

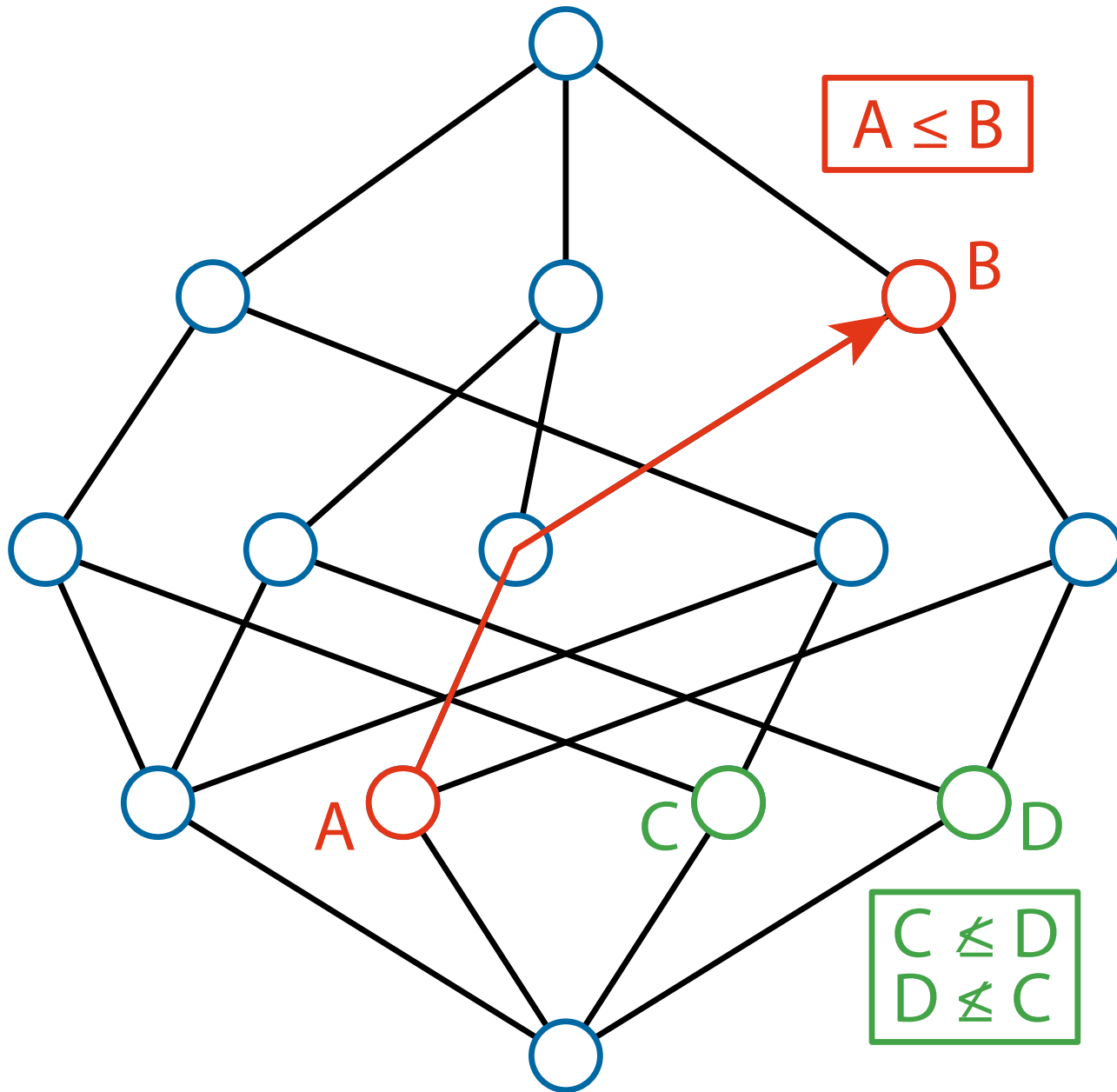
March 17, 2016  
@RIKEN BSI

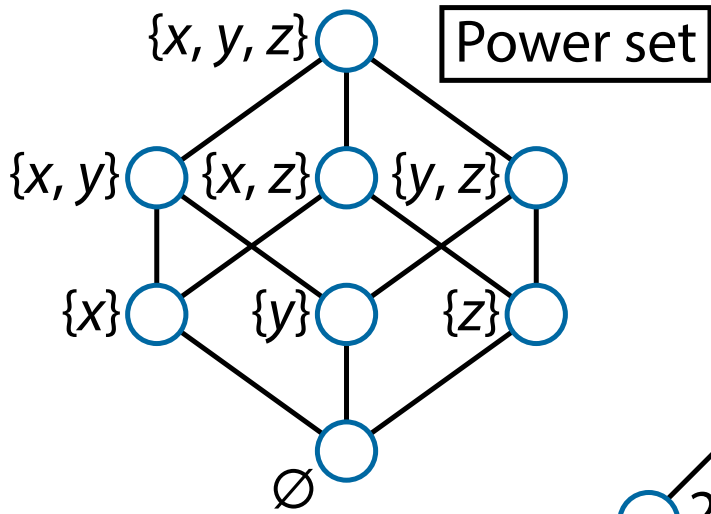


# Statistical Analysis on Order Structures

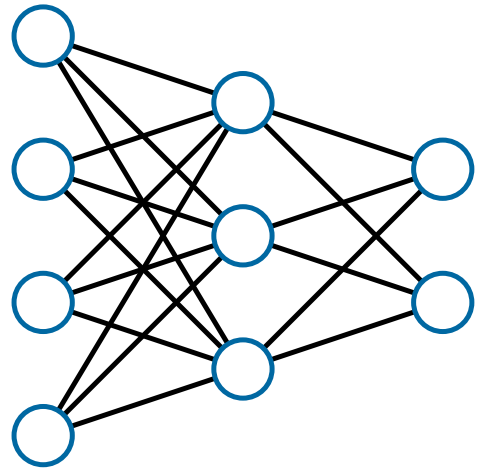
---

Mahito Sugiyama (ISIR, Osaka University)  
(杉山磨人; 大阪大学産業科学研究所)

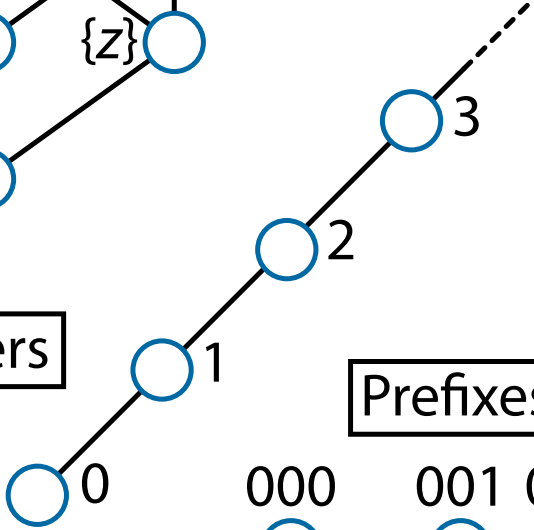




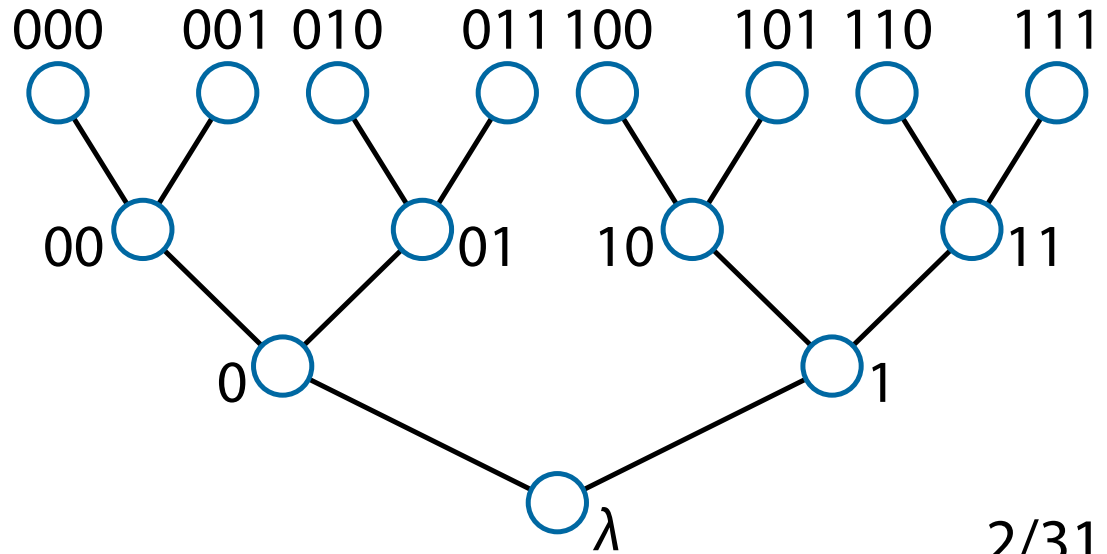
**Neural networks**



**Positive integers**



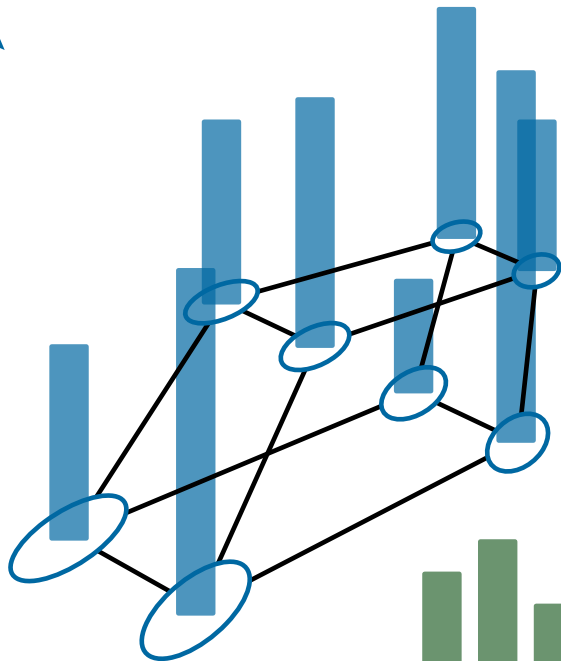
**Prefixes**



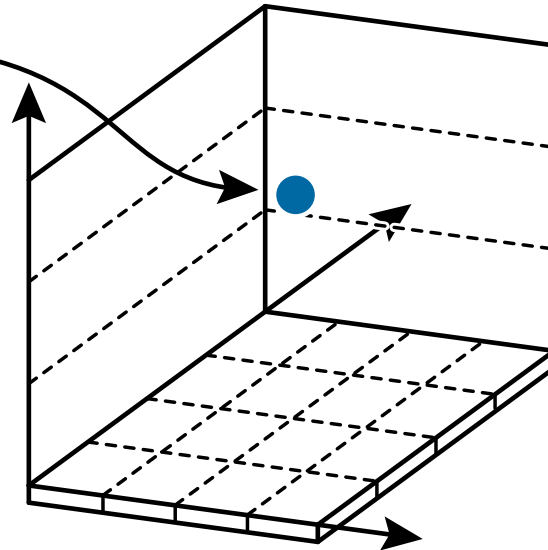
# Summary

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

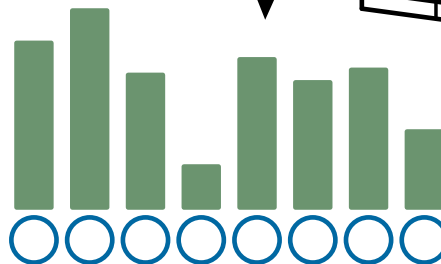
Probability



Information  
geometry



Decomposition in  
the **log-linear model**



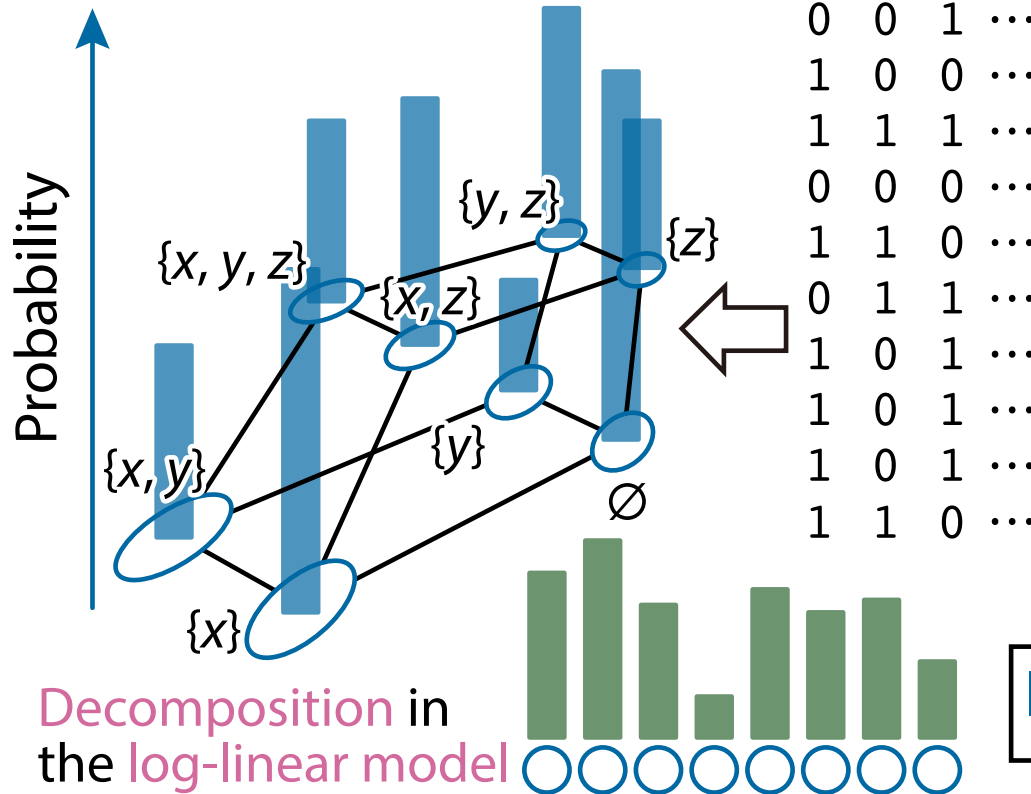
$$\log p(x) = \sum \theta(s)$$

S. Amari, Information geometry on hierarchy of probability distributions, IEEE TIT 2001

M. Sugiyama, H. Nakahara, K. Tsuda, Information Decomposition on Structured Space, arXiv 2016

# Summary

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



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	1	0	...

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

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

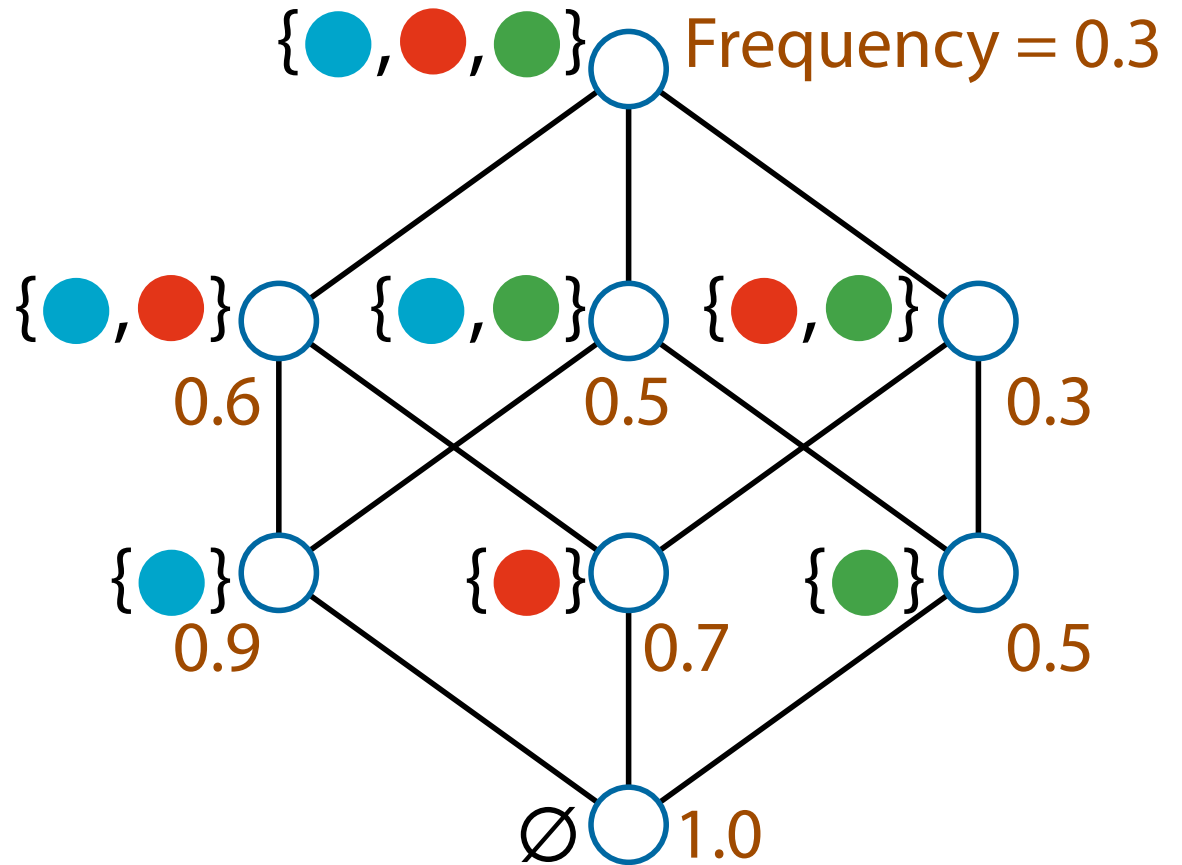
S. Amari, Information geometry on hierarchy of probability distributions, IEEE TIT 2001  
 M. Sugiyama, H. Nakahara, K. Tsuda, Information Decomposition on Structured Space, arXiv 2016

# 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

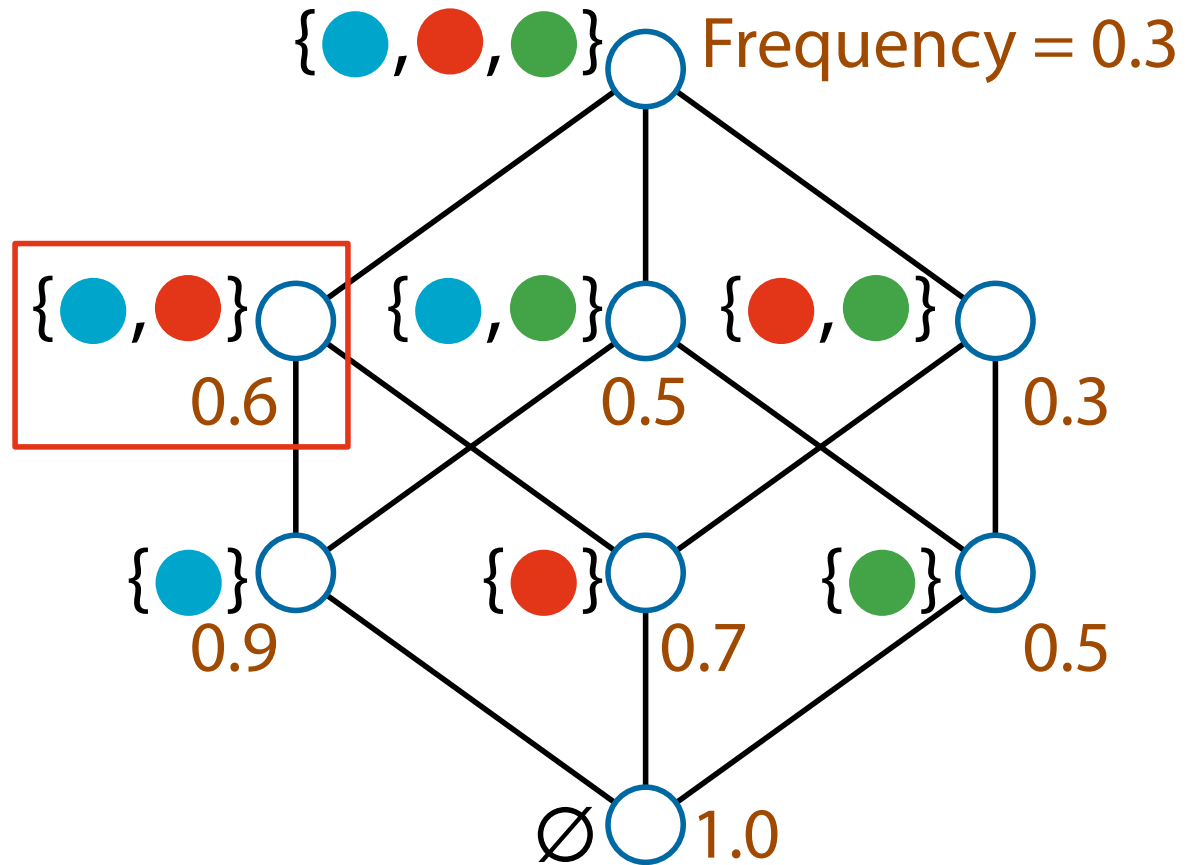
# Itemset lattice



# 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



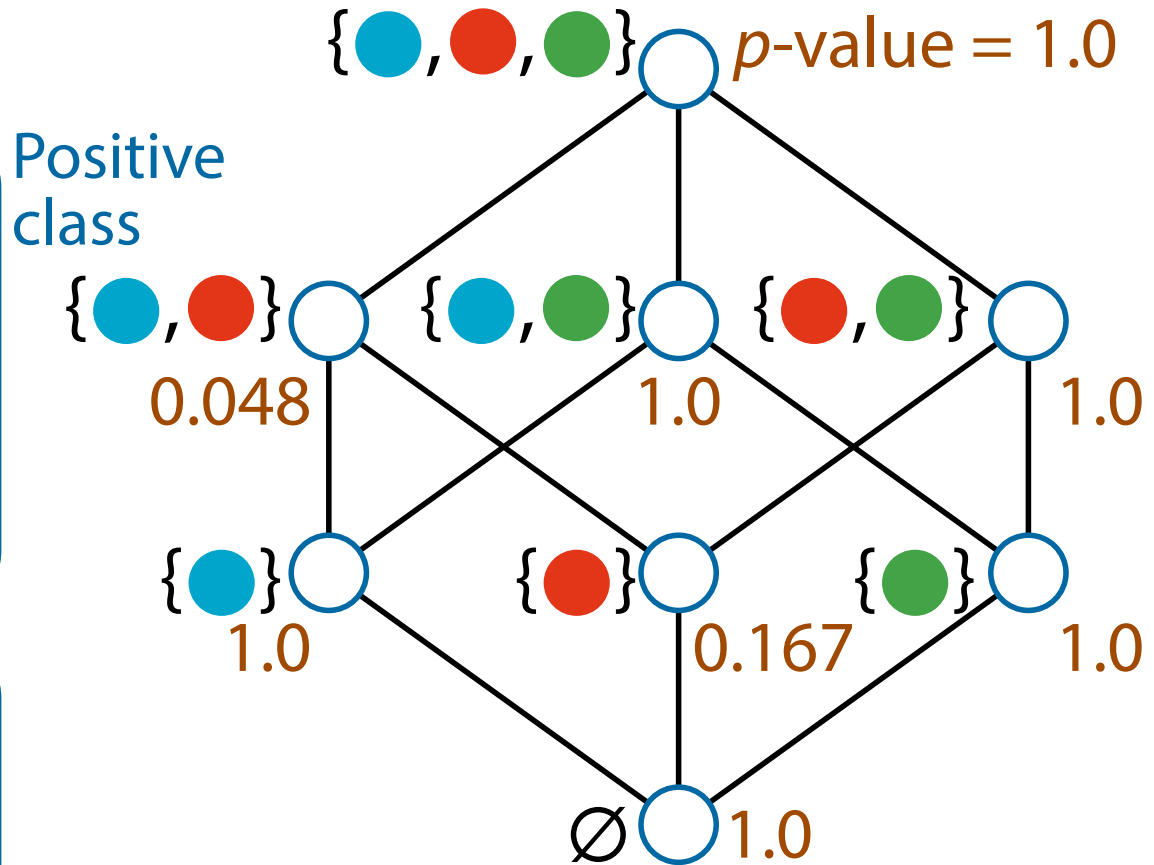
# Transaction database



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
ID 10:	0	1	0

# Itemset lattice



Positive class

Negative class

LAMP (Terada et al. PNAS 2013)

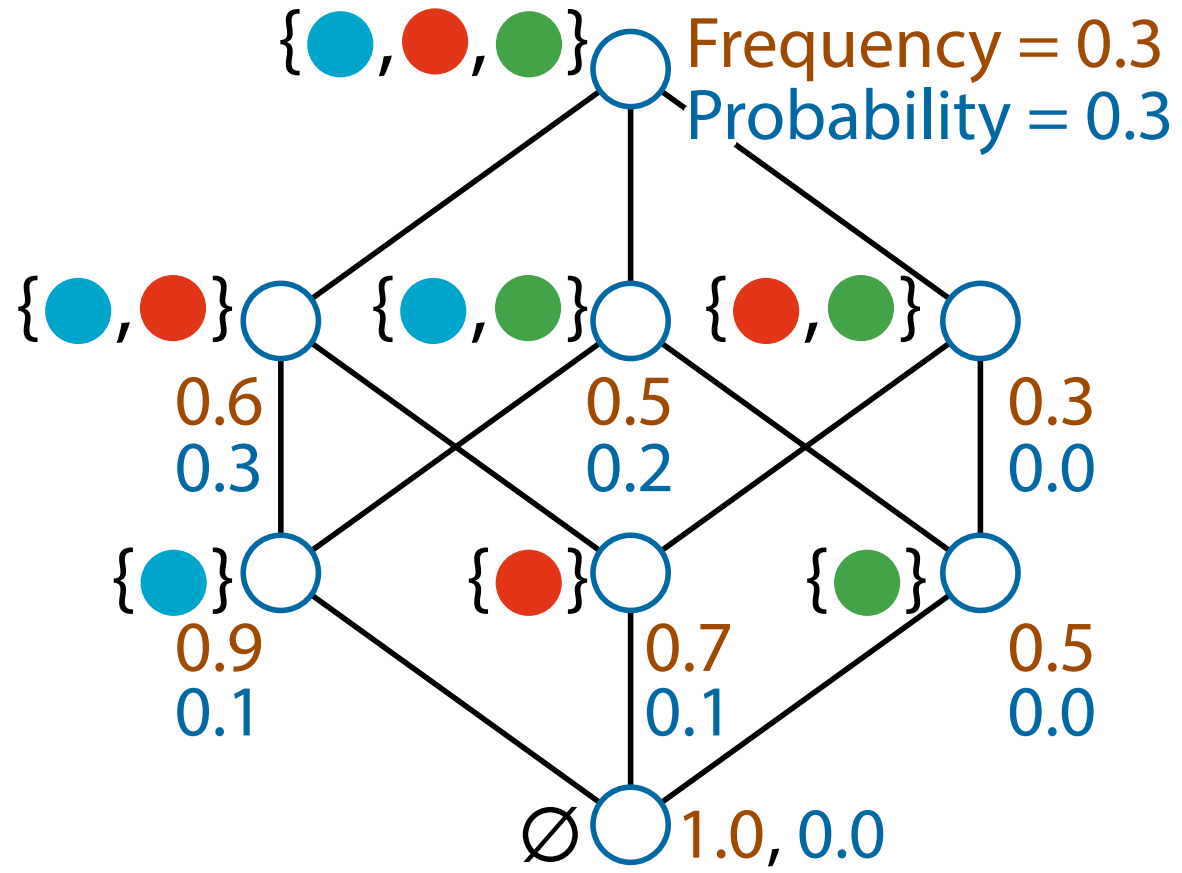
Westfall-Young light (Llinares-López et al. KDD 2015)



# Transaction database

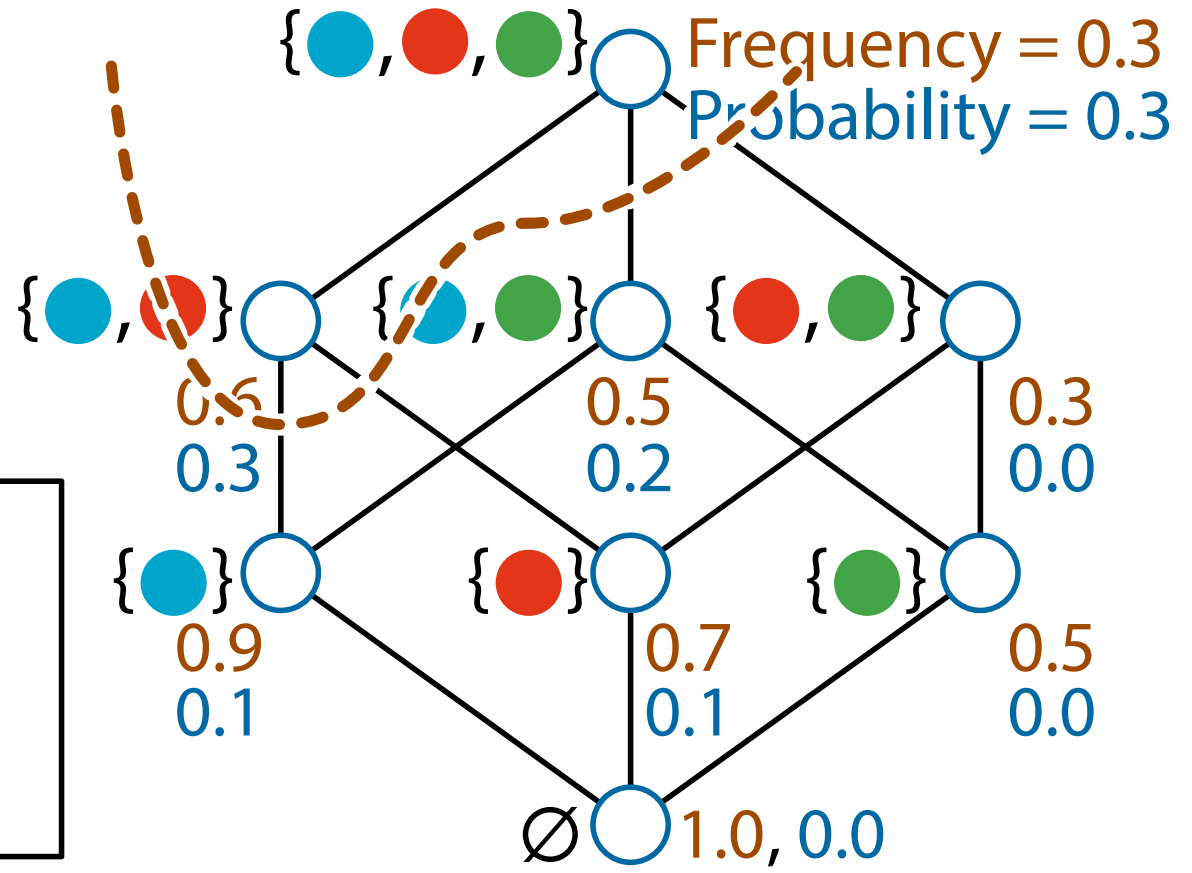
	<span style="color: blue;">●</span>	<span style="color: red;">●</span>	<span style="color: green;">●</span>
ID 1:	1	1	0
ID 2:	1	1	1
ID 3:	1	1	0
ID 4:	1	1	1
ID 5:	1	1	0
ID 6:	1	0	1
ID 7:	1	0	1
ID 8:	1	1	1
ID 9:	1	0	0
ID10:	0	1	0

# Itemset lattice



Upward =  
Pattern mining

Itemset lattice

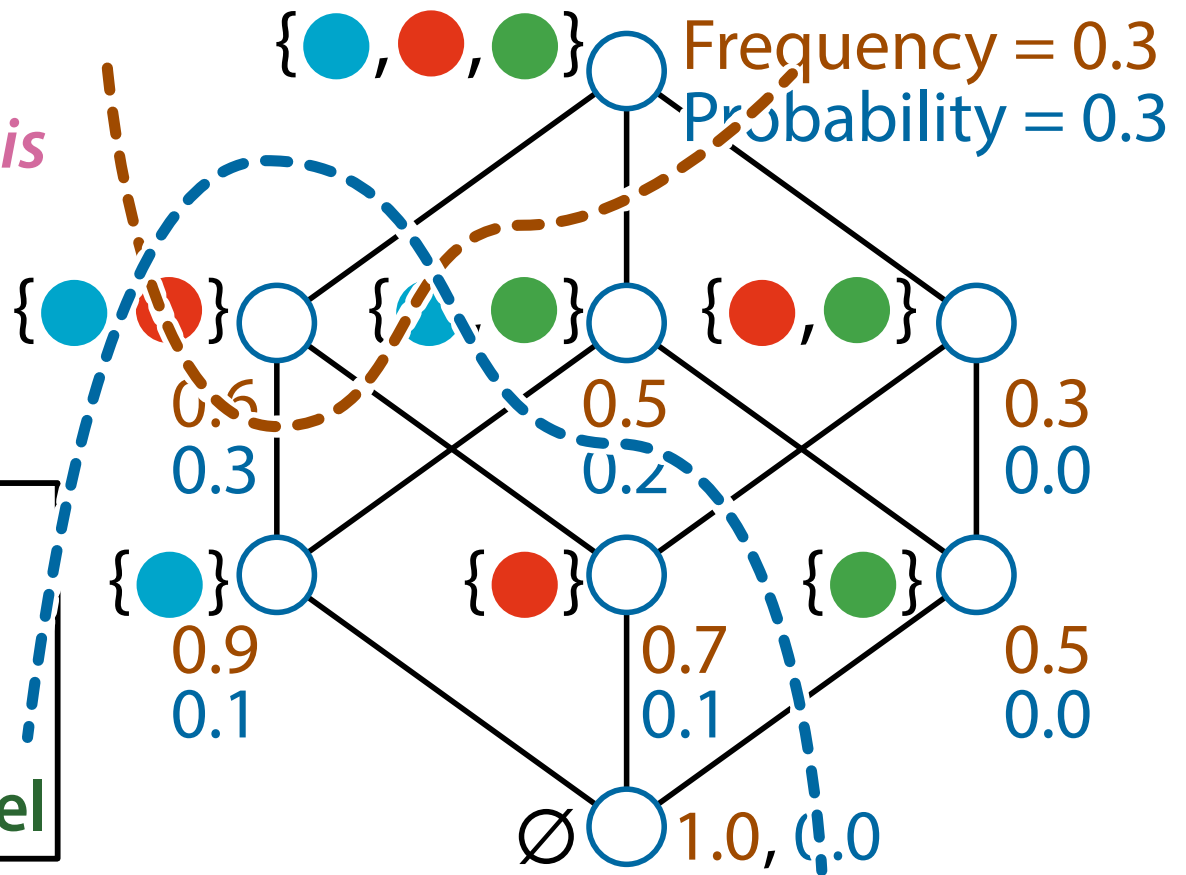


$\eta$ : Frequency  
 $p$ : Probability

$$\eta(\{\text{blue}, \text{red}\}) = p(\{\text{blue}, \text{red}\}) + p(\{\text{blue}, \text{red}, \text{green}\})$$

Upward =  
Pattern mining  
Downward =  
Log-linear analysis

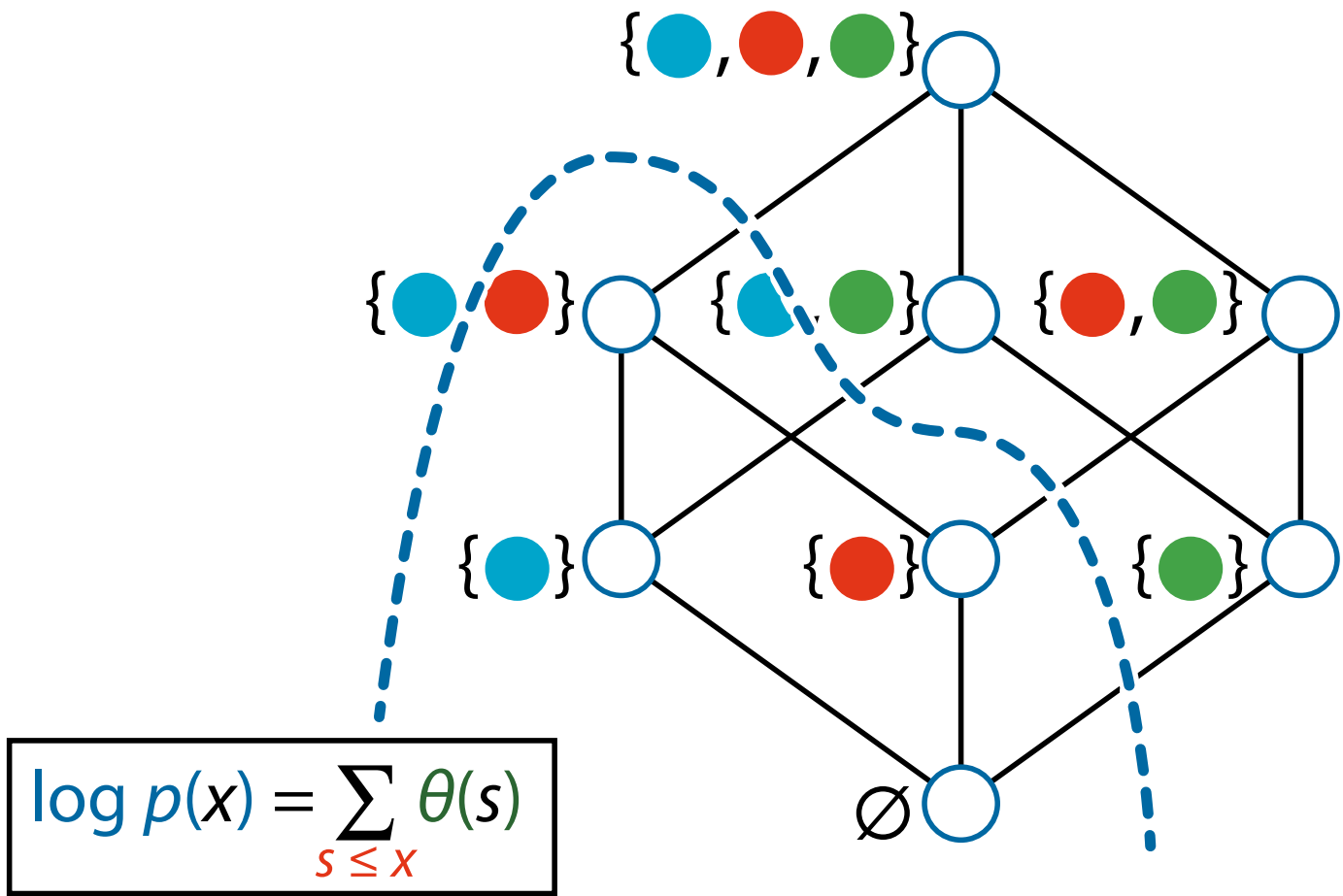
Itemset lattice

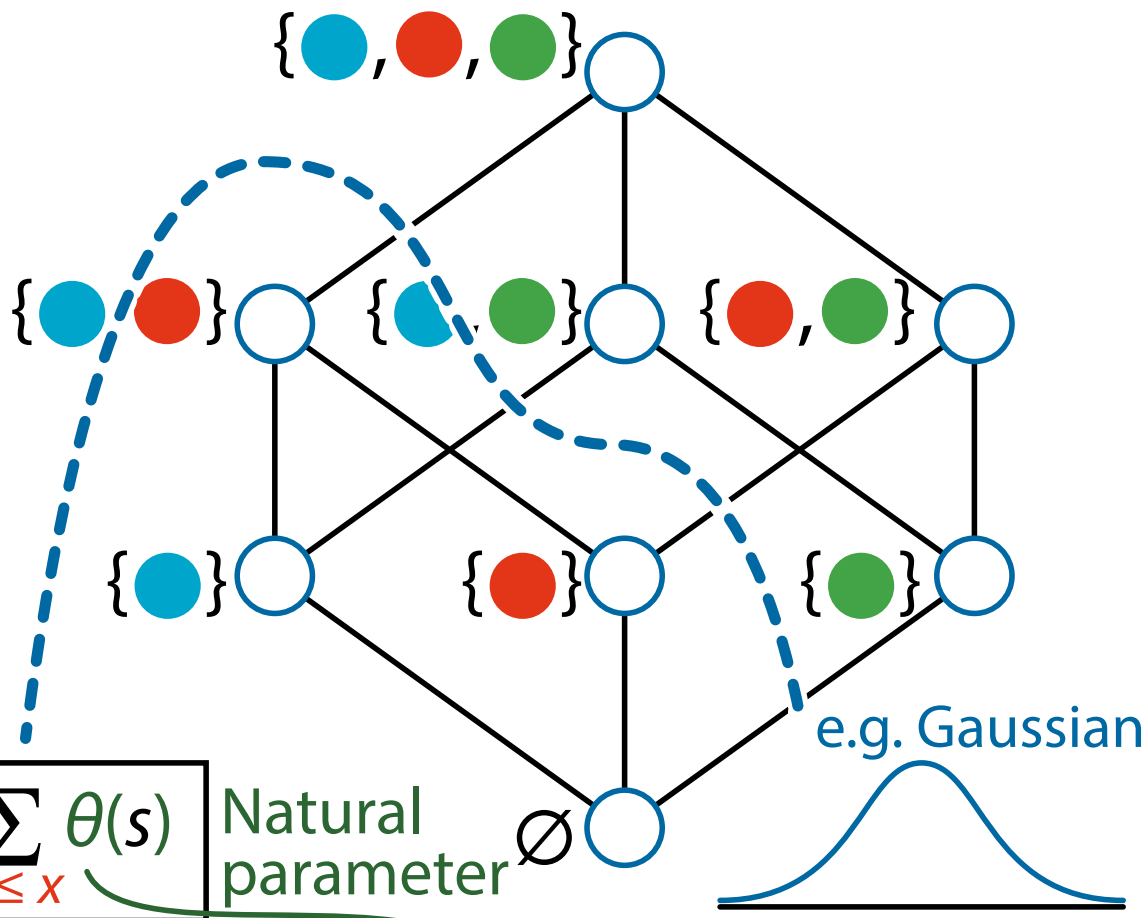


$\eta$ : Frequency  
 $p$ : Probability  
 $\theta$ : Coefficient of  
log-linear model

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

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





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

Natural parameter  $\theta$

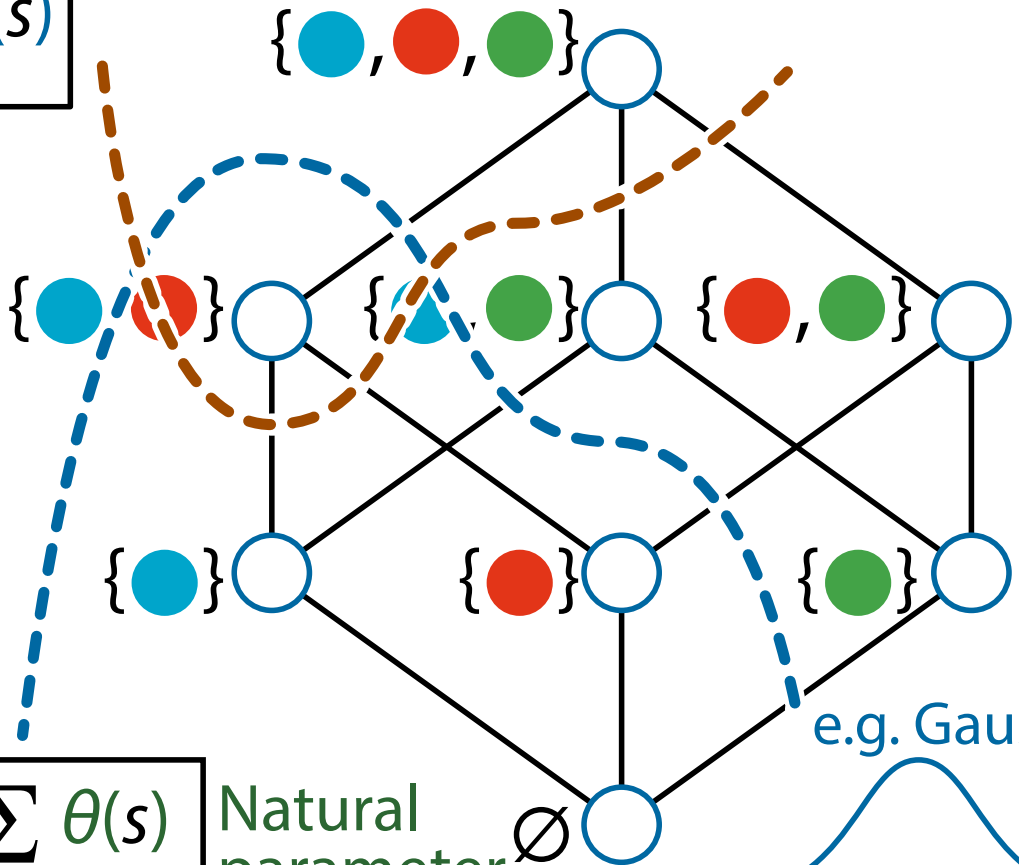
Exponential family:

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

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

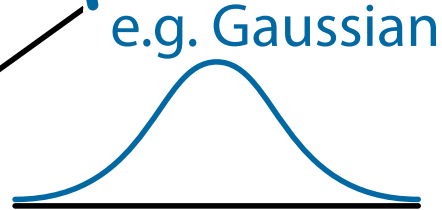
$$\eta(x) = \mathbb{E}[F_x(s)]$$

Sufficient statistics of exponential family



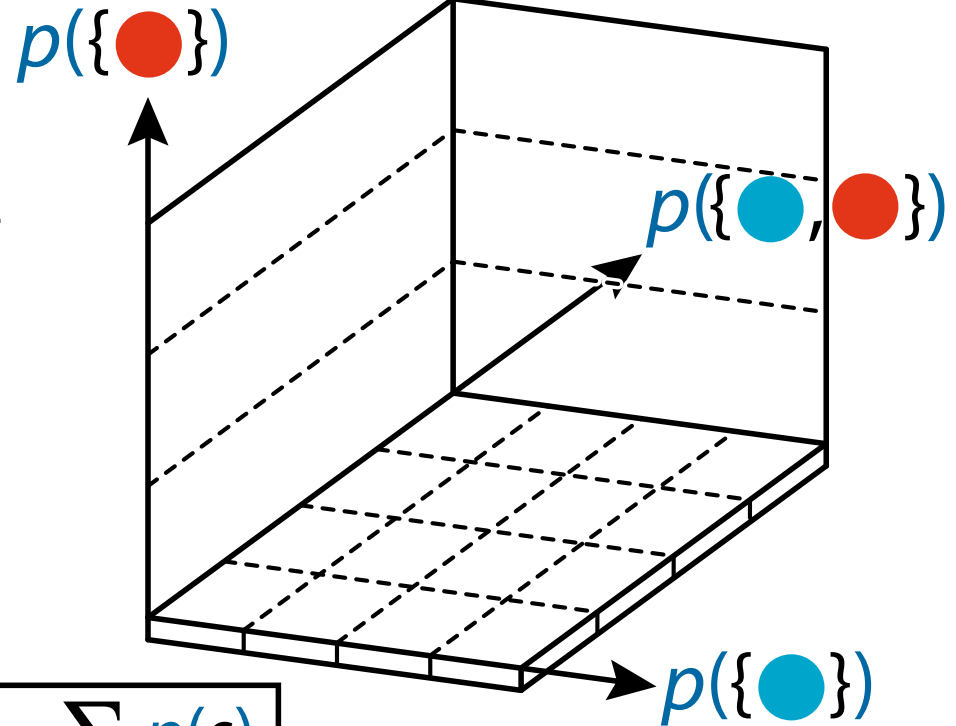
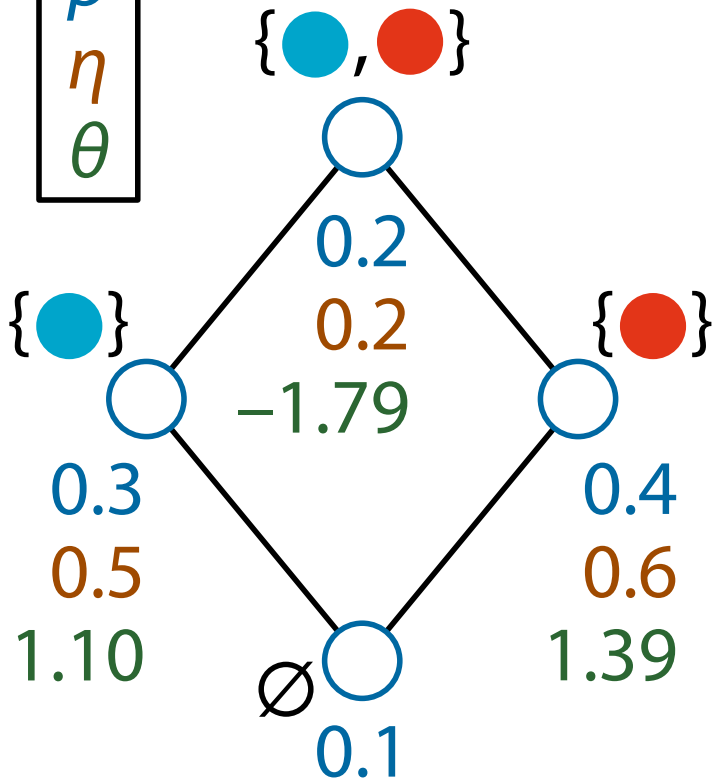
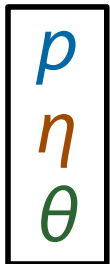
$$\log p(x) = \sum_{s \leq x} \theta(s)$$

Natural parameter  $\theta$



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

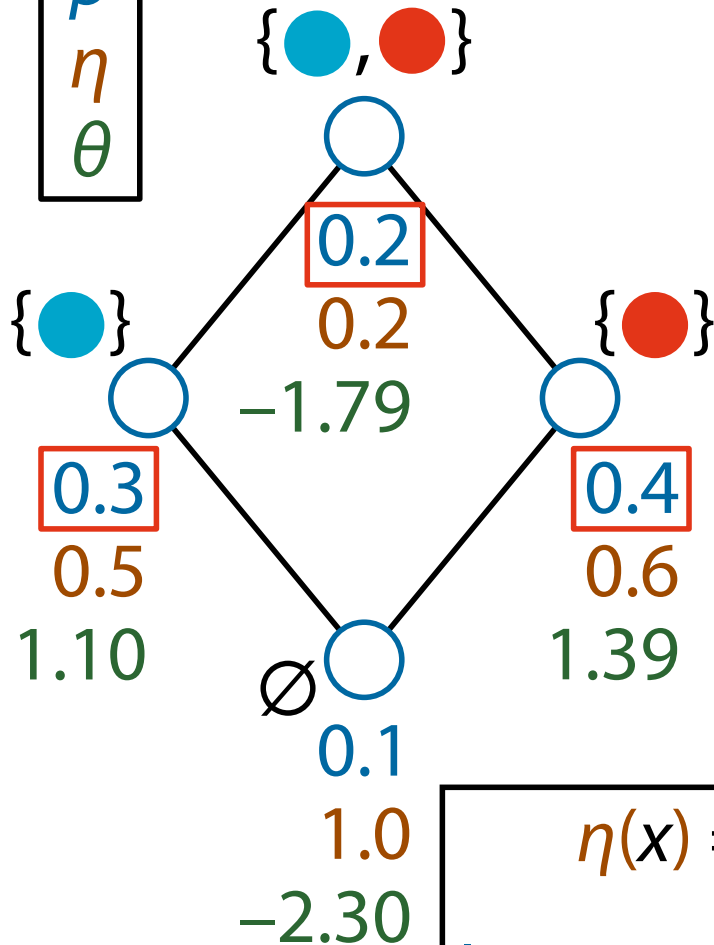
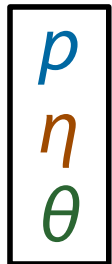
Triple for each node



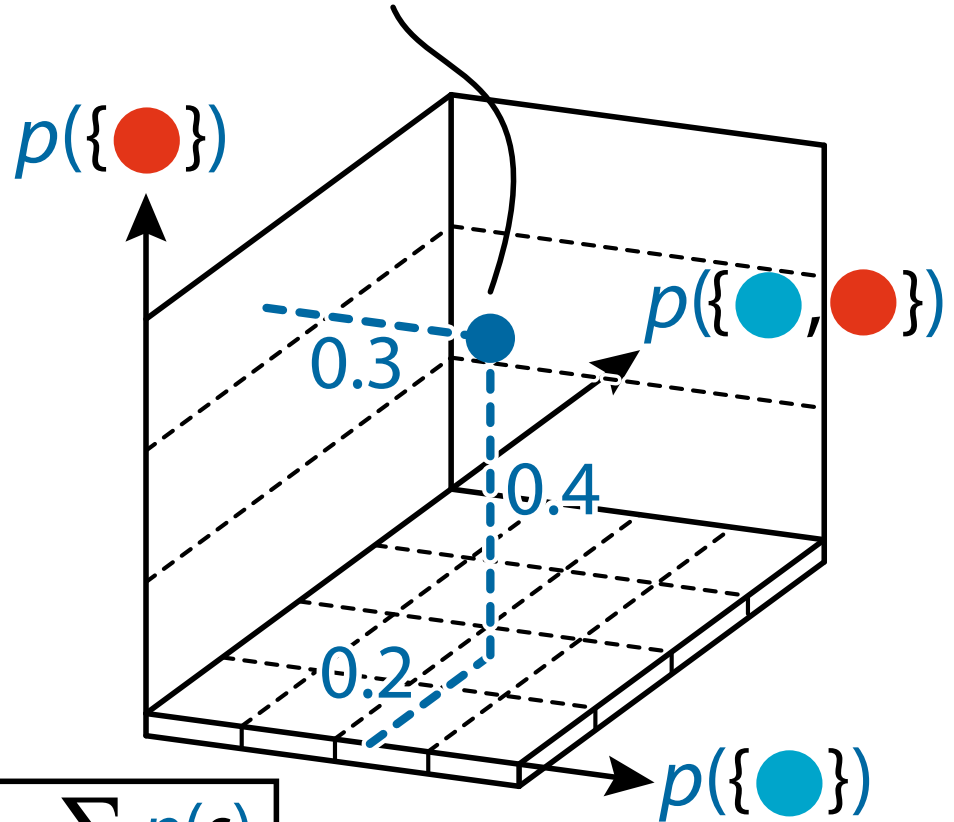
$$\eta(x) = \sum_{s \geq x} p(s)$$

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

Triple for each node



Probability distribution is a "point" in 3D space

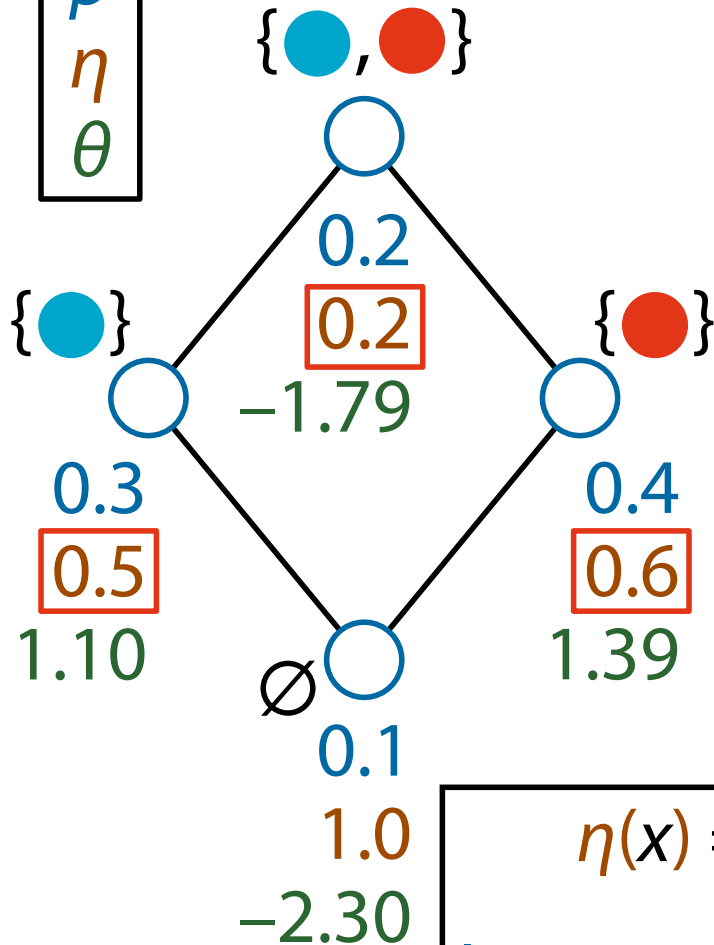
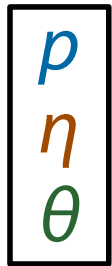


$$\eta(x) = \sum_{s \geq x} p(s)$$

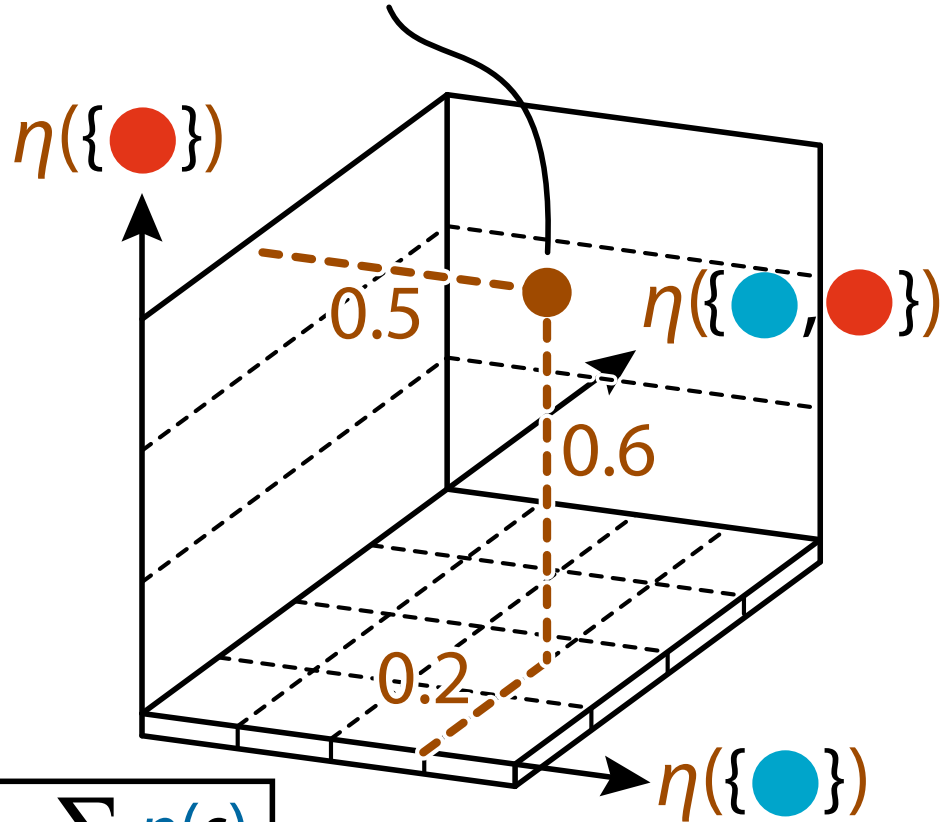
$$\log p(x) = \sum_{s \leq x} \theta(s)$$



Triple for each node



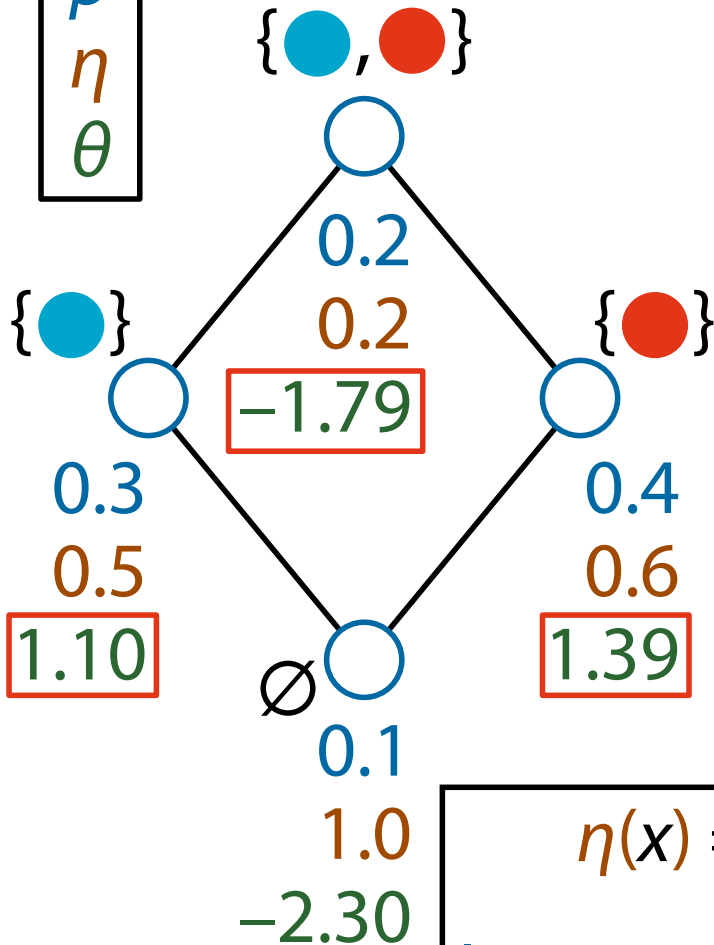
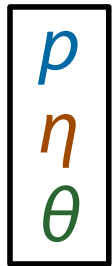
Probability distribution is a "point" in 3D space



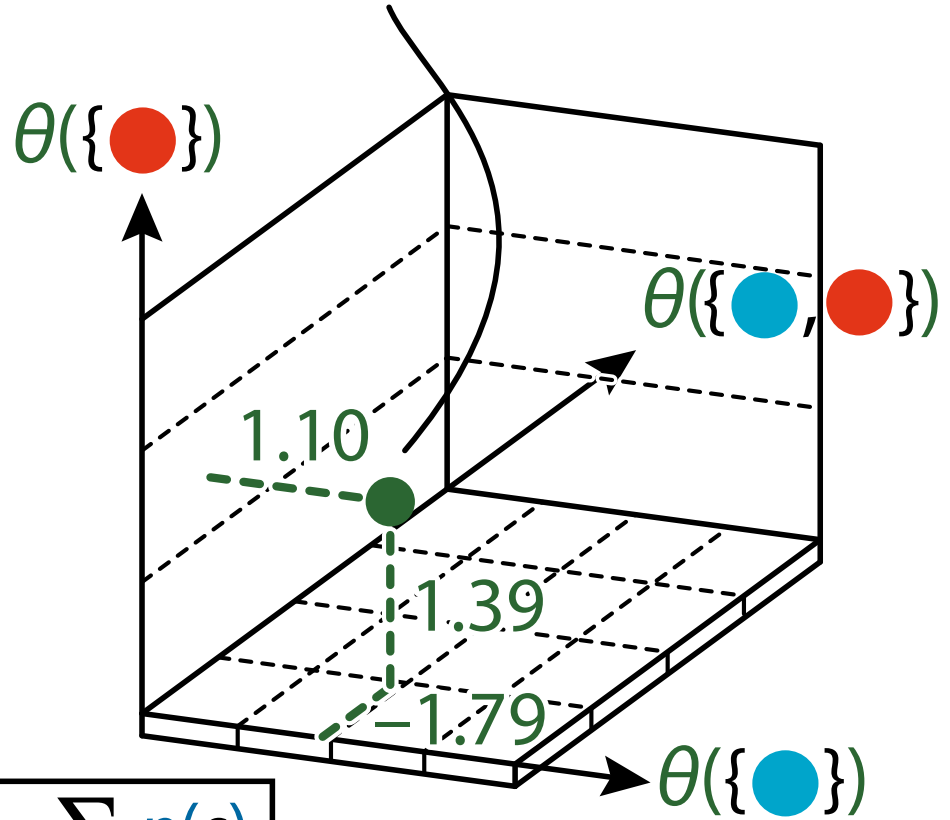
$$\eta(x) = \sum_{s \geq x} p(s)$$

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

Triple for each node



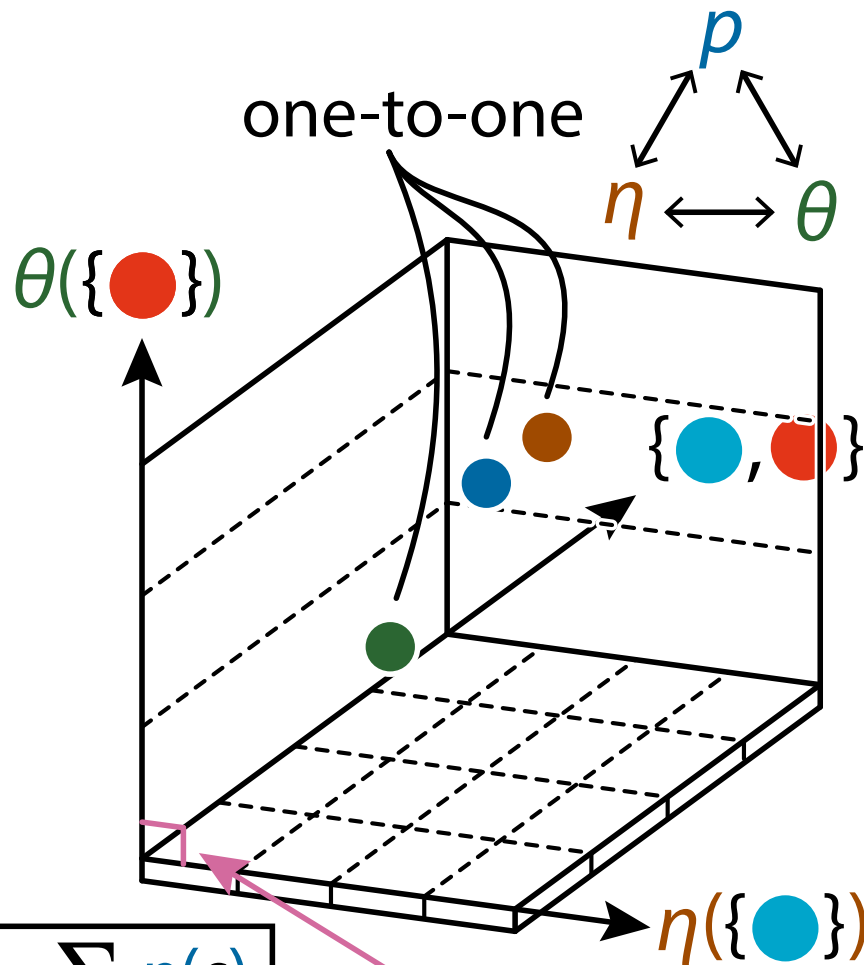
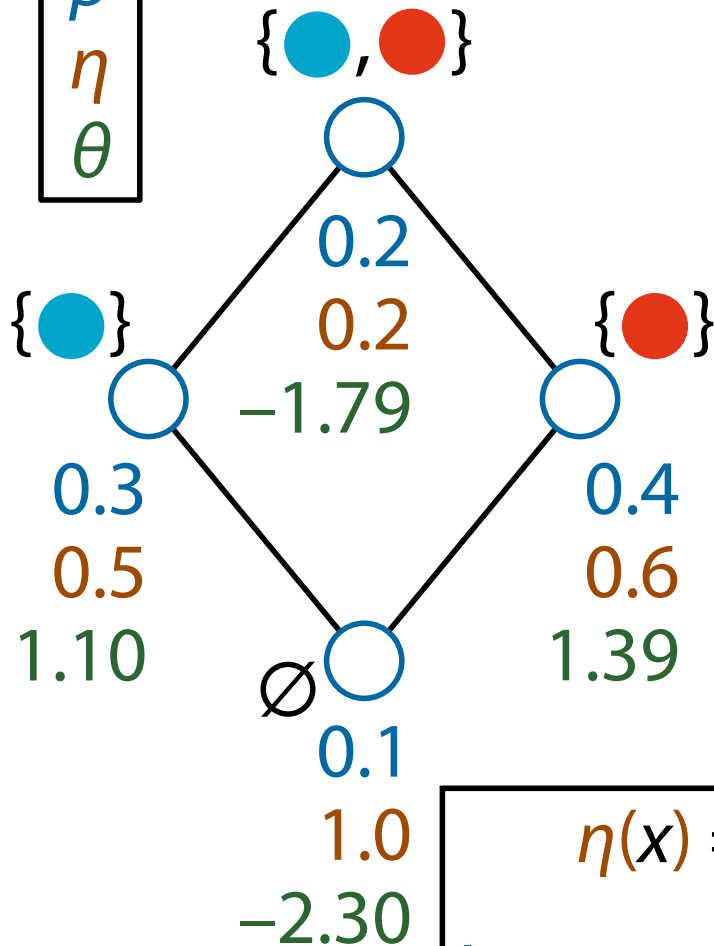
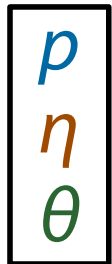
Probability distribution is a "point" in 3D space



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

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

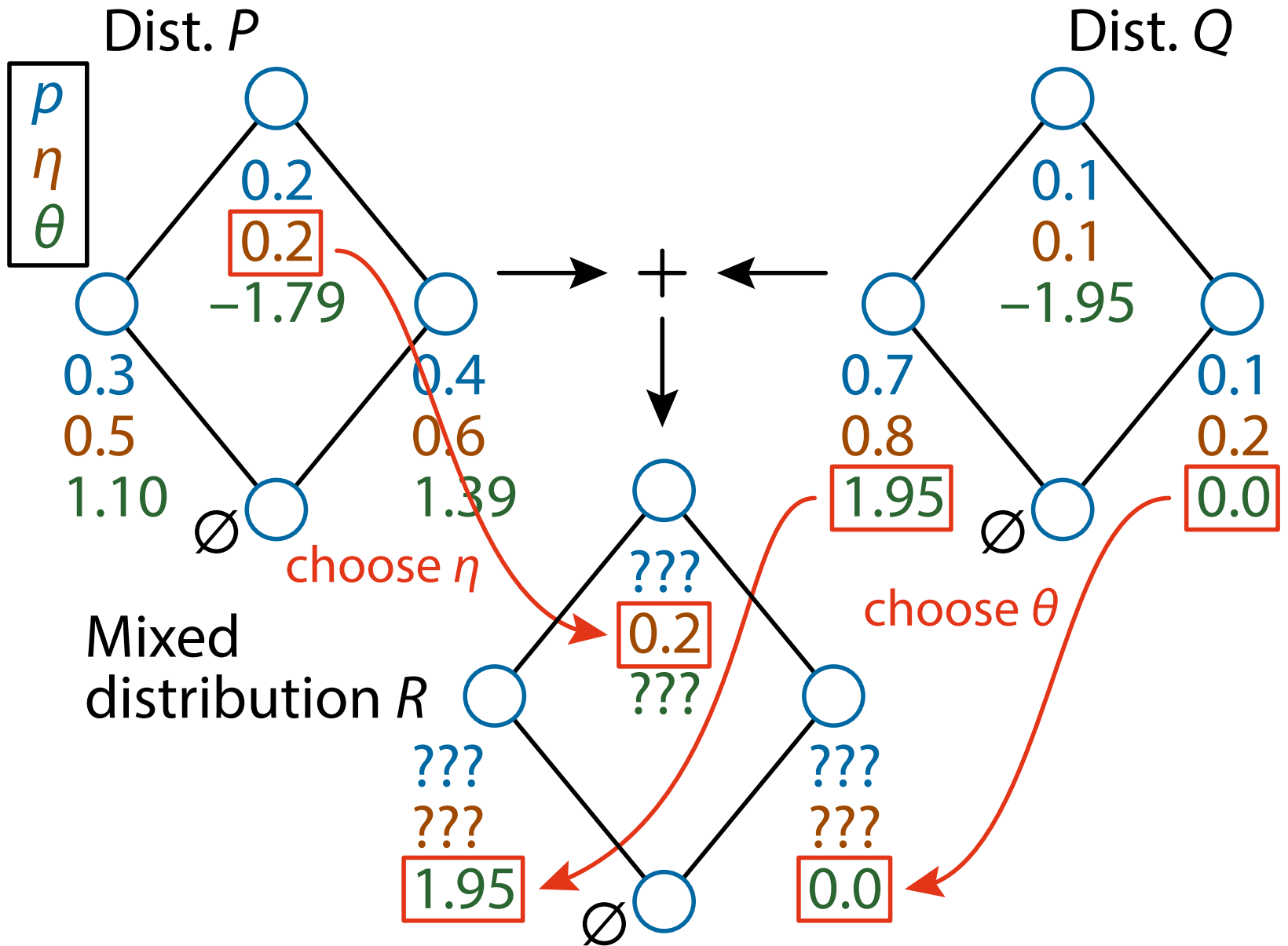
Triple for each node

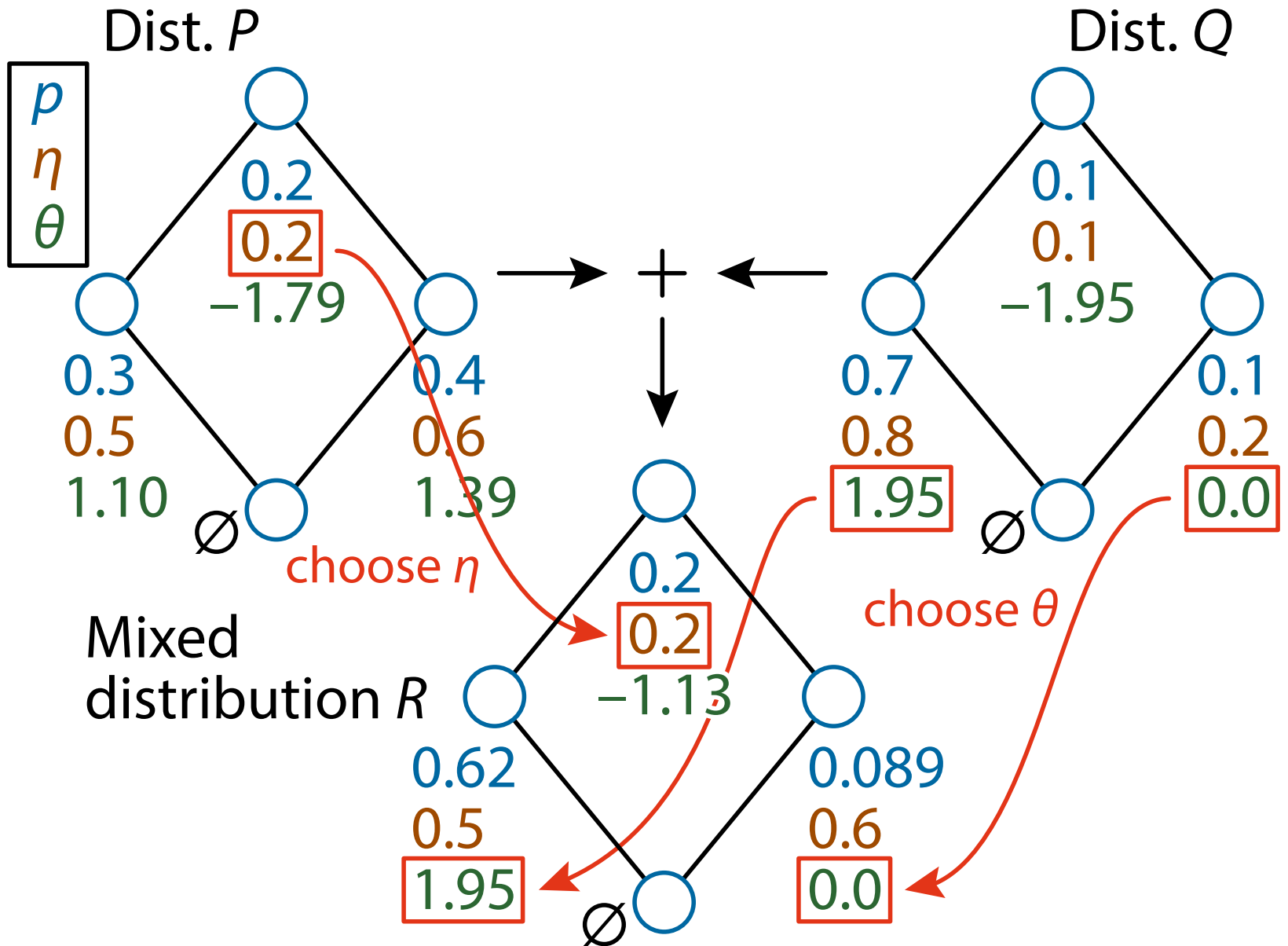


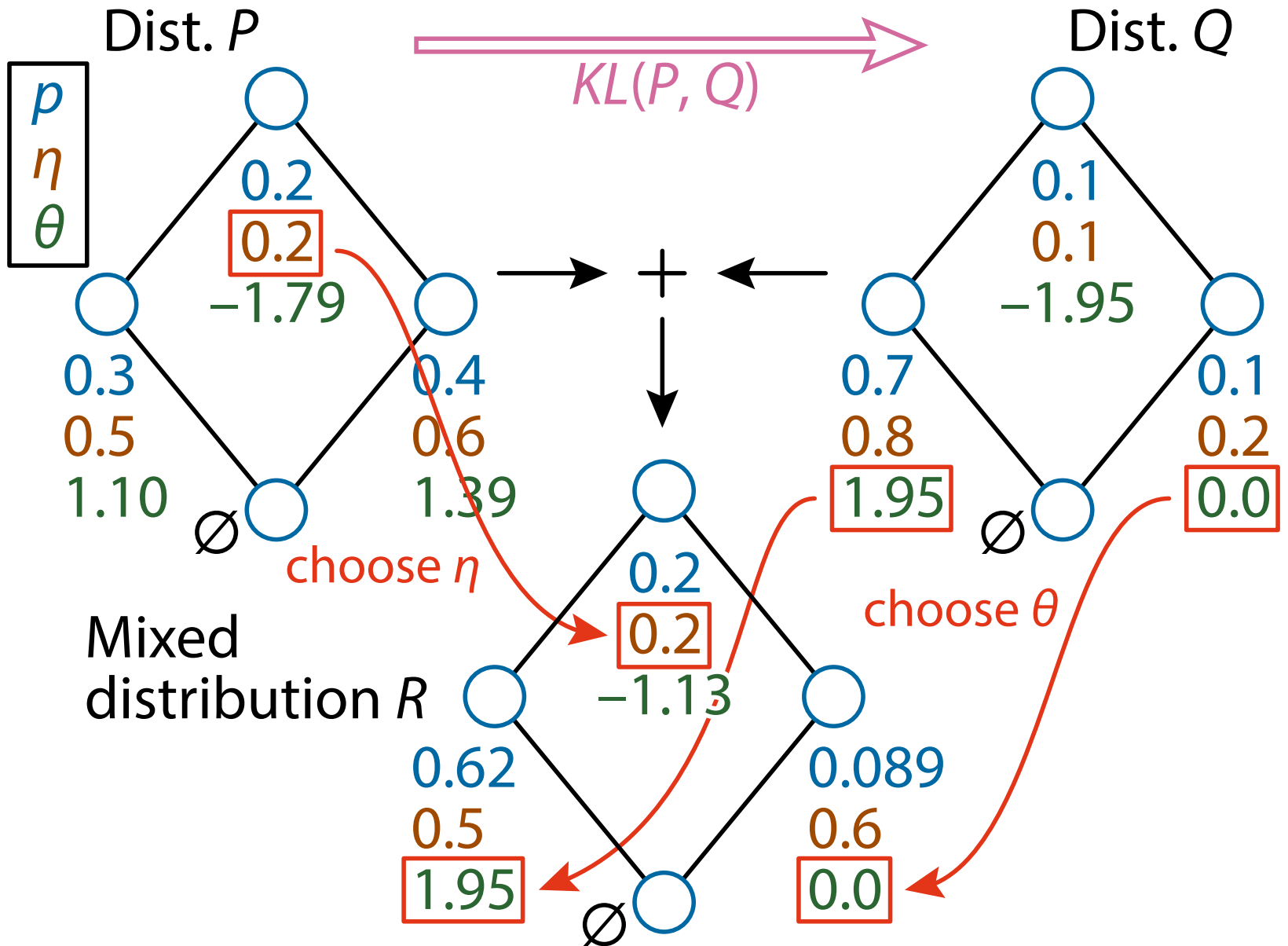
$$\eta(x) = \sum_{s \geq x} p(s)$$

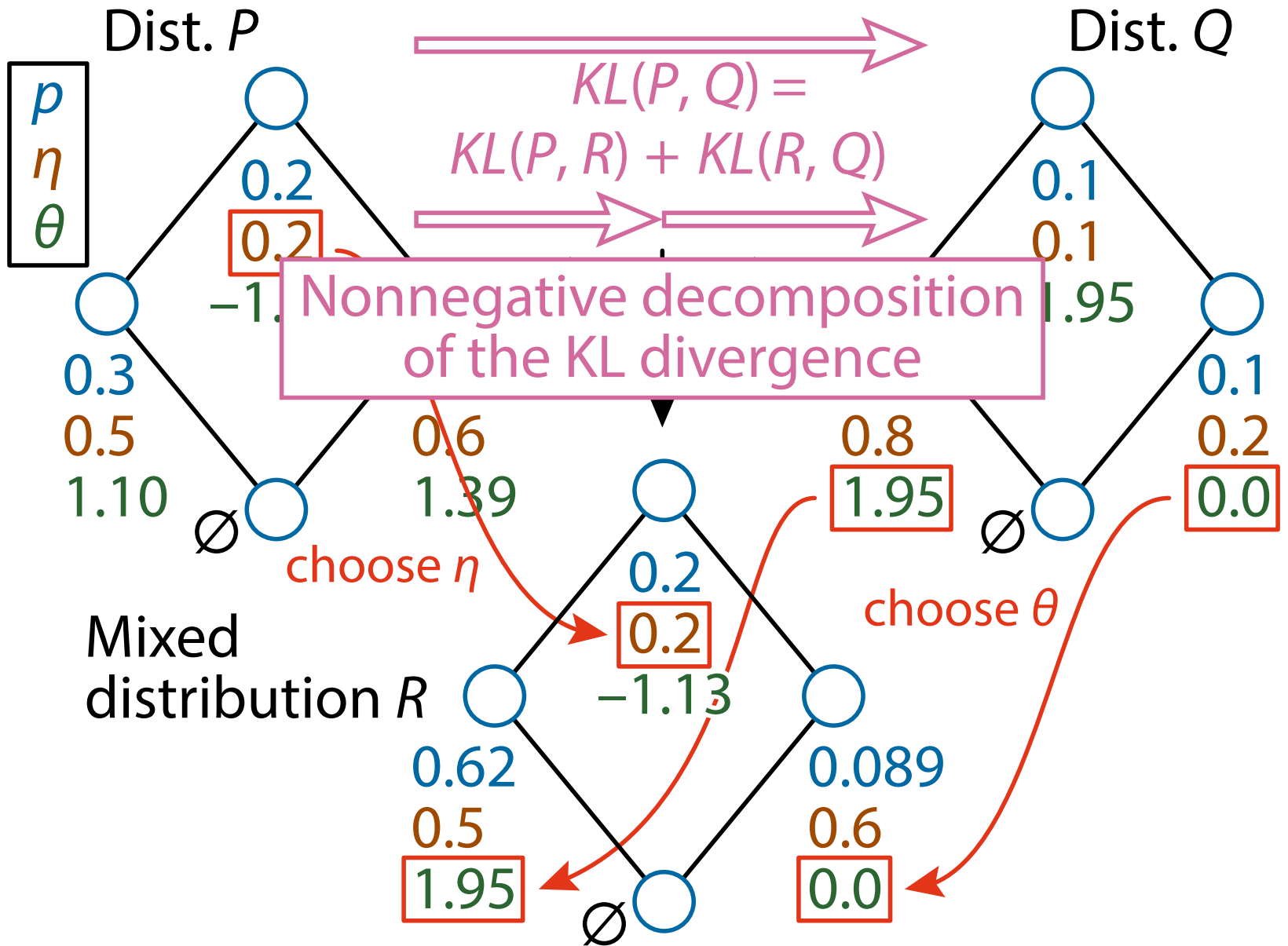
$$\log p(x) = \sum_{s \leq x} \theta(s)$$

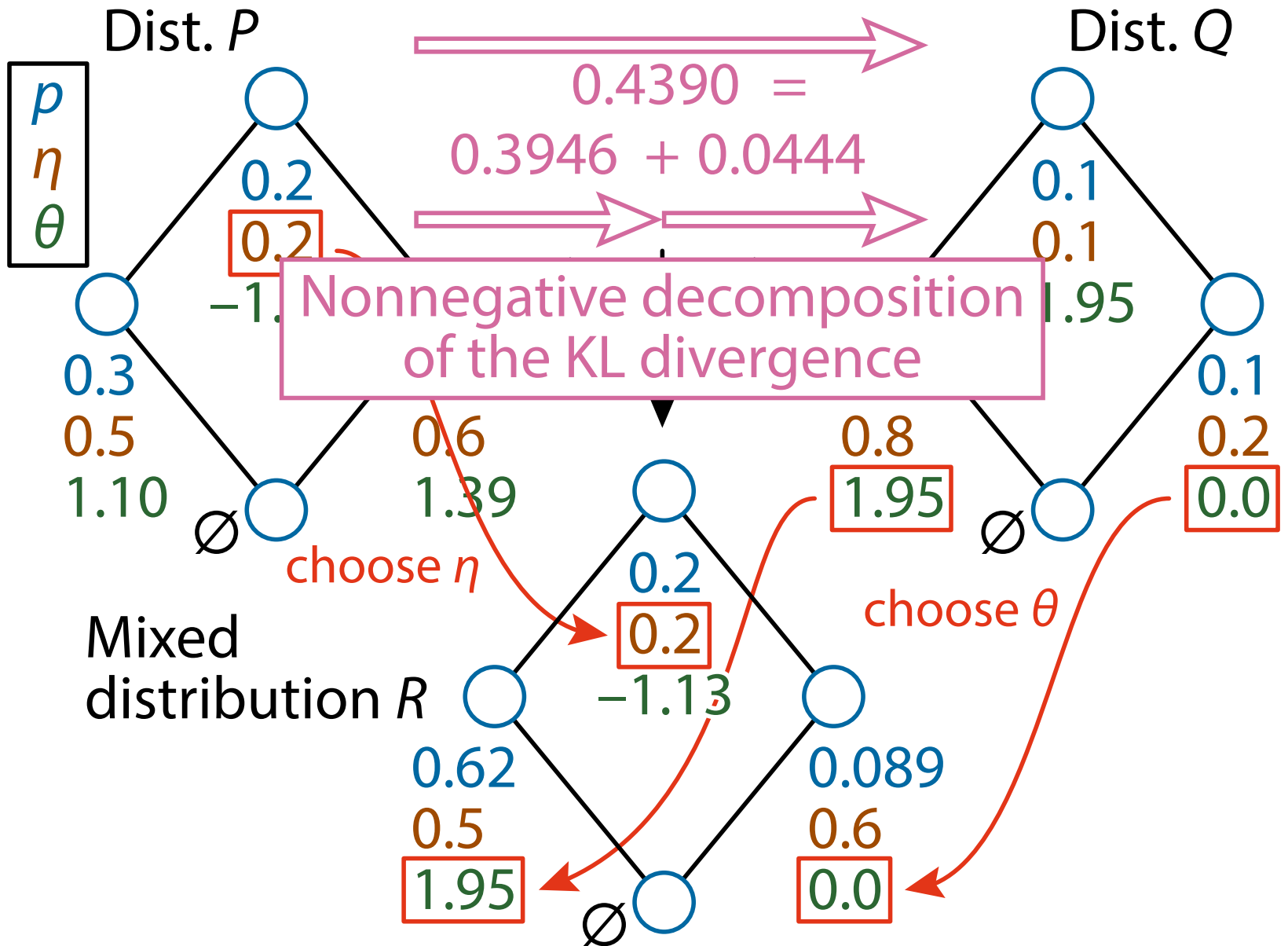
*$\theta$  and  $\eta$  are dually orthogonal*



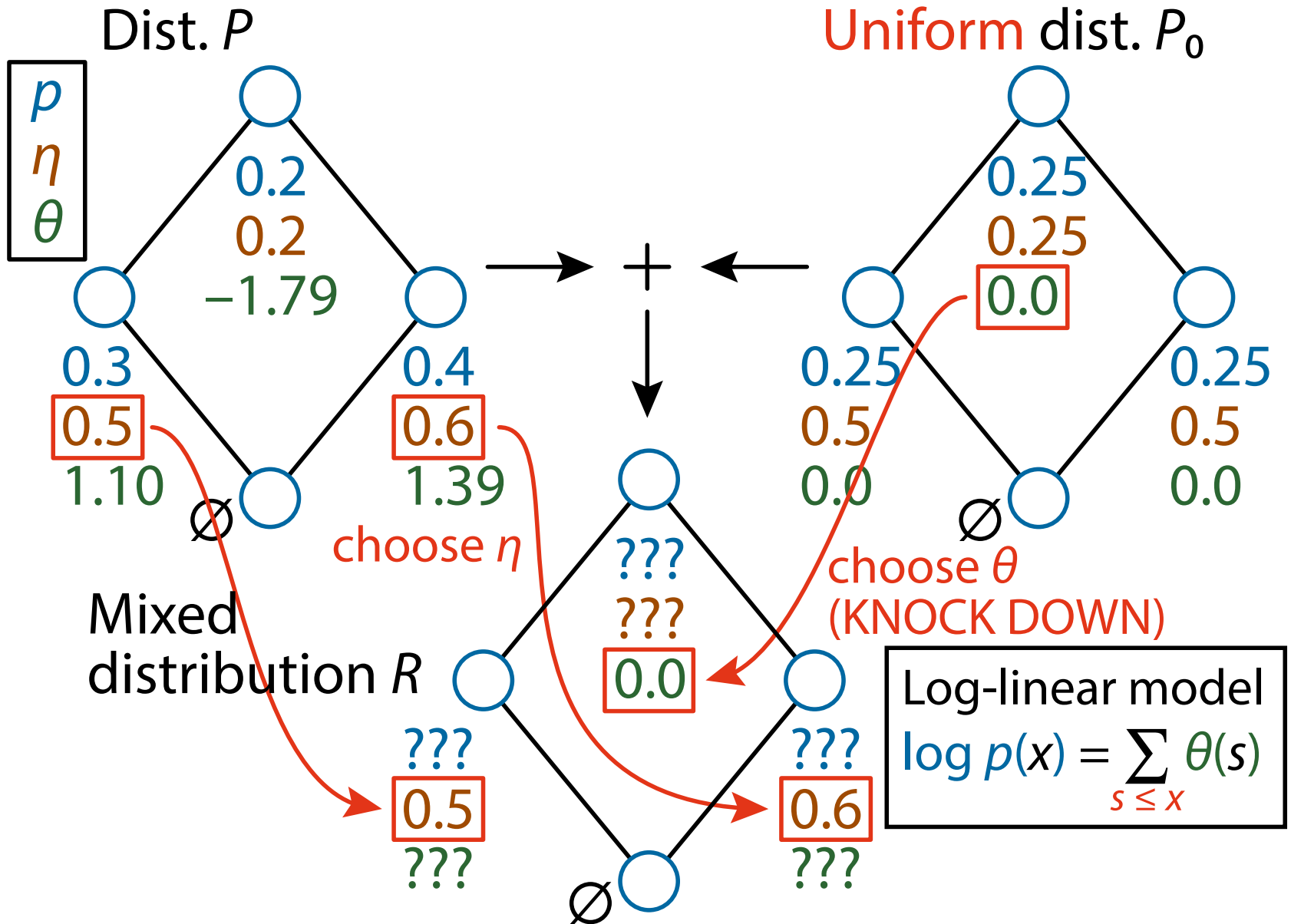


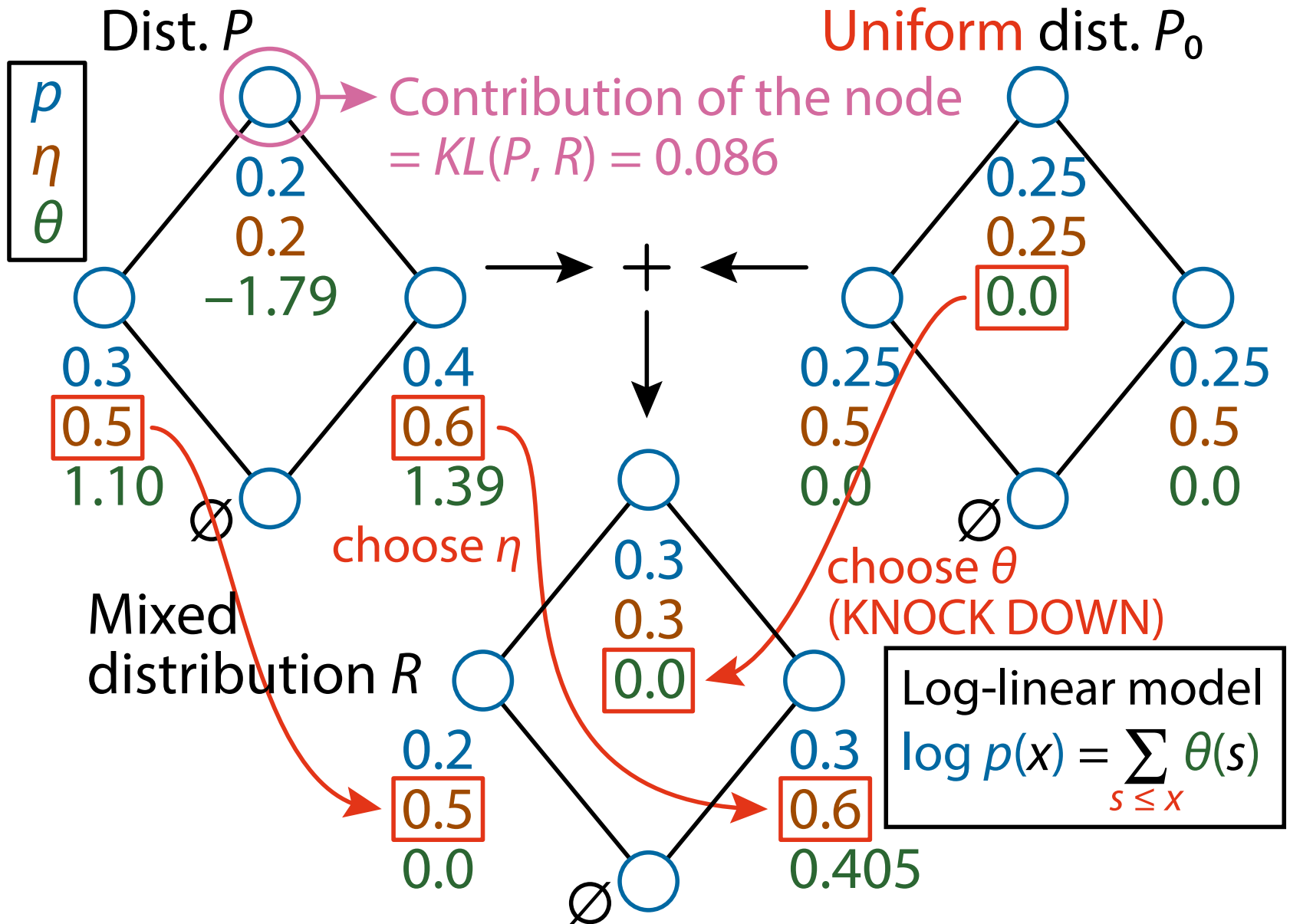


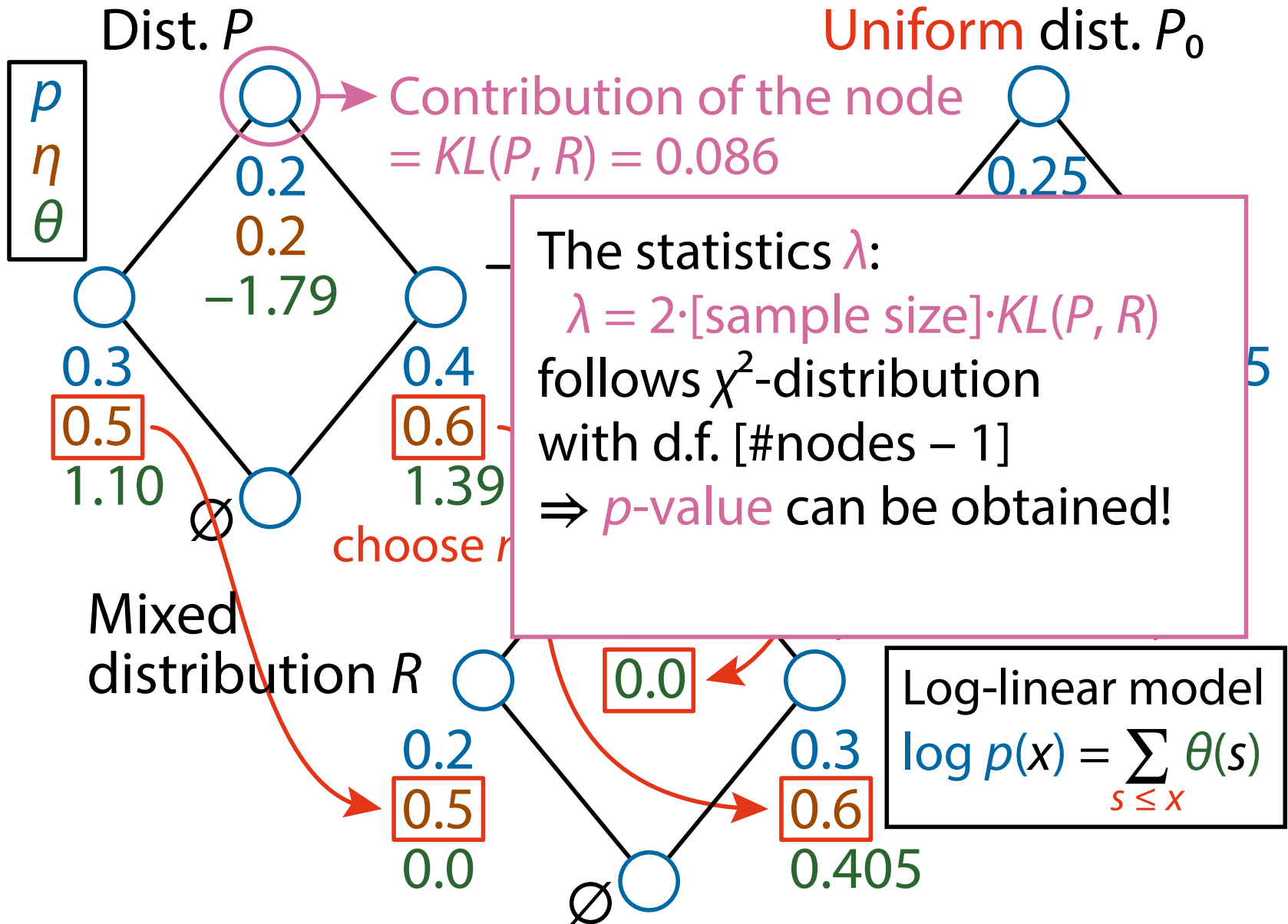


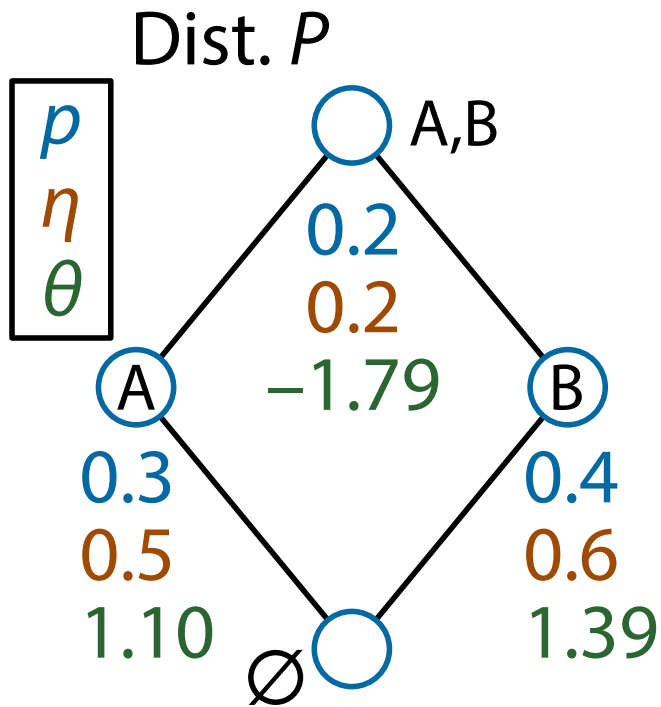




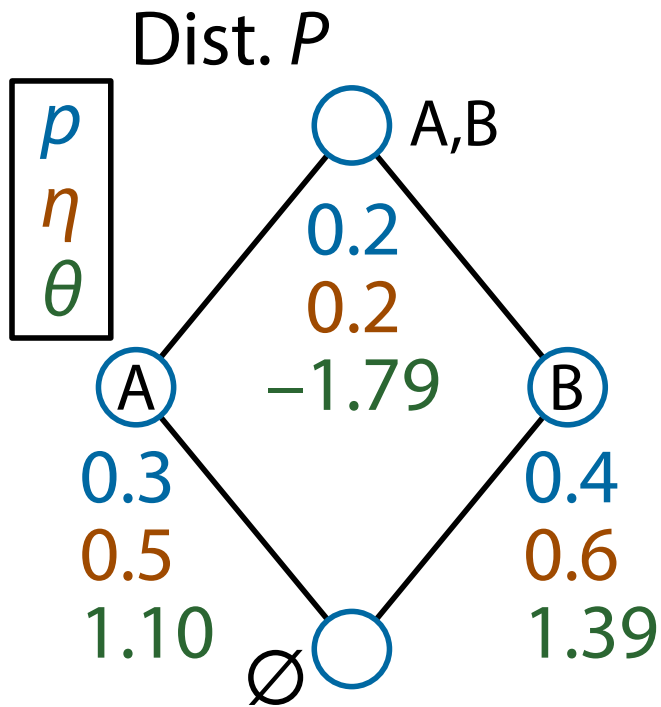








A	B	A, B
0.3	0.4	0.2
0.5	0.6	0.2
1.10	1.39	-1.79
???	???	???
???	0.6	0.2
0.0	???	???
???	???	???
0.5	???	0.2
???	0.0	???
???	???	???
0.5	0.6	???
???	???	0.0



	A	B	A, B
Row 1	0.3	0.4	0.2
Row 2	0.5	0.6	0.2
Row 3	1.10	1.39	-1.79
Row 4	???	???	???
Row 5	???	0.6	0.2
Row 6	0.0	???	???
Row 7	???	???	???
Row 8	0.5	???	0.2
Row 9	???	0.0	???
Row 10	???	???	???
Row 11	0.5	0.6	???
Row 12	???	???	0.0

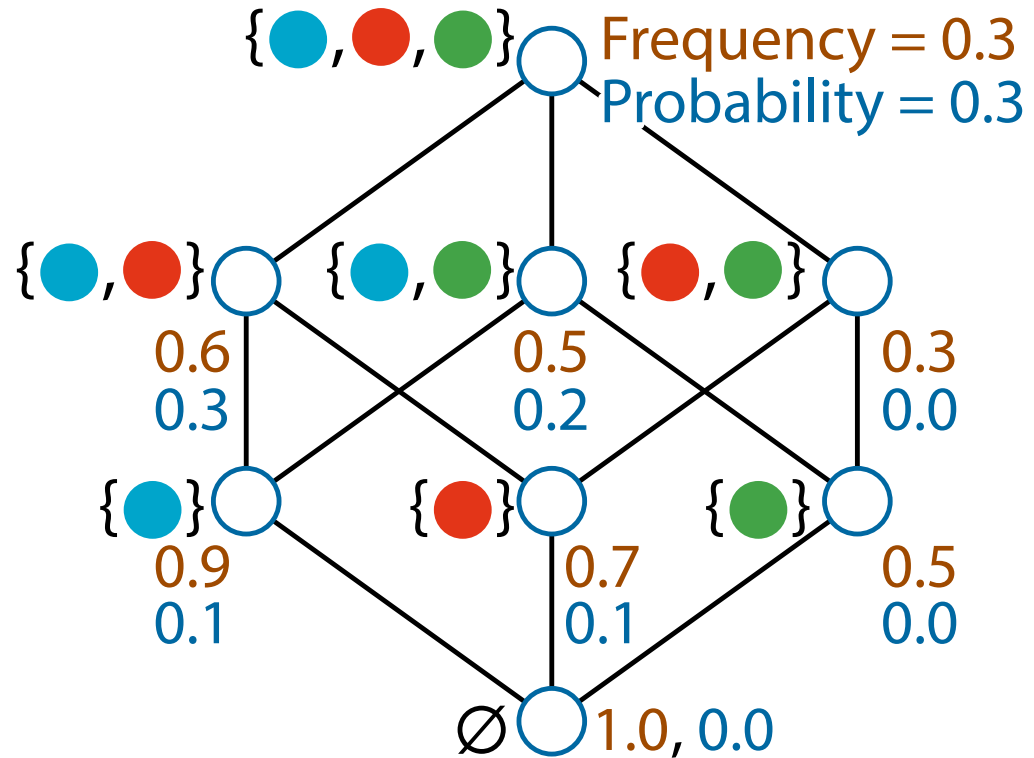
Annotations:

- Row 1: KL = Score of A
- Row 4: KL = Score of B
- Row 7: KL = Score of A, B

# 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

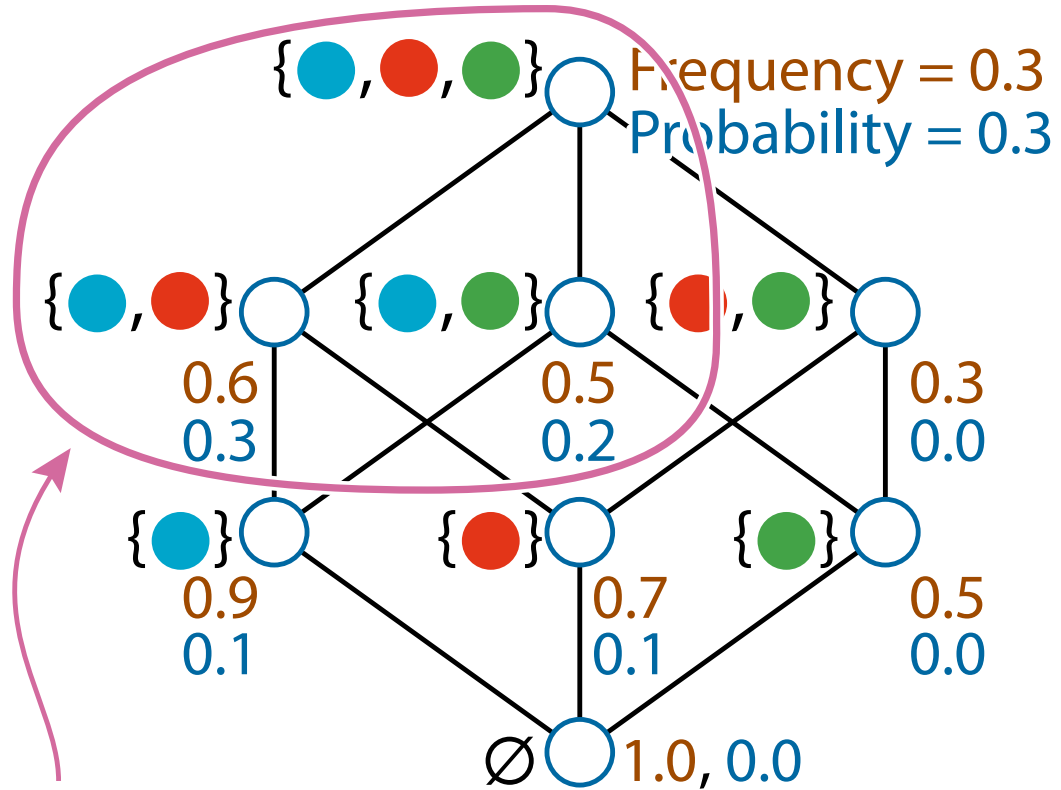


Number of nodes =  $2^{\text{\#features}}$   
 $\Rightarrow$  combinatorial explosion!

# 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

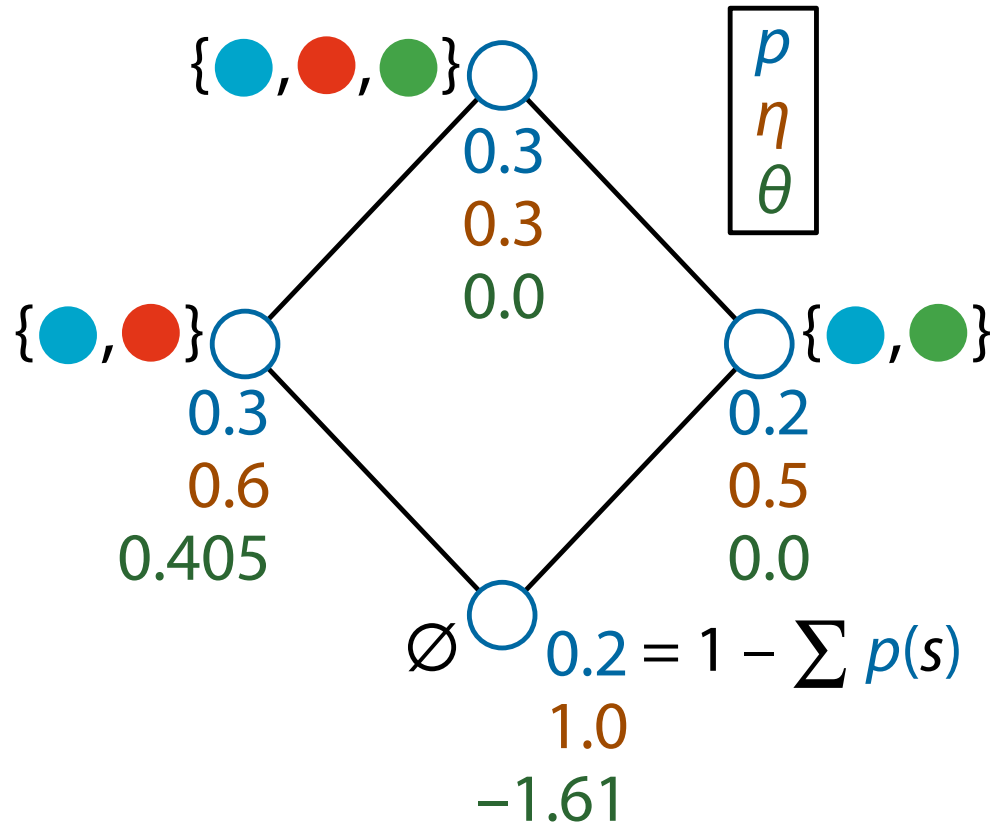


Probability  $\geq 0.2$   
(user specified threshold)

# Remove Nodes with Probability 0

Dataset

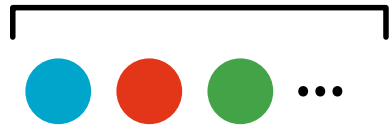
	<span style="color: blue;">●</span>	<span style="color: red;">●</span>	<span style="color: green;">●</span>
ID 1:	1	1	0
ID 2:	1	1	1
ID 3:	1	1	0
ID 4:	1	1	1
ID 5:	1	1	0
ID 6:	1	0	1
ID 7:	1	0	1
ID 8:	1	1	1
ID 9:	1	0	0
ID10:	0	1	0





# Example on Real Data (kosarak)

# features: 41,270

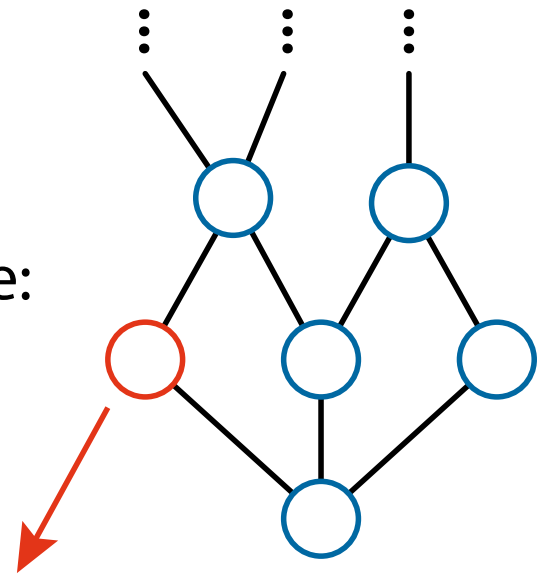


ID 1:	1	1	0
ID 2:	1	1	1
ID 3:	1	1	0 ...
ID 4:	1	1	1
ID 5:	1	1	0
⋮	⋮		

Sample size:  
990,002

Total runtime:  
4.95 seconds

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



# significant interactions: **583**

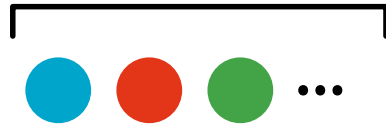
Single feature: 537

Pairwise interactions: 41

Triple interactions: 5

# Example on Real Data (accidents)

# features: 468

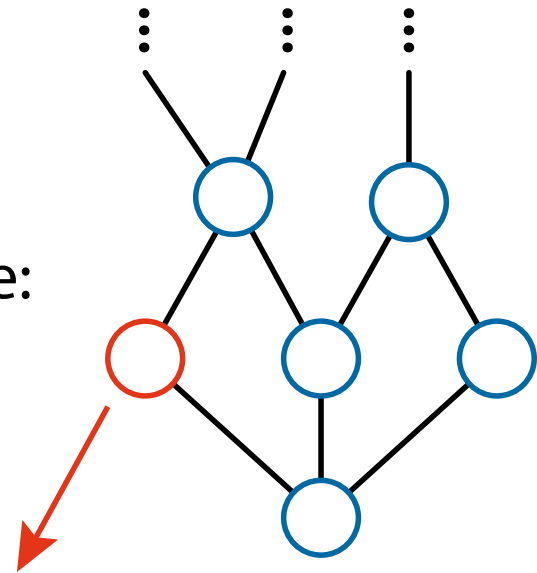


ID 1:	1	1	0
ID 2:	1	1	1
ID 3:	1	1	0 ...
ID 4:	1	1	1
ID 5:	1	1	0
⋮	⋮	⋮	

Total runtime:  
4.95 seconds

Sample size:  
340,183

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



# 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**
    - M. Sugiyama, H. Nakahara, K. Tsuda, *Information Decomposition on Structured Space*, arXiv:1601.05533 (2016)
    - S. Amari, *Information geometry on hierarchy of probability distributions*, IEEE Trans. Info. Theory (2001)
- We can decompose a probability distribution and assess the significance of any-order interactions beyond pairwise interactions