# Lecture 15

# Expectation Maximization (EM), Unsupervised Learning and Clustering

**Tuesday 22 October 2002**

**William H. Hsu**

**Department of Computing and Information Sciences, KSU**

http://www.kddresearch.org

http://www.cis.ksu.edu/~bhsu

Readings:

Section 6.12, Mitchell

Section 3.2.4, Shavlik and Dietterich (Rumelhart and Zipser)

Section 3.2.5, Shavlik and Dietterich (Kohonen)

# Lecture Outline

- **Readings: 6.12, Mitchell; Rumelhart and Zipser**

- **Suggested Reading: Kohonen**

- **This Week's Review: "The Future of Time Series", Gershenfeld and Weigend**

- **Unsupervised Learning and Clustering**

    - **Definitions and framework**

    - **Constructive induction**

        - **Feature construction**

        - **Cluster definition**

    - **EM, *AutoClass*, Principal Components Analysis, Self-Organizing Maps**

- **Expectation-Maximization (EM) Algorithm**

    - **More on EM and Bayesian Learning**

    - **EM and unsupervised learning**

- **Next Lecture: Time Series Learning**

    - **Intro to time series learning, characterization; stochastic processes**

    - **Read Chapter 16, Russell and Norvig (decisions and utility)**

**CIS 732: Machine Learning and Pattern Recognition**

# Unsupervised Learning: Objectives

- **Unsupervised Learning**

  $x \longrightarrow$ [Supervised Learning] $\xrightarrow{\hat{f}(x)}$ $f(x)$ $x \longrightarrow$ [Unsupervised Learning] $\longrightarrow y$

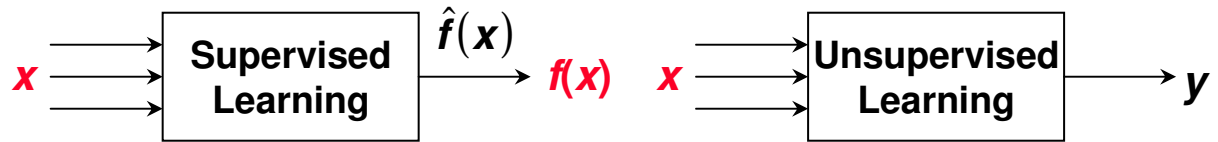  - **Given: data set _D_**

    - **Vectors of attribute values $(x_1, x_2, \ldots, x_n)$**

    - **No distinction between input attributes and output attributes (class label)**

  - **Return: (synthetic) <u>descriptor</u> _y_ of each _x_**

    - **<u>Clustering</u>: _grouping points_ (_x_) into inherent regions of mutual similarity**

    - **<u>Vector quantization</u>: _discretizing continuous space_ with best labels**

    - **<u>Dimensionality reduction</u>: _projecting many attributes_ down to a few**

    - **<u>Feature extraction</u>: _constructing (few) new attributes_ from (many) old ones**
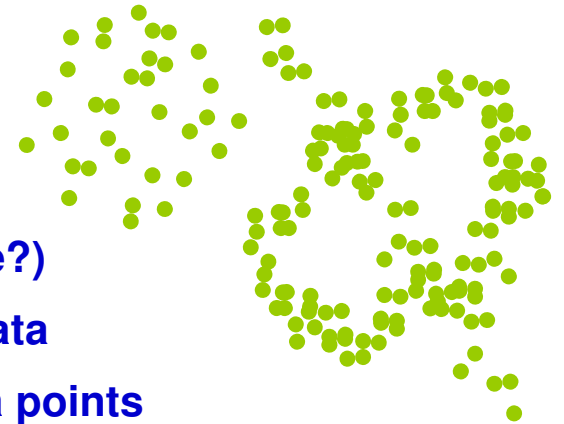
- **Intuitive Idea**

  - **Want to map <u>independent variables</u> (_x_) to <u>dependent variables</u> (_y = f(x)_)**

  - **_Don't always know what "dependent variables" (y) are_**

  - **Need to discover _y_ based on numerical criterion (e.g., <u>distance metric</u>)**

# Clustering

- **A Mode of Unsupervised Learning**

    - **Given: a collection of data points**

    - **Goal: *discover structure* in the data**

        - ***Organize data into sensible groups* (how many here?)**

        - **Criteria: <u>convenient</u> and <u>valid</u> organization of the data**

        - **NB: not necessarily rules for classifying future data points**

    - **<u>Cluster analysis</u>: study of algorithms, methods for discovering this structure**

        - ***Representing* structure: organizing data into clusters (<u>cluster formation</u>)**

        - ***Describing* structure: cluster boundaries, centers (<u>cluster segmentation</u>)**

        - ***Defining* structure: assigning meaningful names to clusters (<u>cluster labeling</u>)**

- **<u>Cluster</u>: Informal and Formal Definitions**

    - **Set whose entities are alike and are different from entities in *other* clusters**

    - **Aggregation of points in the instance space such that distance between any two points in the cluster is less than the distance between any point in the cluster and any point not in it**

**CIS 732: Machine Learning and Pattern Recognition**

# Quick Review:
# Bayesian Learning and EM

- **Problem Definition**

  - <u>Given</u>: data (*n*-tuples) with <u>missing values</u>, *aka* <u>partially observable</u> (PO) data

  - **Want to fill in *?* with <u>expected</u> value**

- **Solution Approaches**

  - **Expected = distribution over possible values**

  - **Use "best guess" Bayesian model (e.g., BBN) to estimate distribution**

  - **<u>E</u>xpectation-<u>M</u>aximization (<u>EM</u>) algorithm can be used here**

- **Intuitive Idea**

  - **Want to find $h_{ML}$ in PO case ($D \equiv$ *unobserved variables* ° *observed variables*)**

  - **<u>Estimation step</u>: calculate *E*[*unobserved variables* | *h*], assuming current *h***

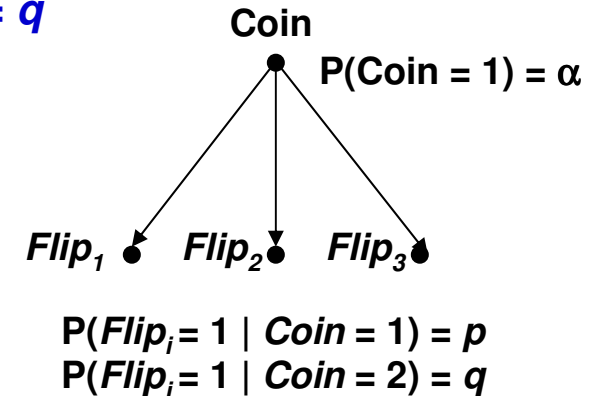  - **<u>Maximization step</u>: update $w_{ijk}$ to maximize *E*[lg *P*(*D* | *h*)], *D* $\equiv$ *all variables***

$$h_{ML} = \arg\max_{h \in H} \frac{\text{\# data cases with } \vec{n}, \vec{e}}{\text{\# data cases with } \vec{e}} = \arg\max_{h \in H} \frac{\sum_j I_{\vec{N}=\vec{n}, \vec{E}=\vec{e}}(\vec{X}_j)}{\sum_j I_{\vec{E}=\vec{e}}(\vec{X}_j)}$$

**CIS 732: Machine Learning and Pattern Recognition**

# EM Algorithm: Example [1]

- **Experiment**
  - Two coins: $P(Head\ on\ Coin\ 1) = p$, $P(Head\ on\ Coin\ 2) = q$
    - **Experimenter first selects a coin: $P(Coin = 1) = \alpha$**
    - **Chosen coin tossed 3 times (per experimental run)**
  - Observe: $D = \{(1\ H\ H\ T),\ (1\ H\ T\ T),\ (2\ T\ H\ T)\}$
  - Want to predict: $\alpha$, $p$, $q$
  - *How to model the problem?*
    - **Simple Bayesian network**
    - **Now, can find most likely values of parameters $\alpha$, $p$, $q$ given data $D$**



Coin

$P(Coin = 1) = \alpha$

$Flip_1$   $Flip_2$   $Flip_3$

$P(Flip_i = 1 \mid Coin = 1) = p$
$P(Flip_i = 1 \mid Coin = 2) = q$

- **Parameter Estimation**
  - Fully observable case: easy to estimate $p$, $q$, and $\alpha$
    - **Suppose $k$ heads are observed out of $n$ coin flips**
    - **Maximum likelihood estimate $v_{ML}$ for $Flip_i$: $p = k/n$**
  - <u>Partially observable</u> case
    - **Don't know which coin the experimenter chose**
    - **Observe: $D = \{(H\ H\ T),\ (H\ T\ T),\ (T\ H\ T)\} \equiv \{(?\ H\ H\ T),\ (?\ H\ T\ T),\ (?\ T\ H\ T)\}$**

# EM Algorithm: Example [2]

- **Problem**
  - When we knew *Coin* = 1 or *Coin* = 2, there was no problem
  - *No known analytical solution* to the partially observable problem
    - i.e., not known how to compute estimates of *p*, *q*, and $\alpha$ to get $v_{ML}$
    - Moreover, not known what the computational complexity is

- **Solution Approach: Iterative Parameter Estimation**
  - Given: a *guess* of *P*(*Coin* = 1 | *x*), *P*(*Coin* = 2 | *x*)
  - Generate "<u>fictional data points</u>", weighted according to this probability
    - *P*(*Coin* = 1 | *x*) = *P*(*x* | *Coin* = 1) *P*(*Coin* = 1) / *P*(*x*) based on our guess of $\alpha$, *p*, *q*
    - <u>Expectation</u> step (the "E" in EM)
  - Now, can find most likely values of parameters $\alpha$, *p*, *q* given "fictional" data
    - Use <u>gradient descent</u> to update our guess of $\alpha$, *p*, *q*
    - <u>Maximization</u> step (the "M" in EM)
  - Repeat until termination condition met (e.g., stopping criterion on validation set)

- **EM Converges to <u>Local Maxima</u> of the Likelihood Function *P*(*D* | $\Theta$)**

**KSU**

# EM Algorithm: Example [3]

- **Expectation Step**

  - **Suppose we observed $m$ actual experiments, each $n$ coin flips long**

    - **Each experiment corresponds to one choice of coin ($\sim\alpha$)**

    - **Let $h$ denote the number of heads in experiment $x_i$ (a single data point)**

  - **Q: How did we simulate the "fictional" data points, $E[\sum \log P(x \mid \hat{\alpha}, \hat{p}, \hat{q})]$?**

  - **A: By estimating (for $1 \le i \le m$, i.e., the real data points)**

$$P(Coin = 1 \mid \vec{x}_i) = \frac{P(\vec{x}_i \mid Coin = 1) \cdot P(Coin = 1)}{P(\vec{x}_i)}$$

$$\approx \frac{\hat{\alpha} \cdot \hat{p}^h (1 - \hat{p})^{n-h}}{\hat{\alpha} \cdot \hat{p}^h (1 - \hat{p})^{n-h} + (1 - \hat{\alpha}) \cdot \hat{q}^h (1 - \hat{q})^{n-h}}$$

- **Maximization Step**

  - **Q: What are we updating? What objective function are we maximizing?**

  - **A: We are updating $\hat{\alpha}, \hat{p}, \hat{q}$ to maximize $\frac{\partial E}{\partial \hat{\alpha}}, \frac{\partial E}{\partial \hat{p}}, \frac{\partial E}{\partial \hat{q}}$ where $E = E\left[\sum_{i=1}^{m} \log P(\vec{x}_i \mid \hat{\alpha}, \hat{p}, \hat{q})\right]$**

$$\hat{\alpha} = \frac{\sum P(Coin = 1 \mid \vec{x}_i)}{m}, \quad \hat{p} = \frac{\sum \frac{h_i}{n} P(Coin = 1 \mid \vec{x}_i)}{\sum P(Coin = 1 \mid \vec{x}_i)}, \quad \hat{q} = \frac{\sum \frac{h_i}{n} [1 - P(Coin = 1 \mid \vec{x}_i)]}{\sum [1 - P(Coin = 1 \mid \vec{x}_i)]}$$

**CIS 732: Machine Learning and Pattern Recognition**

# EM for Unsupervised Learning

- **Unsupervised Learning Problem**

  - **Objective: estimate a probability distribution with unobserved variables**

  - **Use EM to estimate <u>mixture policy</u> (more on this later; see 6.12, Mitchell)**

- **Pattern Recognition Examples**

  - **<u>H</u>uman-<u>c</u>omputer <u>i</u>ntelligent <u>i</u>nteraction (<u>HCII</u>)**

    - **Detecting facial features in emotion recognition**

    - **Gesture recognition in virtual environments**

  - **Computational medicine [Frey, 1998]**

    - **Determining morphology (shapes) of bacteria, viruses in microscopy**

    - **Identifying cell structures (e.g., nucleus) and shapes in microscopy**

  - **Other image processing**

  - **Many other examples (audio, speech, signal processing; motor control; etc.)**

- **Inference Examples**

  - **<u>Plan recognition</u>: mapping from (observed) actions to agent's (hidden) plans**

  - **<u>Hidden changes in context</u>: e.g., aviation; computer security; MUDs**

# Unsupervised Learning:
## *AutoClass* [1]

- **Bayesian Unsupervised Learning**

  - **Given: $D = \{(x_1, x_2, \ldots, x_n)\}$ (vectors of indistingushed attribute values)**

  - **Return: set of class labels that has <u>m</u>aximum <u>*a posteriori*</u> (<u>MAP</u>) probability**

- **Intuitive Idea**

  - **Bayesian learning: $h_{MAP} = \arg\max\limits_{h \in H} P(h \mid D) = \arg\max\limits_{h \in H} P(D \mid h)P(h)$**

  - **MDL/BIC (Occam's Razor): priors $P(h)$ express "cost of coding" each model $h$**

  - ***AutoClass***

    - **Define mutually exclusive, exhaustive clusters (class labels) $y_1, y_2, \ldots, y_J$**

    - **Suppose: each $y_j$ ($1 \le j \le J$) contributes to $x$**

    - **Suppose also: $y_j$'s contribution ~ *known pdf*, e.g., <u>M</u>ixture <u>o</u>f <u>G</u>aussians (<u>MoG</u>)**

    - **<u>Conjugate priors</u>: priors on $y$ of same form as priors on $x$**

- **When to Use for Clustering**

  - **Believe (or can *assume*): clusters generated by known pdf**

  - **Believe (or can *assume*): clusters combined using <u>finite mixture</u> (later)**

**CIS 732: Machine Learning and Pattern Recognition**

# Unsupervised Learning:
## *AutoClass* [2]

- *AutoClass* Algorithm [Cheeseman *et al*, 1988]
  - Based on maximizing $P(x \mid \Theta_j, y_j, J)$
    - $\Theta_j$: class (cluster) parameters (e.g., mean and variance)
    - $y_j$: synthetic classes (can estimate marginal $P(y_j)$ any time)
  - Apply Bayes's Theorem, use numerical BOC estimation techniques (cf. Gibbs)
  - Search objectives
    - Find best $J$ (*ideally*: integrate out $\Theta_j$, $y_j$; *really*: start with big $J$, decrease)
    - Find $\Theta_j$, $y_j$: use MAP estimation, then "integrate in the neighborhood" of $y_{MAP}$
- EM: Find MAP Estimate for $P(x \mid \Theta_j, y_j, J)$ by Iterative Refinement
- Advantages over Symbolic (Non-Numerical) Methods
  - Returns probability distribution over class membership
    - More robust than "best" $y_j$
    - Compare: fuzzy set membership (similar but probabilistically motivated)
  - Can deal with continuous as well as discrete data

# Unsupervised Learning: *AutoClass* [3]

- **AutoClass Resources**

  – **Beginning tutorial (*AutoClass II*): Cheeseman *et al*, 4.2.2 Buchanan and Wilkins**

  – **Project page: http://ic-www.arc.nasa.gov/ic/projects/bayes-group/autoclass/**

- **Applications**

  – **Knowledge discovery in databases (KDD) and data mining**

    - **Infrared astronomical satellite (IRAS): *spectral atlas* (*sky survey*)**

    - **Molecular biology: pre-clustering DNA acceptor, donor sites (mouse, human)**

    - ***LandSat* data from Kansas (30 km$^2$ region, 1024 x 1024 pixels, 7 channels)**

    - **Positive findings: see book chapter by Cheeseman and Stutz, online**

  – **Other typical applications: see KD Nuggets (http://www.kdnuggets.com)**

- **Implementations**

  – **Obtaining source code from project page**

    - ***AutoClass III*: Lisp implementation [Cheeseman, Stutz, Taylor, 1992]**

    - ***AutoClass C*: C implementation [Cheeseman, Stutz, Taylor, 1998]**

  – **These and others at: http://www.recursive-partitioning.com/cluster.html**

**CIS 732: Machine Learning and Pattern Recognition**

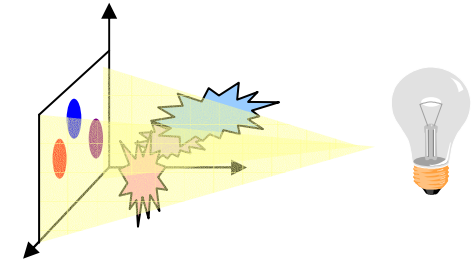# Unsupervised Learning: Competitive Learning for Feature Discovery

- **Intuitive Idea: Competitive Mechanisms for Unsupervised Learning**

  - **Global organization from local, competitive weight update**

    - **Basic principle expressed by Von der Malsburg**

    - **Guiding examples from (neuro)biology: lateral inhibition**

  - **Previous work: Hebb, 1949; Rosenblatt, 1959; Von der Malsburg, 1973; Fukushima, 1975; Grossberg, 1976; Kohonen, 1982**

- **A Procedural Framework for Unsupervised Connectionist Learning**

  - **Start with identical ("neural") processing units, with random initial parameters**

  - **Set limit on "activation strength" of each unit**

  - **Allow units to compete for right to respond to a set of inputs**

- **Feature Discovery**

  - **Identifying (or *constructing*) new features relevant to supervised learning**

  - **Examples: finding distinguishable letter characteristics in handwriten character recognition (HCR), optical character recognition (OCR)**

  - **Competitive learning: transform $X$ into $X'$; train units in $X'$ closest to $x$**

# Unsupervised Learning: Kohonen's Self-Organizing Map (SOM) [1]

- **Another Clustering Algorithm**

  - *aka* Self-Organizing Feature Map (SOFM)

  - Given: vectors of attribute values $(x_1, x_2, ..., x_n)$

  - Returns: vectors of attribute values $(x_1', x_2', ..., x_k')$

    - Typically, $n \gg k$ ($n$ is high, $k$ = 1, 2, or 3; hence "dimensionality reducing")

    - Output: vectors $x'$, the projections of input points $x$; also get $P(x_j' \mid x_i)$

    - Mapping from $x$ to $x'$ is topology preserving

- **Topology Preserving Networks**

  - Intuitive idea: similar input vectors will map to similar clusters

  - Recall: informal definition of cluster (isolated set of mutually similar entities)

  - Restatement: "clusters of $X$ (high-D) will still be clusters of $X'$ (low-D)"

- **Representation of Node Clusters**

  - Group of neighboring artificial neural network units (neighborhood of nodes)

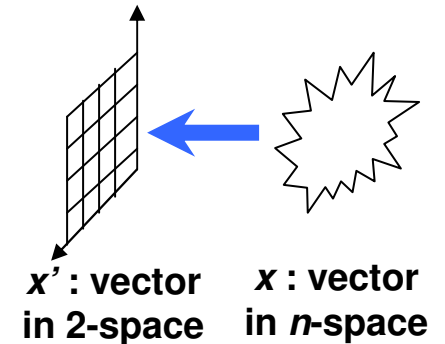  - SOMs: combine ideas of topology-preserving networks, unsupervised learning

- **Implementation: http://www.cis.hut.fi/nnrc/ and *MATLAB* NN Toolkit**

# Unsupervised Learning:
# Kohonen's Self-Organizing Map (SOM) [2]

- **Kohonen Network (SOM) for Clustering**

  - **Training algorithm: unnormalized competitive learning**

  - **Map is organized as a grid (shown here in 2D)**

    - **Each node (grid element) has a weight vector $w_j$**

    - **Dimension of $w_j$ is $n$ (same as input vector)**

    - **Number of trainable parameters (weights): $m \cdot m \cdot n$ for an $m$-by-$m$ SOM**

    - **1999 state-of-the-art: typical small SOMs 5-20, "industrial strength" > 20**

  - **Output found by selecting $j^*$ whose $w_j$ has minimum Euclidean distance from $x$**

    - ***Only one active node, aka Winner-Take-All (WTA): winning node $j^*$***

    - **i.e., $j^* = arg\ min_j \| w_j - x \|_2$**
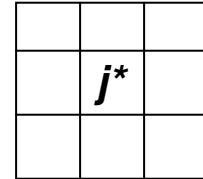
- **Update Rule**

  - **Same as competitive learning algorithm, with one modification**

  - **Neighborhood function associated with $j^*$ spreads the $w_j$ around**

$$\vec{w}_j(t+1) = \begin{cases} \vec{w}_j(t) + r(t)h_{j,j^*}\left(\vec{x} - \vec{w}_j(t)\right) & \textit{if}\ \ j \in \textbf{Neighborhood}(j^*) \\ \vec{w}_j(t) & \textit{otherwise} \end{cases}$$

**$x'$ : vector in 2-space**   **$x$ : vector in $n$-space**

# Unsupervised Learning: Kohonen's Self-Organizing Map (SOM) [3]

- **Traditional Competitive Learning**

    - **Only train $j*$**

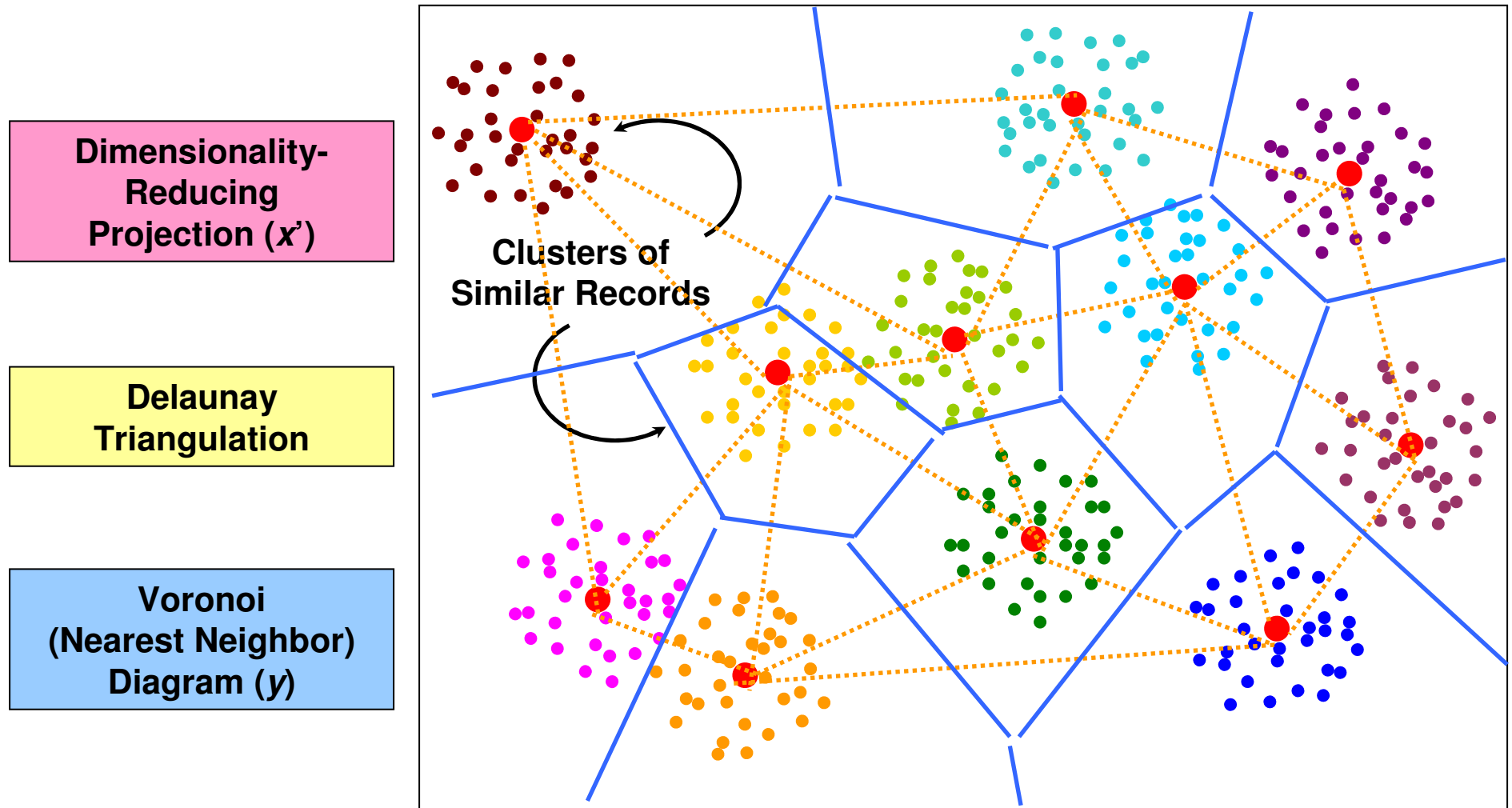    - **Corresponds to neighborhood of 0**

- **Neighborhood Function $h_{j, j*}$**

    - **For 2D Kohonen SOMs, $h$ is typically a square or hexagonal region**

        - **$j*$, the winner, is at the center of *Neighborhood* ($j*$)**

        - **$h_{j*, j*} \equiv 1$**

    - **Nodes in *Neighborhood* ($j$) updated whenever $j$ wins, i.e., $j* = j$**

    - **Strength of information fed back to $w_j$ is inversely proportional to its distance from the $j*$ for each $x$**

    - **Often use exponential or Gaussian (normal) distribution on neighborhood to decay weight delta as distance from $j*$ increases**

- **Annealing of Training Parameters**

    - **Neighborhood must shrink to 0 to achieve convergence**

    - **$r$ (learning rate) must also decrease monotonically**

|  |  |  |
|---|---|---|
|  |  |  |
|  | *j** |  |
|  |  |  |

**Neighborhood of 1**

**CIS 732: Machine Learning and Pattern Recognition**

# Unsupervised Learning: SOM and Other Projections for Clustering



Dimensionality-Reducing Projection ($x'$)

Delaunay Triangulation

Voronoi (Nearest Neighbor) Diagram ($y$)

Clusters of Similar Records

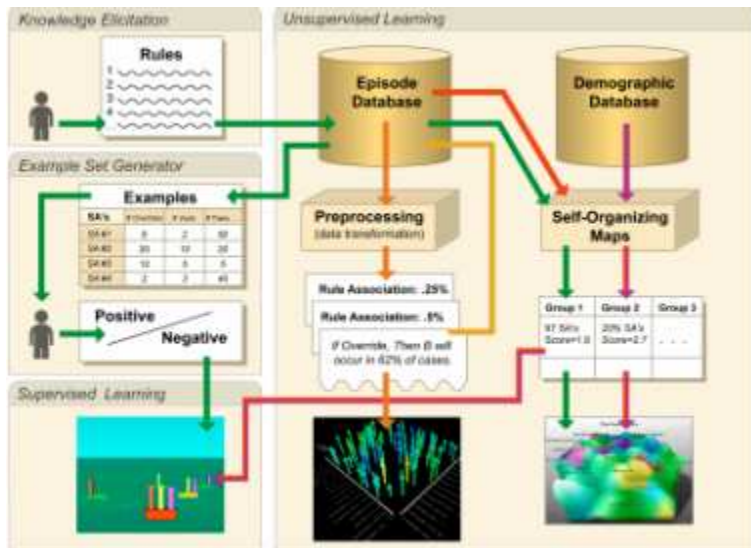**Cluster Formation and Segmentation Algorithm (Sketch)**

# Unsupervised Learning: Other Algorithms (PCA, Factor Analysis)

- **Intuitive Idea**

  - **Q:** *Why are dimensionality-reducing transforms good for supervised learning?*

  - **A: There may be many attributes with undesirable properties, e.g.,**

    - <u>Irrelevance</u>: $x_i$ has little discriminatory power over $c(x) = y_i$

    - <u>Sparseness of information</u>: "feature of interest" spread out over many $x_i$'s (e.g., text document categorization, where $x_i$ is a word position)

    - We want to increase the "information density" by "squeezing $X$ down"

- **<u>P</u>rincipal <u>C</u>omponents <u>A</u>nalysis (<u>PCA</u>)**

  - **Combining redundant variables into a single variable (*aka* <u>component</u>, or <u>factor</u>)**

  - **<u>Example</u>: ratings (e.g., Nielsen) and polls (e.g., Gallup); responses to certain questions may be correlated (e.g., "like fishing?" "time spent boating")**

- **<u>F</u>actor <u>A</u>nalysis (<u>FA</u>)**

  - **General term for a class of algorithms that includes PCA**

  - **Tutorial: http://www.statsoft.com/textbook/stfacan.html**

# Clustering Methods: Design Choices

- **Intuition**
  - **Functional (<u>declarative</u>) definition: easy ("We recognize a cluster when we see it")**
  - **Operational (procedural, <u>constructive</u>) definition: much harder to give**
  - **Possible reason: clustering of objects into groups has <u>taxonomic semantics</u> (e.g., shape, size, time, resolution, etc.)**

- **Possible Assumptions**
  - **Data generated by a particular probabilistic model**
  - **No statistical assumptions**

- **Design Choices**
  - **Distance (similarity) measure: standard metrics, transformation-invariant metrics**
    - **$L_1$ (Manhattan): $\sum |x_i - y_i|$, $L_2$ (Euclidean): $\sqrt{\sum (x_i - y_i)^2}$, $L_\infty$ (Sup): max $|x_i - y_i|$**
    - **Symmetry: Mahalanobis distance**
    - **Shift, scale invariance: covariance matrix**
  - **Transformations (e.g., <u>covariance diagonalization</u>: rotate axes to get rotational invariance, cf. PCA, FA)**
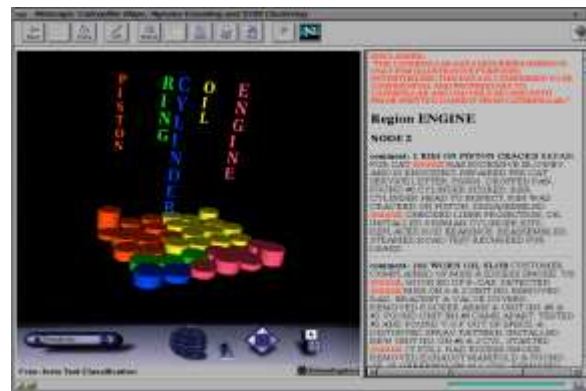
# Clustering: Applications



NCSA *D2K* 1.0 - http://www.ncsa.uiuc.edu/STI/ALG/

## Transactional Database Mining



**Data from T. Mitchell's web site:**
**http://www.cs.cmu.edu/~tom/faces.html**

**http://www.cnl.salk.edu/~wiskott/Bibliographies/**
**FaceFeatureFinding.html**

## Facial Feature Extraction



6500 news stories
from the WWW
in 1997

*ThemeScapes* - http://www.cartia.com



**Confidential and proprietary to Caterpillar; may only be used with prior written consent from Caterpillar.**

## Information Retrieval: Text Document Categorization

NCSA *D2K* 2.0 - http://www.ncsa.uiuc.edu/STI/ALG/

**CIS 732: Machine Learning and Pattern Recognition**

# Unsupervised Learning and Constructive Induction

- **Unsupervised Learning in Support of Supervised Learning**

    - **Given:** $D \equiv$ labeled vectors (*x, y*)

    - **Return:** $D' \equiv$ <u>transformed training examples</u> (*x', y'*)

    - <u>Solution approach</u>: constructive induction

        - **Feature "construction": generic term**

        - **Cluster definition**

- **Feature Construction: Front End**

    - **Synthesizing new attributes**

        - **Logical:** $x_1 \vee \neg x_2$, **arithmetic:** $x_1 + x_5 / x_2$

        - **Other synthetic attributes:** $f(x_1, x_2, \dots, x_n)$, **etc.**

    - **Dimensionality-reducing projection, feature extraction**

    - **Subset selection: finding relevant attributes for a given target *y***

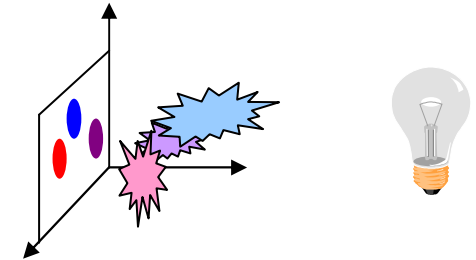    - **Partitioning: finding relevant attributes for given targets $y_1, y_2, \dots, y_p$**

- **Cluster Definition: Back End**

    - **Form, segment, and label clusters to get <u>intermediate</u> targets *y'***

    - <u>Change of representation</u>: **find an (*x', y'*) that is good for learning target *y***

**Constructive Induction**

(x, y)

*Feature (Attribute) Construction and Partitioning*

$x' / (x_1', \dots, x_p')$

*Cluster Definition*

(x', y') or $((x_1', y_1'), \dots, (x_p', y_p'))$

KSU

# Clustering:
# Relation to Constructive Induction

- **Clustering versus Cluster Definition**

  - **Clustering: 3-step process**

  - **Cluster definition: "back end" for feature construction**

- **Clustering: 3-Step Process**

  - <u>**Form**</u>

    - **$(x_1', \ldots, x_k')$ in terms of $(x_1, \ldots, x_n)$**

    - ***NB**: **typically** part of construction step, **sometimes** integrates both*

  - <u>**Segment**</u>

    - **$(y_1', \ldots, y_J')$ in terms of $(x_1', \ldots, x_k')$**

    - ***NB**: number of clusters **J** not necessarily same as number of dimensions **k***

  - <u>**Label**</u>

    - **Assign names (discrete/symbolic labels $(v_1', \ldots, v_J')$) to $(y_1', \ldots, y_J')$**

    - **Important in document categorization (e.g., clustering text for info retrieval)**

- <u>**Hierarchical Clustering**</u>: **Applying Clustering Recursively**

**CIS 732: Machine Learning and Pattern Recognition**

# Terminology

- **Expectation-Maximization (EM) Algorithm**
    - <u>Iterative refinement</u>: repeat until convergence to a locally optimal label
    - <u>Expectation</u> step: estimate parameters with which to simulate data
    - <u>Maximization</u> step: use simulated ("fictitious") data to update parameters

- **Unsupervised Learning and Clustering**
    - <u>Constructive induction</u>: using unsupervised learning *for* supervised learning
        - <u>Feature construction</u>: "front end" - *construct* new *x* values
        - <u>Cluster definition</u>: "back end" - use these to reformulate *y*
    - <u>Clustering</u> problems: <u>formation</u>, <u>segmentation</u>, <u>labeling</u>
    - Key criterion: <u>distance metric</u> (points closer <u>intra</u>-cluster than <u>inter</u>-cluster)
    - Algorithms
        - *<u>AutoClass</u>*: Bayesian clustering
        - <u>P</u>rincipal <u>C</u>omponents <u>A</u>nalysis (<u>PCA</u>), <u>f</u>actor <u>a</u>nalysis (<u>FA</u>)
        - <u>S</u>elf-<u>O</u>rganizing <u>M</u>aps (<u>SOM</u>): <u>topology preserving transform</u> (<u>dimensionality reduction</u>) for <u>competitive</u> unsupervised learning

# Summary Points

- **Expectation-Maximization (EM) Algorithm**

- **Unsupervised Learning and Clustering**

  - **Types of unsupervised learning**
    - **Clustering, vector quantization**
    - **Feature extraction (typically, dimensionality reduction)**

  - **Constructive induction: unsupervised learning *in support of* supervised learning**
    - **Feature construction (*aka* feature extraction)**
    - **Cluster definition**

  - **Algorithms**
    - **EM: mixture parameter estimation (e.g., for *AutoClass*)**
    - ***AutoClass*: Bayesian clustering**
    - **Principal Components Analysis (PCA), factor analysis (FA)**
    - **Self-Organizing Maps (SOM): projection of data; competitive algorithm**

  - **Clustering problems: formation, segmentation, labeling**

- **Next Lecture: Time Series Learning and Characterization**