# Data Mining:
## Concepts and Techniques
### (3rd ed.)

## Classification & Clustering: Other Topics

Jiawei Han, Micheline Kamber, and Jian Pei

University of Illinois at Urbana-Champaign &

Simon Fraser University

# Classification & Clustering: Other Topics

- Additional Topics Regarding Classification
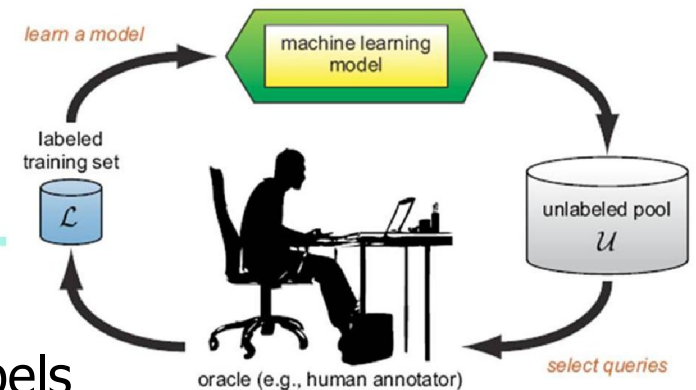- Clustering Categorical Data
- Summary

# Multiclass Classification

- Classification involving more than two classes (i.e., > 2 Classes)
- Method 1. **One-vs.-all** (OVA): Learn a classifier one at a time
    - Given m classes, train m classifiers: one for each class
    - Classifier j: treat tuples in class j as *positive* & all others as *negative*
    - To classify a tuple **X**, the set of classifiers vote as an ensemble
- Method 2. **All-vs.-all** (AVA): Learn a classifier for each pair of classes
    - Given m classes, construct m(m-1)/2 binary classifiers
    - A classifier is trained using tuples of the two classes
    - To classify a tuple **X**, each classifier votes.  X is assigned to the class with maximal vote
- Comparison
    - All-vs.-all tends to be superior to one-vs.-all
    - Problem: Binary classifier is sensitive to errors, and errors affect vote count

# Semi-Supervised Classification

- Semi-supervised: Uses labeled and unlabeled data to build a classifier
- Self-training:
  - Build a classifier using the labeled data
  - Use it to label the unlabeled data, and those with the most confident label prediction are added to the set of labeled data
  - Repeat the above process
  - Adv: easy to understand; disadv: may reinforce errors
- Co-training: Use two or more classifiers to teach each other
  - Each learner uses a mutually independent set of features of each tuple to train a good classifier, say $f_1$
  - Then $f_1$ and $f_2$ are used to predict the class label for unlabeled data X
  - Teach each other: The tuple having the most confident prediction from $f_1$ is added to the set of labeled data for $f_2$, & vice versa
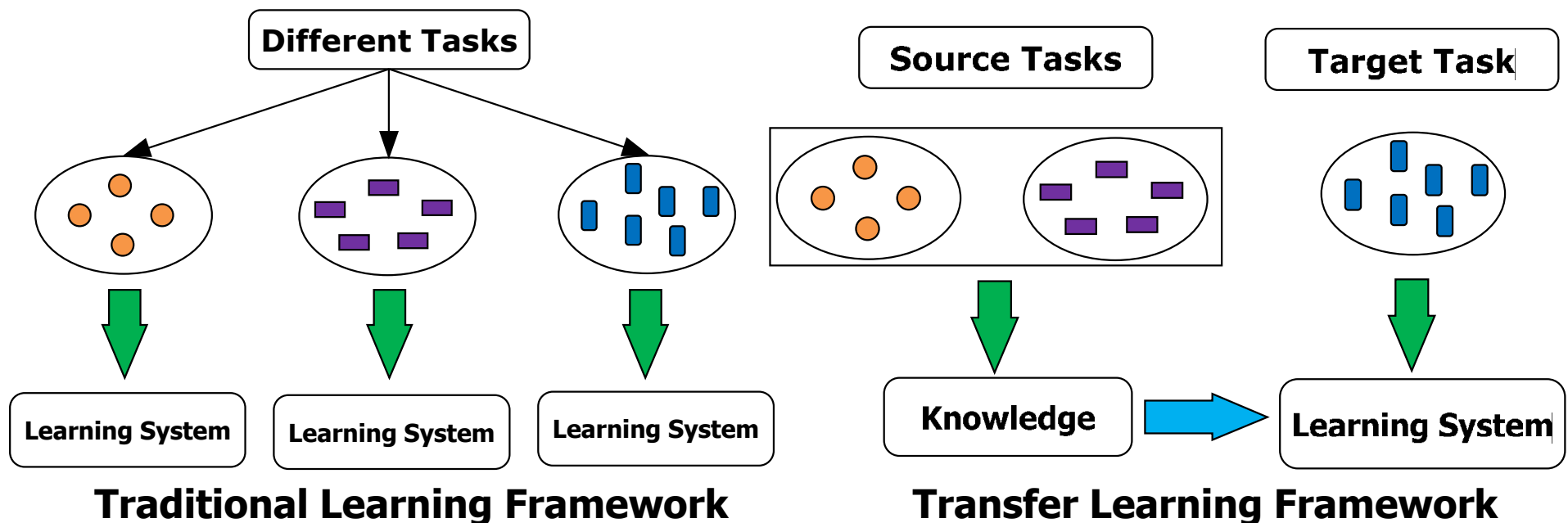- Other methods, e.g., joint probability distribution of features and labels

# Active Learning



learn a model
machine learning model
labeled training set
$\mathcal{L}$
unlabeled pool
$\mathcal{U}$
oracle (e.g., human annotator)
select queries

- Class labels are expensive to obtain
- Active learner: query human (oracle) for labels
- Pool-based approach: Uses a pool of unlabeled data
  - L: a small subset of D is labeled, U: a pool of unlabeled data in D
  - Use a query function to carefully select one or more tuples from U and request labels from an oracle (a human annotator)
  - The newly labeled samples are added to L, and learn a model
  - Goal: Achieve high accuracy using as few labeled data as possible
- Evaluated using *learning curves*: Accuracy as a function of the number of instances queried (# of tuples to be queried should be small)
- Research issue: How to choose the data tuples to be queried?
  - Uncertainty sampling: choose the least certain ones
  - Reduce *version space*, the subset of hypotheses consistent w. the training data
  - Reduce expected entropy over U: Find the greatest reduction in the total number of incorrect predictions

# Transfer Learning: Conceptual Framework

- Transfer learning: Extract knowledge from one or more source tasks and apply the knowledge to a target task

- Traditional learning: Build a new classifier for each new task

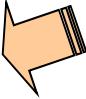- Transfer learning: Build new classifier by applying existing knowledge learned from source tasks



**Traditional Learning Framework**

**Transfer Learning Framework**

# Transfer Learning: Methods and Applications

- Applications: Especially useful when data is outdated or distribution changes, e.g., Web document classification, e-mail spam filtering
- *Instance-based transfer learning*: Reweight some of the data from source tasks and use it to learn the target task
- TrAdaBoost (Transfer AdaBoost)
  - Assume source and target data each described by the same set of attributes (features) & class labels, but rather diff. distributions
  - Require only labeling a small amount of target data
  - Use source data in training: When a source tuple is misclassified, reduce the weight of such tupels so that they will have less effect on the subsequent classifier
- Research issues
  - Negative transfer: When it performs worse than no transfer at all
  - Heterogeneous transfer learning: Transfer knowledge from different feature space or multiple source domains
  - Large-scale transfer learning

# Classification & Clustering: Other Topics

- Additional Topics Regarding Classification

- Clustering Categorical Data

- Summary

# ROCK: Clustering Categorical Data

- ROCK: RObust Clustering using linKs
  - S. Guha, R. Rastogi & K. Shim, ICDE'99
- Major ideas
  - Use links to measure similarity/proximity
  - Not distance-based
- Algorithm: sampling-based clustering
  - Draw random sample
  - Cluster with links
  - Label data in disk
- Experiments
  - Congressional voting, mushroom data

# Similarity Measure in ROCK

- Traditional measures for categorical data may not work well, e.g., Jaccard coefficient
- Example: Two groups (clusters) of transactions
  - $C_1$. <a, b, c, d, e>: {a, b, c}, {a, b, d}, {a, b, e}, {a, c, d}, {a, c, e}, {a, d, e}, {b, c, d}, {b, c, e}, {b, d, e}, {c, d, e}
  - $C_2$. <a, b, f, g>: {a, b, f}, {a, b, g}, {a, f, g}, {b, f, g}
- Jaccard co-efficient may lead to wrong clustering result
  - $C_1$: 0.2 ({a, b, c}, {b, d, e}} to 0.5 ({a, b, c}, {a, b, d})
  - $C_1$ & $C_2$: could be as high as 0.5 ({a, b, c}, {a, b, f})
- Jaccard co-efficient-based similarity function:

$$Sim(T_1, T_2) = \frac{|T_1 \cap T_2|}{|T_1 \cup T_2|}$$

  - Ex. Let $T_1$ = {a, b, c}, $T_2$ = {c, d, e}

$$Sim\ (T_1, T_2) = \frac{|\{c\}|}{|\{a, b, c, d, e\}|} = \frac{1}{5} = 0.2$$

# Link Measure in ROCK

- Clusters
  - $C_1$:<a, b, c, d, e>: {a, b, c}, {a, b, d}, {a, b, e}, {a, c, d}, {a, c, e}, {a, d, e}, {b, c, d}, {b, c, e}, {b, d, e}, {c, d, e}
  - $C_2$: <a, b, f, g>: {a, b, f}, {a, b, g}, {a, f, g}, {b, f, g}
- Neighbors
  - Two transactions are neighbors if sim($T_1$,$T_2$) > threshold (here is to 0.5)
  - Let $T_1$ = {a, b, c}, $T_2$ = {c, d, e}, $T_3$ = {a, b, f}
    - $T_1$ connected to: {a,b,d}, {a,b,e}, {a,c,d}, {a,c,e}, {b,c,d}, {b,c,e}, {a,b,f}, {a,b,g}
    - $T_2$ connected to: {a,c,d}, {a,c,e}, {a,d,e}, {b,c,e}, {b,d,e}, {b,c,d}
    - $T_3$ connected to: {a,b,c}, {a,b,d}, {a,b,e}, {a,b,g}, {a,f,g}, {b,f,g}
- Link Similarity
  - Link similarity between two transactions is the # of common neighbors
  - link($T_1$, $T_2$) = 4, since they have 4 common neighbors
    - {a, c, d}, {a, c, e}, {b, c, d}, {b, c, e}
  - link($T_1$, $T_3$) = 3, since they have 3 common neighbors
    - {a, b, d}, {a, b, e}, {a, b, g}

# Mushroom Data Set

- http://archive.ics.uci.edu/ml/datasets/Mushroom

- Number of Instances: 8124

- Number of Attributes: 22 (all nominally valued) including cap shape, cap color, odor, etc.

- Missing Attribute Values: 2480 of them (denoted by "?")

- Class Distribution:
  - edible: 4208 (51.8%)
  - poisonous: 3916 (48.2%)
  - total: 8124 instances

# Clustering result for mushroom data

| Traditional Hierarchical Algorithm | | | | | |
|---|---|---|---|---|---|
| Cluster No | No of Edible | No of Poisonous | Cluster No | No of Edible | No of Poisonous |
| 1 | 666 | 478 | 11 | 120 | 144 |
| 2 | 283 | 318 | 12 | 128 | 140 |
| 3 | 201 | 188 | 13 | 144 | 163 |
| 4 | 164 | 227 | 14 | 198 | 163 |
| 5 | 194 | 125 | 15 | 131 | 211 |
| 6 | 207 | 150 | 16 | 201 | 156 |
| 7 | 233 | 238 | 17 | 151 | 140 |
| 8 | 181 | 139 | 18 | 190 | 122 |
| 9 | 135 | 78 | 19 | 175 | 150 |
| 10 | 172 | 217 | 20 | 168 | 206 |

| ROCK | | | | | |
|---|---|---|---|---|---|
| Cluster No | No of Edible | No of Poisonous | Cluster No | No of Edible | No of Poisonous |
| 1 | 96 | 0 | 12 | 48 | 0 |
| 2 | 0 | 256 | 13 | 0 | 288 |
| 3 | 704 | 0 | 14 | 192 | 0 |
| 4 | 96 | 0 | 15 | 32 | 72 |
| 5 | 768 | 0 | 16 | 0 | 1728 |
| 6 | 0 | 192 | 17 | 288 | 0 |
| 7 | 1728 | 0 | 18 | 0 | 8 |
| 8 | 0 | 32 | 19 | 192 | 0 |
| 9 | 0 | 1296 | 20 | 16 | 0 |
| 10 | 0 | 8 | 21 | 0 | 36 |
| 11 | 48 | 0 | | | |

ROCK threshold is 0.8

# Classification & Clustering: Other Topics

- Additional Topics Regarding Classification

- Clustering Categorical Data

- Summary

# Summary

- Other Classification methods

  - Multiclass classification

  - Semi-supervised classification

  - Active learning

  - Transfer learning

- Clustering Categorical Data