

July 11, 2016
ISIT 2016



Information Decomposition on Structured Space

Mahito Sugiyama (Osaka Univ.)

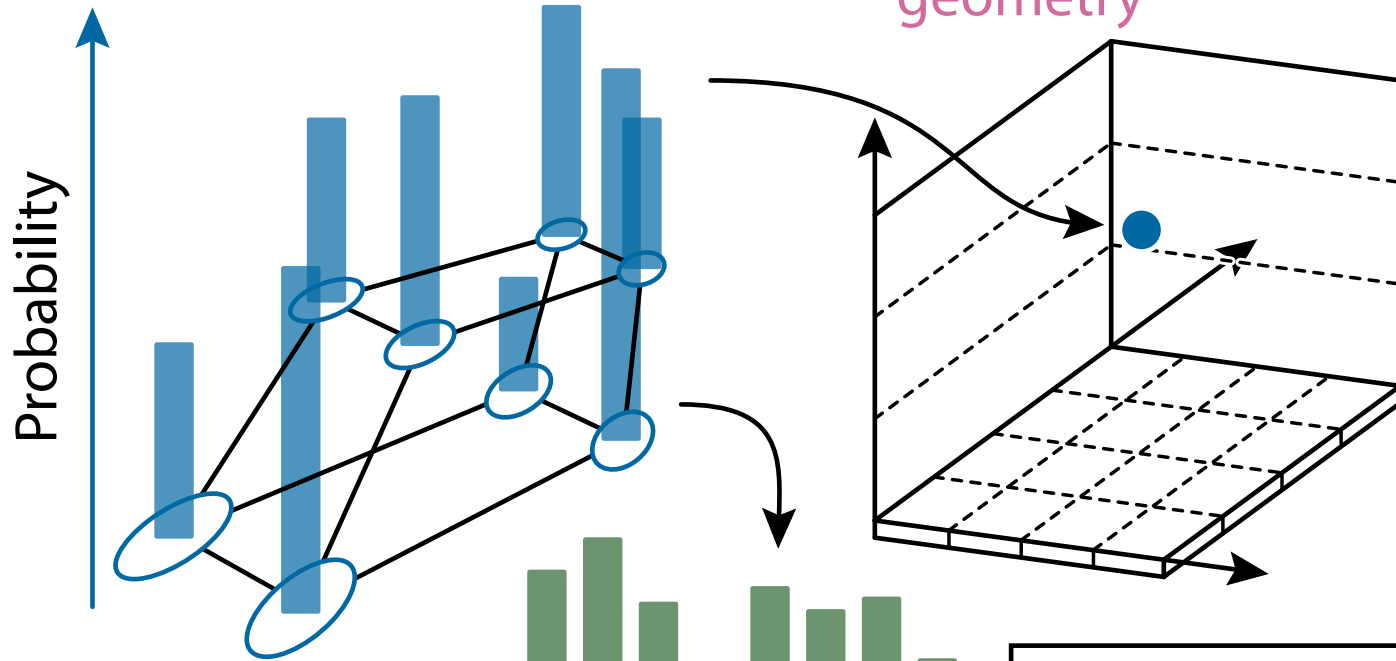
Hiroyuki Nakahara (RIKEN), Koji Tsuda (UTokyo)

Contributions

- We build **information geometry** for **posets** (partially ordered sets)
 - Decomposition of **KL divergence**
- Key observations:
 - θ -coordinate \rightarrow principal **ideals** (lower sets) \rightarrow p -coordinate
 - θ -coordinate: coefficients of a log-linear model
 - p -coordinate: probabilities
 - p -coordinate \rightarrow principal **filters** (upper sets) \rightarrow η -coordinate
 - η -coordinate: frequencies (sufficient statistics)
- Code: <https://git.io/decomp>

Summary

Probability distribution
on **posets** (partially ordered sets)

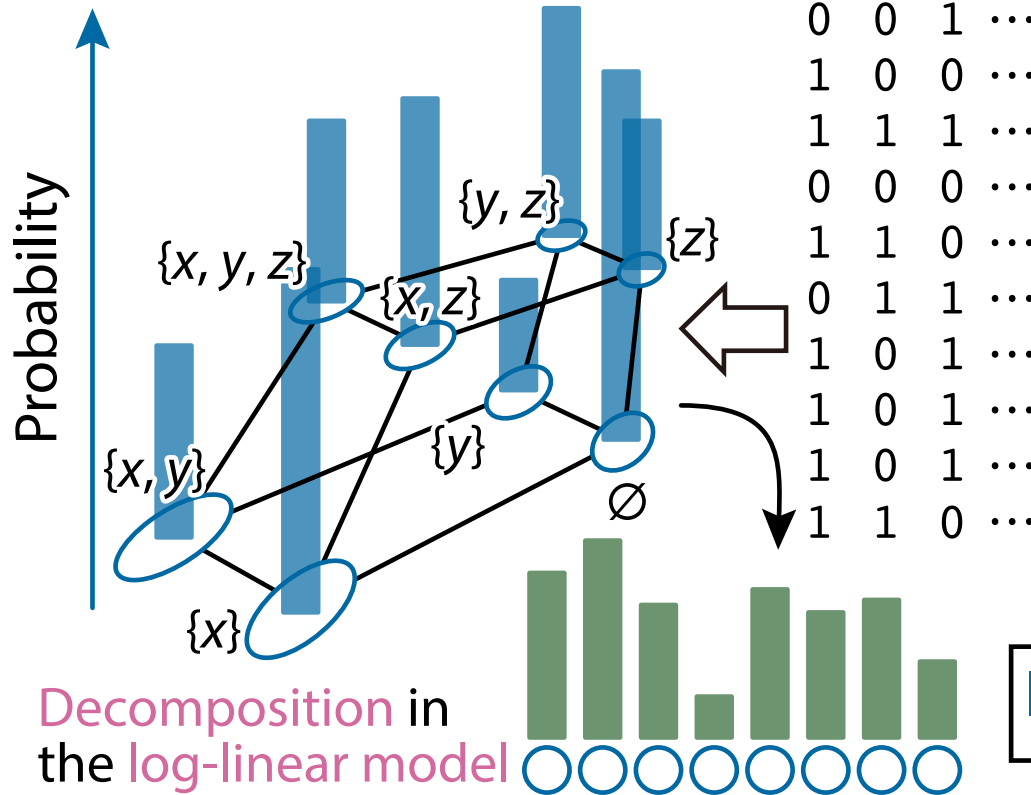


Decomposition in
the **log-linear model**

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

Summary

Probability distribution on **posets** (partially ordered sets)

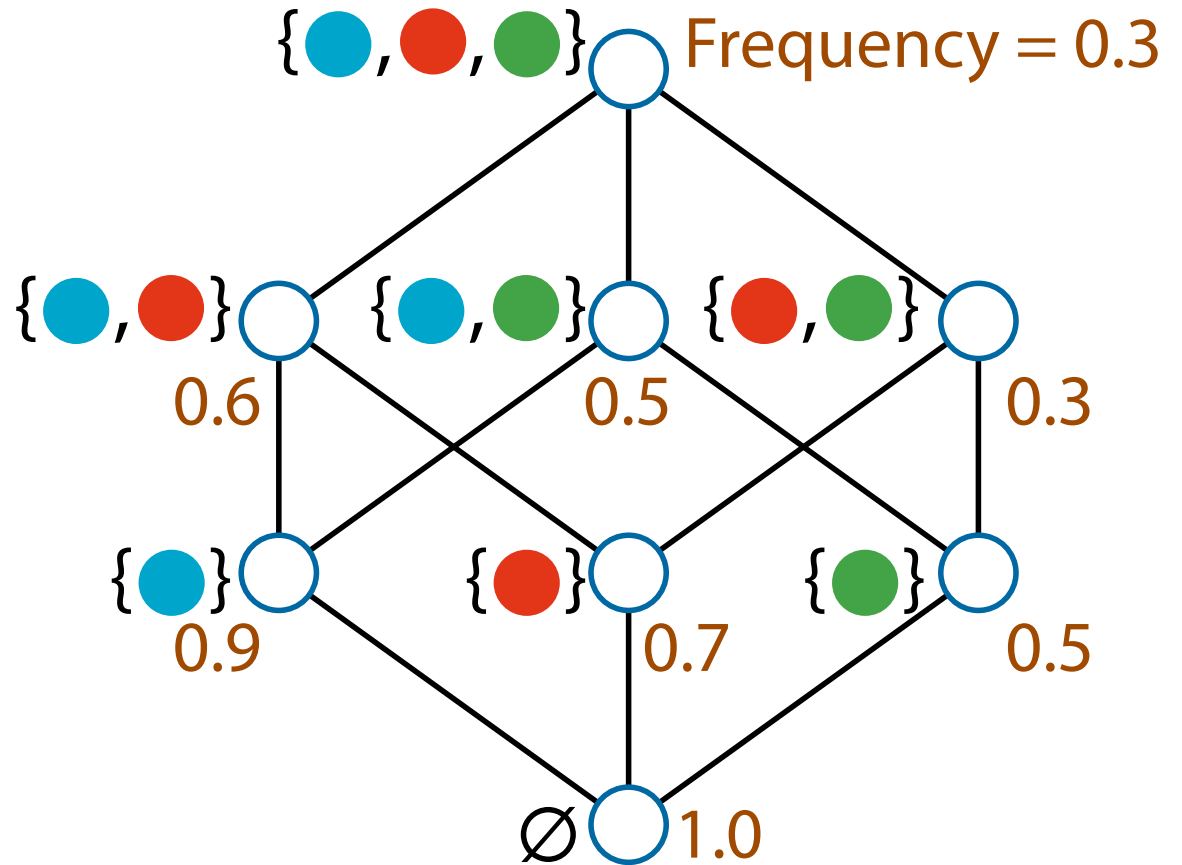


Transaction database






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

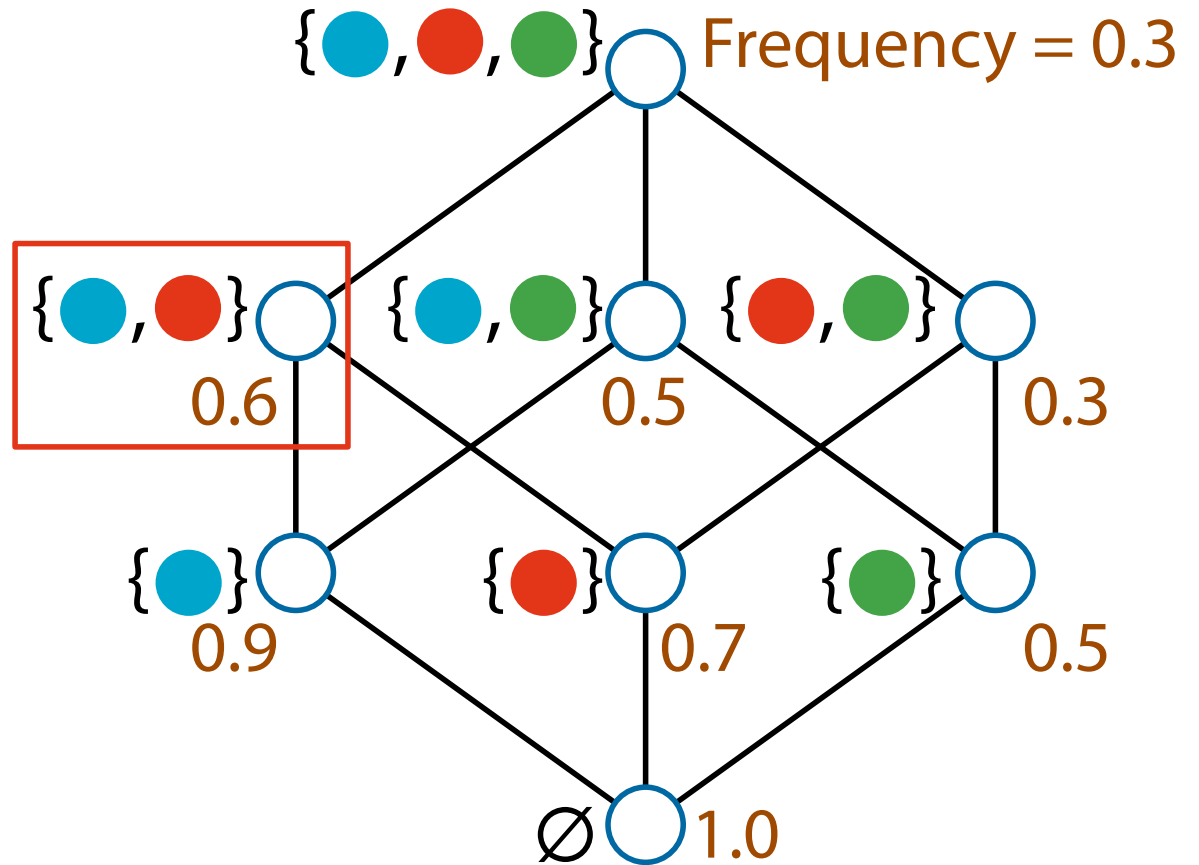
Itemset lattice



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

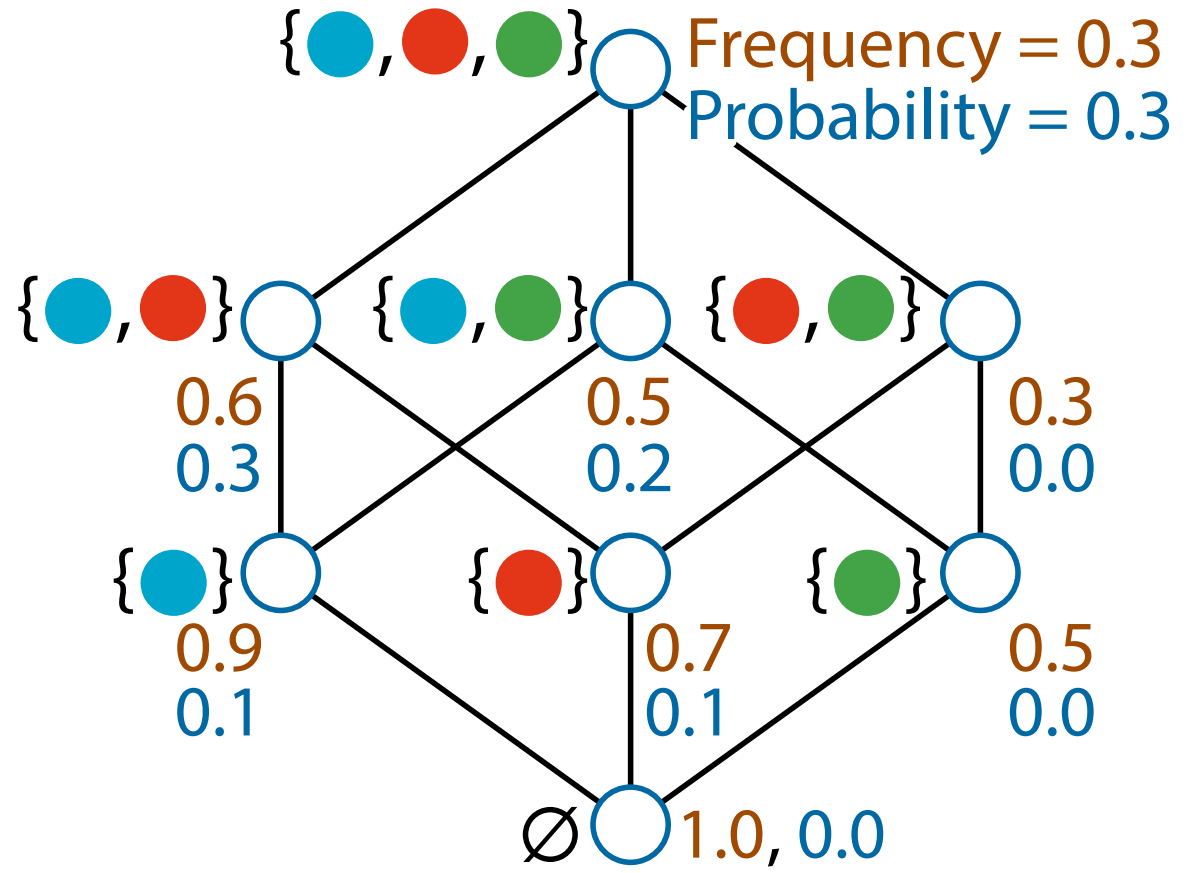
Itemset lattice



Transaction database

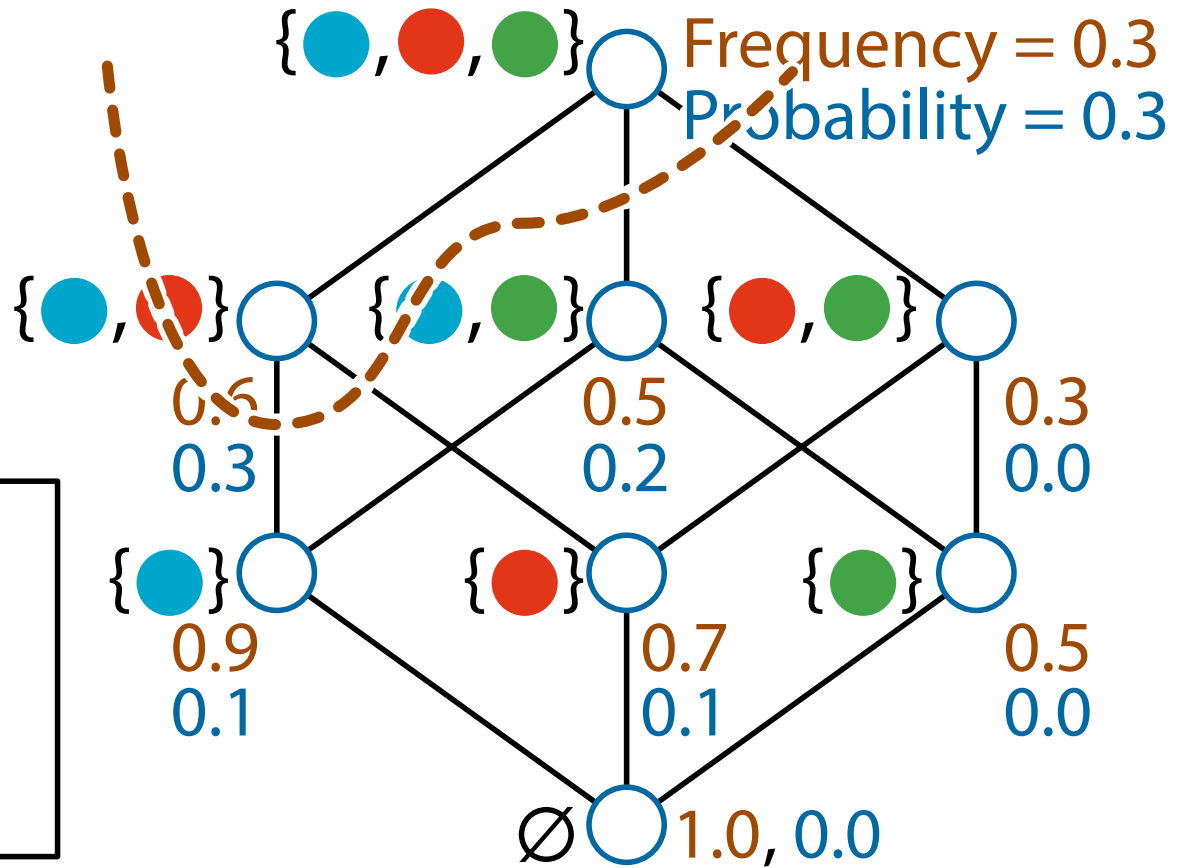
	●	●	●
ID 1:	1	1	0
ID 2:	1	1	1
ID 3:	1	1	0
ID 4:	1	1	1
ID 5:	1	1	0
ID 6:	1	0	1
ID 7:	1	0	1
ID 8:	1	1	1
ID 9:	1	0	0
ID10:	0	1	0

Itemset lattice



Upward =
Pattern mining

Itemset lattice

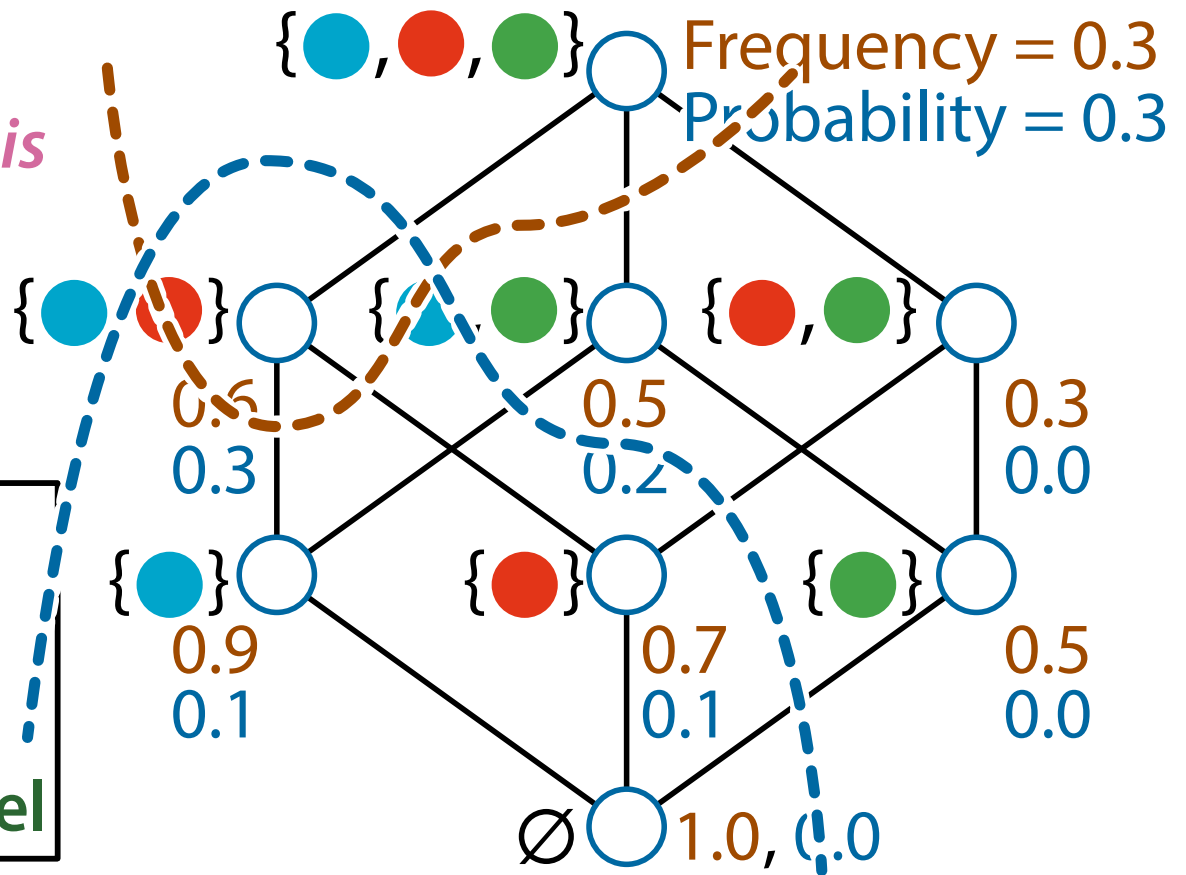


η : Frequency
 p : Probability

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

Upward =
Pattern mining
Downward =
Log-linear analysis

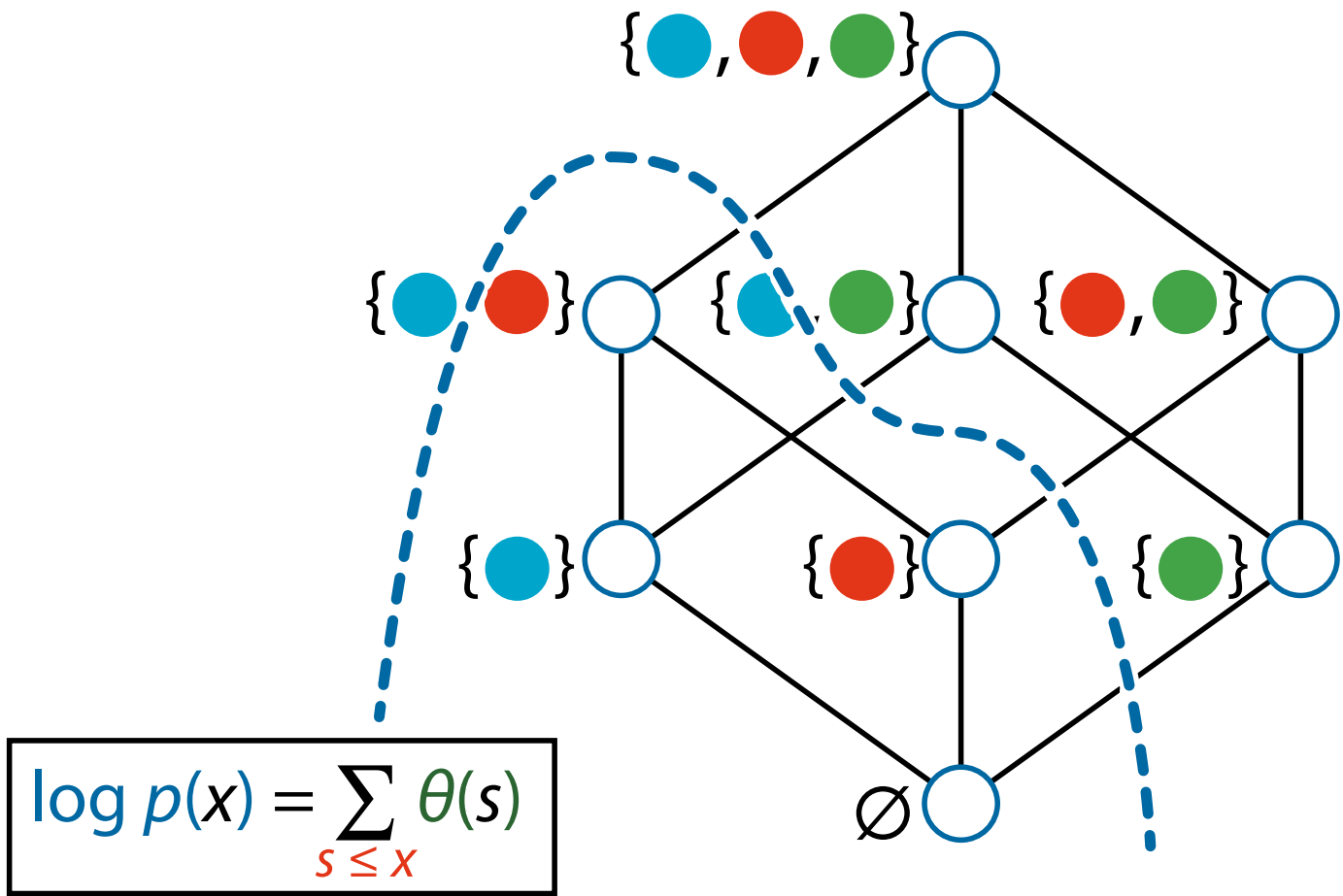
Itemset lattice

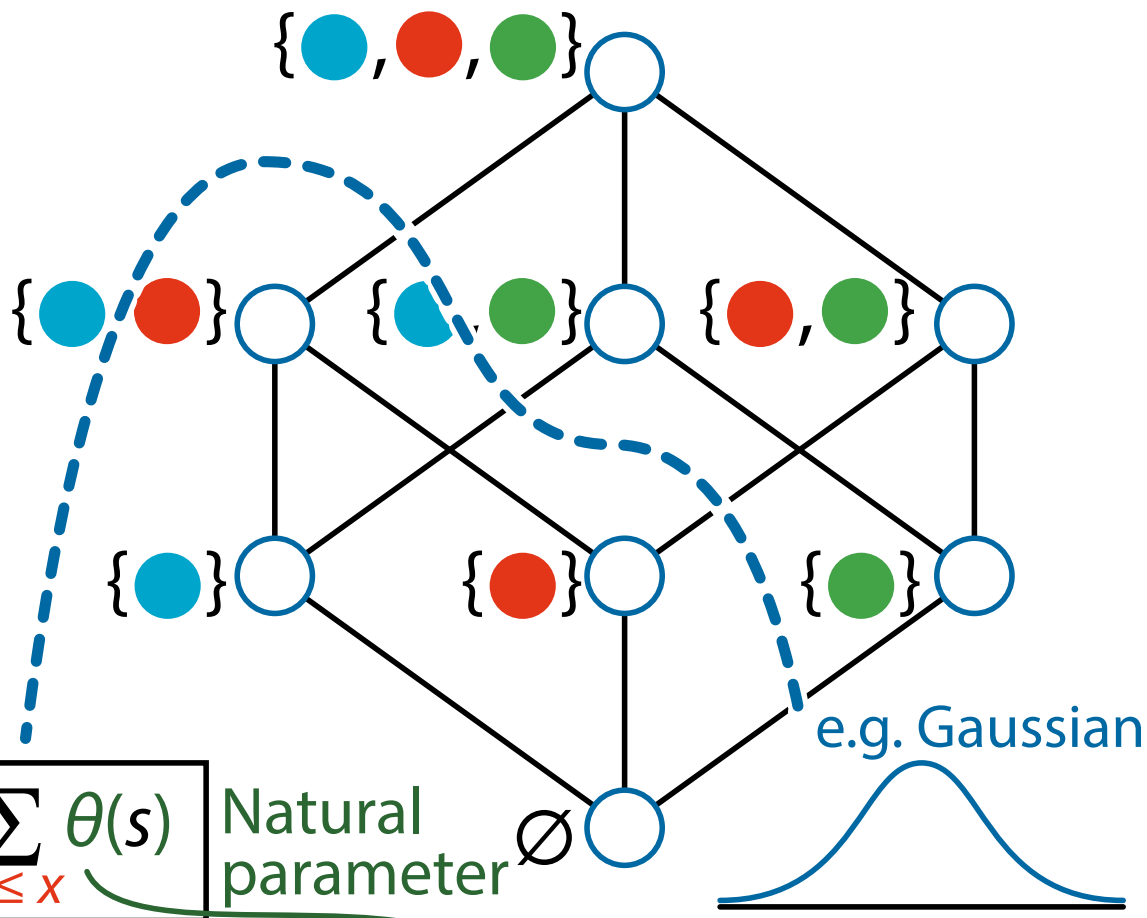


η : Frequency
 p : Probability
 θ : Coefficient of
log-linear model

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

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





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

Natural parameter θ

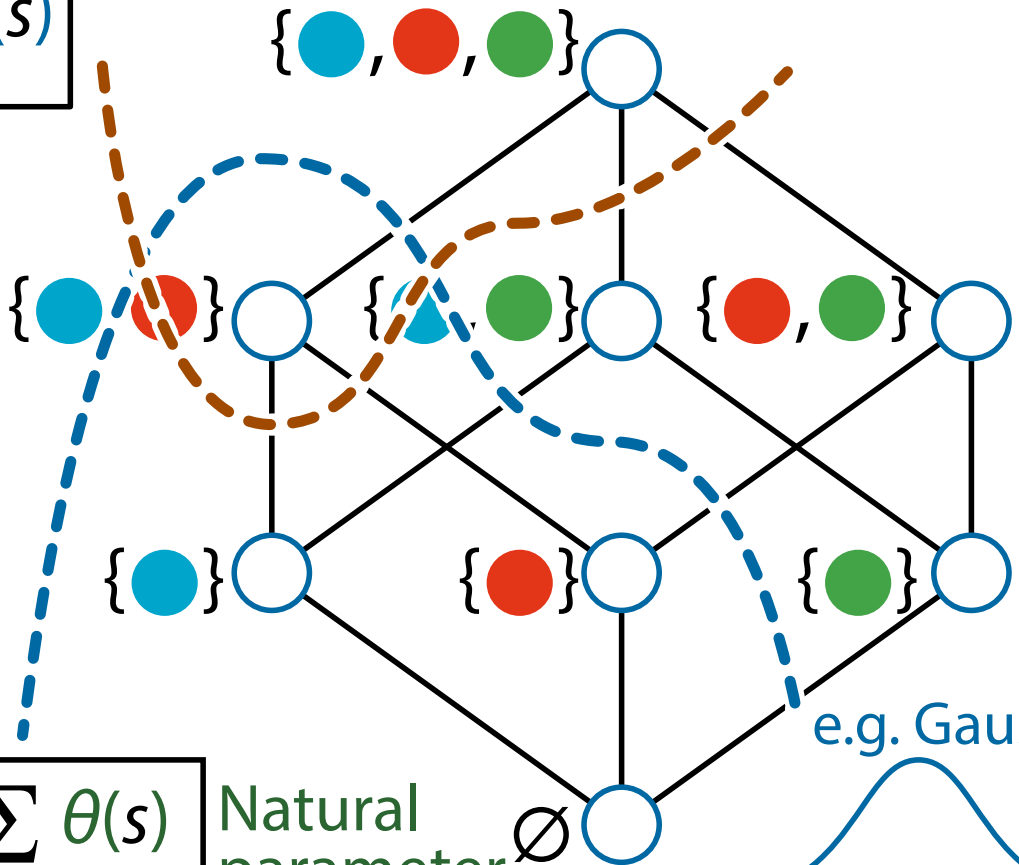
Exponential family:

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

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

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

Sufficient statistics of exponential family

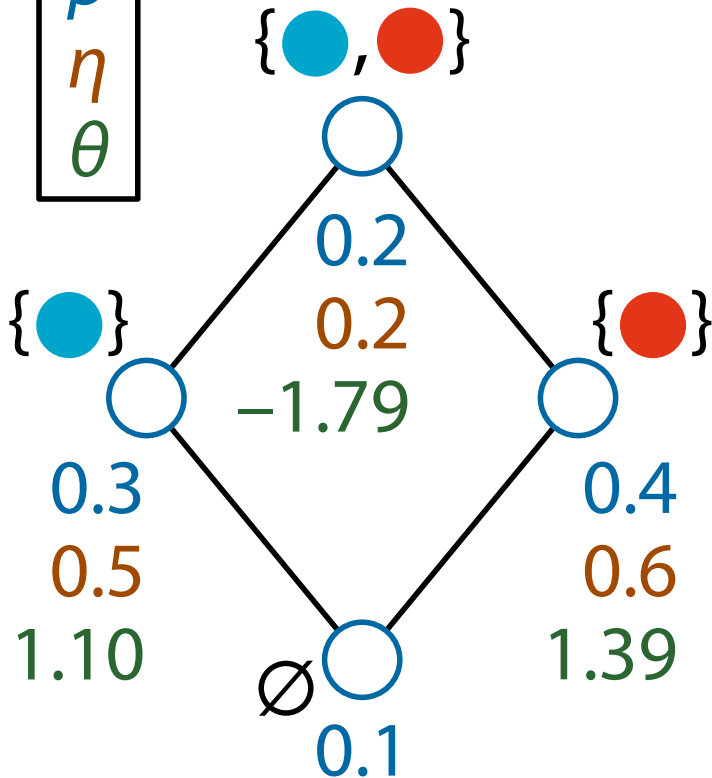
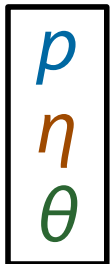


$$\log p(x) = \sum_{s \leq x} \theta(s)$$

Natural parameter θ

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

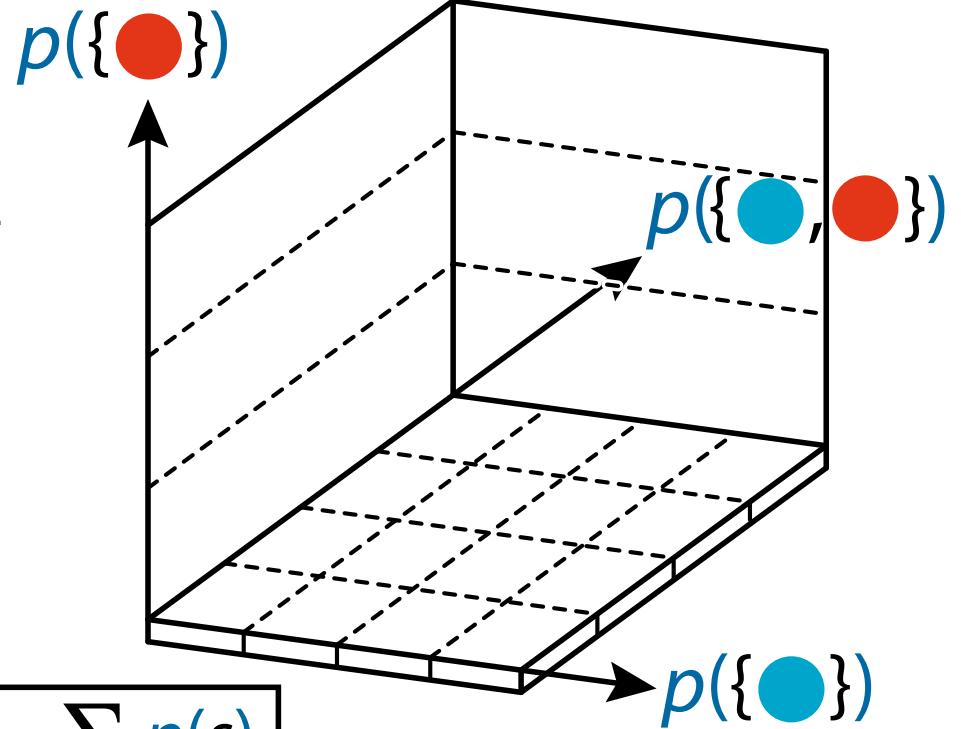
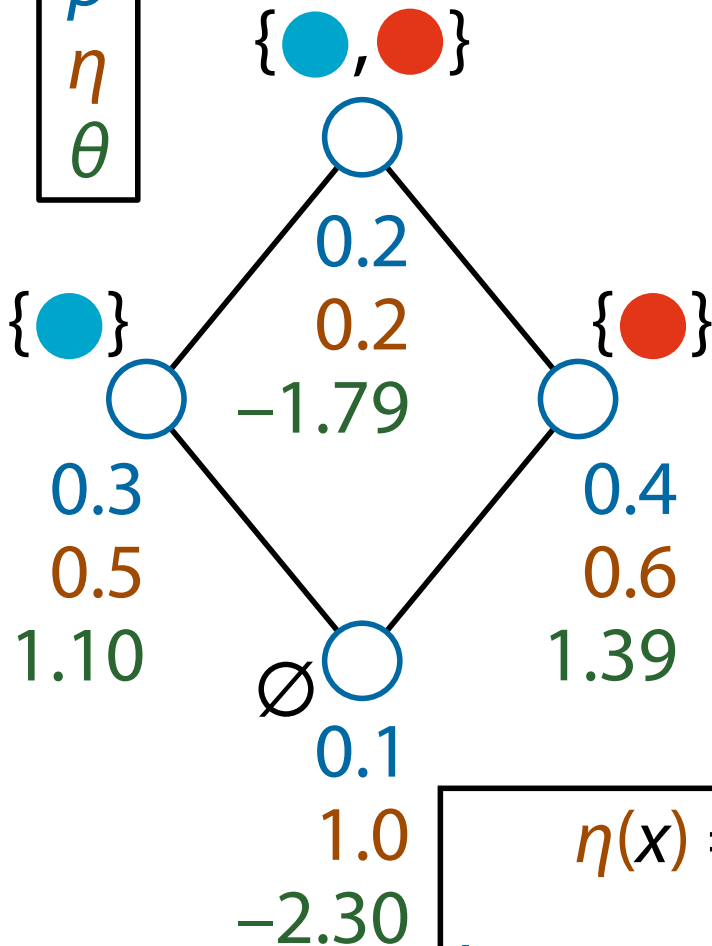
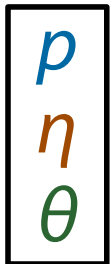
Triple for each node



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

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

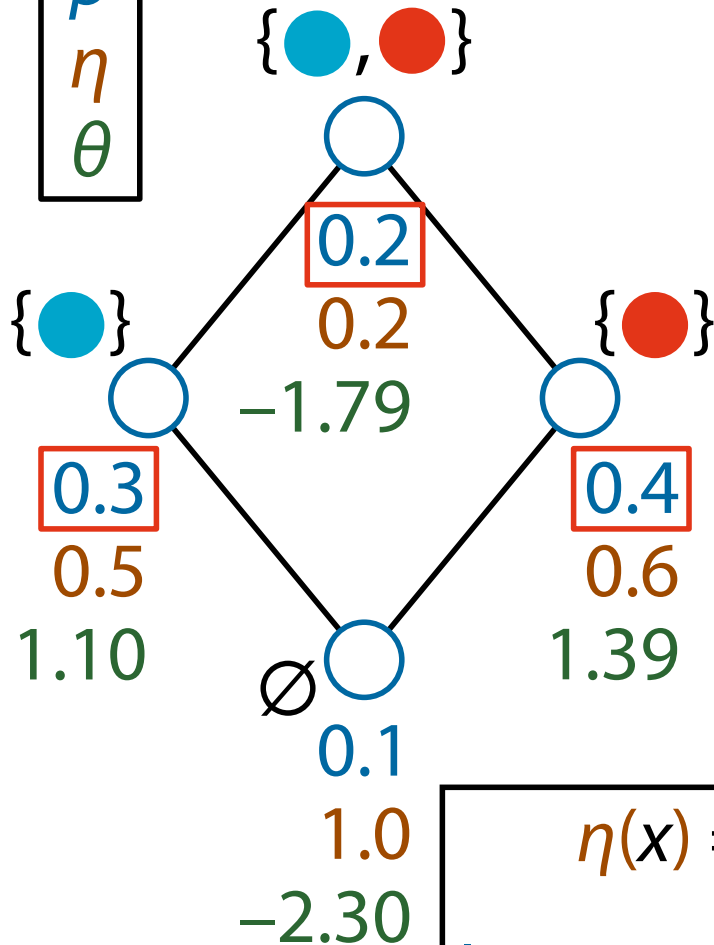
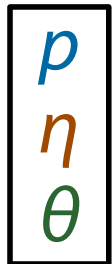
Triple for each node



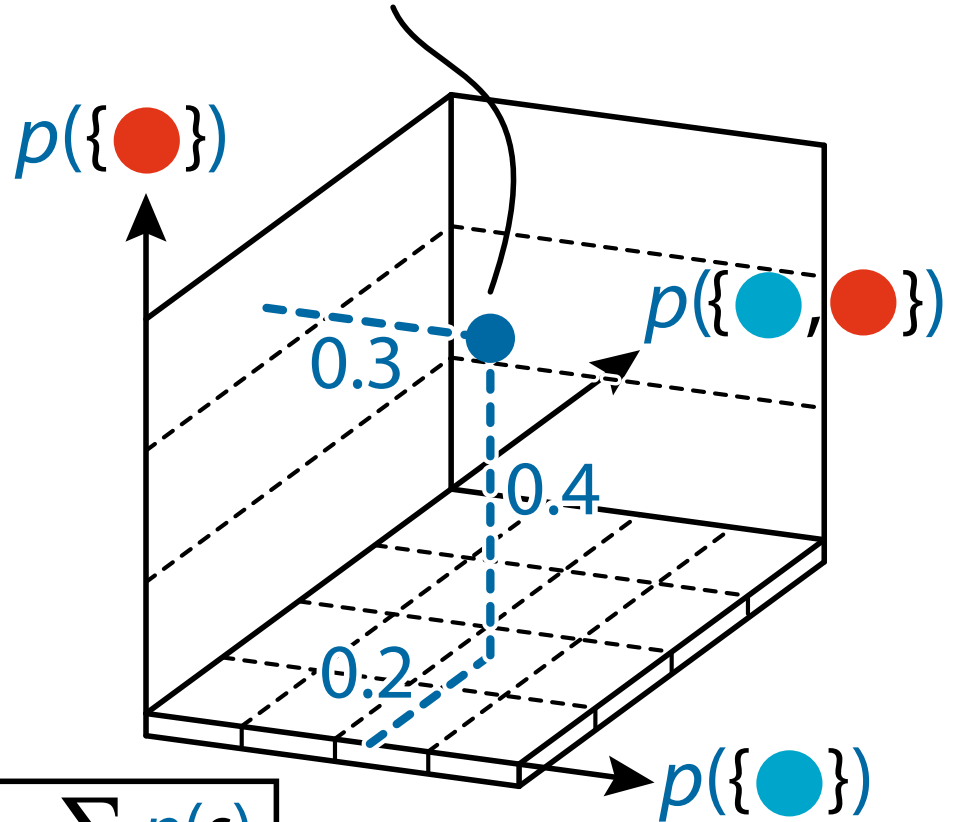
$$\eta(x) = \sum_{s \geq x} p(s)$$

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

Triple for each node



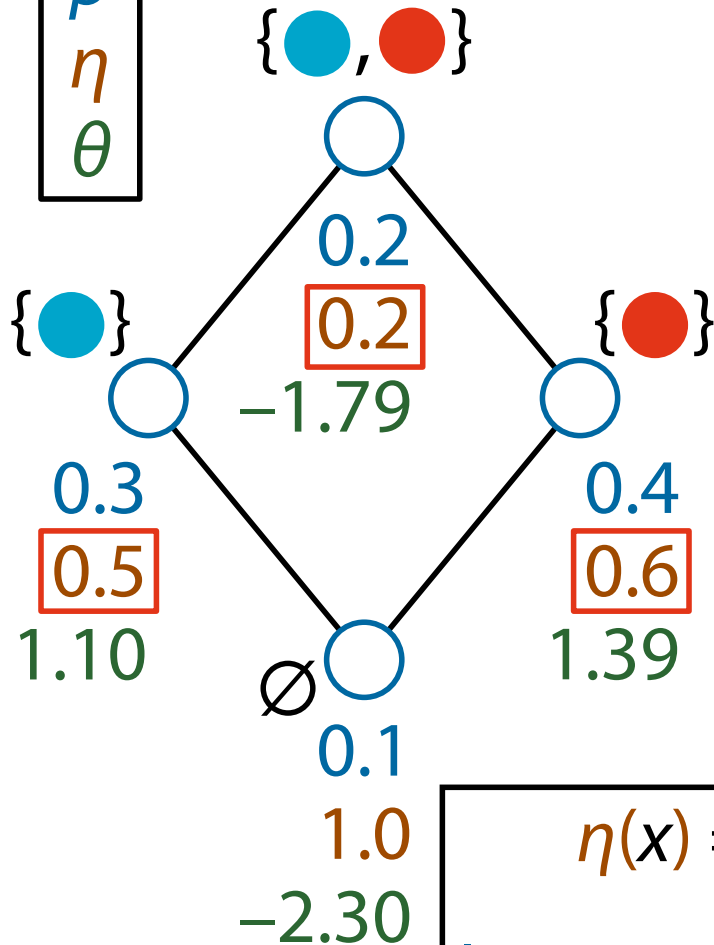
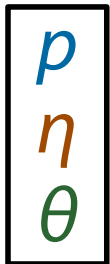
Probability distribution is a "point" in 3D space



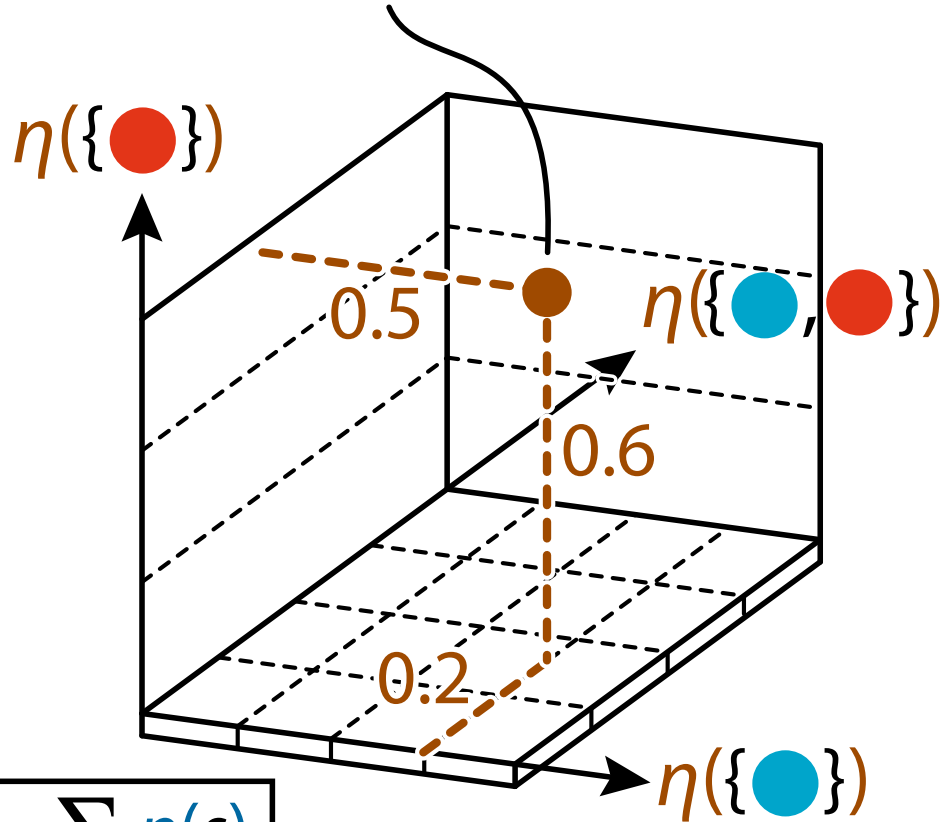
$$\eta(x) = \sum_{s \geq x} p(s)$$

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

Triple for each node



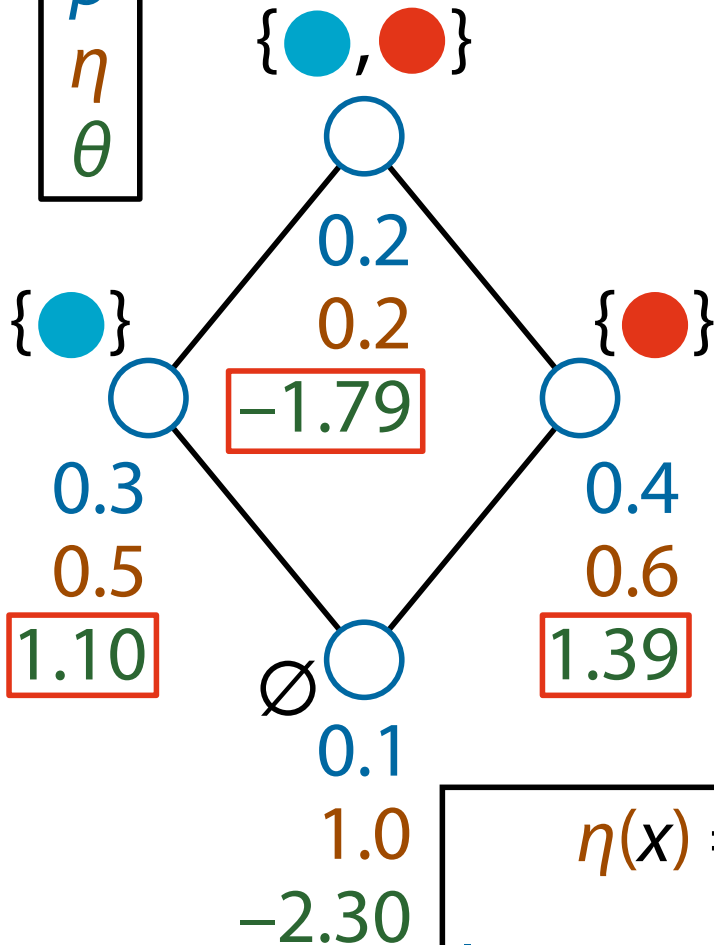
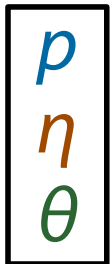
Probability distribution is a "point" in 3D space



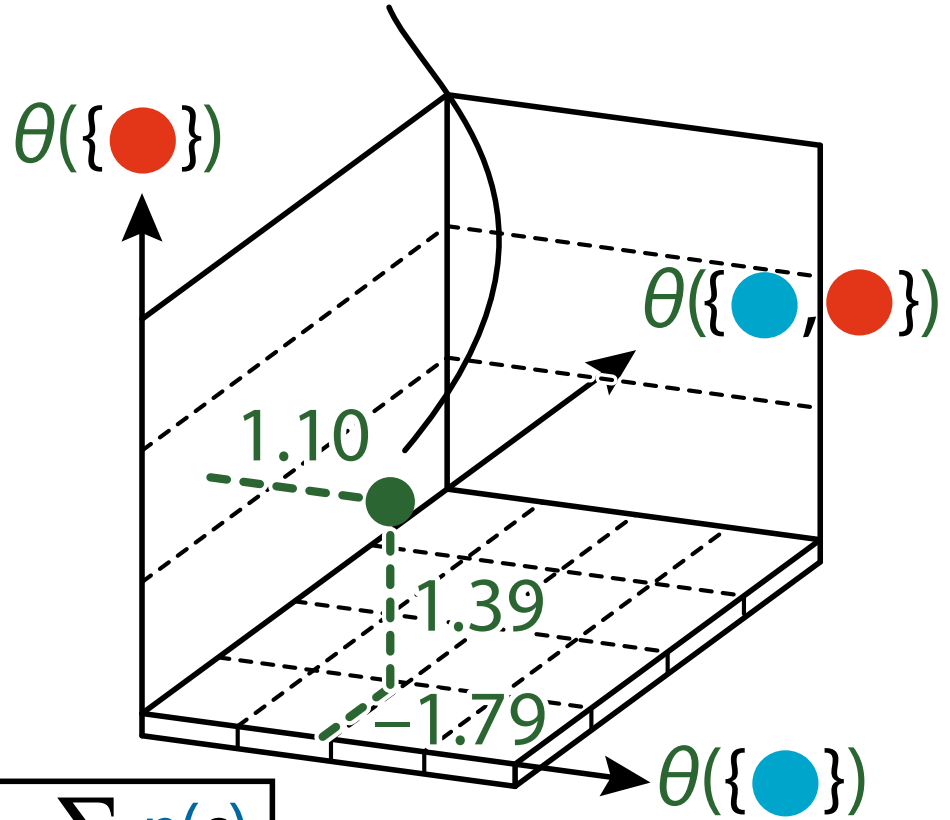
$$\eta(x) = \sum_{s \geq x} p(s)$$

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

Triple for each node



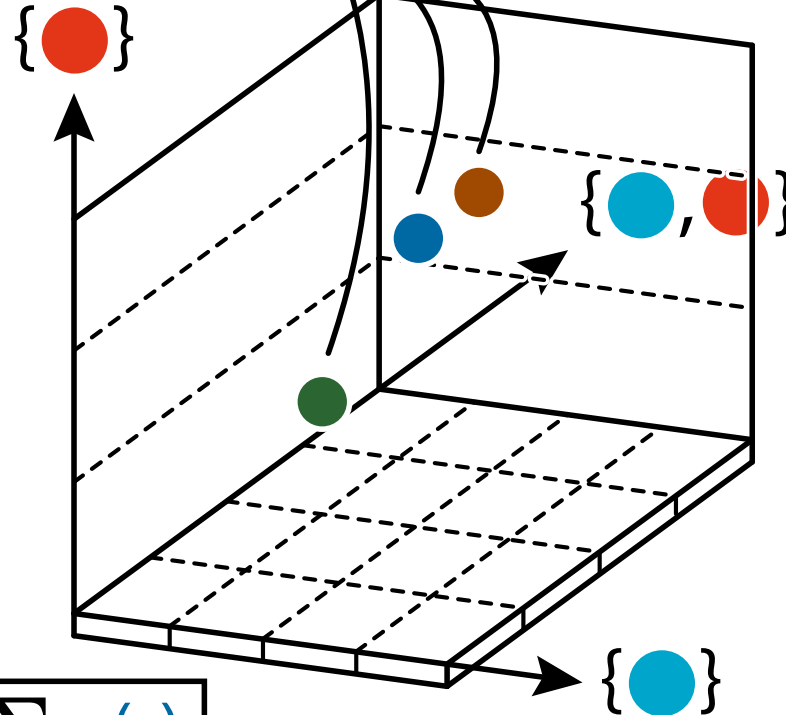
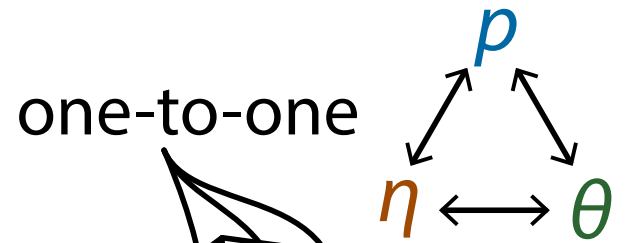
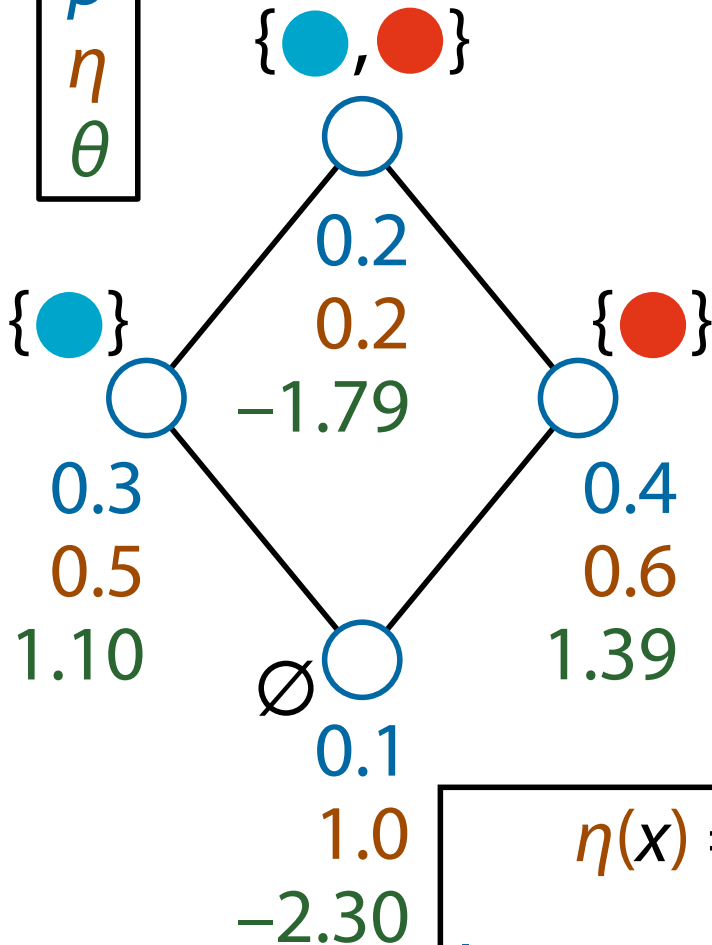
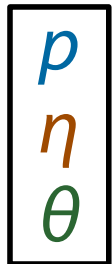
Probability distribution is a "point" in 3D space



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

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

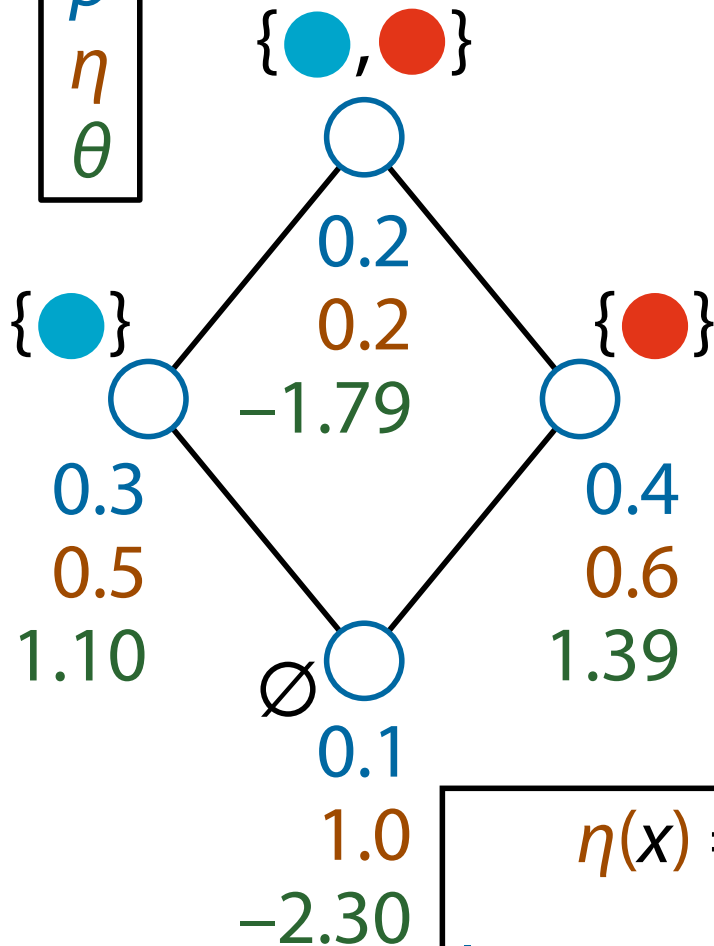
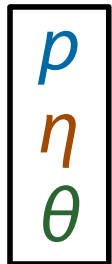
Triple for each node



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

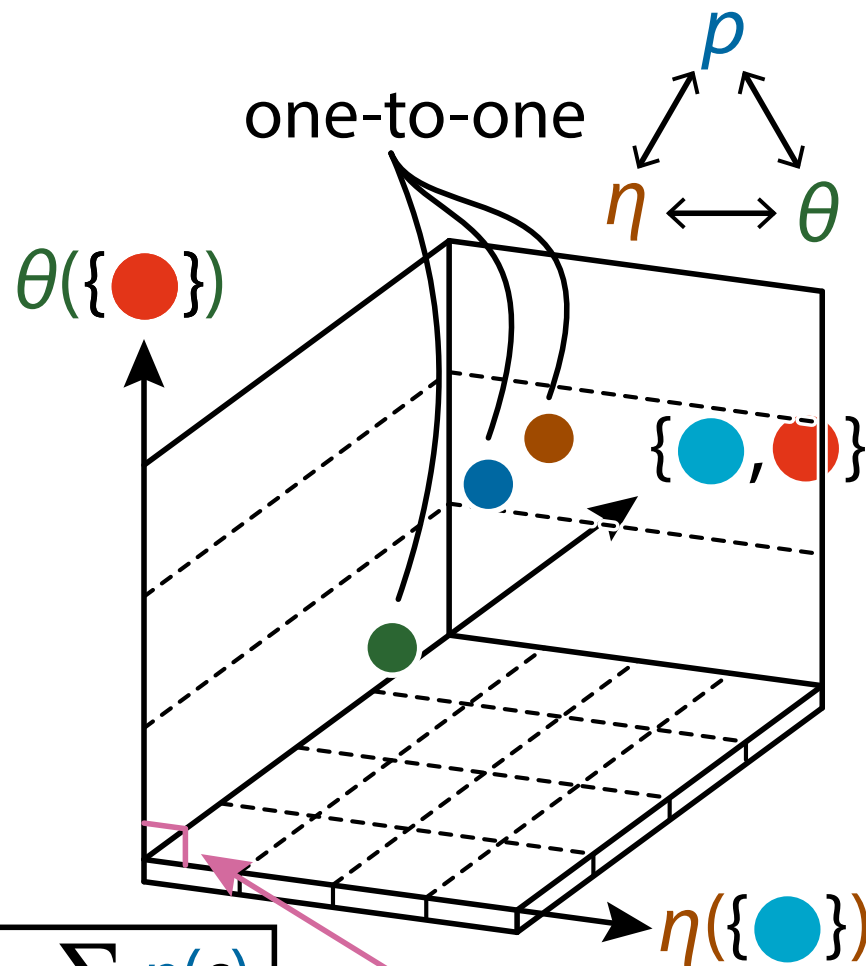
$$\log p(x) = \sum_{s \leq x} \theta(s)$$

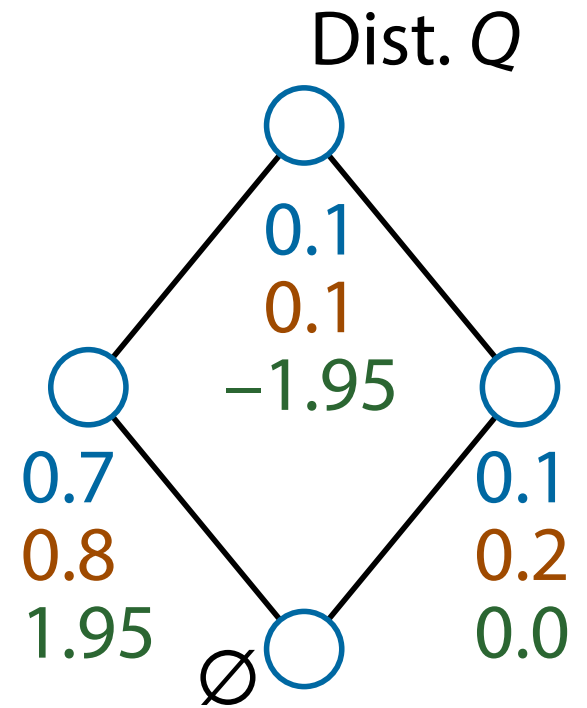
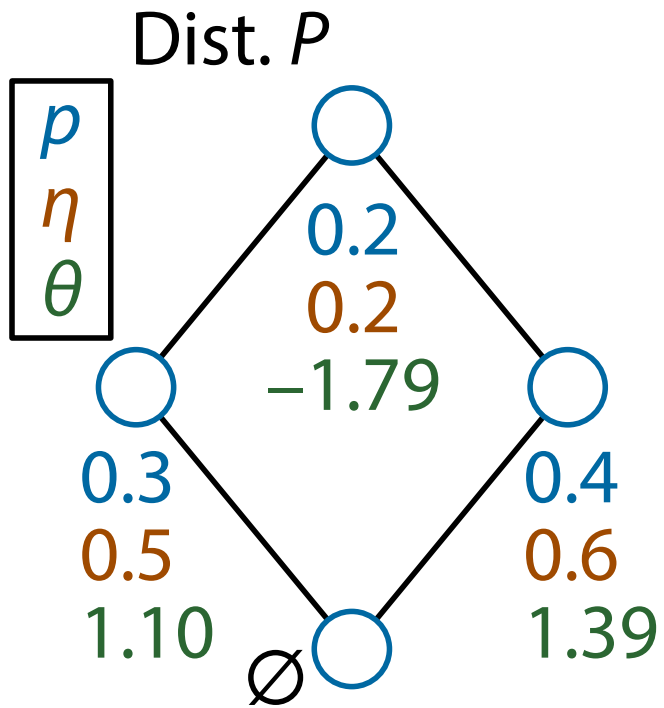
Triple for each node

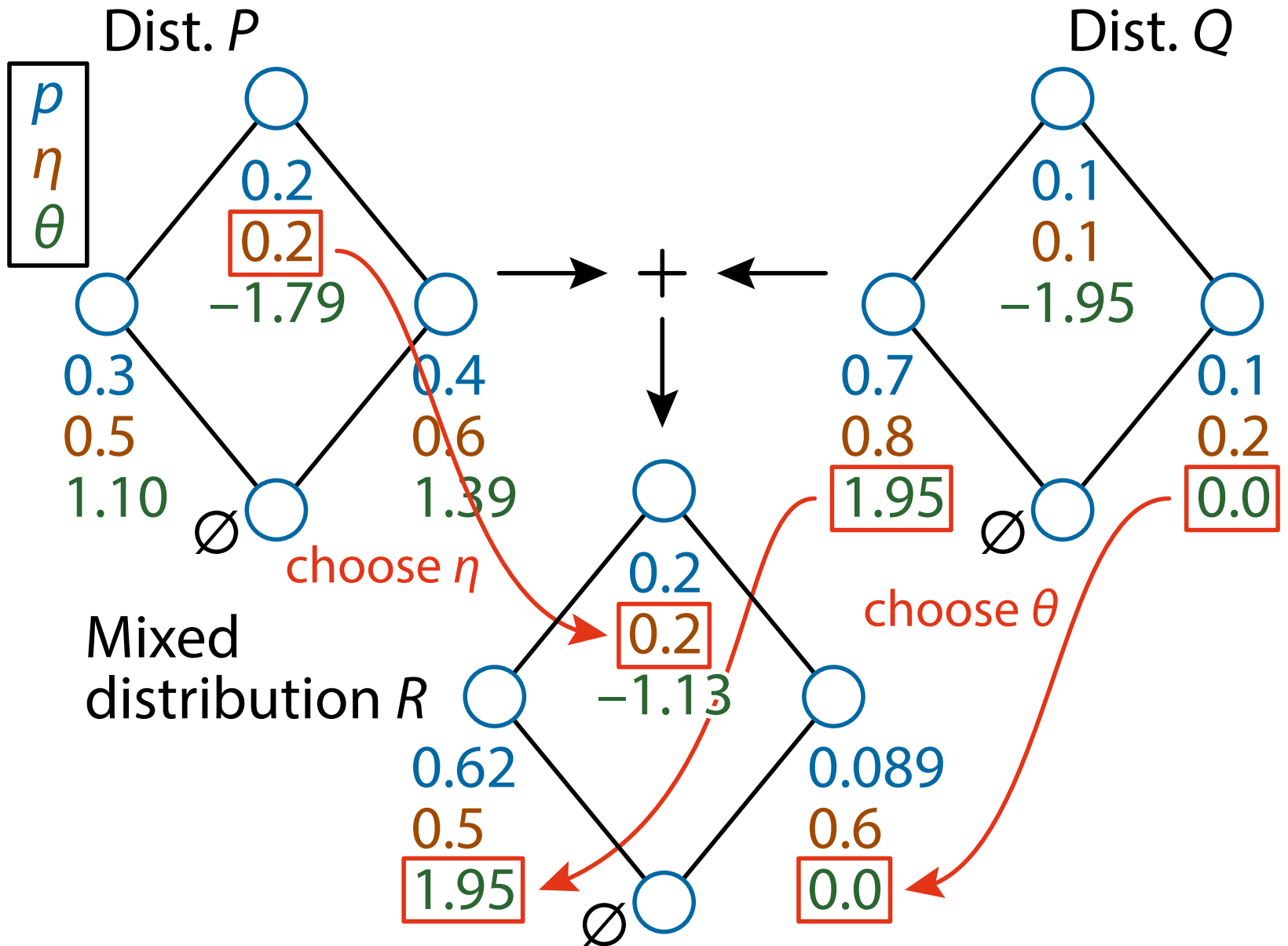


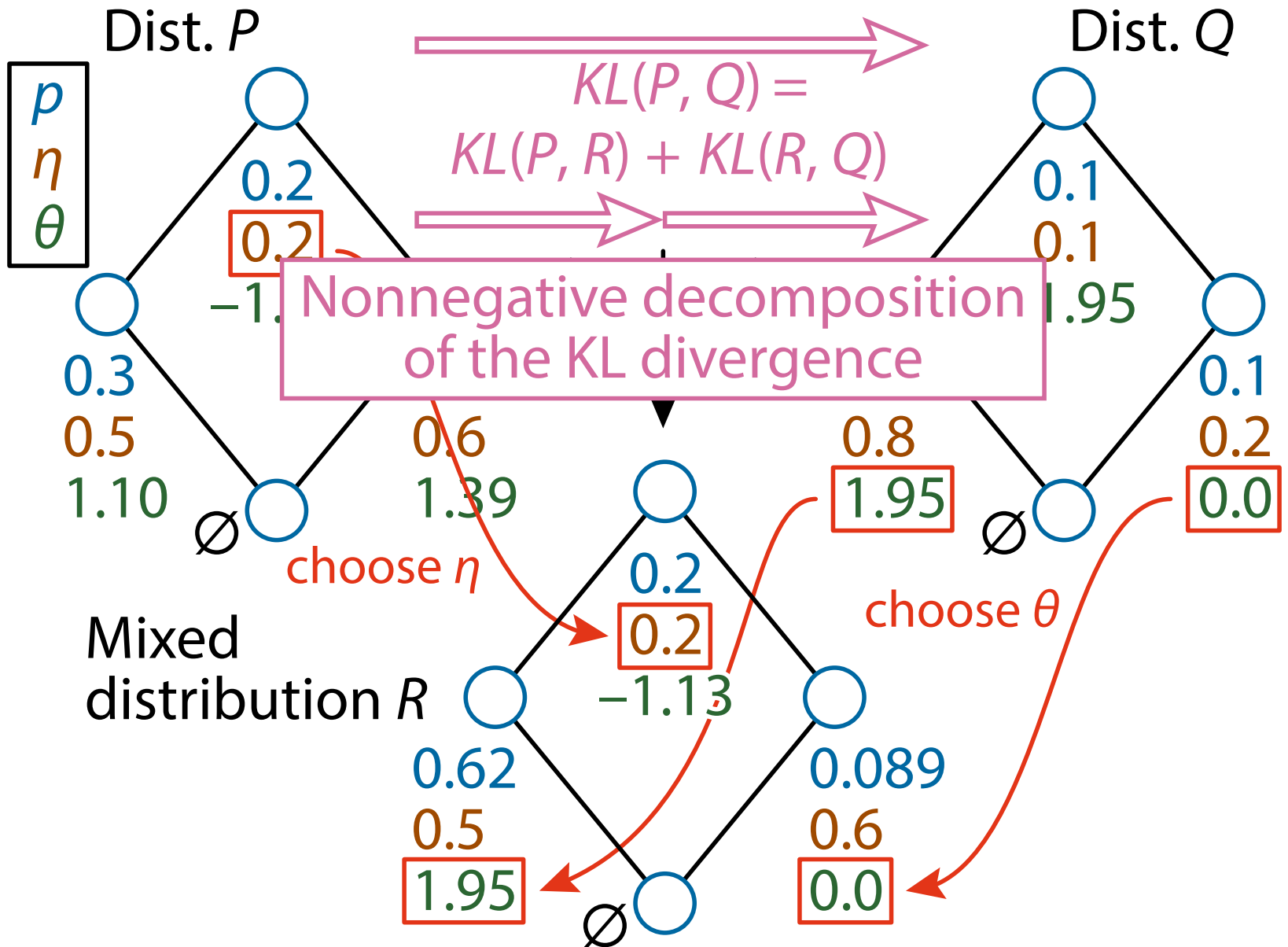
$$\eta(x) = \sum_{s \geq x} p(s)$$

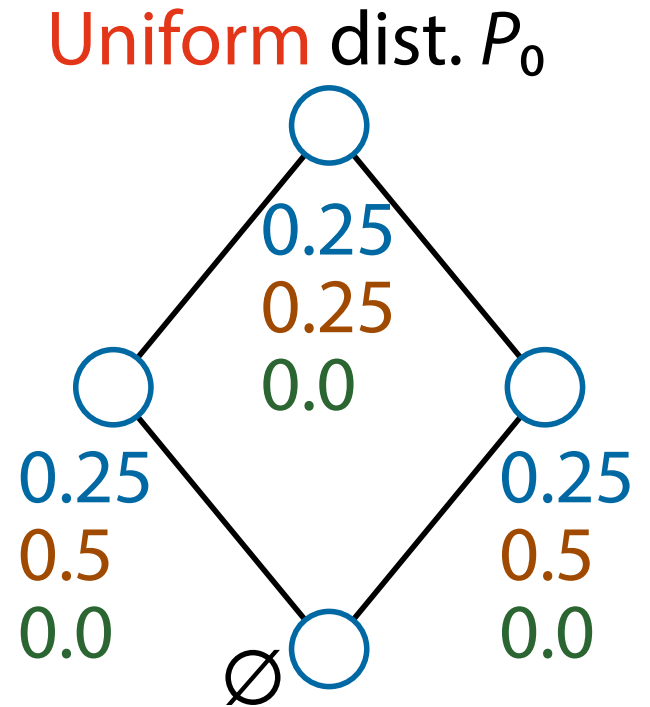
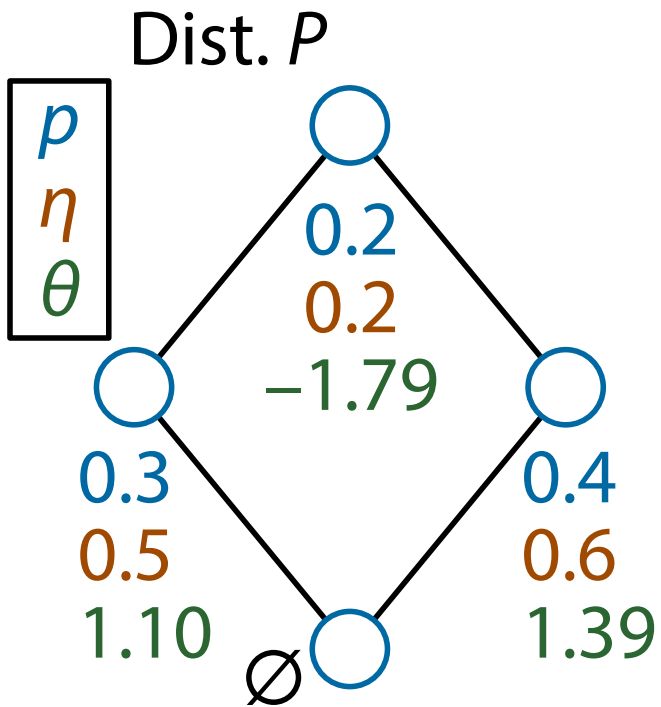
$$\log p(x) = \sum_{s \leq x} \theta(s)$$

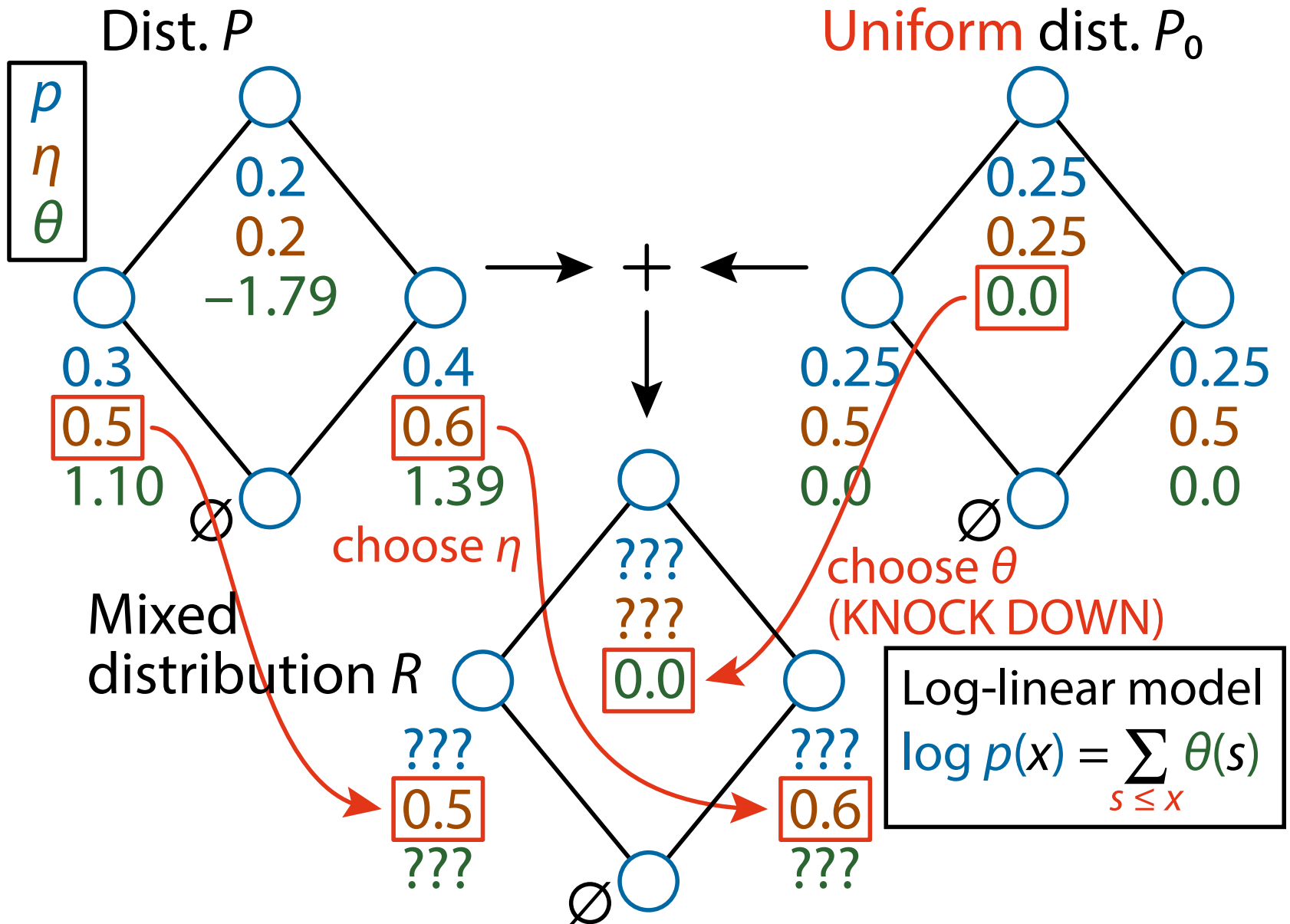


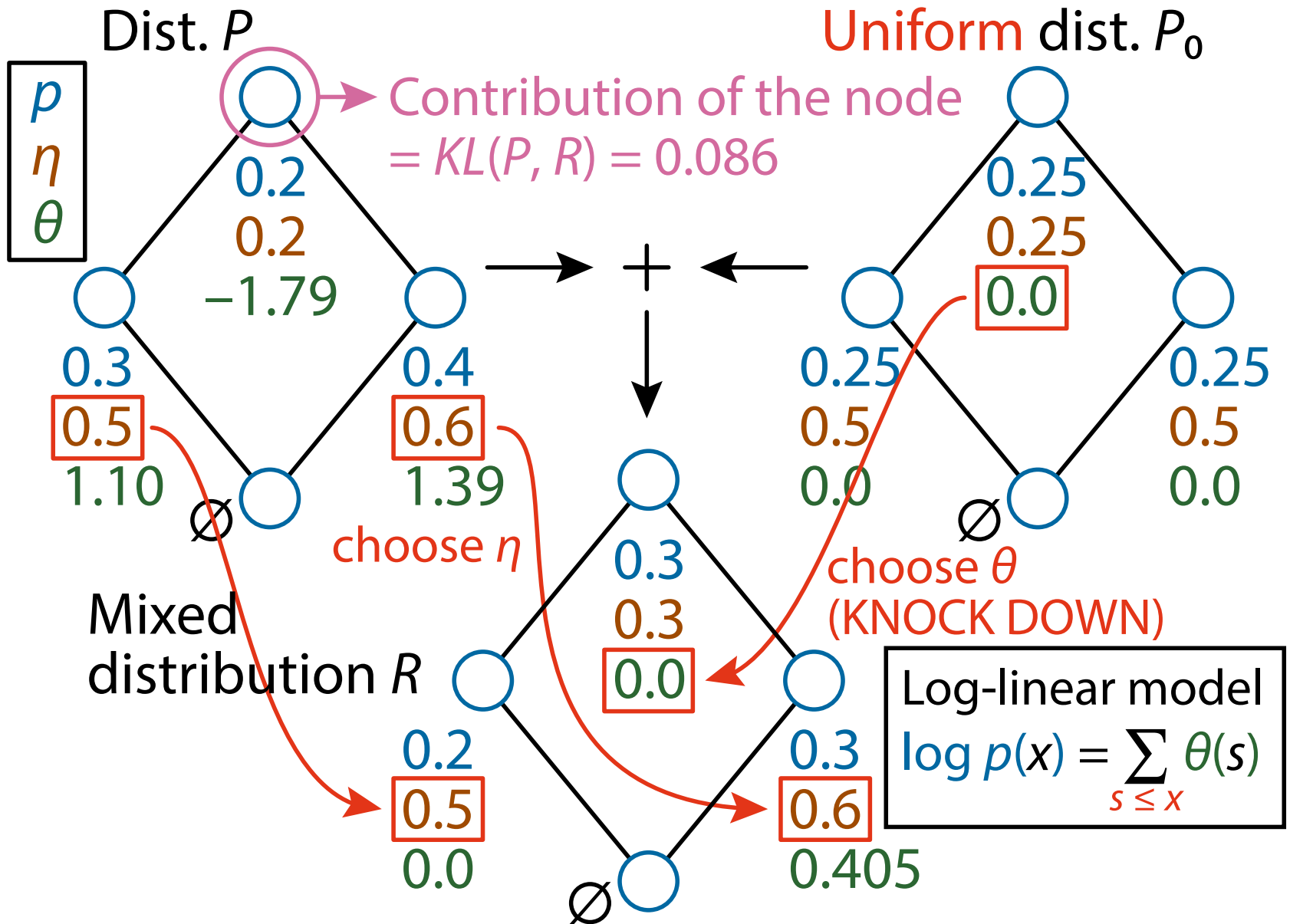


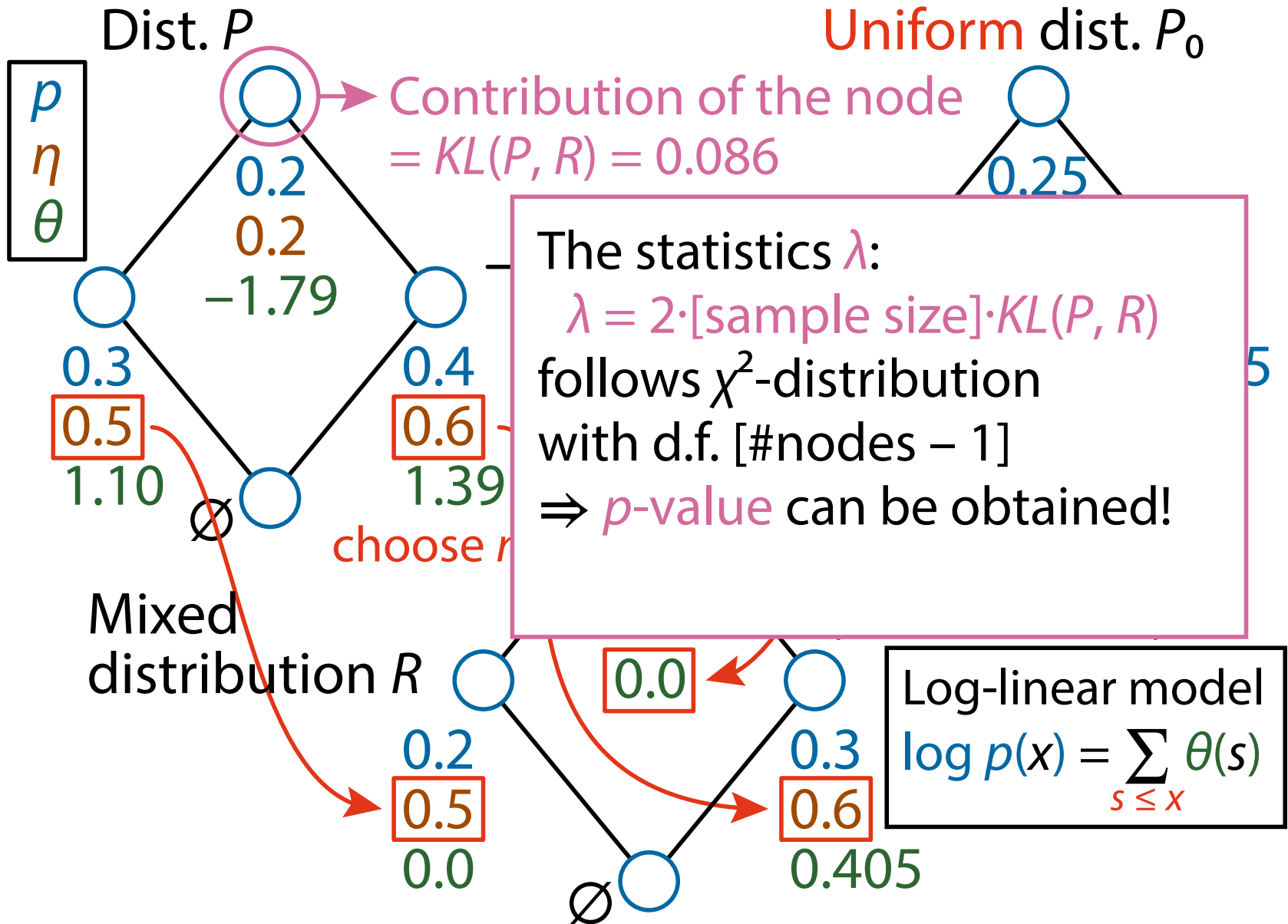








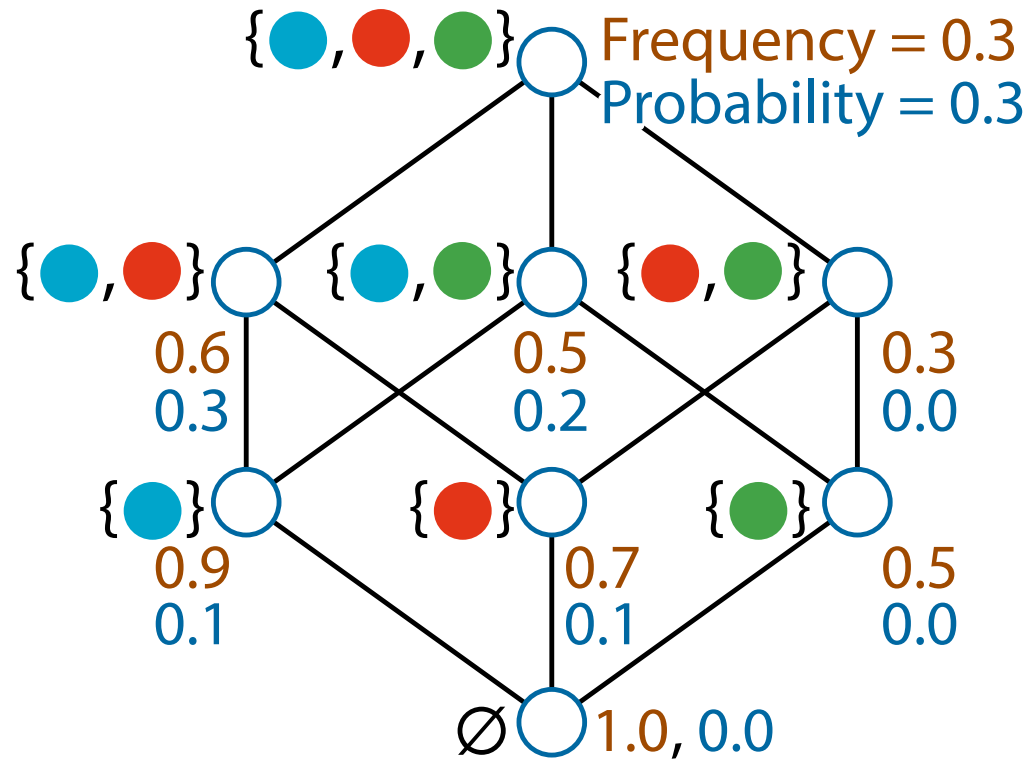




Make a Poset from Data

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

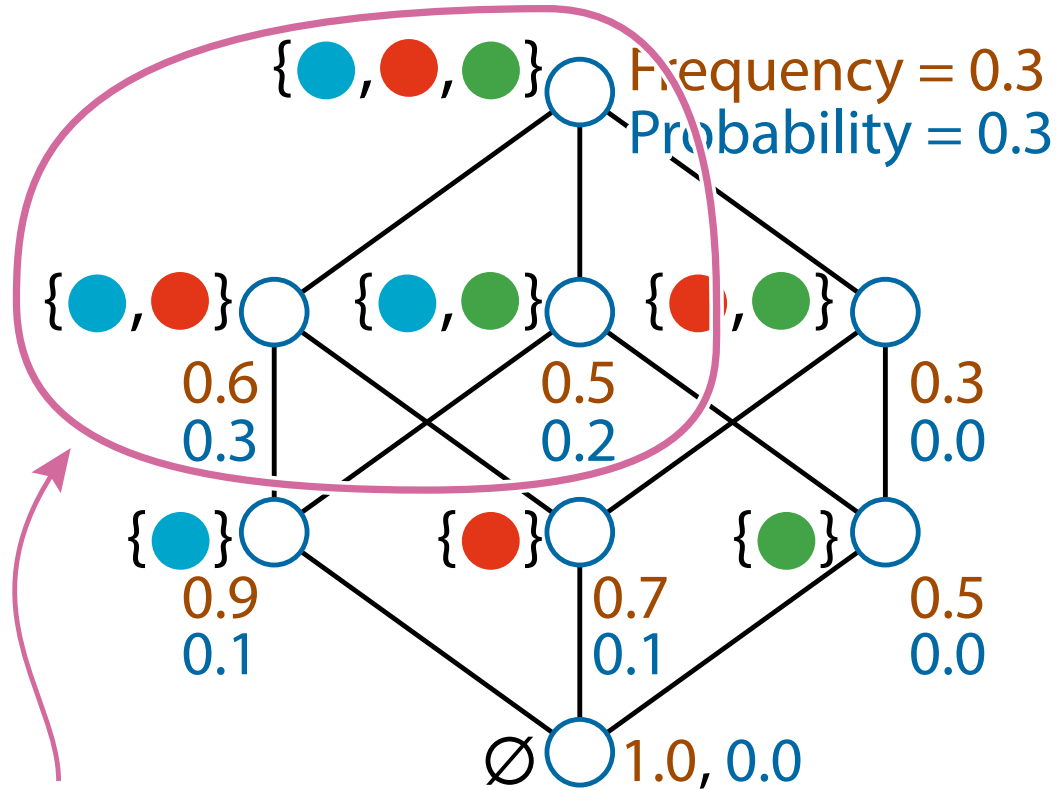


Number of nodes = $2^{\text{\#features}}$
 \Rightarrow combinatorial explosion!

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

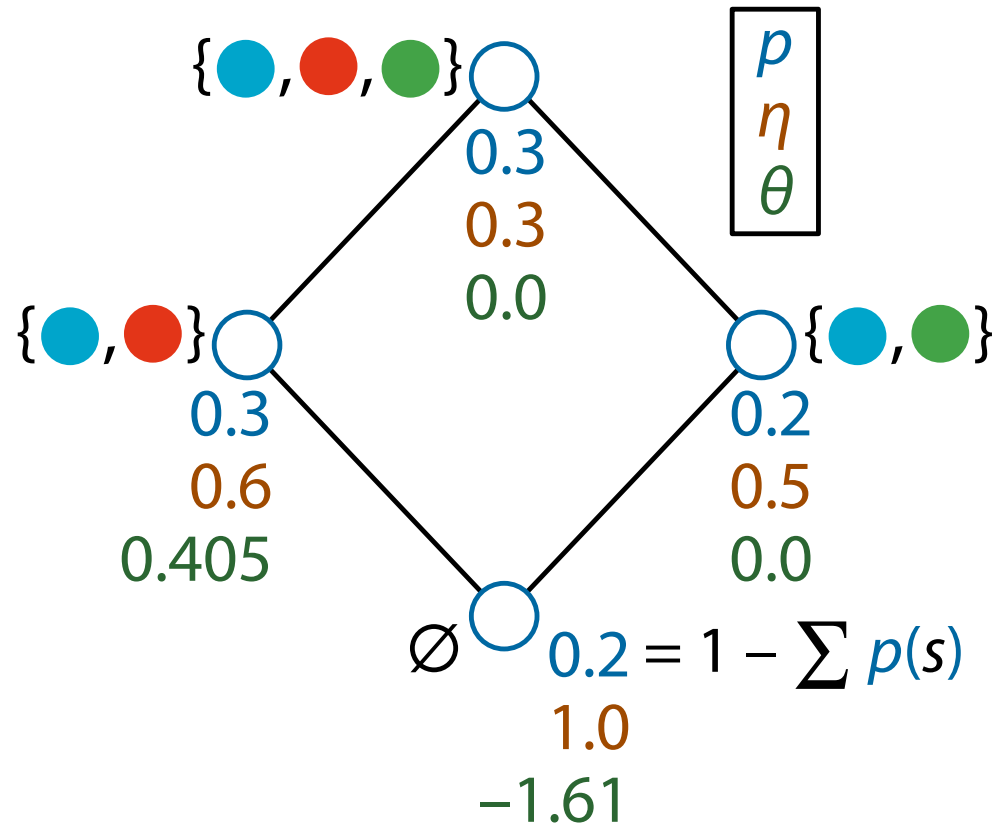


Probability ≥ 0.2
(user specified threshold)

Remove Nodes with Probability 0

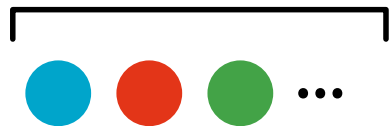
Dataset

	●	●	●
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



Example on Real Data (kosarak)

features: 41,270

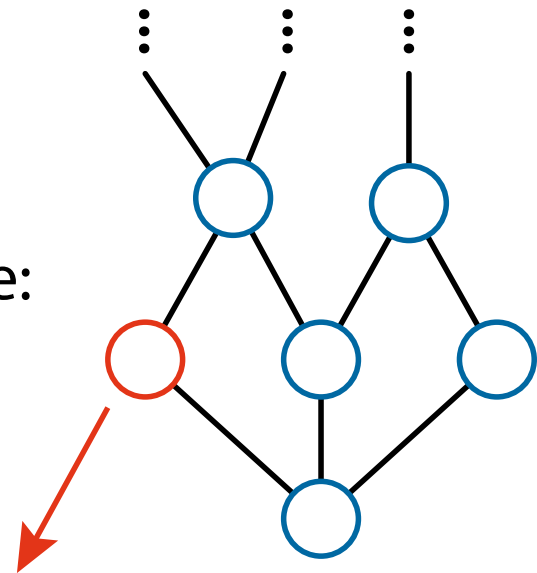


ID 1:	1	1	0
ID 2:	1	1	1
ID 3:	1	1	0 ...
ID 4:	1	1	1
ID 5:	1	1	0
⋮	⋮		

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)

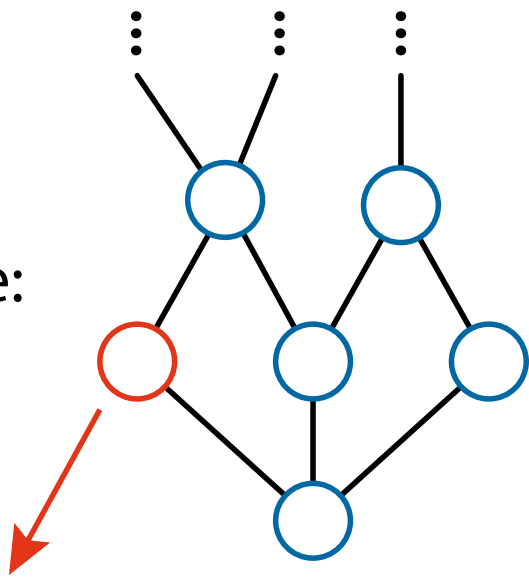
features: 468

				...
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	
⋮	⋮	⋮		

Total runtime:
4.95 seconds

Sample size:
340,183

nodes: 281
(Threshold: 5×10^{-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 assess the significance of any-order interactions
- Related papers:
 - 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)