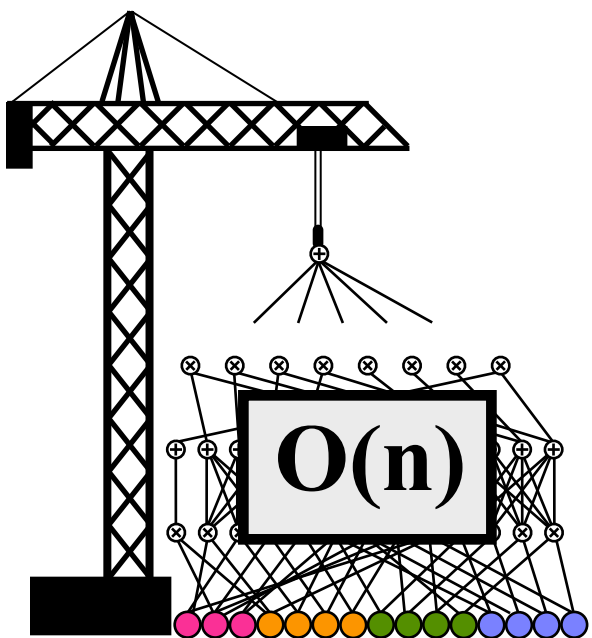


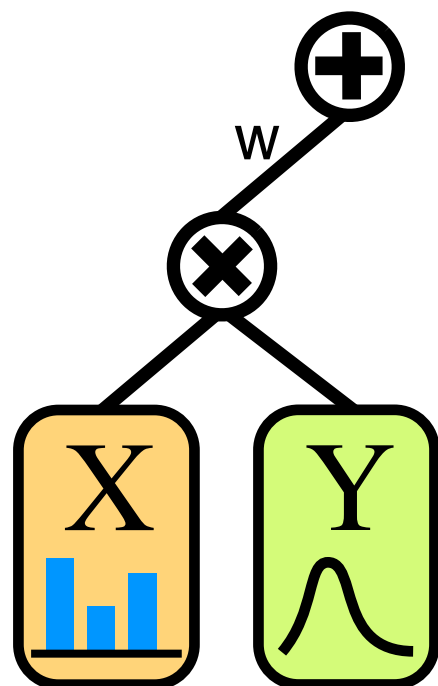
Learning the Structure of Sum-Product Networks

Robert Gens
Pedro Domingos

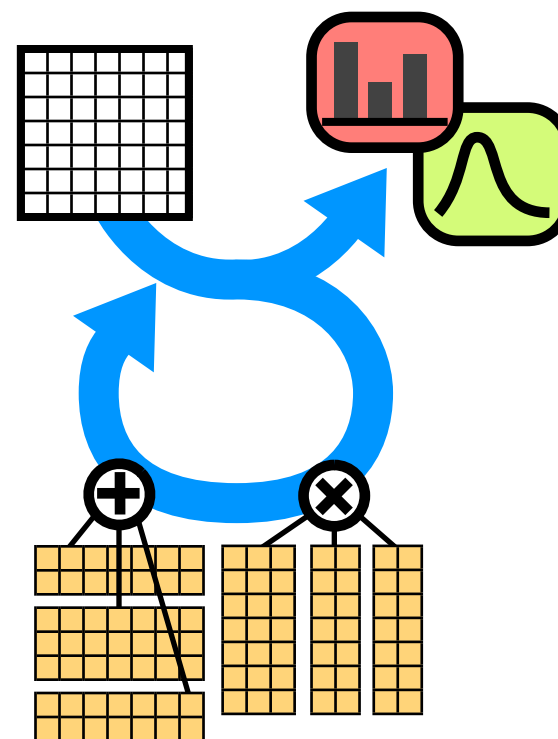




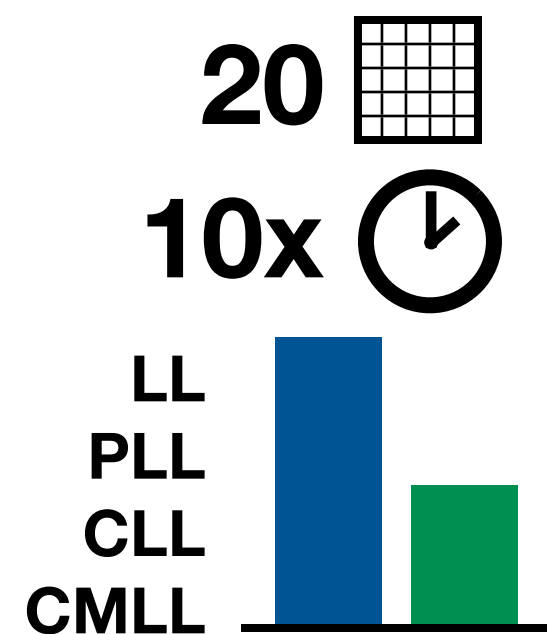
Motivation



SPN Review

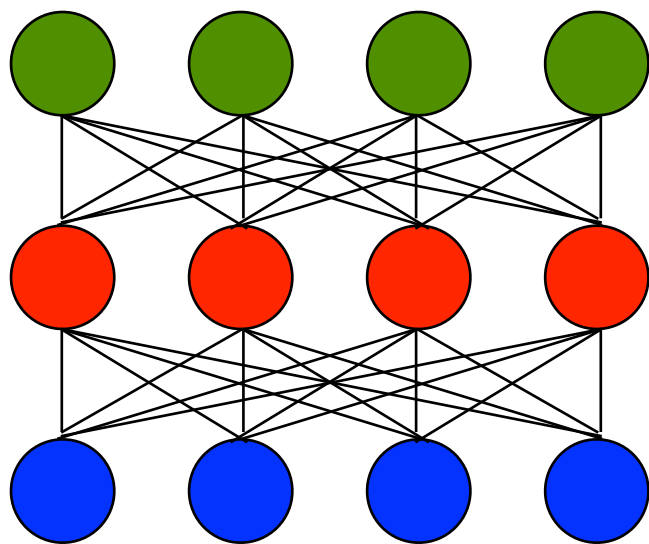


Structure Learning



Experiments

Graphical Models



Representation

Compact and expressive
Global independence

Inference

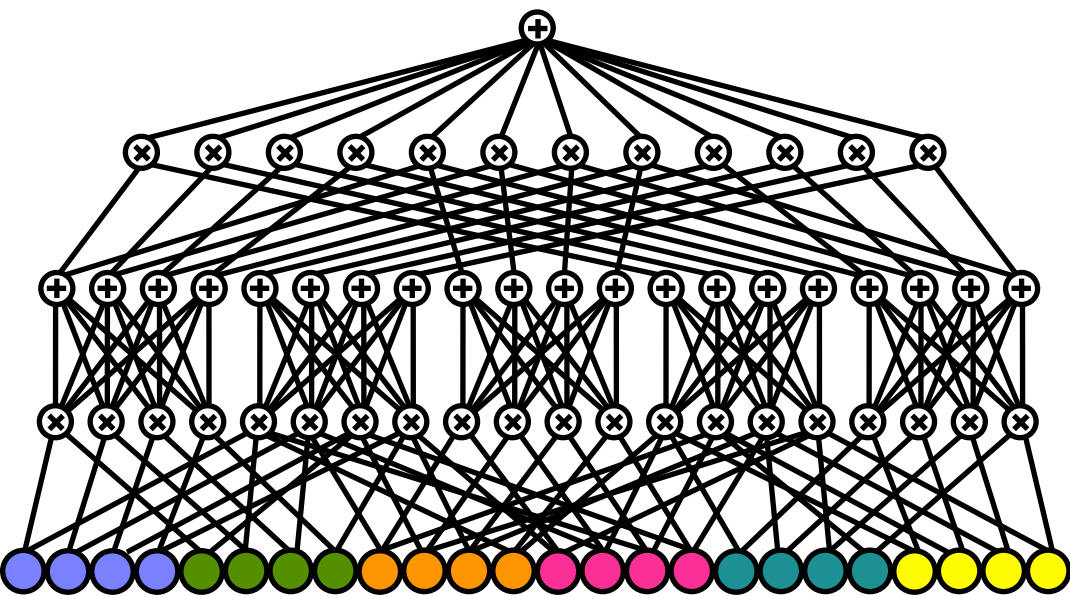
Exponential in treewidth of graph

Learning

Extremely difficult because:

- Learning requires inference
- Approximate inference is unreliable
- Hidden variables → no global optimum

Sum-Product Networks



Representation

Compact and expressive
Local independence

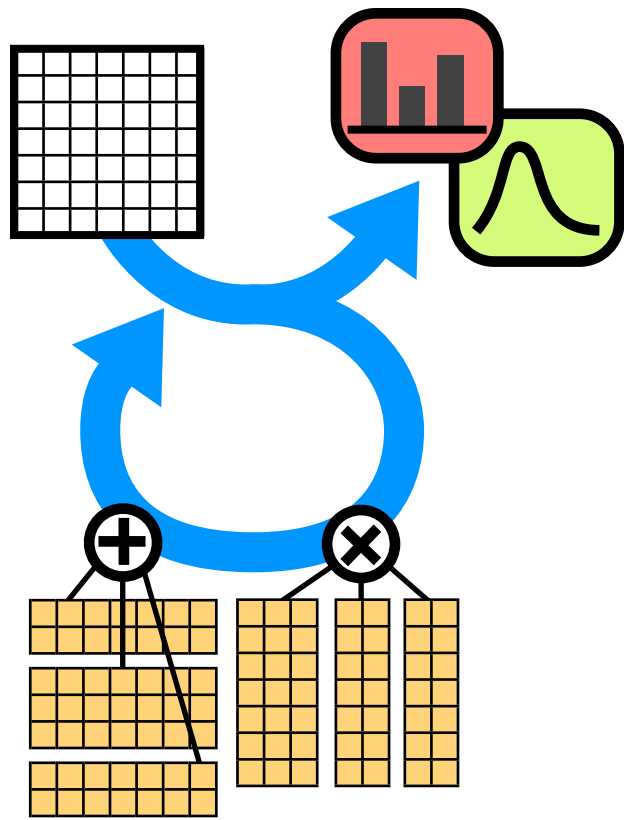
Inference

Linear in number of edges

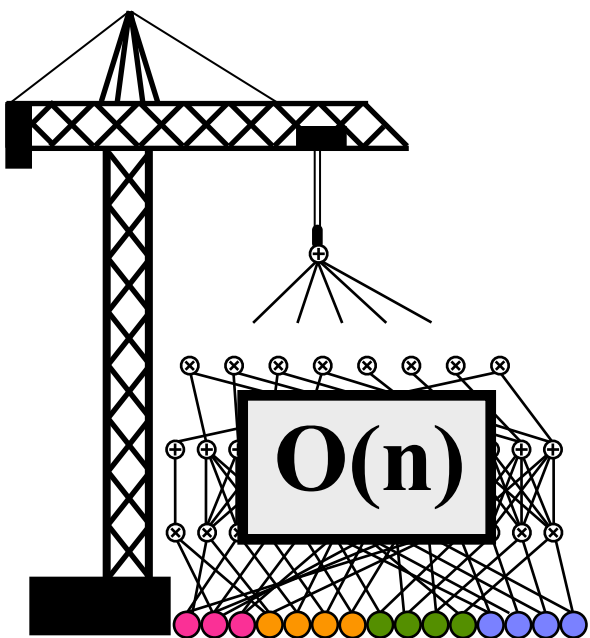
Learning

Much easier because exact
inference
Only weight learning to date

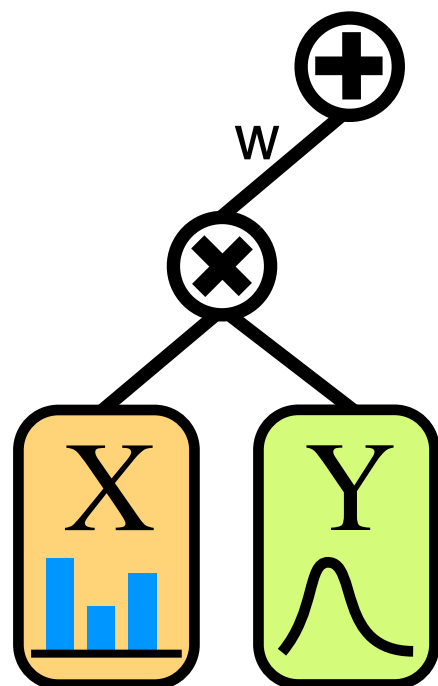
This Paper: SPN Structure Learning



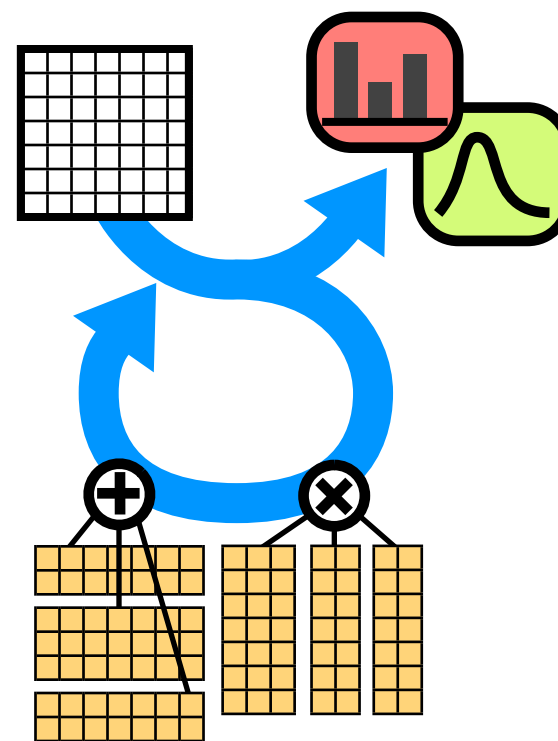
- **General-purpose SPN structure learning**
 - Discrete or continuous
 - Learns layers of hidden variables
 - Fully leverages context-specific independence
- **Simple and intuitive**
- **1-3 orders of magnitude faster, and more accurate at query time**



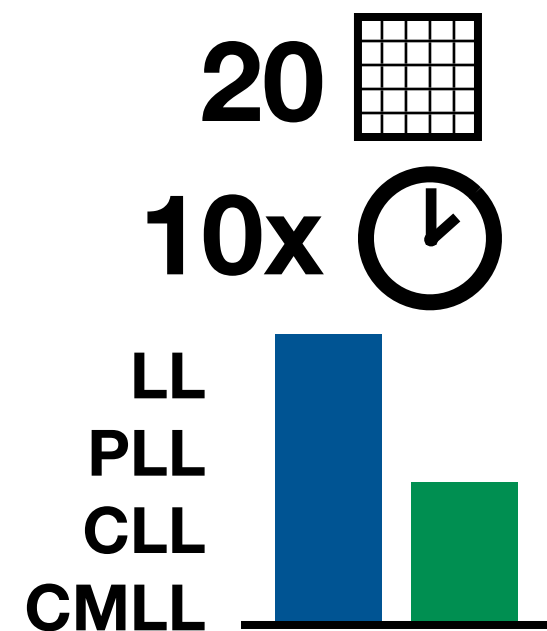
Motivation



SPN Review



Structure Learning

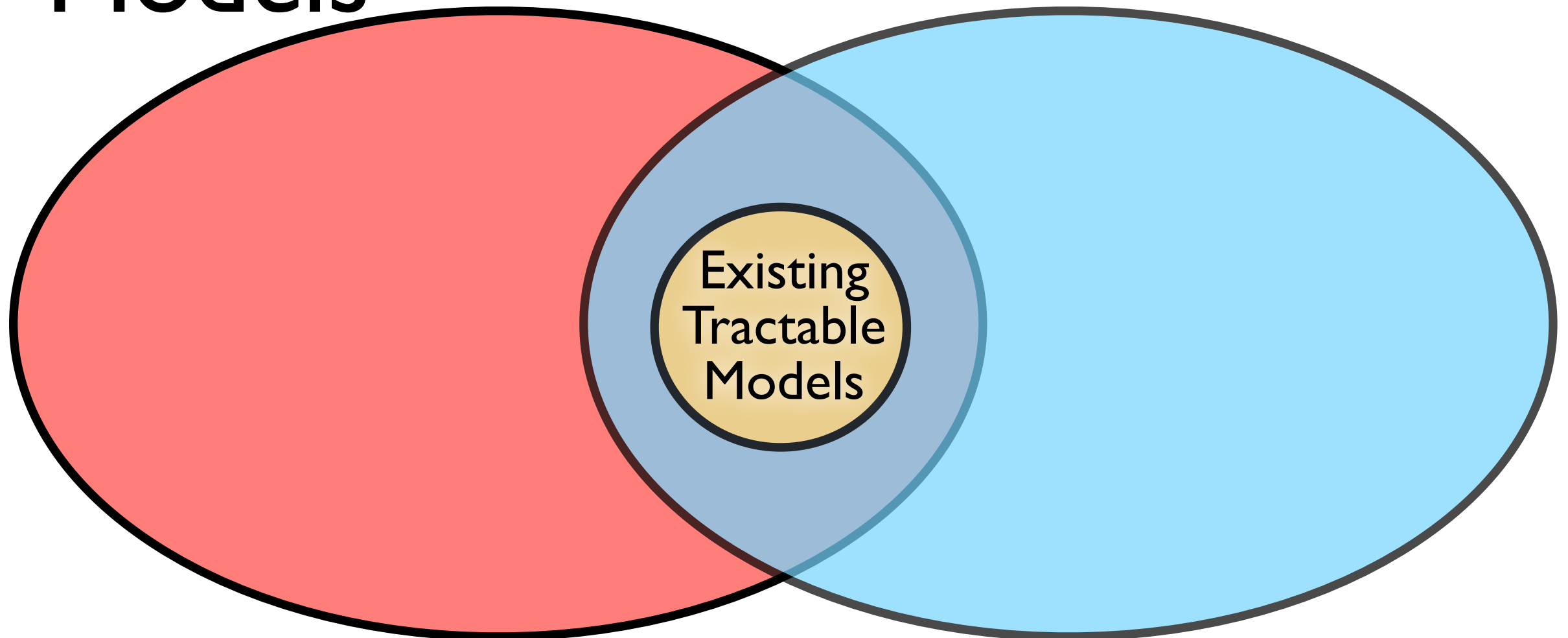


Experiments

Compactly Representable Probability Distributions

Graphical
Models

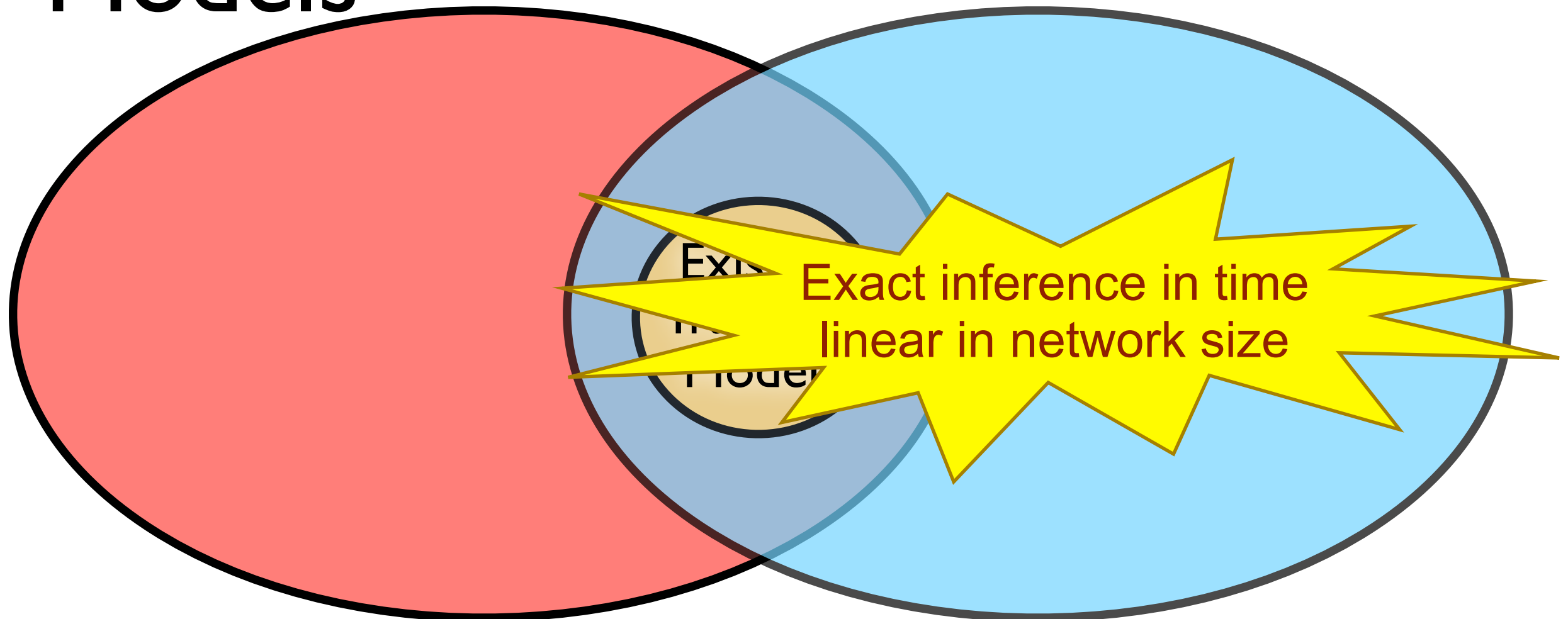
Sum-Product
Networks



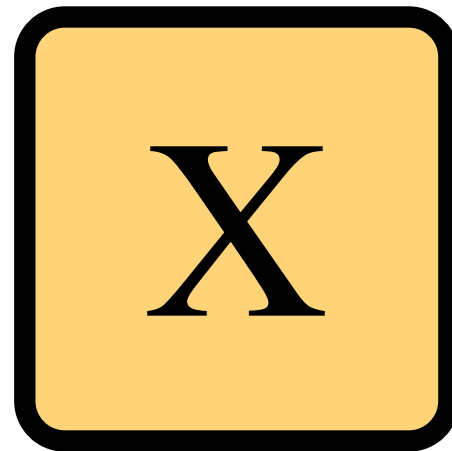
Compactly Representable Probability Distributions

Graphical
Models

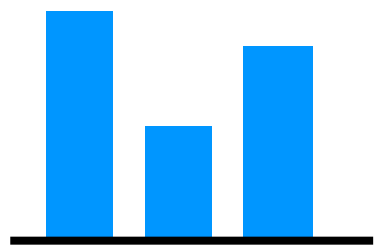
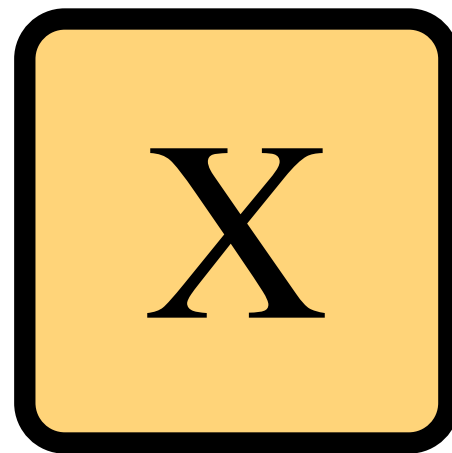
Sum-Product
Networks



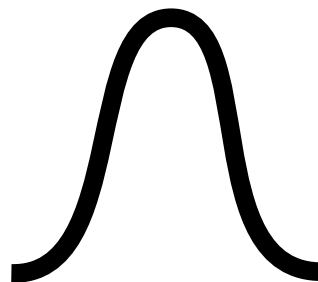
A Univariate Distribution
Is an SPN.



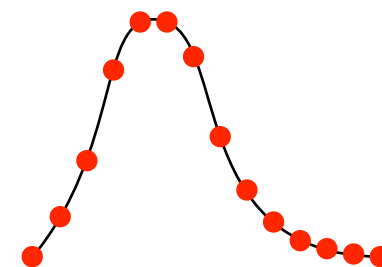
A Univariate Distribution Is an SPN.



Multinomial



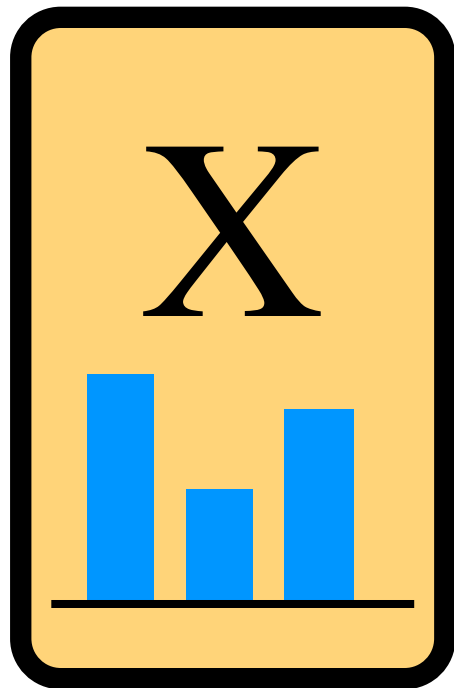
Gaussian

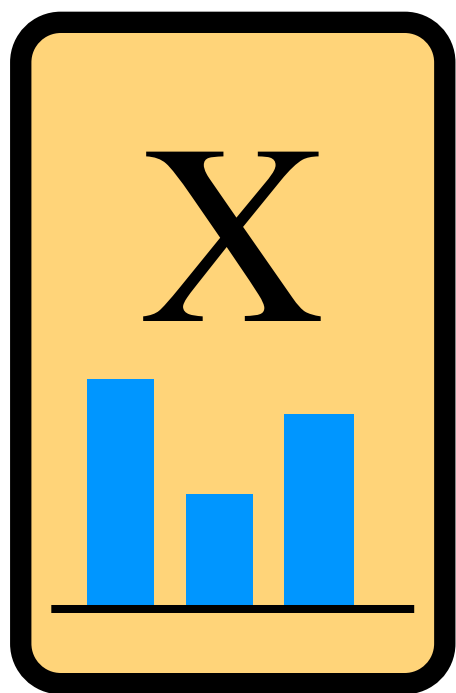


Poisson

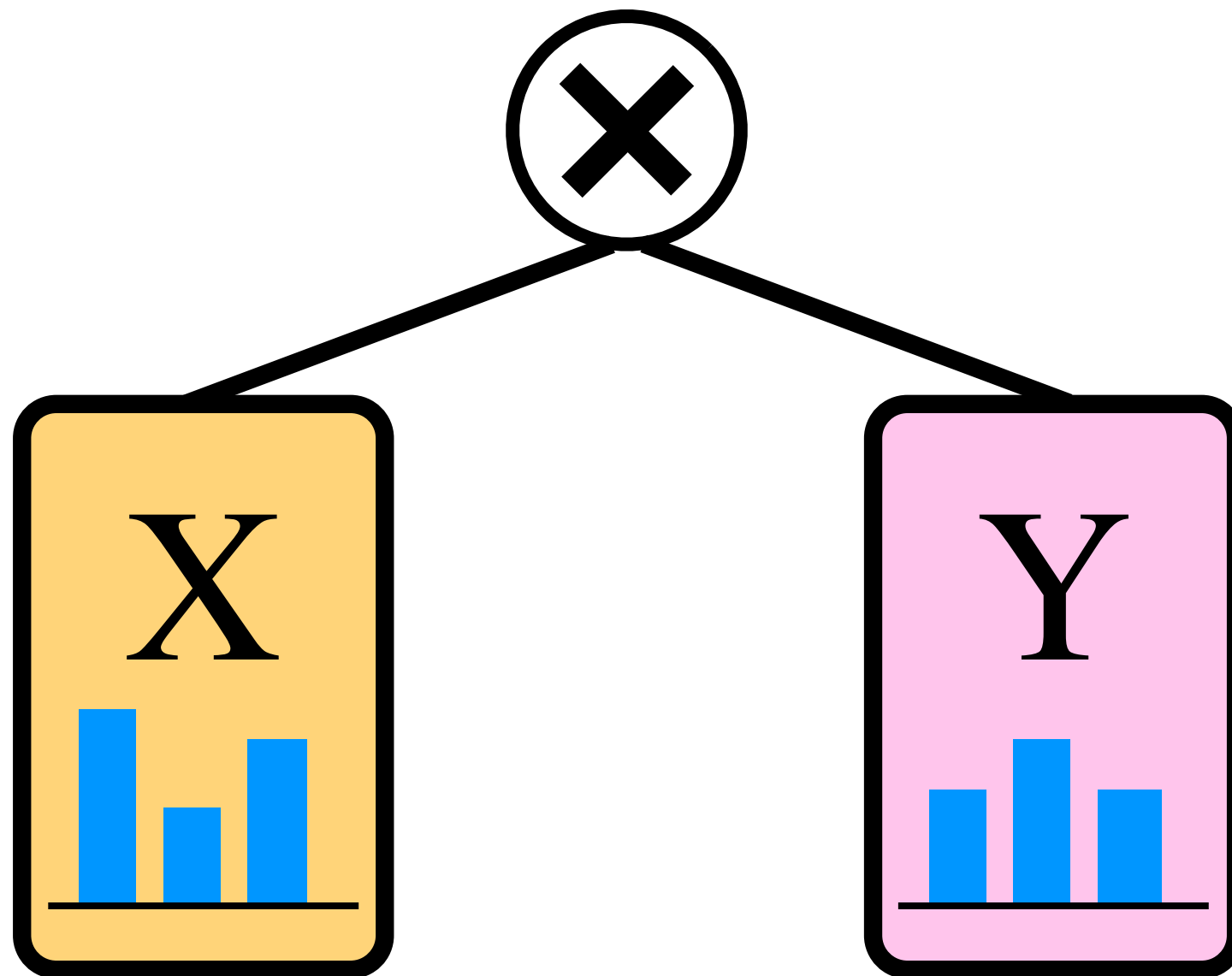
...

A Univariate Distribution Is an SPN.

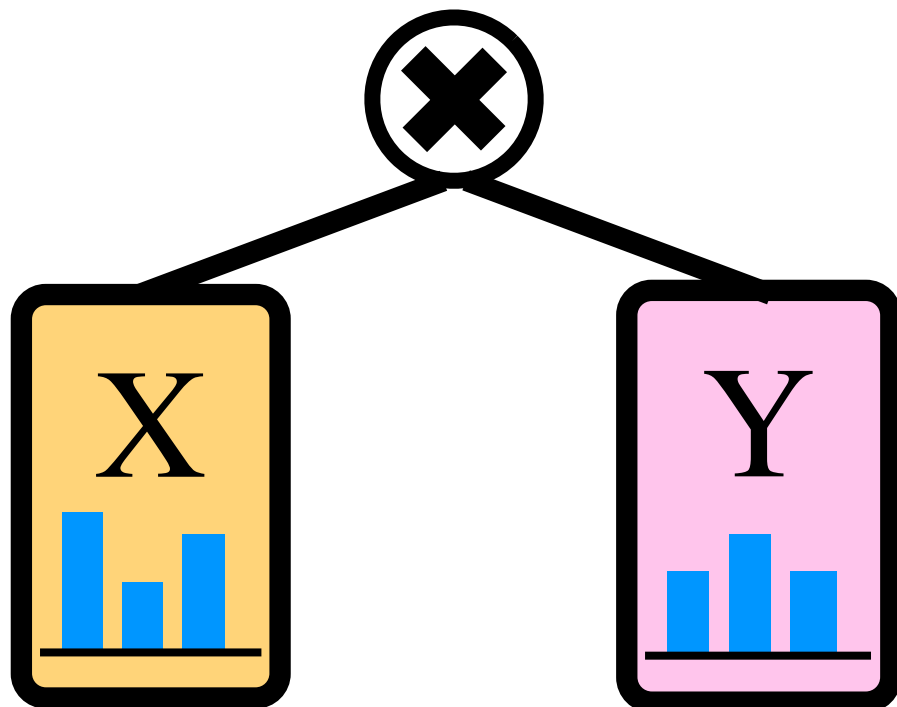


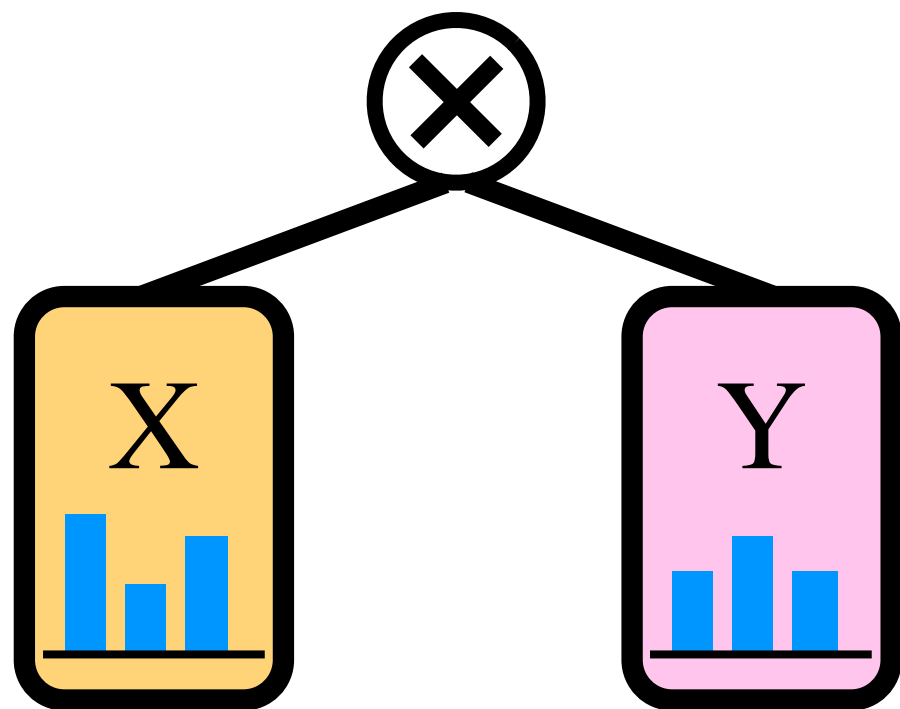


A Product of SPNs over
Disjoint Variables
Is an SPN.

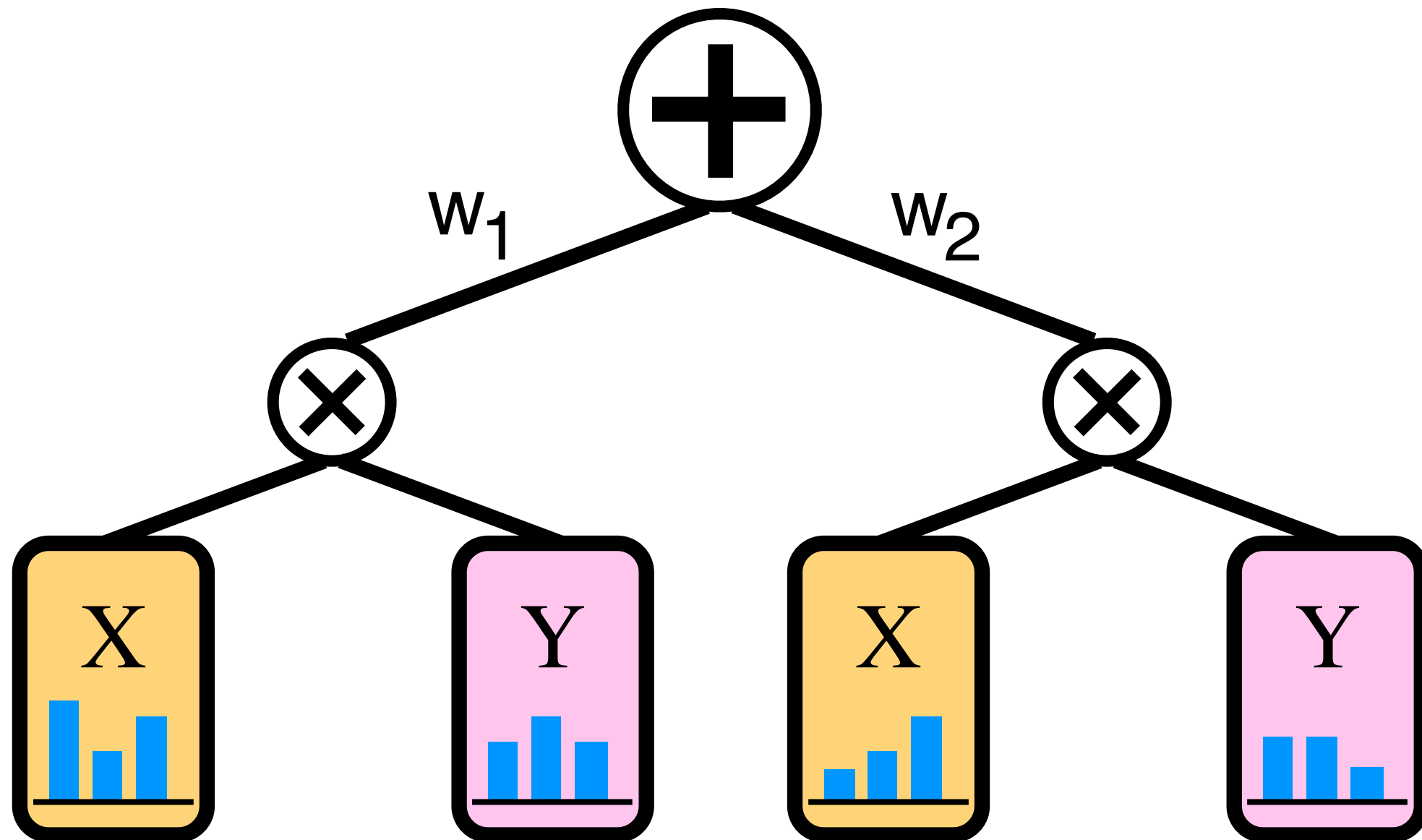


A Product of SPNs over
Disjoint Variables
Is an SPN.

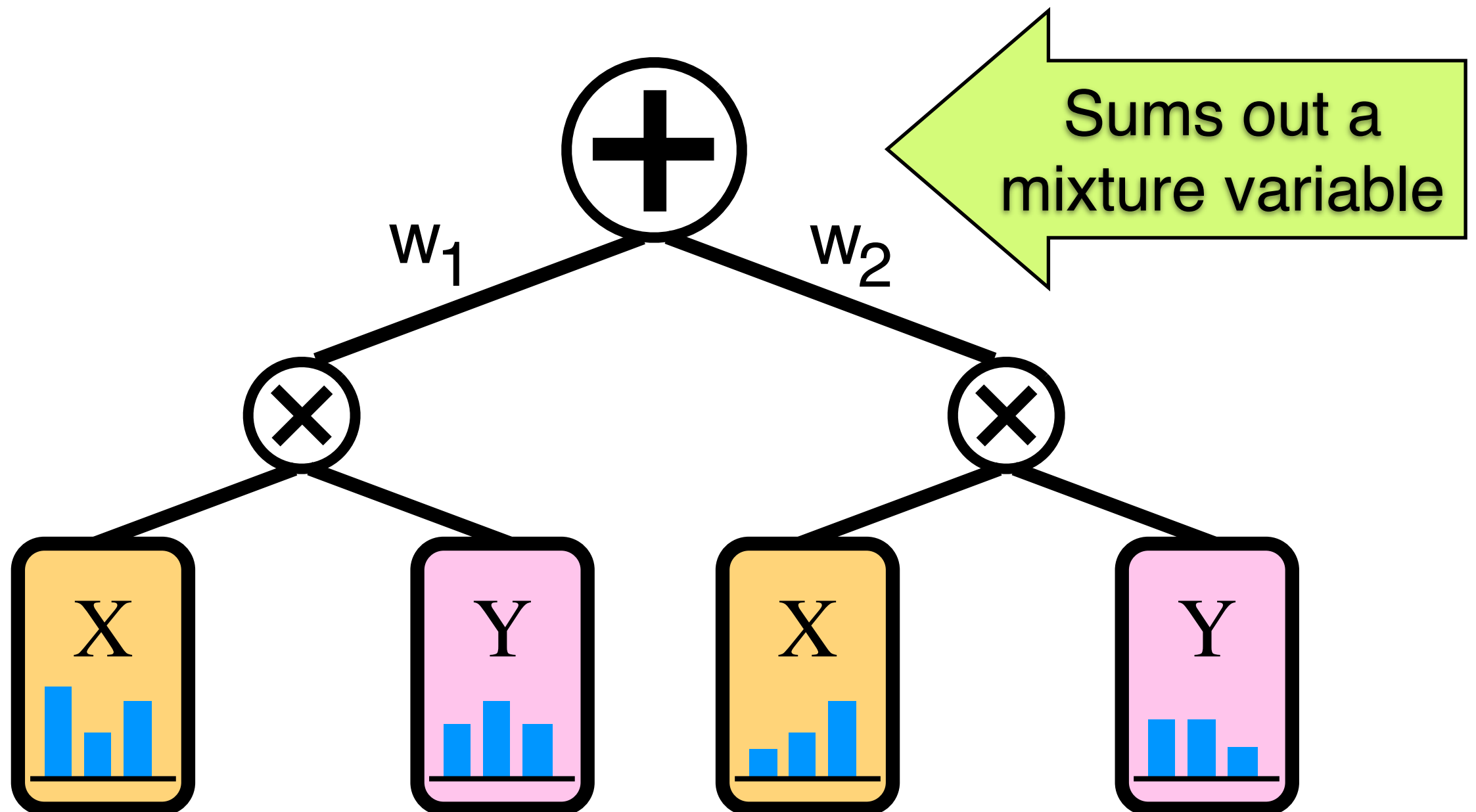


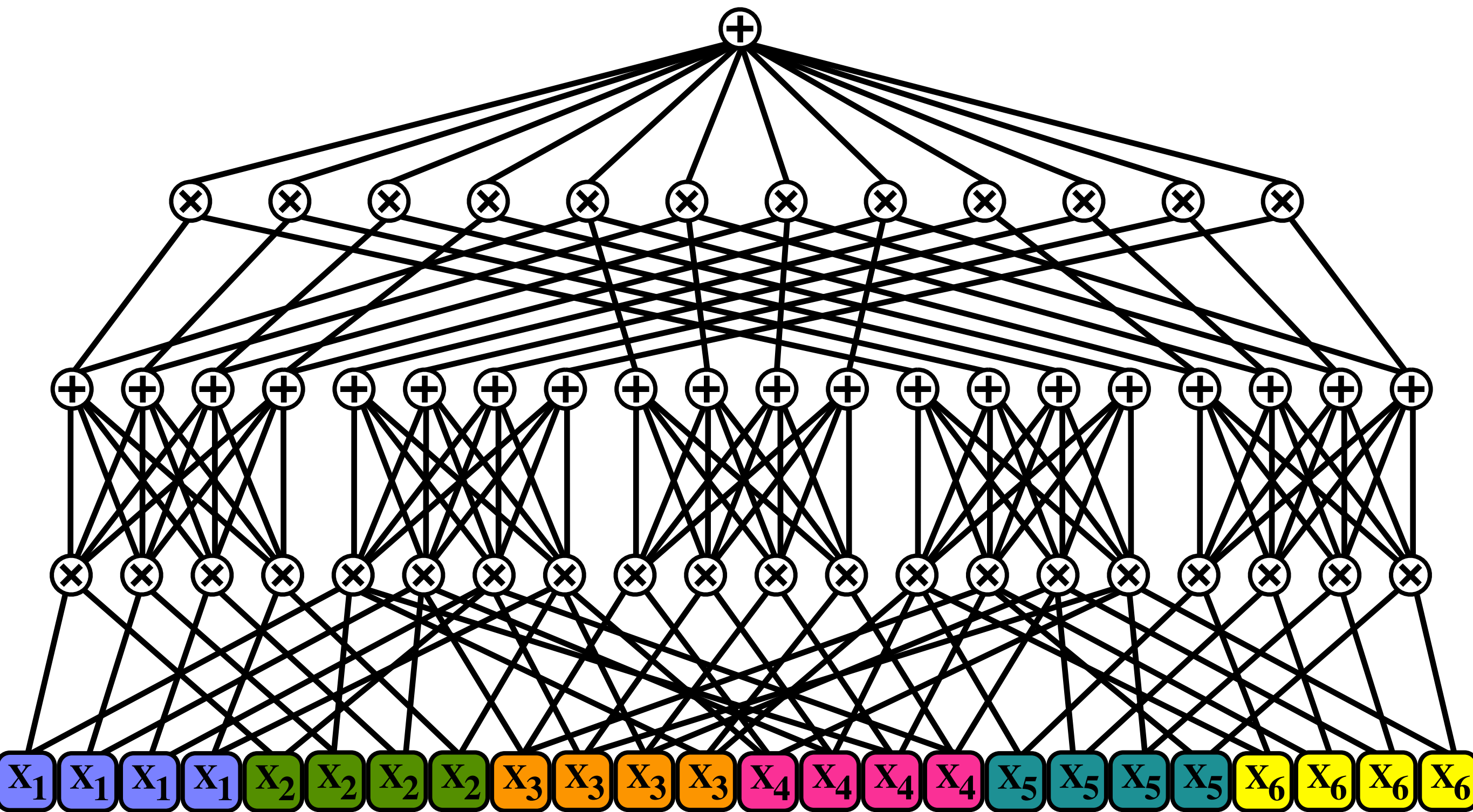


A Weighted Sum of SPNs
over the Same Variables
Is an SPN.



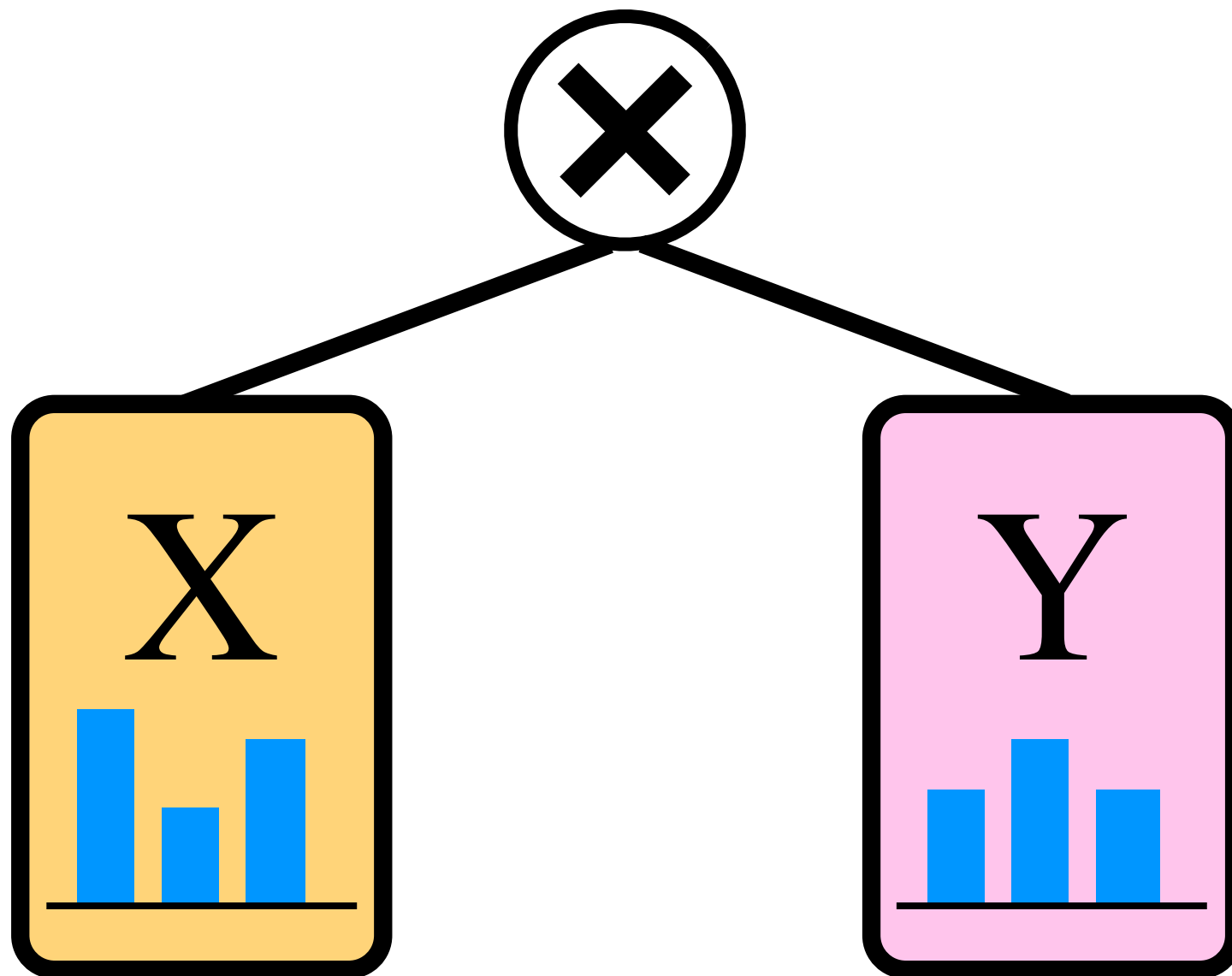
A Weighted Sum of SPNs over the Same Variables Is an SPN.



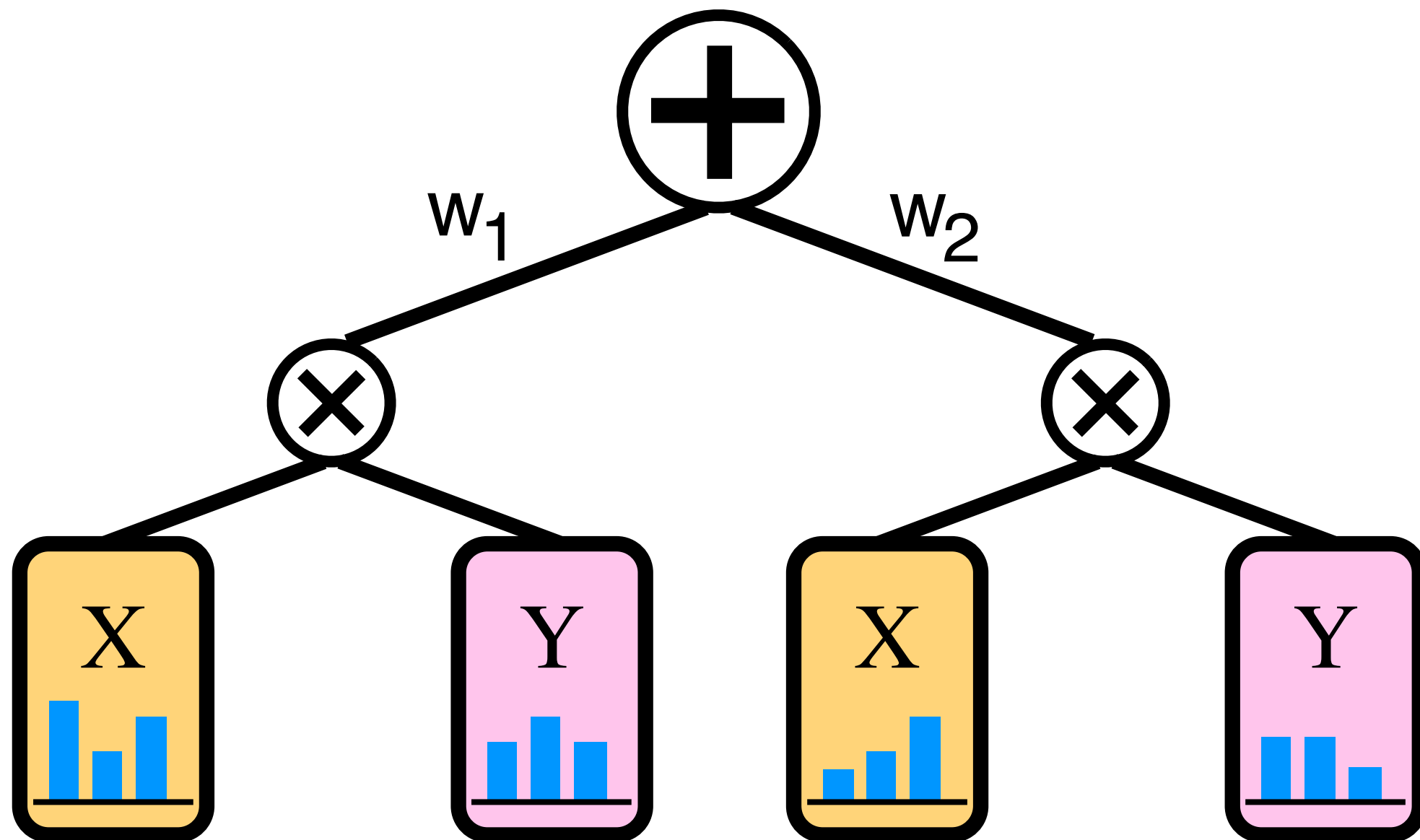


All Marginals Are
Computable in Linear Time

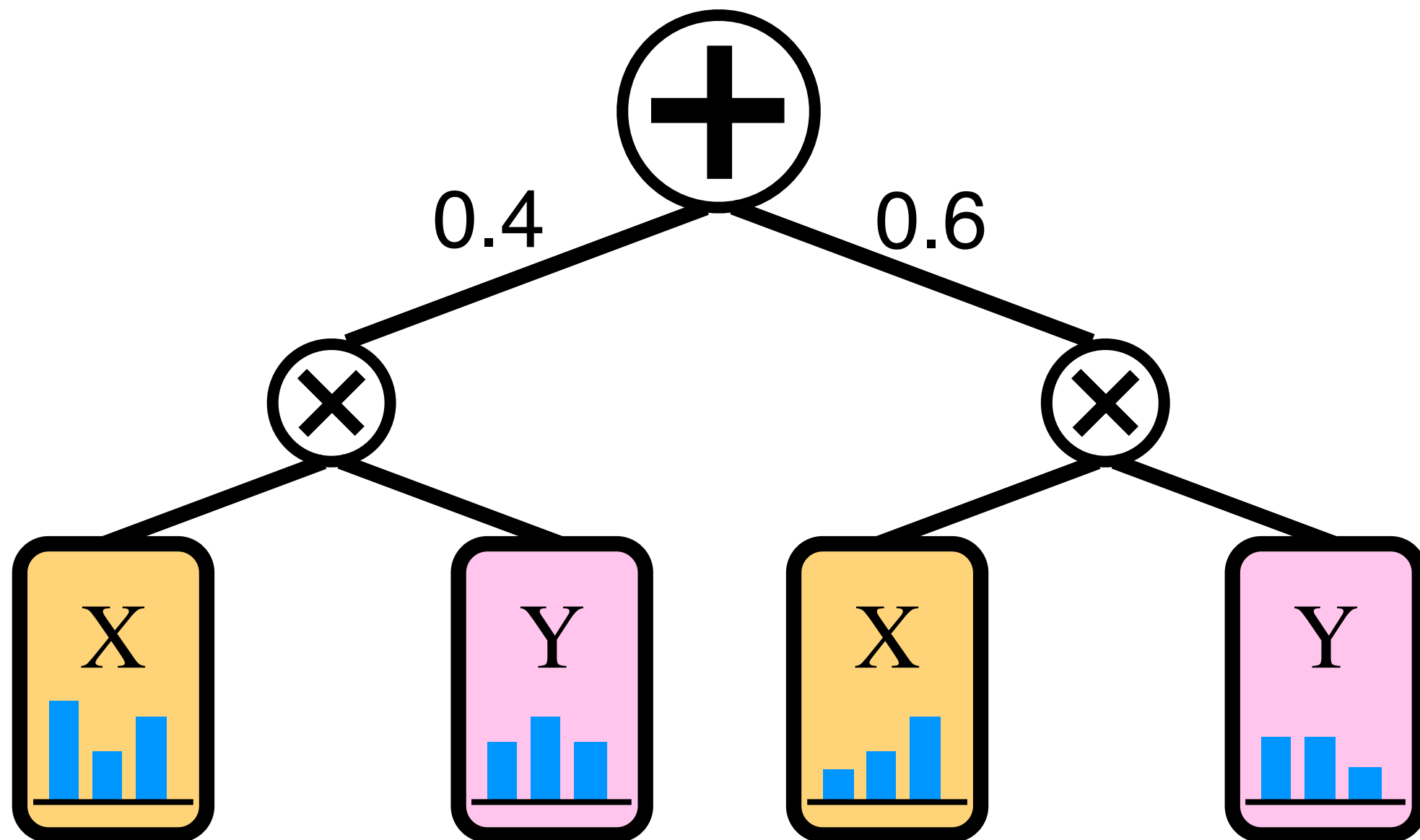
All Marginals Are Computable in Linear Time



All Marginals Are Computable in Linear Time

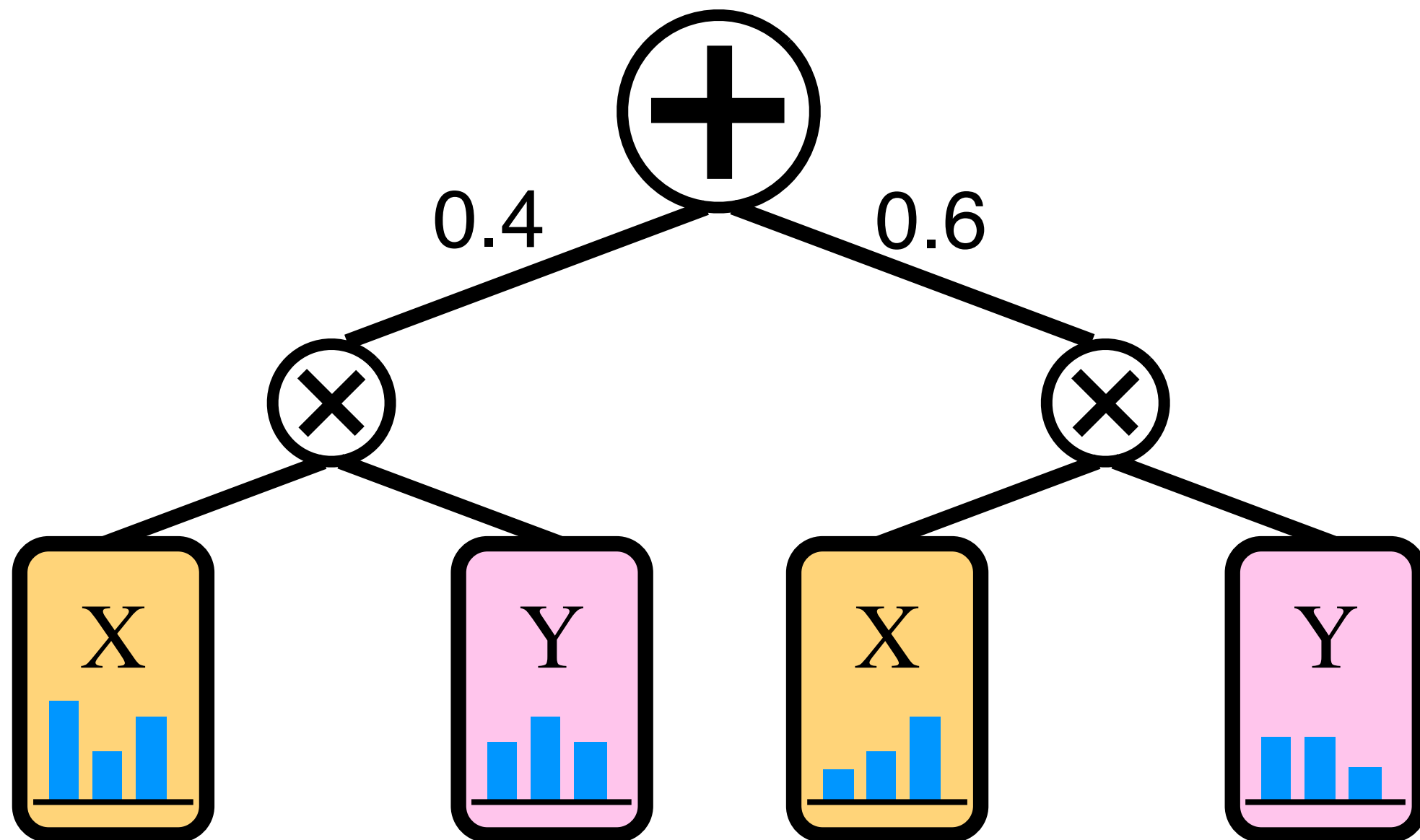


All Marginals Are Computable in Linear Time



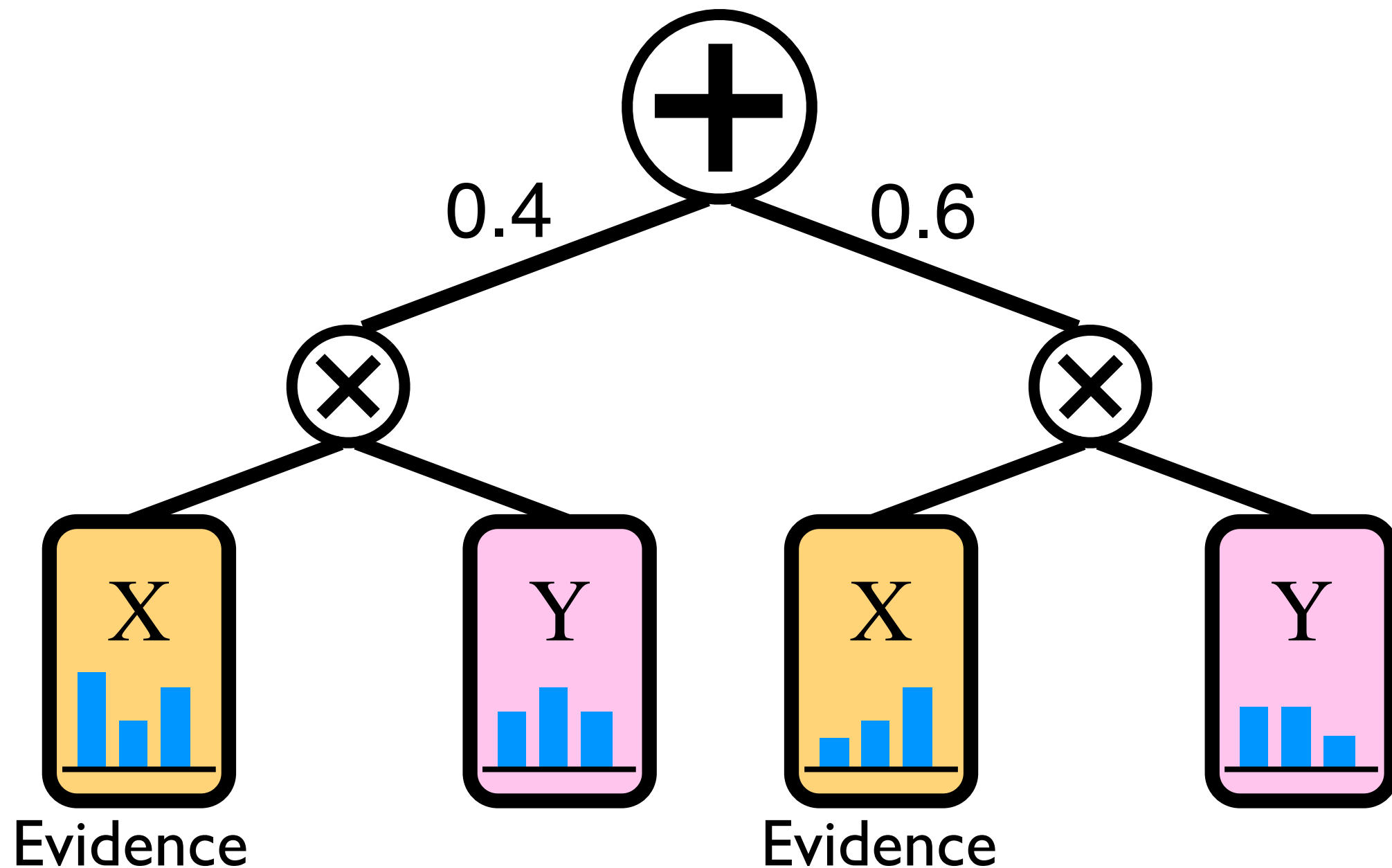
All Marginals Are Computable in Linear Time

$$P(X=0) \text{ ?}$$



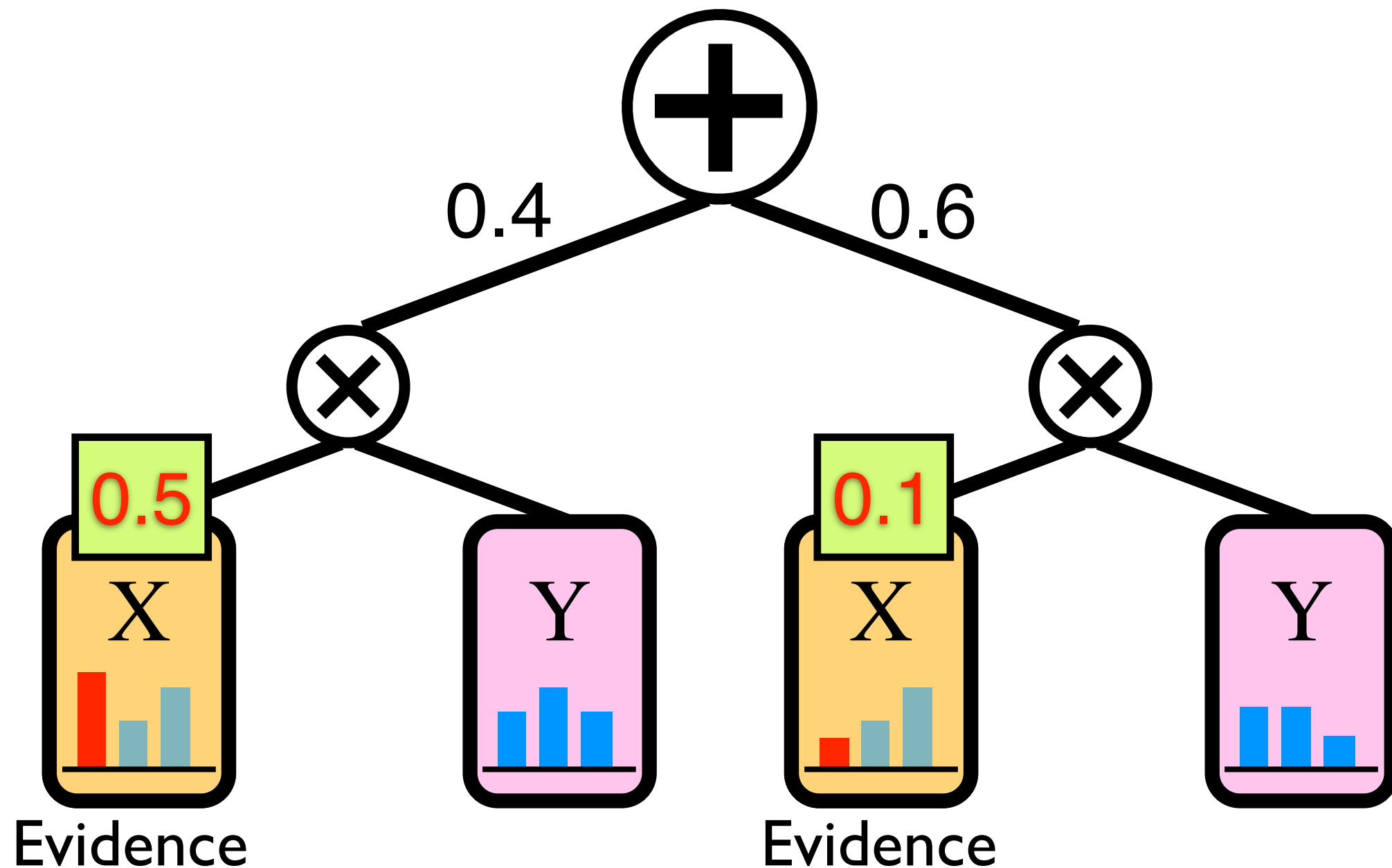
All Marginals Are Computable in Linear Time

$$P(X=0) \text{ ?}$$



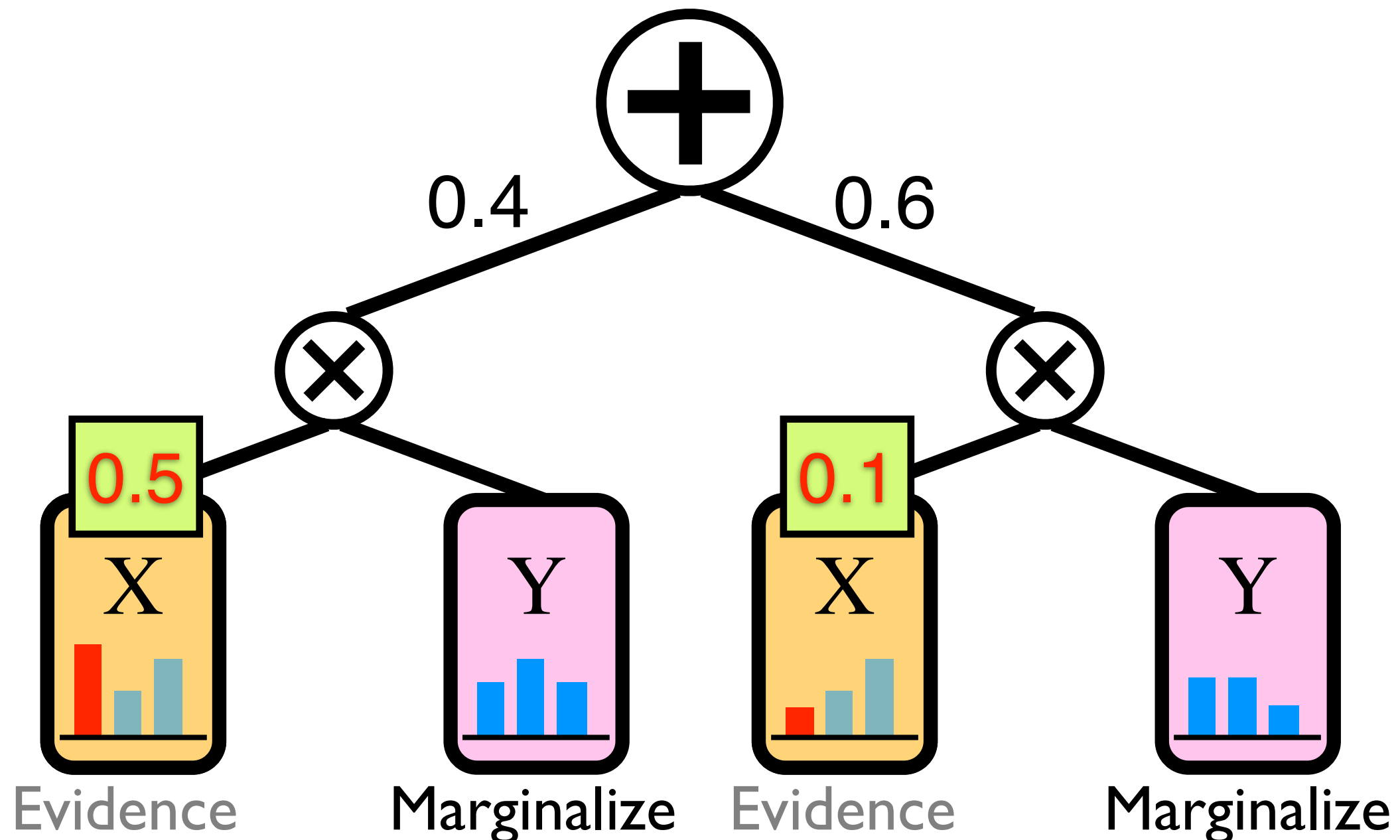
All Marginals Are Computable in Linear Time

$$P(X=0) \text{ ?}$$



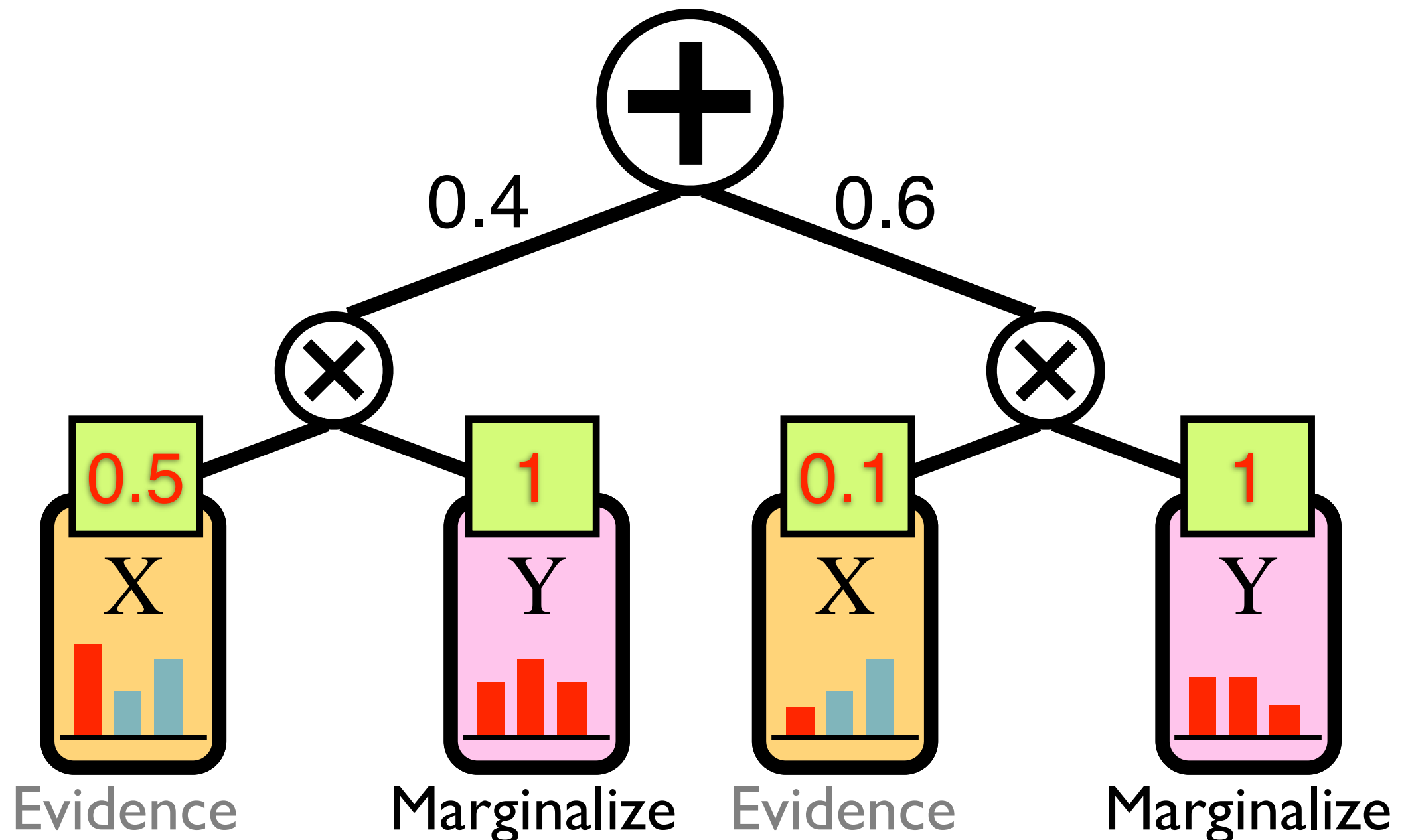
All Marginals Are Computable in Linear Time

$$P(X=0) \text{ ?}$$



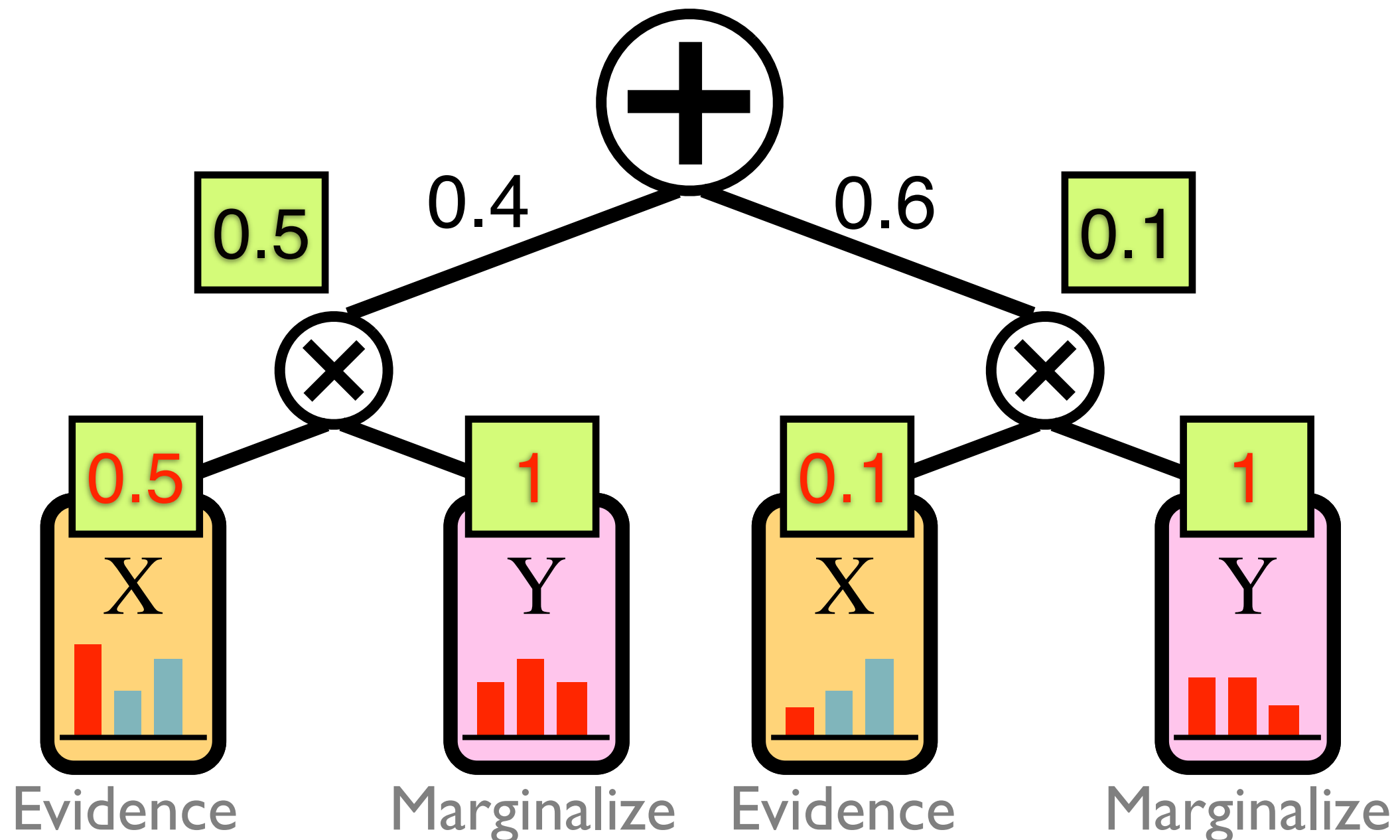
All Marginals Are Computable in Linear Time

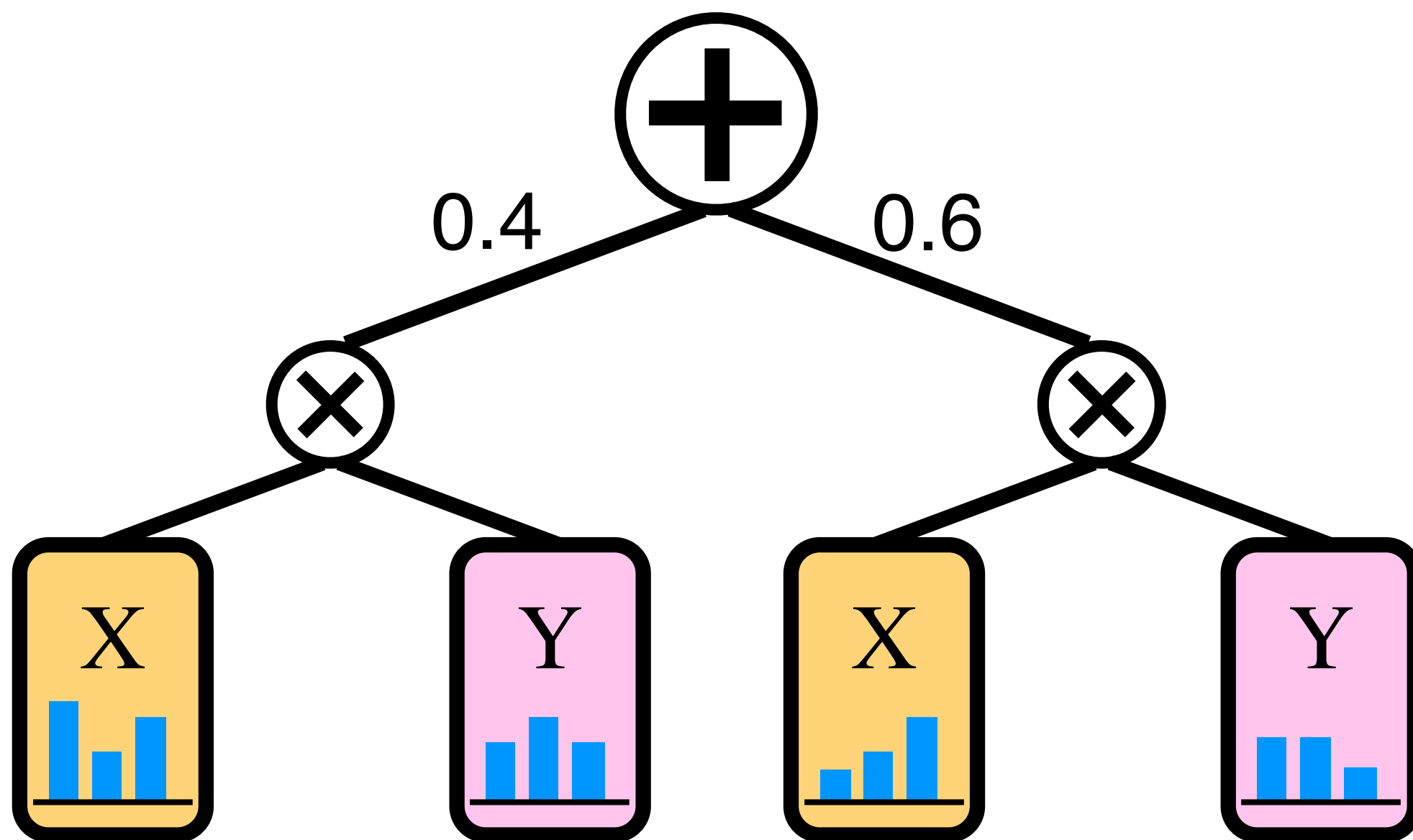
$$P(X=0) ?$$



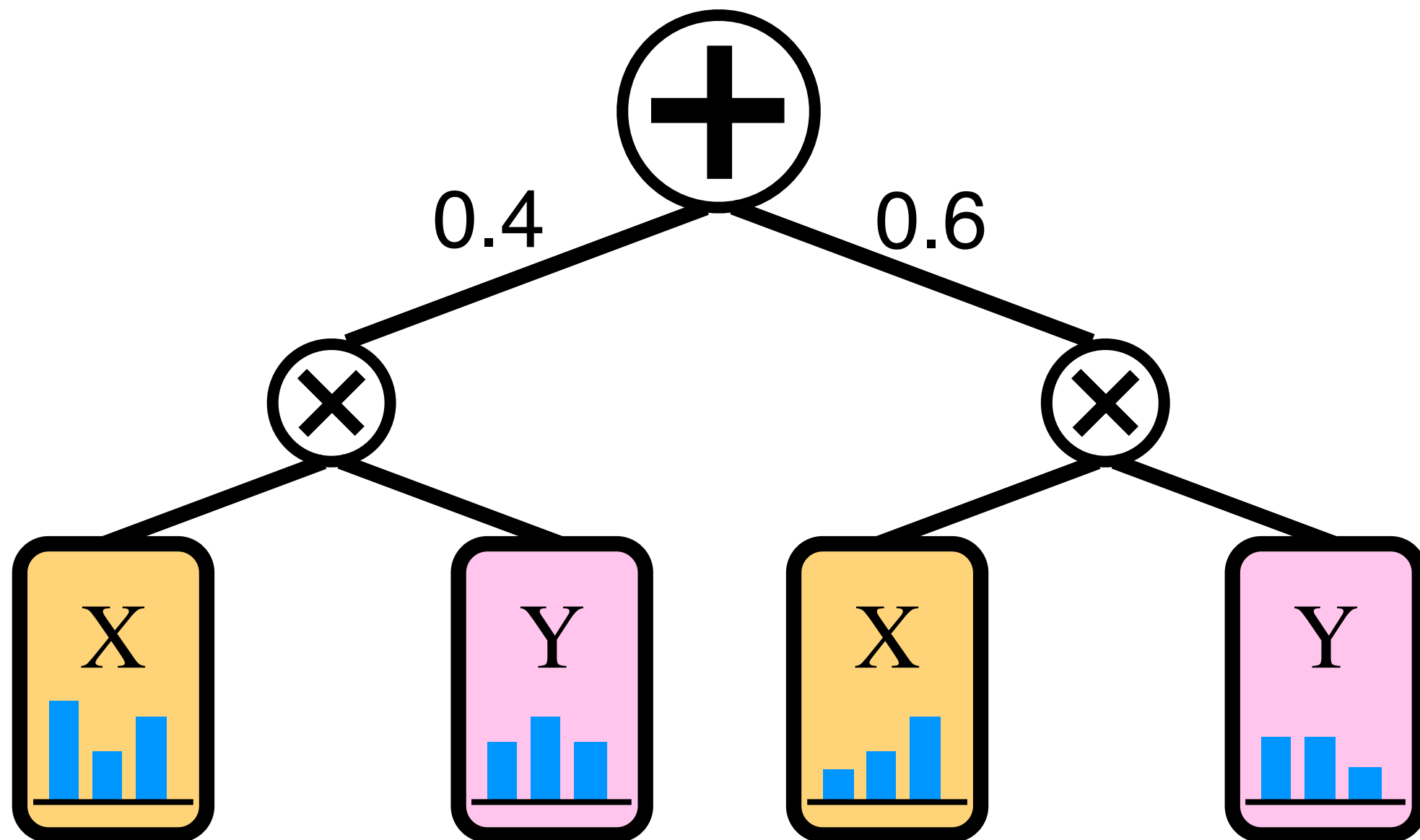
All Marginals Are Computable in Linear Time

$$P(X=0) = 0.26$$

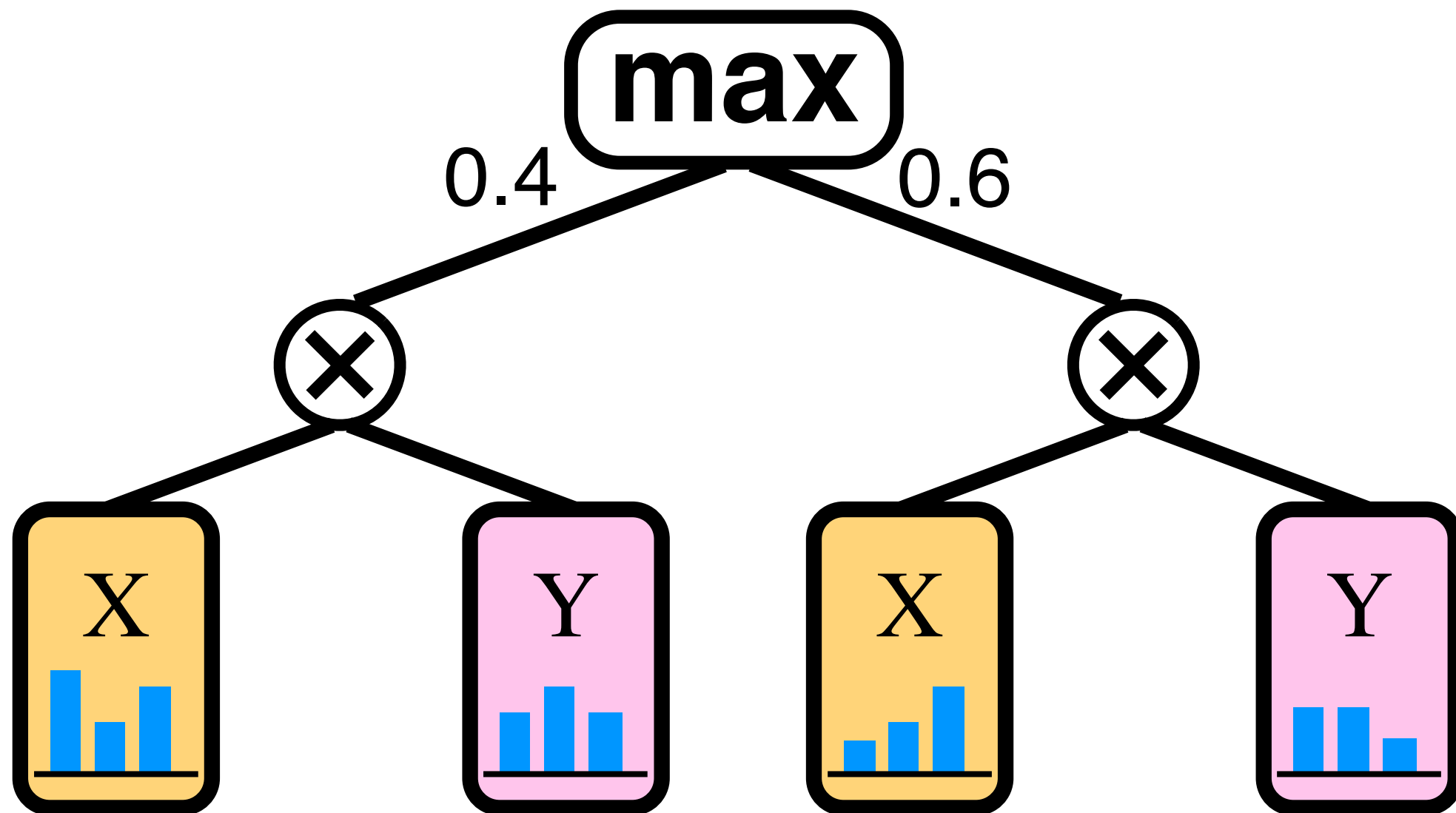




All MAP States Are Computable in Linear Time

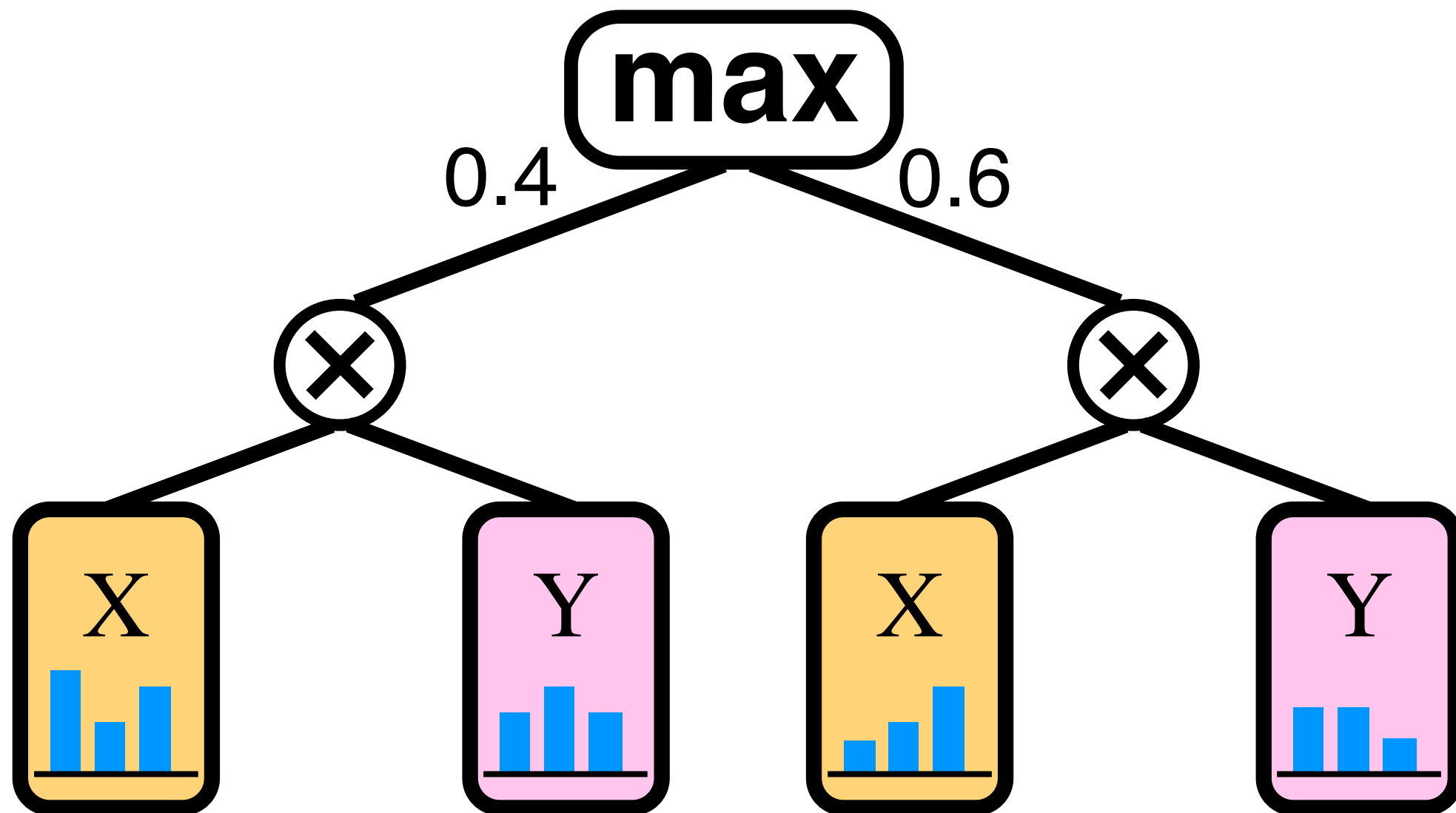


All MAP States Are Computable in Linear Time



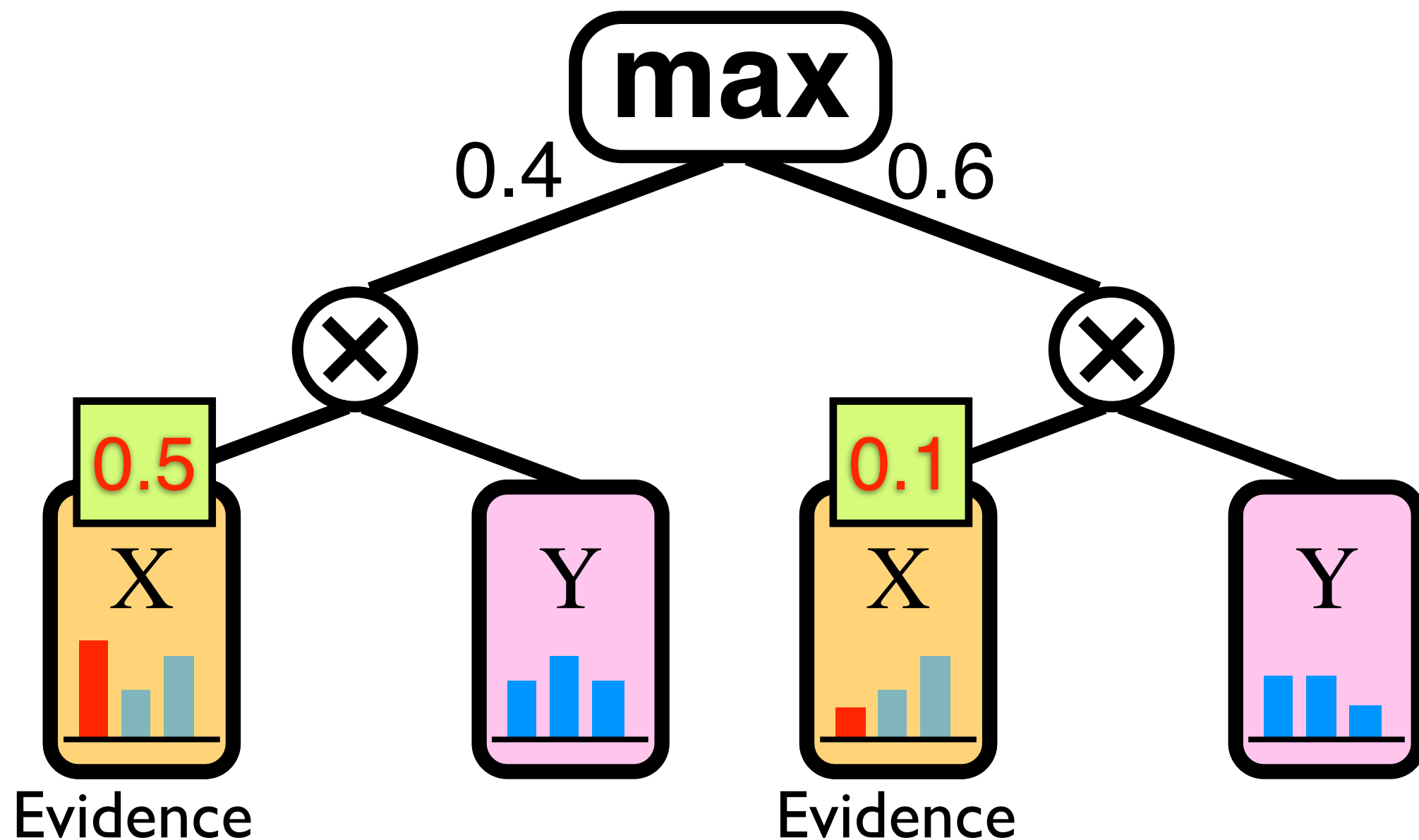
All MAP States Are Computable in Linear Time

$$\max_{y,h} P(X = 0, Y = y, H = h) ?$$



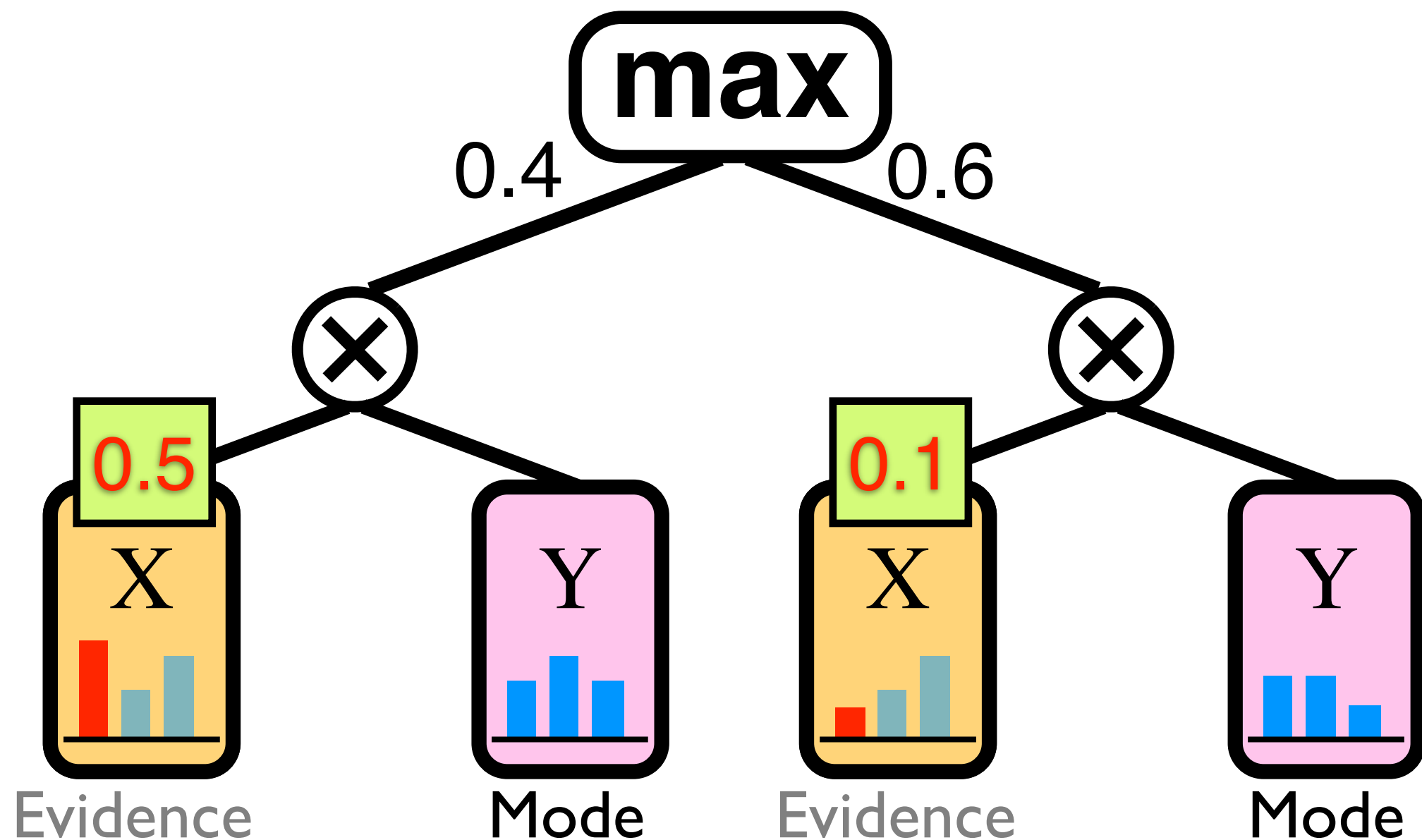
All MAP States Are Computable in Linear Time

$$\max_{y,h} P(X = 0, Y = y, H = h) ?$$



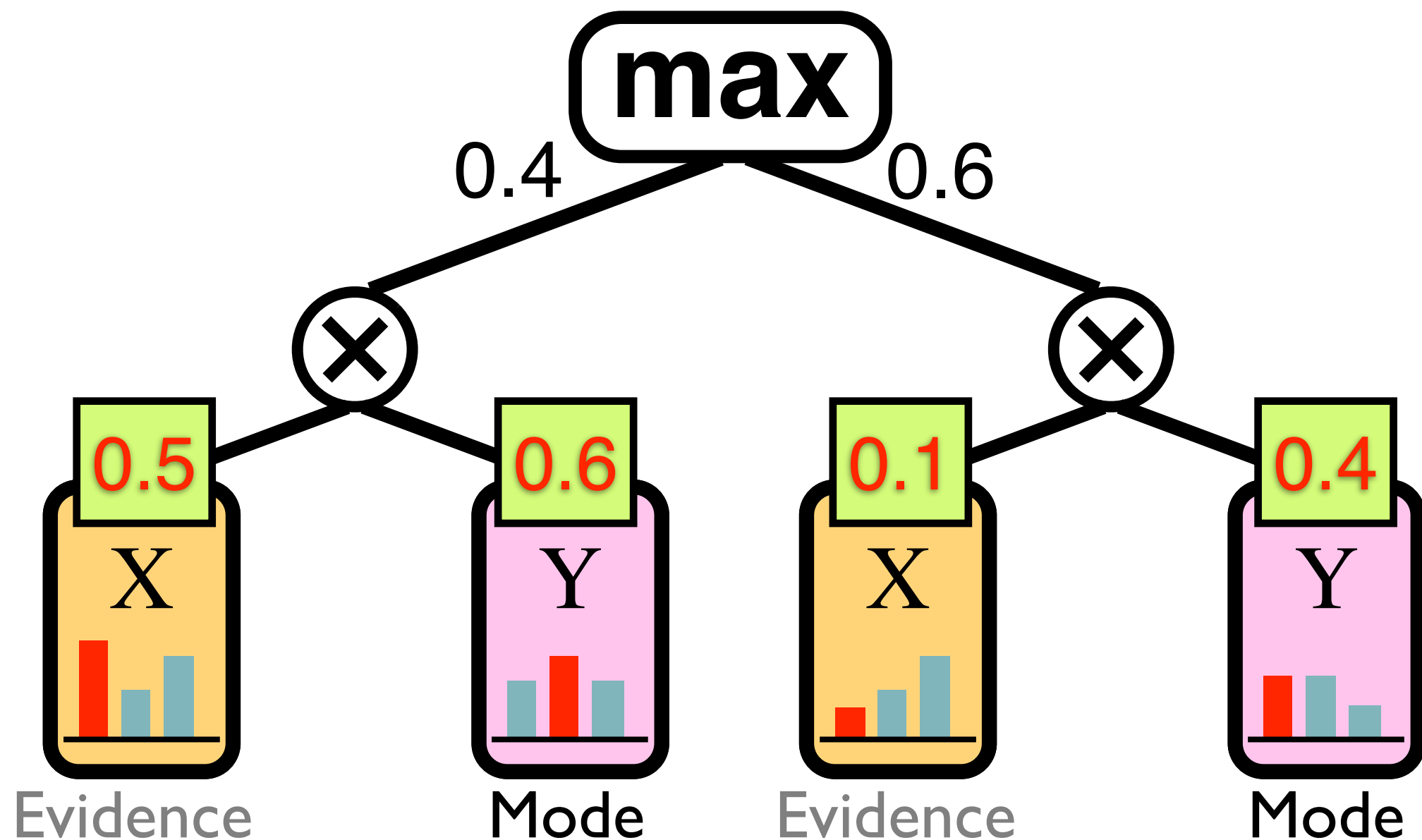
All MAP States Are Computable in Linear Time

$$\max_{y,h} P(X = 0, Y = y, H = h) ?$$



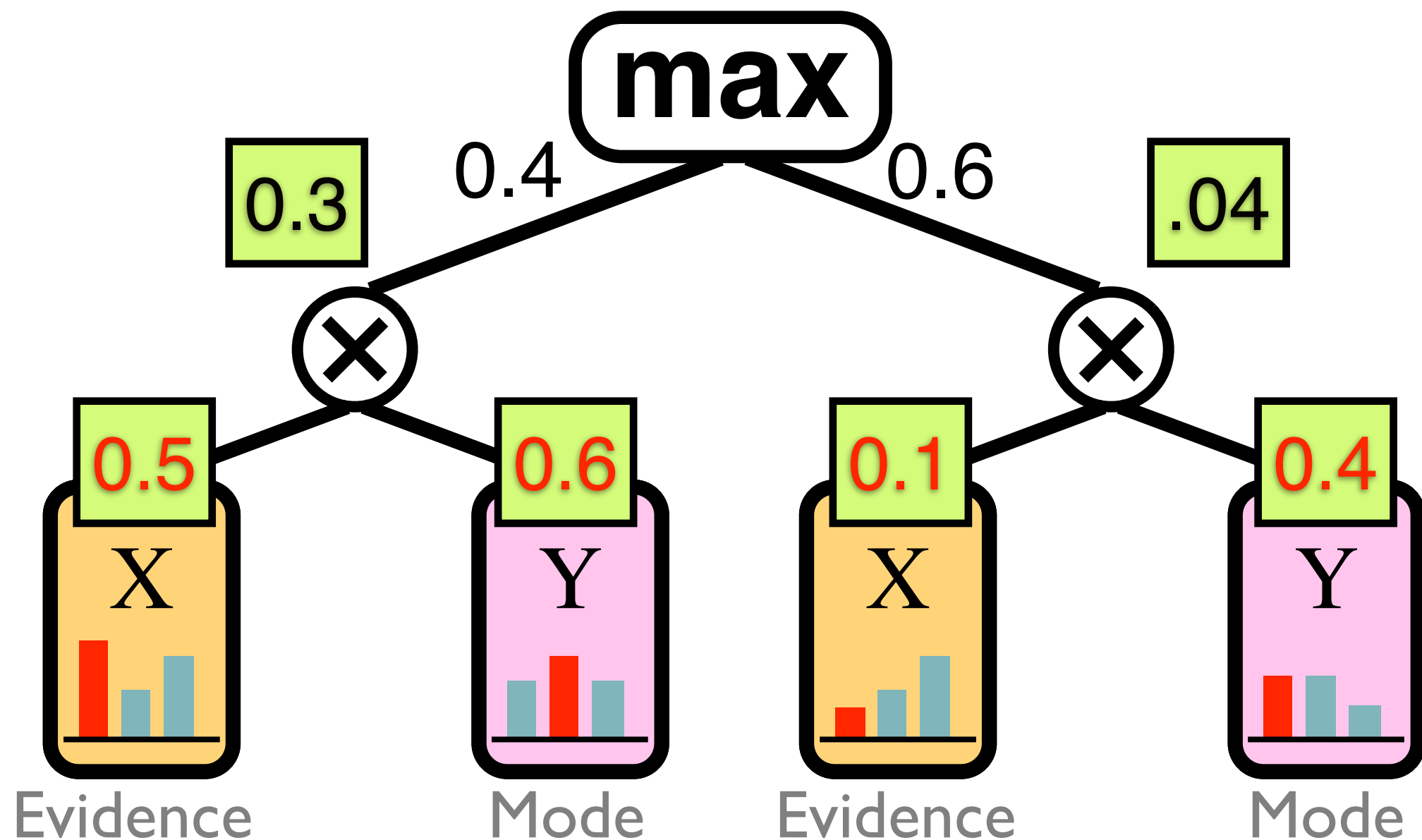
All MAP States Are Computable in Linear Time

$$\max_{y,h} P(X = 0, Y = y, H = h) ?$$



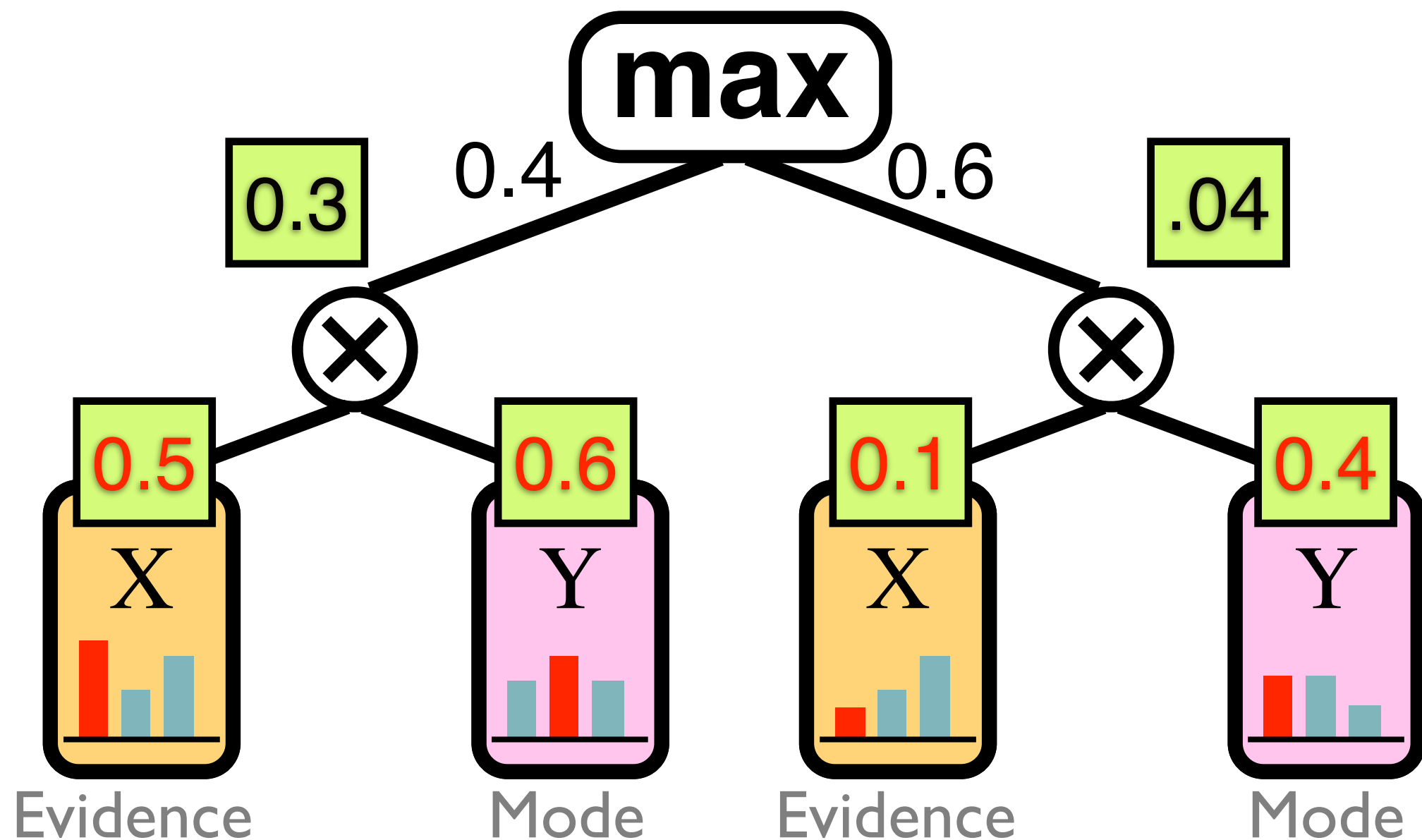
All MAP States Are Computable in Linear Time

$$\max_{y,h} P(X = 0, Y = y, H = h) = \boxed{0.12}$$



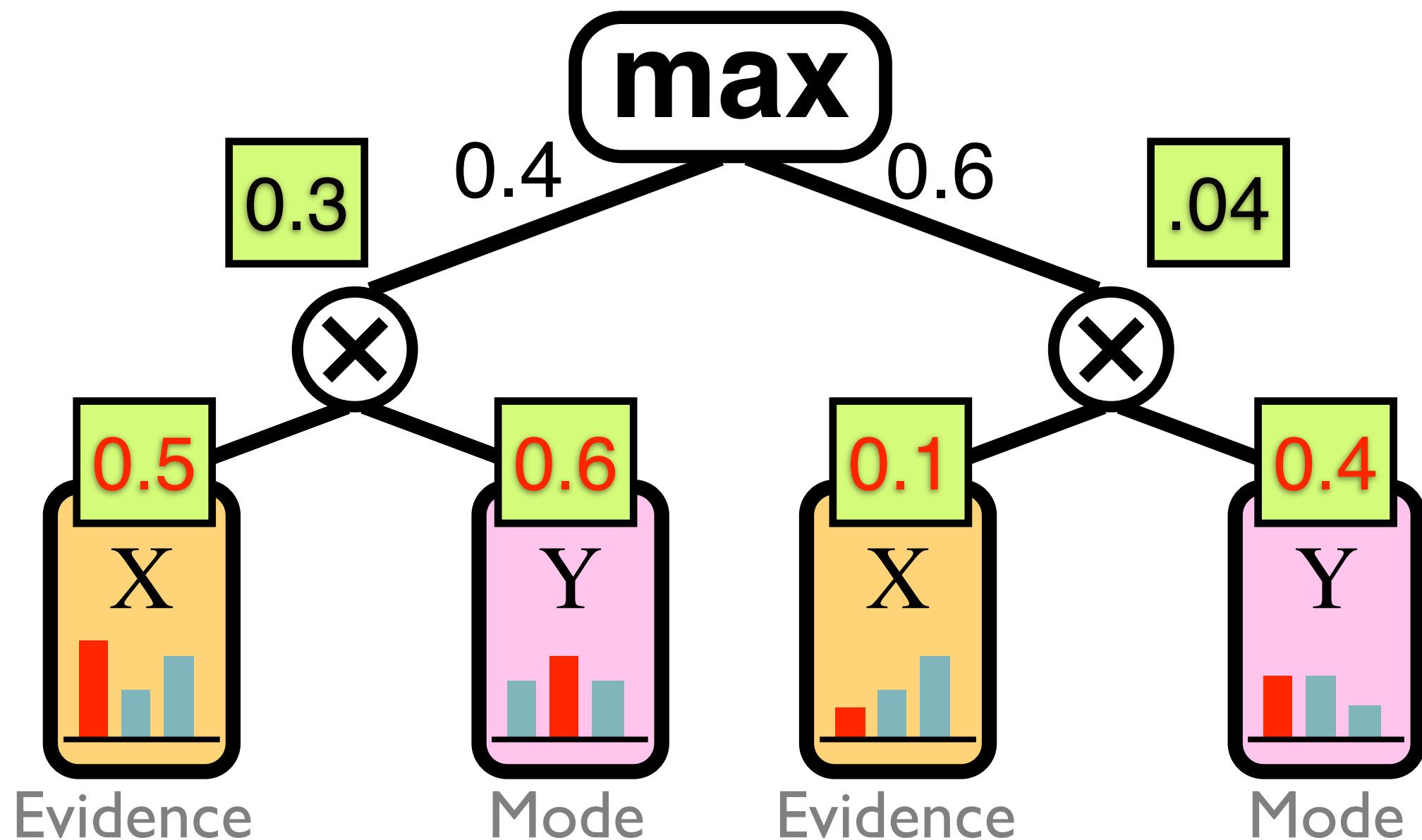
All MAP States Are Computable in Linear Time

$$\max_{y,h} P(X = 0, Y = y, H = h) = \boxed{0.12}$$



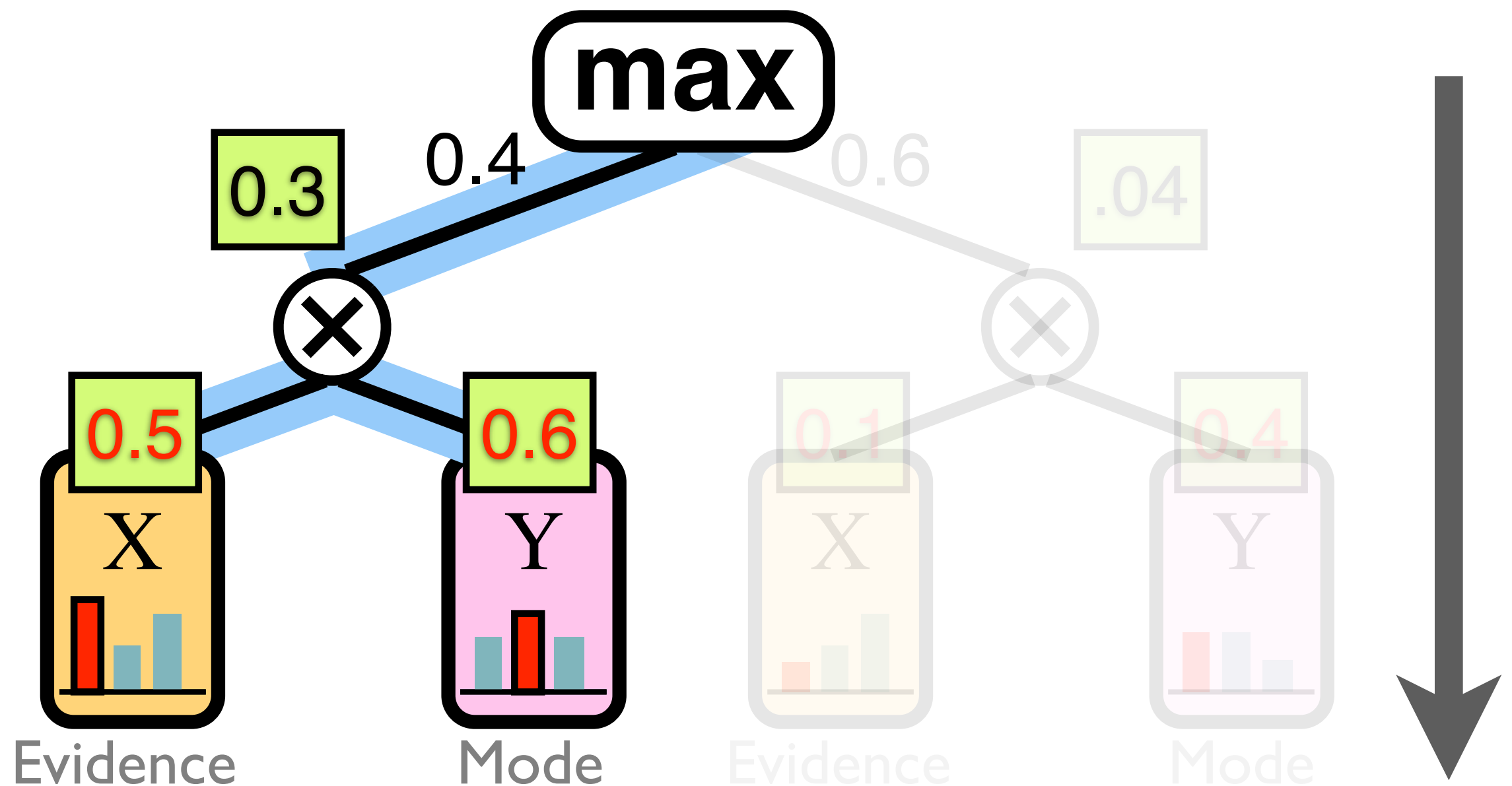
All MAP States Are Computable in Linear Time

$$\max_{y,h} P(X = 0, Y = y, H = h) = \boxed{0.12}$$



All MAP States Are Computable in Linear Time

$$\max_{y,h} P(X = 0, Y = y, H = h) = \boxed{0.12}$$



What Does an SPN Mean?

Products = Features

Sums = Clusterings



What Does an SPN Mean?

Products = Features
Sums = Clusterings



Hay



Gravel



Greenery



Truck



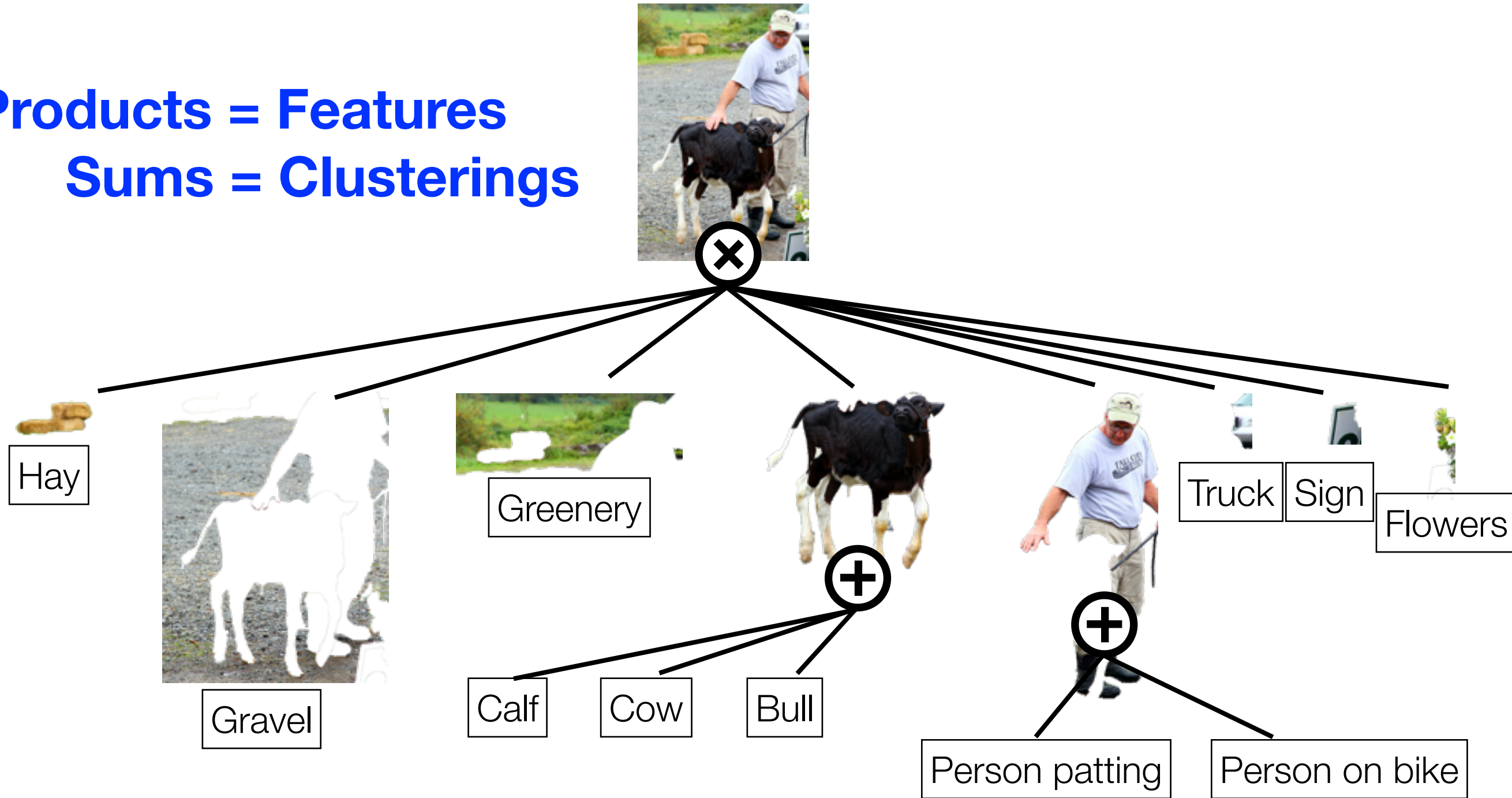
Sign



Flowers

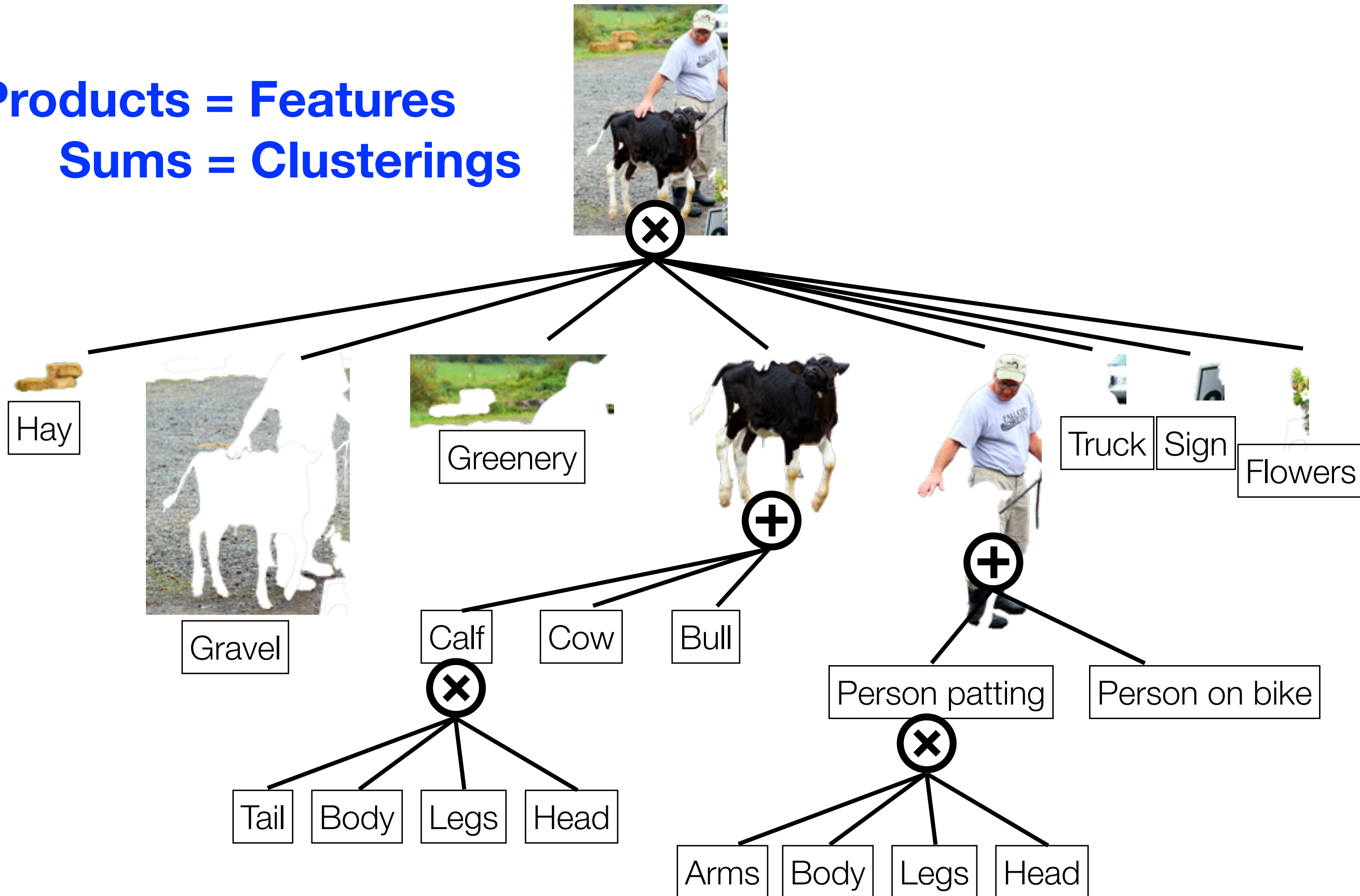
What Does an SPN Mean?

Products = Features
Sums = Clusterings



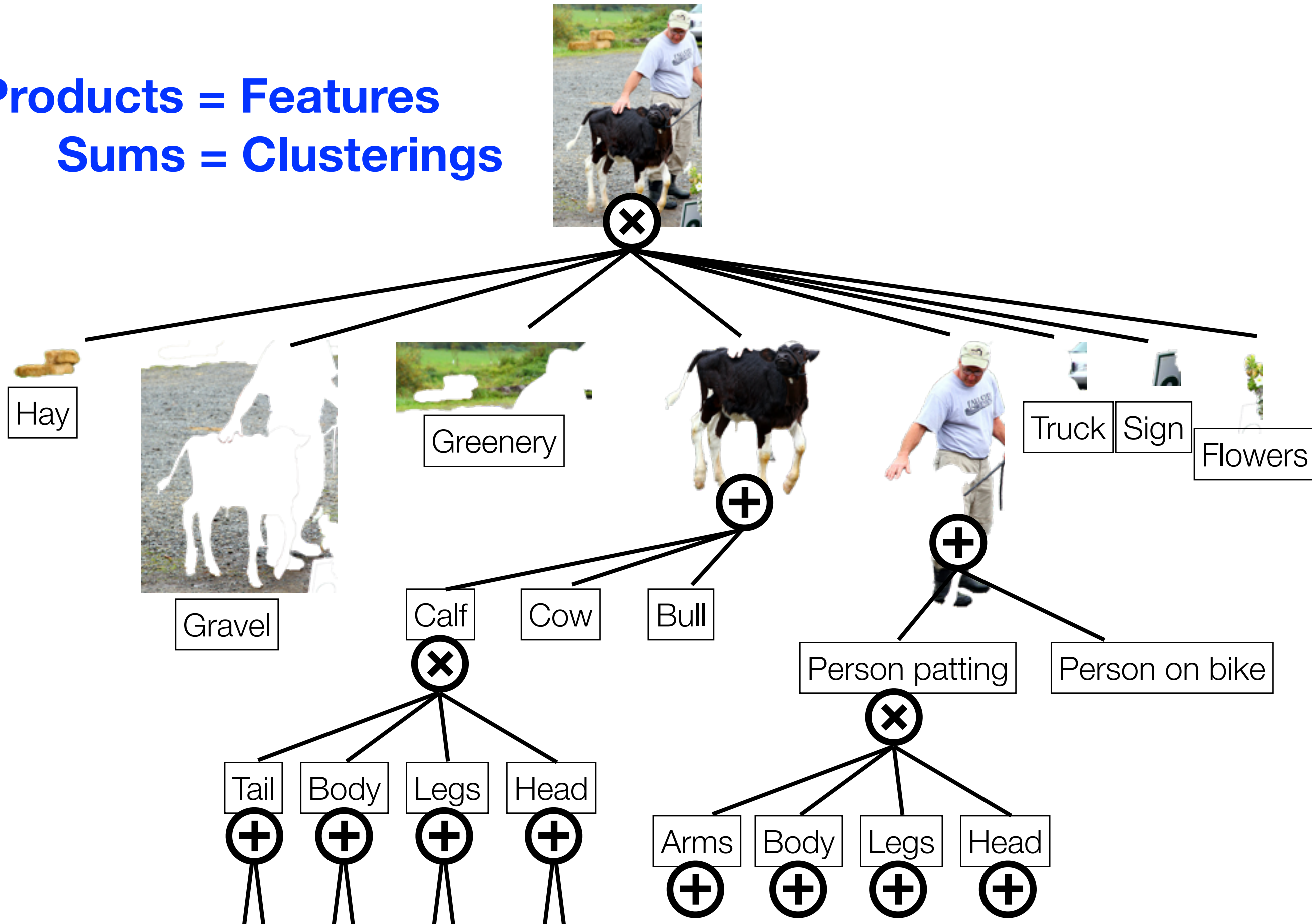
What Does an SPN Mean?

Products = Features
Sums = Clusterings



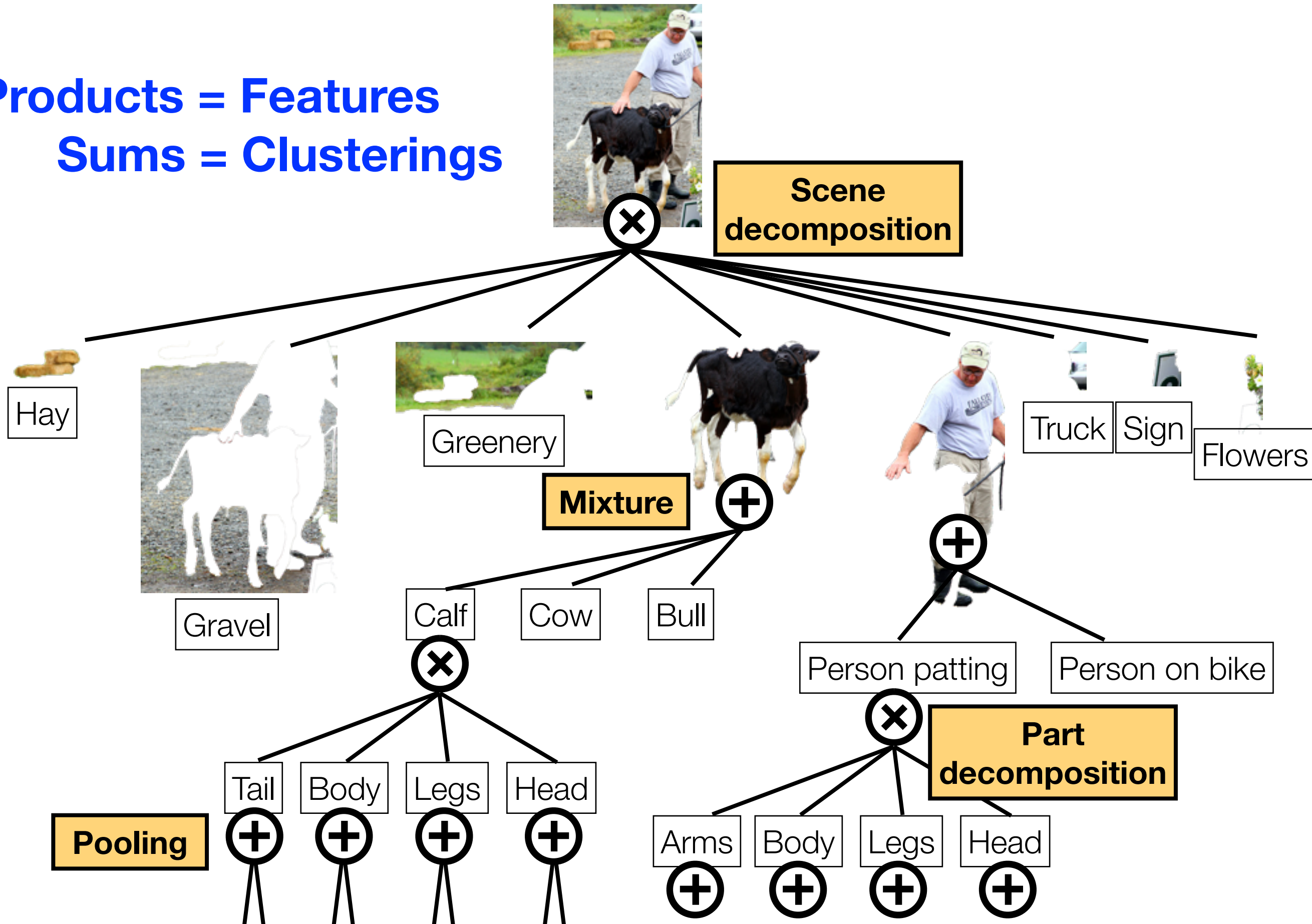
What Does an SPN Mean?

Products = Features
Sums = Clusterings



What Does an SPN Mean?

Products = Features
Sums = Clusterings



Special Cases

- Hierarchical mixture models
- Thin junction trees (e.g.: hidden Markov models)
- Non-recursive probabilistic context-free grammars
- Etc.

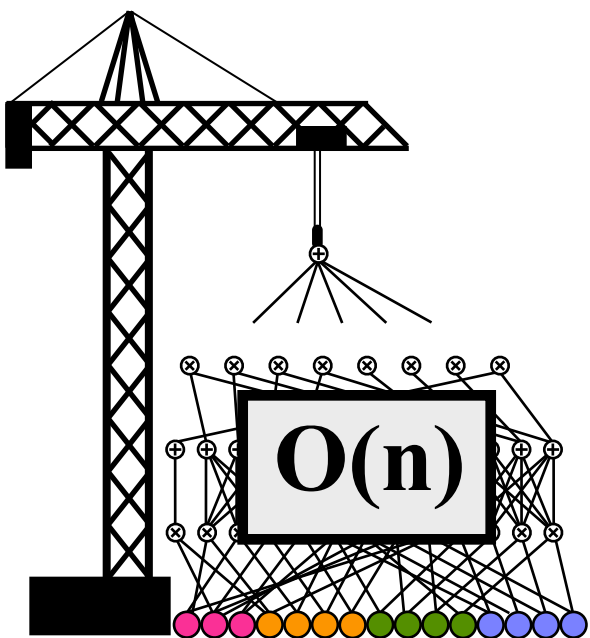
Weight Learning

- **Generative** (Poon & Domingos, UAI 2011)
UAI 2011 Best Paper Award
- **Discriminative** (Gens & Domingos, NIPS 2012)
NIPS 2012 Outstanding Student Paper Award

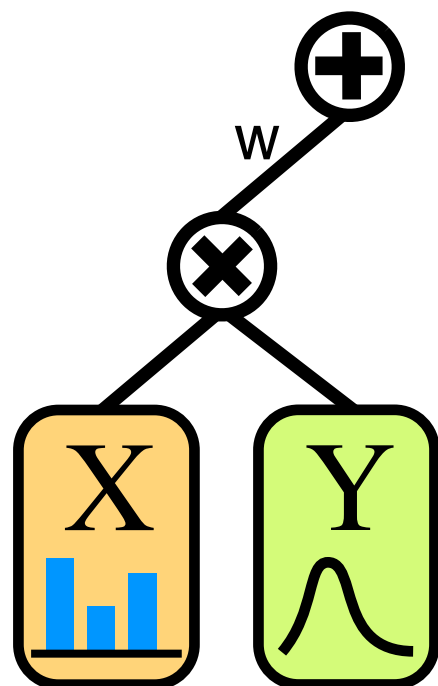
Weight Learning

- **Generative** (Poon & Domingos, UAI 2011)
UAI 2011 Best Paper Award
- **Discriminative** (Gens & Domingos, NIPS 2012)
NIPS 2012 Outstanding Student Paper Award

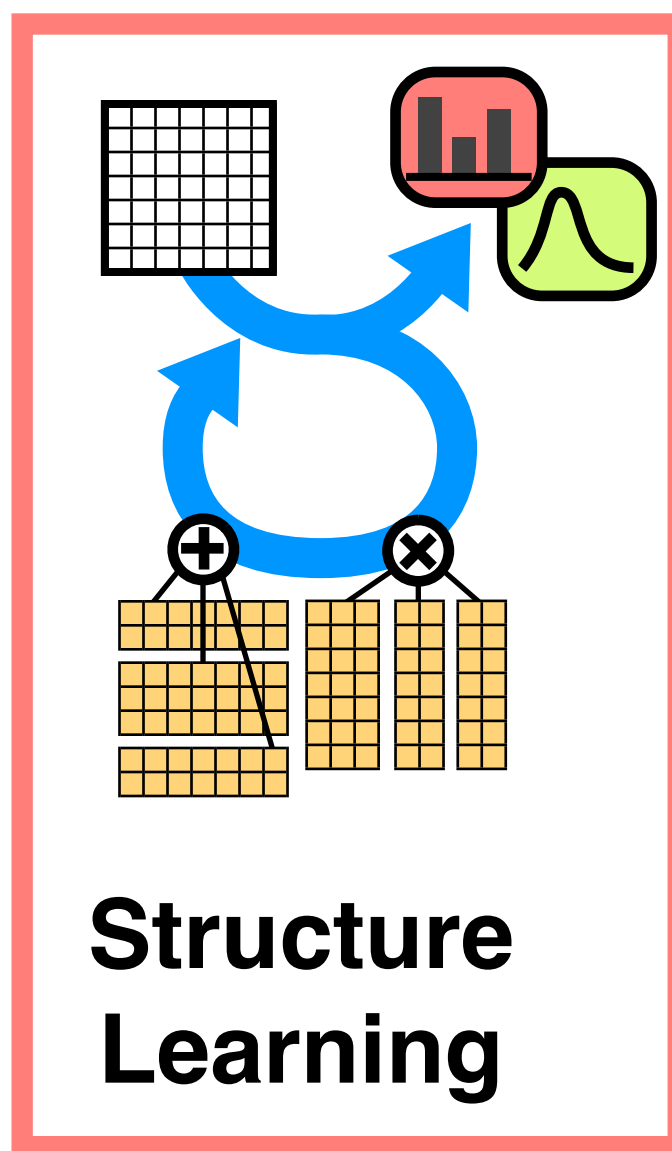
Key limitation: Requires a structure



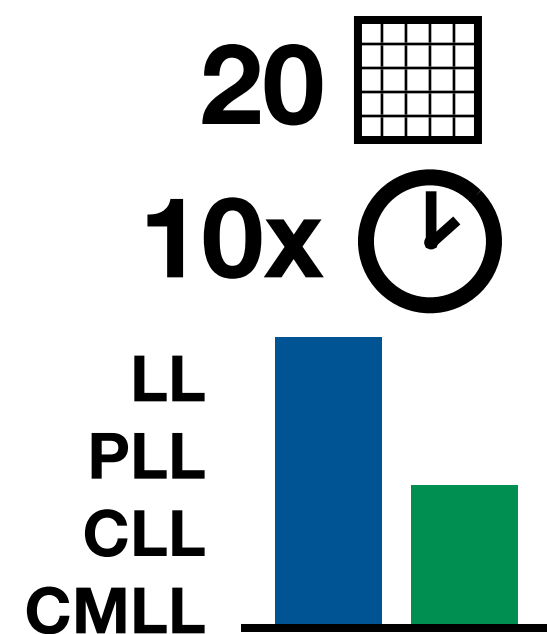
Motivation



SPN Review



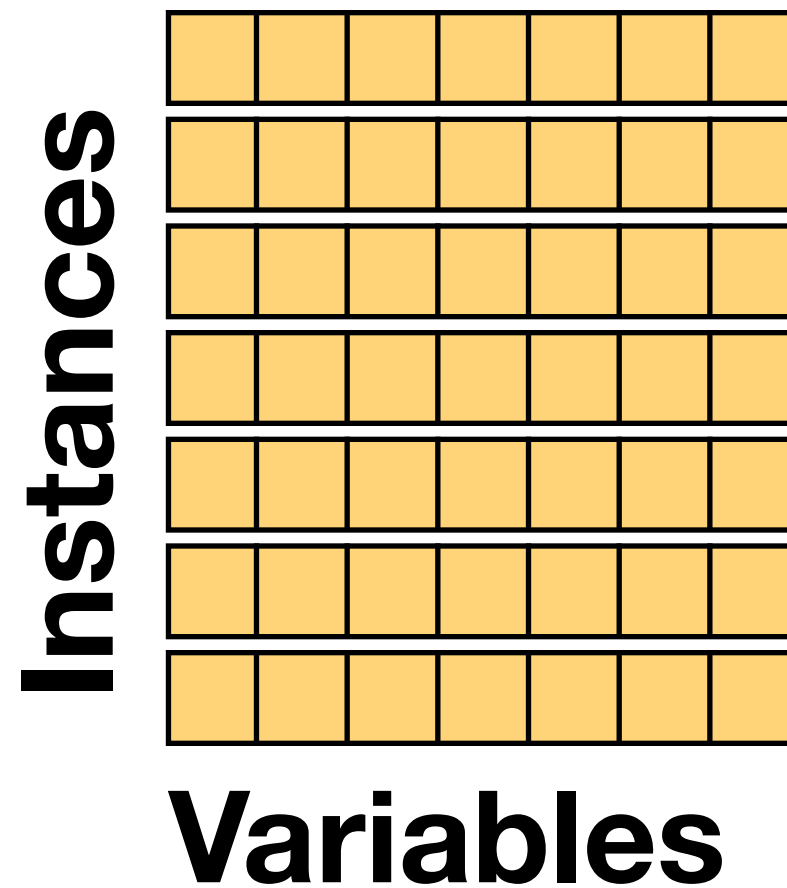
Structure Learning



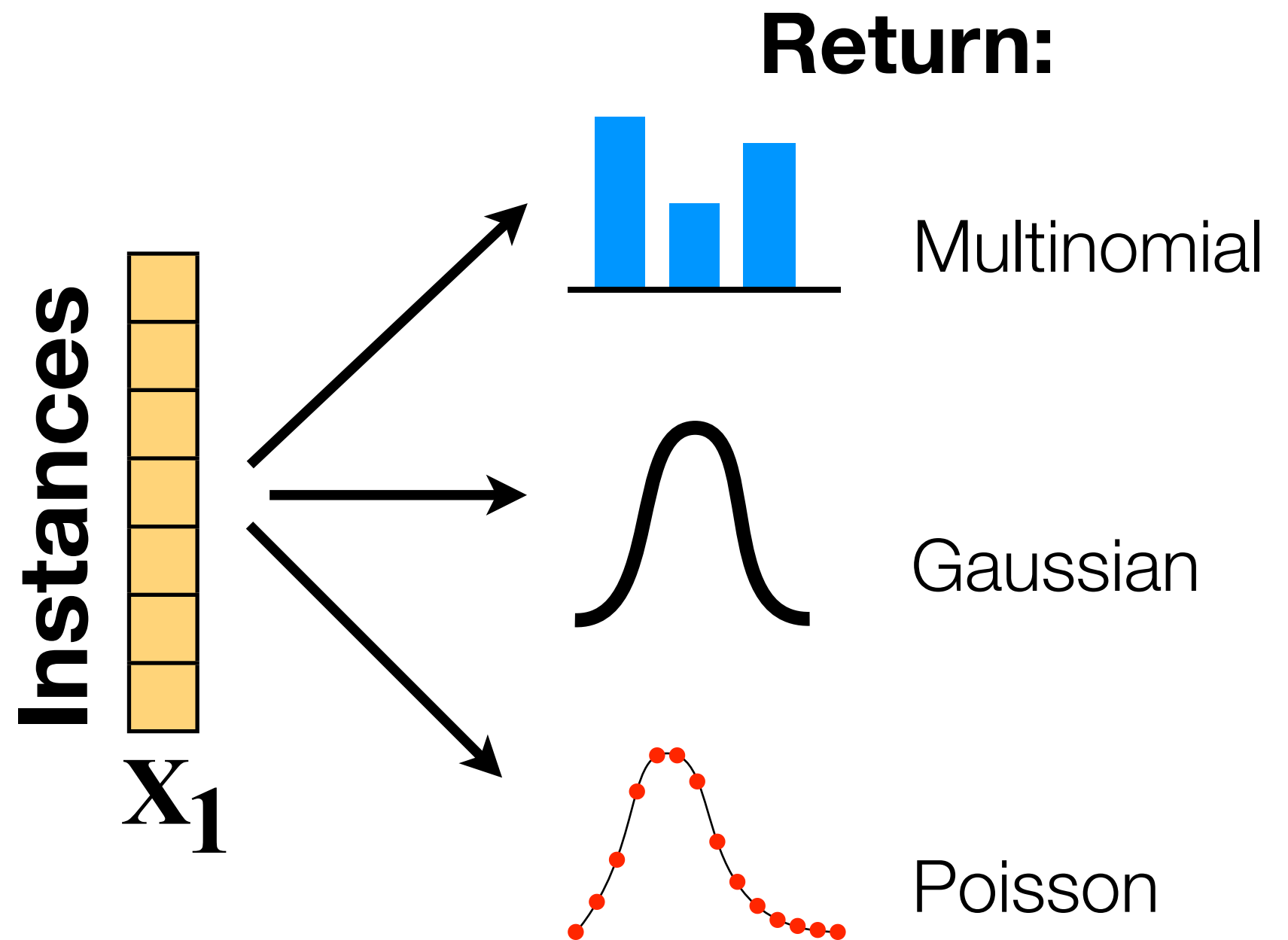
Experiments

LearnSPN

Recursive algorithm



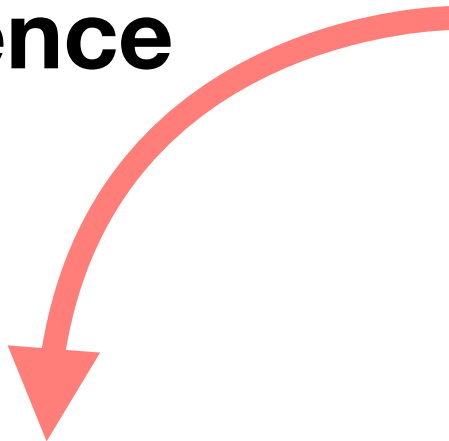
LearnSPN



Instances

Variables

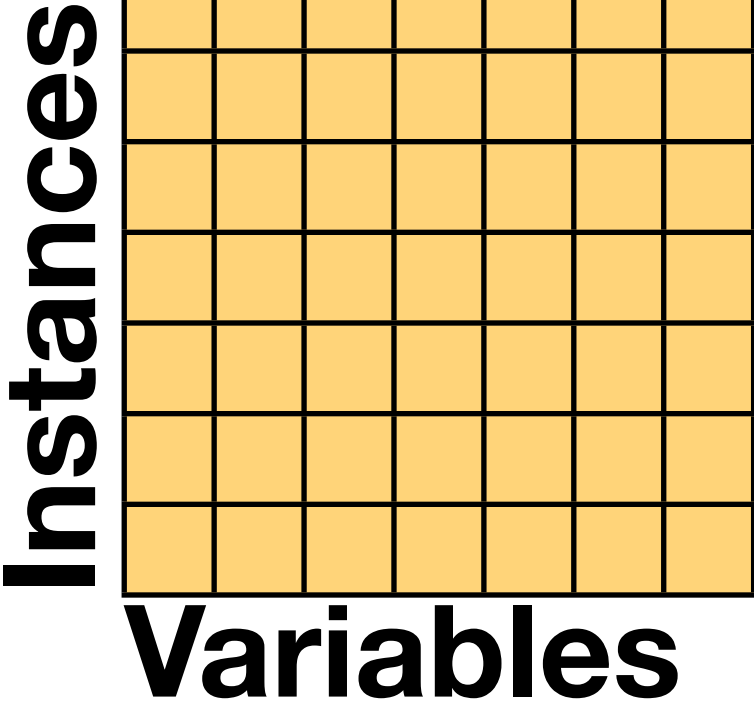
**Establish
approximate
independence**



Instances

Variables

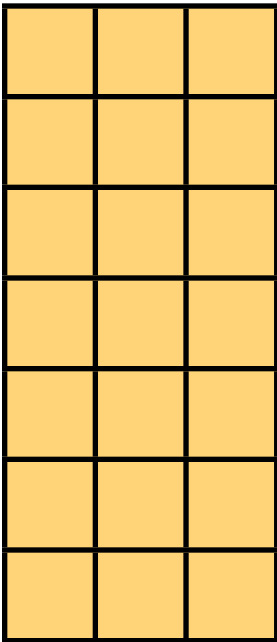
**Establish
approximate
independence**



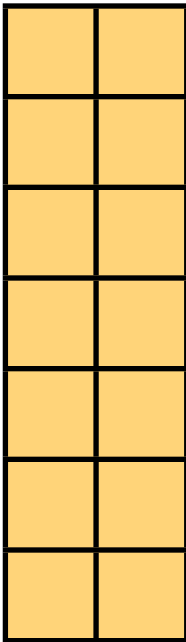
Return:



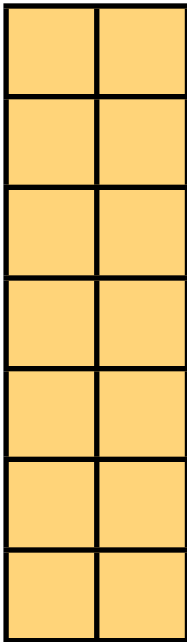
Recurse



Recurse



Recurse



Establish approximate independence

Instances

Variables

**If no
independence,
cluster similar
instances**

Return:



Recurse

[illegible]

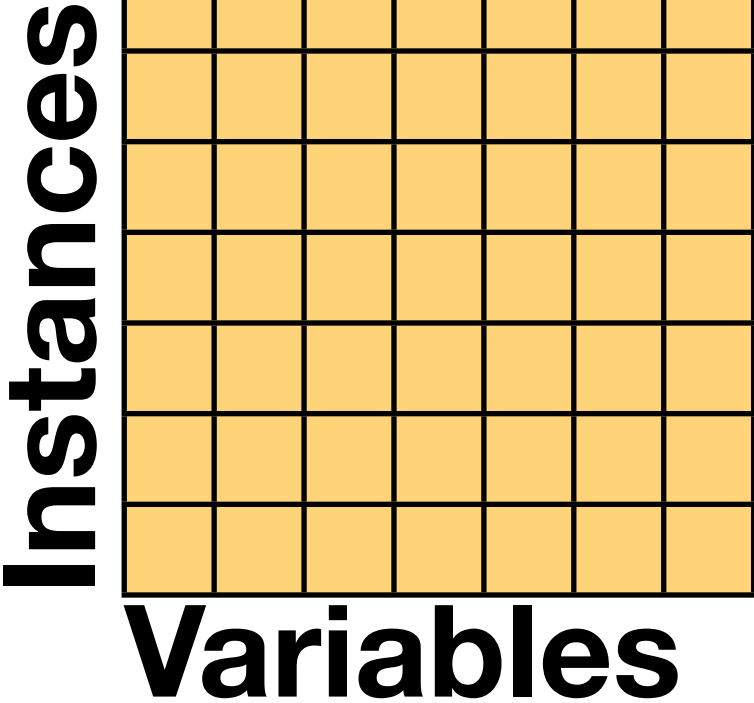
Recurse

[illegible]

Recurse

[illegible]

**Establish
approximate
independence**

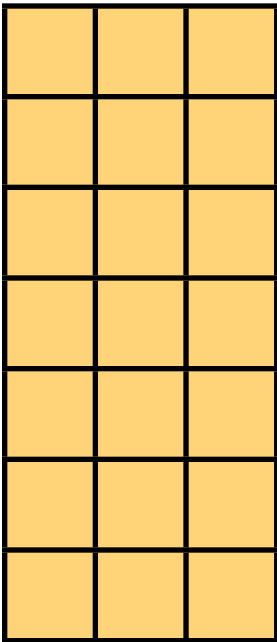


**If no
independence,
cluster similar
instances**

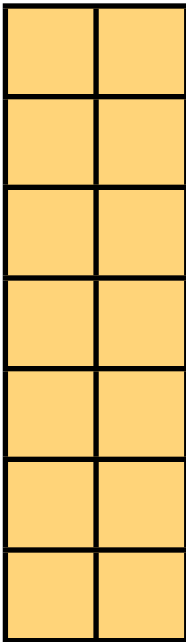
Return:



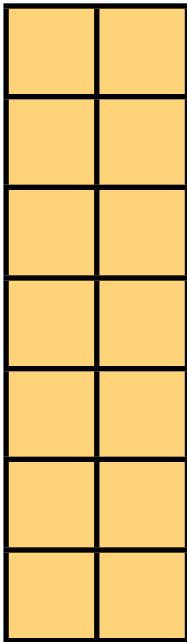
Recurse



Recurse



Recurse



Return:

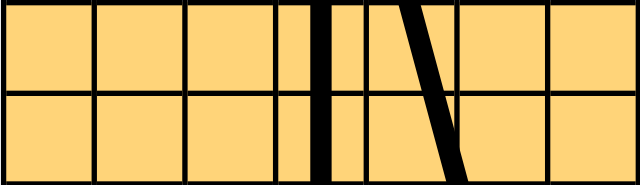


w_1

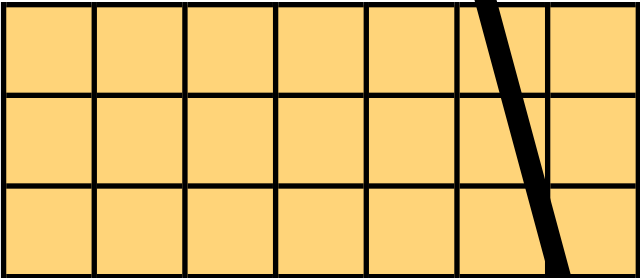
w_2

w_3

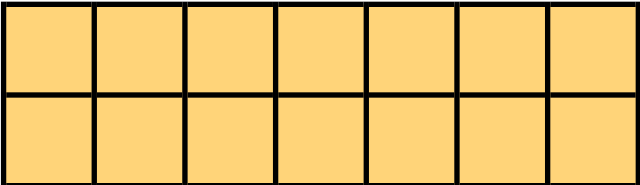
Recurse



Recurse

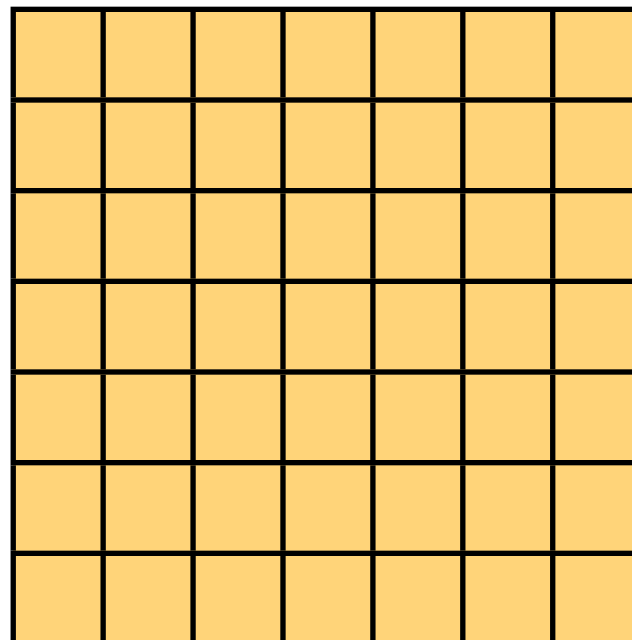


Recurse

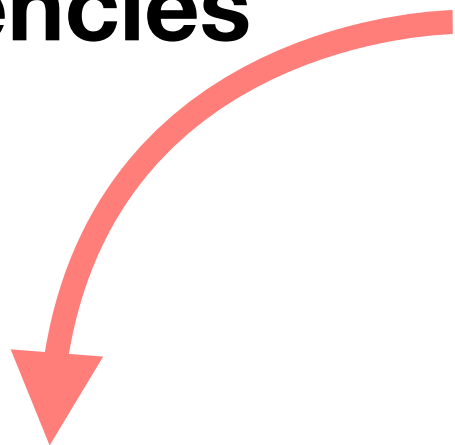


If no detectable dependencies

Instances



Variables



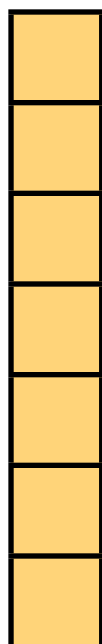
Return:



R...



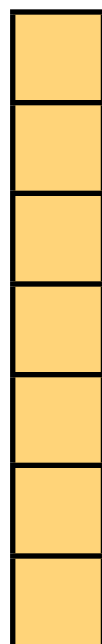
R...



R...

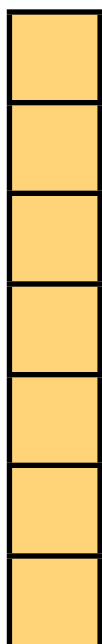


R...



...

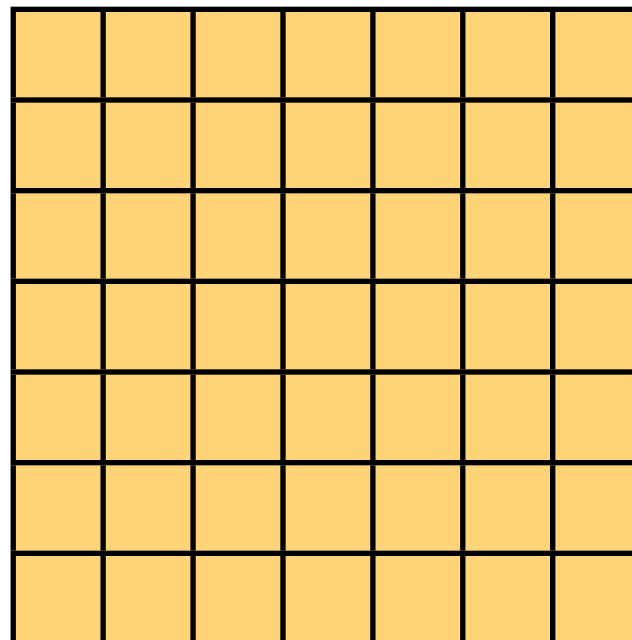
R...



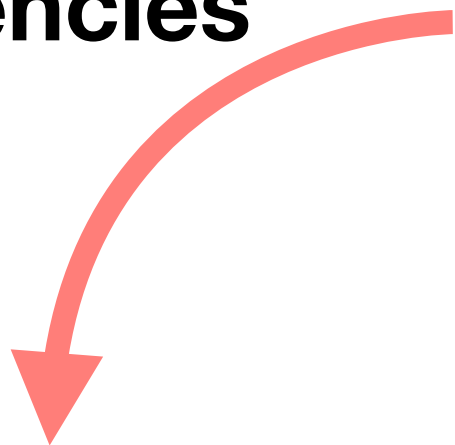
**Fully factorized
distribution**

If no detectable dependencies

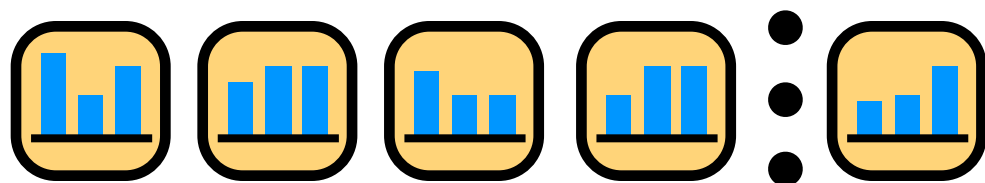
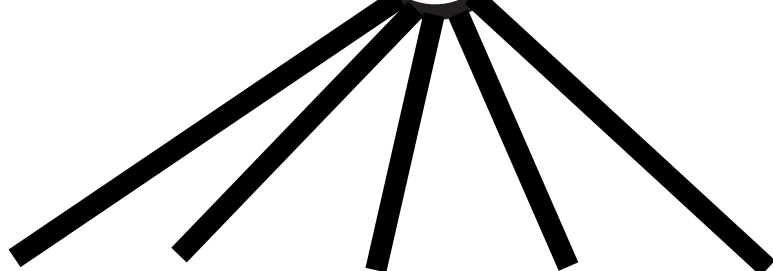
Instances



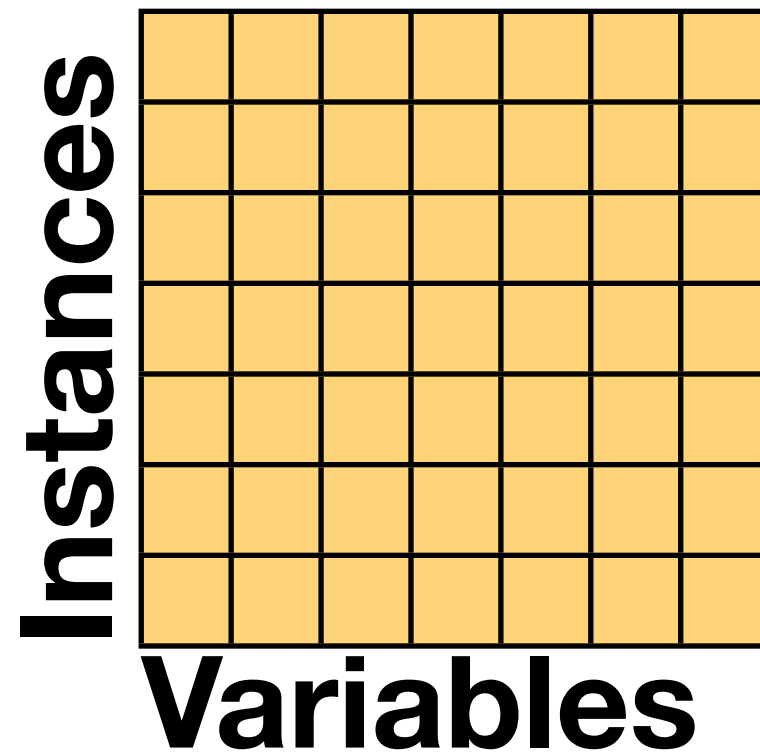
Variables



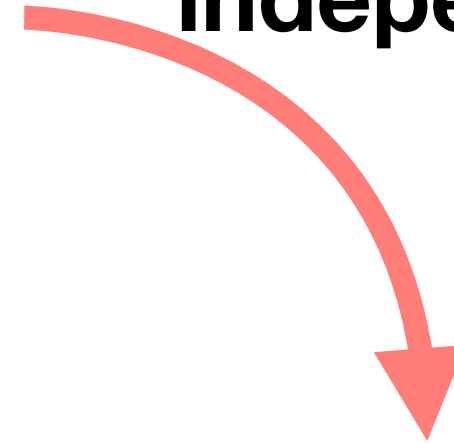
Return:



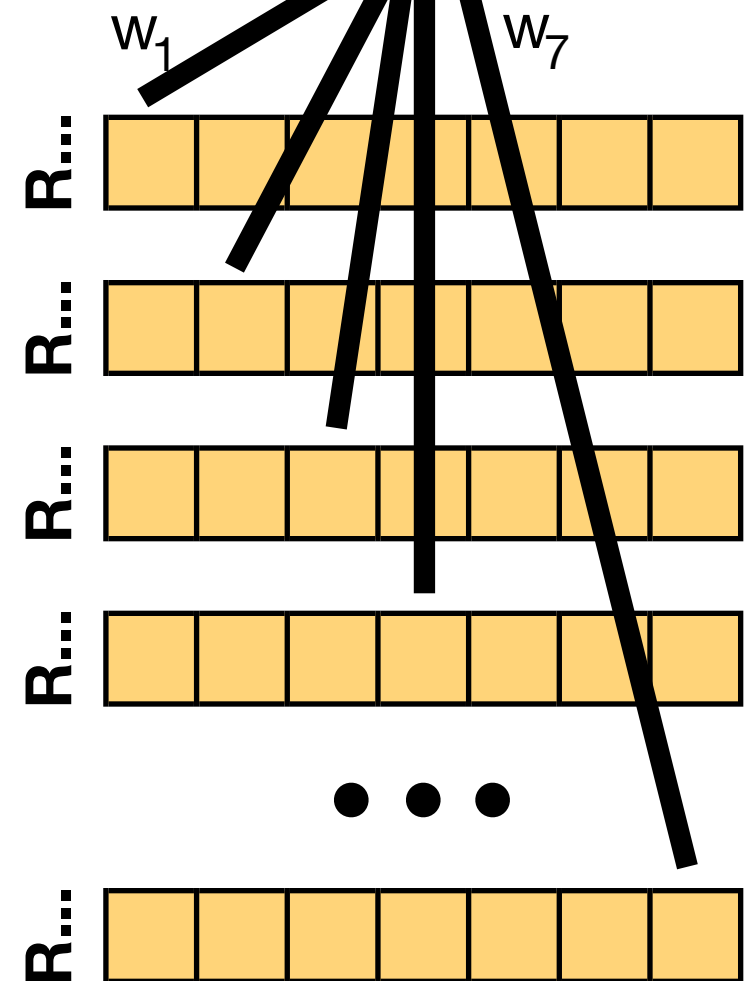
Fully factorized distribution



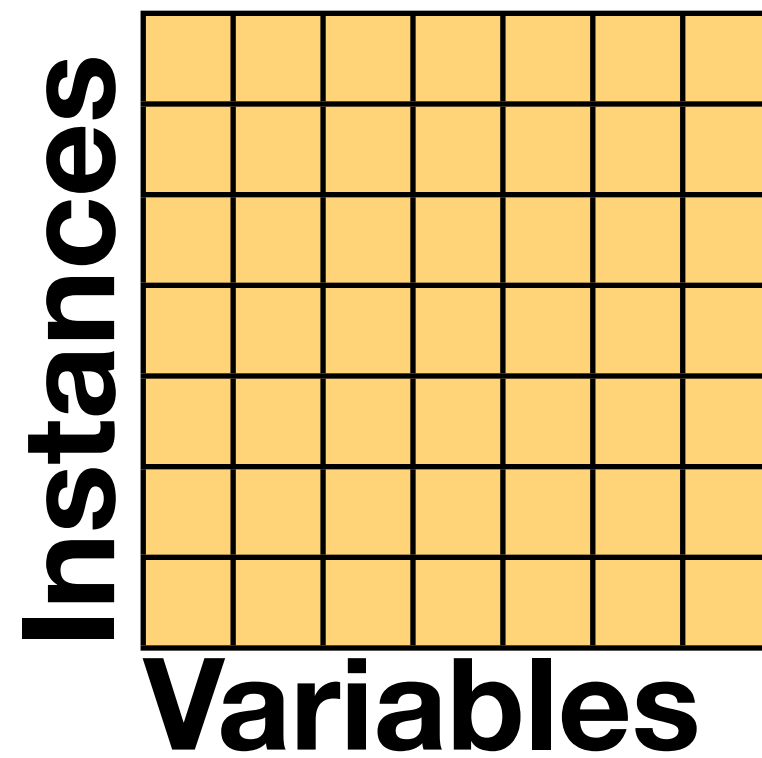
If never finds
independence



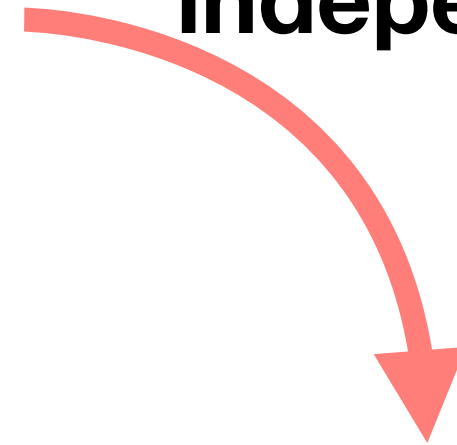
Return: \oplus



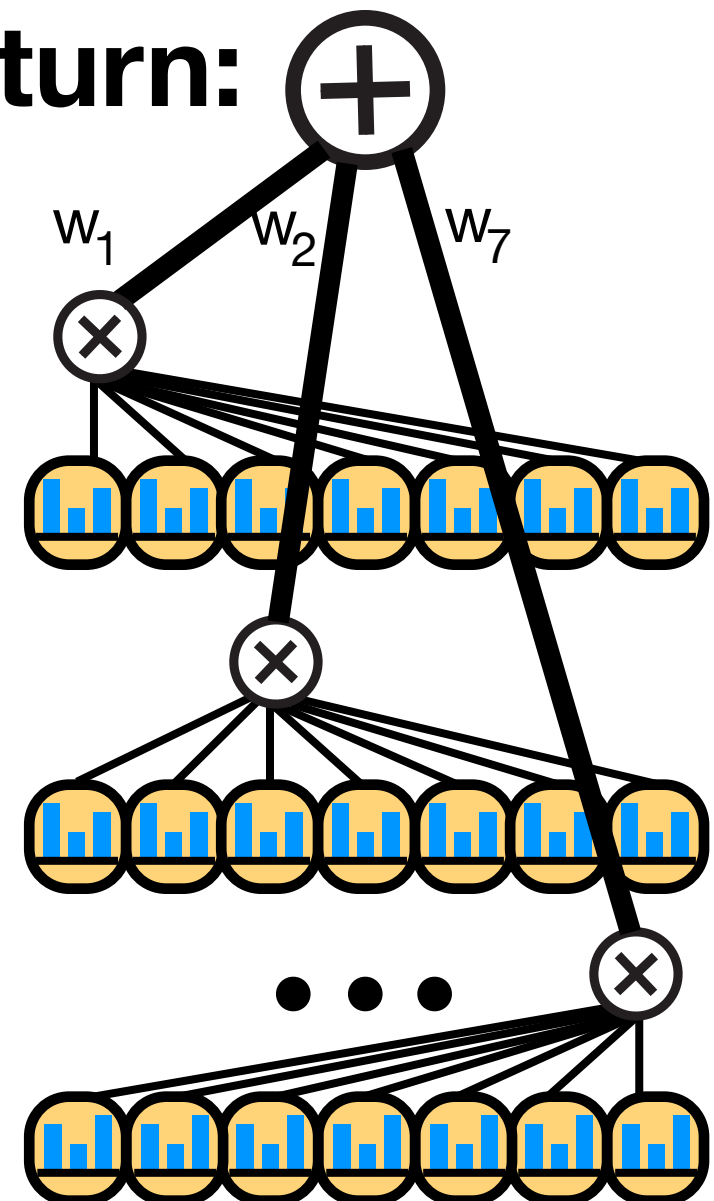
Kernel density estimate



If never finds
independence

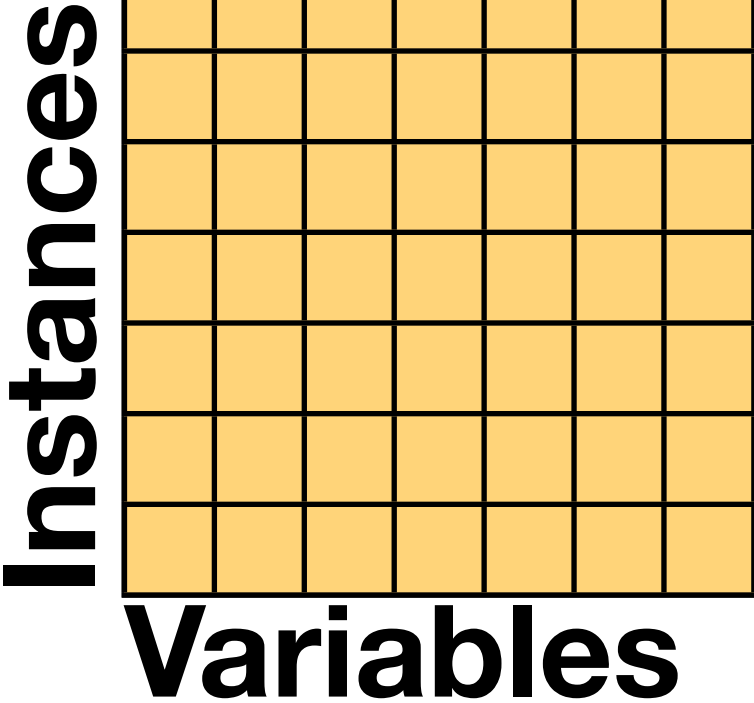


Return:



Kernel density estimate

**Establish
approximate
independence**

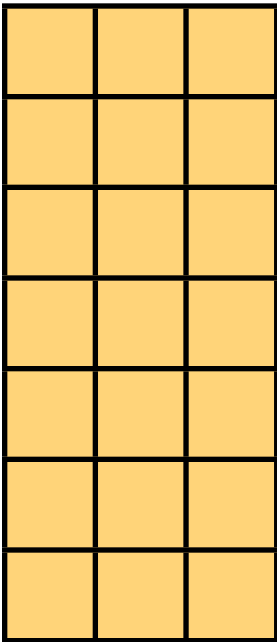


**If no
independence,
cluster similar
instances**

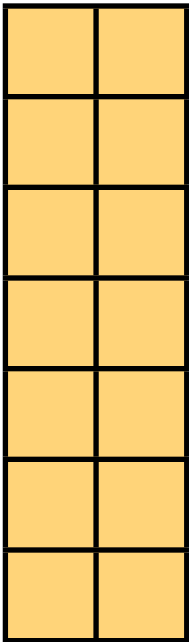
Return:



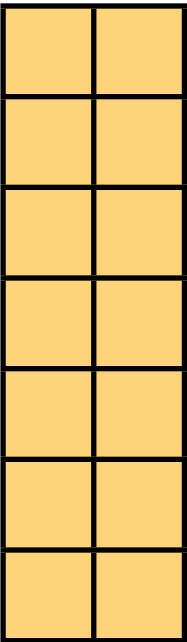
Recurse



Recurse



Recurse



Return:

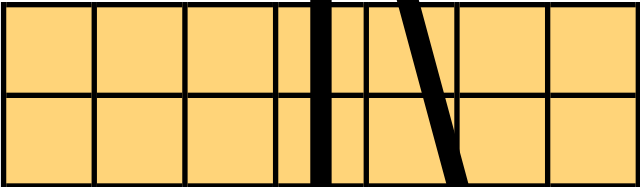


w_1

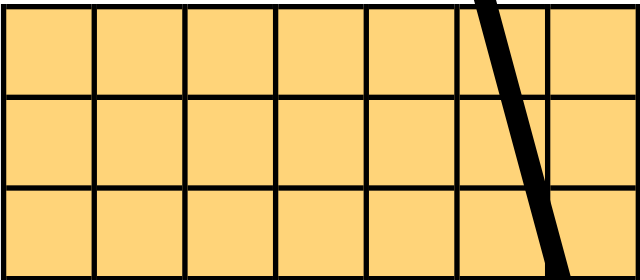
w_2

w_3

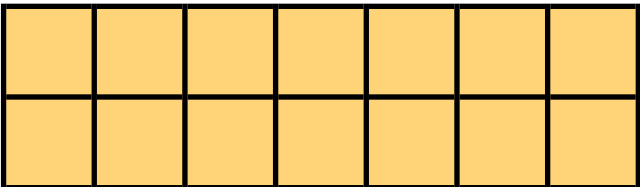
Recurse



Recurse



Recurse

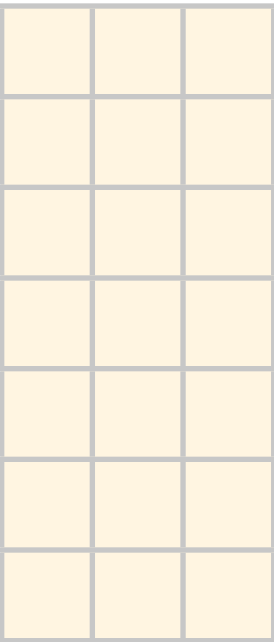


Establish
approximate
independence

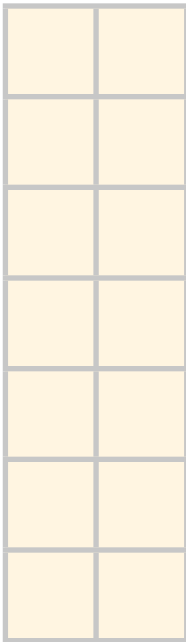
Return:



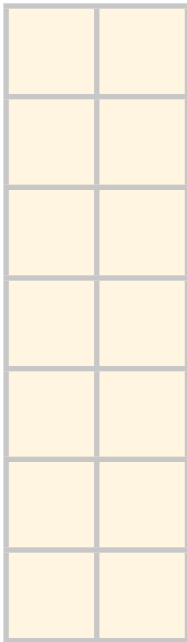
Recurse



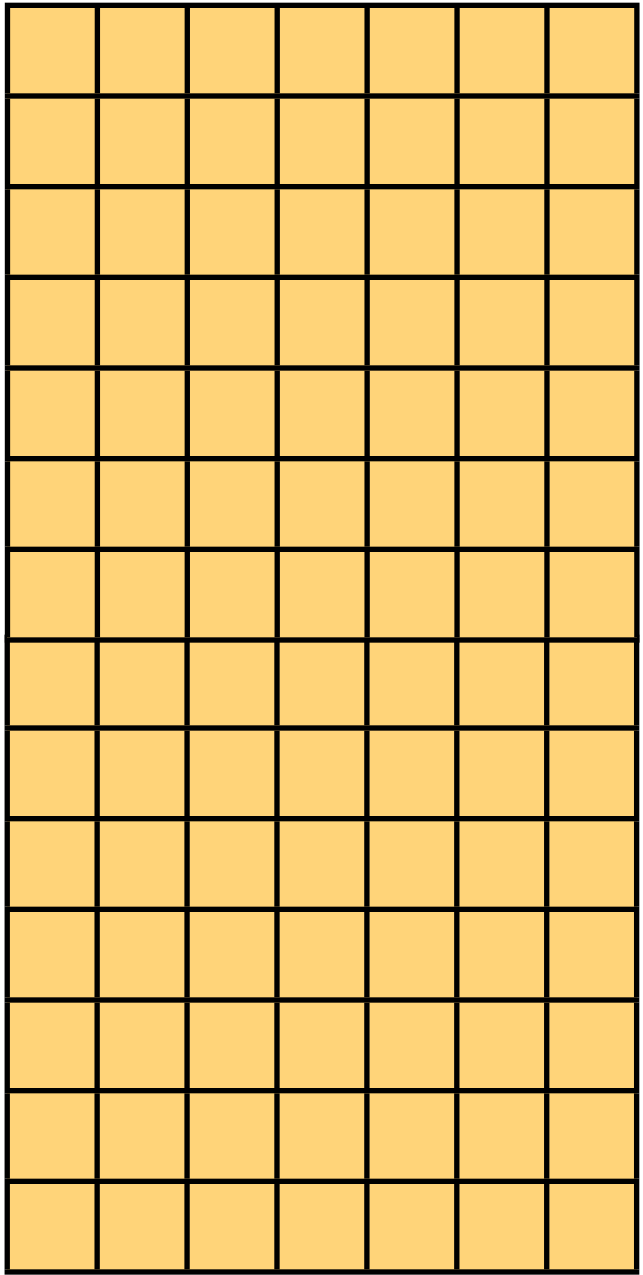
Recurse



Recurse



Instances



Variables

If no
independence,
cluster similar
instances

Return:



w_1

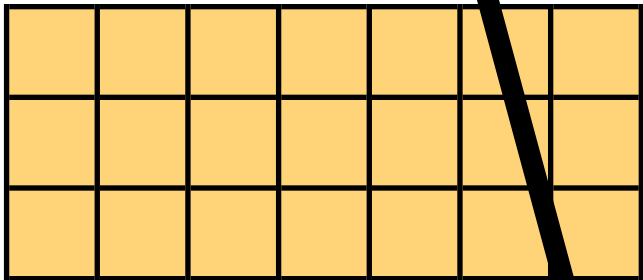
w_2

w_3

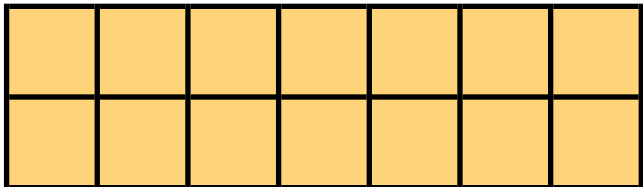
Recurse



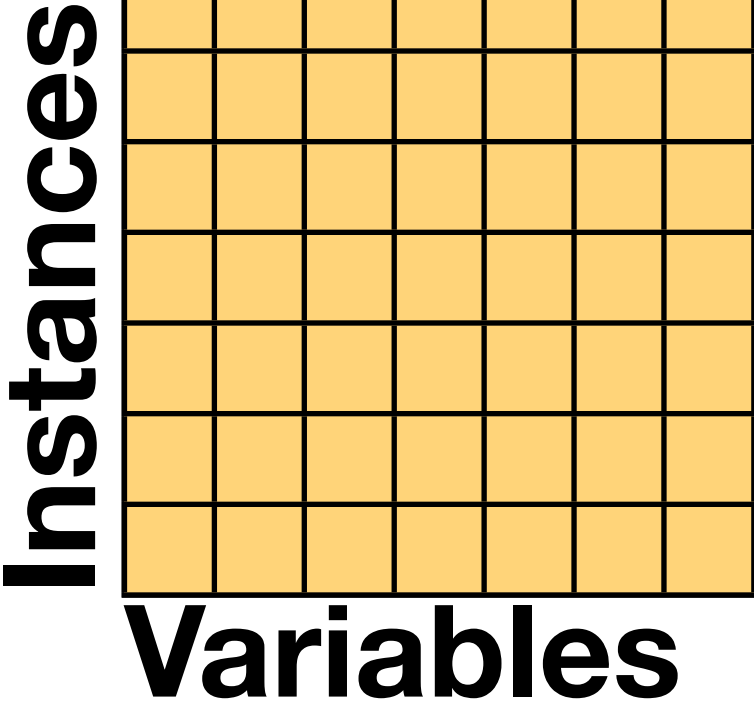
Recurse



Recurse



**Establish
approximate
independence**

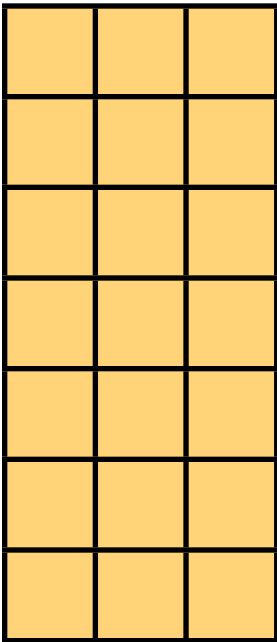


**If no
independence,
cluster similar
instances**

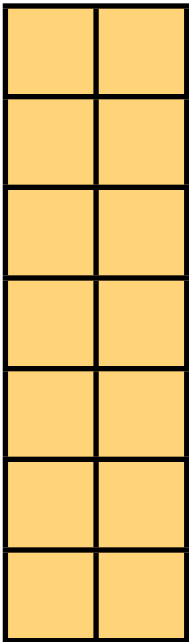
Return:



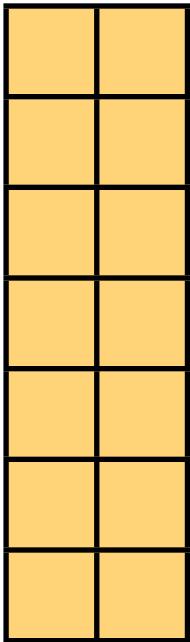
Recurse



Recurse



Recurse



Return:

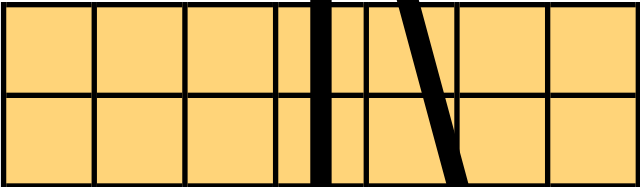


w_1

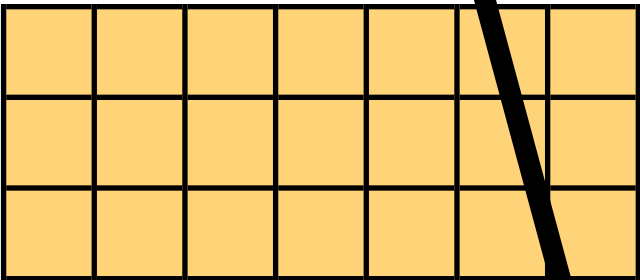
w_2

w_3

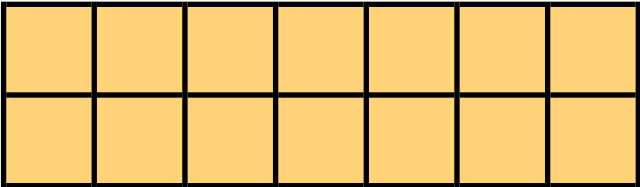
Recurse



Recurse

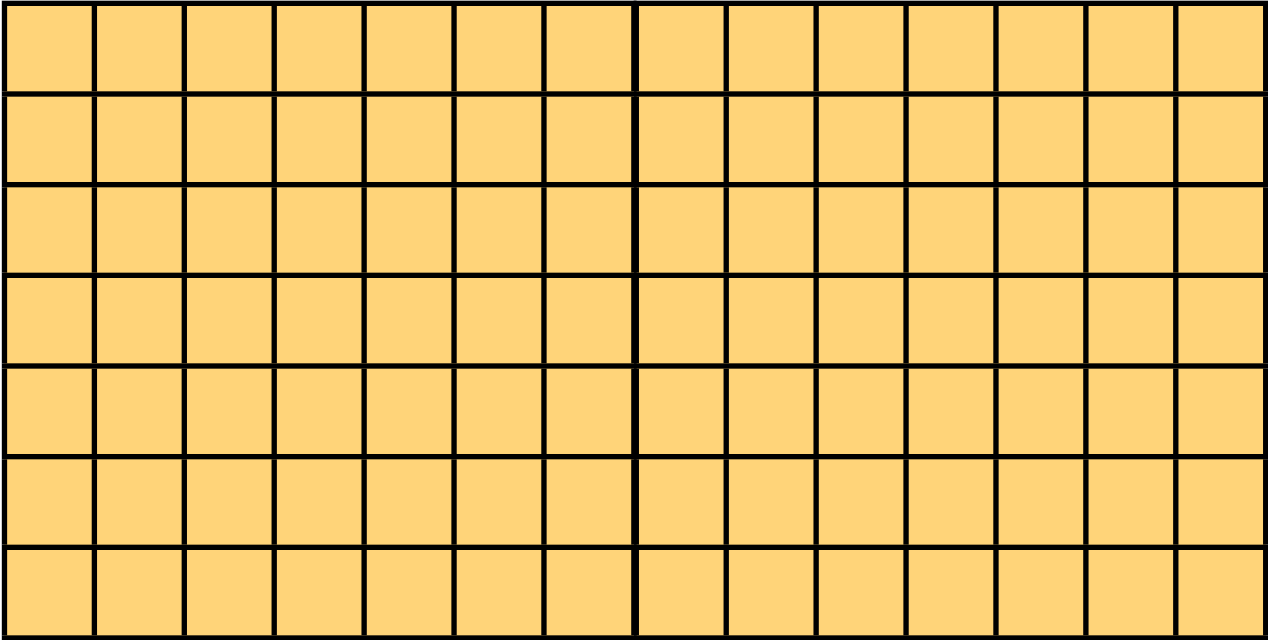


Recurse



**Establish
approximate
independence**

Instances



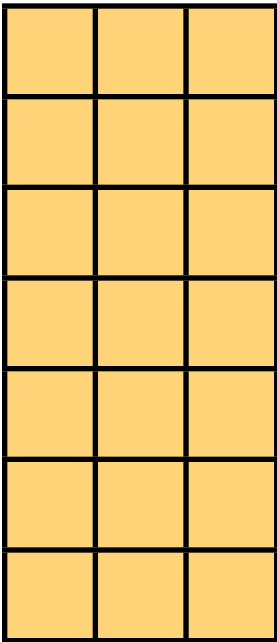
dependence,
similar
instances

Variables

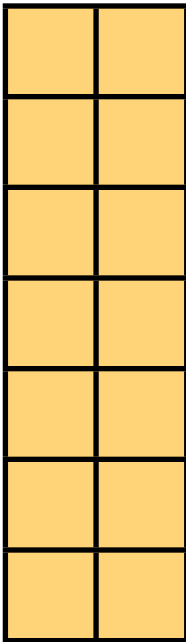
Return:



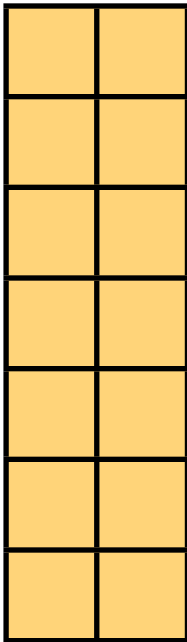
Recurse



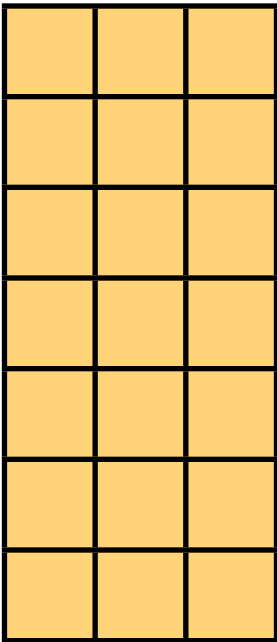
Recurse



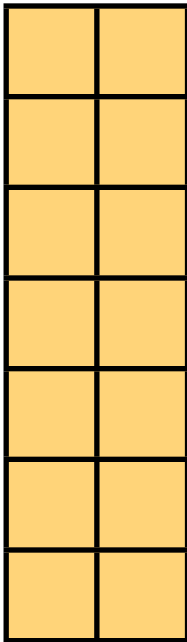
Recurse



Recurse



Recurse



Return:

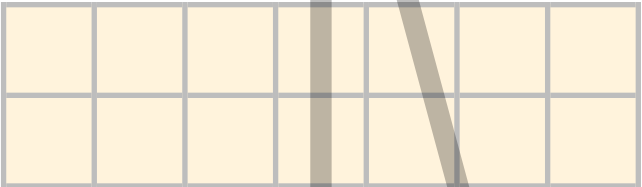


w_1

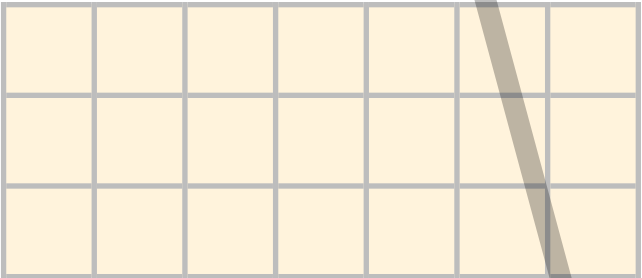
w_2

w_3

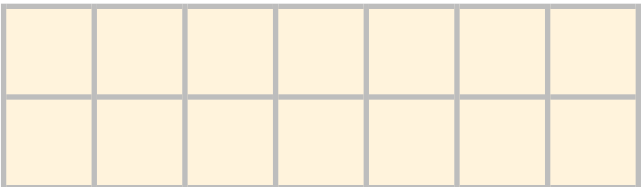
Recurse



Recurse



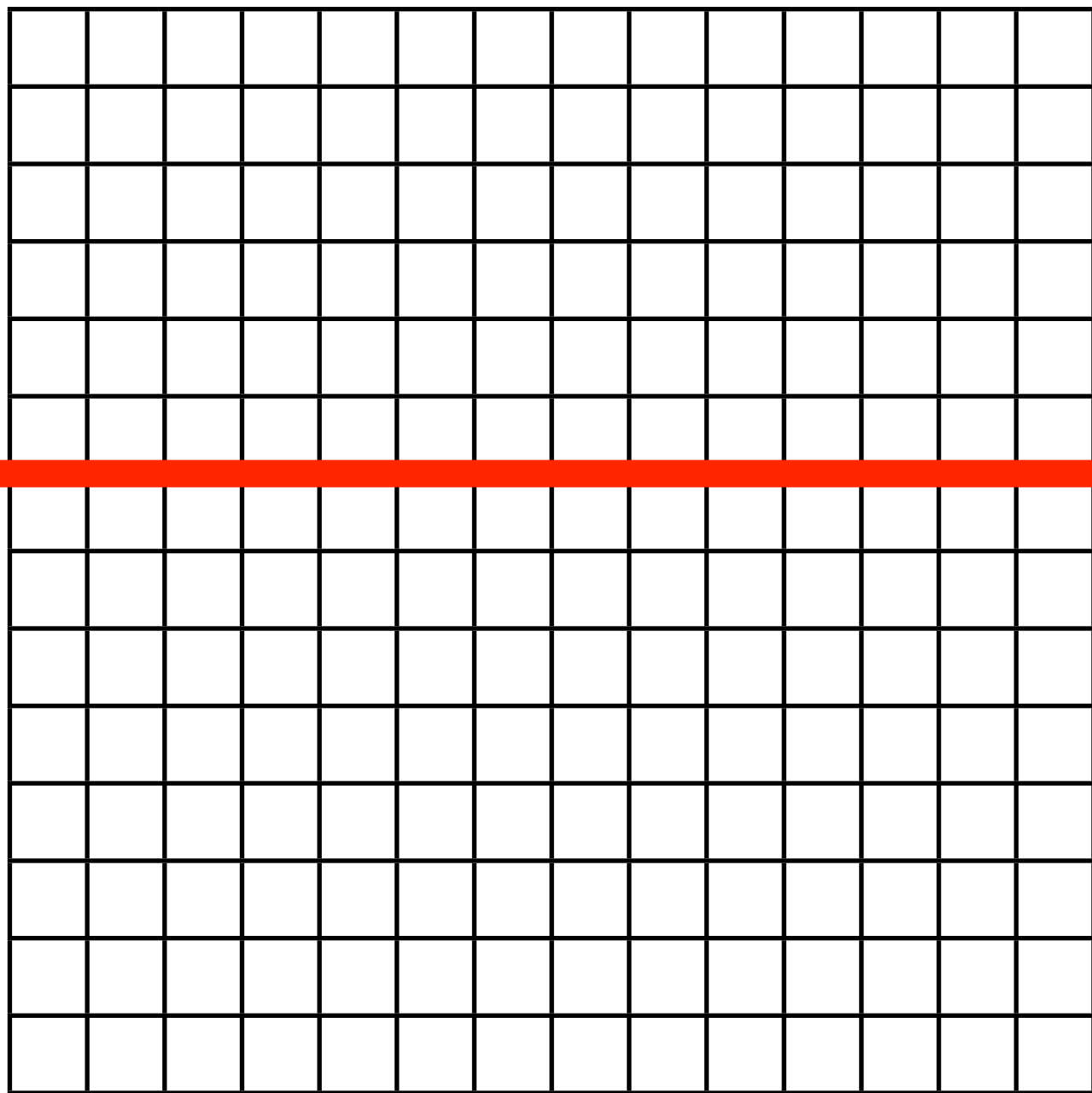
Recurse



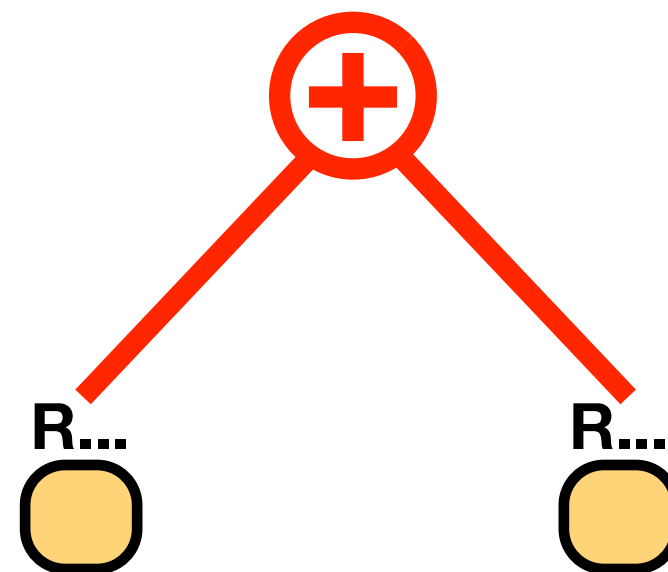
Instances

Variables

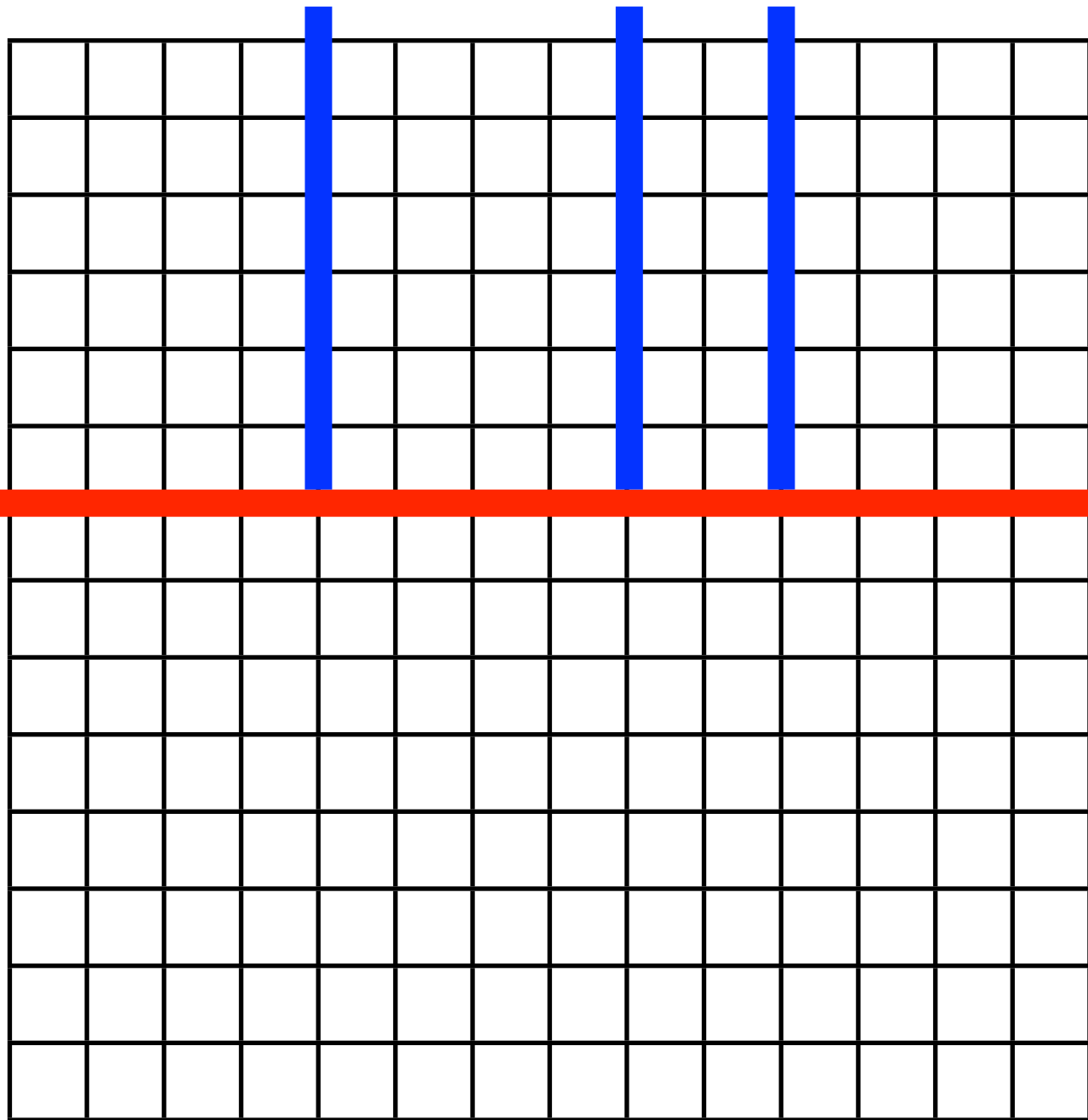
Instances



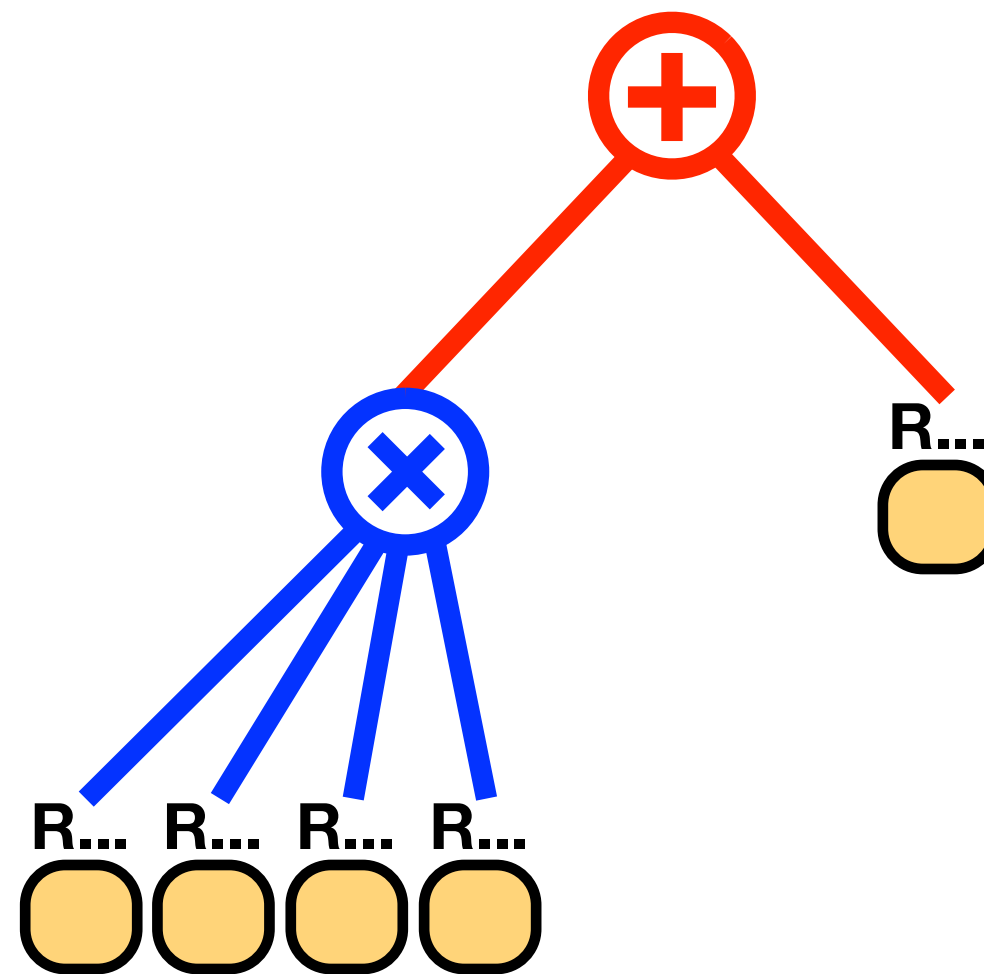
Variables



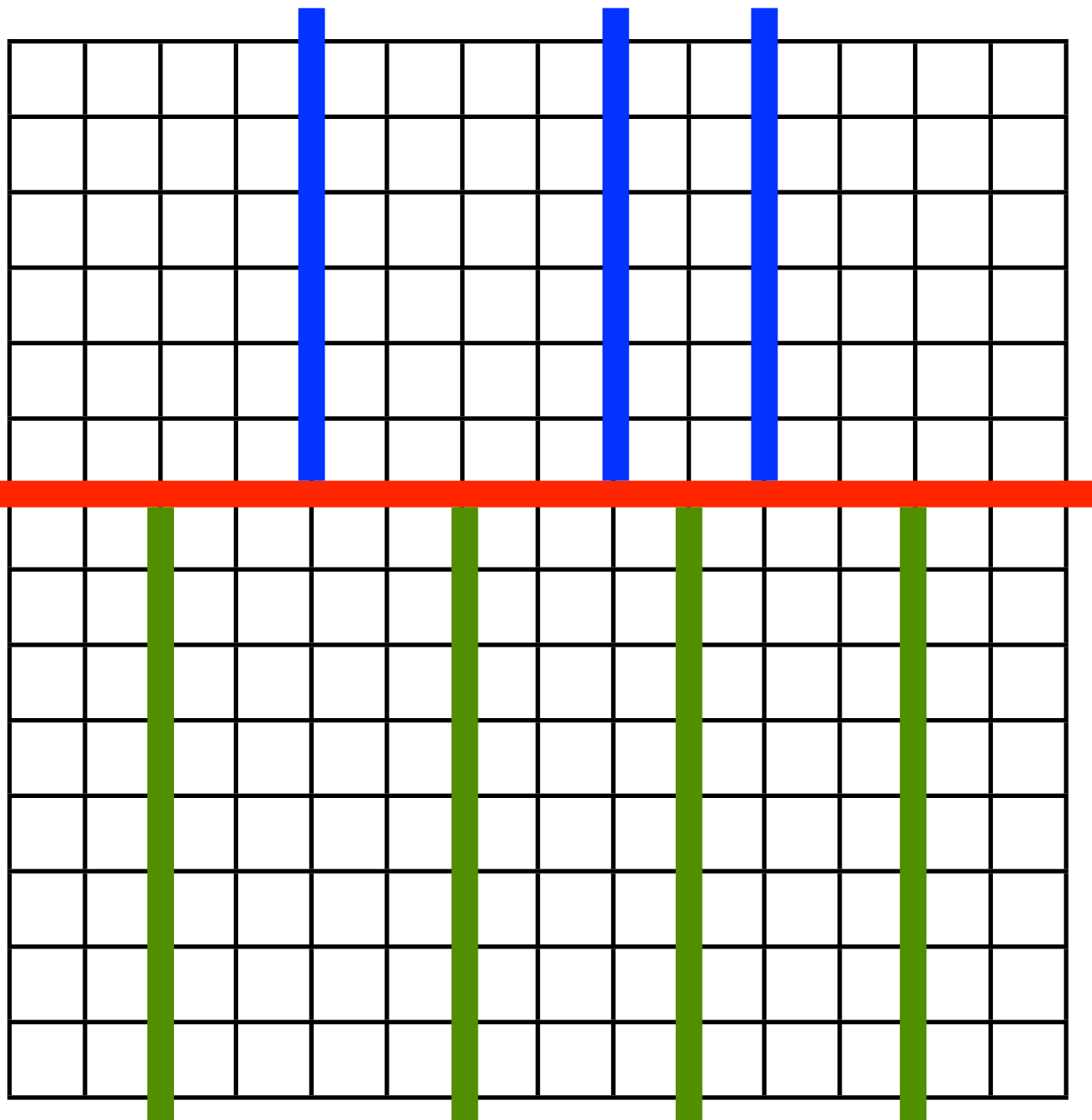
Instances



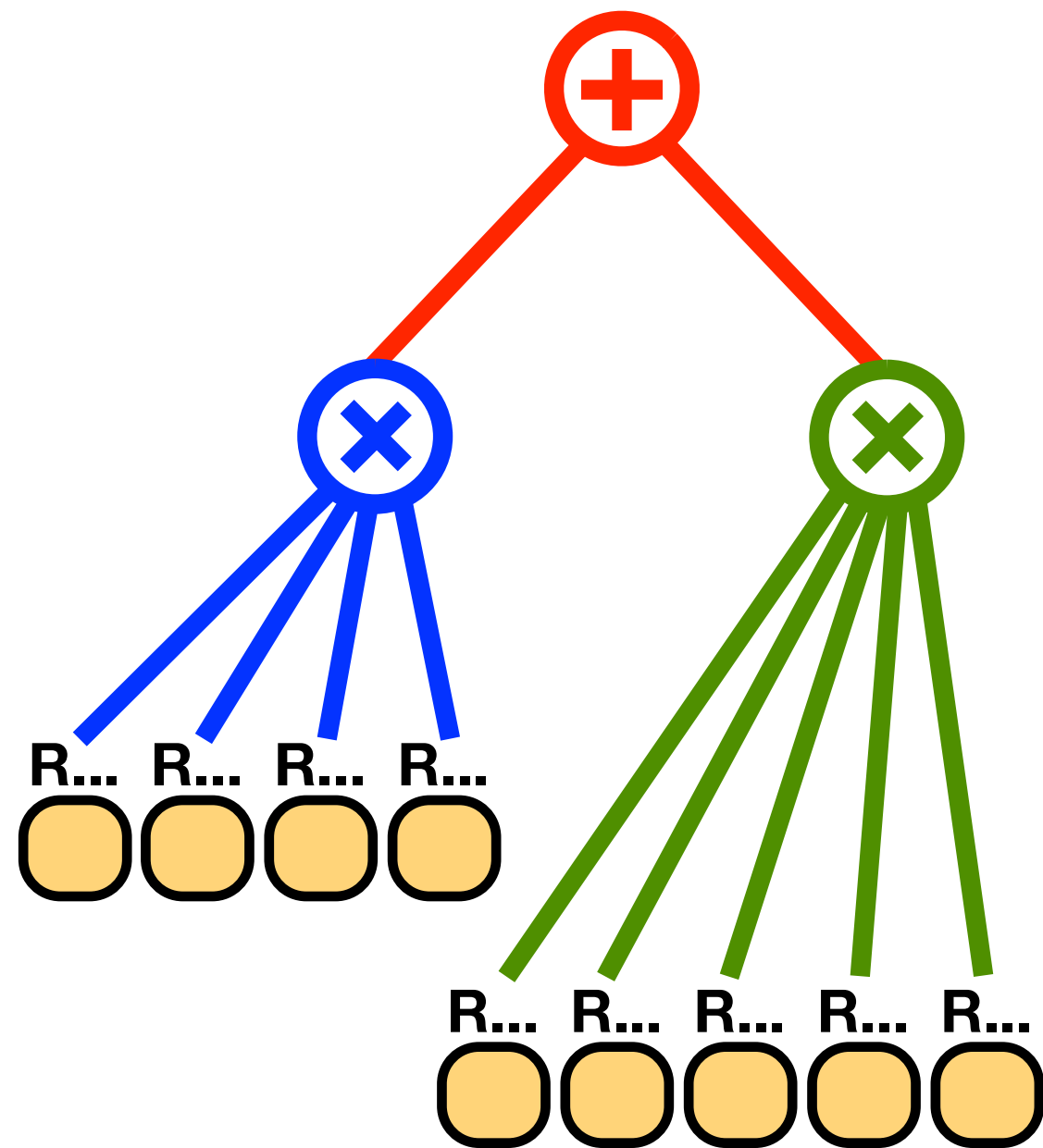
Variables



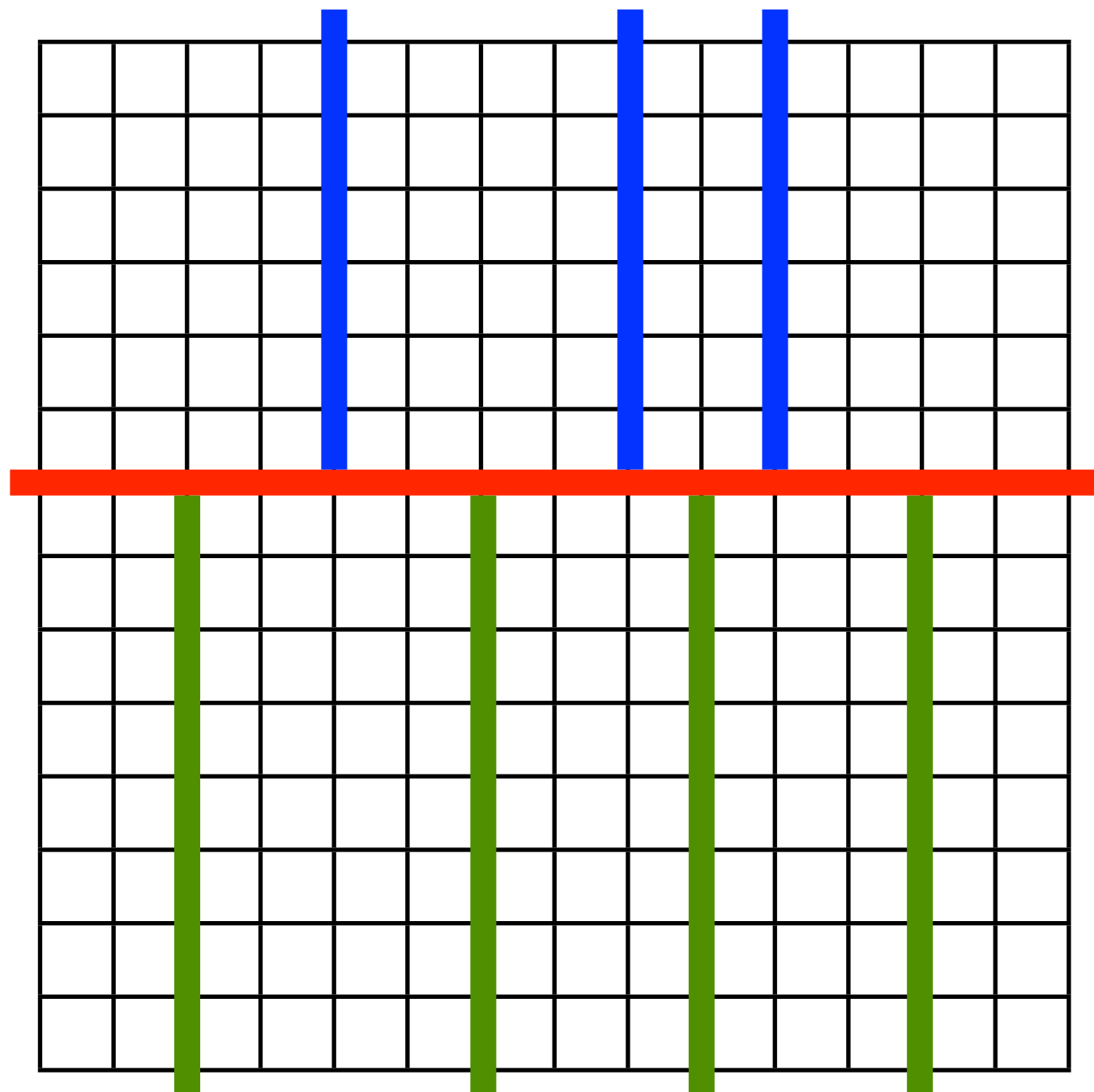
Instances



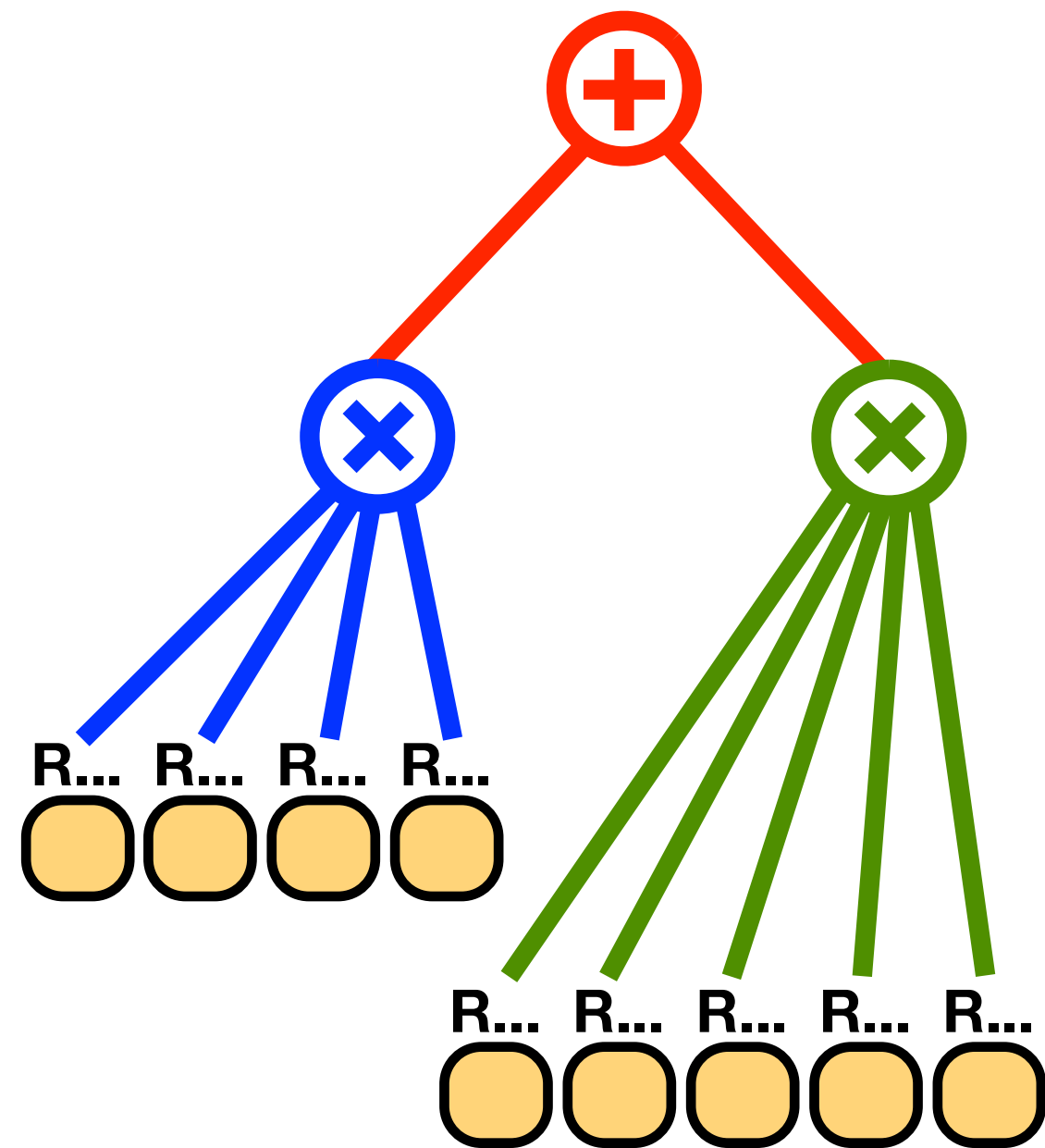
Variables



Instances

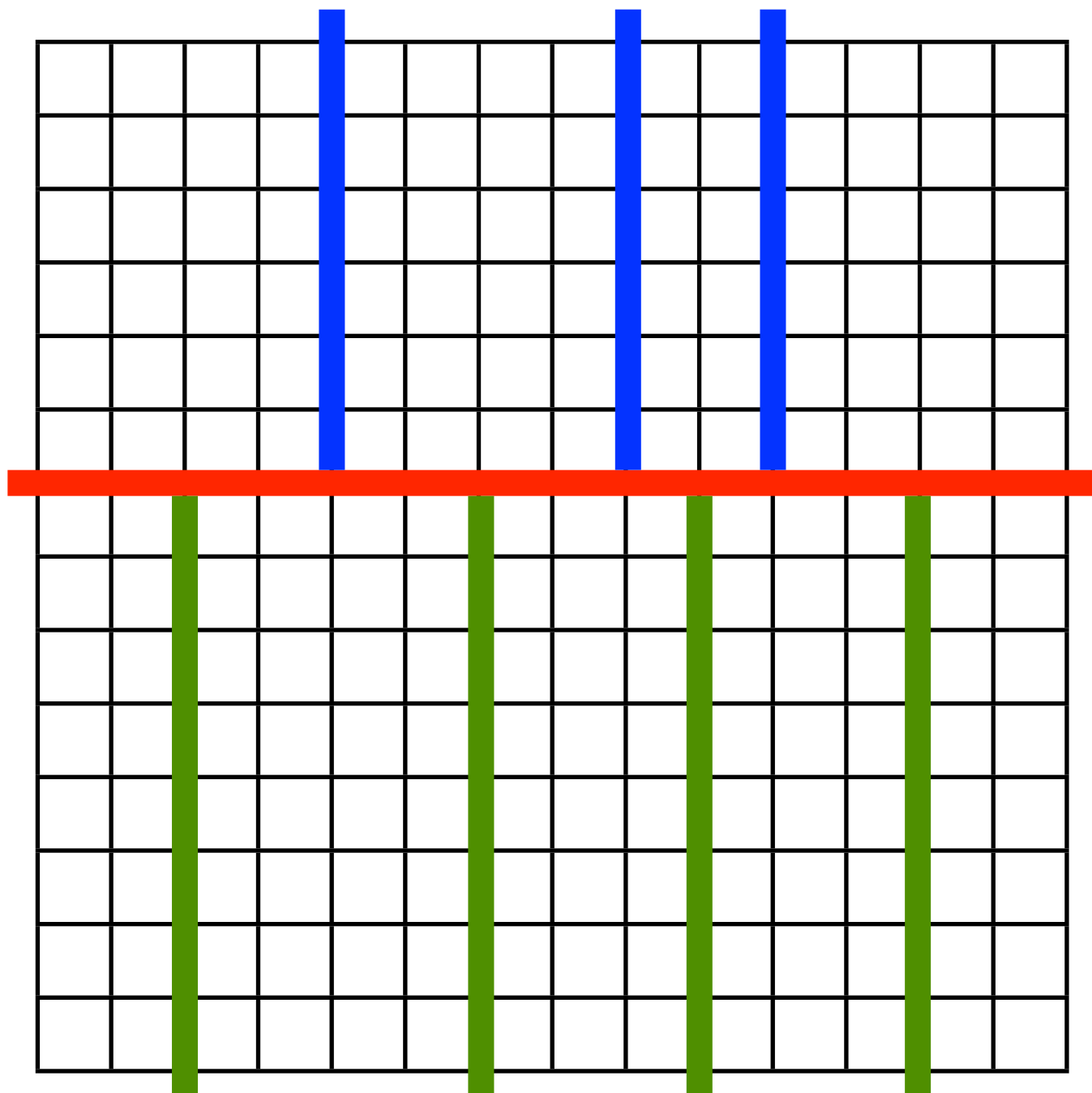


Variables

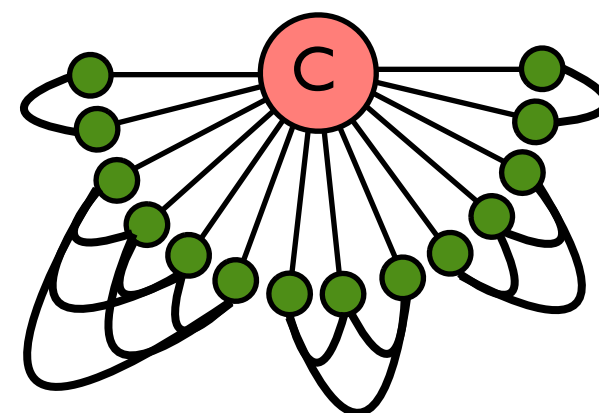
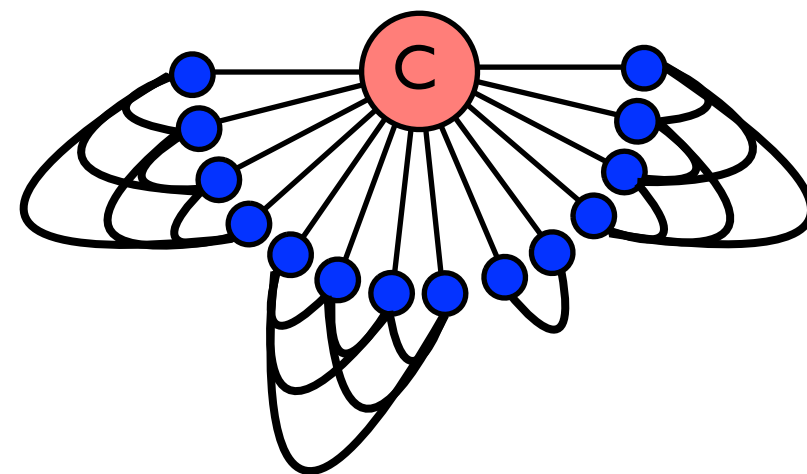


High treewidth

Instances

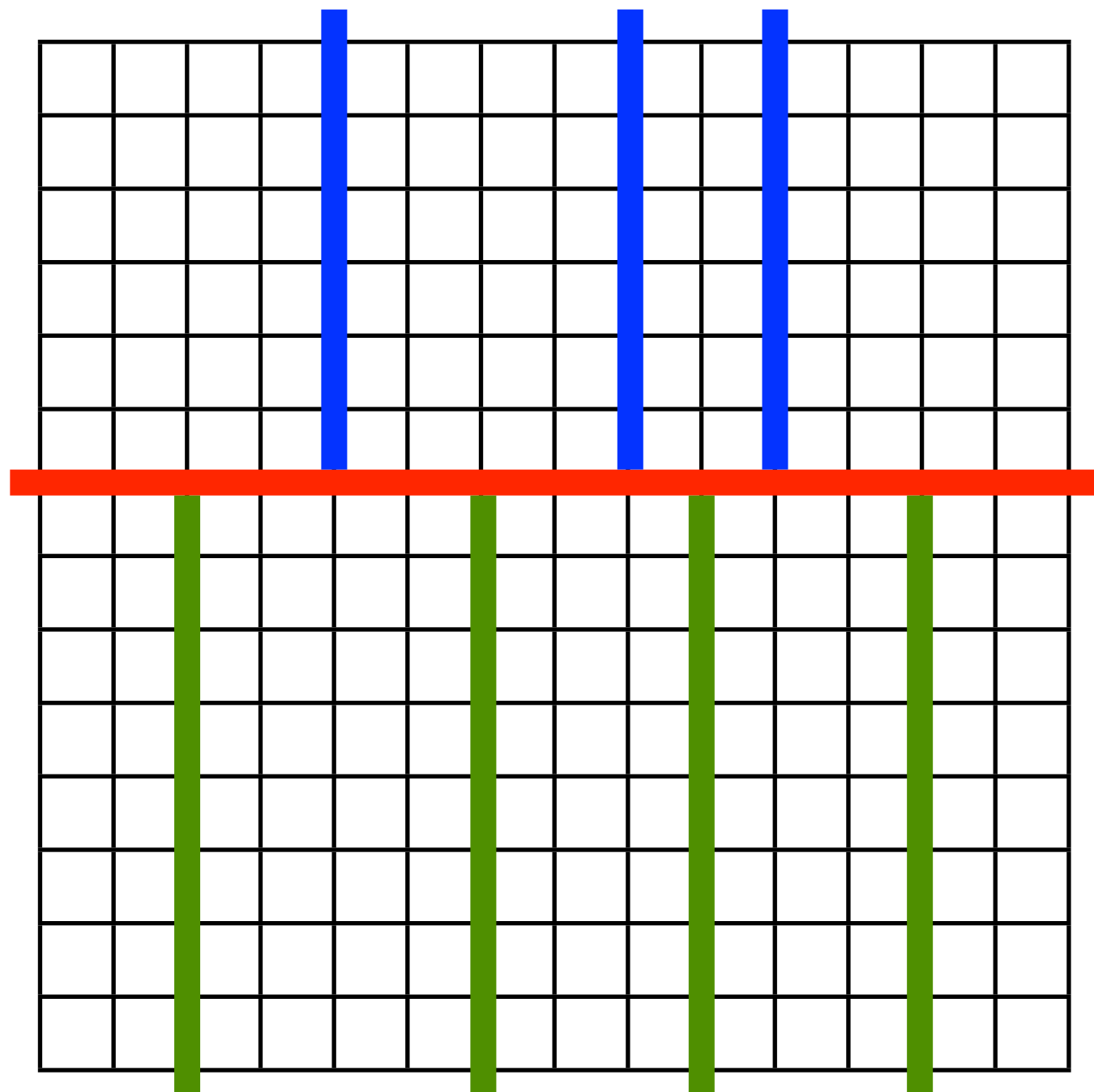


Variables

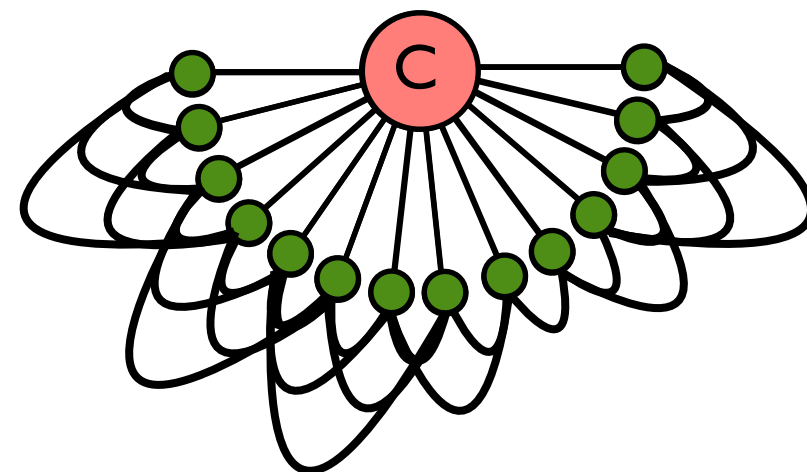


High treewidth

Instances



Variables



High treewidth

Training
set

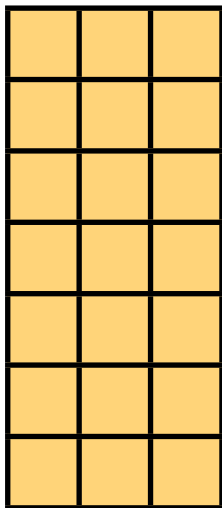
LearnSPN

Establish
approximate
independence

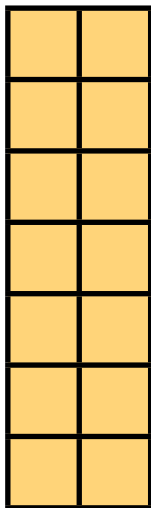
Return:



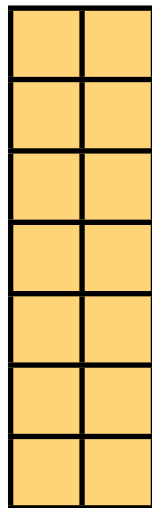
Recurse



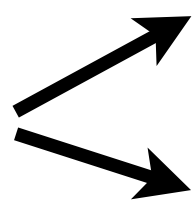
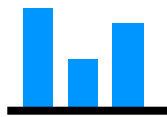
Recurse



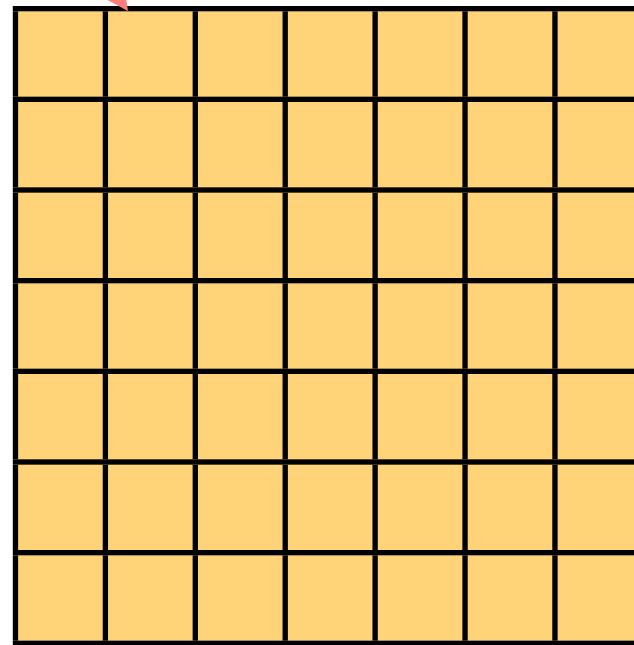
Recurse



If $|V|=1$,
Return:



Instances



Variables

If no
independence,
cluster similar
instances

Return:

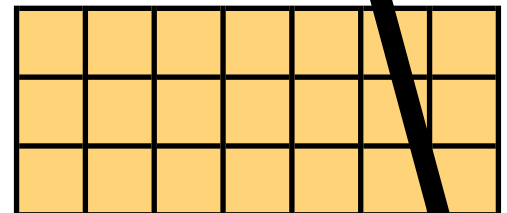


w_1 w_2 w_3

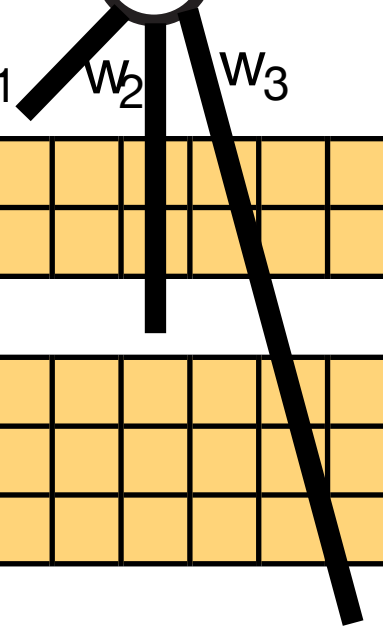
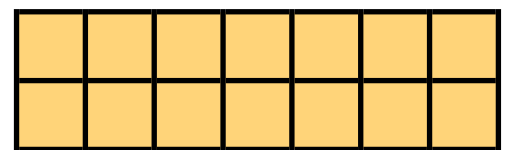
Recurse



Recurse



Recurse



Variable Splits (×)

- Pairwise independence test
- p-value (validation set)

X_6

X_7

X_5

X_1

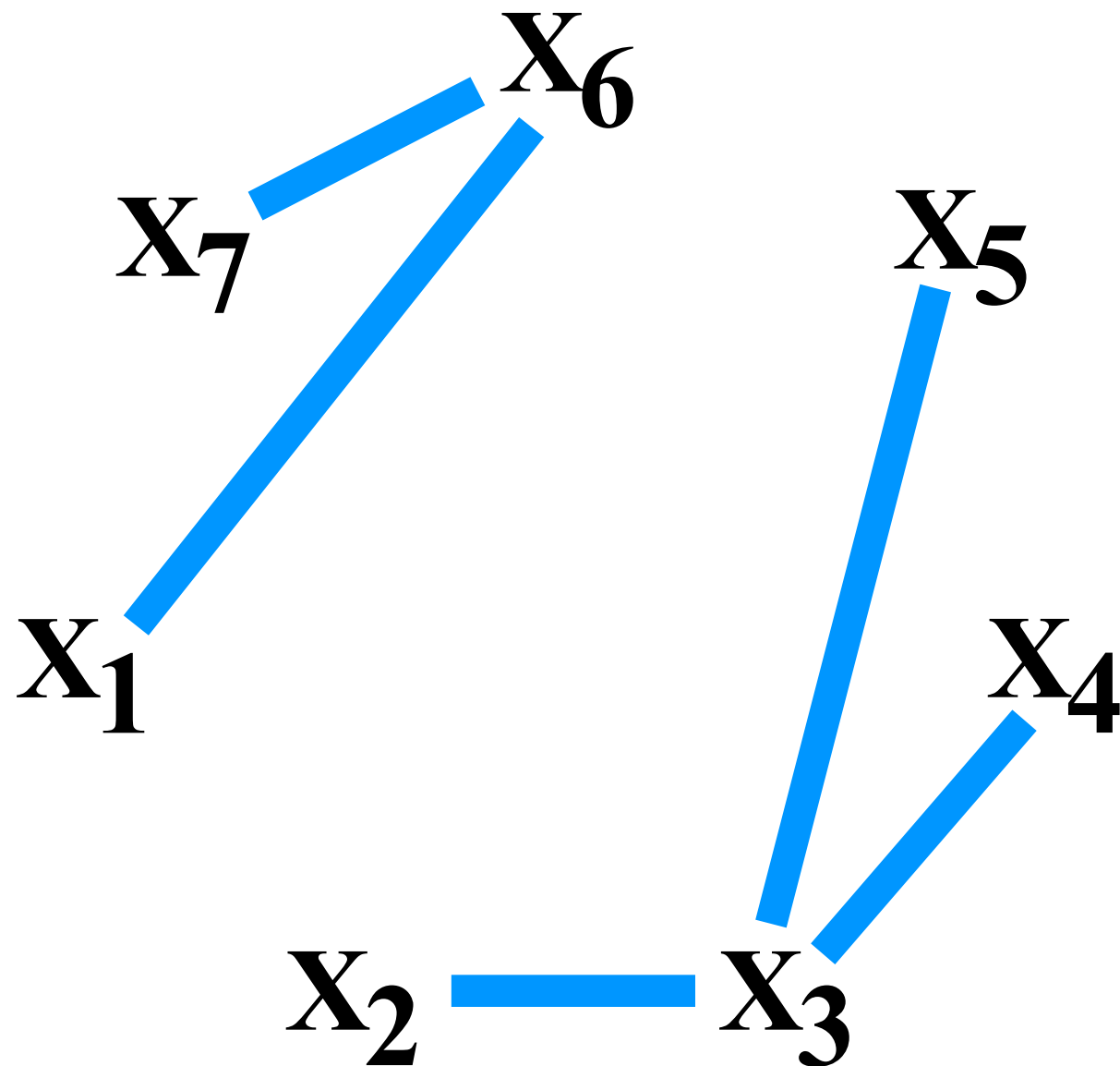
X_4

X_2

X_3

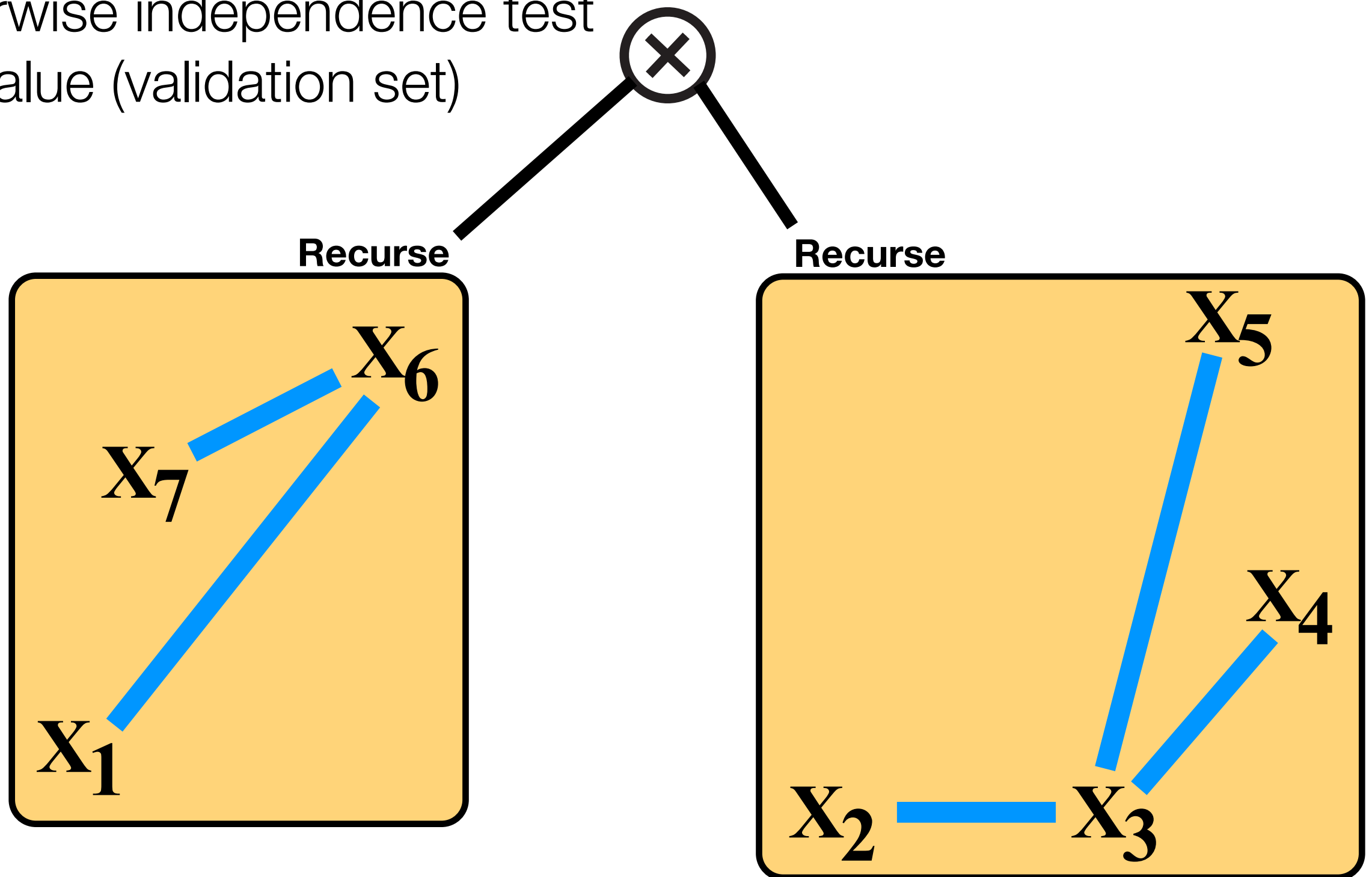
Variable Splits (×)

- Pairwise independence test
- p-value (validation set)



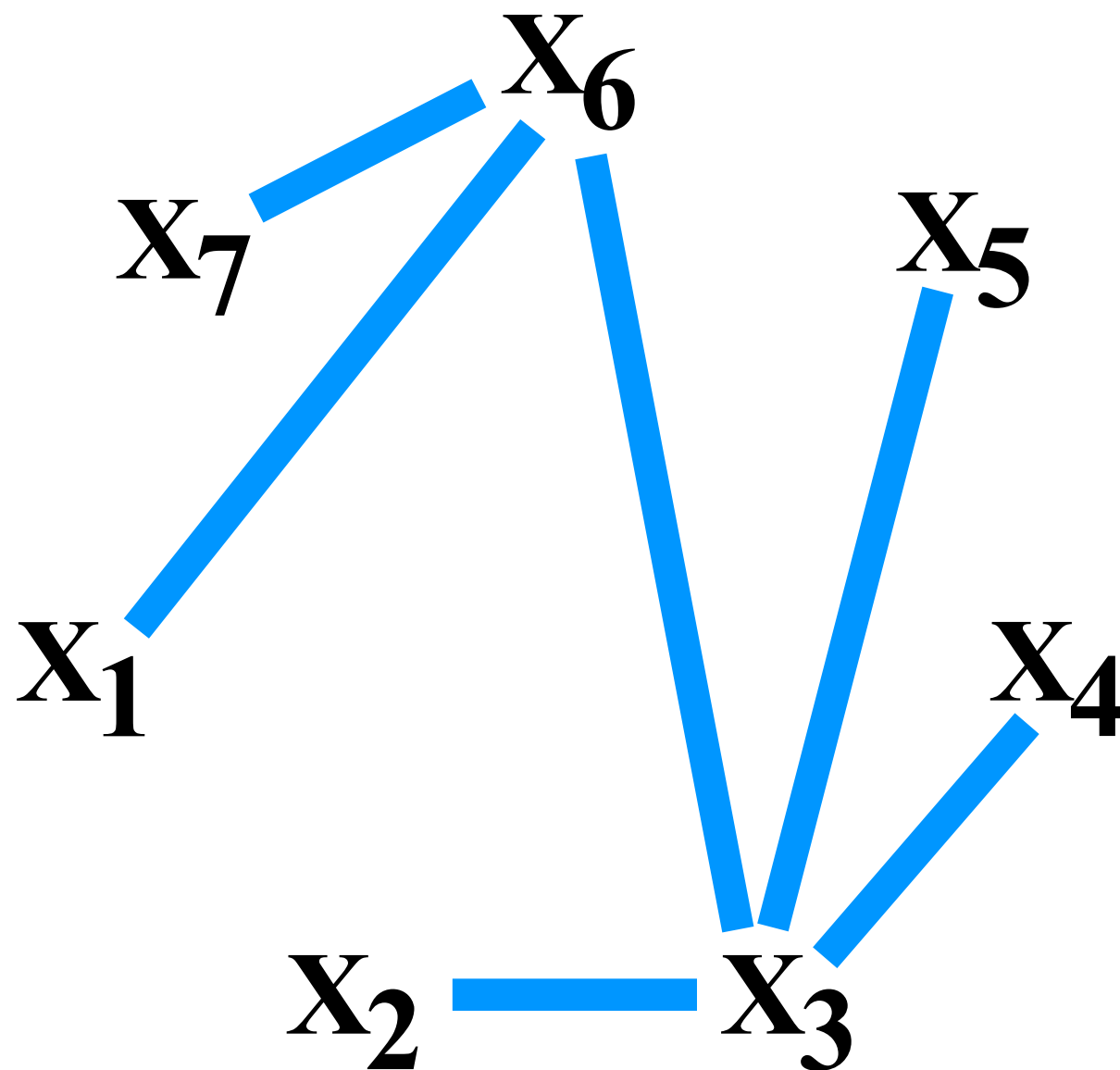
Variable Splits (×)

- Pairwise independence test
- p-value (validation set)

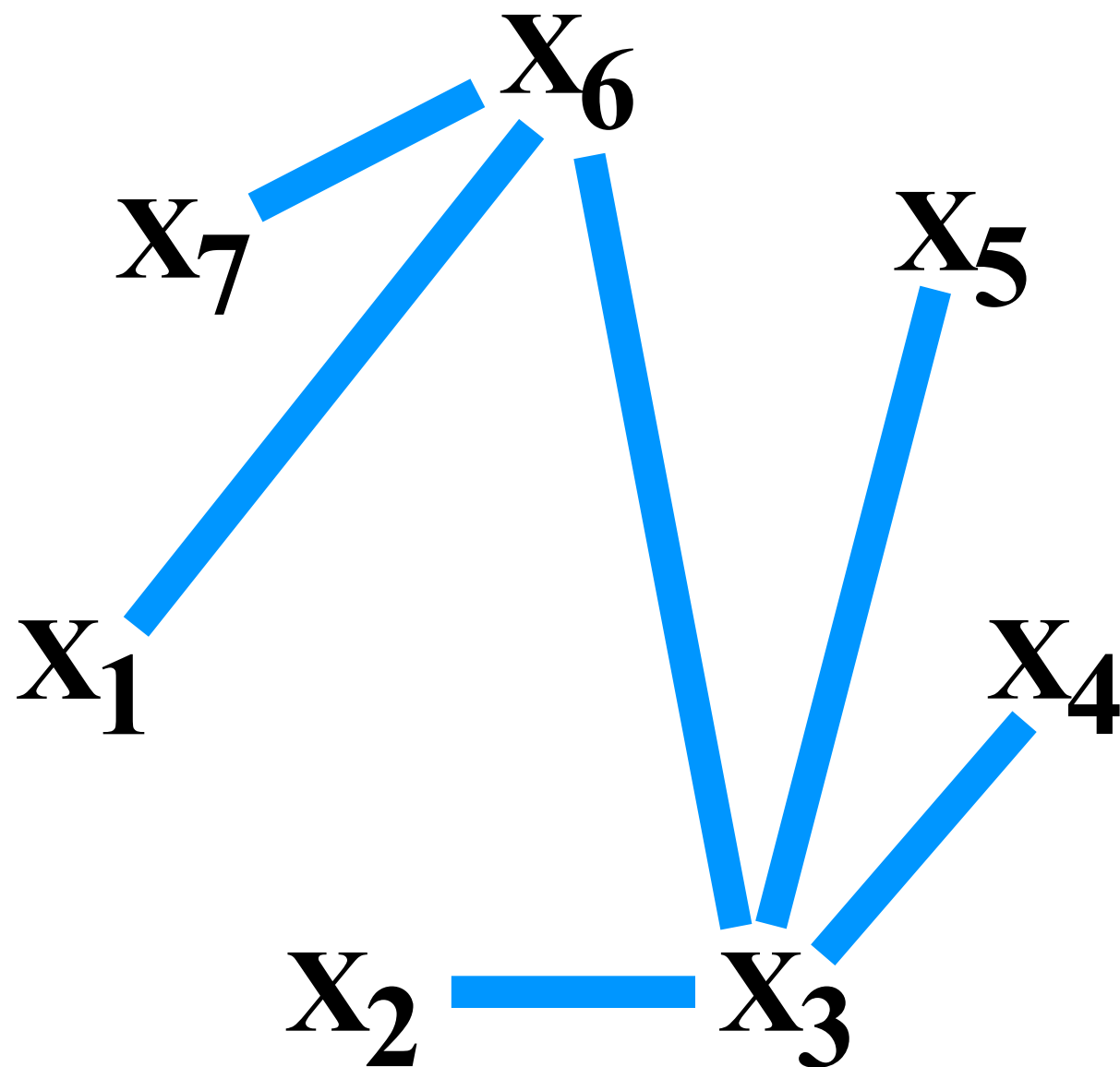


Variable Splits (×)

- Pairwise independence test
- p-value (validation set)



Instance Clustering (+)



Instance Clustering (+)

x_1 x_2 x_3 x_4 x_5 x_6 x_7

Instance Clustering (+)

$\mathbf{m_1}$							
$\mathbf{m_2}$							
$\mathbf{m_3}$							
$\mathbf{m_4}$							
$\mathbf{m_5}$							
$\mathbf{m_6}$							
$\mathbf{m_7}$							
	$\mathbf{X_1}$	$\mathbf{X_2}$	$\mathbf{X_3}$	$\mathbf{X_4}$	$\mathbf{X_5}$	$\mathbf{X_6}$	$\mathbf{X_7}$

Instance Clustering (+)

m_1

m_2

m_3

m_4

m_5

m_6

m_7

Instance Clustering (+)

- Online hard EM
- Naive Bayes mixture model
- Cluster penalty (validation set)
- Learned weights are the mixture priors

$$P(V) = \sum_i P(C_i) \prod_j P(X_j | C_i)$$

m₁

m₂

m₃

m₄

m₅

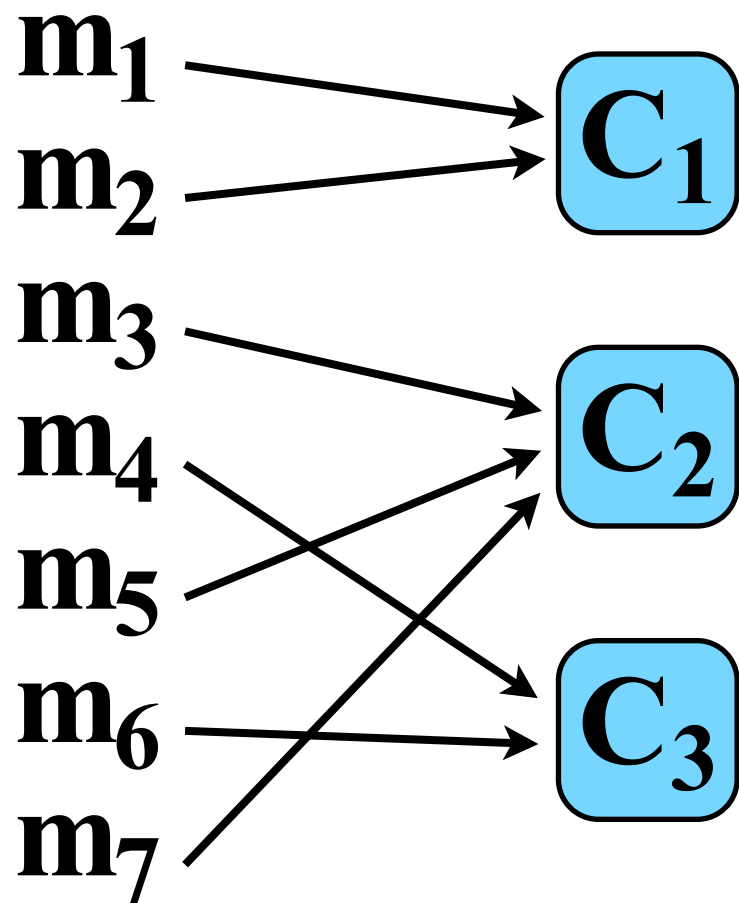
m₆

m₇

Instance Clustering (+)

- Online hard EM
- Naive Bayes mixture model
- Cluster penalty (validation set)
- Learned weights are the mixture priors

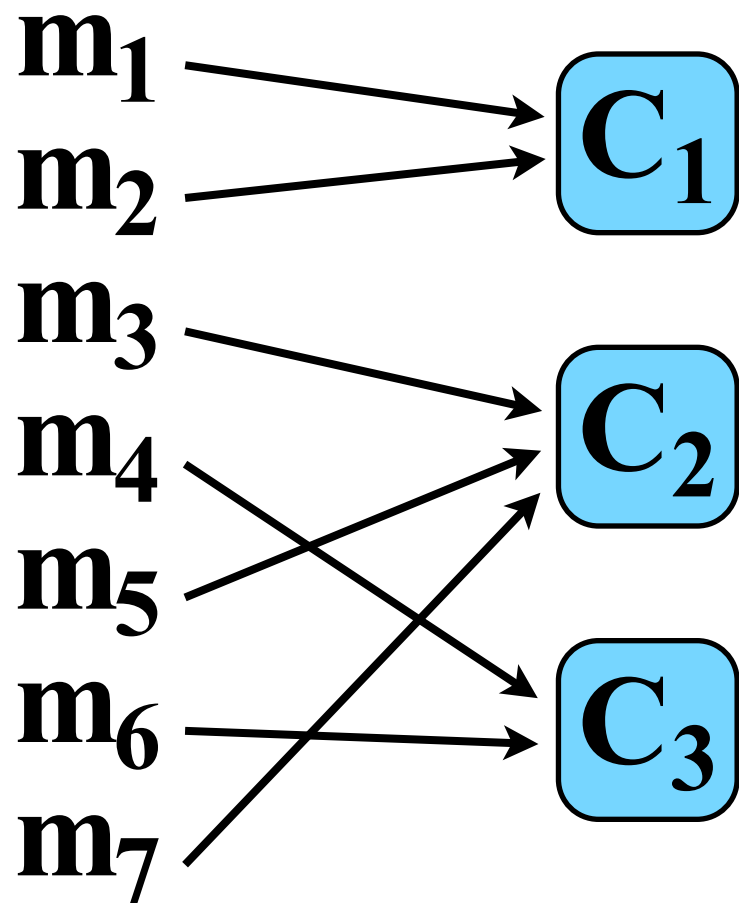
$$P(V) = \sum_i P(C_i) \prod_j P(X_j | C_i)$$



Instance Clustering (+)

- Online hard EM
- Naive Bayes mixture model
- Cluster penalty (validation set)
- Learned weights are the mixture priors

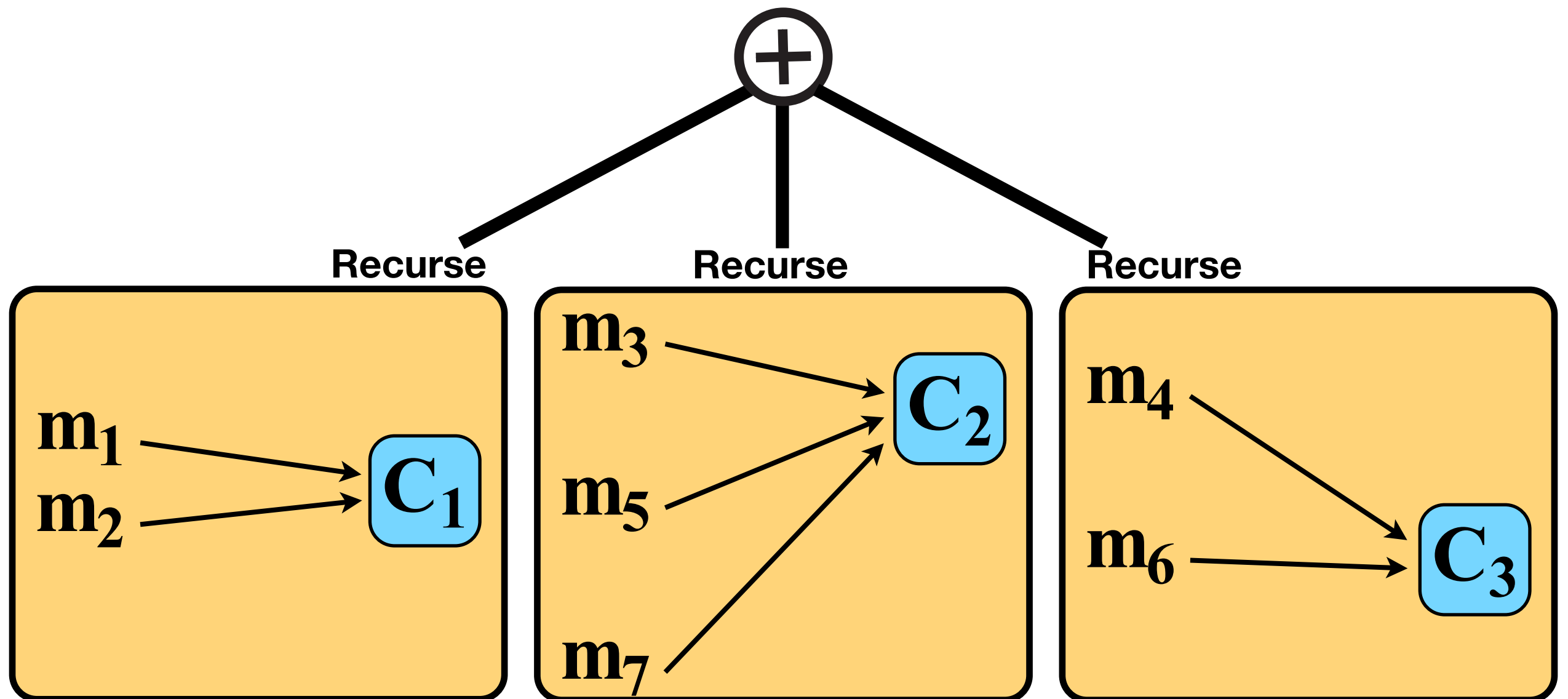
$$P(V) = \sum_i P(C_i) \prod_j P(X_j | C_i)$$



Instance Clustering (+)

- Online hard EM
- Naive Bayes mixture model
- Cluster penalty (validation set)
- Learned weights are the mixture priors

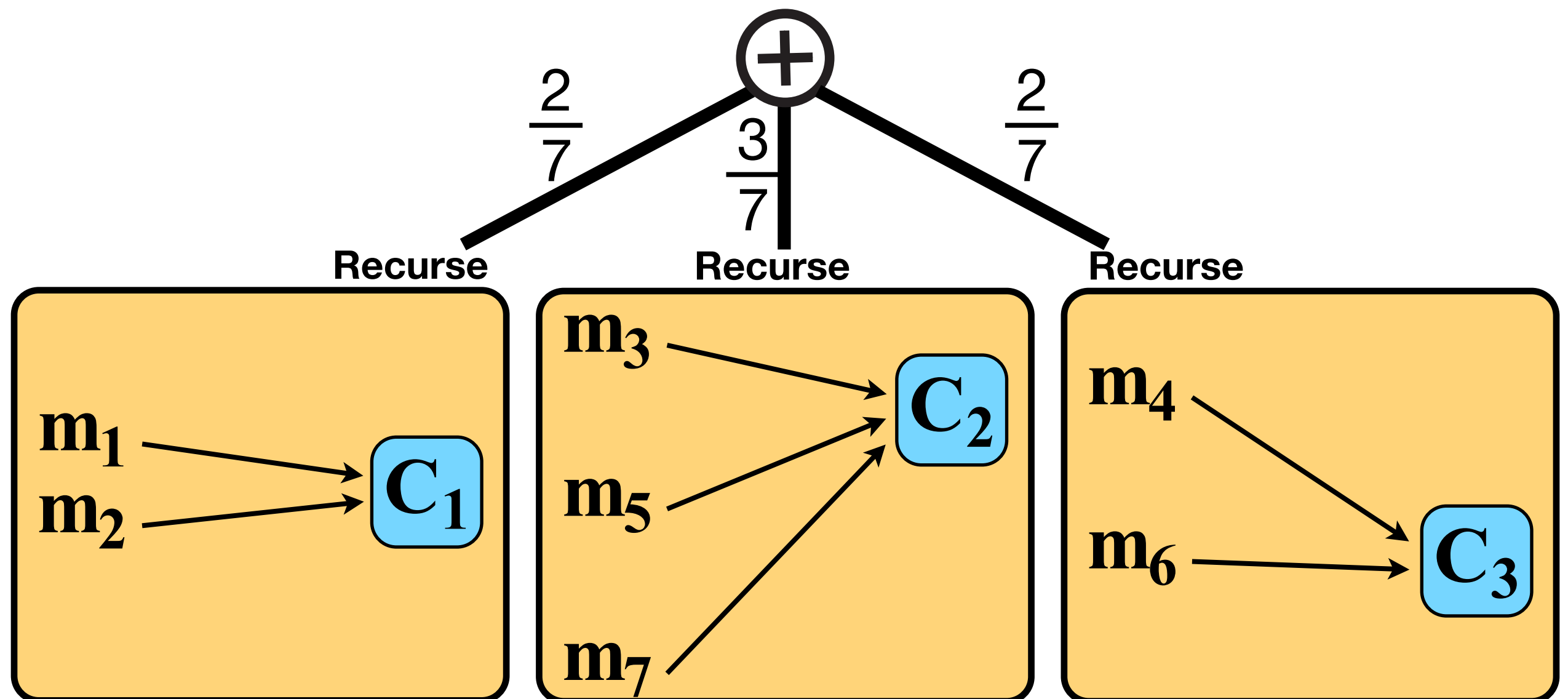
$$P(V) = \sum_i P(C_i) \prod_j P(X_j | C_i)$$



Instance Clustering (+)

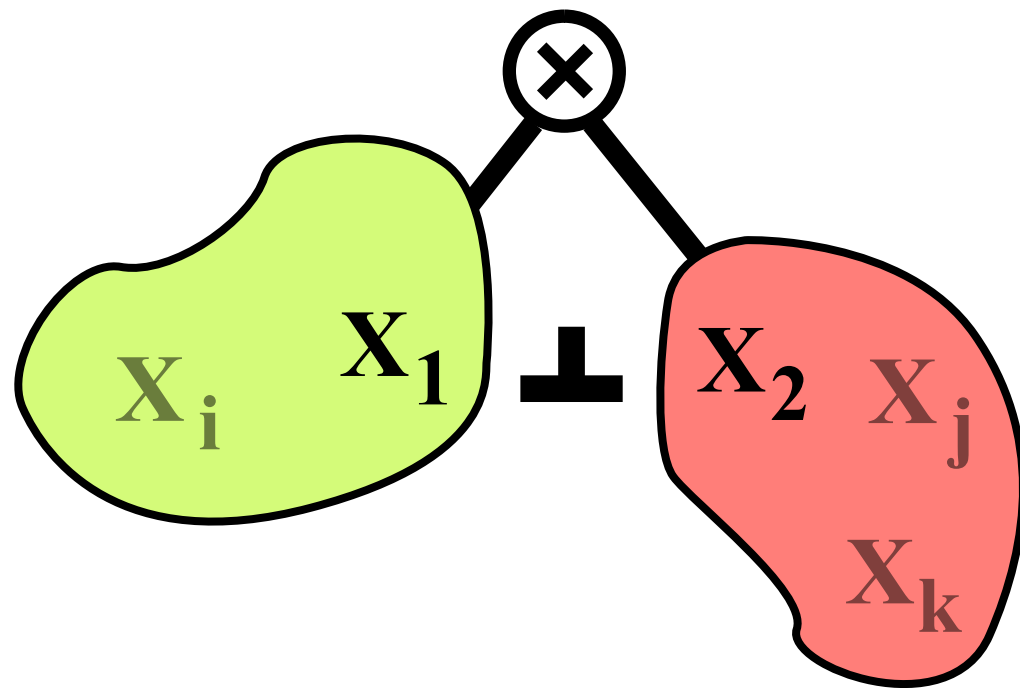
- Online hard EM
- Naive Bayes mixture model
- Cluster penalty (validation set)
- Learned weights are the mixture priors

$$P(V) = \sum_i P(C_i) \prod_j P(X_j|C_i)$$



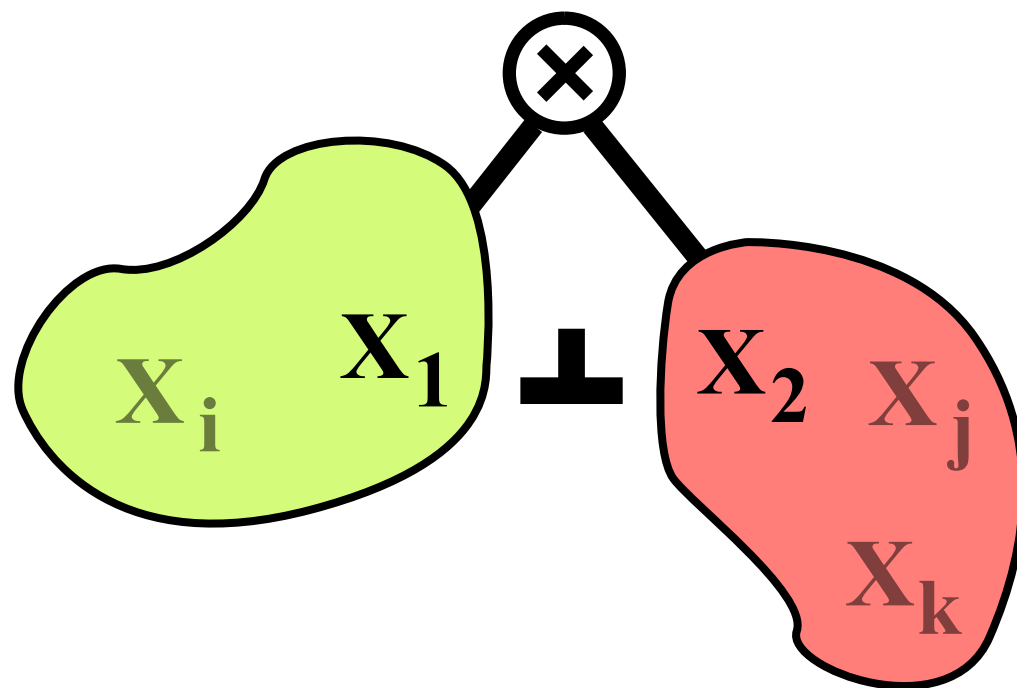
LearnSPN Locally Optimizes Likelihood

LearnSPN Locally Optimizes Likelihood



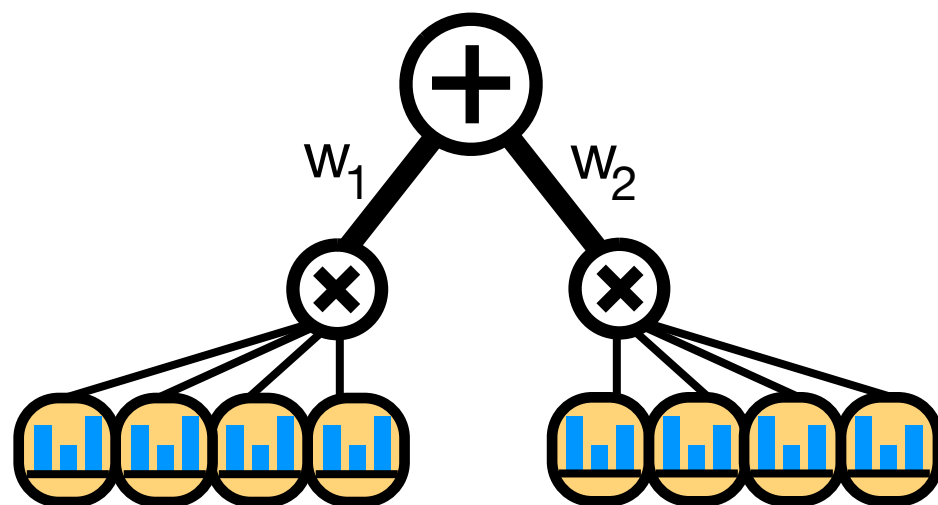
No loss of likelihood
if truly independent

LearnSPN Locally Optimizes Likelihood



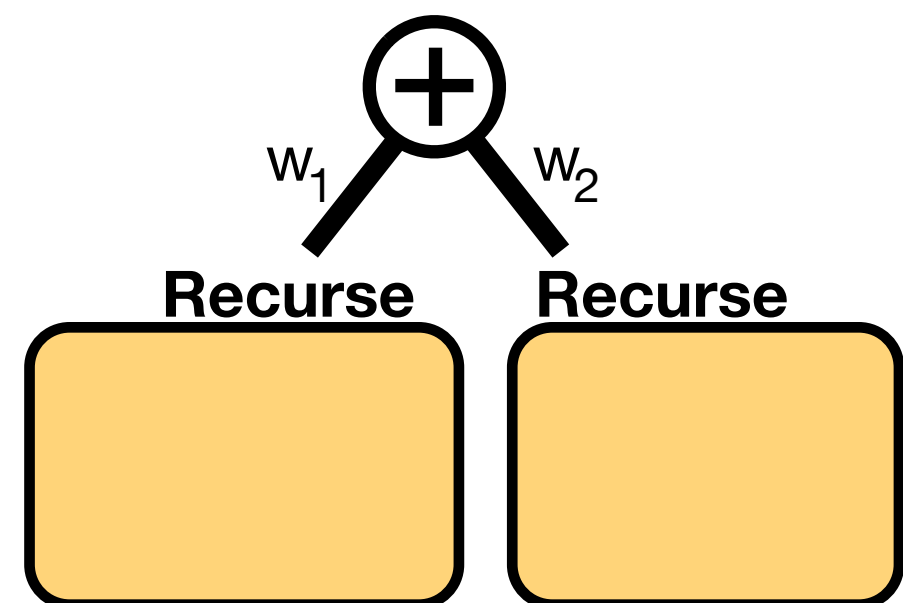
No loss of likelihood
if truly independent

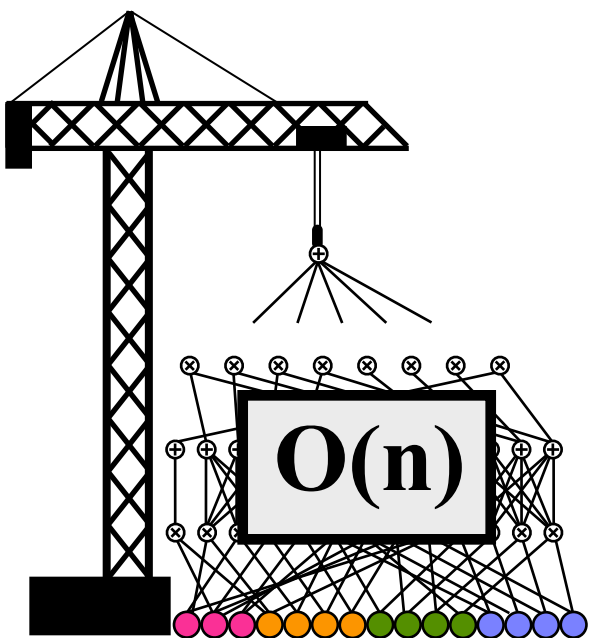
Naive Bayes likelihood



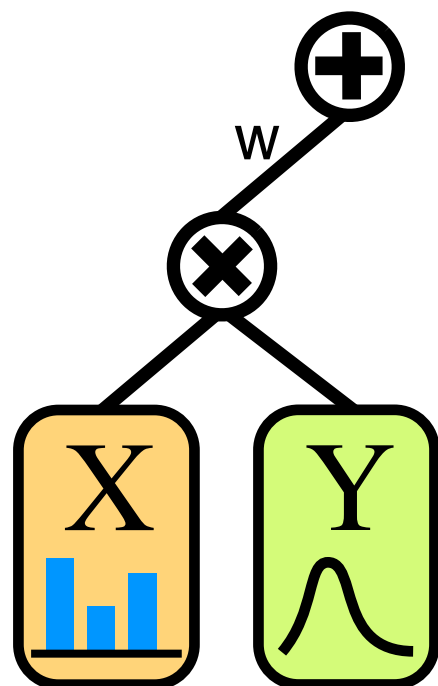
\geq

LearnSPN likelihood

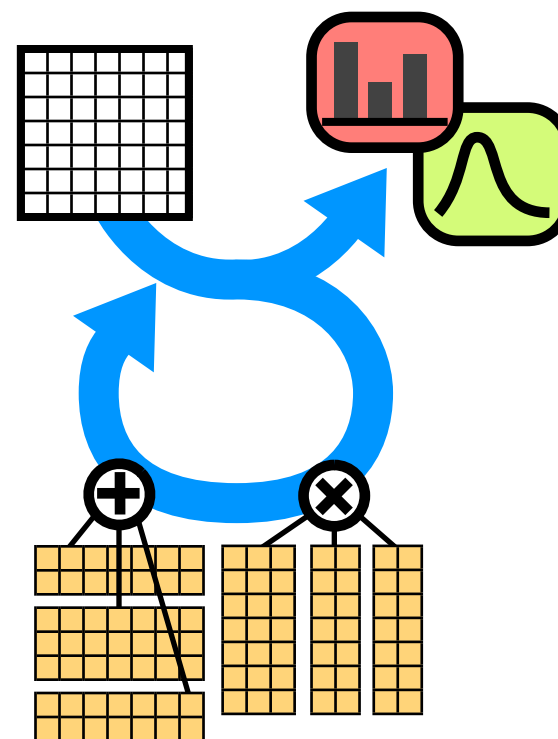




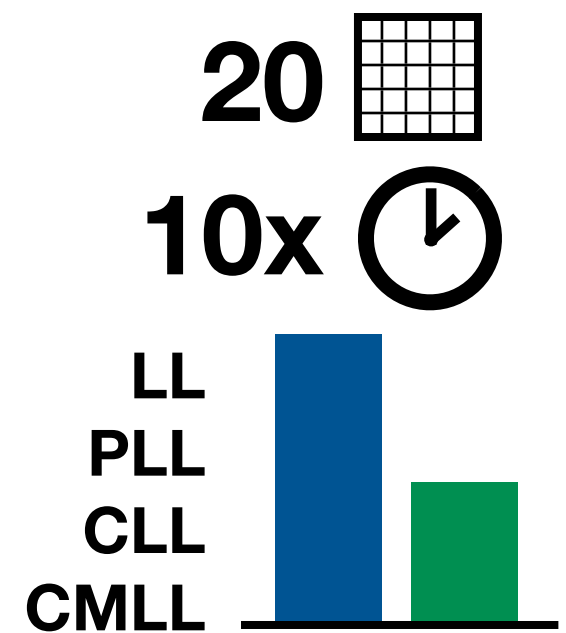
Motivation



**SPN
Review**



**Structure
Learning**



Experiments

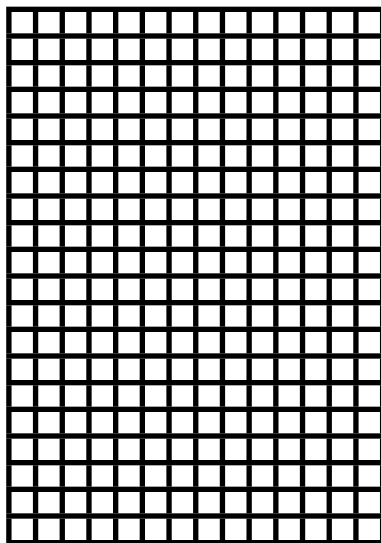
Experiments

20 Datasets

collaborative filtering
click-through logs
nucleic acid sequences
...

Instances

2k-388k



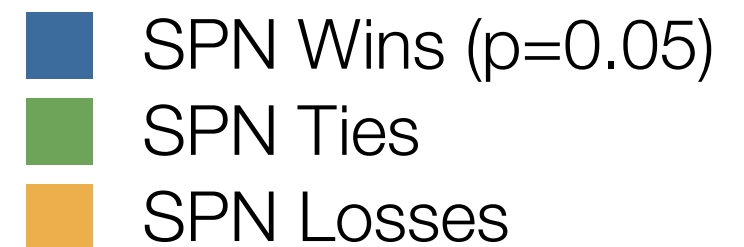
Variables

16-1556

Representation	Learning	Inference
SPN	LearnSPN	Exact
Bayesian Network	WinMine	Gibbs
		Loopy BP
Markov Network	Della Pietra	Gibbs
		Loopy BP
	L1	Gibbs
		Loopy BP

Total: 1680 Experiments

Learning Results



WinMine

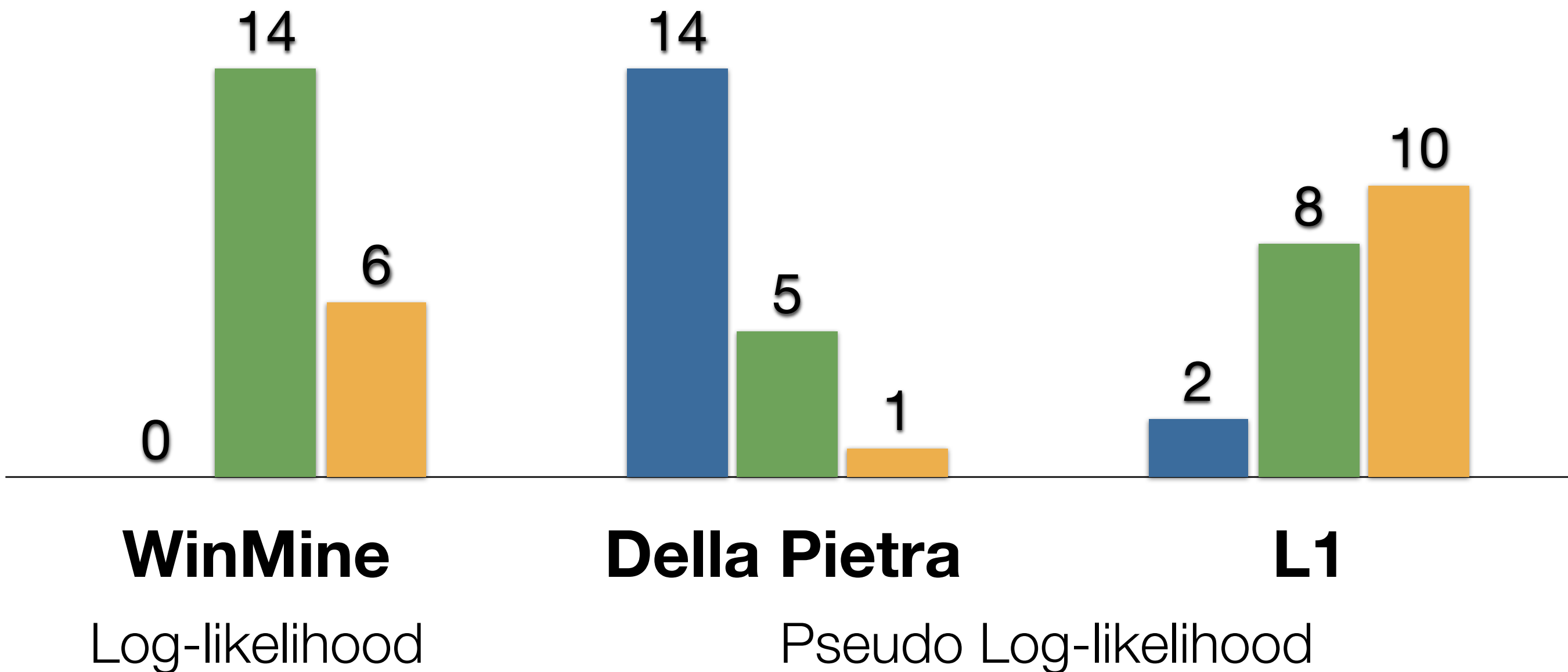
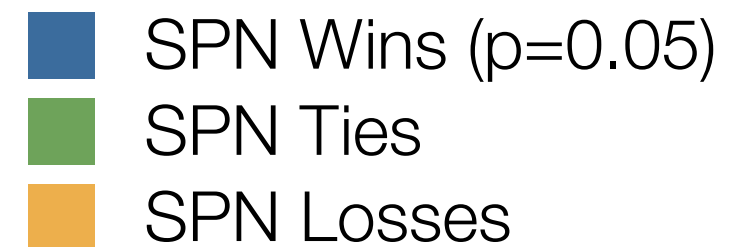
Log-likelihood

Della Pietra

Pseudo Log-likelihood

L1

Learning Results



Inference

Inference

$$P(\textit{Query} \mid \textit{Evidence})$$

Inference

$$P(\textit{Query} \mid \textit{Evidence})$$

Q	E
10%	30%
20%	30%
30%	30%
40%	30%
50%	30%
30%	0%
30%	10%
30%	20%
30%	40%
30%	50%

For each proportion,
generate 1000
queries from test set

Inference

$$P(\textit{Query} \mid \textit{Evidence})$$

Q	E
10%	30%
20%	30%
30%	30%
40%	30%
50%	30%
30%	0%
30%	10%
30%	20%
30%	40%
30%	50%

$\left\{ \begin{array}{c} \text{WinMine} \\ \text{Della Pietra} \\ \text{L1} \end{array} \right\} \times 20 \text{ Datasets} \times 10 \text{ Variable proportions}$

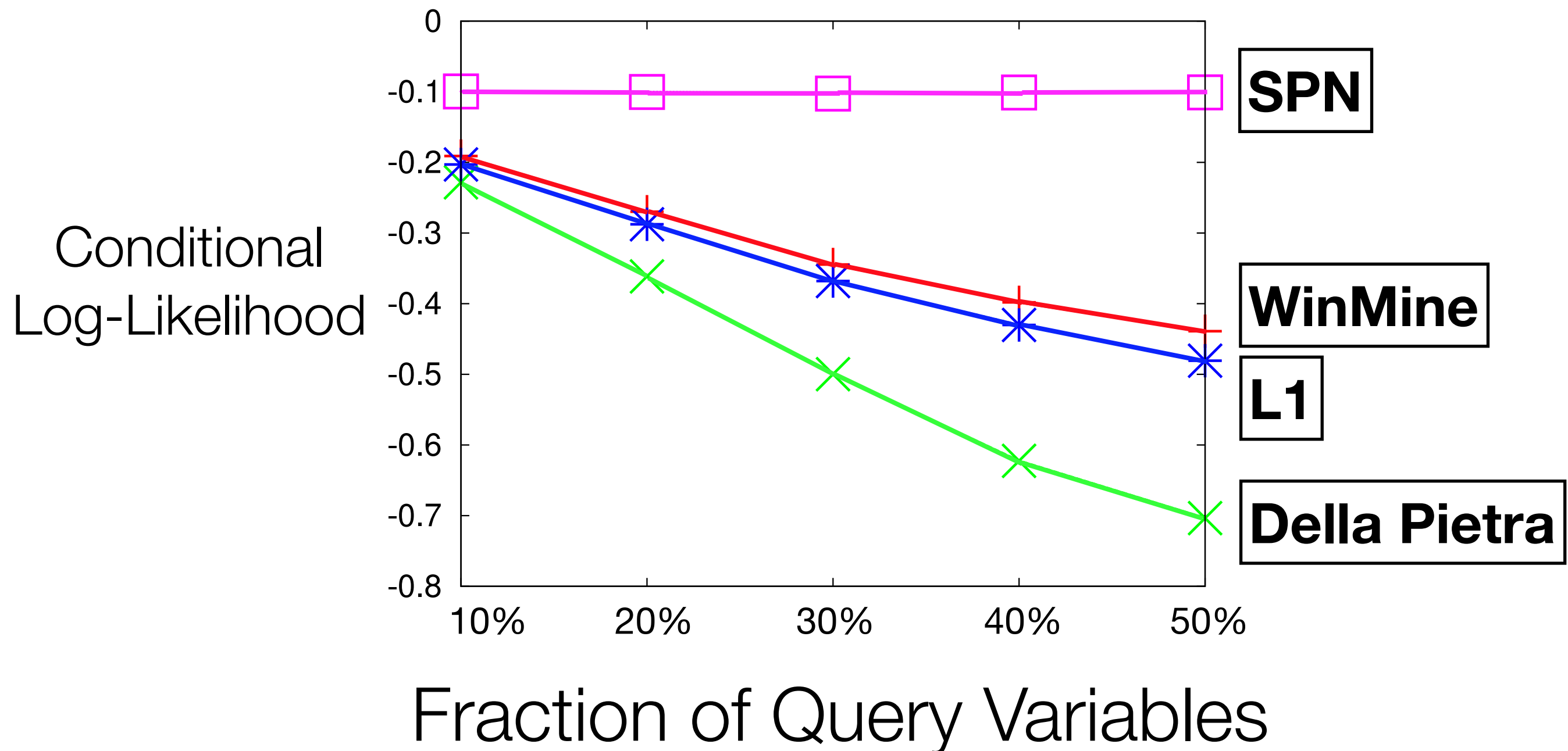
$\times \left\{ \begin{array}{c} \text{Gibbs} \\ \text{Loopy BP} \end{array} \right\} = 1200 \text{ Experiments}$

(In addition to 400 for SPNs)

For each proportion,
generate 1000
queries from test set

Inference Accuracy

“EachMovie”, 500 variables



Inference Time

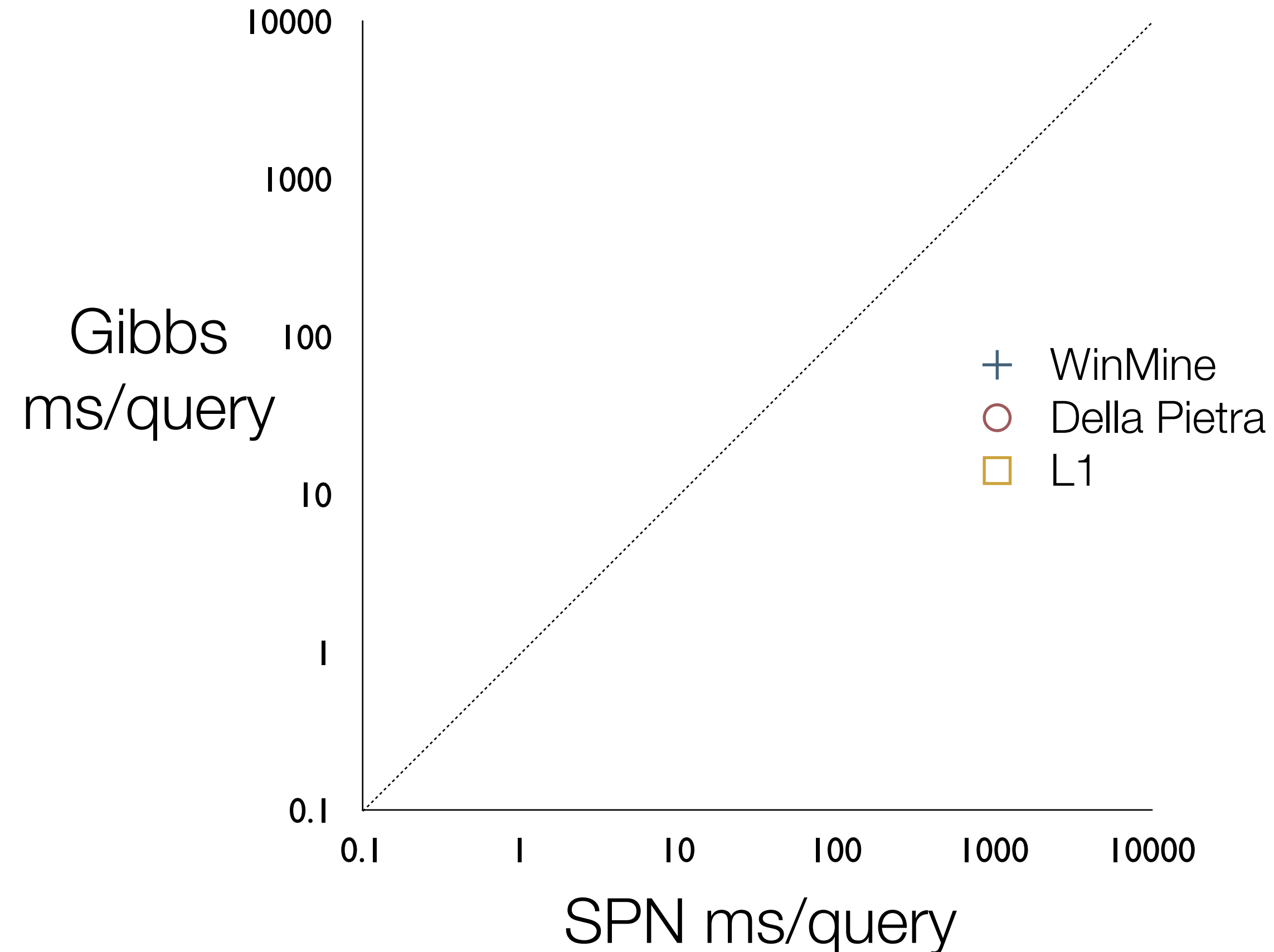
Gibbs
ms/query

+ WinMine
○ Della Pietra
□ L1

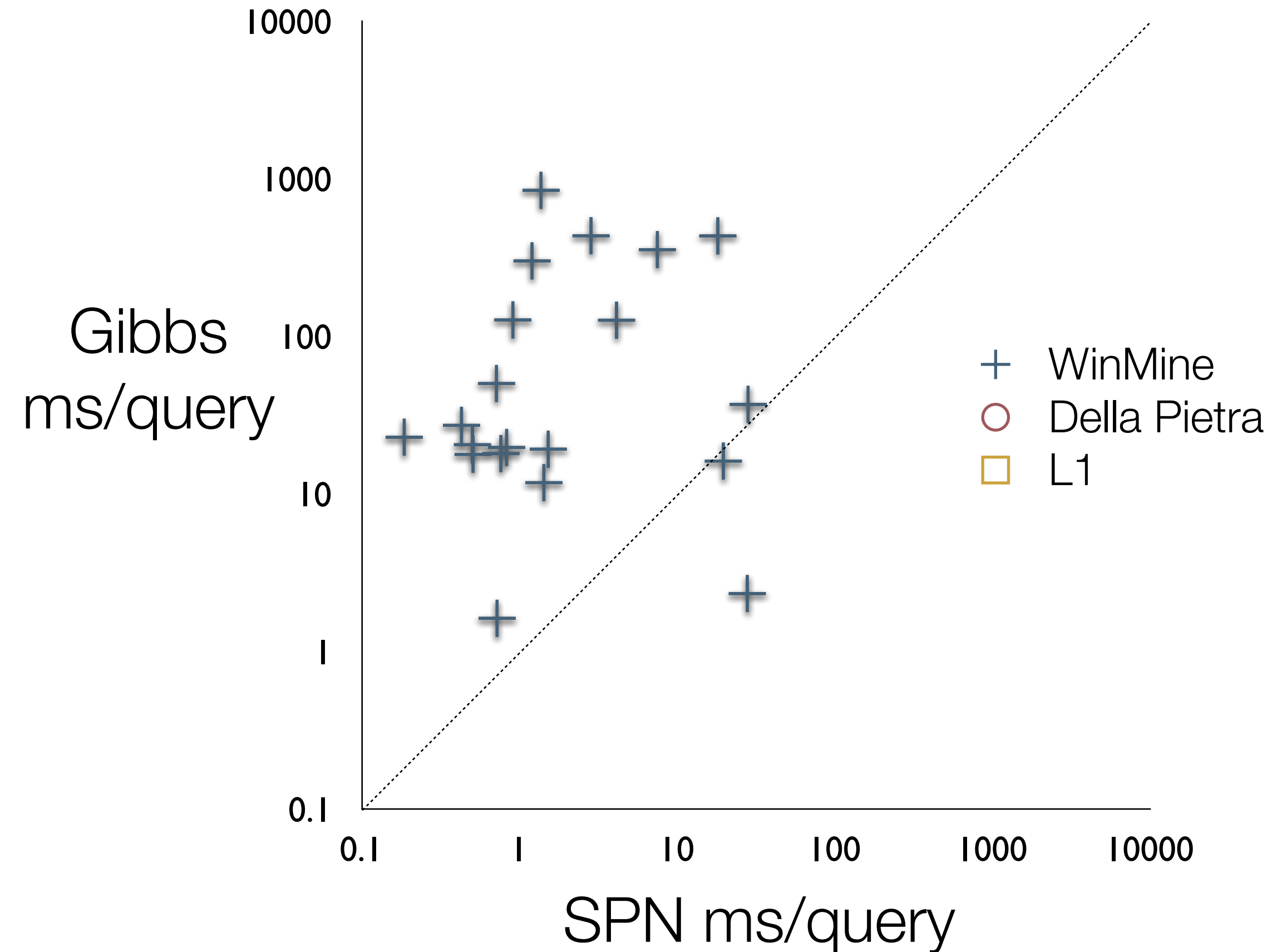
SPN ms/query



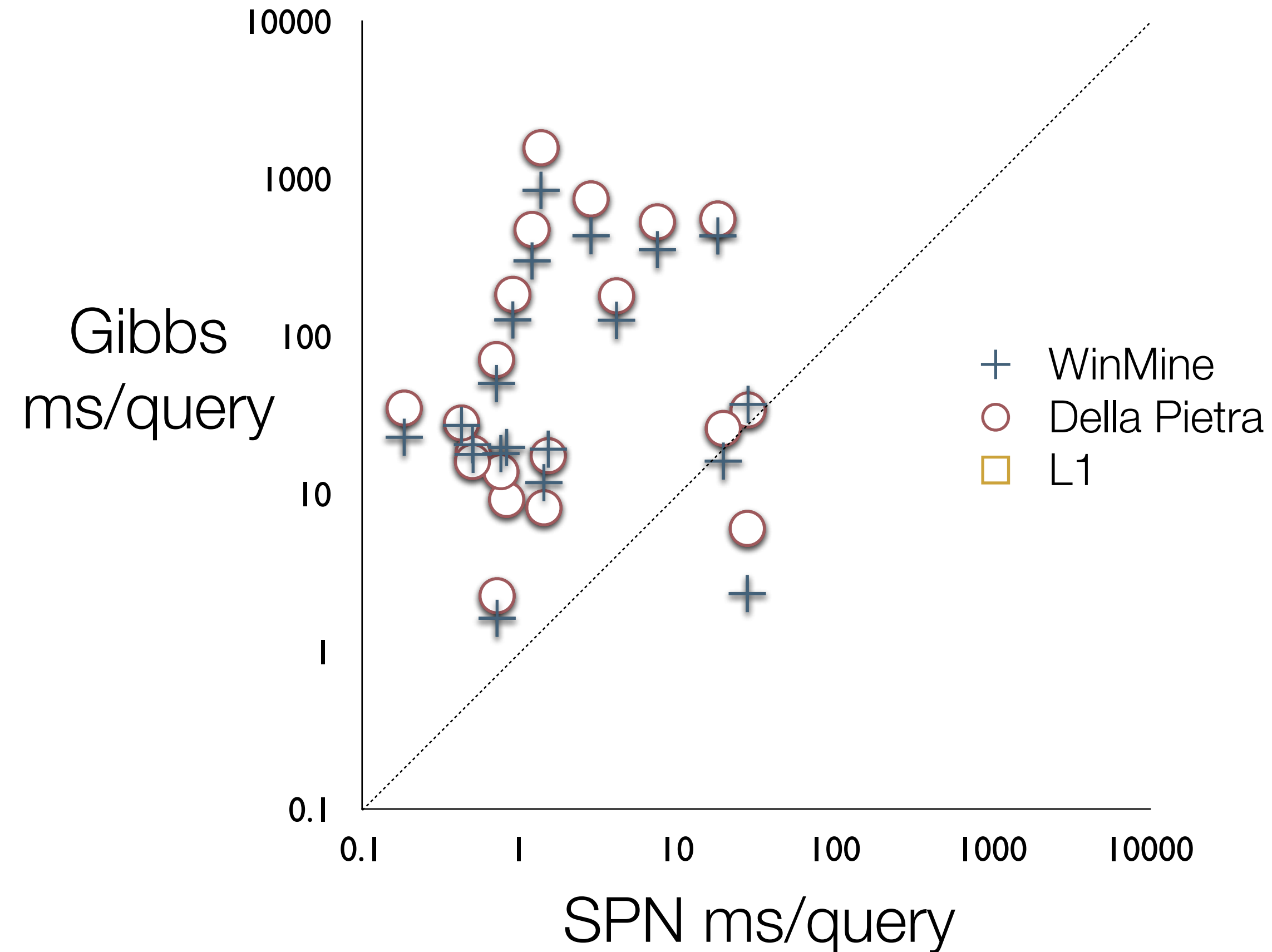
Inference Time



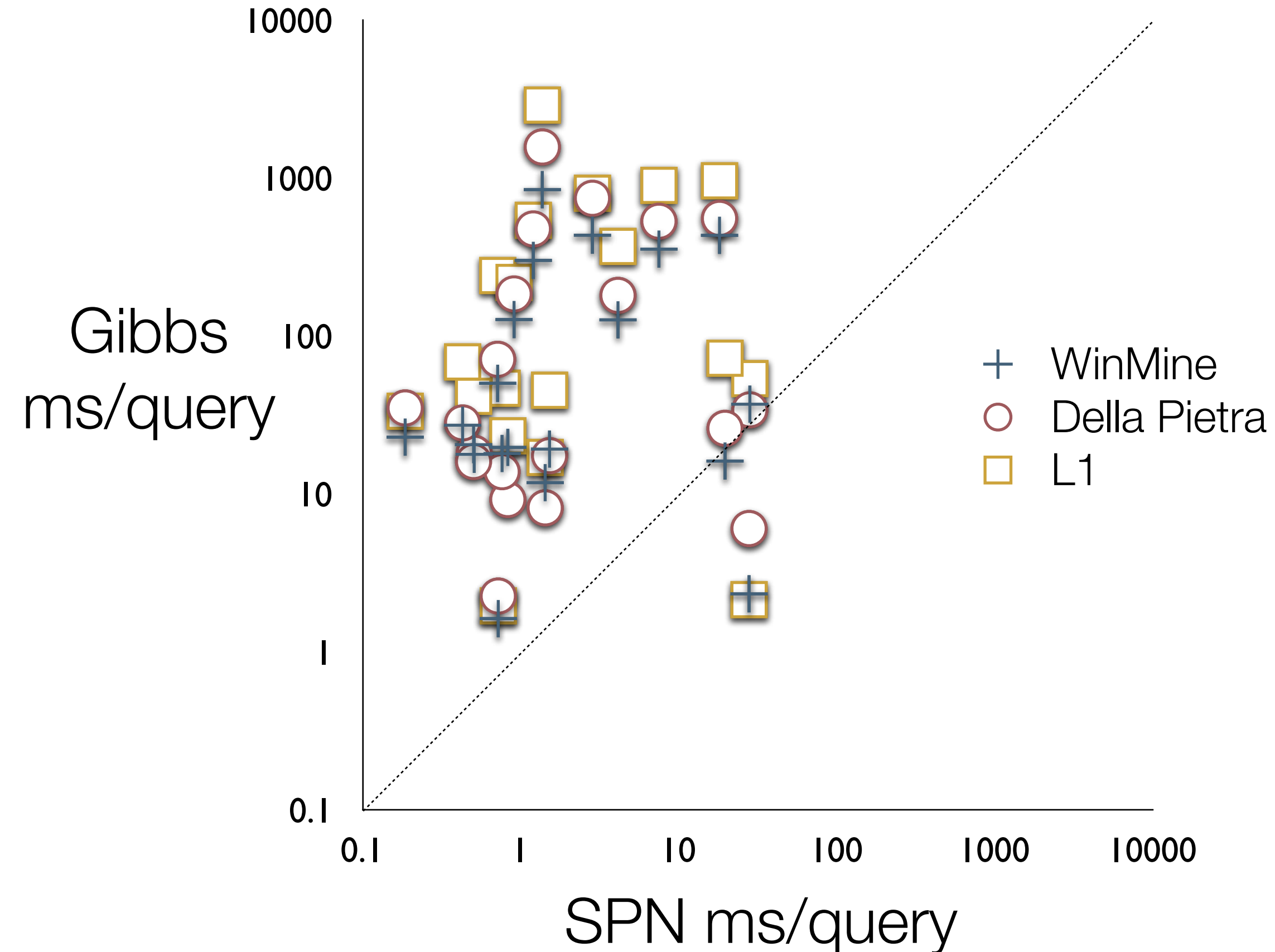
Inference Time



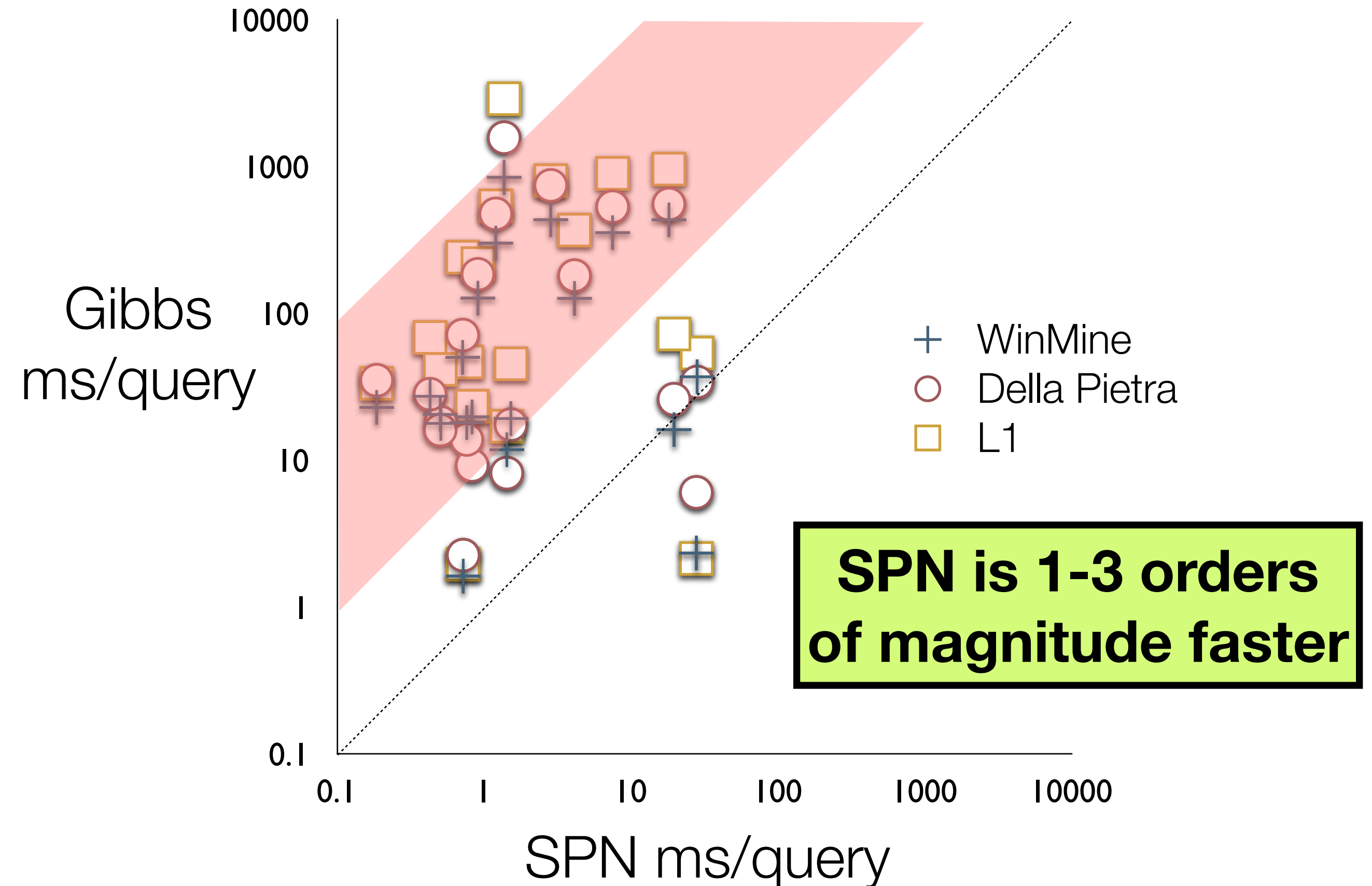
Inference Time



Inference Time



Inference Time



Inference Time

Loopy Belief Propagation

Single variable marginals

$$\left\{ \begin{array}{l} \text{WinMine} \\ \text{Della Pietra} \\ \text{L1} \end{array} \right\} \times 20 \text{ Datasets} \times 10 \text{ Variable proportions} \\ = 600 \text{ Experiments}$$

SPNs had higher conditional marginal log-likelihood on 78% of experiments

(95% higher CLL vs. Gibbs)

SPN ms/query

Experiment Conclusions

- SPN learning accuracy is comparable
- SPN exact inference is 1-3 orders of magnitude faster, and more accurate
- Inference does not involve tuning settings or diagnosing convergence
- Inference takes a predictable amount of time
- Can now apply SPNs to many domains

Future Work

- Other SPN structure and weight learning algorithms
- Approximating intractable distributions with SPNs
- Parallelizing SPN learning and inference

Code and supplemental results available at
spn.cs.washington.edu/learnspn/