

Naive Bayesian Classifier (20.2.2)

Given a vector $V = \{v_{11}, v_{22}, v_{31}, \dots, v_{m2}, c?\}$

Compute:

$P(c_1 | v_{11}, v_{22}, v_{31}, \dots, v_{m2})$ proportional to

$P(v_{11} | v_{22}, v_{31}, \dots, v_{m2}, c_1) P(v_{22} | v_{31}, \dots, v_{m2}, c_1) \dots P(v_{m2} | c_1) P(c_1)$
equals (under assumption that V_i 's are independent conditioned
on C)

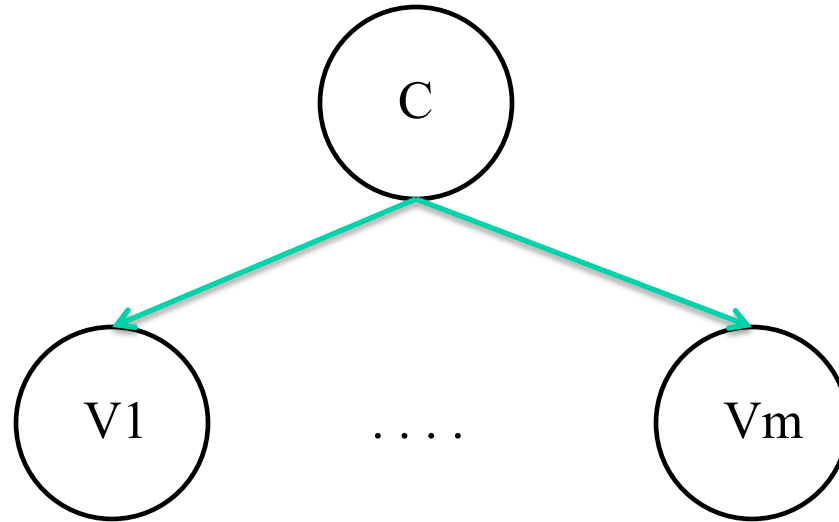
$$P(v_{11} | c_1) P(v_{22} | c_1) P(v_{31} | c_1) \dots P(v_{m2} | c_1) P(c_1)$$

$P(c_2 | v_{11}, v_{22}, v_{31}, \dots, v_{m2})$ proportional to

$P(v_{11} | v_{22}, v_{31}, \dots, v_{m2}, c_2) P(v_{22} | v_{31}, \dots, v_{m2}, c_2) \dots P(v_{m2} | c_2) P(c_2)$
equals (under assumption that V_i 's are independent conditioned
on C)

$$P(v_{11} | c_2) P(v_{22} | c_2) P(v_{31} | c_2) \dots P(v_{m2} | c_2) P(c_2)$$

Classify V as in c_1 or c_2 , whichever yields higher probability



$P(C | V_1, V_2, V_3, \dots, V_m)$ proportional to

$$P(V_1, V_2, V_3, V_m, C) =$$

$$P(V_1|C)P(V_2|C)P(V_3|C)\dots P(V_m|C)P(C)$$

Given a vector $V = \{1, -1, 0, \dots, 1\}$

Compute:

$P(-1 | 1, -1, 0, \dots, 1)$ as ~~$P(1|-1)P(-1|-1)P(0|-1)\dots P(1|-1)P(-1)$~~

$P(1 | 1, -1, 0, \dots, 1)$ as ~~$P(1|1)P(-1|1)P(0|1)\dots P(1|1)P(1)$~~

Classify V as in $c1$ or $c2$, whichever yields higher probability

Learning a Naïve Bayesian Classifier.

View probabilities as proportions computed over training set.

$$P(v_{11}|c_1)P(v_{22}|c_1)P(v_{31}|c_1)\dots P(v_{m2}|c_1)P(c_1)$$

$$= \frac{[v_{11},c_1]}{[c_1]} * \frac{[v_{22},c_1]}{[c_1]} * \frac{[v_{31},c_1]}{[c_1]} * \dots * \frac{[v_{m2},c_1]}{[c_1]} * \frac{[c_1]}{[]}$$

where $[conditions]$ is the number of objects/rows in the training set that satisfy all the $conditions$. So $[v_{11},c_1]$ is the number of training data that are members of c_1 and have $V_1=v_{11}$, $[c_1]$ is the number of training objects in c_1 , $[]$ is the total number of training objects.

Learning in this case, is a matter of counting the number of rows in training data in which various conditions satisfied. What conditions? Each class/variable-value pair, each class, total number of rows.

$$\begin{array}{cccccc}
 & V1 & V2 & V3 & \dots & Vm \\
 c1 & \left[\begin{array}{ccccc} [v11,c1] & [v21,c1] & [v31,c1] & \dots & [vm1,c1] \\ [v12,c1] & [v22,c1] & [v32,c1] & \dots & [vm2,c1] \end{array} \right] & & & & [c1] \\
 c2 & \left[\begin{array}{ccccc} [v11,c2] & [v21,c2] & [v31,c2] & \dots & [vm1,c2] \\ [v12,c2] & [v22,c2] & [v32,c2] & \dots & [vm2,c2] \end{array} \right] & & & & [c2] \\
 & [v11] & [v21] & [v31] & & [vm1] \\
 & [v12] & [v22] & [v32] & & [vm2] \\
 & & & & & []
 \end{array}$$

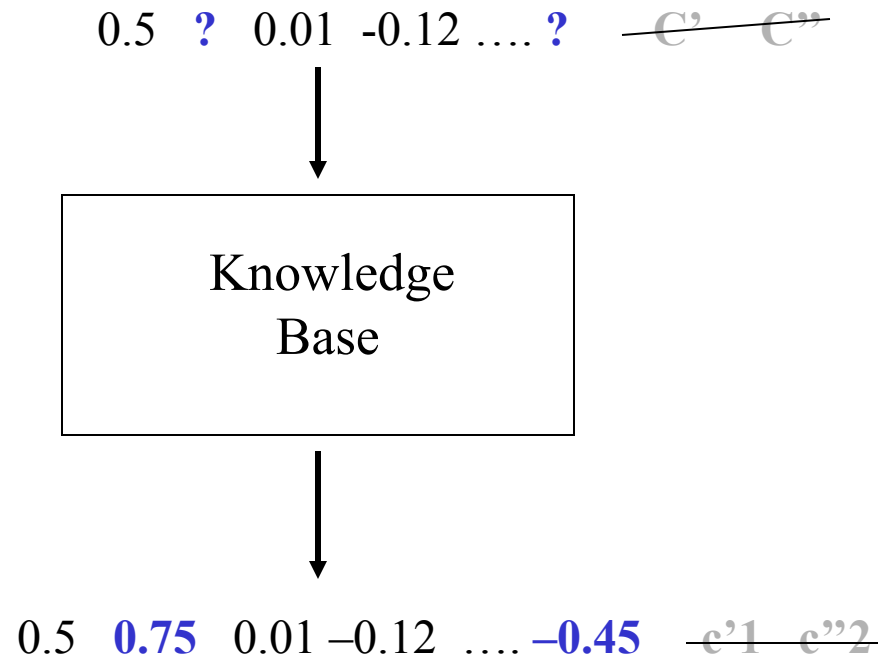
Consider an (multidimensional) array implementation of int, and estimate $P(v_{ij}|c_k)$ as $([v_{ij},c_k]+1) / ([c_k]+2)$, and $P(c_k)$ as $([c_k]+1) / ([]+2)$

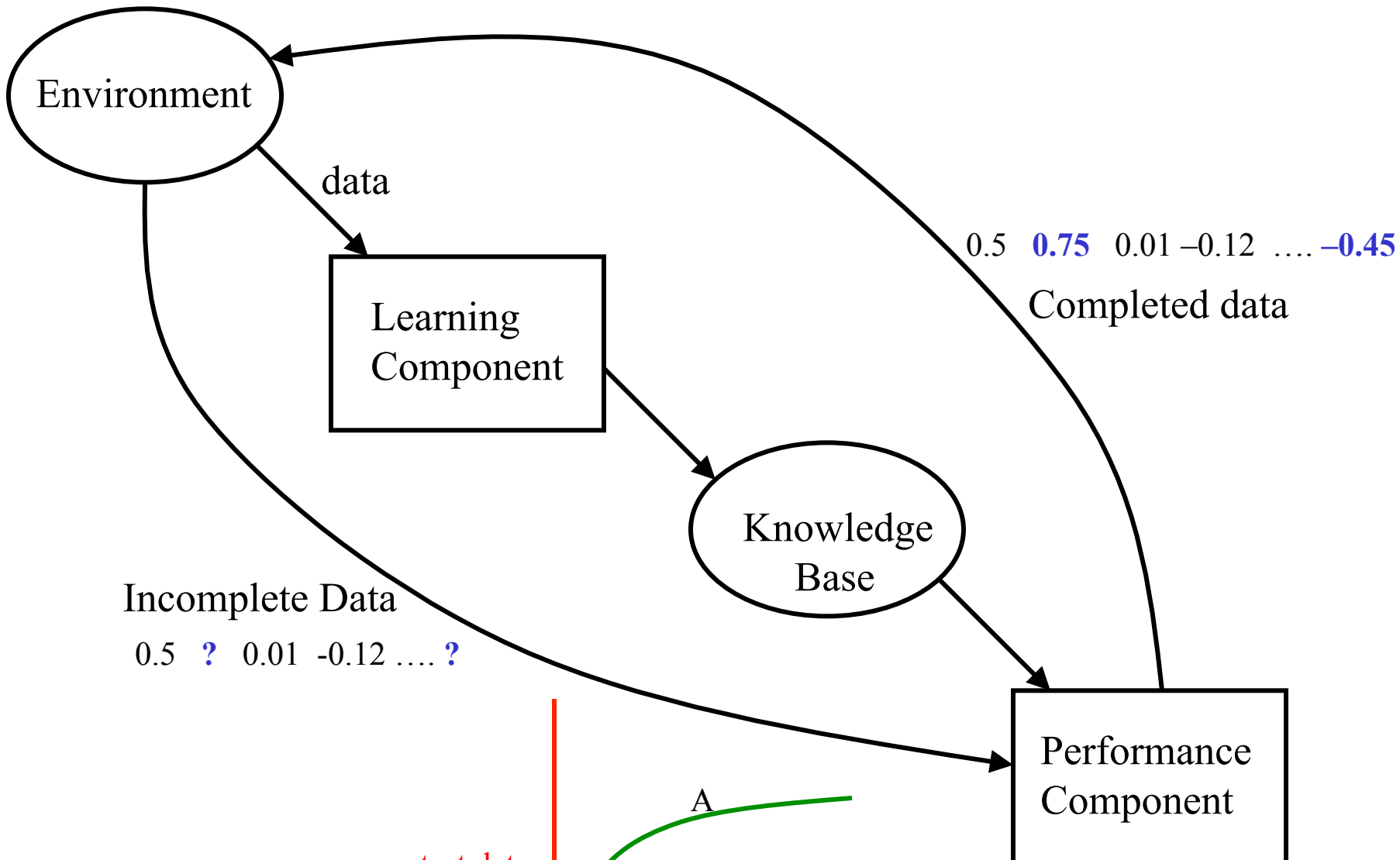
Number of classes

Number of V_i values

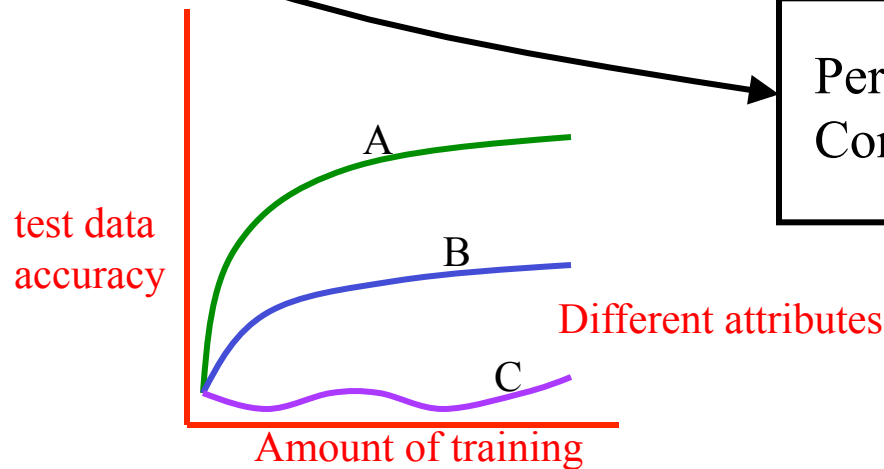
“various tricks are used” to avoid probabilities of 0

Unsupervised Performance Task: Pattern Completion





Incomplete Data
0.5 ? 0.01 -0.12 ?



Example: Unsupervised rule induction of Association Rules (market-basket analysis)

In a nutshell: run “brute force” rule discovery for all possible consequents, not simply single variable values (e.g., $V1=v12$), but consequents that are conjunctions of variable values (e.g., $V1=v12 \ \& \ V4=v42 \ \& \ V5=v51$).

Retain rules $A \rightarrow C$ such that $P(A \ \& \ C) \geq T1$ and $P(C|A) \geq T2$. These thresholds enable pruning of the search space (A and C are themselves conjunctions).

Problem: a plethora of rules, most uninteresting, are produced.

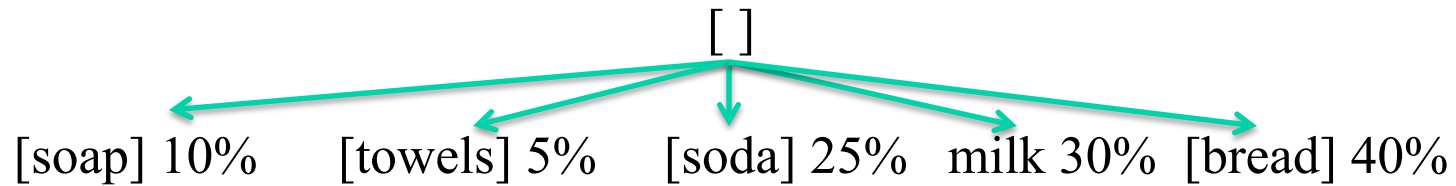
Solutions: Organize/prune rules by

- a) Interestingness (e.g., $A \rightarrow C$ interesting if $P(A, C) \gg P(A)P(C)$ or $\ll P(A)P(C)$)
- b) confidence (a confidence interval around coverage and/or accuracy)
- c) support for top-level goal

A Priori Algorithm

Prune conjunctions of features that fall below threshold support

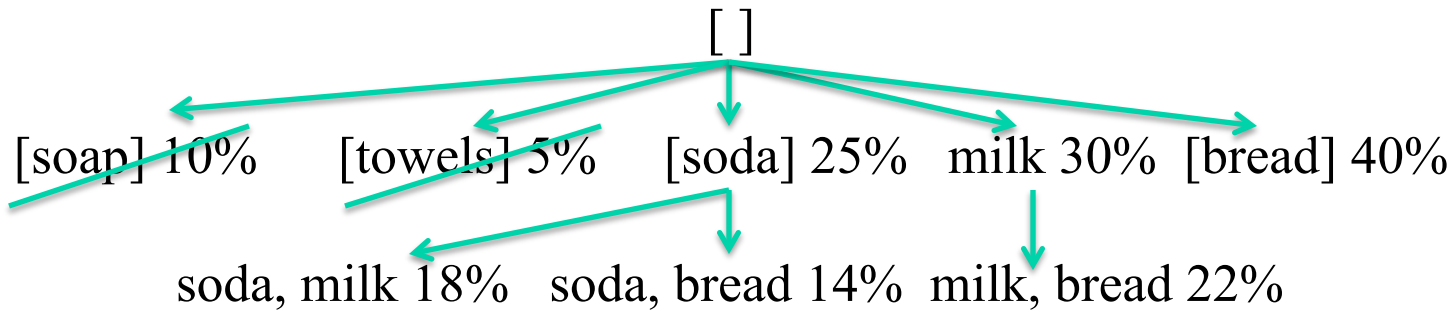
T1 = 20%



A Priori Algorithm

Prune conjunctions of features that fall below threshold support

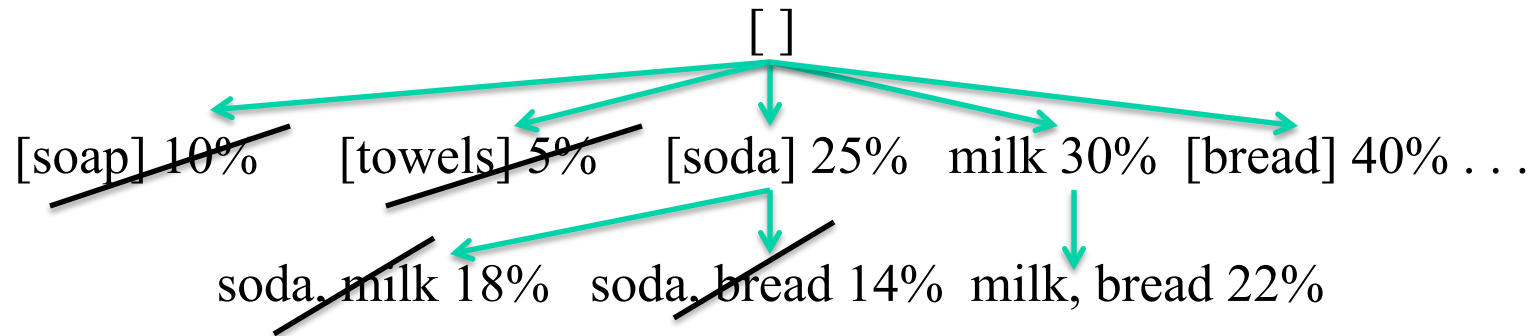
T1 = 20%



A Priori Algorithm

Prune conjunctions of features that fall below threshold support

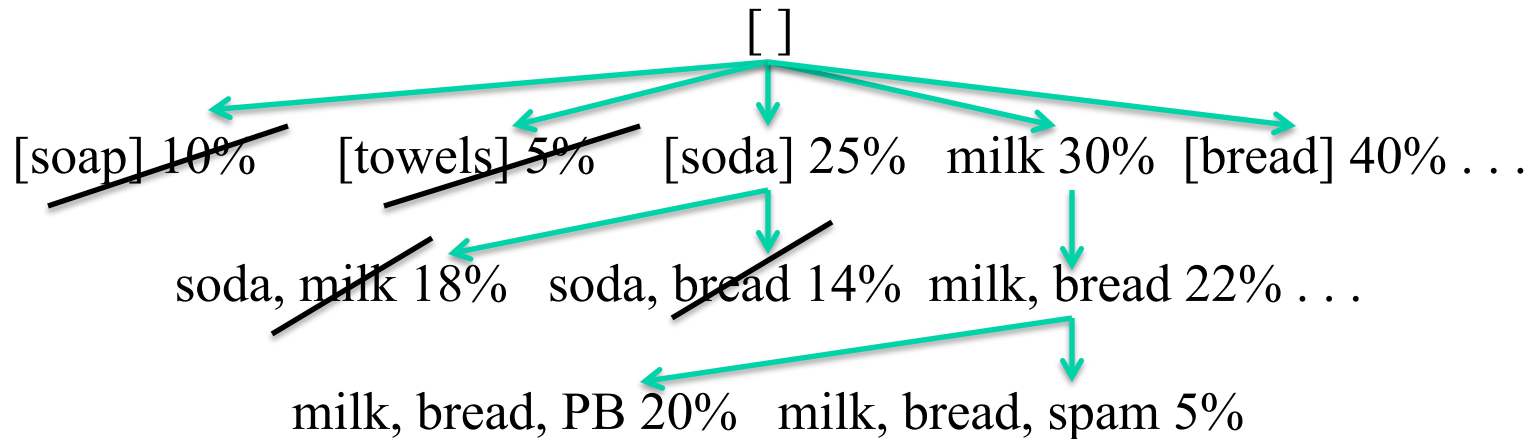
T1 = 20%



A Priori Algorithm

Prune conjunctions of features that fall below threshold support

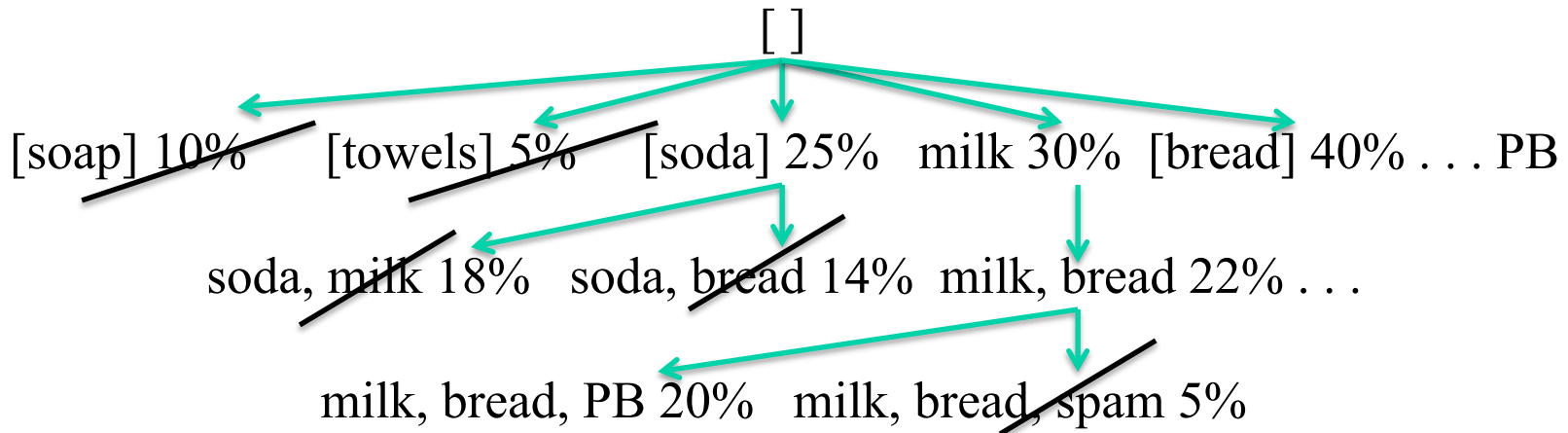
T1 = 20%



A Priori Algorithm

Prune conjunctions of features that fall below threshold support

T1 = 20%



And search for implication rules within each frequent item set

e.g., milk \rightarrow bread $P(\text{bread} | \text{milk})$, and bread \rightarrow milk $P(\text{milk} | \text{bread})$

[bread, milk] ; bread \rightarrow milk (40%, 80%)

[wheat medium bread, low fat milk]; wheat medium bread \rightarrow low fat milk (16%, 85%)

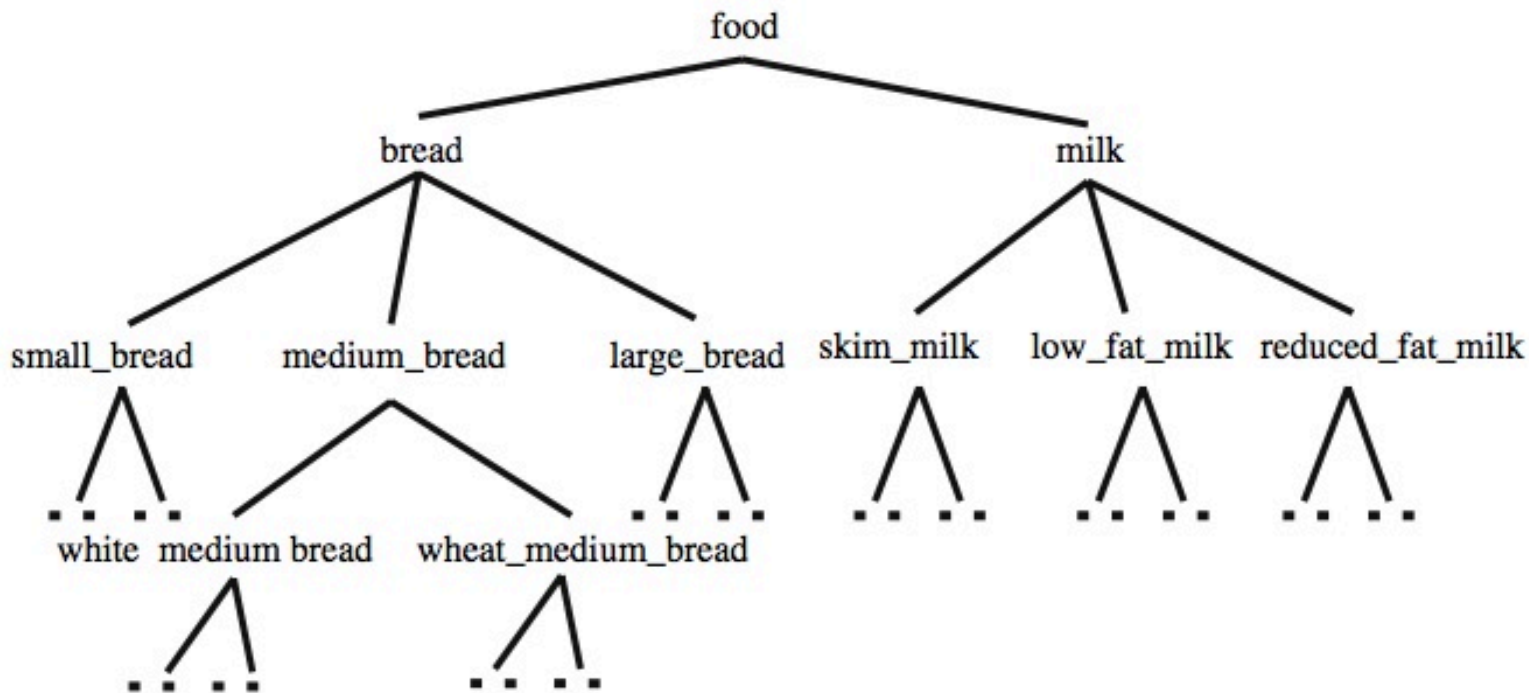


Figure 10.11. Example of a concept hierarchy concerning food.

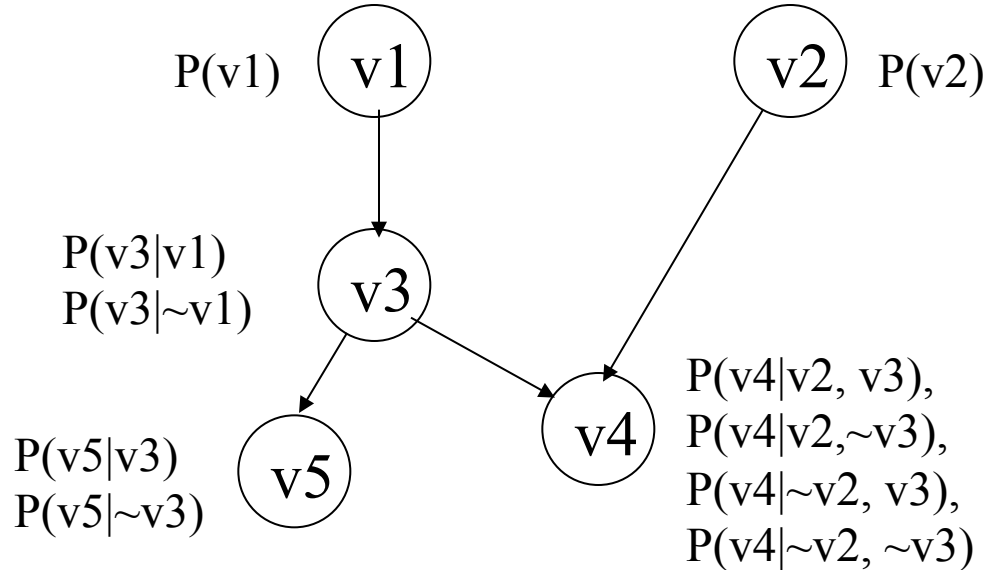
Pairing weak supervision with unsupervised learning to derive a characteristic learner

e.g., Biswas and Kinnebrew

Separate high learners with low learners

Learn association rules for each group

Example (Empirical, Unsupervised): Learning Bayesian Networks



Components of a Bayesian Network: a **topology (graph)** that qualitatively indicates displays the conditional independencies, and **probability tables** at each node

Semantics of graphical component: for each variable, v , v is independent of all of its non-descendants conditioned on its parents

A Bayesian Network is a graphical representation of a joint probability distribution with (conditional) independence relationships made explicit

Recall the chain rule:

Assume V_i a binary valued variable (T or F)

$P(v_1 \text{ and } v_2 \text{ and } \sim v_3 \text{ and } v_4 \text{ and } \sim v_5)$

A factorization ordering

$$= P(v_1)P(v_2|v_1)P(\sim v_3|v_1,v_2)P(v_4|v_1,v_2,\sim v_3)P(\sim v_5|v_1,v_2,\sim v_3,v_4)$$

$P(v_1, v_2)$

$P(v_1, v_2)$

$P(v_1, v_2, \sim v_3)$

$P(v_1, v_2, \sim v_3, v_4)$

$P(v_1, v_2, \sim v_3, v_4, \sim v_5)$


$P(v_1 \text{ and } v_2 \text{ and } \sim v_3 \text{ and } v_4 \text{ and } \sim v_5)$

An alternative ordering

$$= P(v_4)P(v_2|v_4)P(\sim v_3|v_4,v_2)P(v_1|v_4,v_2,\sim v_3)P(\sim v_5|v_4,v_2,\sim v_3,v_1)$$

$P(v_1 \text{ and } v_2 \text{ and } \sim v_3 \text{ and } v_4 \text{ and } \sim v_5)$

A factorization ordering


$$= P(v_1)P(v_2|v_1)P(\sim v_3|v_1, v_2)P(v_4|v_1, v_2, \sim v_3)P(\sim v_5|v_1, v_2, \sim v_3, v_4)$$

Assume the following conditional independencies:

$P(v_1)$ v_2 independent of v_1

$P(v_2|v_1) = P(v_2)$ and $P(v_2|\sim v_1) = P(v_2)$, $P(\sim v_2|v_1) = P(\sim v_2)$, $P(\sim v_2|\sim v_1) = P(\sim v_2)$

$P(\sim v_3|v_1, v_2) = P(\sim v_3|v_1)$ v_3 independent of v_2 conditioned on v_1

and $P(\sim v_3|v_1, \sim v_2) = P(\sim v_3|v_1)$, $P(\sim v_3|\sim v_1, v_2) = P(\sim v_3|\sim v_1)$, $P(\sim v_3|\sim v_1, \sim v_2) = P(\sim v_3|\sim v_1)$,

$P(v_3|v_1, v_2) = P(v_3|v_1)$, $P(v_3|v_1, \sim v_2) = P(v_3|v_1)$, $P(v_3|\sim v_1, v_2) = P(v_3|\sim v_1)$,

$P(v_3|\sim v_1, \sim v_2) = P(v_3|\sim v_1)$

$P(v_4|v_1, v_2, \sim v_3) = P(v_4|v_2, \sim v_3)$ and

$P(\sim v_5|v_1, v_2, \sim v_3, v_4) = P(\sim v_5|\sim v_3)$ and

$$\begin{aligned}
& P(v_1 \text{ and } v_2 \text{ and } \sim v_3 \text{ and } v_4 \text{ and } \sim v_5) \\
&= P(v_1)P(v_2|v_1)P(\sim v_3|v_1,v_2)P(v_4|v_1,v_2,\sim v_3)P(\sim v_5|v_1,v_2,\sim v_3,v_4) \\
&= P(v_1)P(v_2)P(\sim v_3|v_1)P(v_4|v_2,\sim v_3)P(\sim v_5|\sim v_3)
\end{aligned}$$

How many probabilities need be stored?

$P(v_1), P(\sim v_1)$ 2 probabilities (actually only one, since $P(\sim v_1) = 1 - P(v_1)$)

2 probabilities (or 1) instead of 4 (or 2)

$P(v_2|v_1) = P(v_2)$ and $P(v_2|\sim v_1) = P(v_2), P(\sim v_2|v_1) = P(\sim v_2), P(\sim v_2|\sim v_1) = P(\sim v_2)$

$P(\sim v_3|v_1) = 1 - P(v_3|v_1)$

$P(\sim v_3|v_1,v_2) = P(\sim v_3|v_1)$ 4 probabilities (or 2) instead of 8 (or 4)

and $P(\sim v_3|v_1,\sim v_2) = P(\sim v_3|v_1), P(\sim v_3|\sim v_1,v_2) = P(\sim v_3|\sim v_1), P(\sim v_3|\sim v_1,\sim v_2) = P(\sim v_3|\sim v_1),$

$P(v_3|v_1,v_2) = P(v_3|v_1), P(v_3|v_1,\sim v_2) = P(v_3|v_1), P(v_3|\sim v_1,v_2) = P(v_3|\sim v_1),$

$P(v_3|\sim v_1,\sim v_2) = P(v_3|\sim v_1)$

8 probabilities (or 4) instead of 16 (or 8)

$P(v_4|v_1,v_2,\sim v_3) = P(v_4|v_2, \sim v_3)$ and

4 probabilities (or 2) instead of 32 (or 16)

$P(\sim v_5|v_1,v_2,\sim v_3,v_4) = P(\sim v_5|\sim v_3)$ and

For a *particular factorization ordering*, construct a Bayesian network as follows:

v_1 a “root”

$$P(v_1), P(\sim v_1) \quad \textcircled{v_1} \quad P(v_1) = 0.75$$

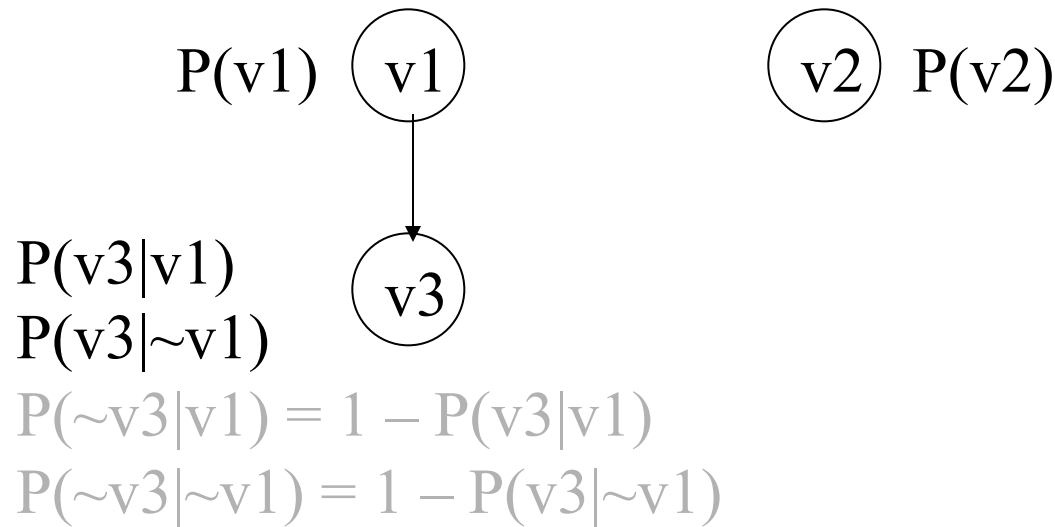
$$P(\sim v_1) = 0.25 = 1 - P(v_1)$$

v_2 is second variable in ordering. If v_2 independent of a subset of its predecessors (possibly the empty set) in ordering conditioned on a disjoint subset of predecessors (including possibly all its predecessors), then the latter subset is its **parents**, else if latter subset is empty then v_2 is a “root”

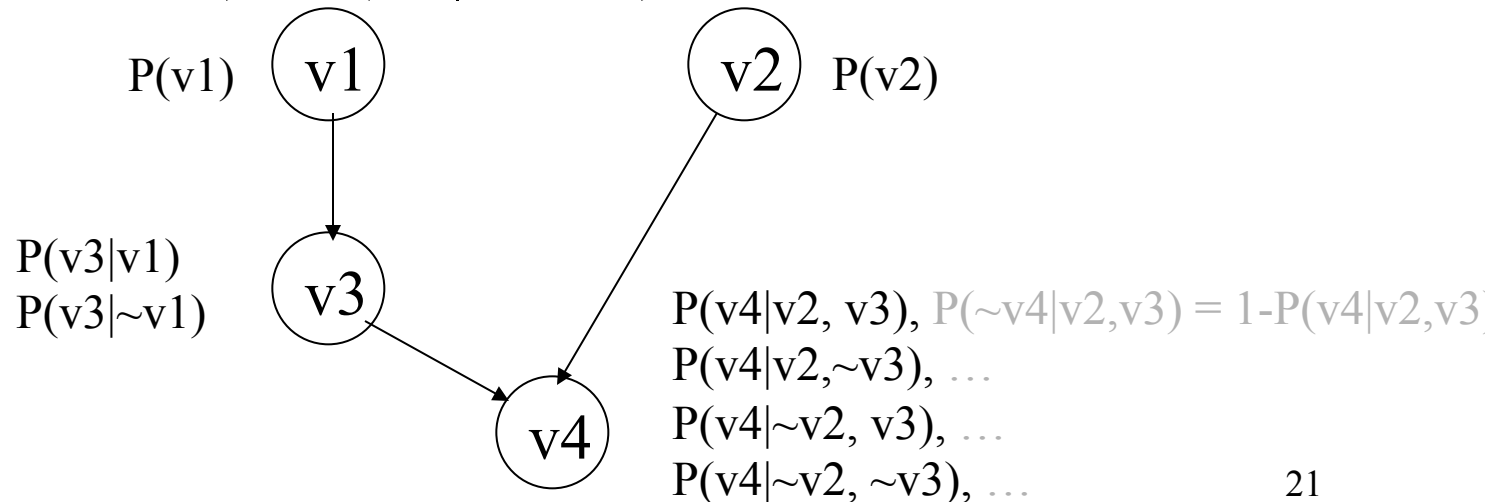
Since $P(v_2|v_1) = P(v_2) \dots$

$$P(v_1) \quad \textcircled{v_1} \quad \textcircled{v_2} \quad P(v_2)$$

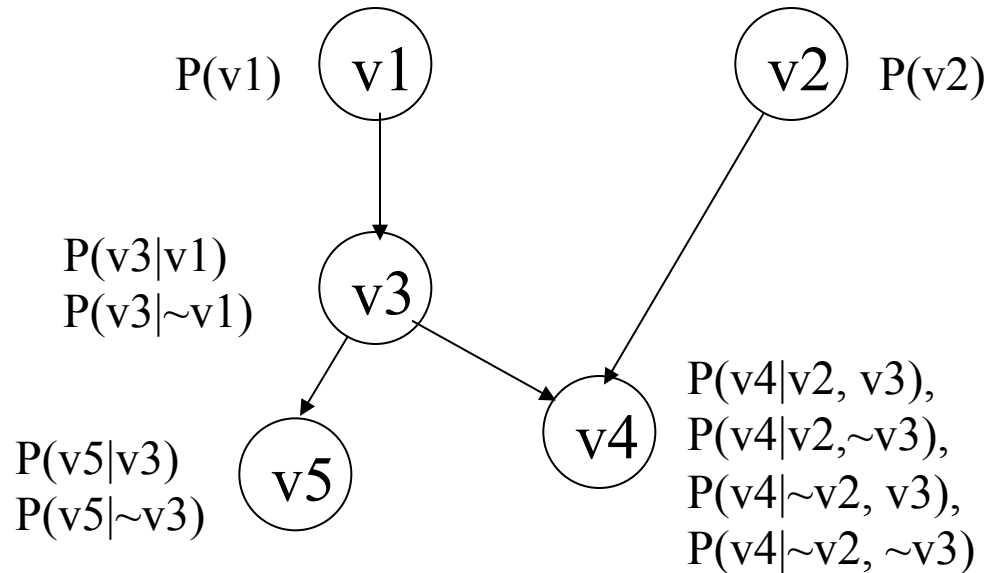
v_3 is third variable in ordering. Since $P(v_3|v_1, v_2) = P(v_3|v_1), \dots$:



Since $P(v_4|v_1, v_2, v_3) = P(v_4|v_2, v_3), \dots$



Since $P(v_5|v_1, v_2, v_3, v_4) = P(v_5|v_3), \dots$:



Components of a Bayesian Network: a **topology (graph)** that qualitatively indicates displays the conditional independencies, and **probability tables** at each node

Semantics of graphical component: for each variable, v , v is independent of all of its non-descendants conditioned on its parents

Where does knowledge of conditional independence come from?

a) **From data.** Consider congressional voting records. Suppose that we have data on House votes (and political party). Suppose variables are ordered Party, Immigration, StarWars,

Party $P(\text{Republican}) = 0.52$ (226/435 Republicans
209/435 Democrats)

To determine relationship between Party and Immigration, we count

	Actual Counts		Predicted Counts (if Immigration and Party independent)	
	Yes	No	Yes	No
Republican	17	209	92	134
Democrat	160	49	85	124

Very different distributions – conclude **dependent**

$$P(\text{Rep}) * P(\text{Yes}) * 435 = 0.52 * (17+160)/435 * 435$$

17/226



$P(\text{Republican}) = 0.52$ (226/435 Republicans
209/435 Democrats)

$P(\text{Yes} | \text{Rep}) = 0.075$

$P(\text{Yes} | \text{Dem}) = 0.765$



Actual Counts

	Immigration	
	Yes	No
Republican	17	209
Democrat	160	49

Consider StarWars

Is StarWars independent of Party and Immigration?

(i.e., is $P(\text{StarWars} | \text{Party}, \text{Immigration})$ approx equal $P(\text{StarWars})$ for all combinations of variable values?)

if yes, then stop and make StarWars a “root”, else continue

Is StarWars independent of Immigration conditioned on Party?

if yes, then stop and make StarWars a child of Party, else continue

Is StarWars independent of Party conditioned on Immigration?

if yes, then stop and make StarWars a child of Immigration, else continue

Make StarWars a child of both Party and Immigration

17/226

Party

$P(\text{Republican}) = 0.52$ (226/435 Republicans
209/435 Democrats)

$P(\text{Yes} | \text{Rep}) = 0.075$

$P(\text{Yes} | \text{Dem}) = 0.765$

Immigration

	Actual Counts		Actual Counts	
	Immigration		StarWars	
	Yes	No	Yes	No
Republican	17	209	219	7
Democrat	160	49	24	185

Consider StarWars

Is StarWars independent of Party and Immigration?

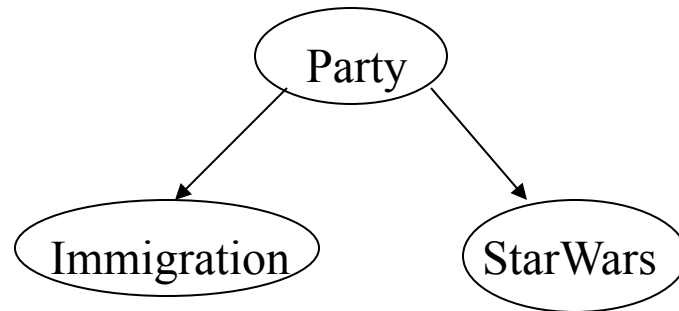
Actual Counts

Predicted Counts

	Actual Counts					Predicted Counts			
	Immigration					Immigration			
	Yes	No	Yes	No		Yes	No	Yes	No
Republican	14	3	205	4	Republican	9.5	7.5	117	92
Democrat	8	152	16	33	Democrat	89	71	27	22
	Yes	No	Yes	No		Yes	No	Yes	No
	StarWars					StarWars			

different – not independent

Further tests might indicate

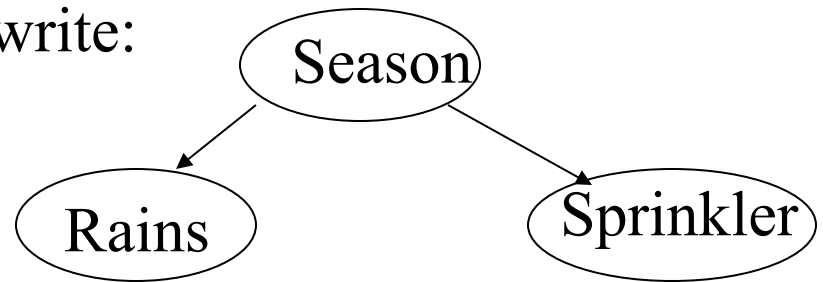


i.e., Immigration and StarWars are independent conditioned on Party

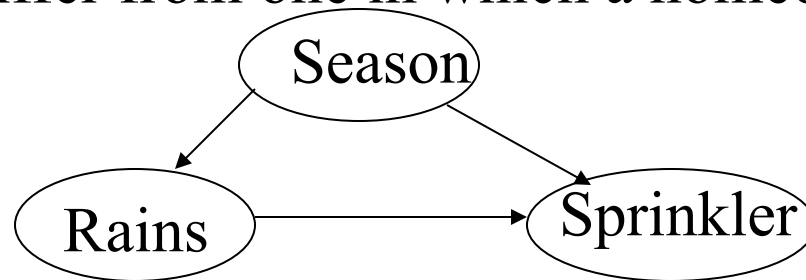
Where does knowledge of conditional independence come from?

b) “First principles”

For example, suppose that the grounds keeper sets sprinkler timers to a fixed schedule that depends on the season (Summer, Winter, Spring, Fall), and suppose that the probability that it rains or not is dependent on season. We might write:



This model might differ from one in which a homeowner manually turns on a sprinkler



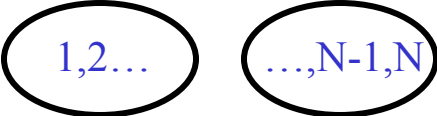
Example (Empirical, Unsupervised): Clustering

Given data (vectors of variable values)

Compute a partition (clusters) of the vectors, such that vectors within a cluster tend to be similar, and vectors across clusters tend to be dissimilar

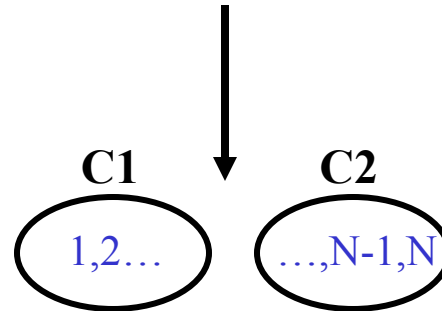
For example,

	V1	V2	V3	V4	VM
1	0.3	0.7	0.1	-0.2	-0.5
2	0.4	0.8	0.01	0.1	-0.4
.....
N-1	-0.3	0.1	1.01	0.8	1.3
N	-0.5	0.03	1.1	0.9	0.9

→ 

Cluster summary representations (e.g., the centroid)

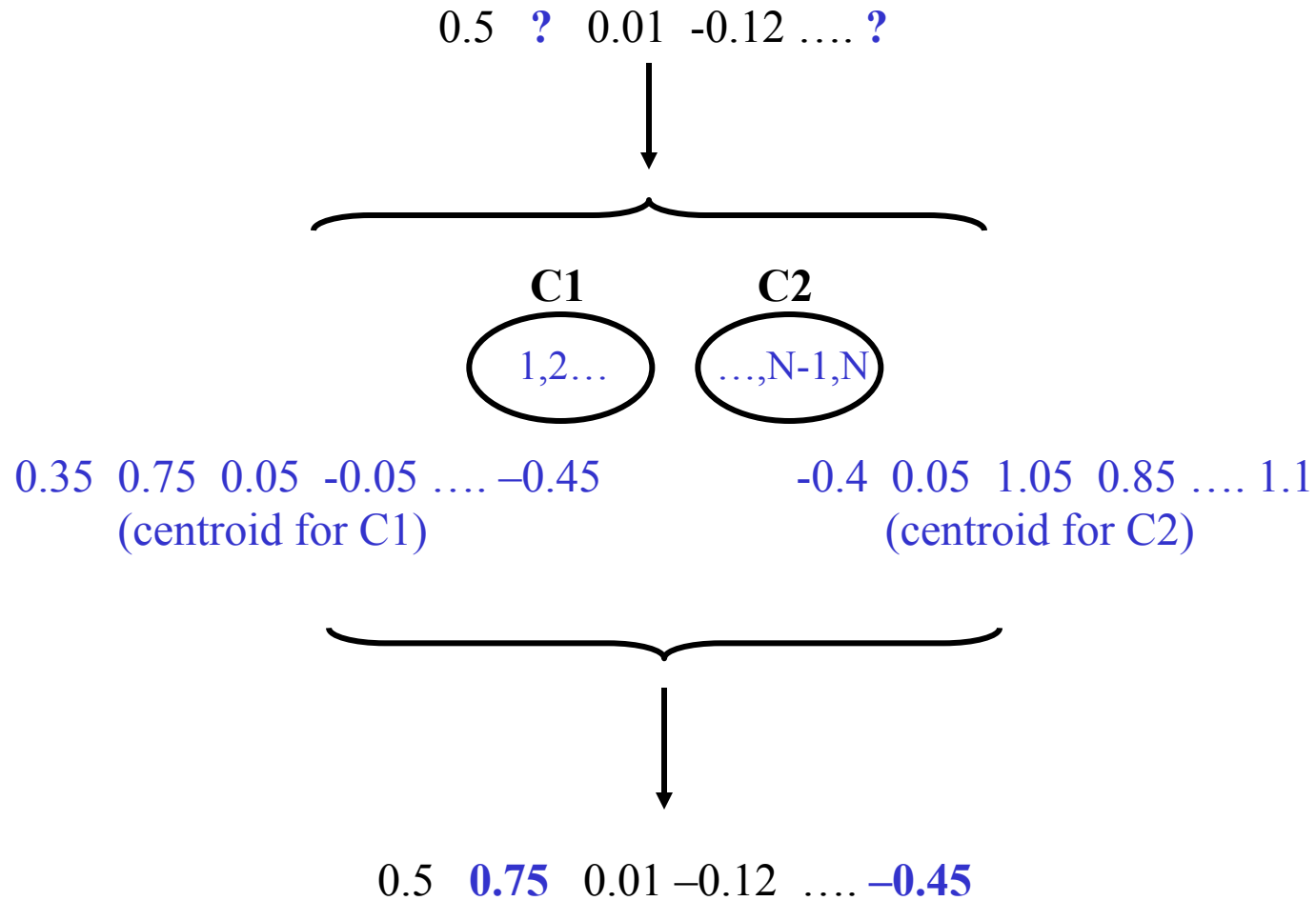
	V1	V2	V3	V4	VM
1	0.3	0.7	0.1	-0.2	-0.5
2	0.4	0.8	0.01	0.1	-0.4
.....						
N-1	-0.3	0.1	1.01	0.8	1.3
N	-0.5	0.03	1.1	0.9	0.9



0.35 0.75 0.05 -0.05 -0.45
 (centroid for C1)

-0.4 0.05 1.05 0.85 1.1
 (centroid for C2)

Using summary representations for inference



K-means

```
Clustering K-Means (Data, K) {  
  ClusterCentroids = K randomly selected vectors from Data  
  for each d in Data  
    assign d to cluster with closest centroid  
  do {  
    compute new cluster centroids  
    for each d in Data  
      assign d to cluster with closest centroid  
    } while NOT termination condition  
}
```

“closest”: Euclidean distance