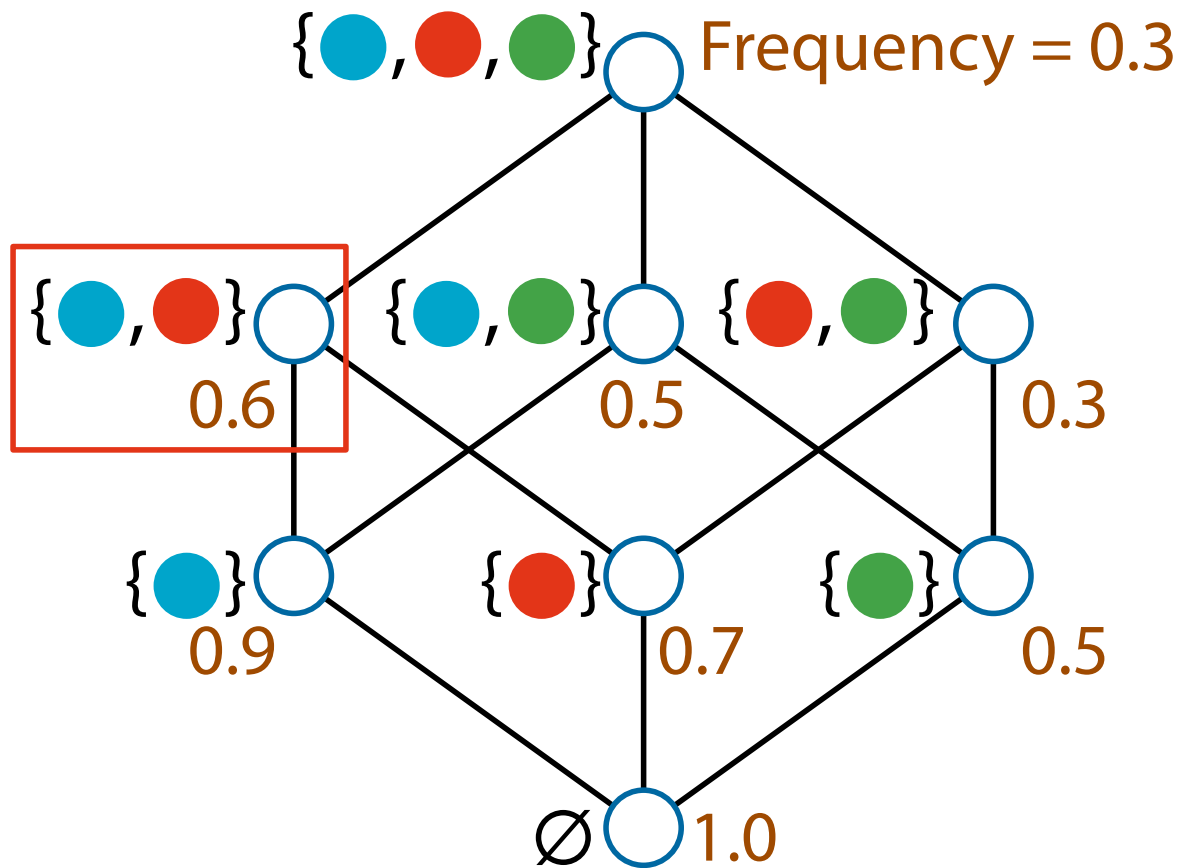| | 🔵 | 🔴 | 🟢 |
|---|---|---|---|
| ID 1: | 1 | 1 | 0 |
| ID 2: | 1 | 1 | 1 |
| ID 3: | 1 | 1 | 0 |
| ID 4: | 1 | 1 | 1 |
| ID 5: | 1 | 1 | 0 |
| ID 6: | 1 | 0 | 1 |
| ID 7: | 1 | 0 | 1 |
| ID 8: | 1 | 1 | 1 |
| ID 9: | 1 | 0 | 0 |
| ID10: | 0 | 1 | 0 |

Itemset lattice

{🔵,🔴,🟢}  Frequency = 0.3

{🔵,🔴} 0.6   {🔵,🟢} 0.5   {🔴,🟢} 0.3

{🔵} 0.9   {🔴} 0.7   {🟢} 0.5

∅ 1.0

Transaction database

Itemset lattice

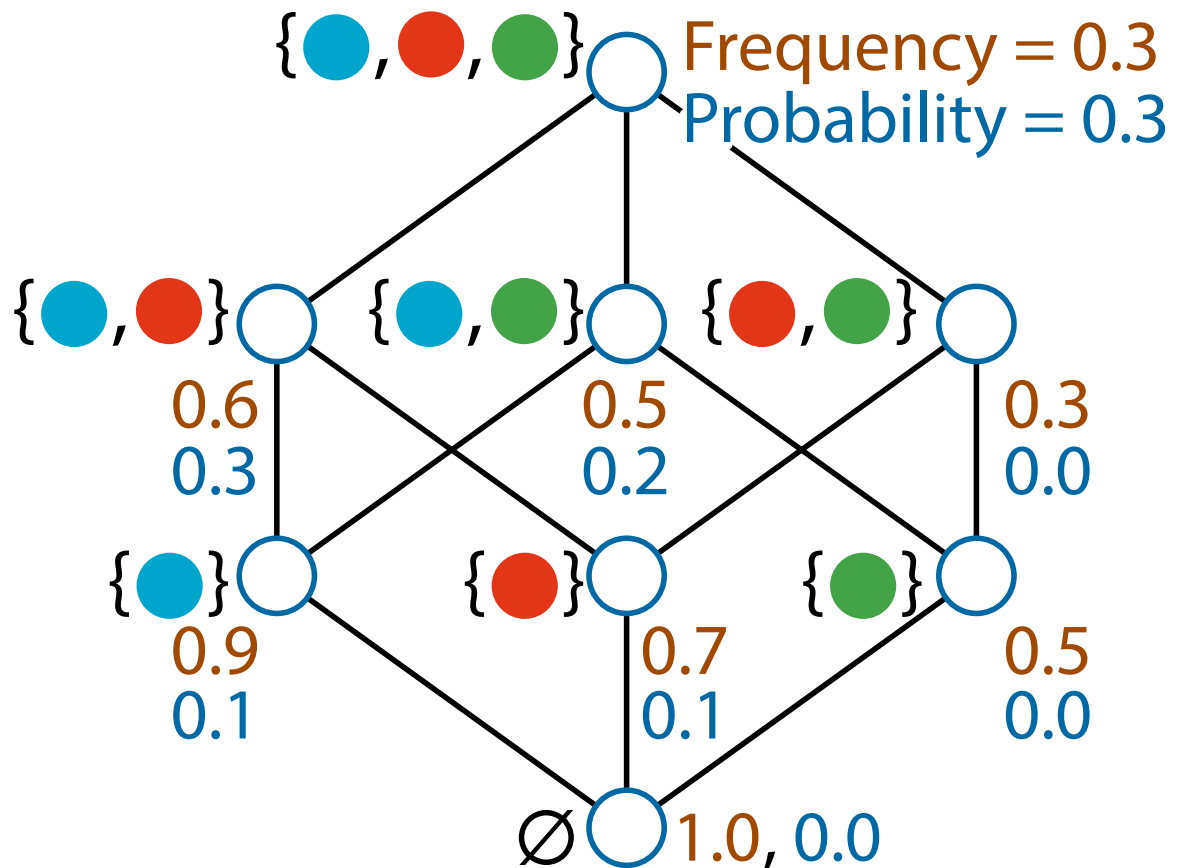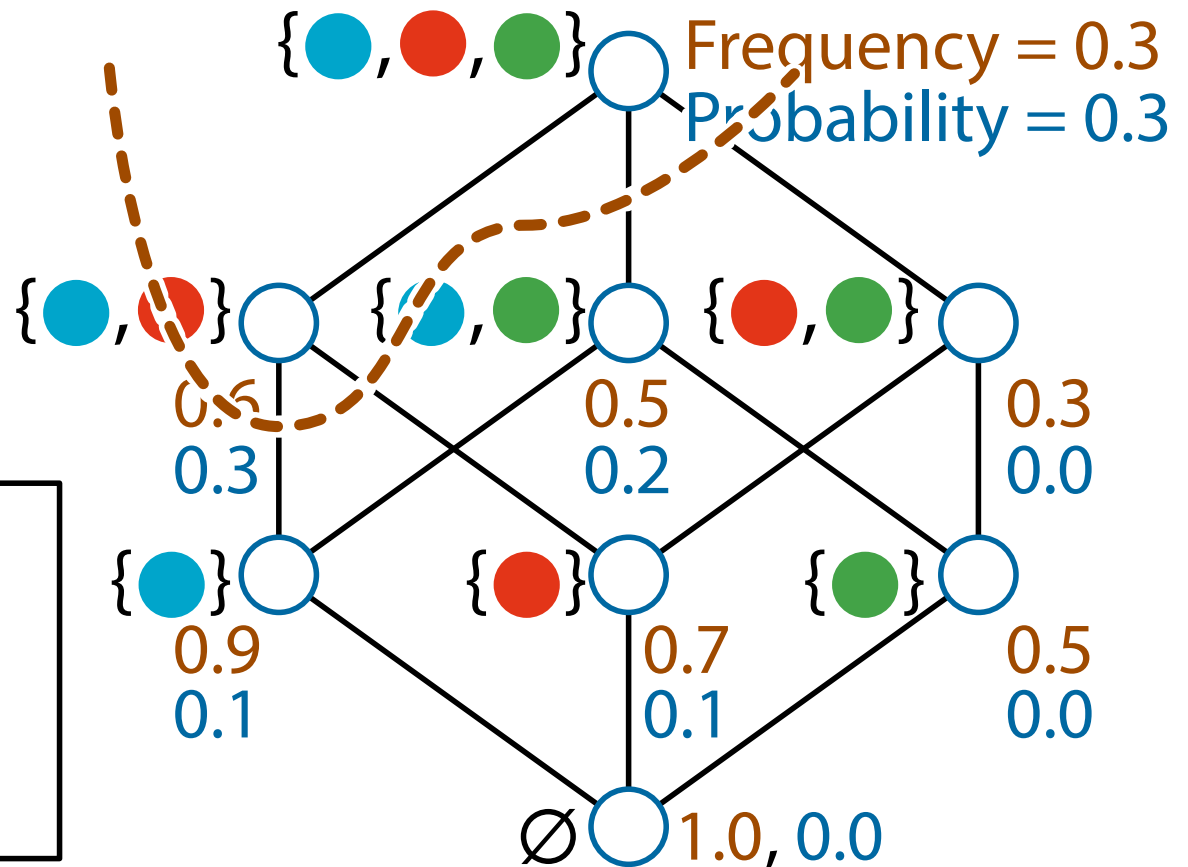| ID 1: | 1 | 1 | 0 |
| ID 2: | 1 | 1 | 1 |
| ID 3: | 1 | 1 | 0 |
| ID 4: | 1 | 1 | 1 |
| ID 5: | 1 | 1 | 0 |
| ID 6: | 1 | 0 | 1 |
| ID 7: | 1 | 0 | 1 |
| ID 8: | 1 | 1 | 1 |
| ID 9: | 1 | 0 | 0 |
| ID10: | 0 | 1 | 0 |

Frequency = 0.3
Probability = 0.3

0.6  0.3
0.5  0.3
0.3  0.2
0.0

0.9  0.7  0.5
0.1  0.1  0.0

∅  1.0, 0.0

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

Itemset lattice

Frequency = 0.3
Probability = 0.3

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

$\eta(\{\bullet,\bullet\}) = p(\{\bullet,\bullet\}) + p(\{\bullet,\bullet,\bullet\})$

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

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

$\{$ 🔵 , 🔴 , 🟢 $\}$

$\{$ 🔵 🔴 $\}$    $\{$ 🔵 🟢 $\}$    $\{$ 🔴 , 🟢 $\}$

$\{$ 🔵 $\}$    $\{$ 🔴 $\}$    $\{$ 🟢 $\}$

$\varnothing$

e.g. Gaussian

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

Natural parameter

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

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

$$\eta(x) = \mathbb{E}[\, F_x(s) \,]$$
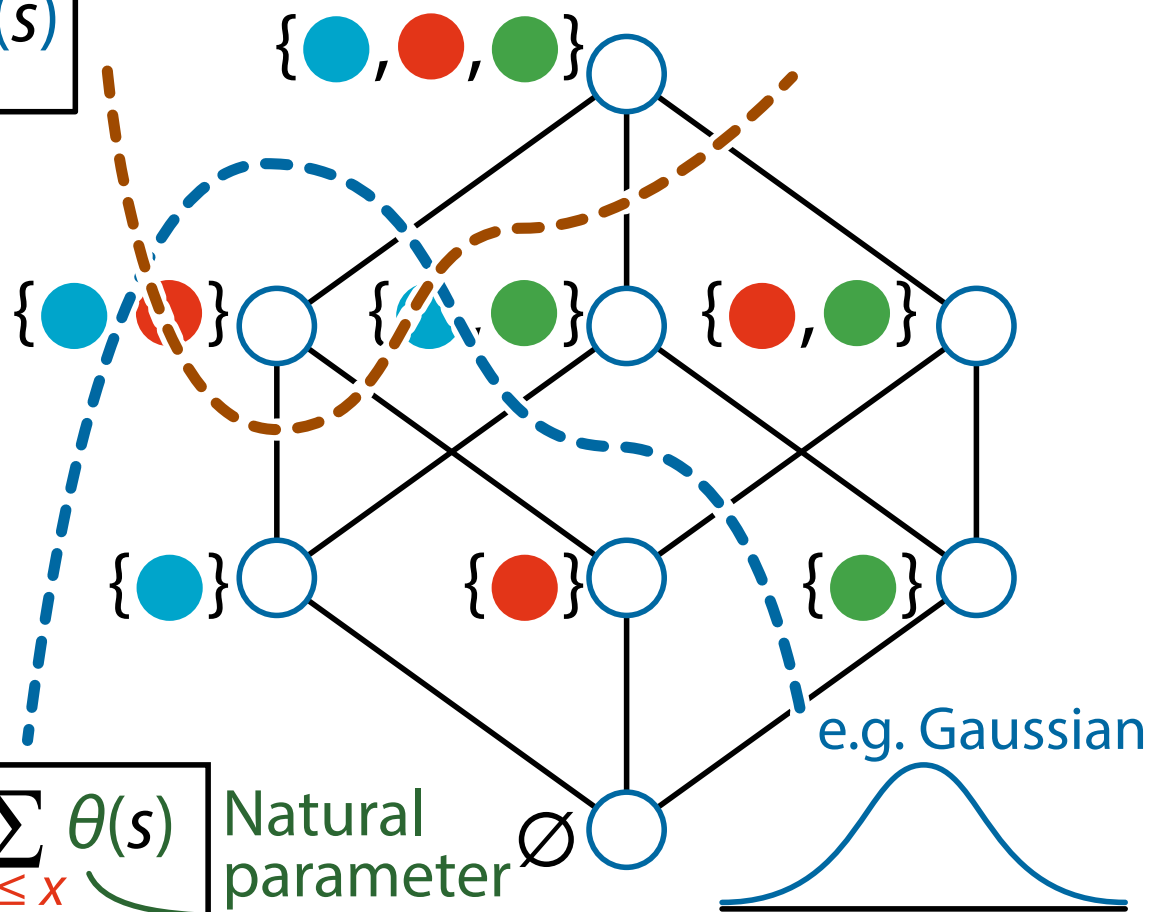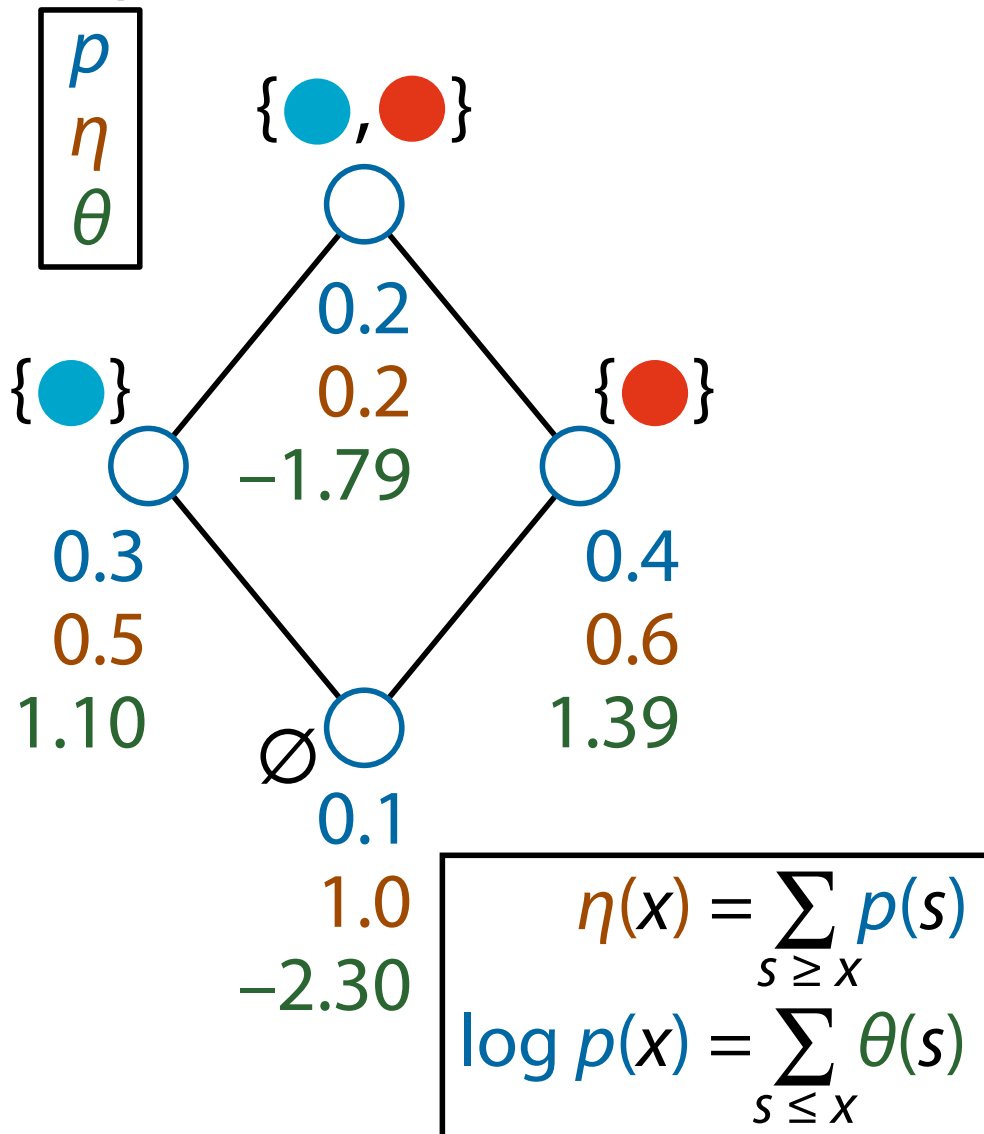
Sufficient statistics of exponential family

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

Natural parameter

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

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

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

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

$\varnothing$

e.g. Gaussian

Triple for each node

$p$
$\eta$
$\theta$

$\{ \textcolor{teal}{\bullet} , \textcolor{red}{\bullet} \}$

$\{ \textcolor{teal}{\bullet} \}$

$\{ \textcolor{red}{\bullet} \}$

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

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

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

$p(\{$🔴$\})$

$p(\{$🔵,🔴$\})$

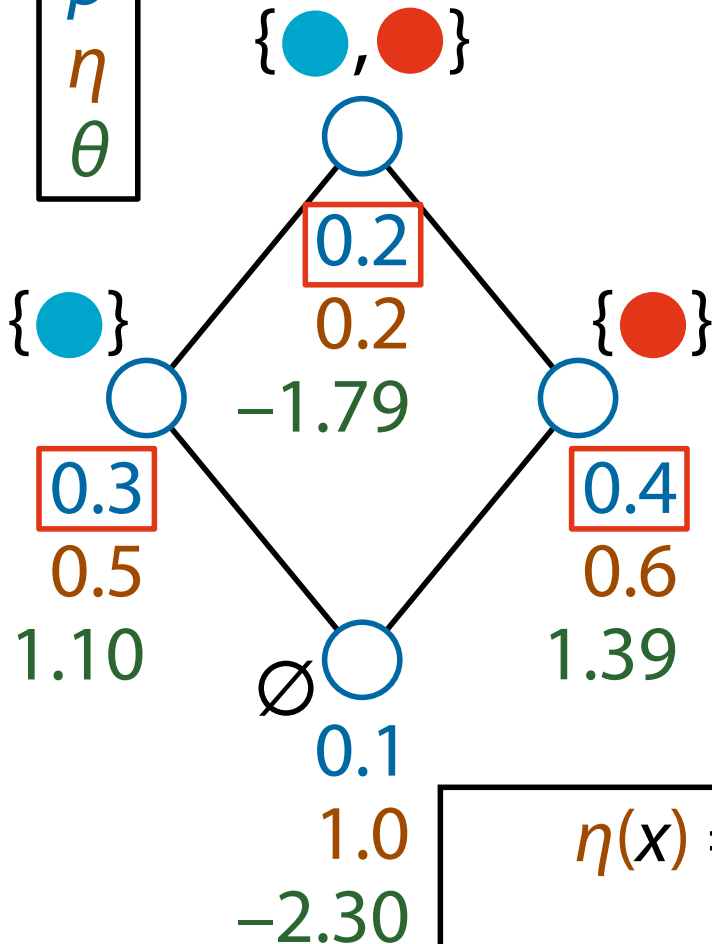$p(\{$🔵$\})$

Triple for each node

$p$
$\eta$
$\theta$

$\{$●,●$\}$

0.2
0.2
−1.79

$\{$●$\}$

$\{$●$\}$

0.3
0.5
1.10

0.4
0.6
1.39

$\varnothing$

0.1
1.0
−2.30

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

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

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
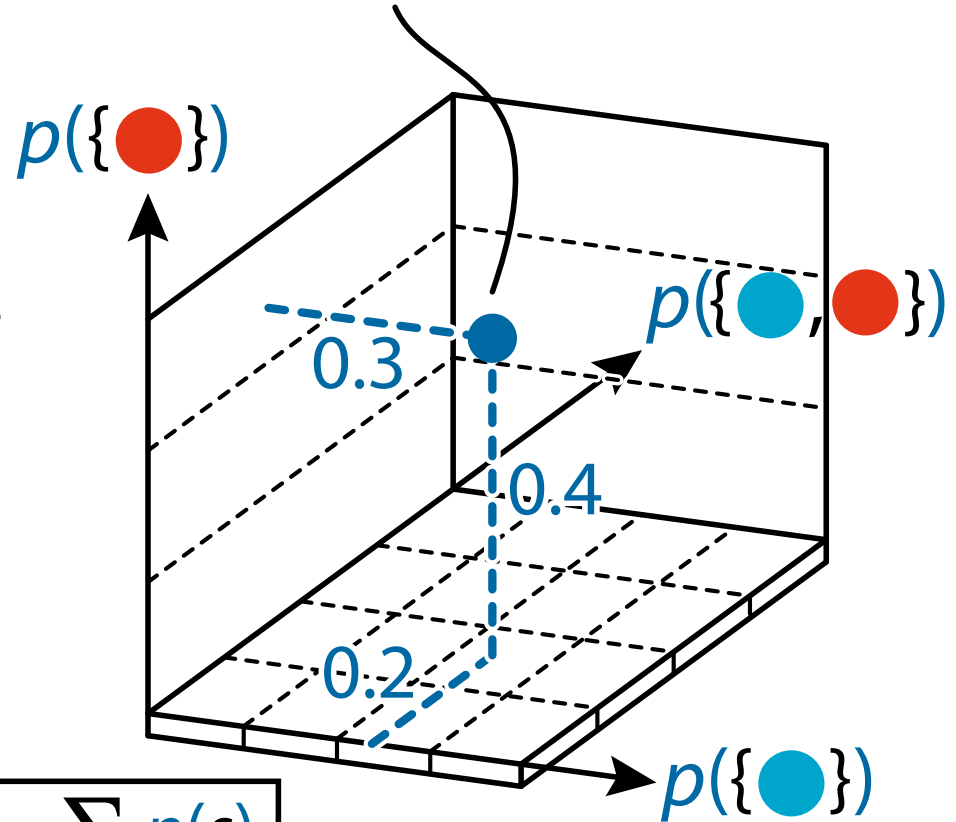       0.5          0.6
       1.10         1.39

∅

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

$\begin{bmatrix} p \\ \eta \\ \theta \end{bmatrix}$

$\{ \bullet, \bullet \}$

$\{ \bullet \}$

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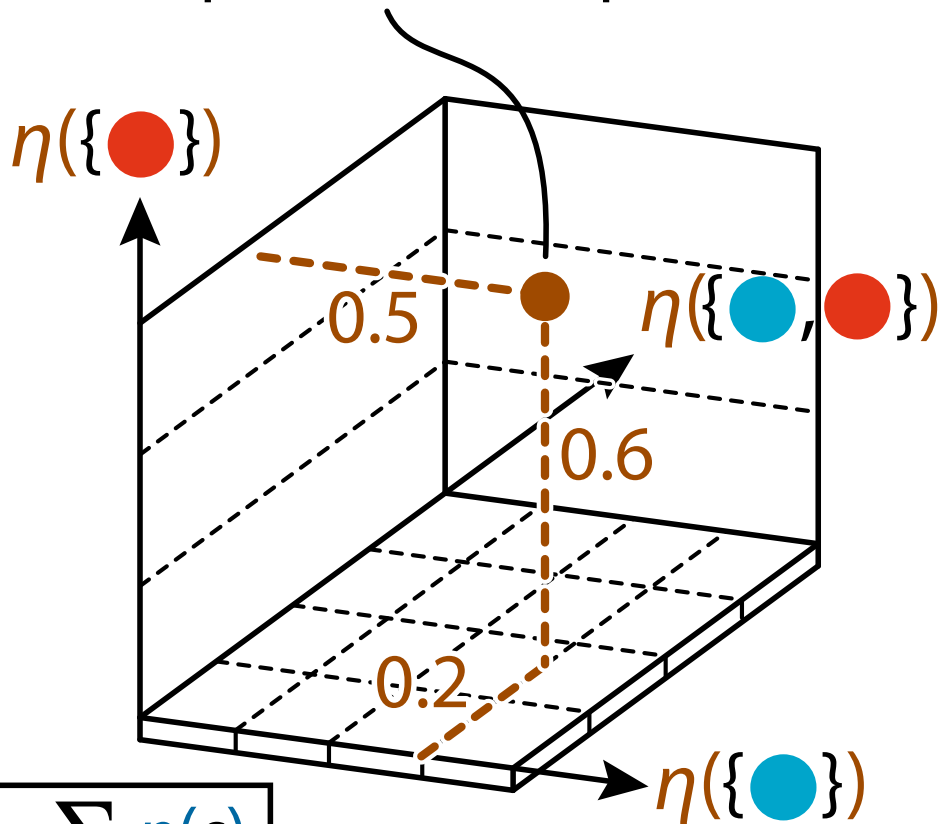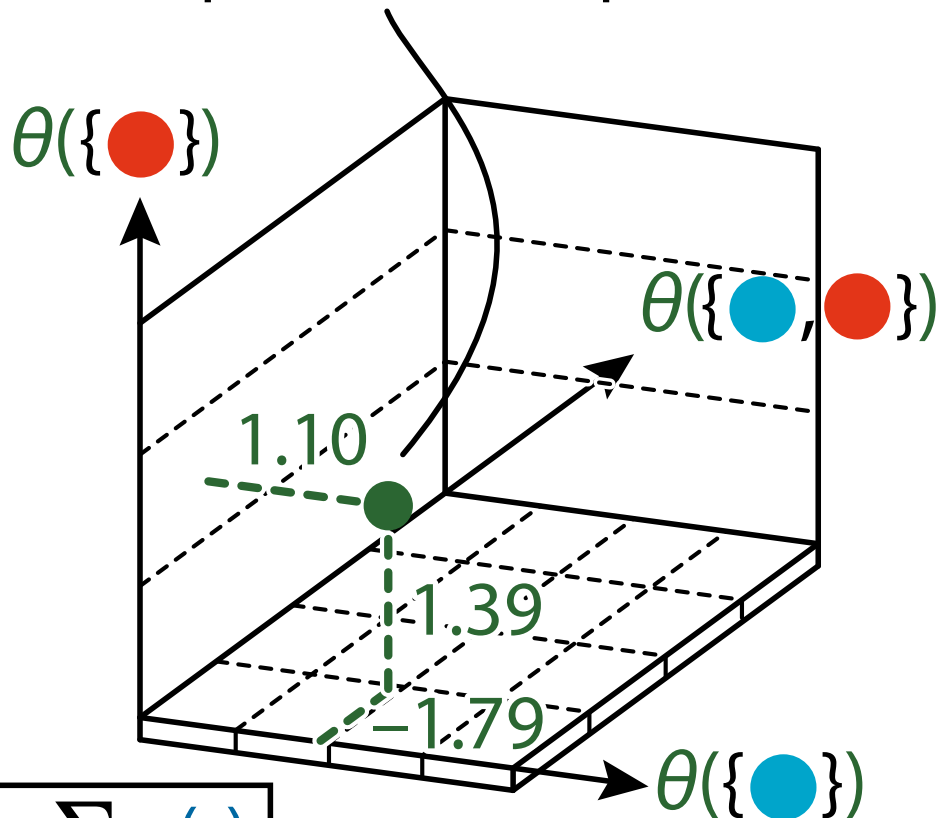
$$\eta(x) = \sum_{s \geq x} p(s)$$

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

# 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

$p$
$\eta$
$\theta$

$\{$🔵$,$🔴$\}$

0.2
0.2
−1.79

$\{$🔵$\}$

0.3
0.5
1.10

$\{$🔴$\}$

0.4
0.6
1.39

∅

0.1
1.0
−2.30

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

one-to-one

$p$
$\eta \longleftrightarrow \theta$

$\{$🔴$\}$

$\{$🔵$,$🔴$\}$

$\{$🔵$\}$

Triple for each node

$p$
$\eta$
$\theta$

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

0.2
0.2
−1.79

$\{\textcolor{cyan}{\bullet}\}$

0.3
0.5
1.10

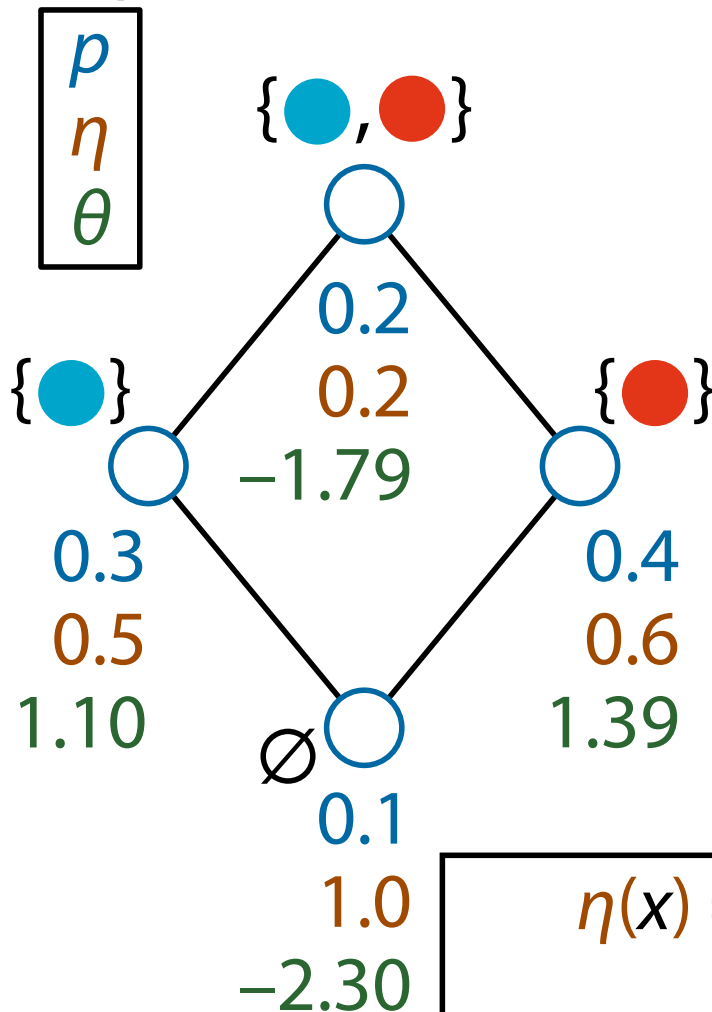$\{\textcolor{red}{\bullet}\}$

0.4
0.6
1.39

$\varnothing$

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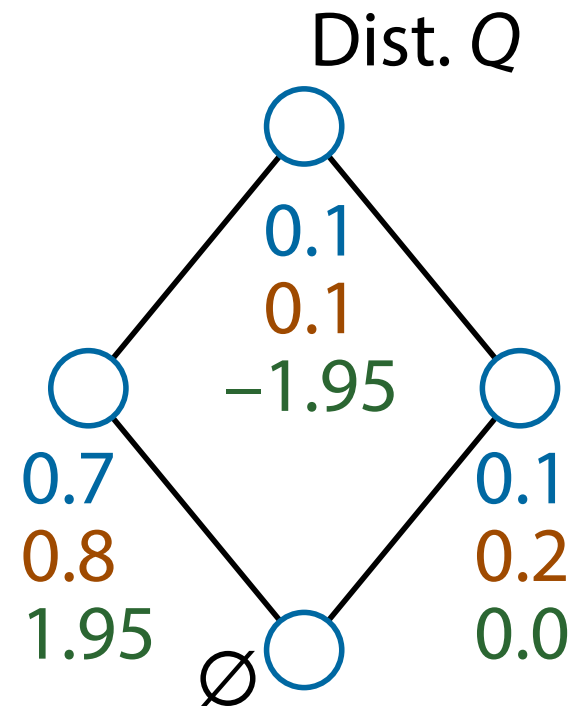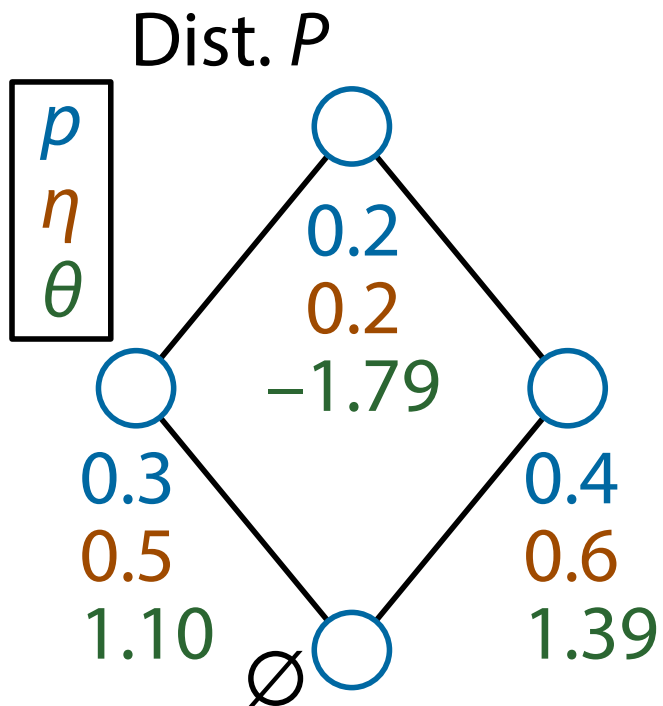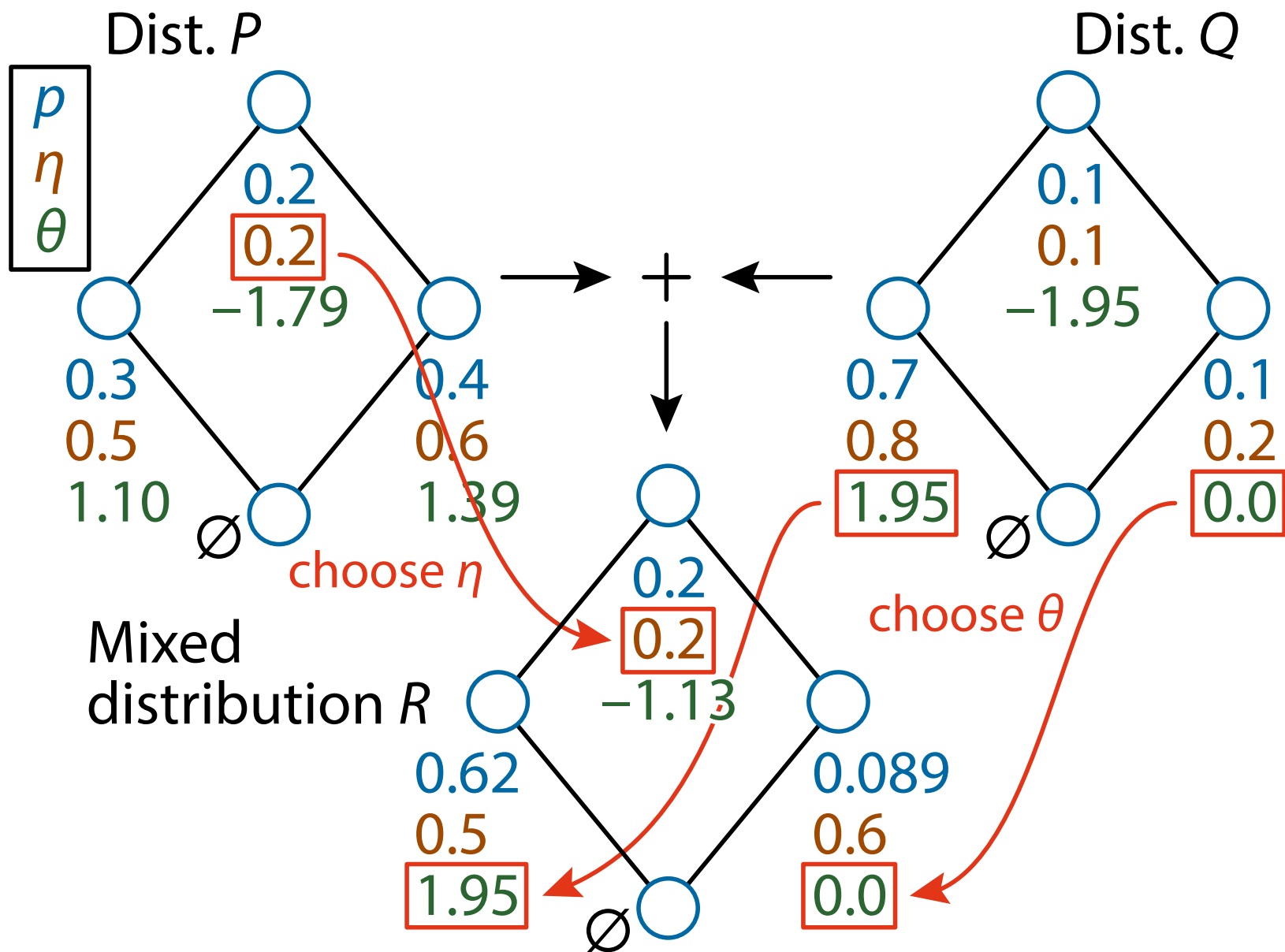
$$\eta(x) = \sum_{s \geq x} p(s)$$
$$\log p(x) = \sum_{s \leq x} \theta(s)$$

one-to-one

$p$
$\eta \longleftrightarrow \theta$

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

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

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

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

8/15

Dist. $P$

$p$
$\eta$
$\theta$

0.2
0.2
−1.79

0.3          0.4
0.5          0.6
1.10  ∅   1.39

Dist. $Q$

0.1
0.1
−1.95

0.7          0.1
0.8          0.2
1.95  ∅   0.0

Dist. $P$

$p$
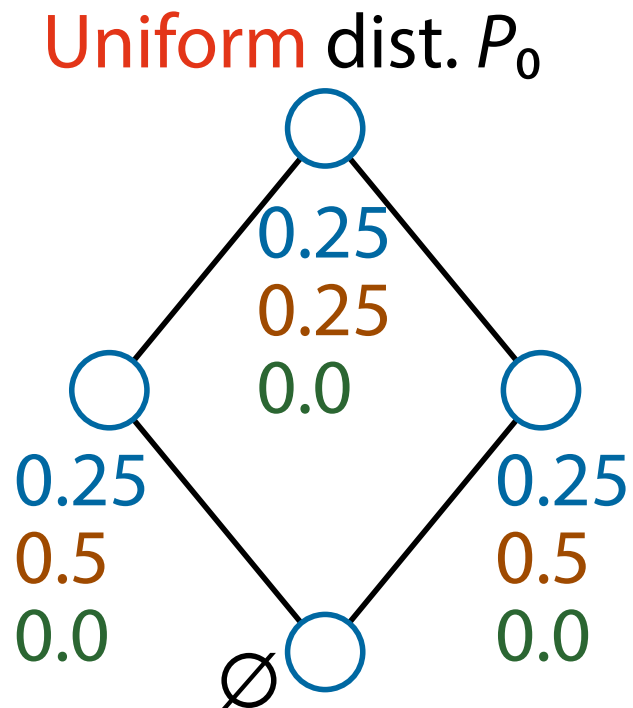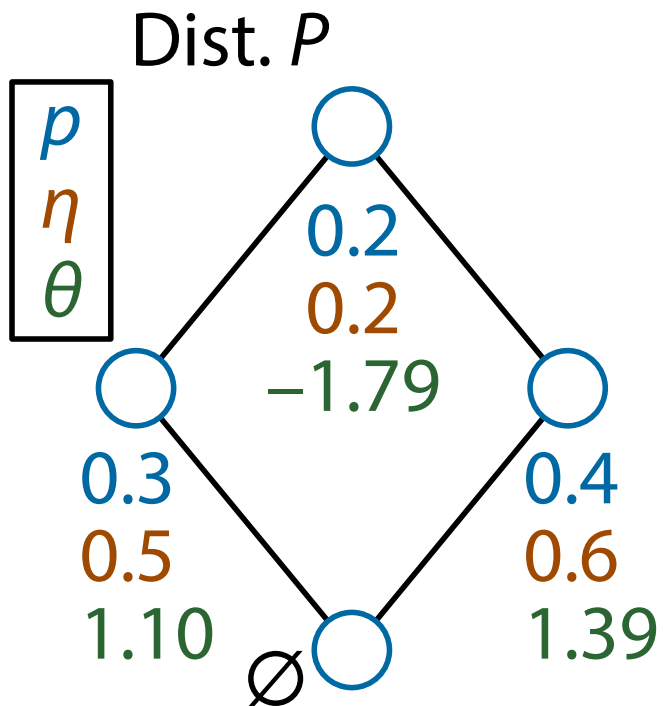$\eta$
$\theta$

0.2
0.2
−1.79

0.3
0.5
1.10

$\varnothing$

0.4
0.6
1.39

choose $\eta$

Mixed
distribution $R$

0.2
0.2
−1.13

0.62
0.5
1.95

$\varnothing$

0.089
0.6
0.0

Dist. $Q$

0.1
0.1
−1.95

0.7
0.8
1.95

$\varnothing$

0.1
0.2
0.0

choose $\theta$

Dist. *P*

$p$
$\eta$
$\theta$

0.2
0.2
−1.

0.3
0.5
1.10

∅

0.6
1.39

$KL(P, Q) =$
$KL(P, R) + KL(R, Q)$

Nonnegative decomposition
of the KL divergence

Dist. *Q*

0.1
0.1
1.95

0.1
0.2
0.0

∅

0.8
1.95

choose $\eta$

Mixed
distribution *R*

0.2
0.2
−1.13

0.62
0.5
1.95

∅

0.089
0.6
0.0

choose $\theta$

Dist. $P$

$p$
$\eta$
$\theta$

0.2
0.2
−1.79

0.3
0.5
1.10

0.4
0.6
1.39

$\emptyset$

Uniform dist. $P_0$

0.25
0.25
0.0

0.25
0.5
0.0

0.25
0.5
0.0

$\emptyset$

Dist. $P$

Uniform dist. $P_0$

$p$
$\eta$
$\theta$

0.2
0.2
−1.79

0.3
0.5
1.10

0.4
0.6
1.39

$\varnothing$

choose $\eta$

0.25
0.25
0.0

0.25
0.5
0.0

0.25
0.5
0.0

$\varnothing$

choose $\theta$
(KNOCK DOWN)

Mixed
distribution $R$

???
???
0.0

???
0.5
???

???
0.6
???

$\varnothing$

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

10/15

Dist. $P$

$p$
$\eta$
$\theta$

Contribution of the node
$= KL(P, R) = 0.086$

0.2
0.2
−1.79

0.3
0.5
1.10

$\varnothing$

0.4
0.6
1.39

choose $\eta$

Mixed
distribution $R$

Uniform dist. $P_0$

0.25
0.25
0.0

0.25
0.5
0.0

$\varnothing$

0.25
0.5
0.0

0.3
0.3
0.0

0.2
0.5
0.0

$\varnothing$

0.3
0.6
0.405

choose $\theta$
(KNOCK DOWN)

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
= $KL(P, R)$ = 0.086

0.2
0.2
−1.79

0.25

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

0.3
0.5
1.10

0.4
0.6
1.39

∅

choose *r*

Mixed
distribution *R*

0.0

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

0.2
0.5
0.0

0.3
0.6
0.405

∅

# Make a Poset from Data

Dataset

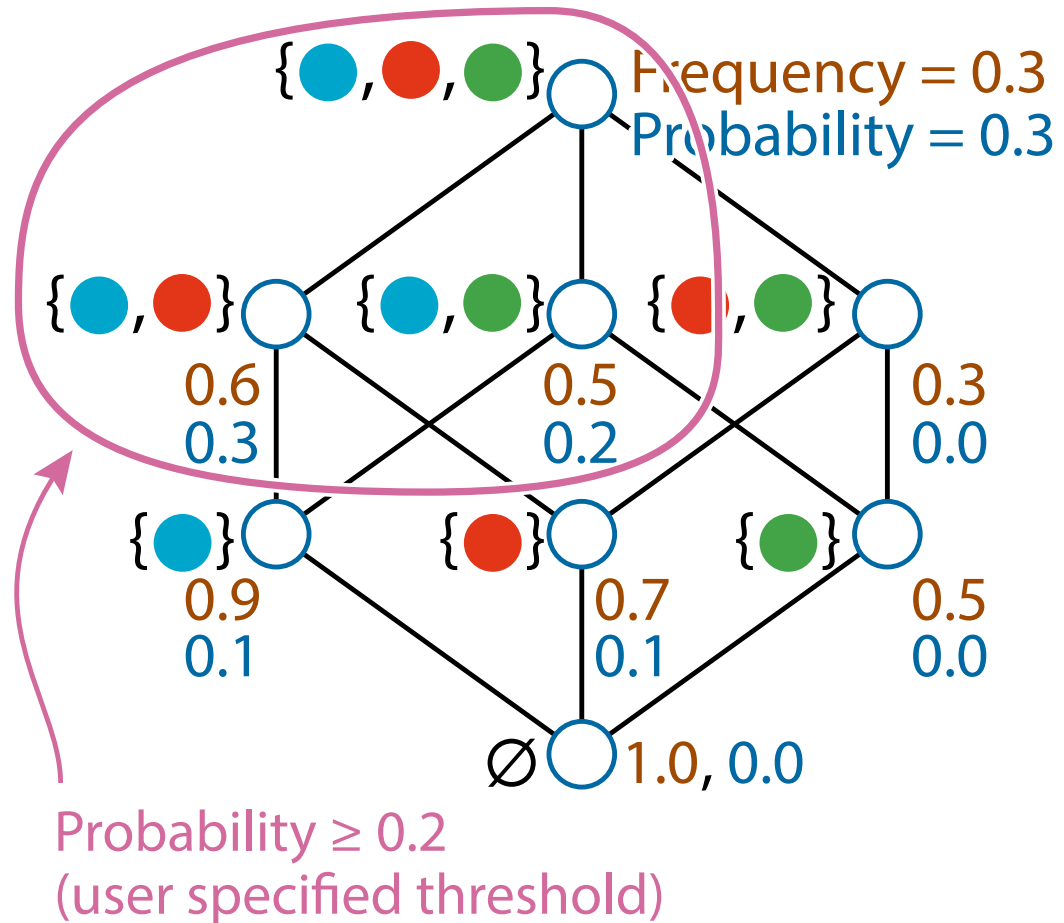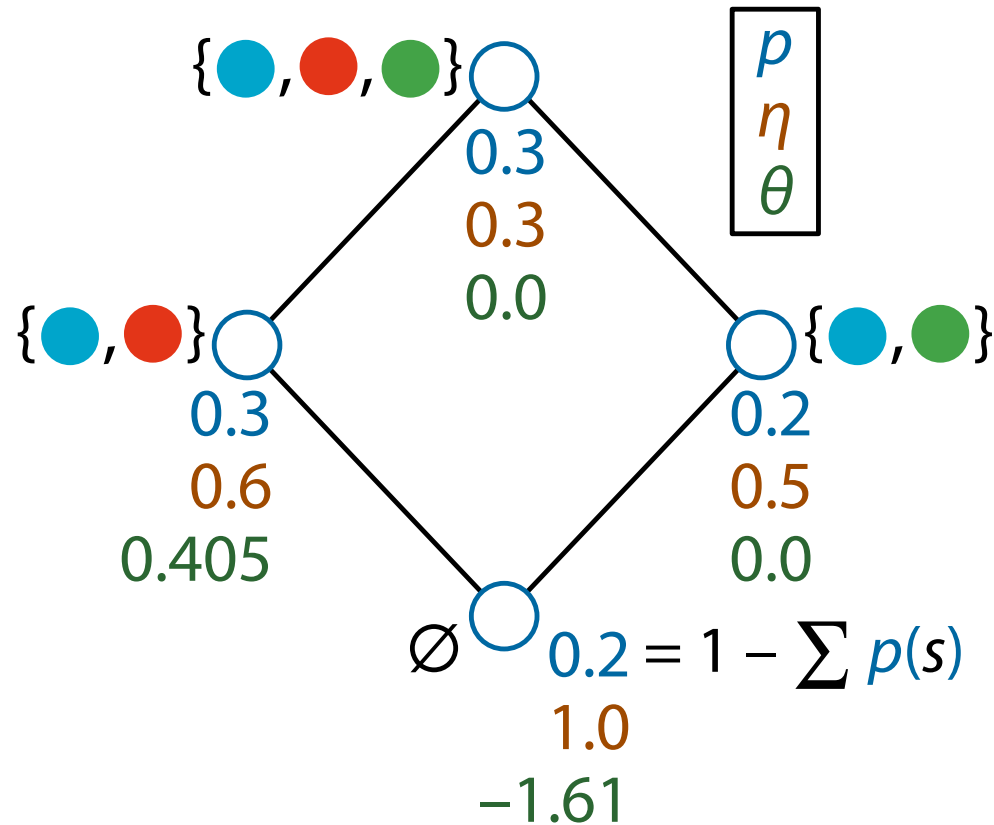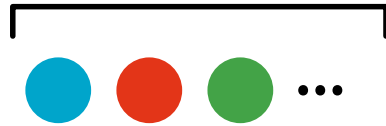| | 🔵 | 🔴 | 🟢 |
|---|---|---|---|
| ID 1: | 1 | 1 | 0 |
| ID 2: | 1 | 1 | 1 |
| ID 3: | 1 | 1 | 0 |
| ID 4: | 1 | 1 | 1 |
| ID 5: | 1 | 1 | 0 |
| ID 6: | 1 | 0 | 1 |
| ID 7: | 1 | 0 | 1 |
| ID 8: | 1 | 1 | 1 |
| ID 9: | 1 | 0 | 0 |
| ID10: | 0 | 1 | 0 |

{🔵,🔴,🟢}  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

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        0.5         0.3
0.3        0.2         0.0

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

∅ ⬡ 1.0, 0.0

Probability ≥ 0.2
(user specified threshold)

11/15

# Remove Nodes with Probability 0

# 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

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

Total runtime: 4.95 seconds

# significant interactions: 583

Single feature: 537
Pairwise interactions: 41
Triple interactions: 5

# Example on Real Data (accidents)

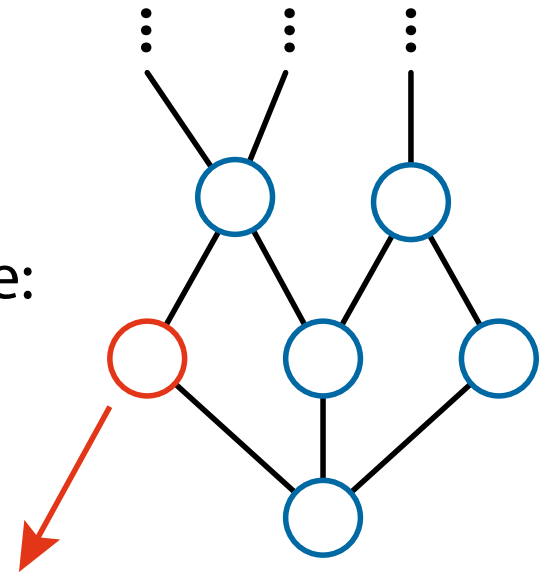# features: 468

🔵 🔴 🟢 ...

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

- We build information geometry for posets (partially ordered sets)
  - Natural connection between the information geometric dual coordinates and the partial order structure
  - Code: `https://git.io/decomp`

- We can decompose a probability distribution and asses the significance of any-order interactions

- Our results generalizes the following:
  - S. Amari, *Information geometry on hierarchy of probability distributions*, IEEE Trans. on Information Theory (2001)
  - H. Nakahara, S. Amari, *Information-geometric measure for neural spikes*, Neural Computation (2002)