

Data Mining

Practical Machine Learning Tools and Techniques

Slides for Chapter 6 of *Data Mining* by I. H. Witten, E. Frank and M. A. Hall

Part 2

Implementation:

Real machine learning schemes

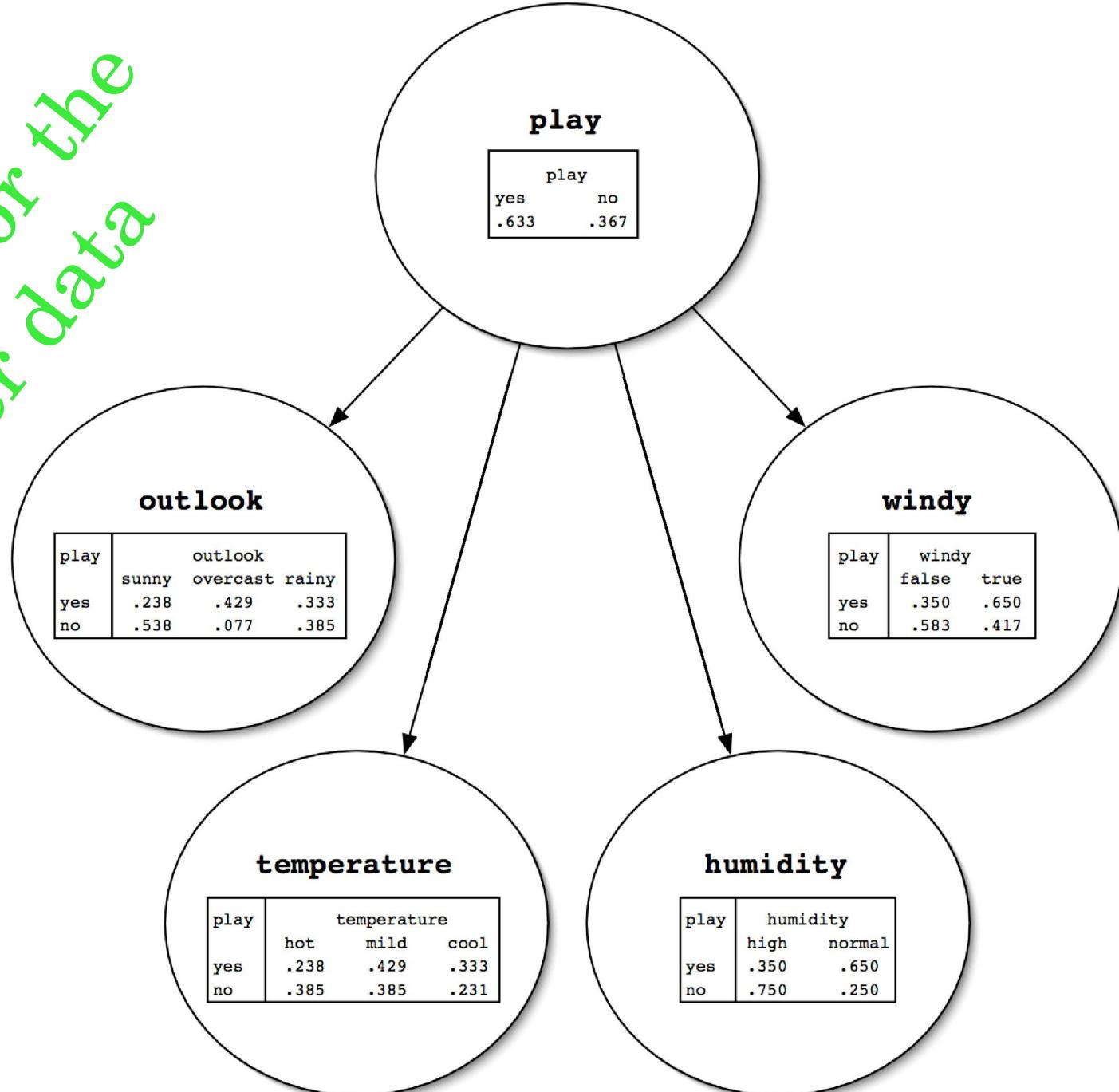
- 6.7 Bayesian networks
 - ◆ Learning and prediction, fast data structures for learning
- 6.8 Clustering: hierarchical, incremental, probabilistic
 - ◆ Hierarchical, incremental, probabilistic, Bayesian
- 6.9 Semisupervised learning
 - ◆ Clustering for classification, co-training

- Naïve Bayes assumes:
attributes conditionally independent
given the class
- Doesn't hold in practice but classification
accuracy often high
- However: sometimes performance much
worse than e.g. decision tree
- Can we eliminate the assumption?

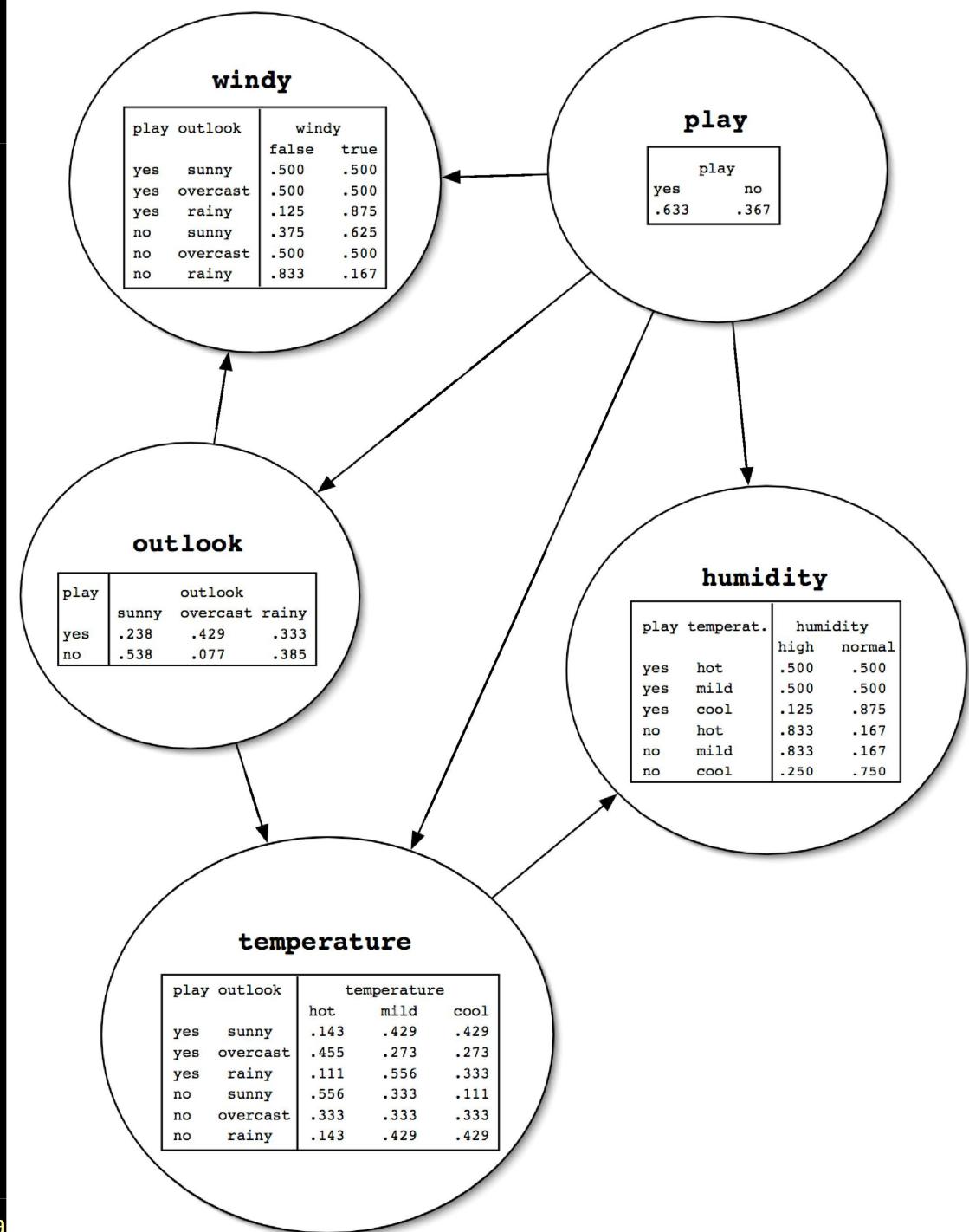
Enter Bayesian networks

- Graphical models that can represent any probability distribution
- Graphical representation: directed acyclic graph, one node for each attribute
- Overall probability distribution factorized into component distributions
- Graph's nodes hold component distributions (conditional distributions)

Network for the weather data



Network for the weather data



Computing the class probabilities

- Two steps: computing a product of probabilities for each class and normalization
 - For each class value
 - Take all attribute values and class value
 - Look up corresponding entries in conditional probability distribution tables
 - Take the product of all probabilities
 - Divide the product for each class by the sum of the products (normalization)

Why can we do this? (Part I)

- Single assumption: values of a node's parents completely determine probability distribution for current node

$$Pr[\text{node}|\text{ancestors}] = Pr[\text{node}|\text{parents}]$$

Means that node/attribute is conditionally independent of other ancestors given parents

Why can we do this? (Part II)

- Chain rule from probability theory:

$$Pr[a_1, a_2, \dots, a_n] = \prod_{i=1}^n Pr[a_i | a_{i-1}, \dots, a_1]$$

Because of our assumption from the previous slide:

$$\begin{aligned} Pr[a_1, a_2, \dots, a_n] &= \prod_{i=1}^n Pr[a_i | a_{i-1}, \dots, a_1] = \\ &\quad \prod_{i=1}^n Pr[a_i | a_i \text{'s parents}] \end{aligned}$$

Learning Bayes nets

- Basic components of algorithms for learning Bayes nets:
 - Method for evaluating the goodness of a given network
 - Measure based on probability of training data given the network (or the logarithm thereof)
 - Method for searching through space of possible networks
 - Amounts to searching through sets of edges because nodes are fixed

Problem: overfitting

- Can't just maximize probability of the training data
 - Because then it's always better to add more edges (fit the training data more closely)
- Need to use cross-validation or some penalty for complexity of the network

AIC measure: $AIC score = -LL + K$

MDL measure: $MDL score = -LL + \frac{K}{2} \log N$

LL : log-likelihood (log of probability of data), K : number of free parameters, N : #instances

Another possibility: Bayesian approach with prior distribution over networks

Searching for a good structure

- Task can be simplified: can optimize each node separately
 - Because probability of an instance is product of individual nodes' probabilities
 - Also works for AIC and MDL criterion because penalties just add up
- Can optimize node by adding or removing edges from other nodes
- Must not introduce cycles!

The K2 algorithm

- Starts with given ordering of nodes (attributes)
- Processes each node in turn
- Greedily tries adding edges from previous nodes to current node
- Moves to next node when current node can't be optimized further
- Result depends on initial order

Some tricks

- Sometimes it helps to start the search with a naïve Bayes network
- It can also help to ensure that every node is in Markov blanket of class node
 - Markov blanket of a node includes all parents, children, and children's parents of that node
 - Given values for Markov blanket, node is conditionally independent of nodes outside blanket
 - I.e. node is irrelevant to classification if not in Markov blanket of class node

Other algorithms

- Extending K2 to consider greedily adding or deleting edges between any pair of nodes
 - Further step: considering inverting the direction of edges
- TAN (Tree Augmented Naïve Bayes):
 - Starts with naïve Bayes
 - Considers adding second parent to each node (apart from class node)
 - Efficient algorithm exists

Likelihood vs. conditional likelihood

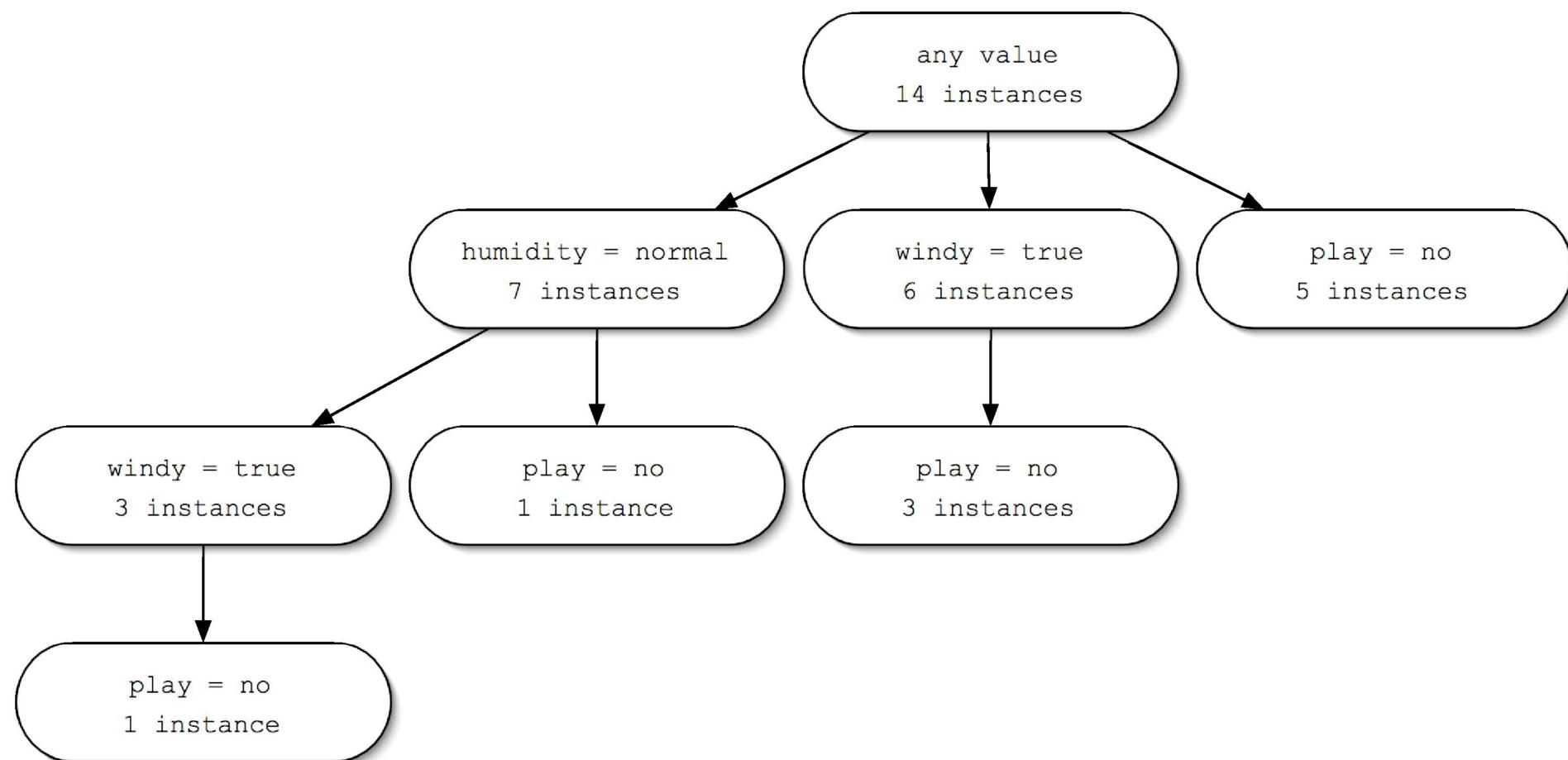
- In classification what we really want is to maximize probability of class given other attributes
 - Not* probability of the instances
- But: no closed-form solution for probabilities in nodes' tables that maximize this
- However: can easily compute conditional probability of data based on given network
- Seems to work well when used for network scoring

Data structures for fast learning

- Learning Bayes nets involves a lot of counting for computing conditional probabilities
- Naïve strategy for storing counts: hash table
 - ◆ Runs into memory problems very quickly
- More sophisticated strategy: *all-dimensions (AD) tree*
 - ◆ Analogous to k D-tree for numeric data
 - ◆ Stores counts in a tree but in a clever way such that redundancy is eliminated
 - ◆ Only makes sense to use it for large datasets

AD tree example

humidity	windy	play	count
high	true	yes	1
high	true	no	2
high	false	yes	2
high	false	no	2
normal	true	yes	2
normal	true	no	1
normal	false	yes	4
normal	false	no	0



Building an AD tree

- Assume each attribute in the data has been assigned an index
- Then, expand node for attribute i with the values of all attributes $j > i$
 - ◆ Two important restrictions:
 - Most populous expansion for each attribute is omitted (breaking ties arbitrarily)
 - Expansions with counts that are zero are also omitted
- The root node is given index zero

Discussion

- We have assumed: discrete data, no missing values, no new nodes
- Different method of using Bayes nets for classification: *Bayesian multinets*
 - I.e. build one network for each class and make prediction using Bayes' rule
- Different class of learning methods for Bayes nets: testing conditional independence assertions
- Can also build Bayes nets for regression tasks

Clustering: how many clusters?

- How to choose k in k -means? Possibilities:
 - Choose k that minimizes cross-validated squared distance to cluster centers
 - Use penalized squared distance on the training data (eg. using an MDL criterion)
 - Apply k -means recursively with $k = 2$ and use stopping criterion (eg. based on MDL)
 - Seeds for subclusters can be chosen by seeding along direction of greatest variance in cluster (one standard deviation away in each direction from cluster center of parent cluster)
 - Implemented in algorithm called X-means (using Bayesian Information Criterion instead of MDL)

Hierarchical clustering

- Recursively splitting clusters produces a hierarchy that can be represented as a *dendrogram*
 - ◆ Could also be represented as a Venn diagram of sets and subsets (without intersections)
 - ◆ Height of each node in the dendrogram can be made proportional to the dissimilarity between its children

Agglomerative clustering

- Bottom-up approach
- Simple algorithm
 - ◆ Requires a distance/similarity measure
 - ◆ Start by considering each instance to be a cluster
 - ◆ Find the two closest clusters and merge them
 - ◆ Continue merging until only one cluster is left
 - ◆ The record of mergings forms a hierarchical clustering structure – a *binary dendrogram*

Distance measures

- *Single-linkage*
 - ◆ Minimum distance between the two clusters
 - ◆ Distance between the clusters closest two members
 - ◆ Can be sensitive to outliers
- *Complete-linkage*
 - ◆ Maximum distance between the two clusters
 - ◆ Two clusters are considered close only if all instances in their union are relatively similar
 - ◆ Also sensitive to outliers
 - ◆ Seeks compact clusters

Distance measures cont.

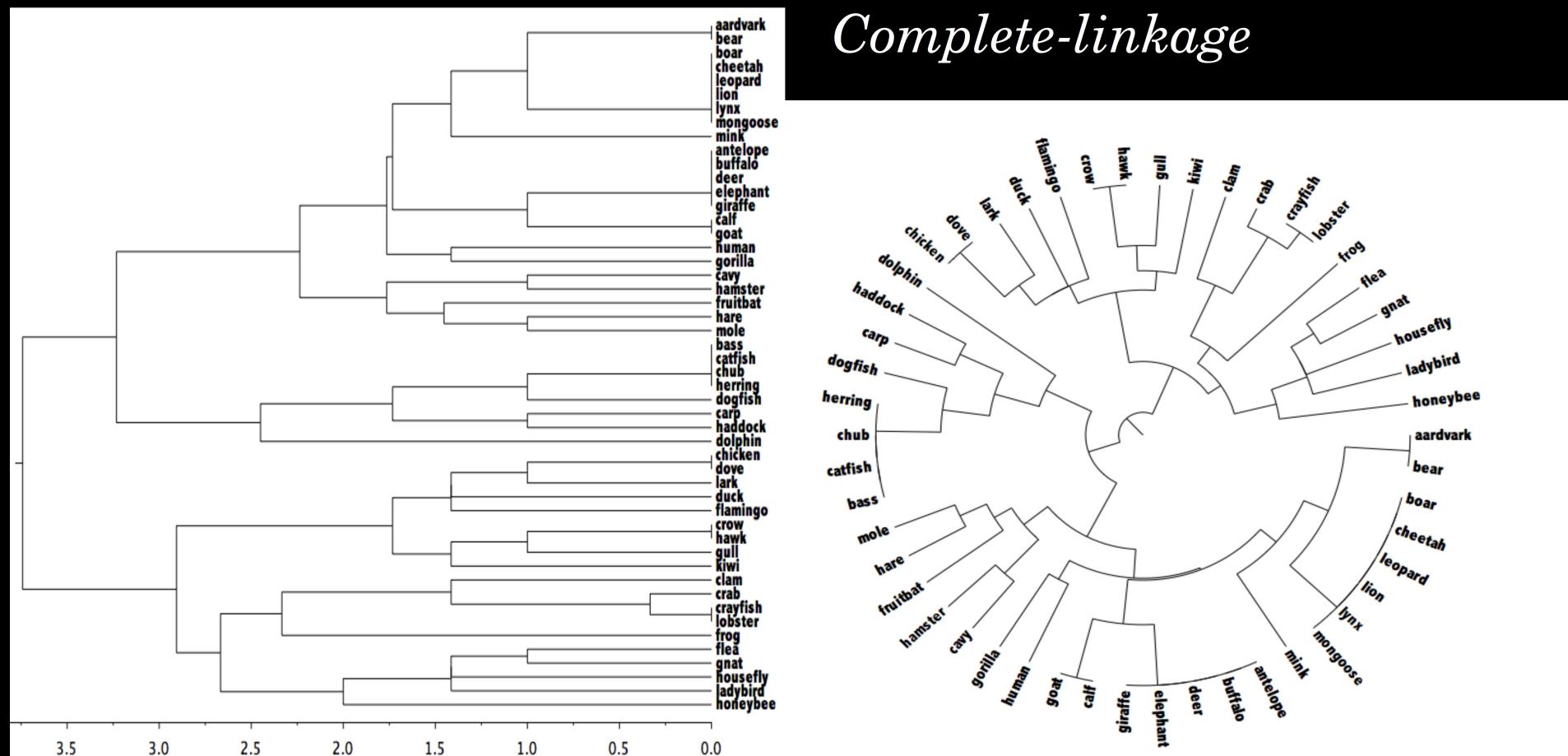
- Compromise between the extremes of minimum and maximum distance
- Represent clusters by their centroid, and use distance between centroids – *centroid linkage*
 - Works well for instances in multidimensional Euclidean space
 - Not so good if all we have is pairwise similarity between instances
- Calculate average distance between each pair of members of the two clusters – *average-linkage*
- Technical deficiency of both: results depend on the numerical scale on which distances are measured

More distance measures

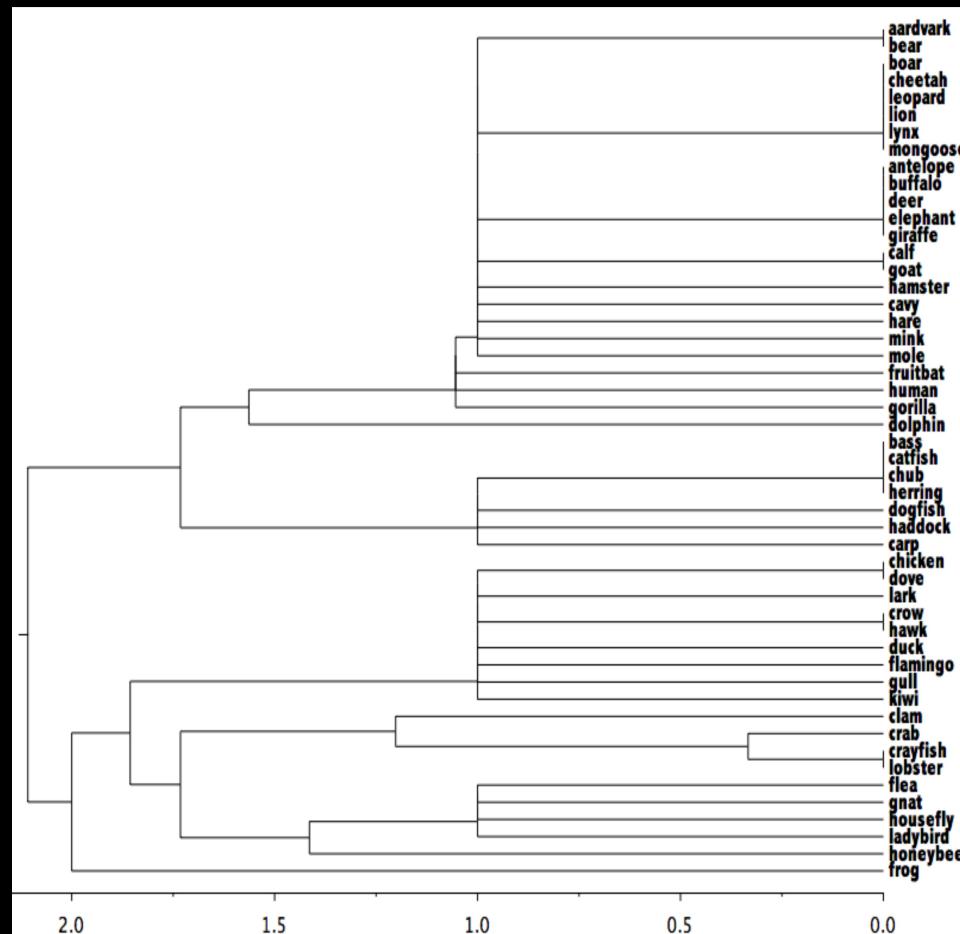
- *Group-average* clustering
 - Uses the average distance between all members of the merged cluster
 - Differs from average-linkage because it includes pairs from the same original cluster
- *Ward's* clustering method
 - Calculates the increase in the sum of squares of the distances of the instances from the centroid before and after fusing two clusters
 - Minimize the increase in this squared distance at each clustering step
- **All** measures will produce the same result if the clusters are compact and well separated

Example hierarchical clustering

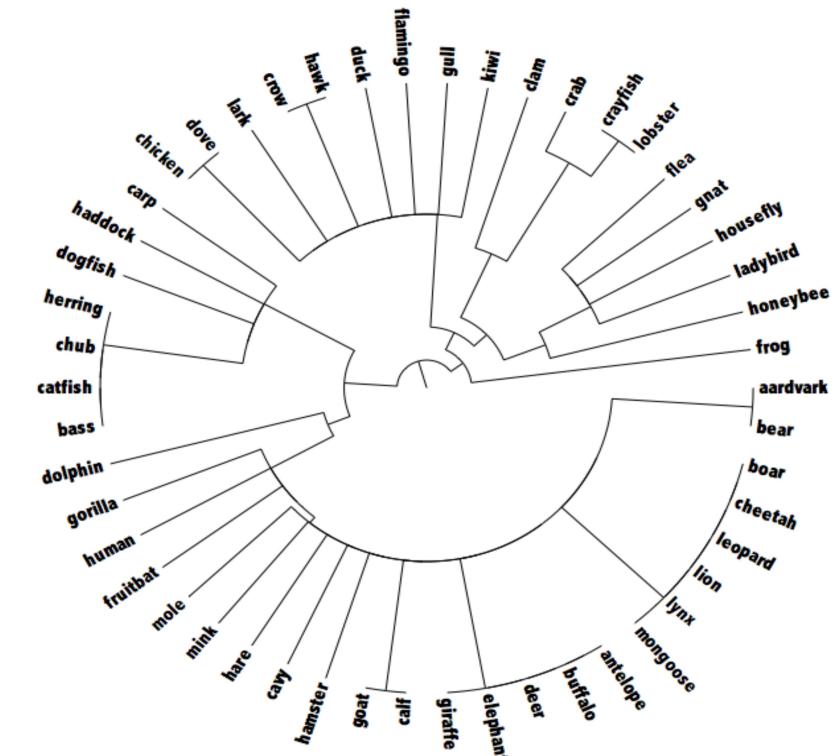
- 50 examples of different creatures from the zoo data



Example hierarchical clustering 2



Single-linkage



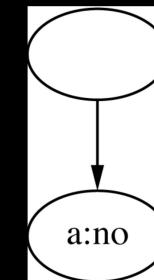
Incremental clustering

- Heuristic approach (COBWEB/CLASSIT)
- Form a hierarchy of clusters incrementally
- Start:
 - tree consists of empty root node
- Then:
 - add instances one by one
 - update tree appropriately at each stage
 - to update, find the right leaf for an instance
 - May involve restructuring the tree
- Base update decisions on *category utility*

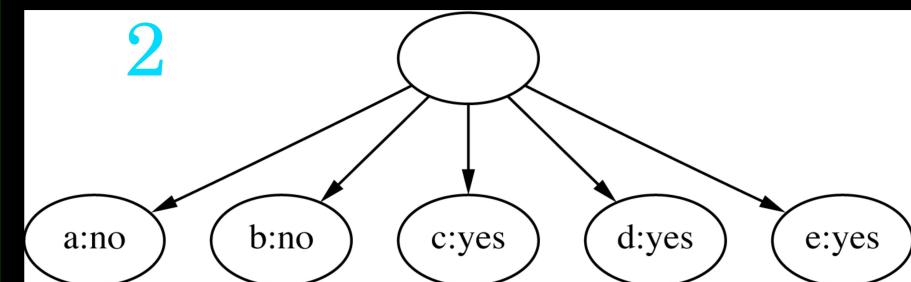
Clustering weather data

ID	Outlook	Temp.	Humidity	Windy
A	Sunny	Hot	High	False
B	Sunny	Hot	High	True
C	Overcast	Hot	High	False
D	Rainy	Mild	High	False
E	Rainy	Cool	Normal	False
F	Rainy	Cool	Normal	True
G	Overcast	Cool	Normal	True
H	Sunny	Mild	High	False
I	Sunny	Cool	Normal	False
J	Rainy	Mild	Normal	False
K	Sunny	Mild	Normal	True
L	Overcast	Mild	High	True
M	Overcast	Hot	Normal	False
N	Rainy	Mild	High	True

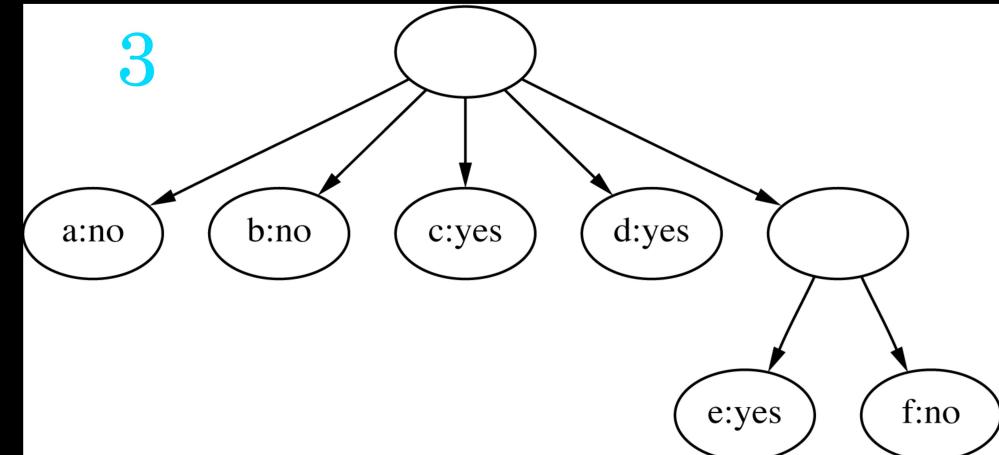
1



2



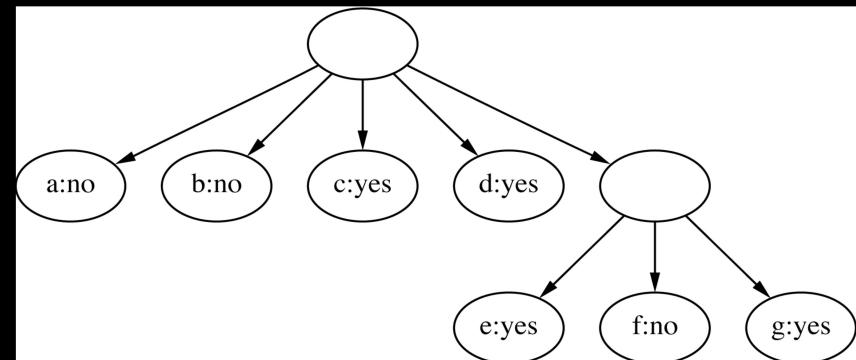
3



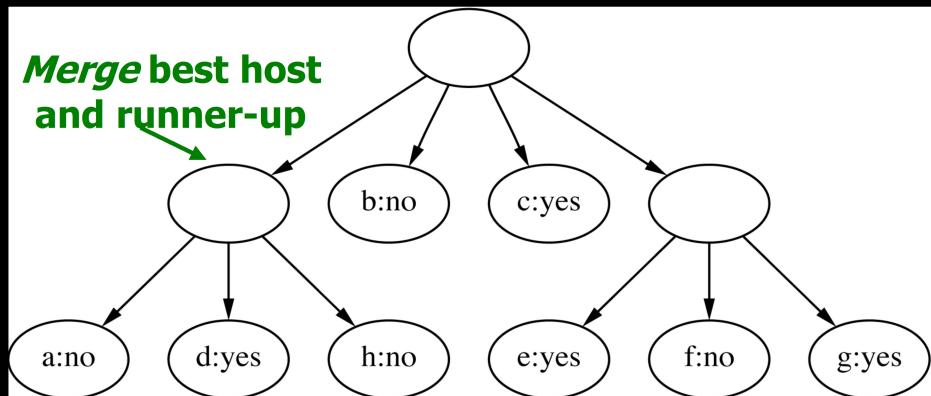
Clustering weather data

ID	Outlook	Temp.	Humidity	Windy
A	Sunny	Hot	High	False
B	Sunny	Hot	High	True
C	Overcast	Hot	High	False
D	Rainy	Mild	High	False
E	Rainy	Cool	Normal	False
F	Rainy	Cool	Normal	True
G	Overcast	Cool	Normal	True
H	Sunny	Mild	High	False
I	Sunny	Cool	Normal	False
J	Rainy	Mild	Normal	False
K	Sunny	Mild	Normal	True
L	Overcast	Mild	High	True
M	Overcast	Hot	Normal	False
N	Rainy	Mild	High	True

4

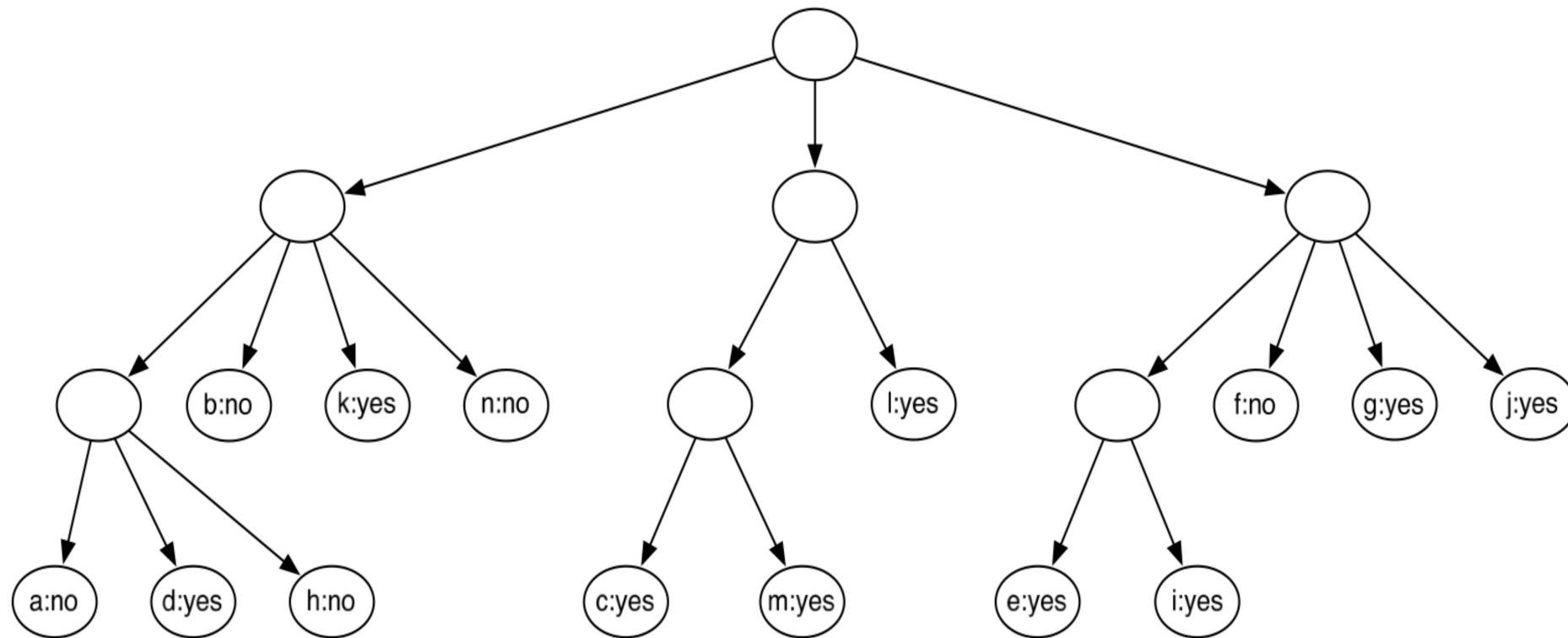


5

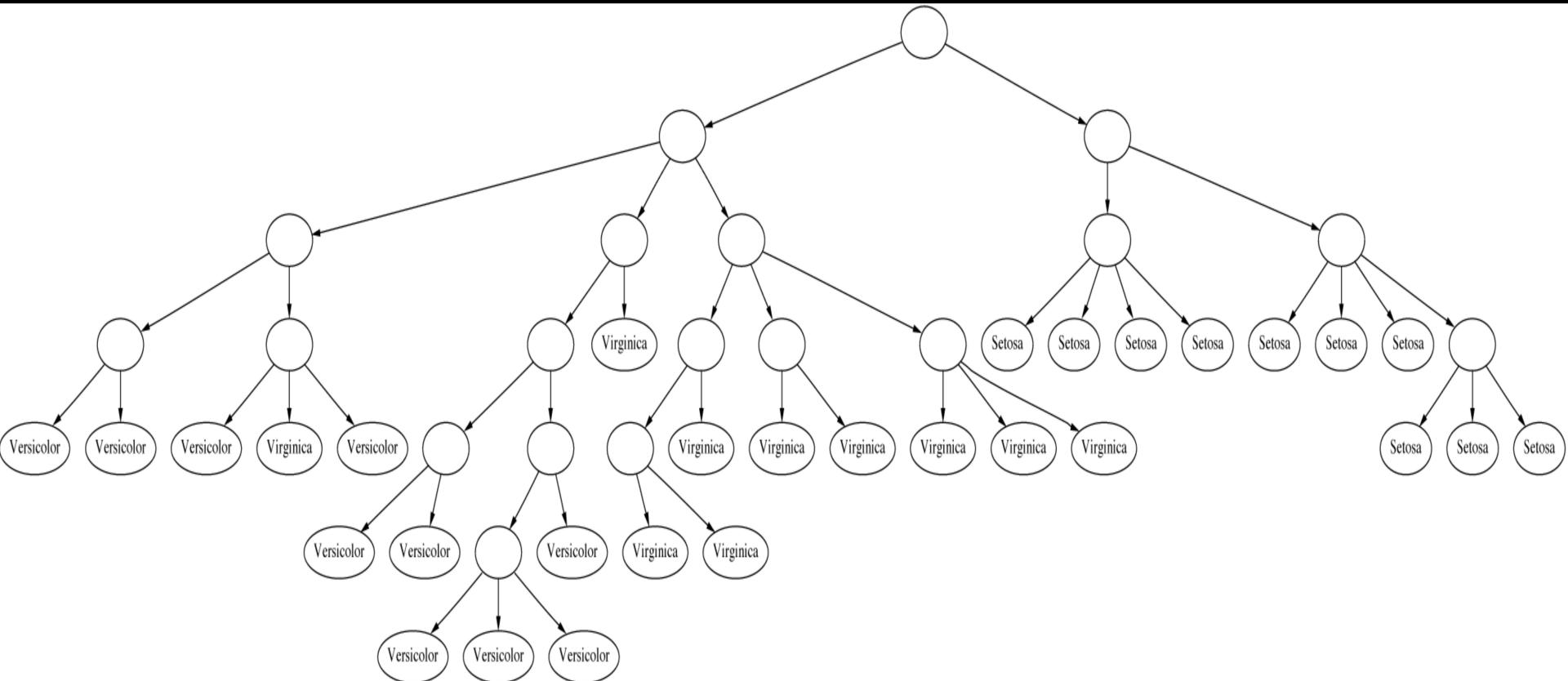


Consider *splitting* the best host if merging doesn't help

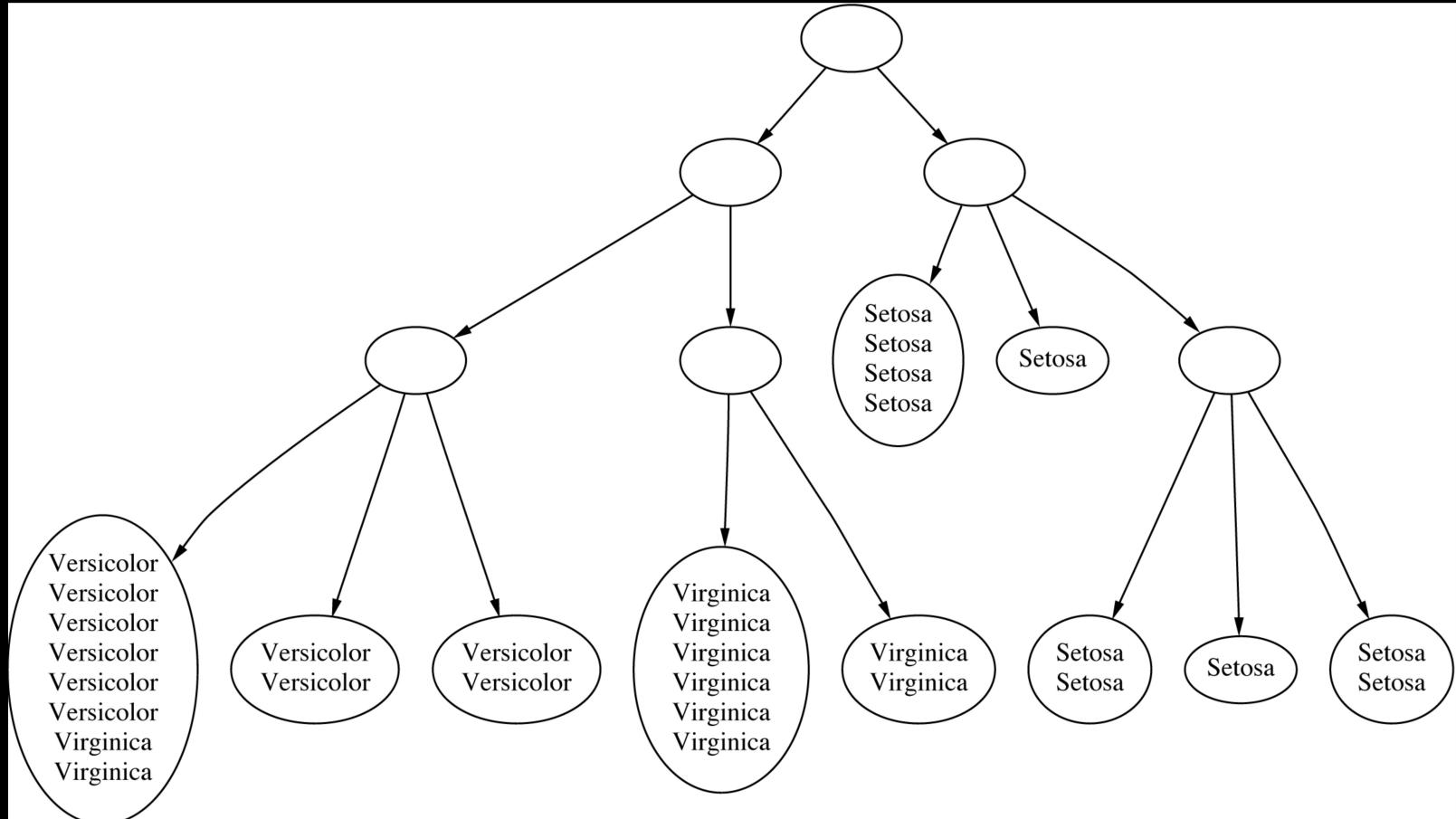
Final hierarchy



Example: the iris data (subset)



Clustering with cutoff



Category utility

- Category utility: quadratic loss function defined on conditional probabilities:

$$CU(C_1, C_2, \dots, C_k) = \frac{\sum_l Pr[C_l] \sum_i \sum_j (Pr[a_i=v_{ij}|C_l]^2 - Pr[a_i=v_{ij}]^2)}{k}$$

- Every instance in different category \Rightarrow numerator becomes

$$n - \sum_i \sum_j Pr[a_i=v_{ij}]^2$$

maximum

↑
number of attributes

Numeric attributes

- Assume normal distribution:

$$f(a) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(a-\mu)^2}{2\sigma^2}\right)$$

- Then:

$$\sum_j Pr[a_i=v_{ij}]^2 \equiv \int f(a_i)^2 da_i = \frac{1}{2\sqrt{\pi}\sigma_i}$$

- Thus $CU(C_1, C_2, \dots, C_k) = \frac{\sum_l Pr[C_l] \sum_i \sum_j (Pr[a_i=v_{ij}|C_l]^2 - Pr[a_i=v_{ij}]^2)}{k}$

becomes $CU(C_1, C_2, \dots, C_k) = \frac{\sum_l Pr[C_l] \frac{1}{2\sqrt{\pi}} \sum_i \left(\frac{1}{\sigma_u} - \frac{1}{\sigma_i}\right)}{k}$

- Prespecified minimum variance
 - ◆ *acuity* parameter

Probability-based clustering

- Problems with heuristic approach:
 - Division by k ?
 - Order of examples?
 - Are restructuring operations sufficient?
 - Is result at least *local* minimum of category utility?
- Probabilistic perspective \Rightarrow seek the *most likely* clusters given the data
- Also: instance belongs to a particular cluster *with a certain probability*

Finite mixtures

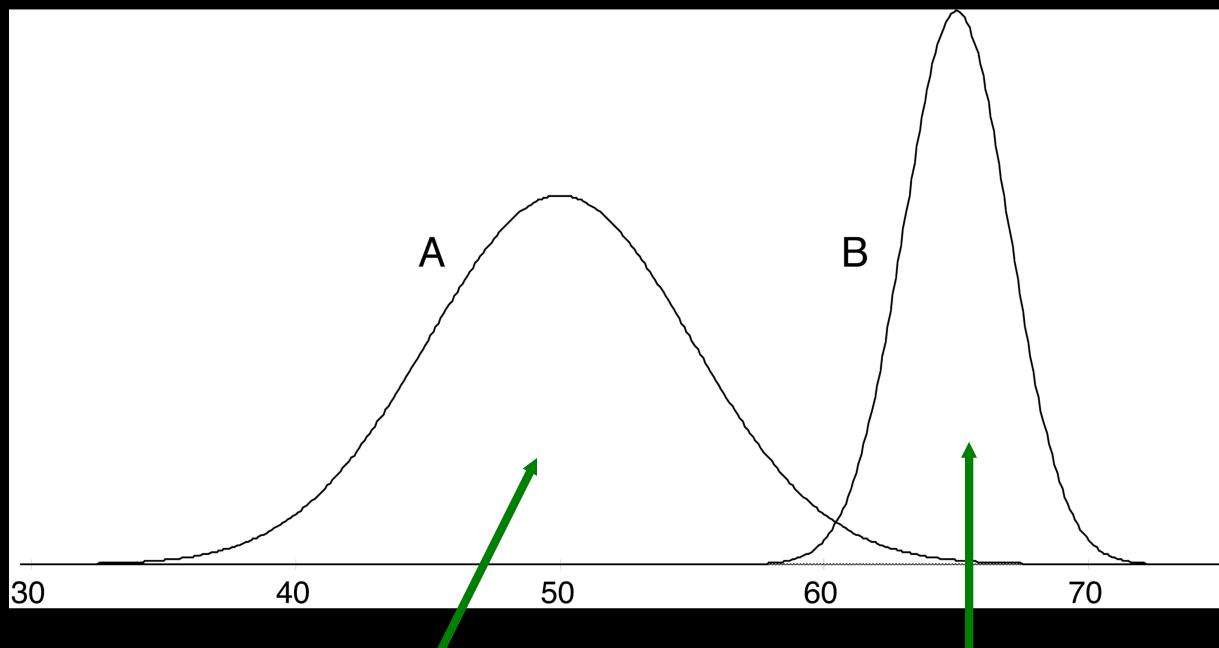
- Model data using a *mixture* of distributions
- One cluster, one distribution
 - governs probabilities of attribute values in that cluster
- *Finite mixtures* : finite number of clusters
- Individual distributions are normal (usually)
- Combine distributions using cluster weights

Two-class mixture model

data

A	51	B	62	B	64	A	48	A	39	A	51
A	43	A	47	A	51	B	64	B	62	A	48
B	62	A	52	A	52	A	51	B	64	B	64
B	64	B	64	B	62	B	63	A	52	A	42
A	45	A	51	A	49	A	43	B	63	A	48
A	42	B	65	A	48	B	65	B	64	A	41
A	46	A	48	B	62	B	66	A	48		
A	45	A	49	A	43	B	65	B	64		
A	45	A	46	A	40	A	46	A	48		

model



$$\mu_A = 50, \sigma_A = 5, \rho_A = 0.6$$

$$\mu_B = 65, \sigma_B = 2, \rho_B = 0.4$$

Using the mixture model

- Probability that instance x belongs to cluster A:

$$Pr[A|x] = \frac{Pr[x|A]Pr[A]}{Pr[x]} = \frac{f(x; \mu_A, \sigma_A)p_A}{Pr[x]}$$

with

$$f(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

- Probability of an instance given the clusters:

$$Pr[x|\text{the_clusters}] = \sum_i Pr[x|\text{cluster}_i] Pr[\text{cluster}_i]$$

Learning the clusters

- Assume:
 - ◆ we know there are k clusters
- Learn the clusters \Rightarrow
 - ◆ determine their parameters
 - ◆ I.e. means and standard deviations
- Performance criterion:
 - ◆ *probability of training data given the clusters*
- EM algorithm
 - ◆ finds a local maximum of the likelihood

EM algorithm

- EM = Expectation-Maximization
 - Generalize k -means to probabilistic setting
- Iterative procedure:
 - E “expectation” step:
Calculate cluster probability for each instance
 - M “maximization” step:
Estimate distribution parameters from cluster probabilities
- Store cluster probabilities as instance weights
- Stop when improvement is negligible

More on EM

- Estimate parameters from weighted instances

$$\mu_A = \frac{w_1 x_1 + w_2 x_2 + \dots + w_n x_n}{w_1 + w_2 + \dots + w_n}$$

$$\sigma_A = \frac{w_1(x_1 - \mu)^2 + w_2(x_2 - \mu)^2 + \dots + w_n(x_n - \mu)^2}{w_1 + w_2 + \dots + w_n}$$

- Stop when log-likelihood saturates
- Log-likelihood:

$$\sum_i \log(p_A Pr[x_i | A] + p_B Pr[x_i | B])$$

Extending the mixture model

- More than two distributions: easy
- Several attributes: easy—assuming independence!
- Correlated attributes: difficult
 - ◆ Joint model: bivariate normal distribution with a (symmetric) covariance matrix
 - ◆ n attributes: need to estimate $n + n(n+1)/2$ parameters

More mixture model extensions

- Nominal attributes: easy if independent
- Correlated nominal attributes: difficult
 - Two correlated attributes $\Rightarrow v_1 v_2$ parameters
- Missing values: easy
- Can use other distributions than normal:
 - “log-normal” if predetermined minimum is given
 - “log-odds” if bounded from above and below
 - Poisson for attributes that are integer counts
- Use cross-validation to estimate k !

Bayesian clustering

- Problem: many parameters \Rightarrow EM overfits
- *Bayesian approach* : give every parameter a prior probability distribution
 - Incorporate prior into overall likelihood figure
 - Penalizes introduction of parameters
- Eg: Laplace estimator for nominal attributes
- Can also have prior on number of clusters!
- Implementation: NASA's AUTOCLASS

Discussion

- Can interpret clusters by using supervised learning
 - ◆ post-processing step
- Decrease dependence between attributes?
 - ◆ pre-processing step
 - ◆ E.g. use *principal component analysis*
- Can be used to fill in missing values
- Key advantage of probabilistic clustering:
 - ◆ Can estimate likelihood of data
 - ◆ Use it to compare different models objectively

Semisupervised learning

- *Semisupervised learning*: attempts to use unlabeled data as well as labeled data
 - ◆ The aim is to improve classification performance
- Why try to do this? Unlabeled data is often plentiful and labeling data can be expensive
 - ◆ Web mining: classifying web pages
 - ◆ Text mining: identifying names in text
 - ◆ Video mining: classifying people in the news
- Leveraging the large pool of unlabeled examples would be very attractive

Clustering for classification

- Idea: use naïve Bayes on labeled examples and then apply EM
 - First, build naïve Bayes model on labeled data
 - Second, label unlabeled data based on class probabilities (“expectation” step)
 - Third, train new naïve Bayes model based on all the data (“maximization” step)
 - Fourth, repeat 2nd and 3rd step until convergence
- Essentially the same as EM for clustering with fixed cluster membership probabilities for labeled data and #clusters = #classes

Comments

- Has been applied successfully to document classification
 - Certain phrases are indicative of classes
 - Some of these phrases occur only in the unlabeled data, some in both sets
 - EM can generalize the model by taking advantage of co-occurrence of these phrases
- Refinement 1: reduce weight of unlabeled data
- Refinement 2: allow multiple clusters per class

Co-training

- Method for learning from *multiple views* (multiple sets of attributes), eg:
 - First set of attributes describes content of web page
 - Second set of attributes describes links that link to the web page
- Step 1: build model from each view
- Step 2: use models to assign labels to unlabeled data
- Step 3: select those unlabeled examples that were most confidently predicted (ideally, preserving ratio of classes)
- Step 4: add those examples to the training set
- Step 5: go to Step 1 until data exhausted
- Assumption: views are independent

EM and co-training

- Like EM for semisupervised learning, but view is switched in each iteration of EM
 - Uses all the unlabeled data (probabilistically labeled) for training
- Has also been used successfully with support vector machines
 - Using logistic models fit to output of SVMs
- Co-training also seems to work when views are chosen randomly!
 - Why? Possibly because co-trained classifier is more robust