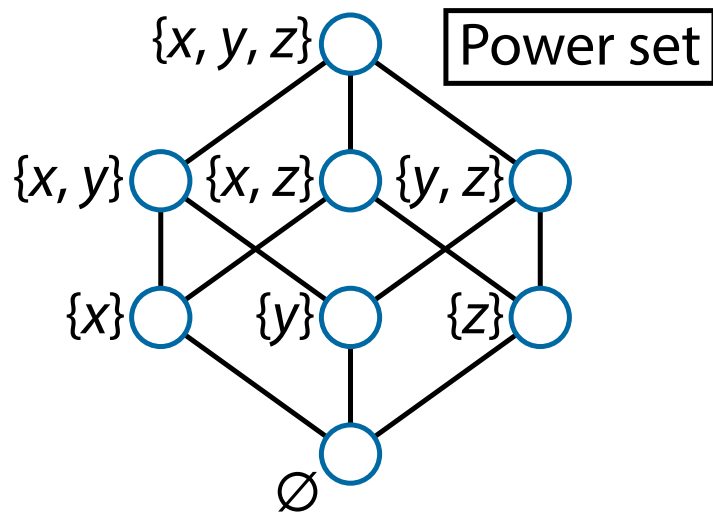


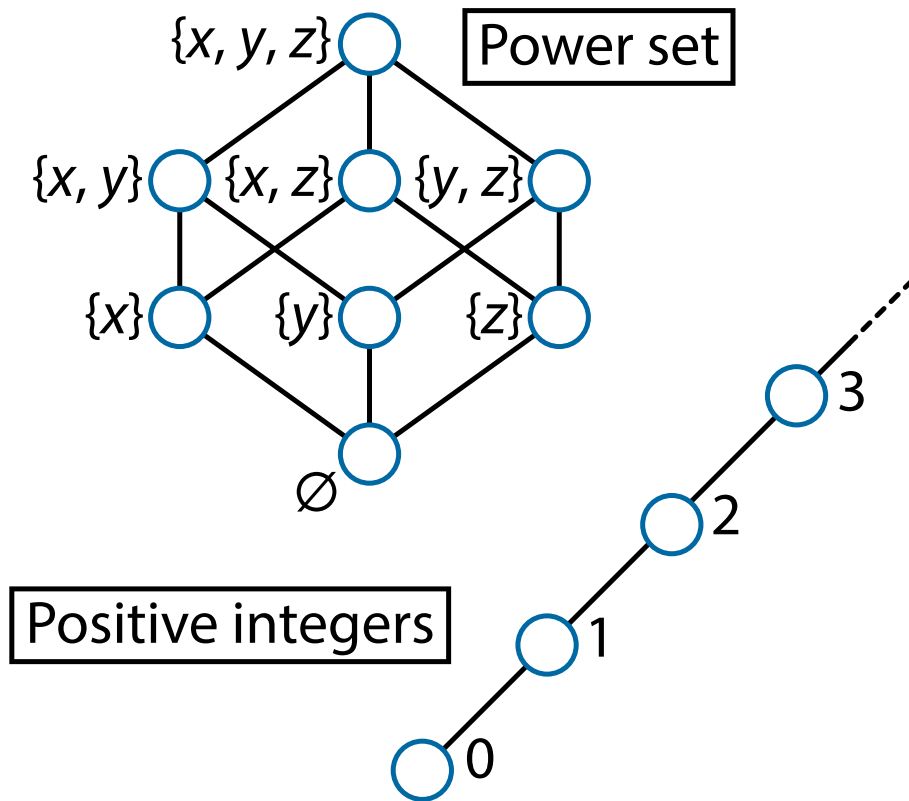
July 11, 2016
ISIT 2016

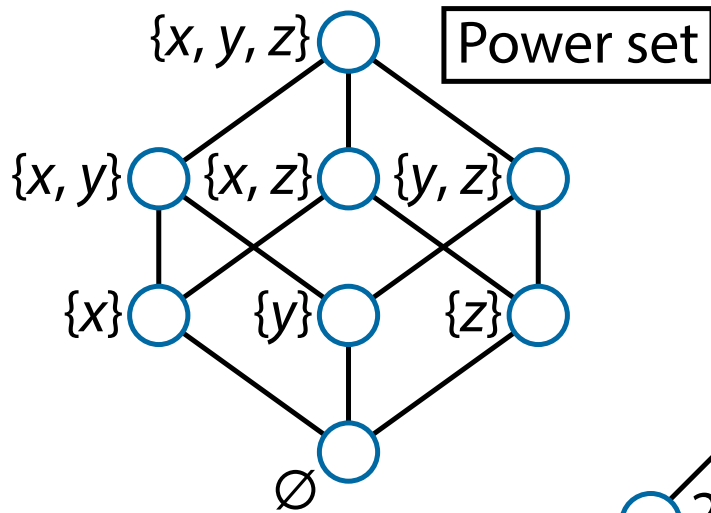


Information Decomposition on Structured Space

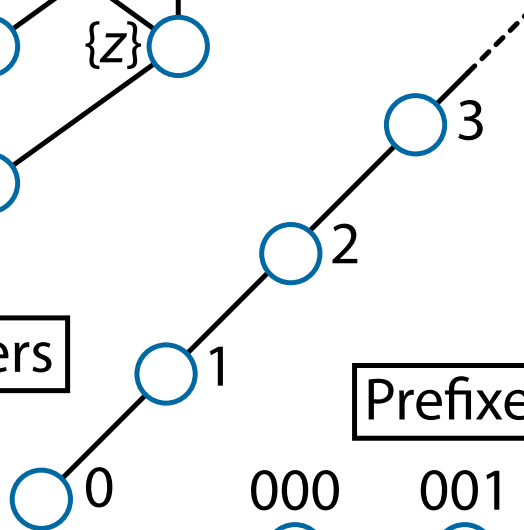
Mahito Sugiyama (Osaka Univ.)
Hiroyuki Nakahara (RIKEN), Koji Tsuda (UTokyo)



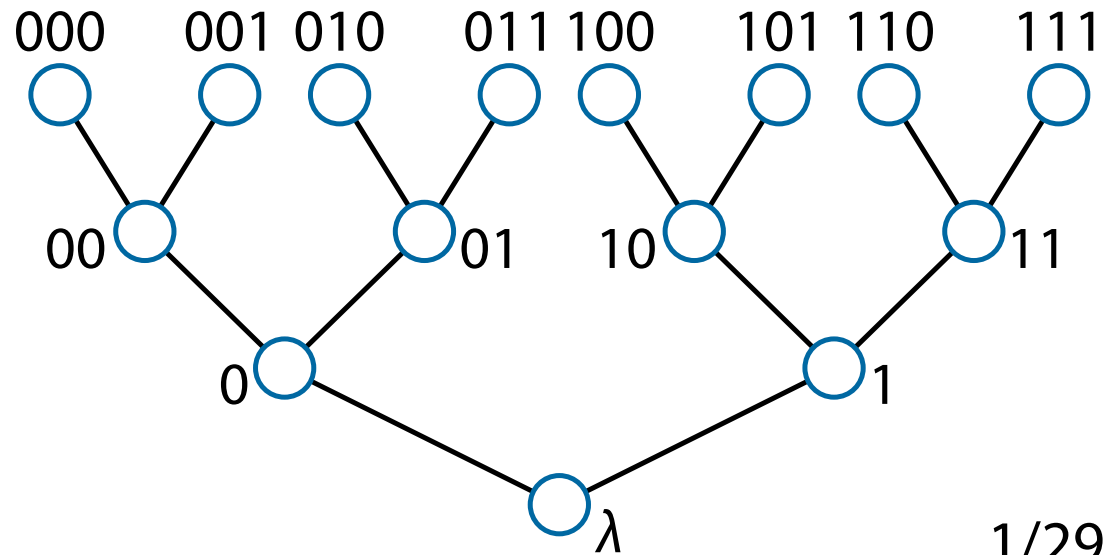


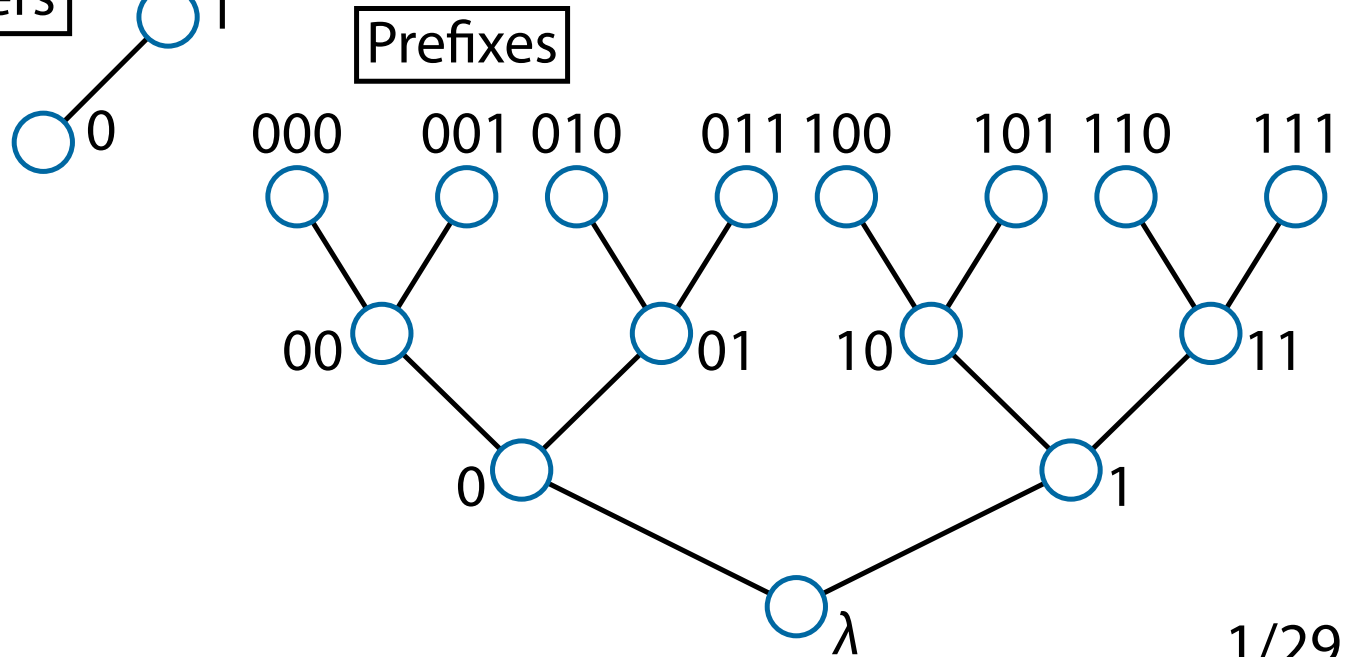
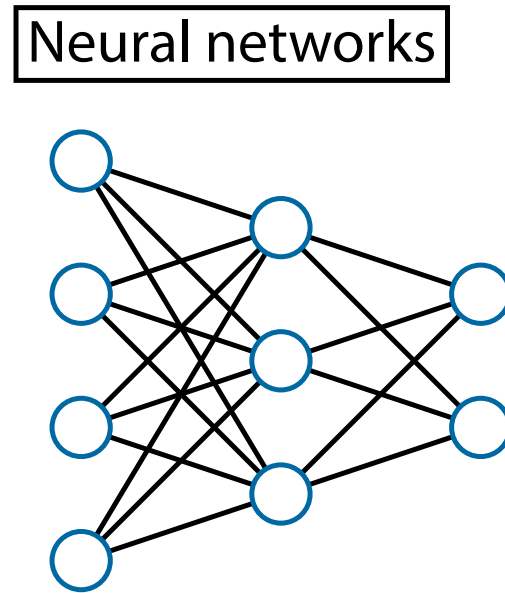
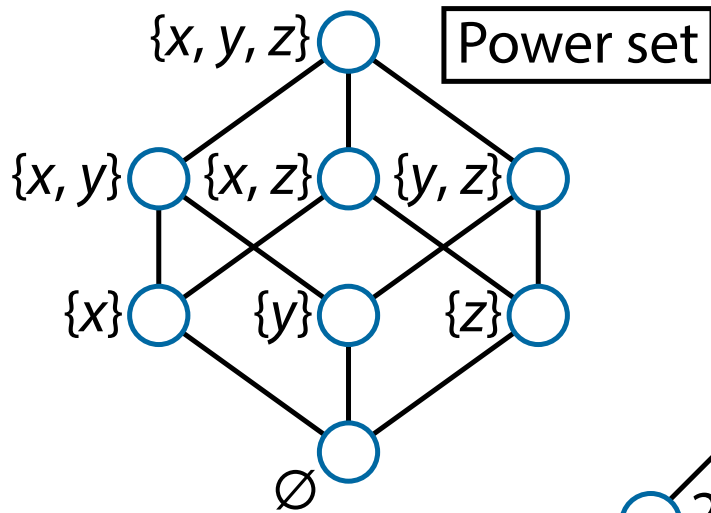


Positive integers



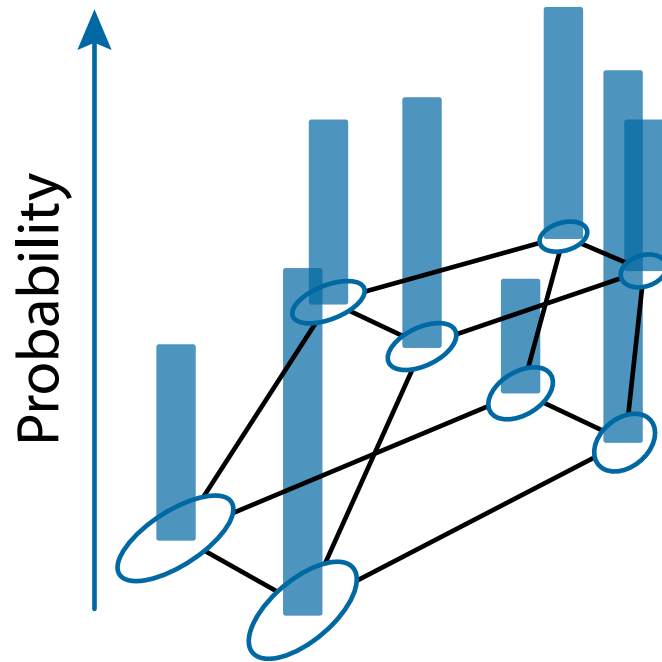
Prefixes





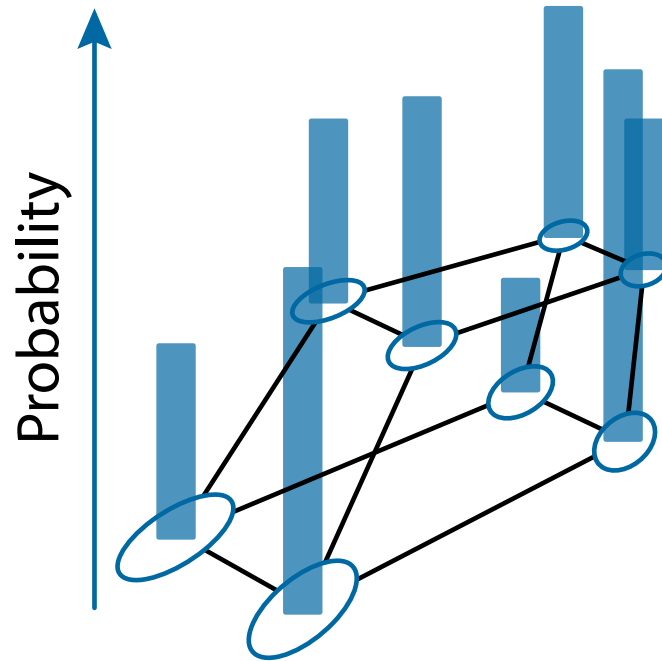
Summary

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

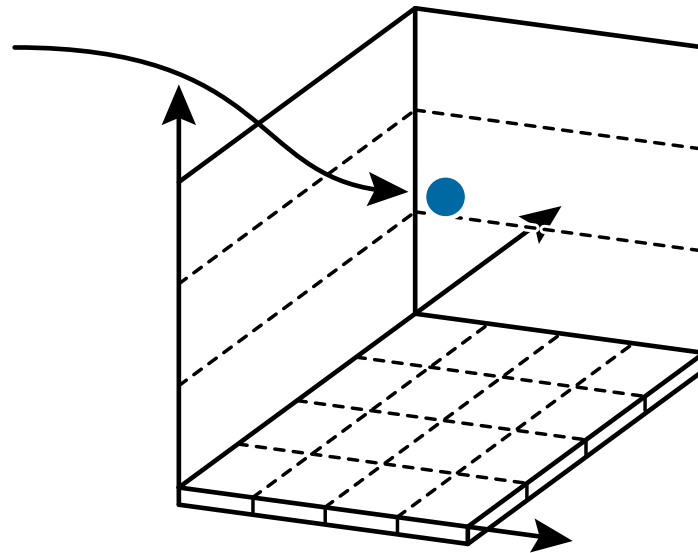


Summary

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

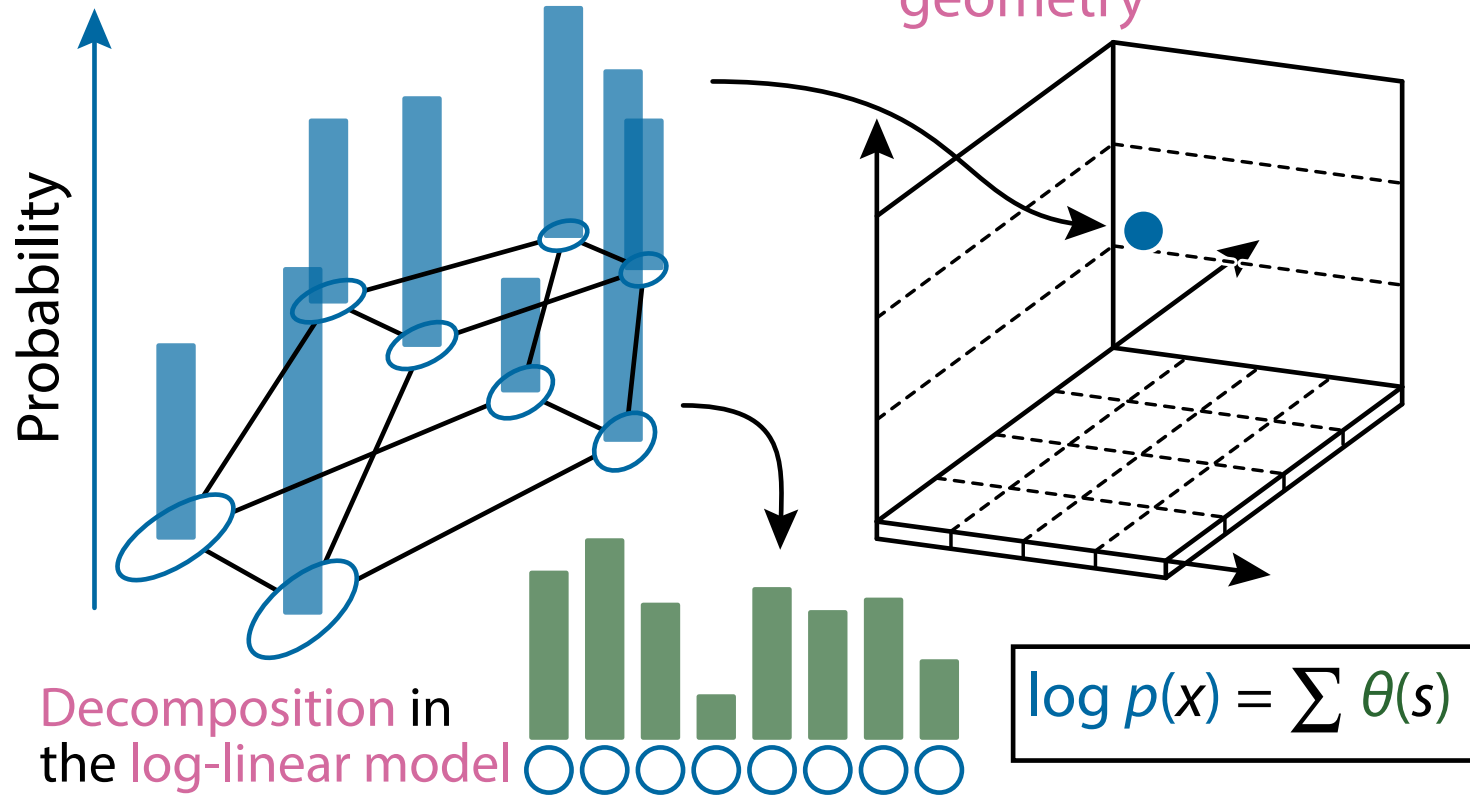


**Information
geometry**



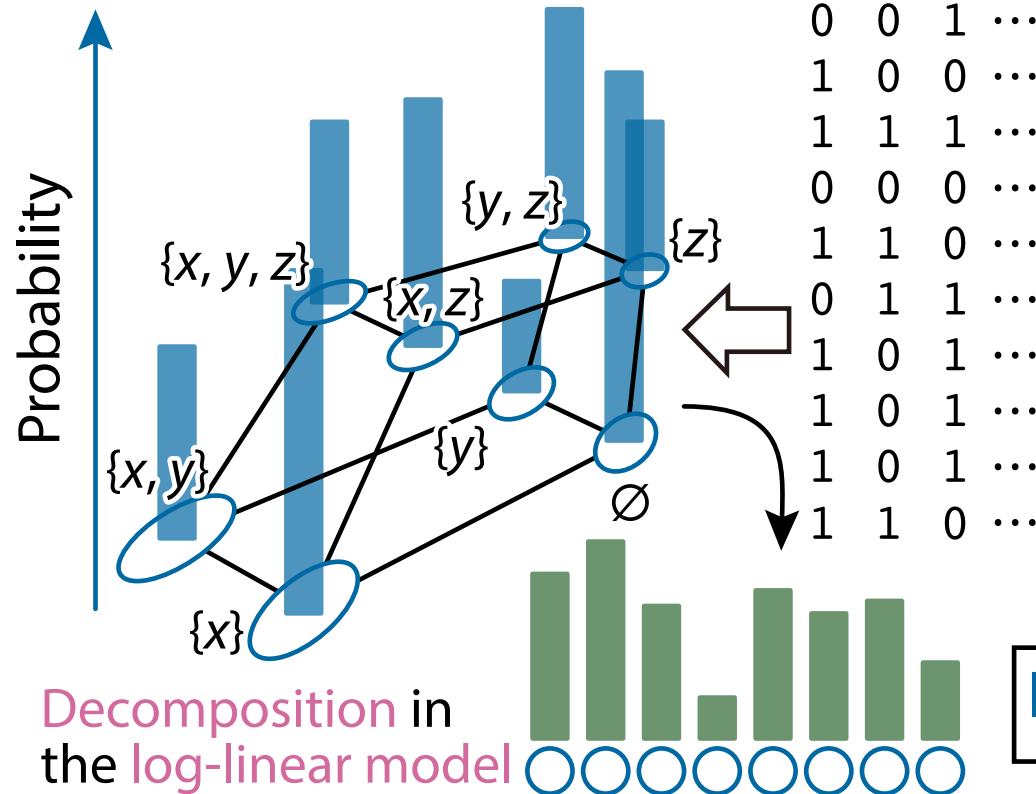
Summary

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



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	0	1	...
1	1	0	...

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

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

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

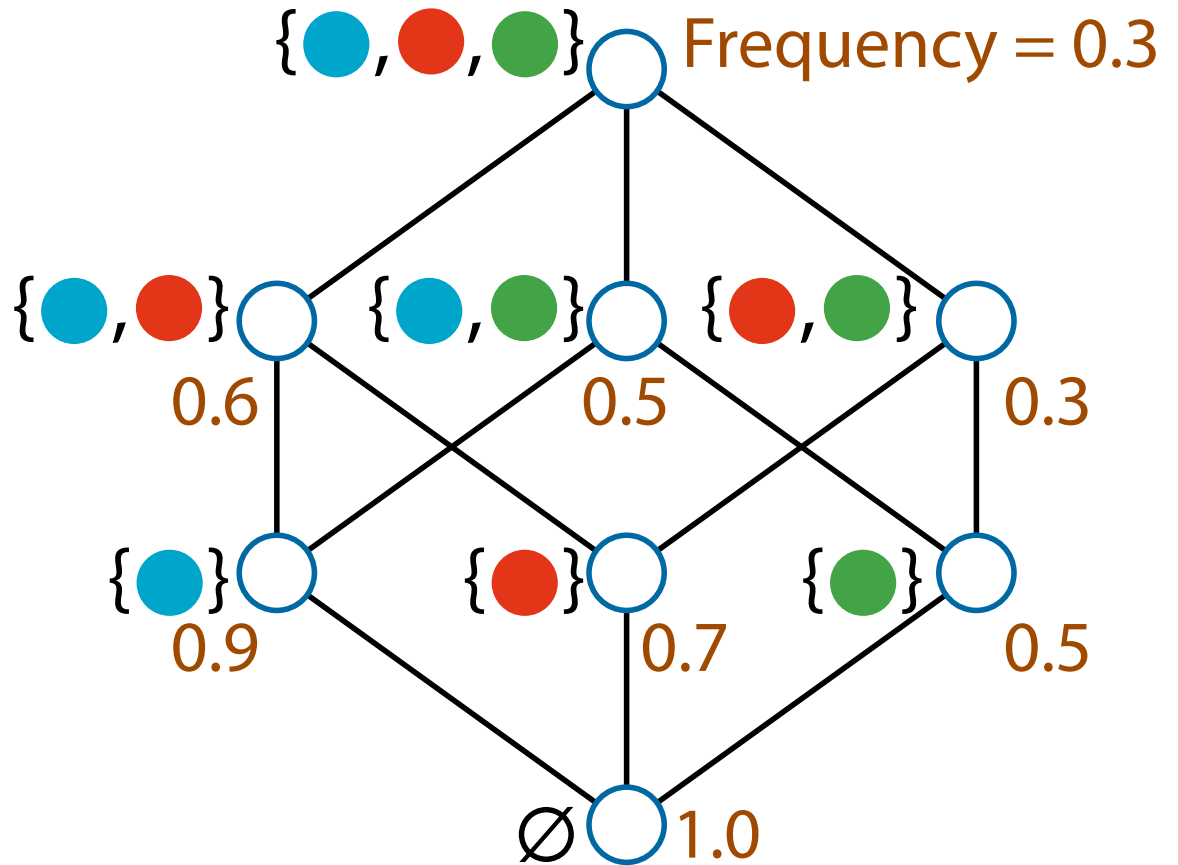
ID 9: 1 0 0

ID10: 0 1 0




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

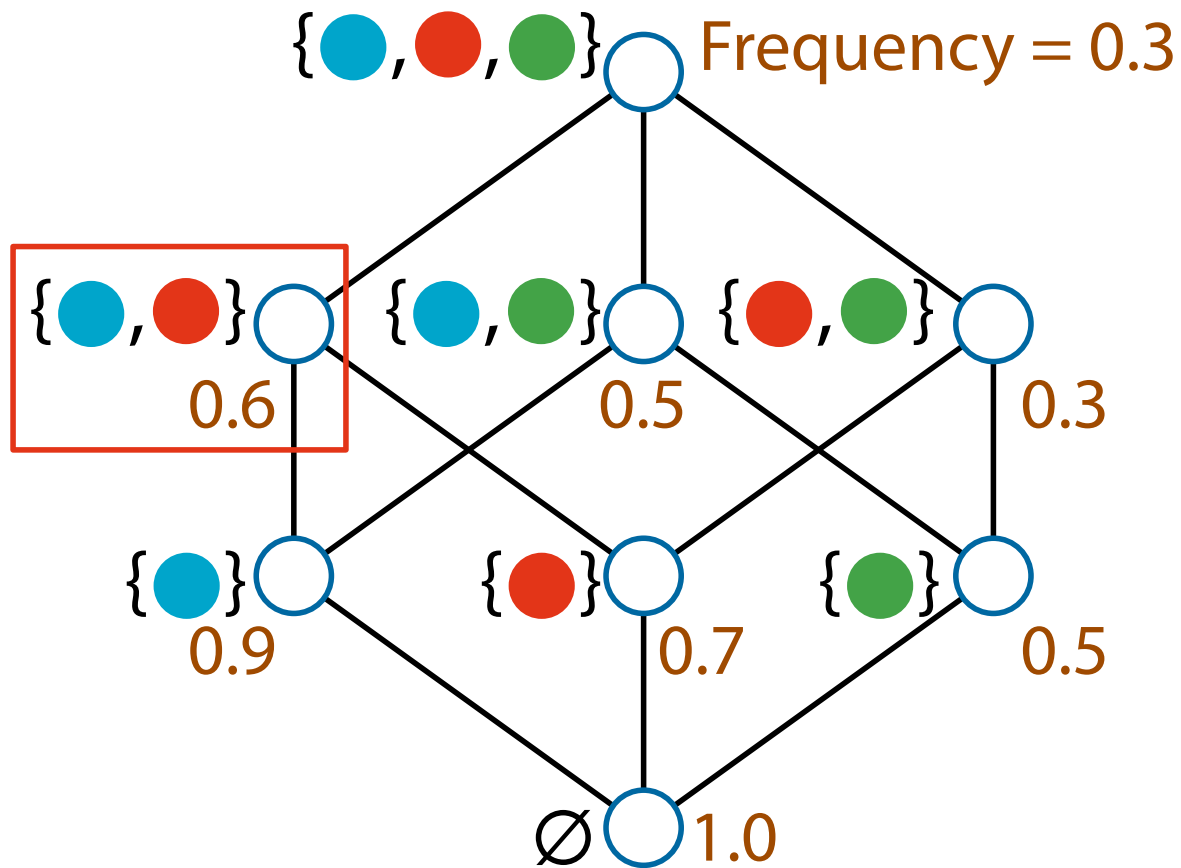
Itemset lattice






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
ID10:	0	1	0

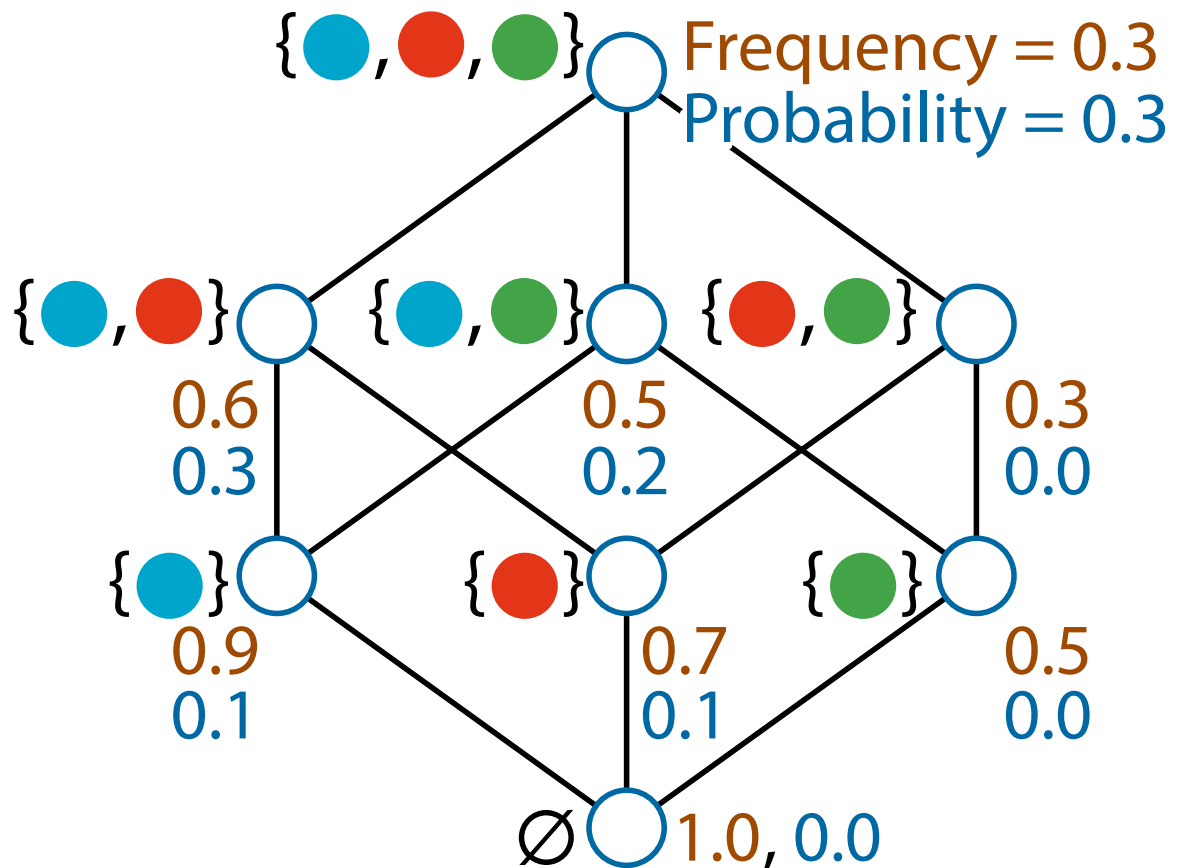
Itemset lattice



Transaction database

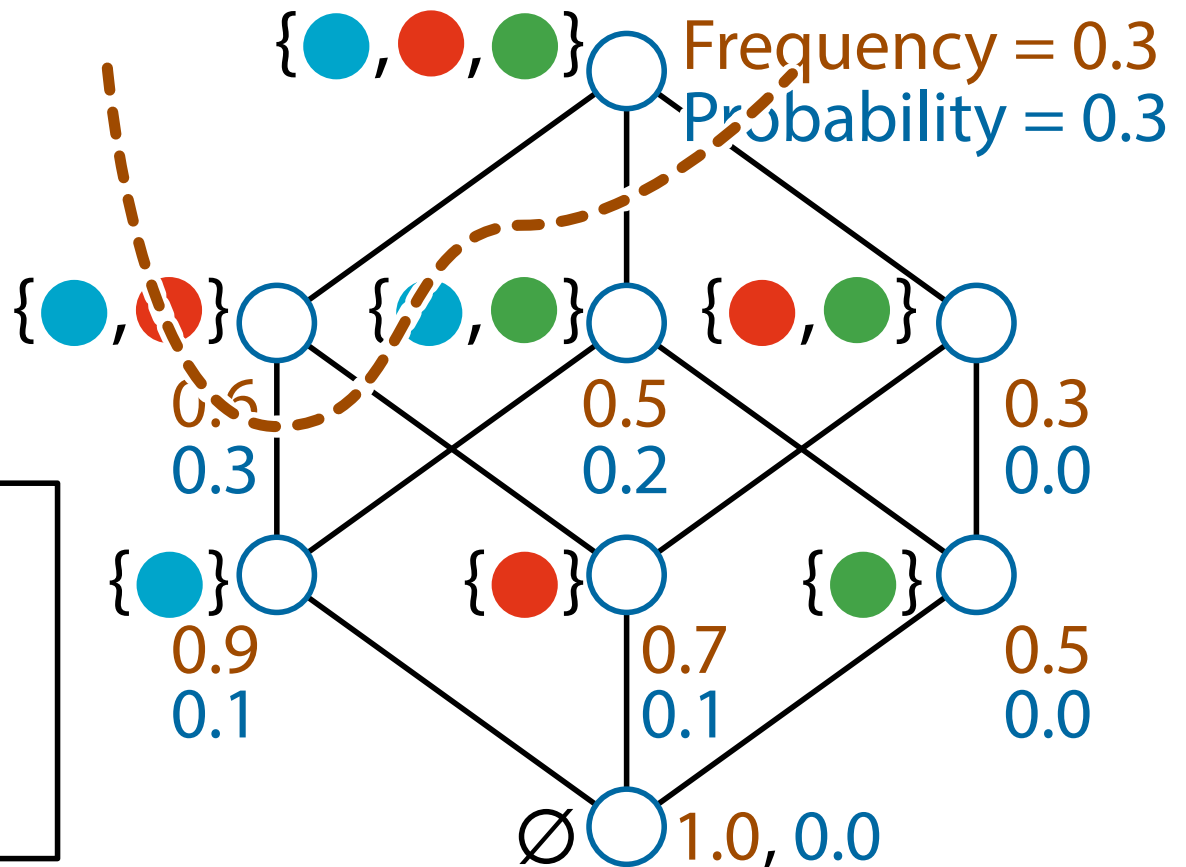
			
ID 1:	1	1	0
ID 2:	1	1	1
ID 3:	1	1	0
ID 4:	1	1	1
ID 5:	1	1	0
ID 6:	1	0	1
ID 7:	1	0	1
ID 8:	1	1	1
ID 9:	1	0	0
ID10:	0	1	0

Itemset lattice



Upward =
Pattern mining

Itemset lattice

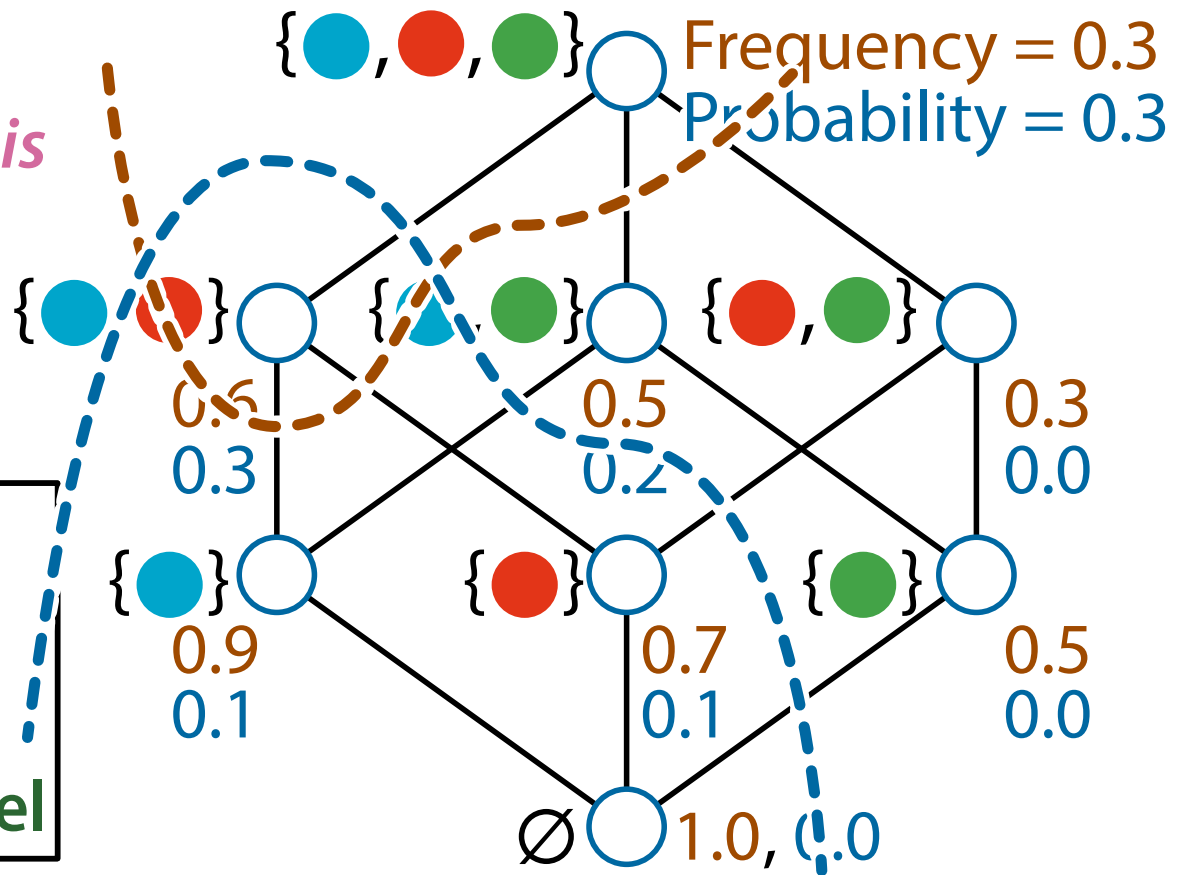


η : Frequency
 p : Probability

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

Upward =
Pattern mining
Downward =
Log-linear analysis

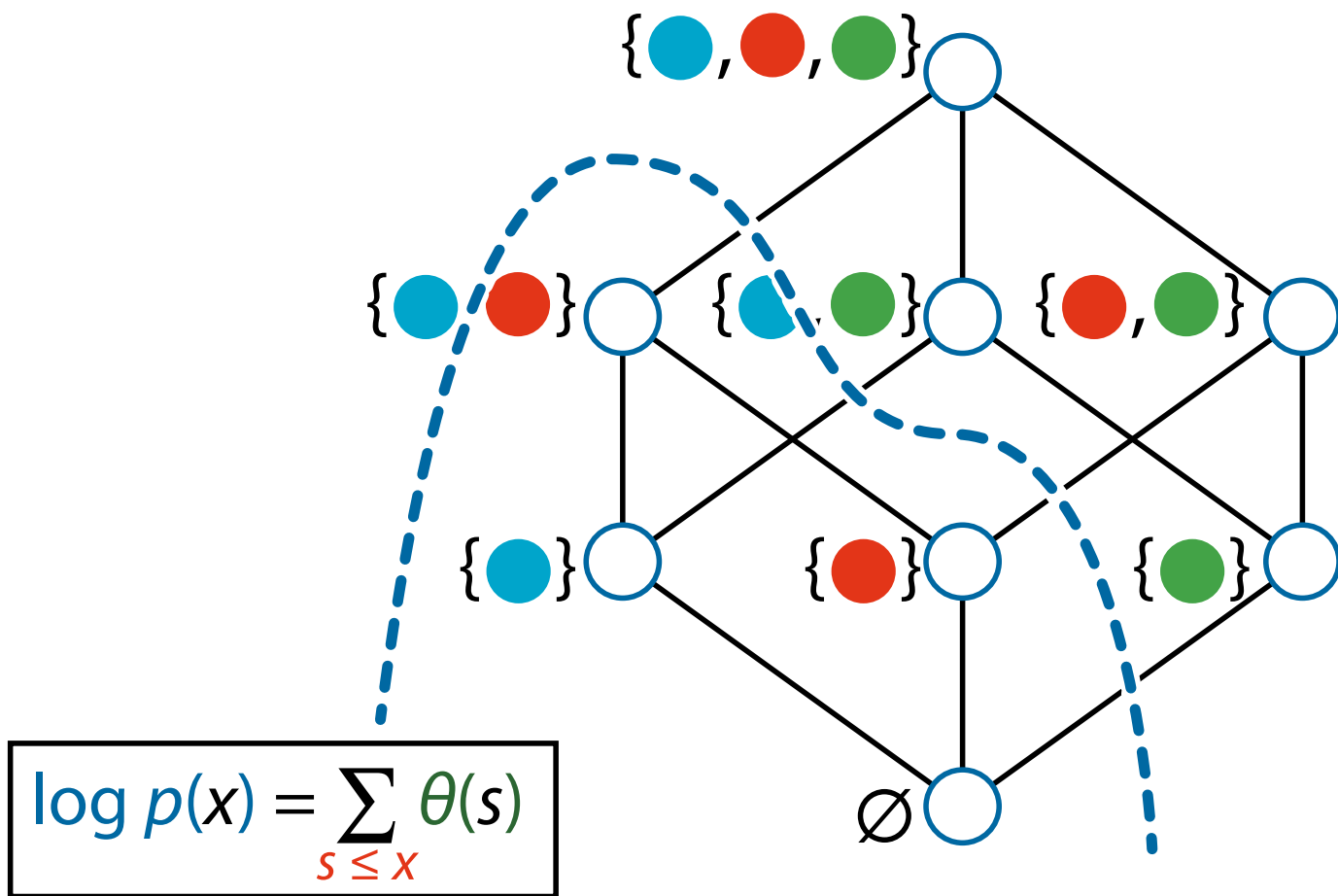
Itemset lattice

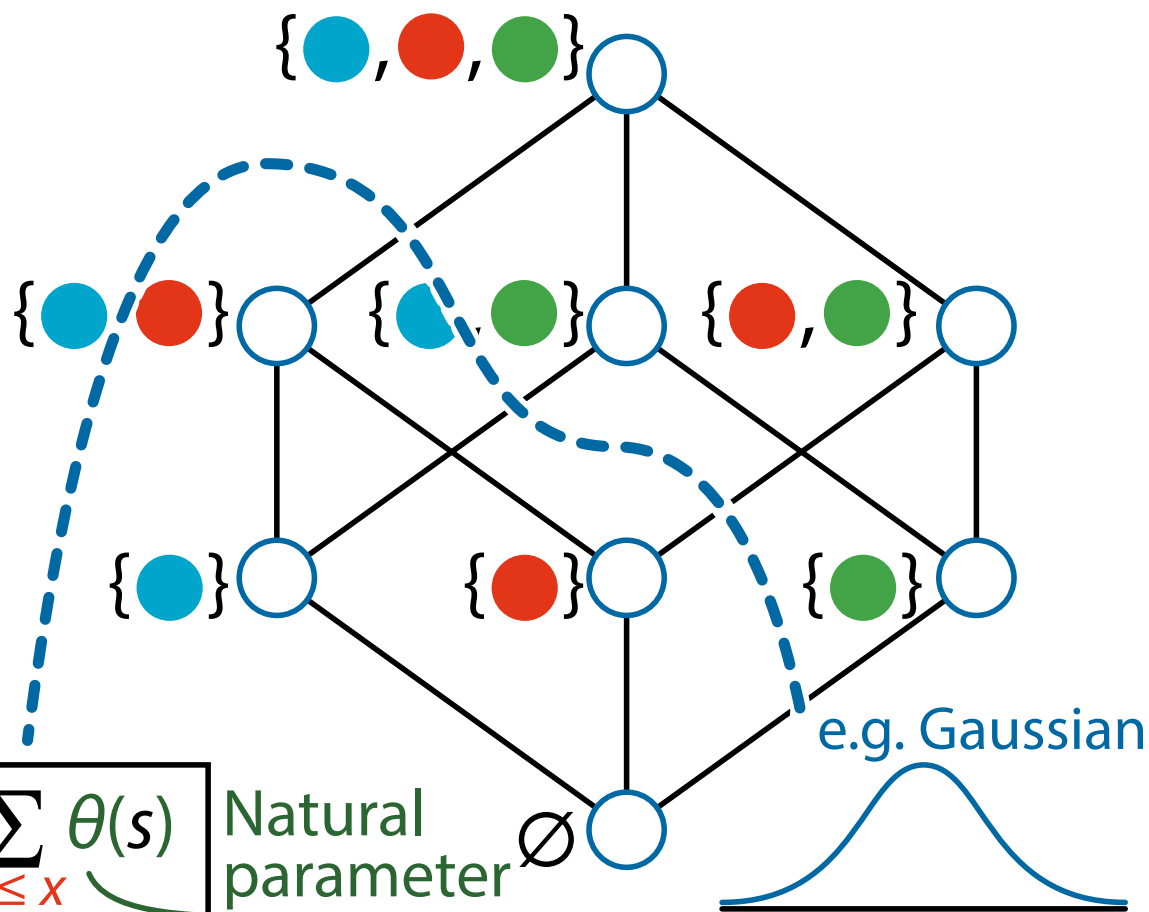


η : Frequency
 p : Probability
 θ : Coefficient of
log-linear model

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

$$\log p(\{\text{blue}, \text{red}\}) = \theta(\{\text{blue}, \text{red}\}) + \theta(\{\text{blue}\}) + \theta(\{\text{red}\}) + \theta(\emptyset)$$





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

Natural parameter θ

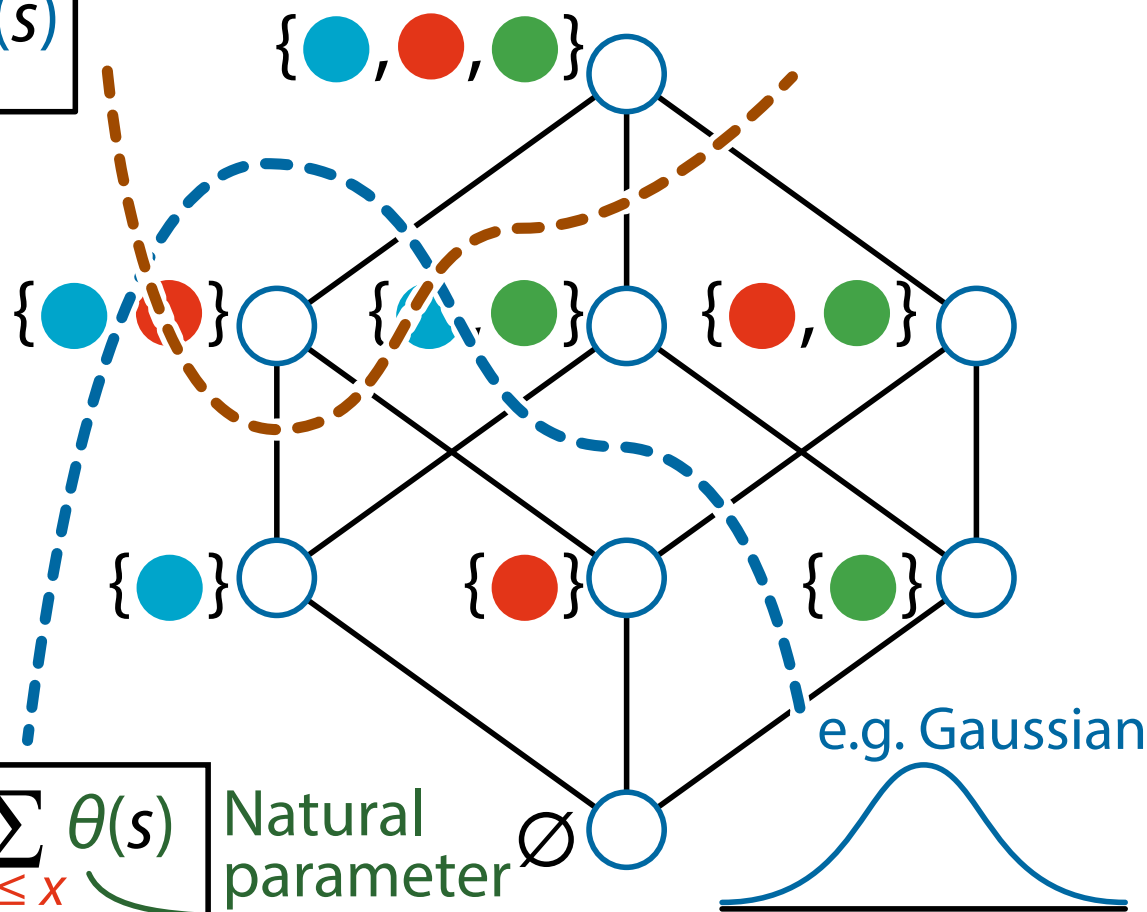
Exponential family:

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

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

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

Sufficient
statistics of
exponential
family



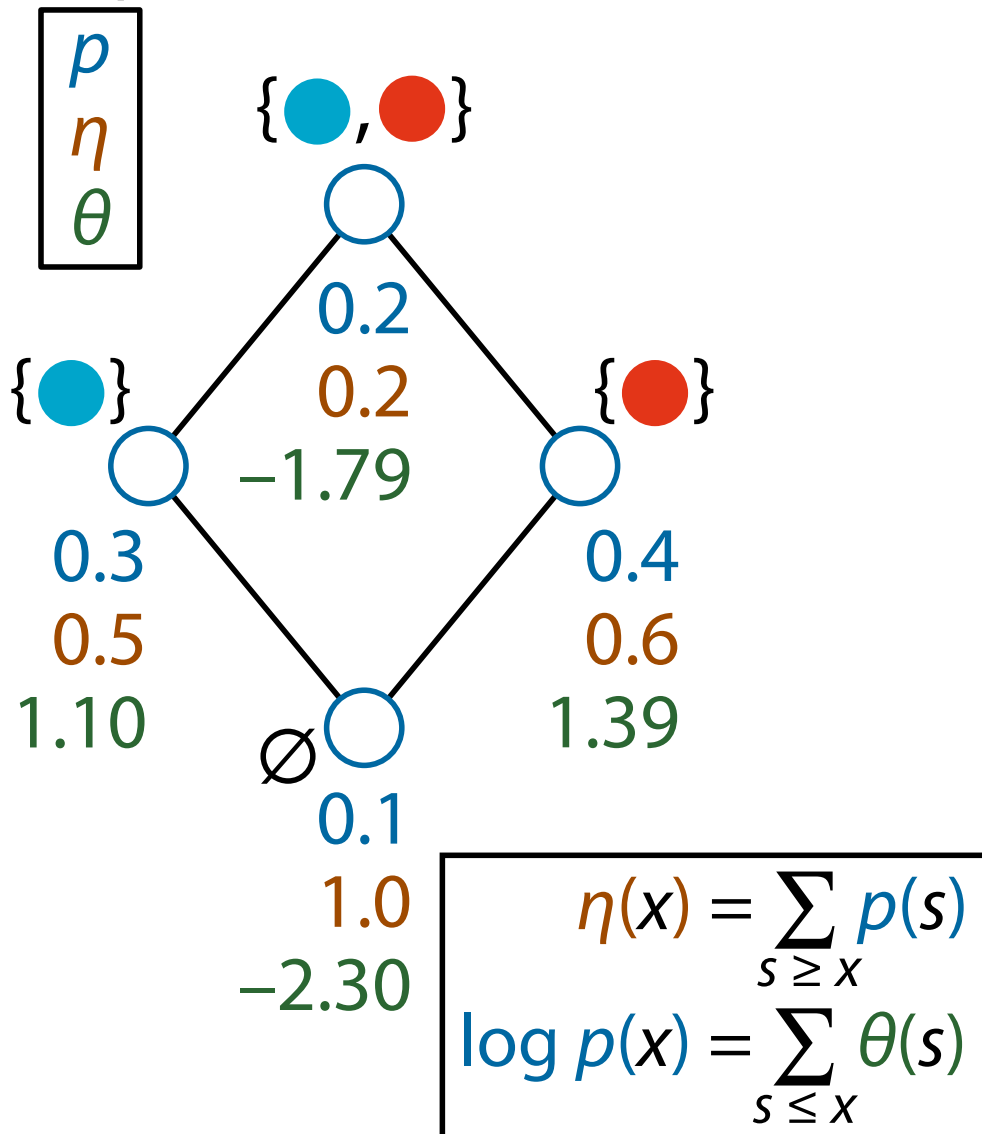
$$\log p(x) = \sum_{s \leq x} \theta(s)$$

Natural parameter θ

Exponential
family:

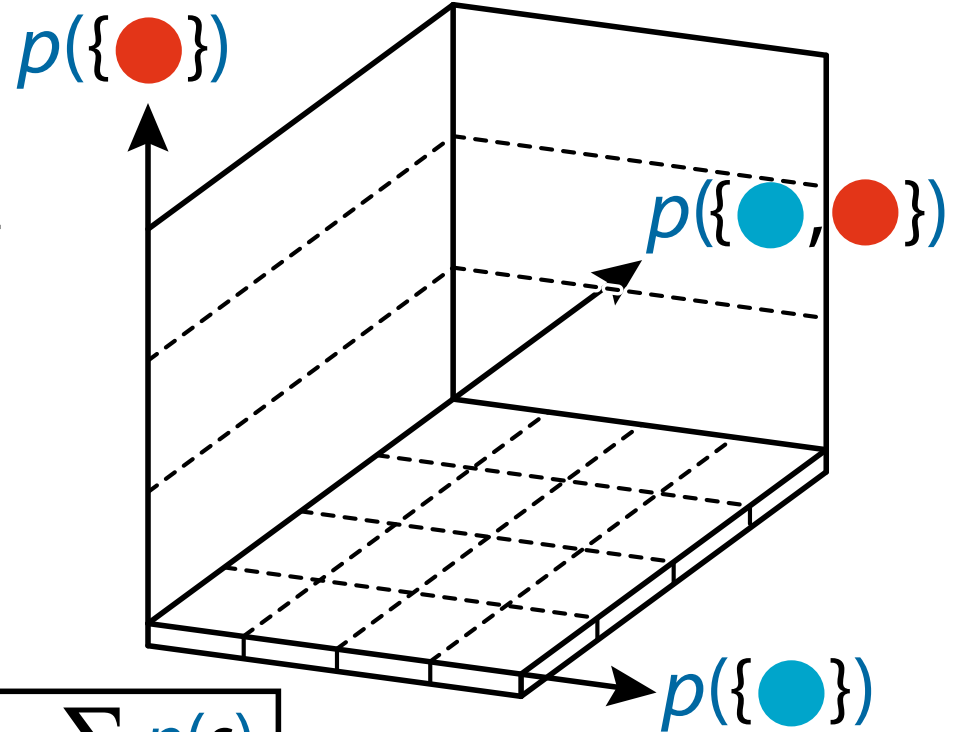
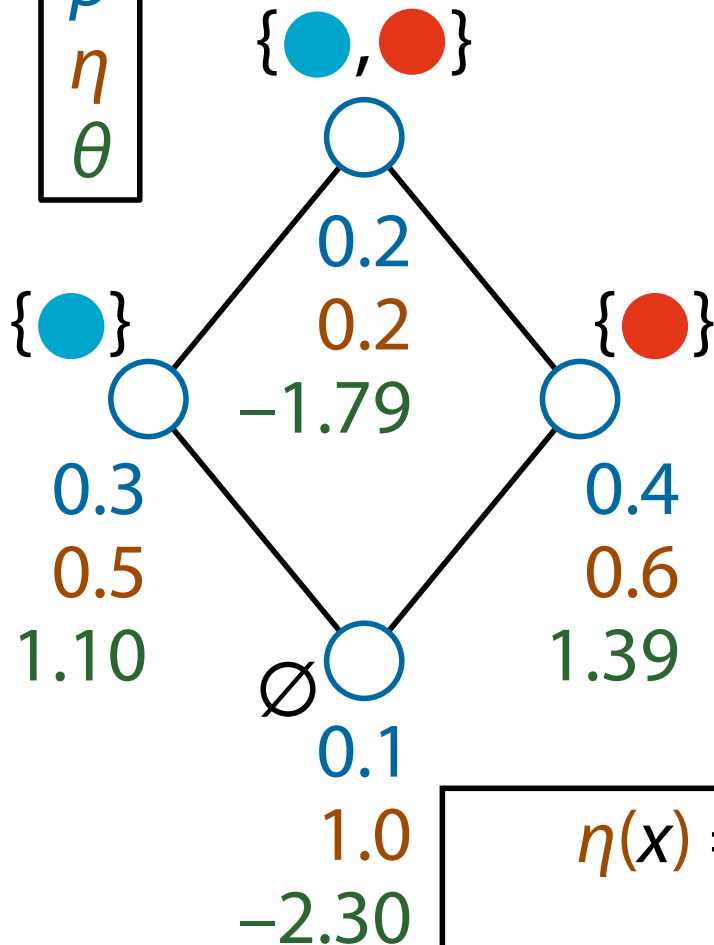
$$p(x) = \exp\left(\sum \theta(s) F_s(x) - \psi(\theta)\right)$$

Triple for each node



Triple for each node

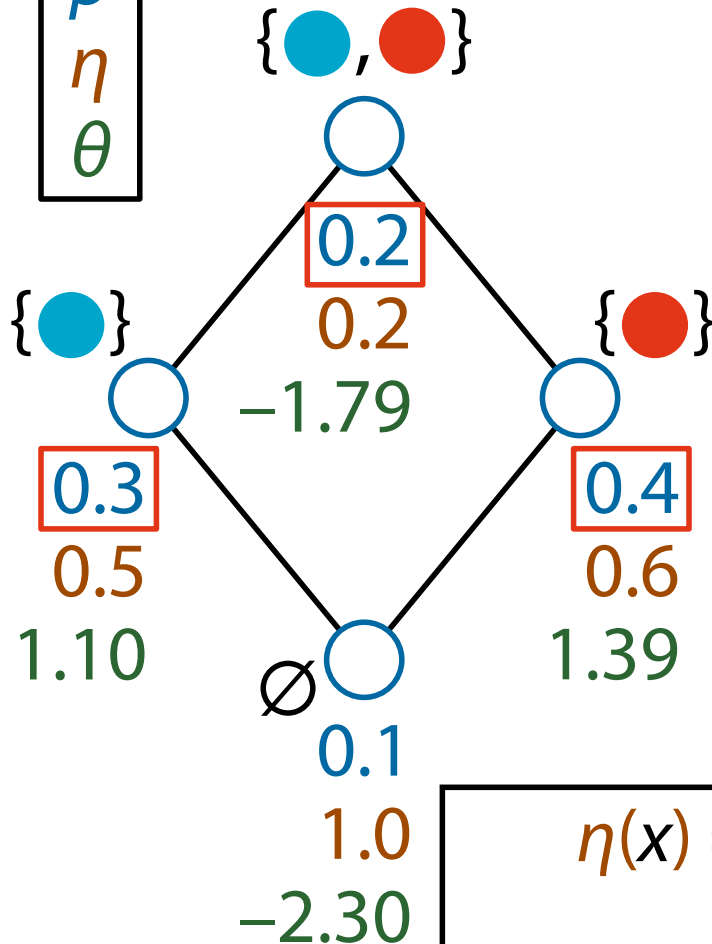
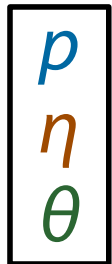
p
η
θ



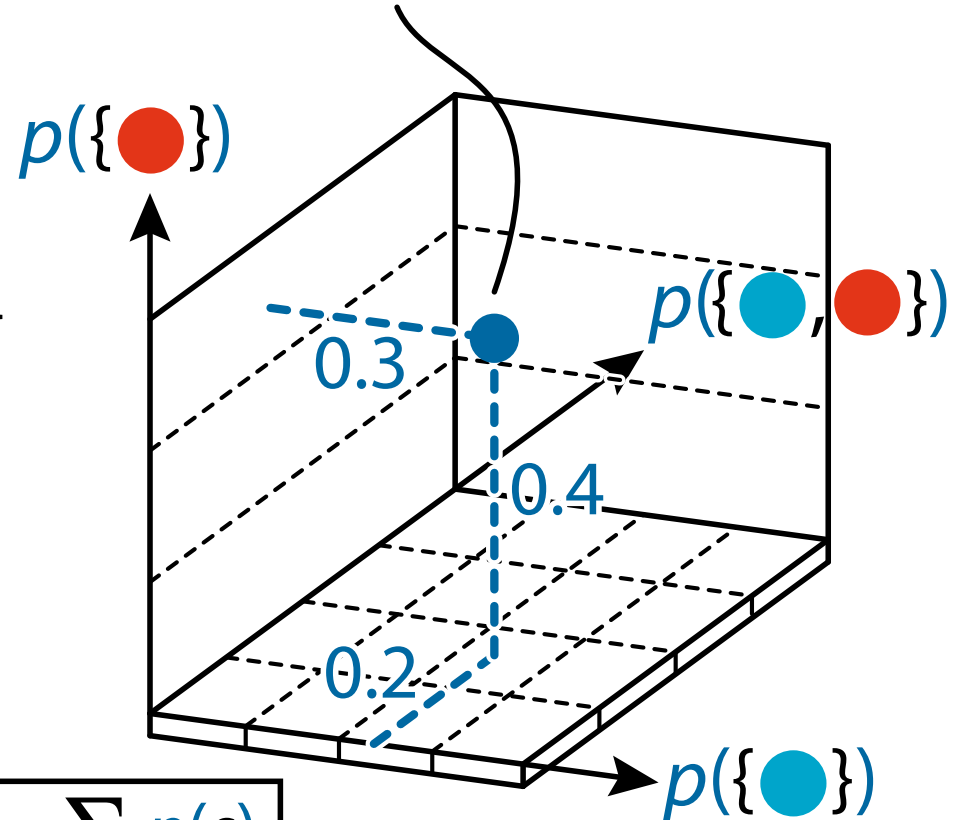
$$\eta(x) = \sum_{s \geq x} p(s)$$

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

Triple for each node



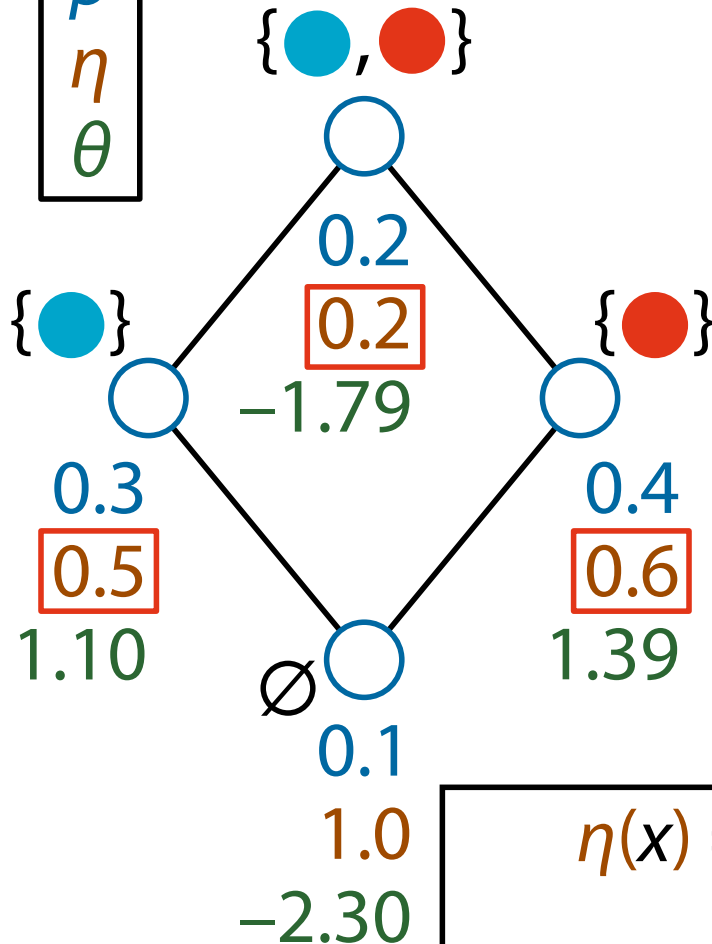
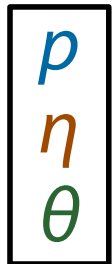
Probability distribution is a "point" in 3D space



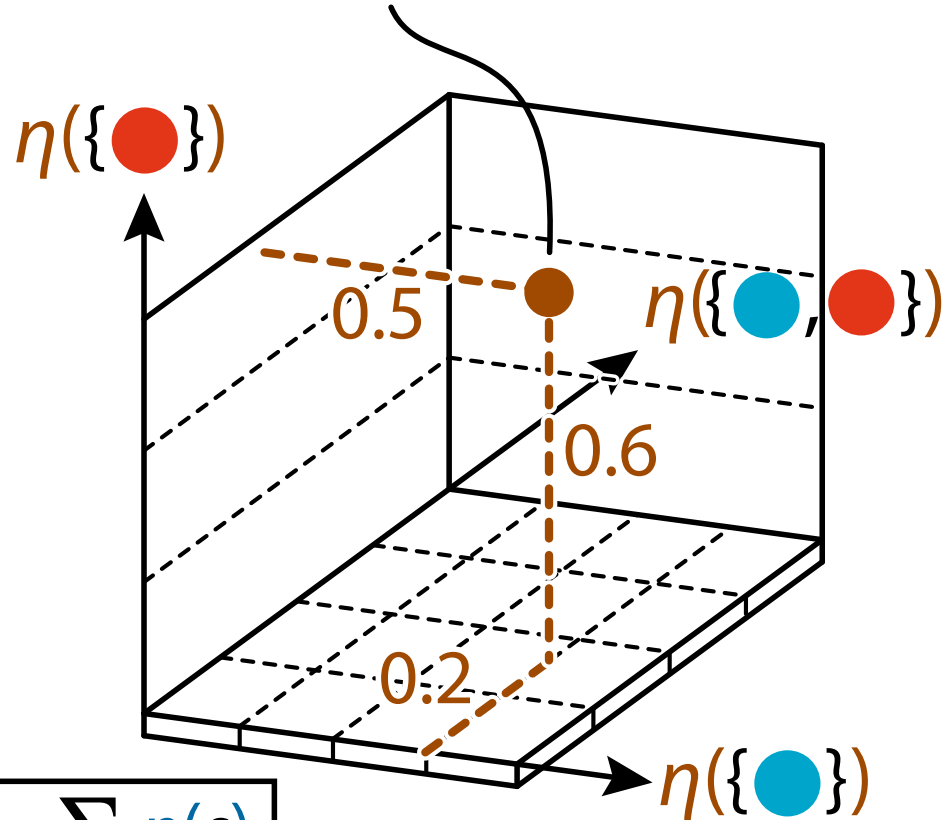
$$\eta(x) = \sum_{s \geq x} p(s)$$

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

Triple for each node



Probability distribution is a "point" in 3D space

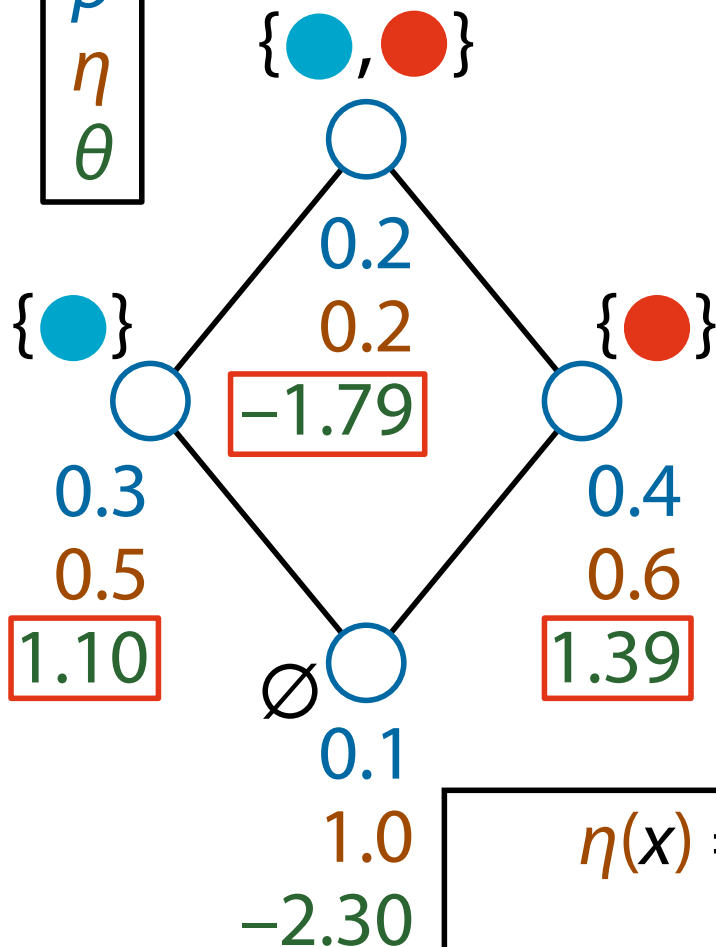


$$\eta(x) = \sum_{s \geq x} p(s)$$

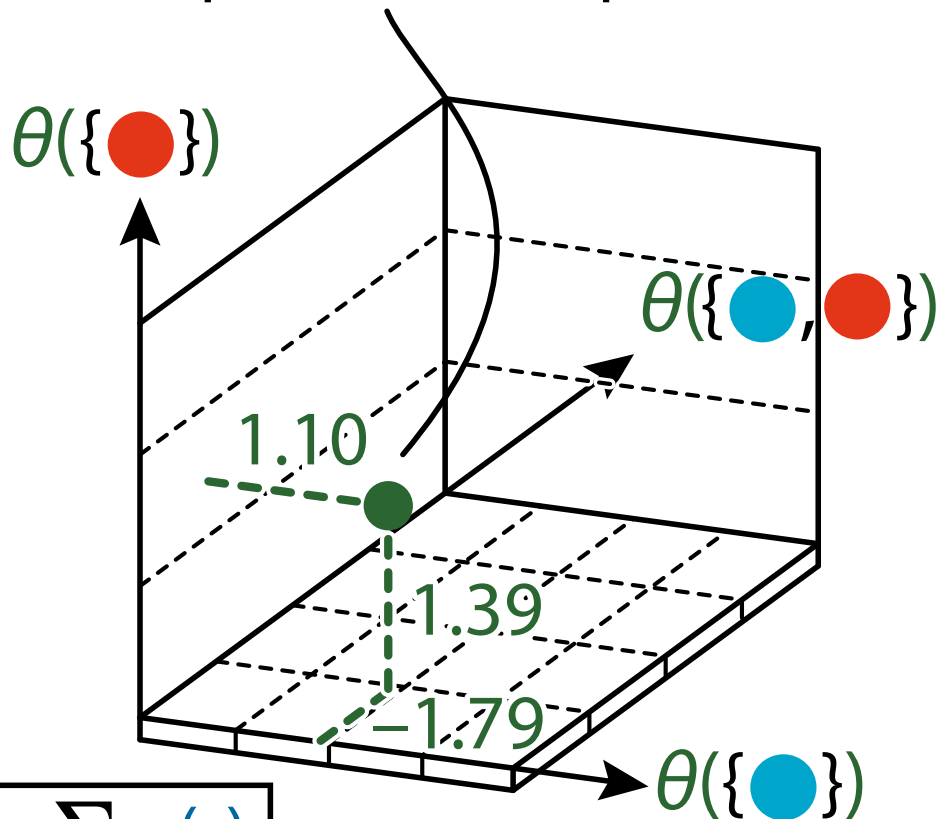
$$\log p(x) = \sum_{s \leq x} \theta(s)$$

Triple for each node

p
 η
 θ



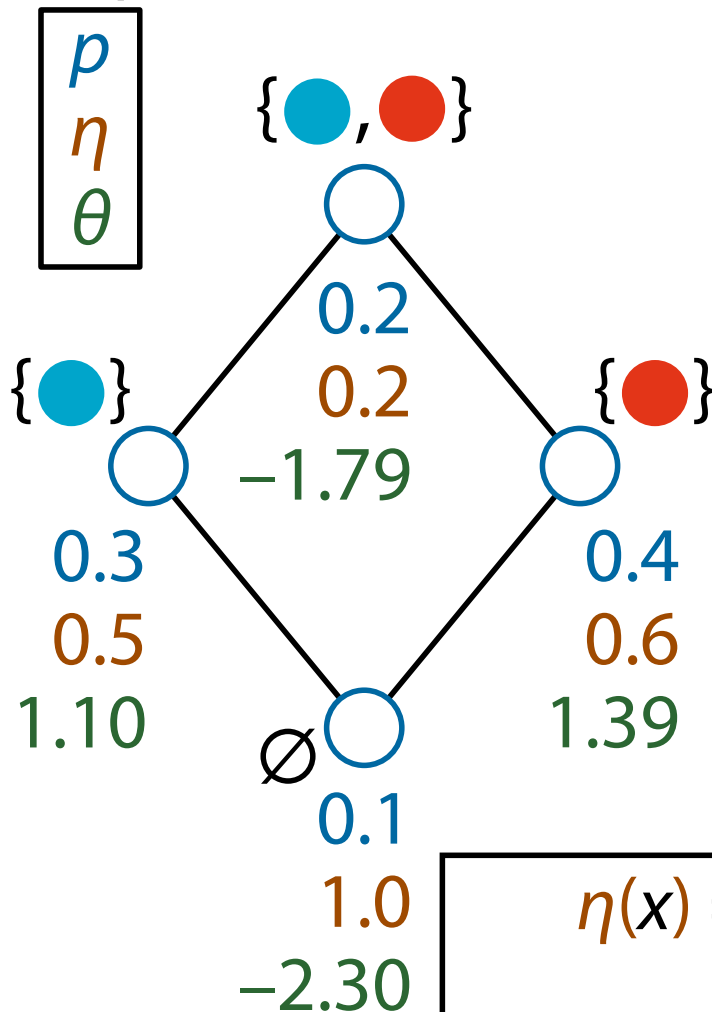
Probability distribution is a "point" in 3D space



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

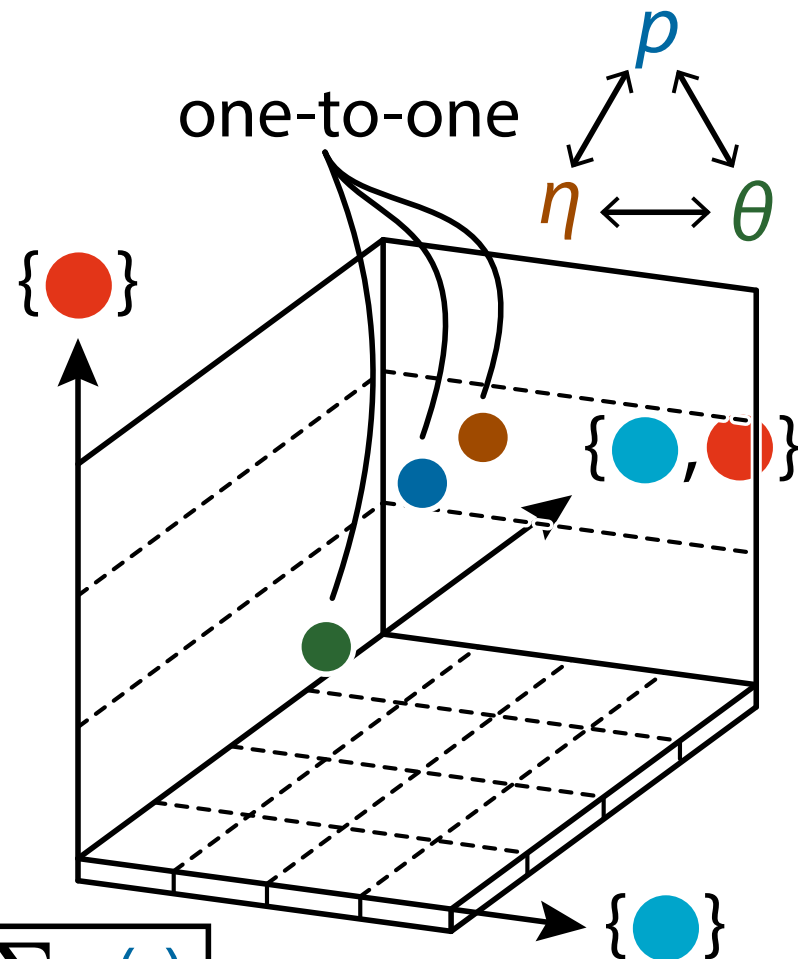
$$\log p(x) = \sum_{s \leq x} \theta(s)$$

Triple for each node

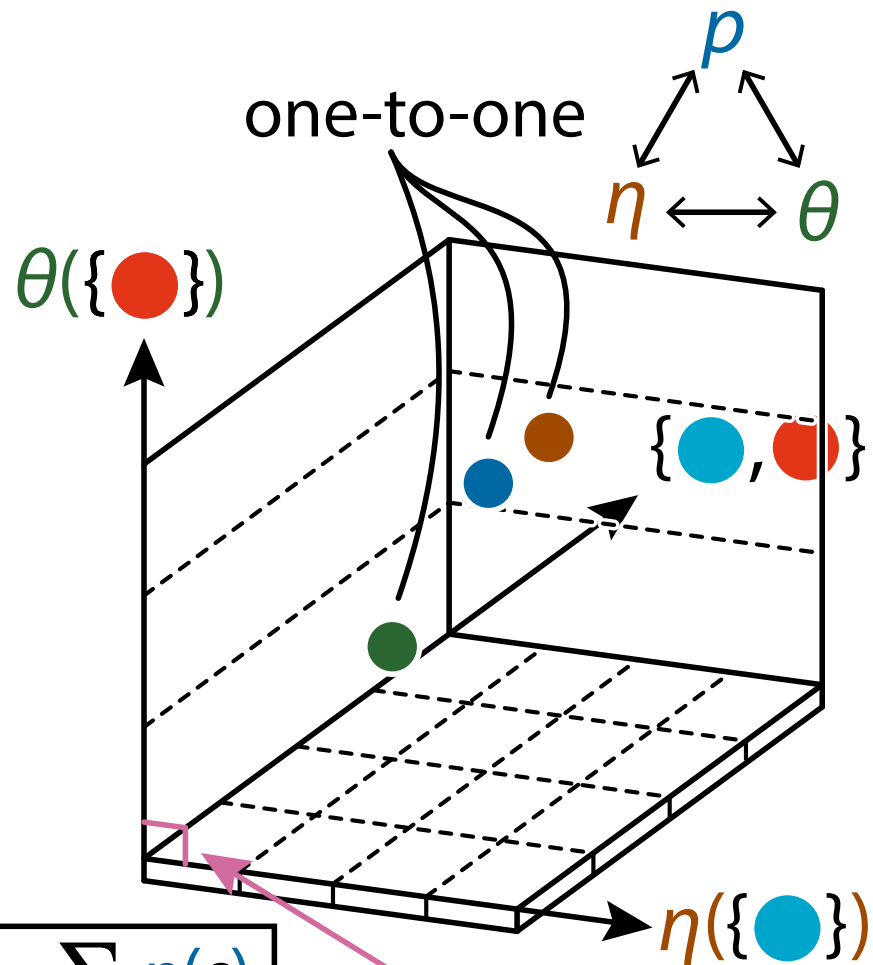
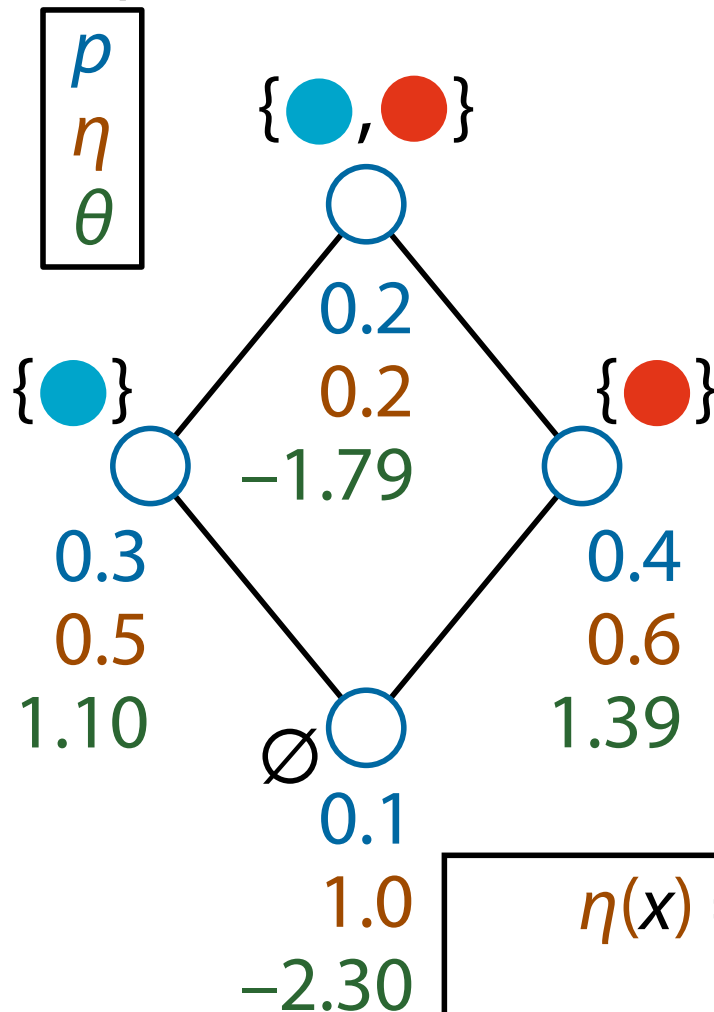


$$\eta(x) = \sum_{s \geq x} p(s)$$

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



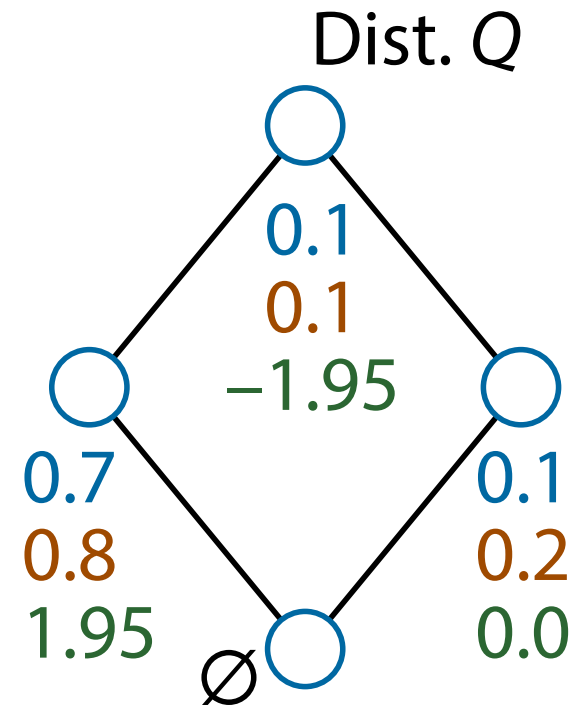
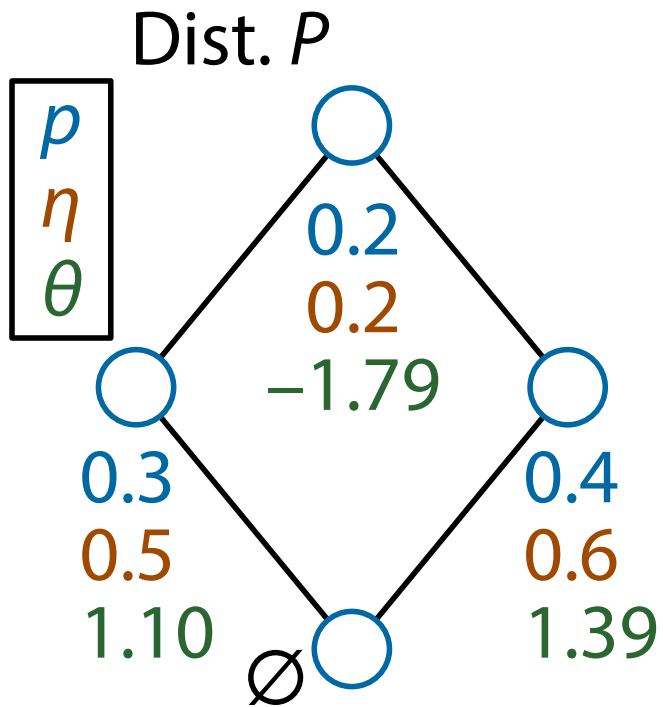
Triple for each node

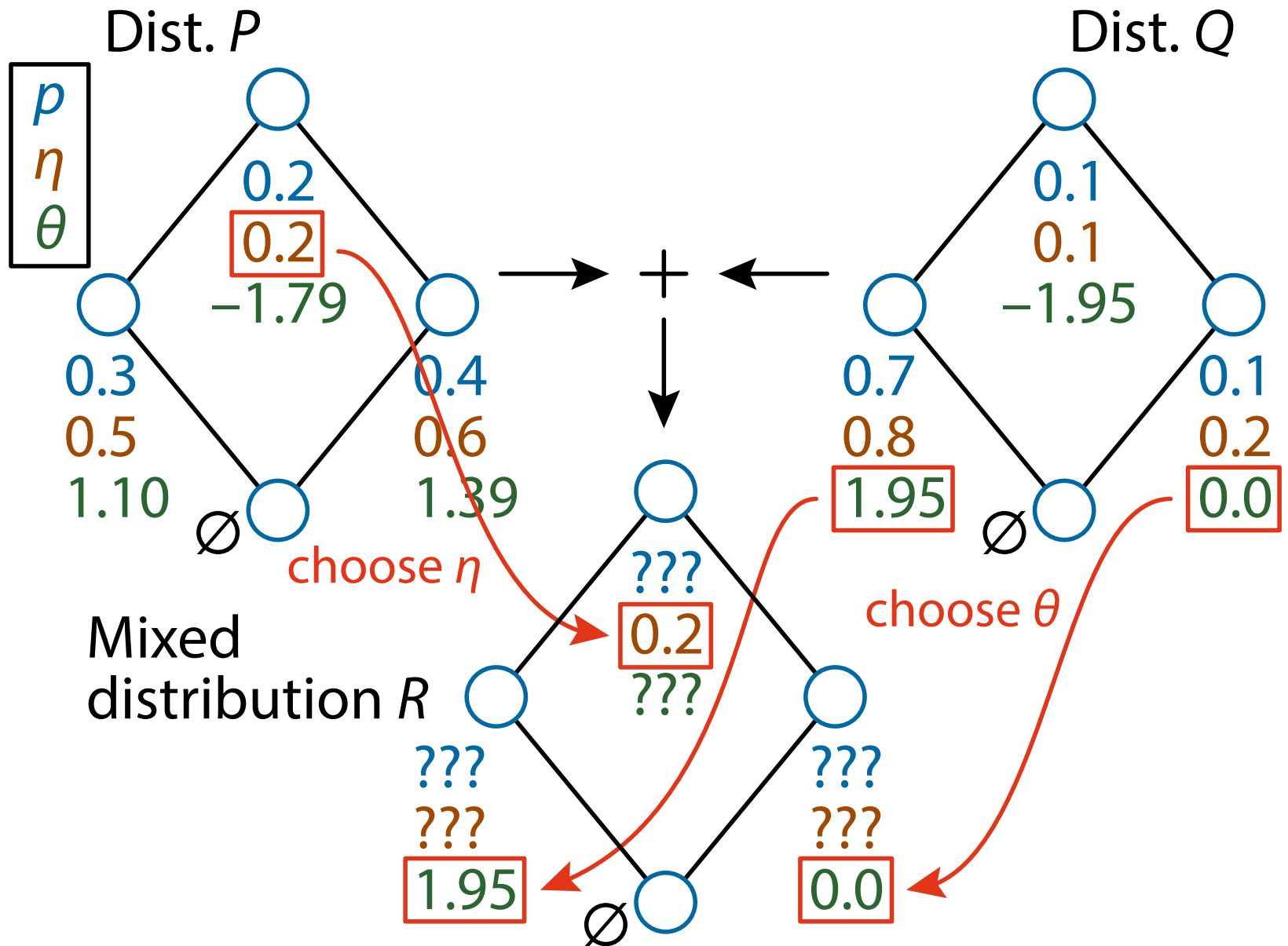


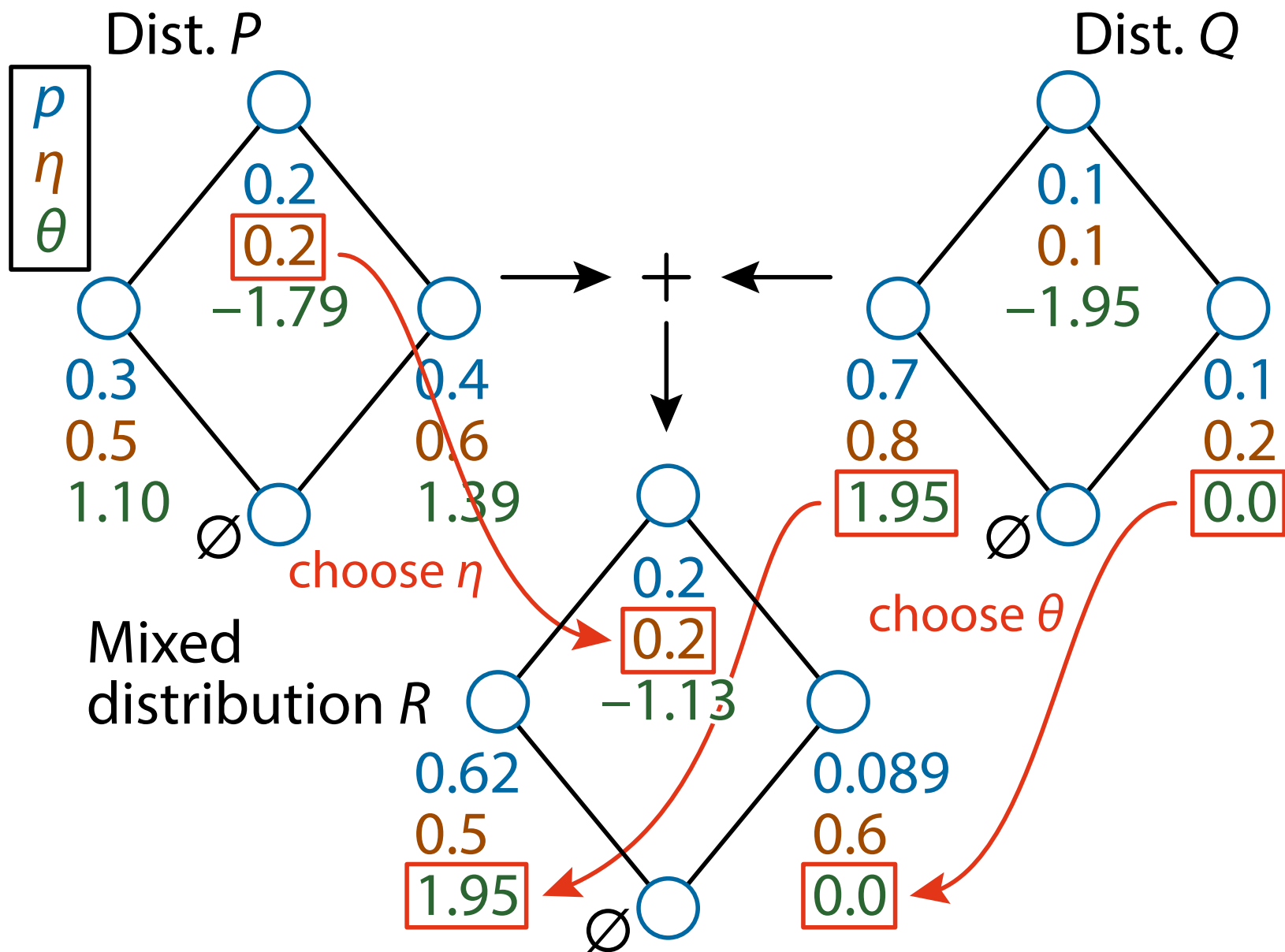
$$\eta(x) = \sum_{s \geq x} p(s)$$

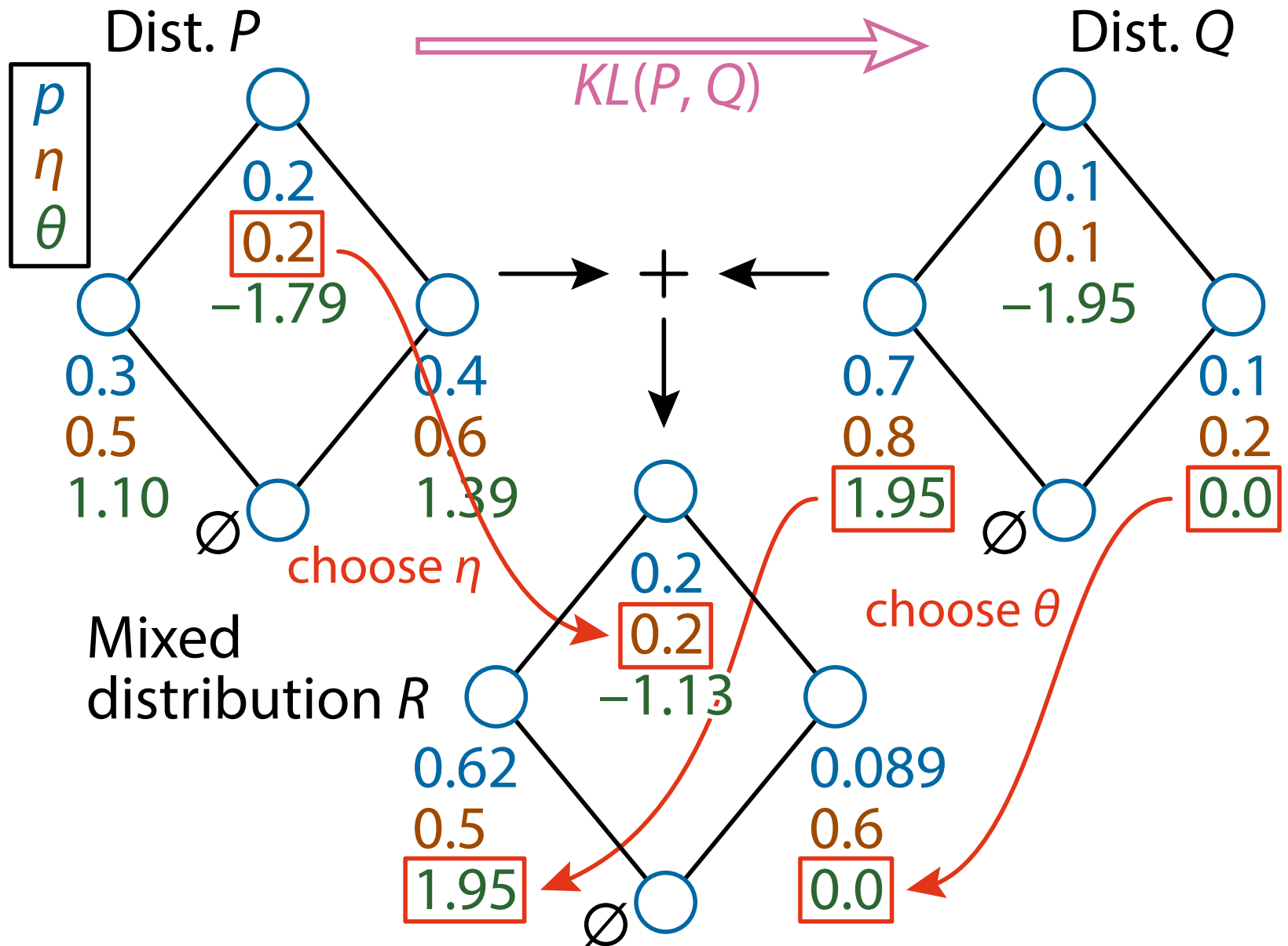
$$\log p(x) = \sum_{s \leq x} \theta(s)$$

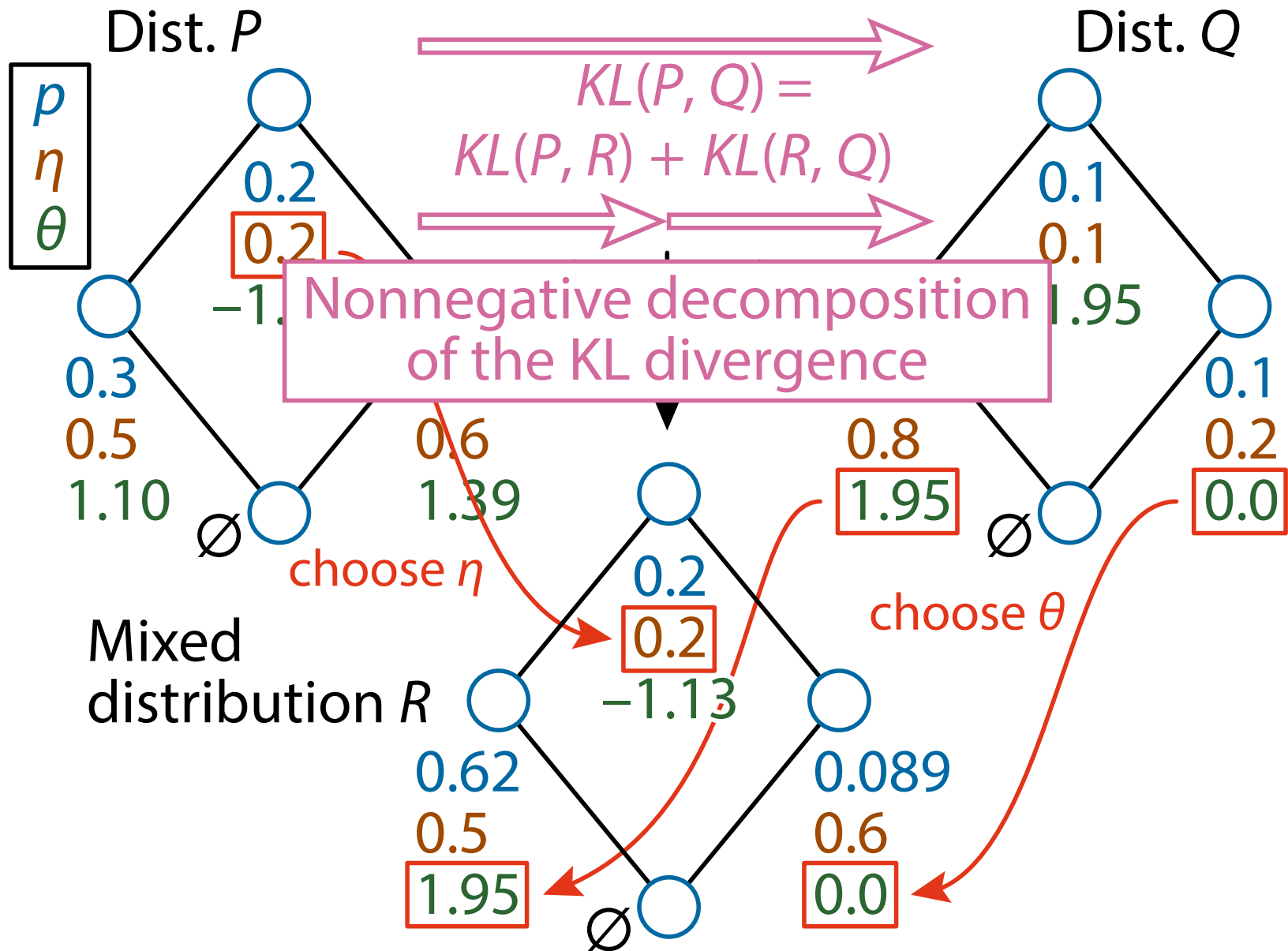
θ and η are dually orthogonal

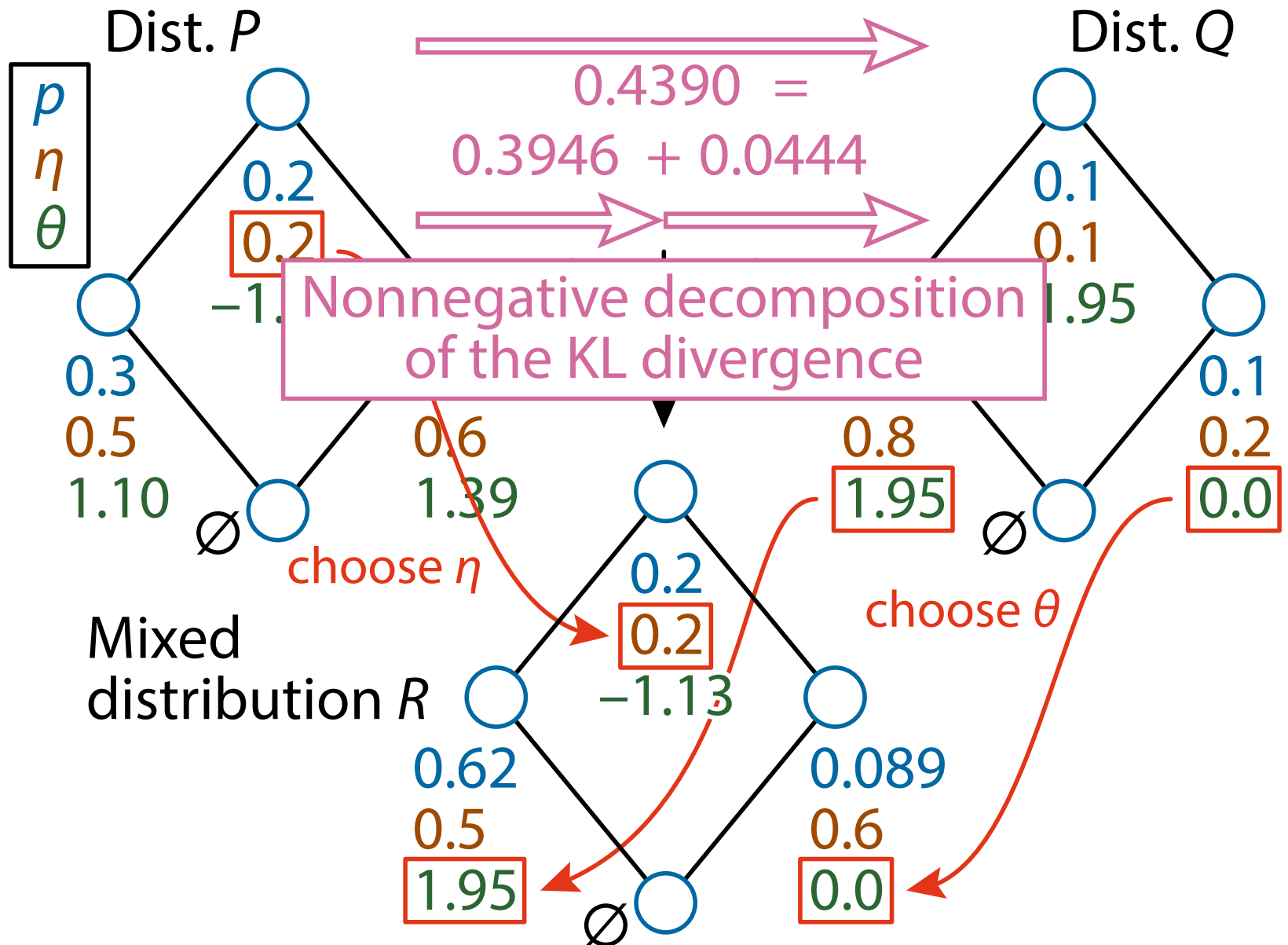


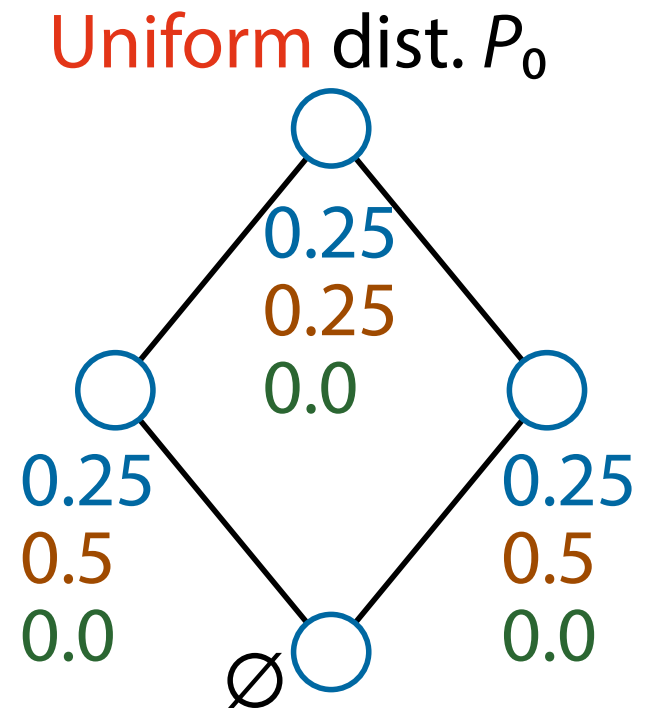
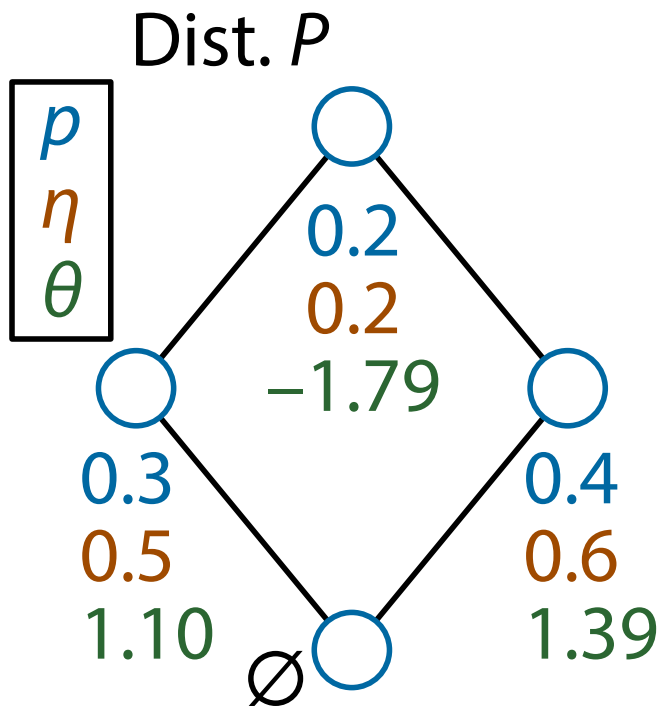


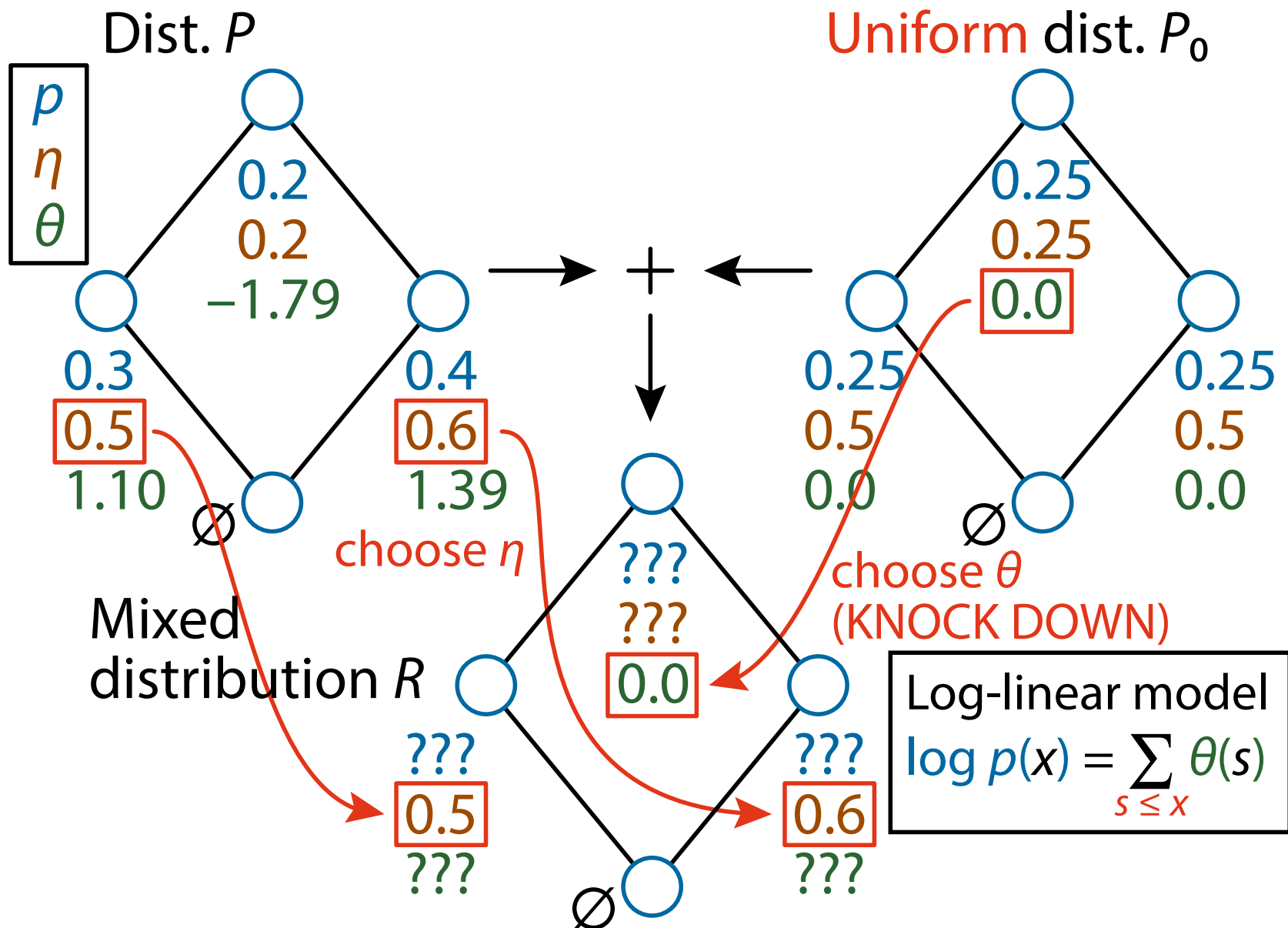


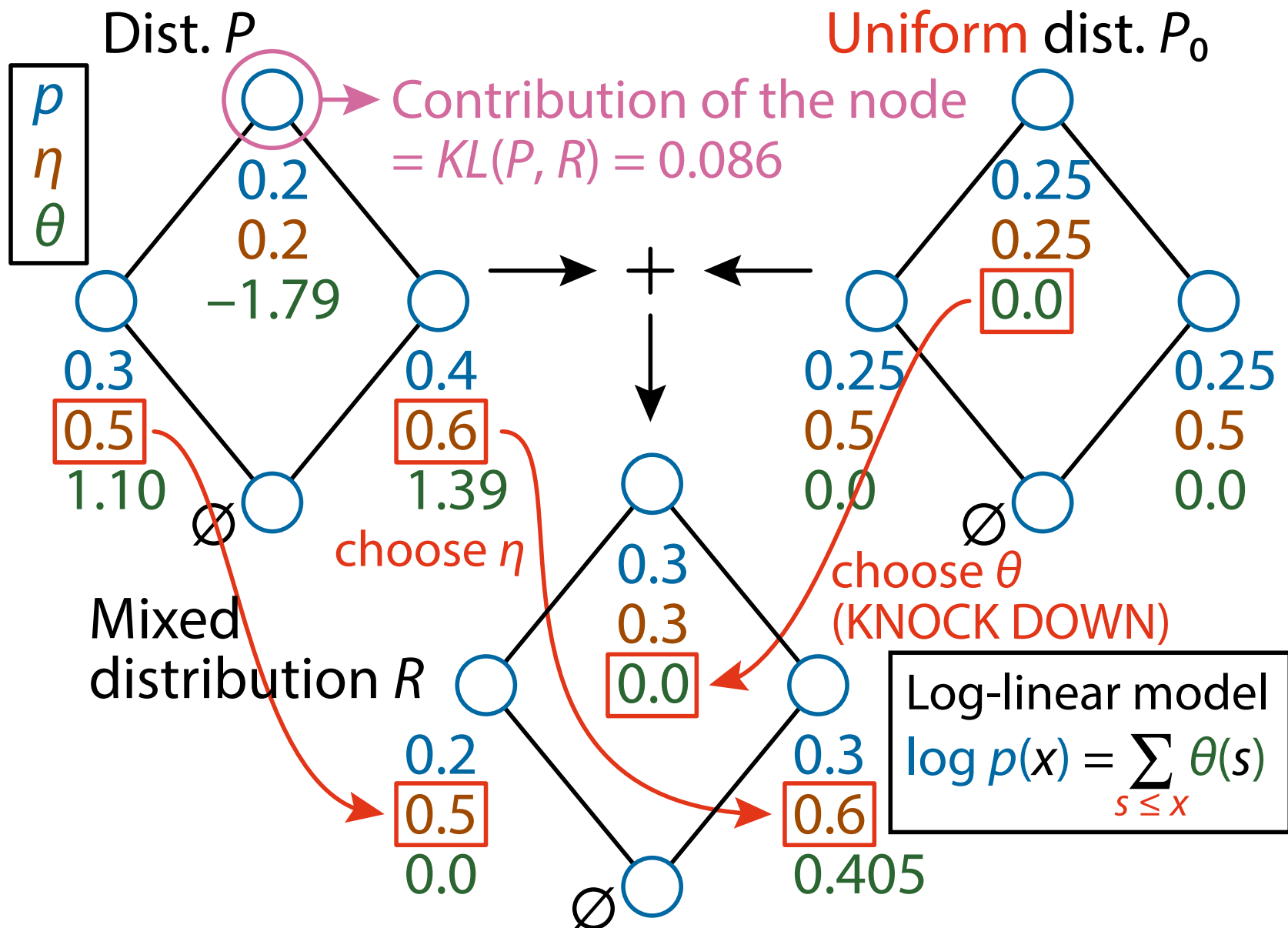


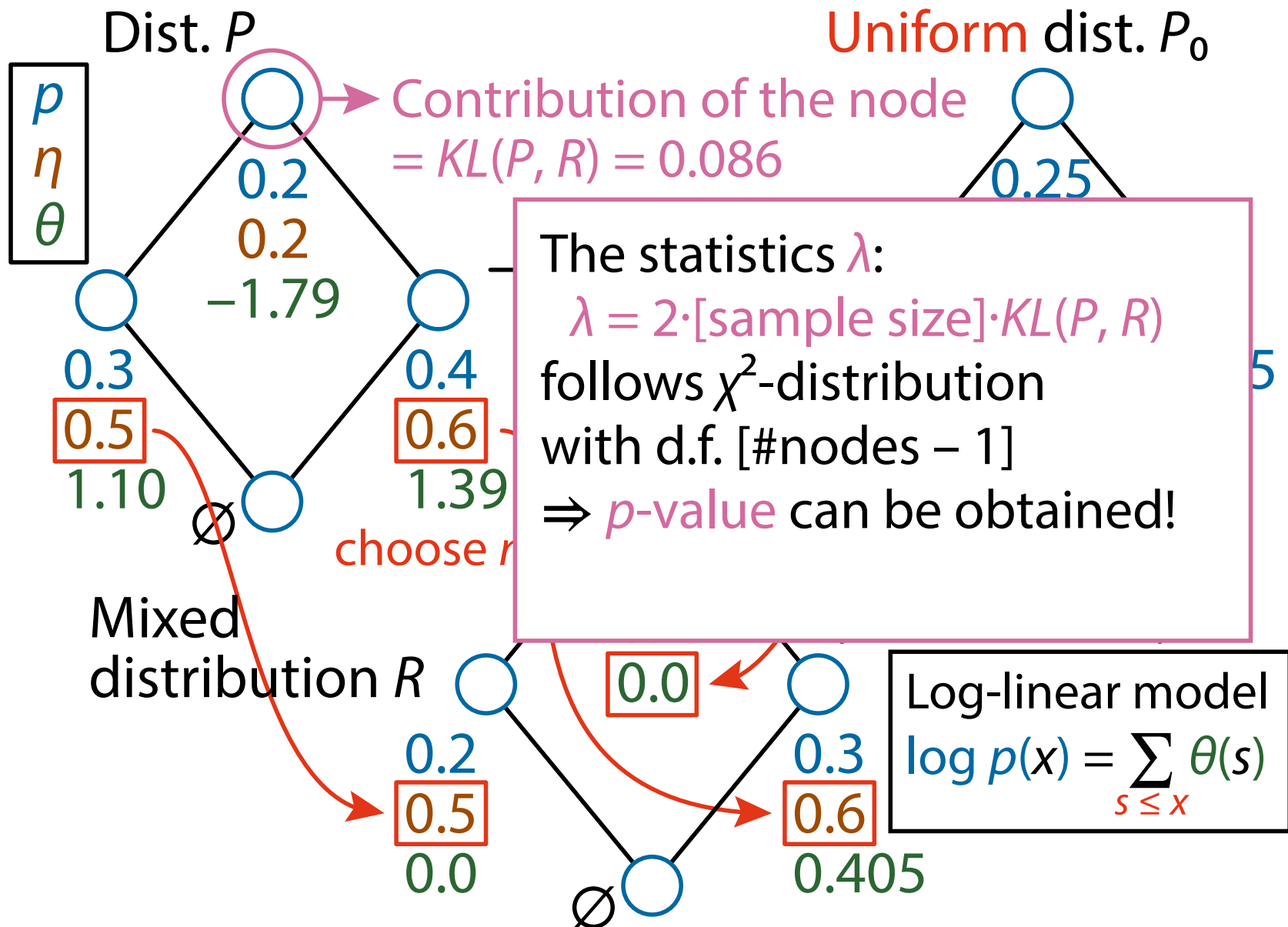


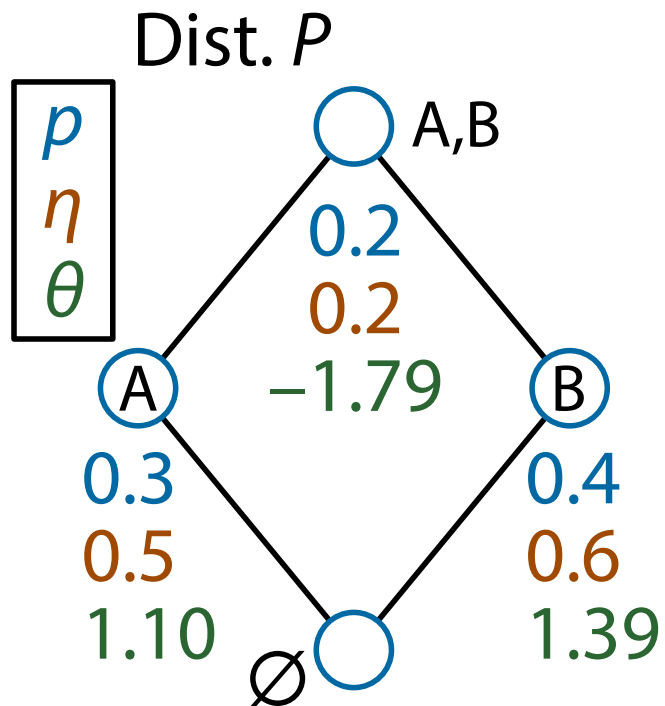




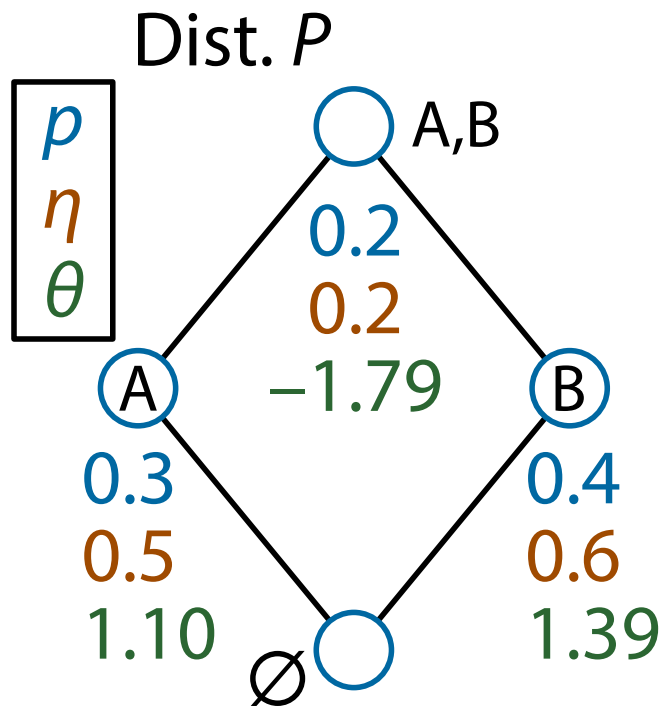




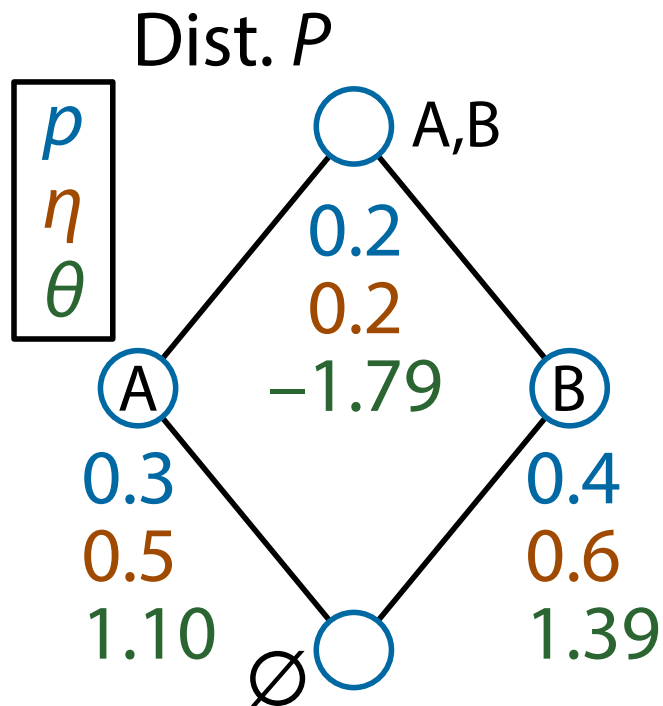




A	B	A, B
0.3	0.4	0.2
0.5	0.6	0.2
1.10	1.39	-1.79



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
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




KL = Score of A

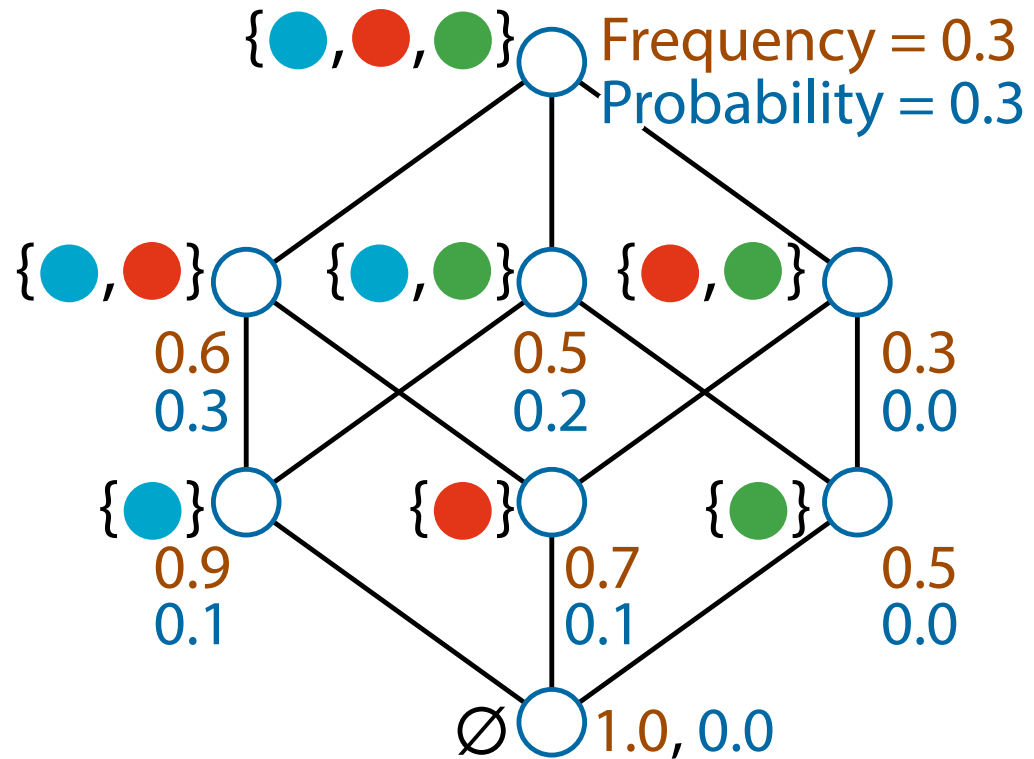
KL = Score of B

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

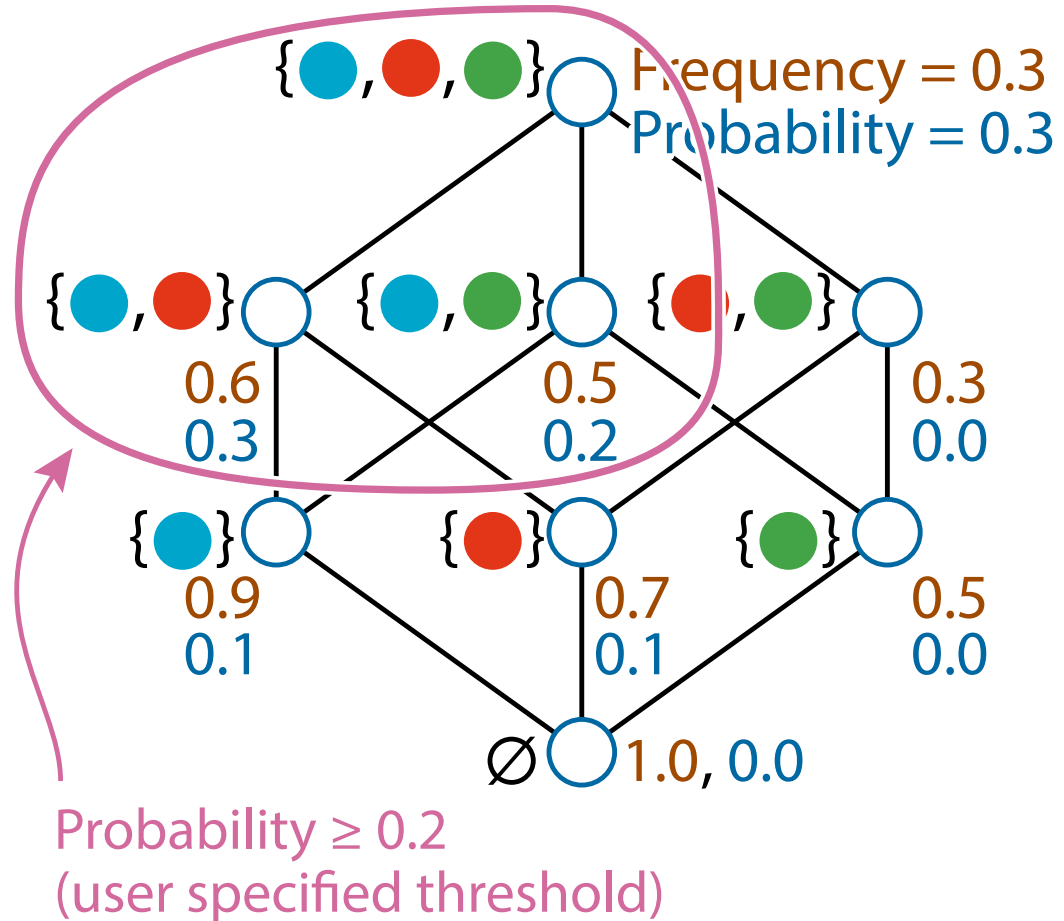


Number of nodes = $2^{\text{\#features}}$
⇒ combinatorial explosion!

Make a Poset from Data




Dataset

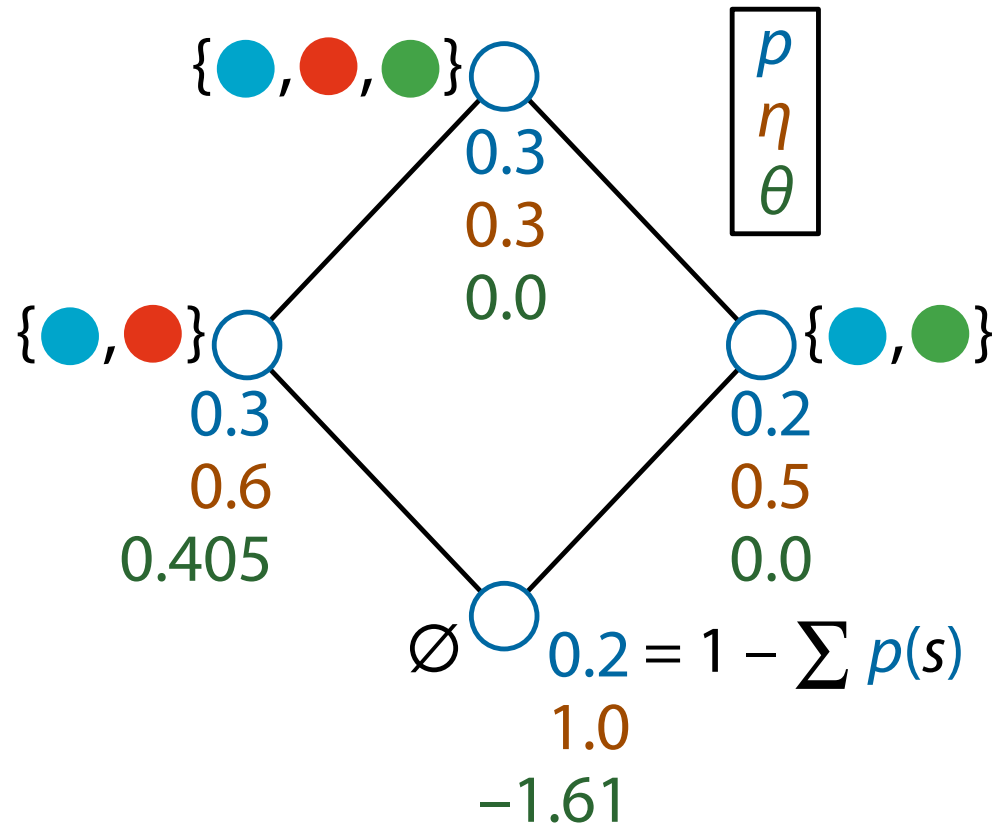
			
ID 1:	1	1	0
ID 2:	1	1	1
ID 3:	1	1	0
ID 4:	1	1	1
ID 5:	1	1	0
ID 6:	1	0	1
ID 7:	1	0	1
ID 8:	1	1	1
ID 9:	1	0	0
ID10:	0	1	0



Remove Nodes with Probability 0

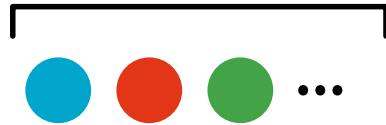
Dataset

			
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

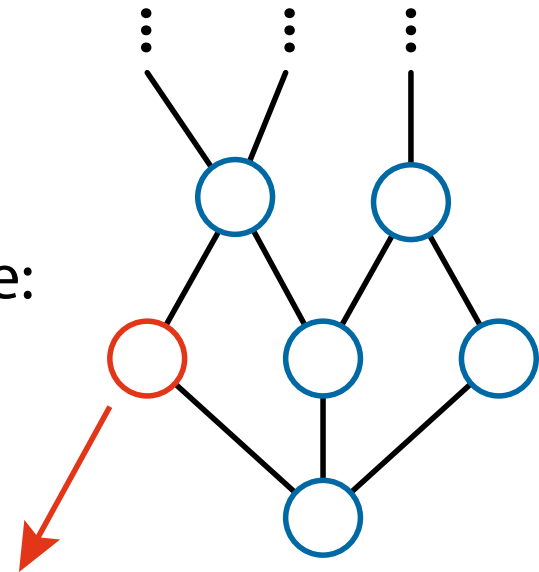


ID 1:	1	1	0
ID 2:	1	1	1
ID 3:	1	1	0 ...
ID 4:	1	1	1
ID 5:	1	1	0
⋮	⋮		

Total runtime:
4.95 seconds

Sample size:
990,002

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



significant interactions: **583**



Single feature: 537

Pairwise interactions: 41

Triple interactions: 5

Example on Real Data (accidents)

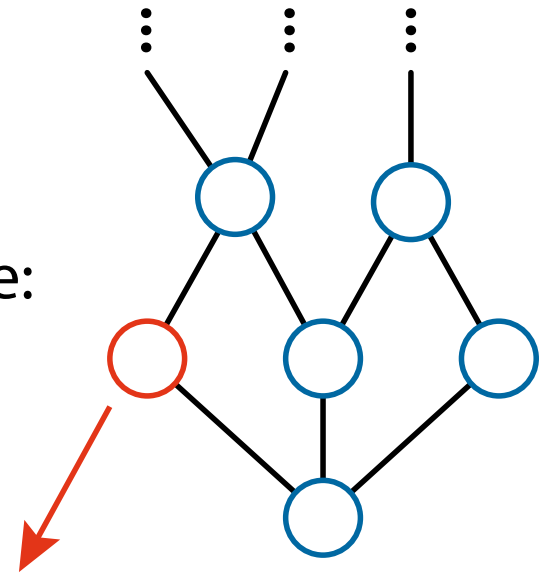
features: 468

				...
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
- We can decompose a probability distribution and assess the significance of any-order interactions beyond pairwise interactions