

SDM'2010
Columbus, OH

On the Power of Ensemble: Supervised and Unsupervised Methods Reconciled*

Jing Gao¹, Wei Fan², Jiawei Han¹

1 Department of Computer Science
University of Illinois

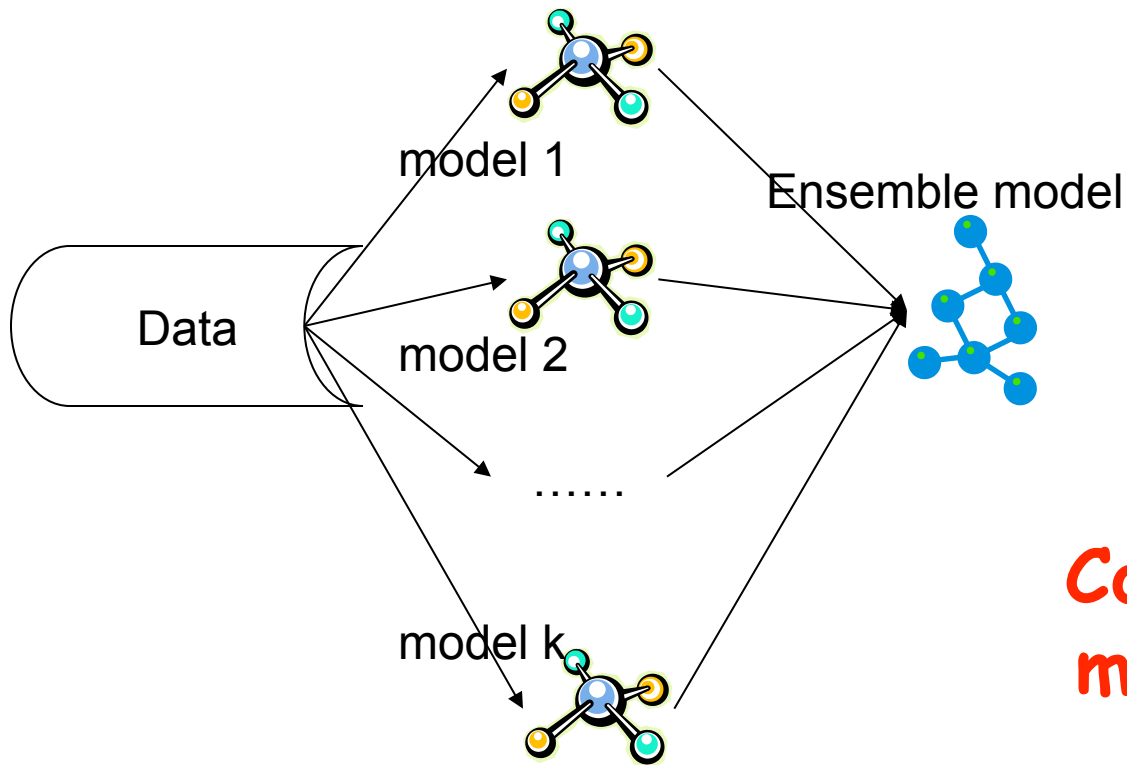
2 IBM TJ Watson Research Center

*Slides and references available at <http://ews.uiuc.edu/~jinggao3/sdm10ensemble.htm>

Outline

- An overview of ensemble methods
 - Motivations
 - Tutorial overview
- Supervised ensemble
- Unsupervised ensemble
- Semi-supervised ensemble
 - Multi-view learning
 - Consensus maximization among supervised and unsupervised models
- Applications
 - Transfer learning, stream classification, anomaly detection

Ensemble



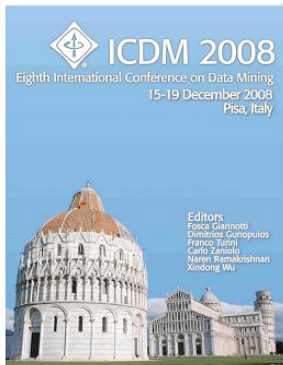
**Combine multiple
models into one!**

Applications: classification, clustering,
collaborative filtering, anomaly detection.....

Stories of Success



- **Million-dollar prize**
 - Improve the baseline movie recommendation approach of Netflix by 10% in accuracy
 - The top submissions all combine several teams and algorithms as an ensemble



- **Data mining competitions**
 - Classification problems
 - Winning teams employ an ensemble of classifiers

Netflix Prize

- **Supervised learning task**

- Training data is a set of users and ratings (1,2,3,4,5 stars) those users have given to movies.
- Construct a classifier that given a user and an unrated movie, correctly classifies that movie as either 1, 2, 3, 4, or 5 stars
- \$1 million prize for a 10% improvement over Netflix's current movie recommender

- **Competition**

- At first, single-model methods are developed, and performances are improved
- However, improvements slowed down
- Later, individuals and teams merged their results, and significant improvements are observed

Leaderboard

Rank	Team Name	Best Test Score	% Improvement	Best Submit Time
Grand Prize - RMSE = 0.8567 - Winning Team: BellKor's Pragmatic Chaos				
1	BellKor's Pragmatic Chaos	0.8567	10.06	2009-07-26 18:18:28
2	The Ensemble	0.8567	10.06	2009-07-26 18:38:22
3	Grand Prize Team	0.8582	9.90	2009-07-10 21:24:40
4	Opera Solutions and Vandelay United	0.8588	9.84	2009-07-10 01:12:31
5	Vandelay Industries !	0.8591	9.81	2009-07-10 00:32:20
6	PragmaticTheory	0.8594	9.77	2009-06-24 12:06:56
7	BellKor in BigChaos	0.8601	9.70	2009-05-13 08:14:09
8	Dace	0.8612	9.59	2009-07-24 17:18:43
9	Feeds2	0.8622	9.48	2009-07-12 13:11:51
10	BigChaos	0.8623	9.47	2009-04-07 12:33:59
12	BellKor	0.8624	9.46	2009-07-26 17:19:11
Progress Prize 2008 - RMSE = 0.8627 - Winning Team: BellKor in BigChaos				
13	xianqiang	0.8642	9.27	2009-07-15 14:53:22
14	Gravity	0.8643	9.26	2009-04-22 18:31:32
15	Ces	0.8651	9.18	2009-06-21 19:24:53

“Our final solution (RMSE=0.8712) consists of blending 107 individual results. “

“Predictive accuracy is substantially improved when blending multiple predictors. Our experience is that most efforts should be concentrated in deriving substantially different approaches, rather than refining a single technique. “

Progress Prize 2007 - RMSE = 0.8725 - Winning Team: Korben

Cinematch score - RMSE = 0.9525

Motivations

- **Motivations of ensemble methods**
 - Ensemble model improves accuracy and robustness over single model methods
 - Applications:
 - distributed computing
 - privacy-preserving applications
 - large-scale data with reusable models
 - multiple sources of data
 - Efficiency: a complex problem can be decomposed into multiple sub-problems that are easier to understand and solve (divide-and-conquer approach)

Relationship with Related Studies (1)

- **Multi-task learning**
 - Learn **multiple** tasks simultaneously
 - Ensemble methods: use multiple models to learn **one** task
- **Data integration**
 - Integrate raw data
 - Ensemble methods: integrate information at the **model** level

Relationship with Related Studies (2)

- **Meta learning**
 - **Learn** on meta-data (include base model output)
 - Ensemble methods: besides learn a joint model based on model output, we can also combine the output by **consensus**
- **Non-redundant clustering**
 - Give **multiple** non-redundant clustering solutions to users
 - Ensemble methods: give **one** solution to users which represents the consensus among all the base models

Why Ensemble Works? (1)

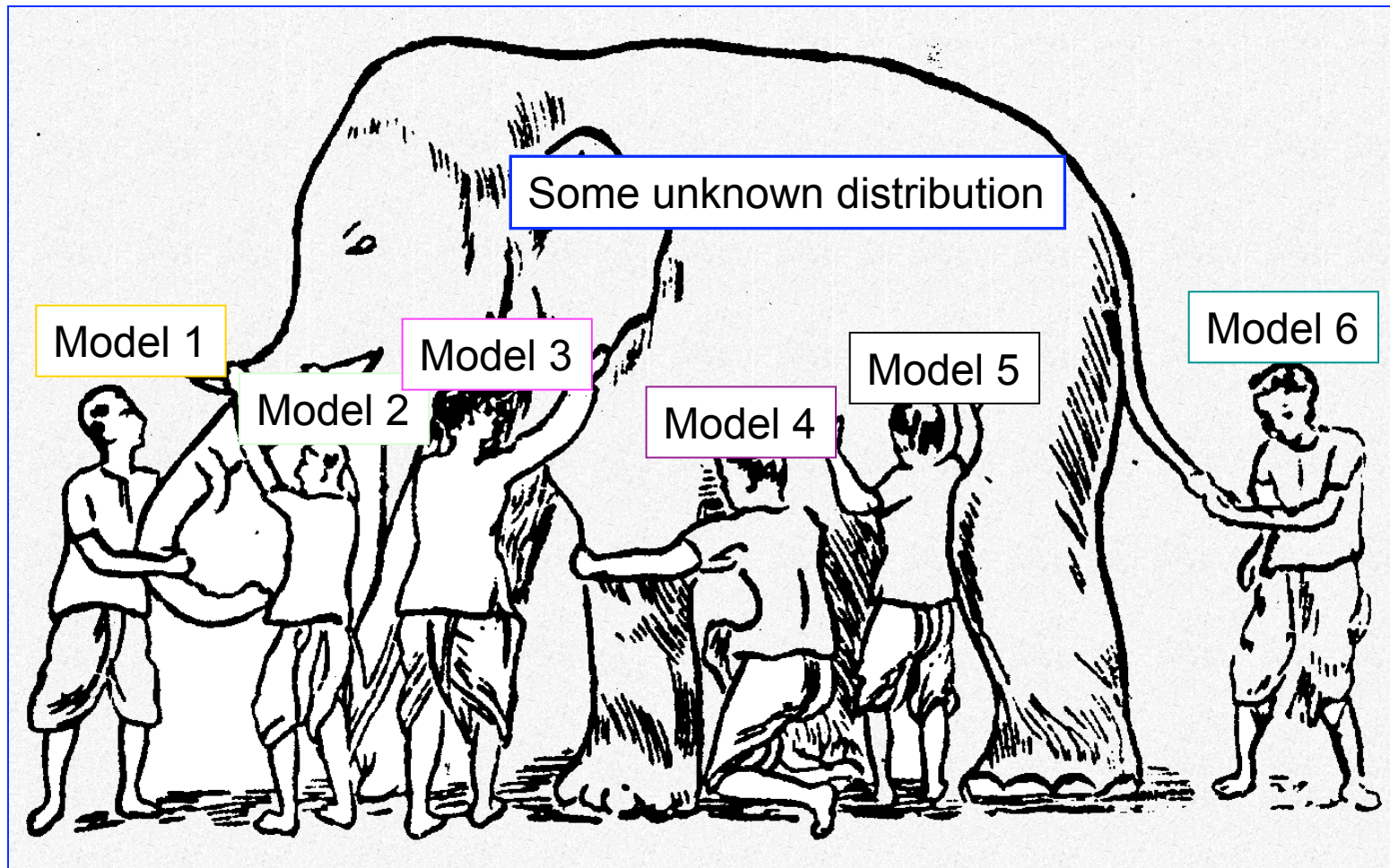
- **Intuition**

- combining diverse, independent opinions in human decision-making as a protective mechanism (e.g. stock portfolio)

- **Uncorrelated error reduction**

- Suppose we have 5 completely independent classifiers for majority voting
- If accuracy is 70% for each
 - $10 (.7^3)(.3^2)+5(.7^4)(.3)+(.7^5)$
 - **83.7% majority vote accuracy**
- 101 such classifiers
 - **99.9% majority vote accuracy**

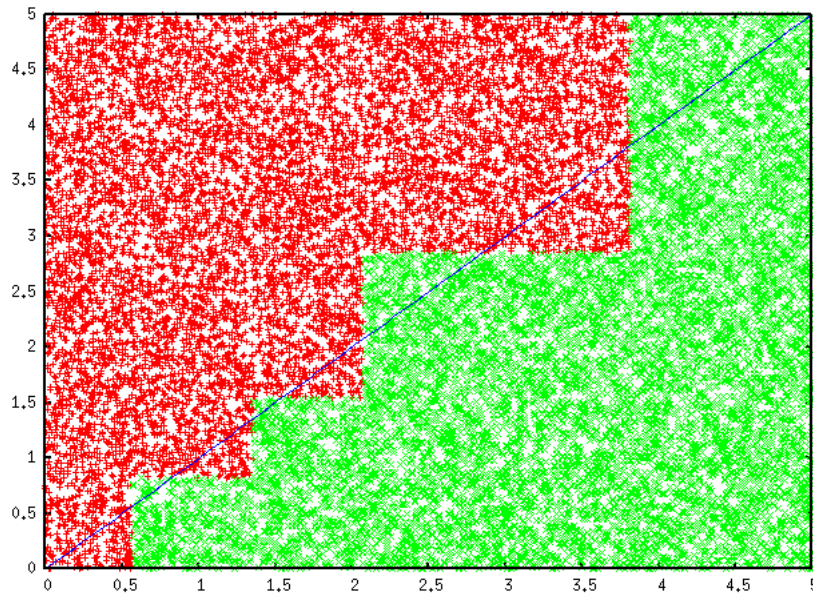
Why Ensemble Works? (2)



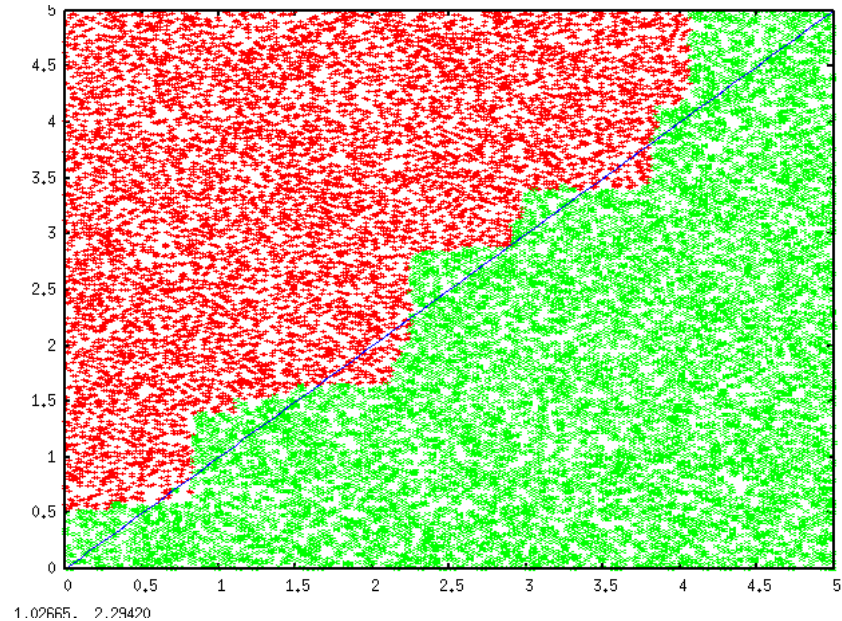
Ensemble gives the global picture!

Why Ensemble Works? (3)

- **Overcome limitations of single hypothesis**
 - The target function may not be implementable with individual classifiers, but may be approximated by model averaging



Decision Tree

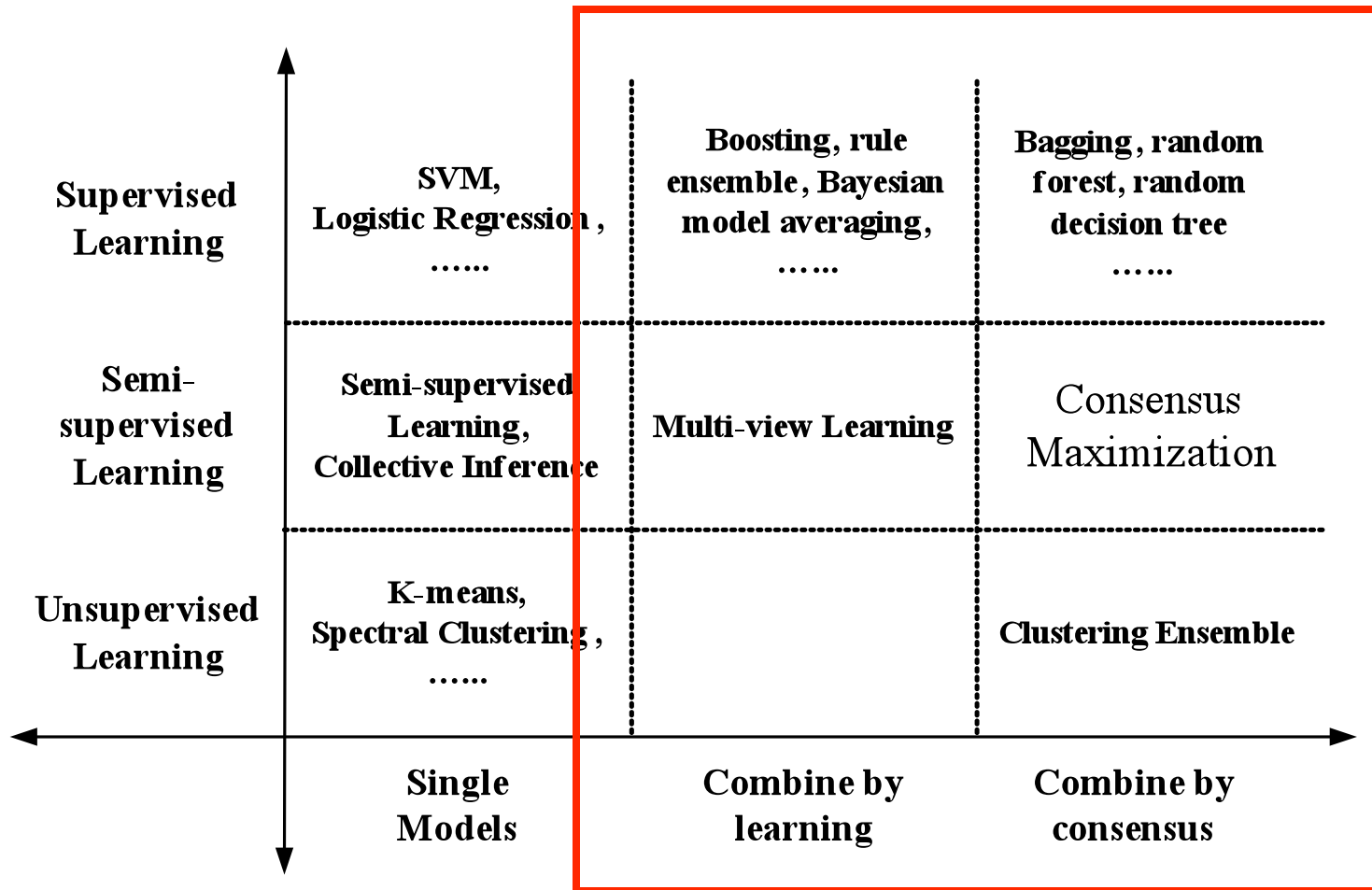


Model Averaging

Research Focus

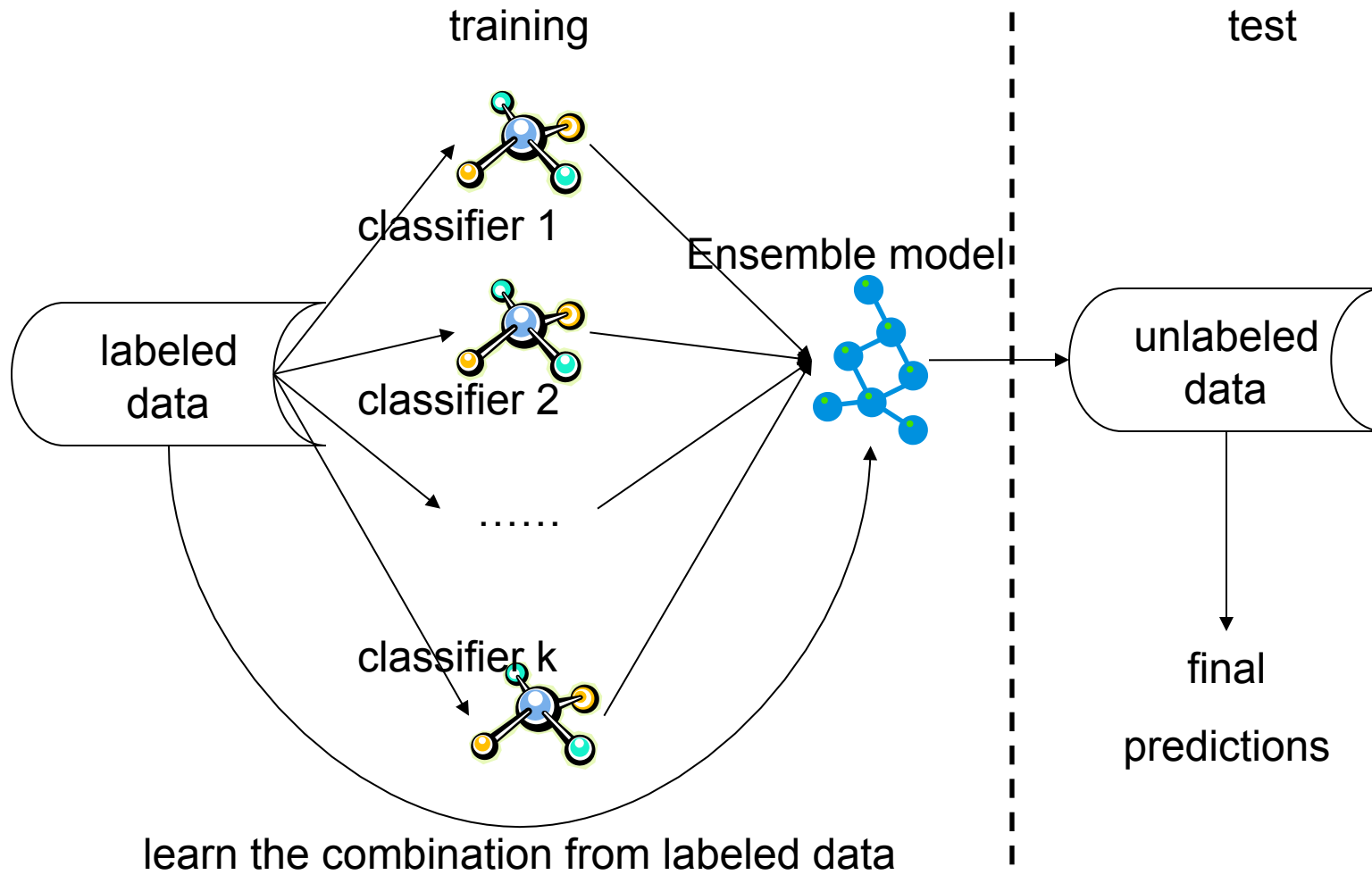
- **Base models**
 - Improve diversity!
- **Combination scheme**
 - Consensus (unsupervised)
 - Learn to combine (supervised)
- **Tasks**
 - Classification (supervised or semi-supervised ensemble)
 - Clustering (unsupervised ensemble)

Summary



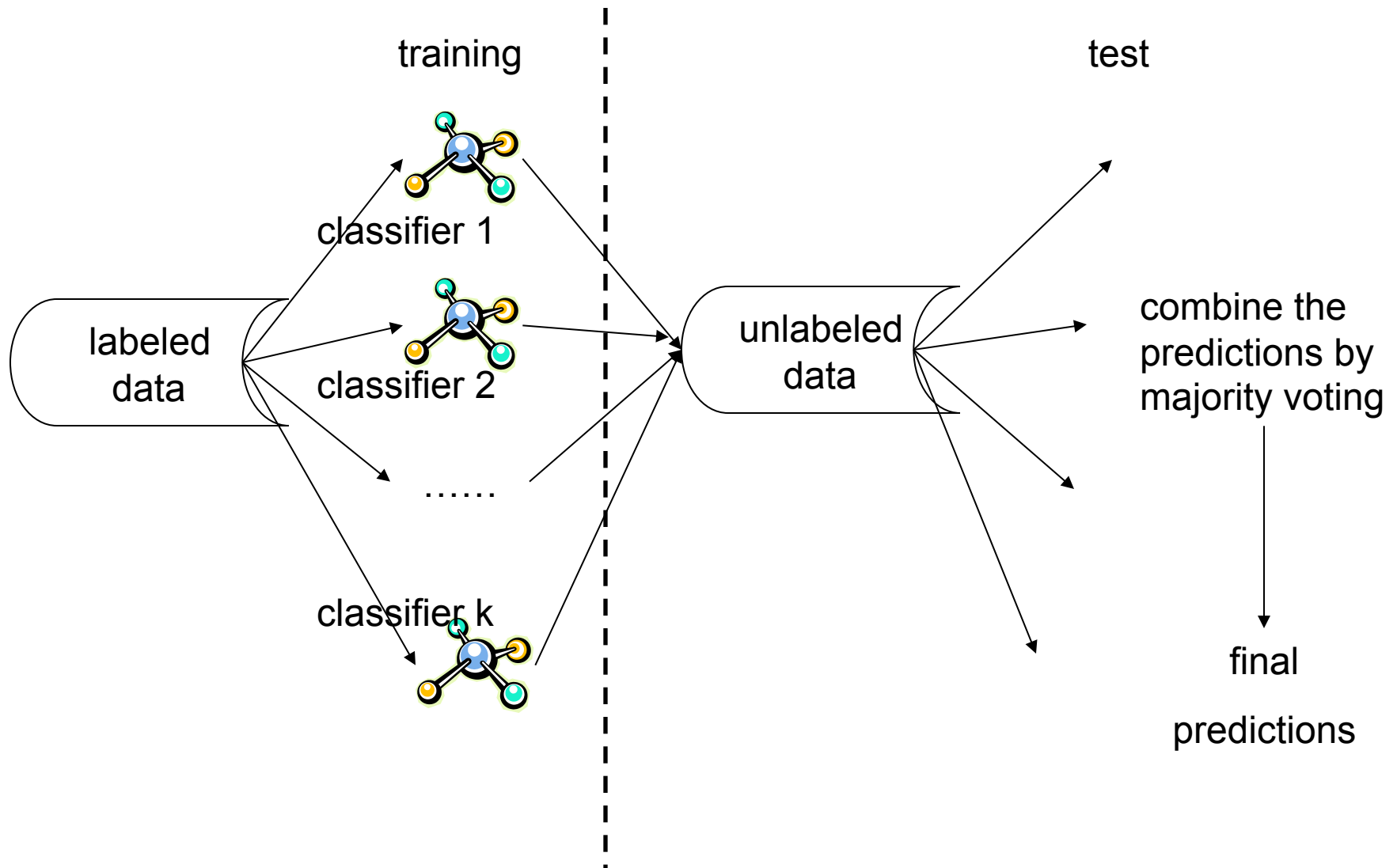
Review the ensemble methods in the tutorial

Ensemble of Classifiers—Learn to Combine



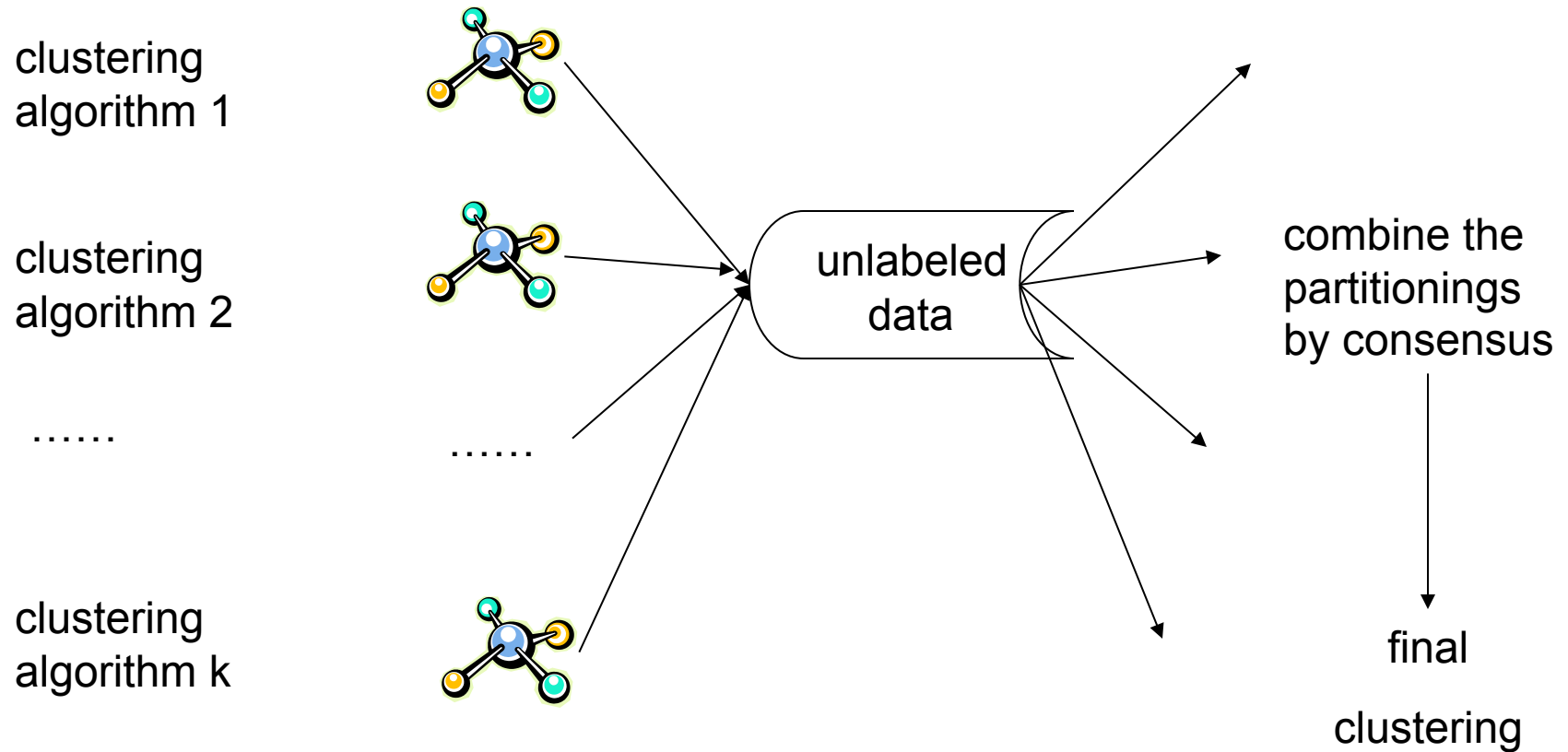
Algorithms: boosting, stacked generalization, rule ensemble, Bayesian model averaging.....

Ensemble of Classifiers—Consensus



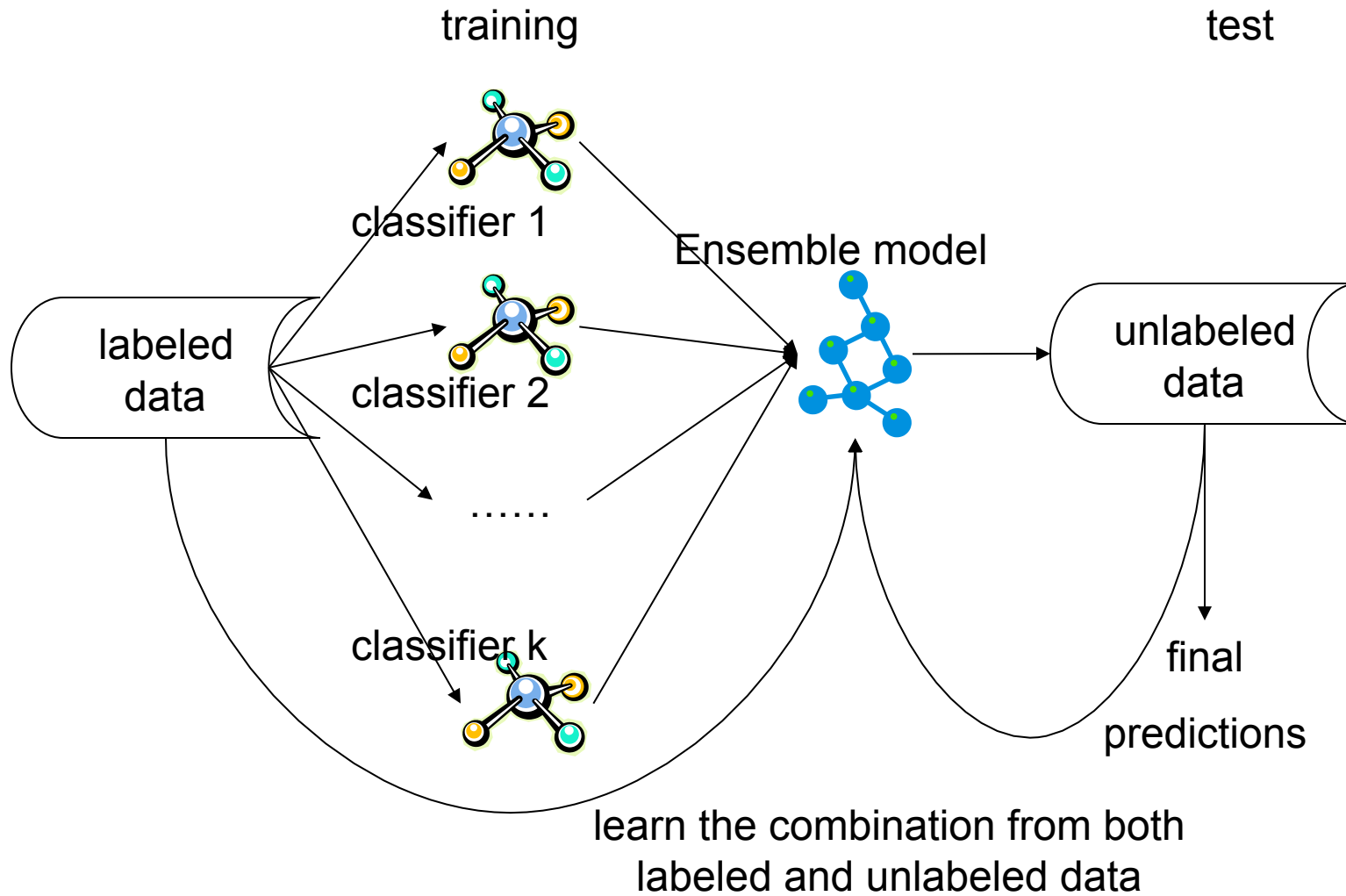
Algorithms: bagging, random forest, random decision tree, model averaging of probabilities.....

Clustering Ensemble—Consensus



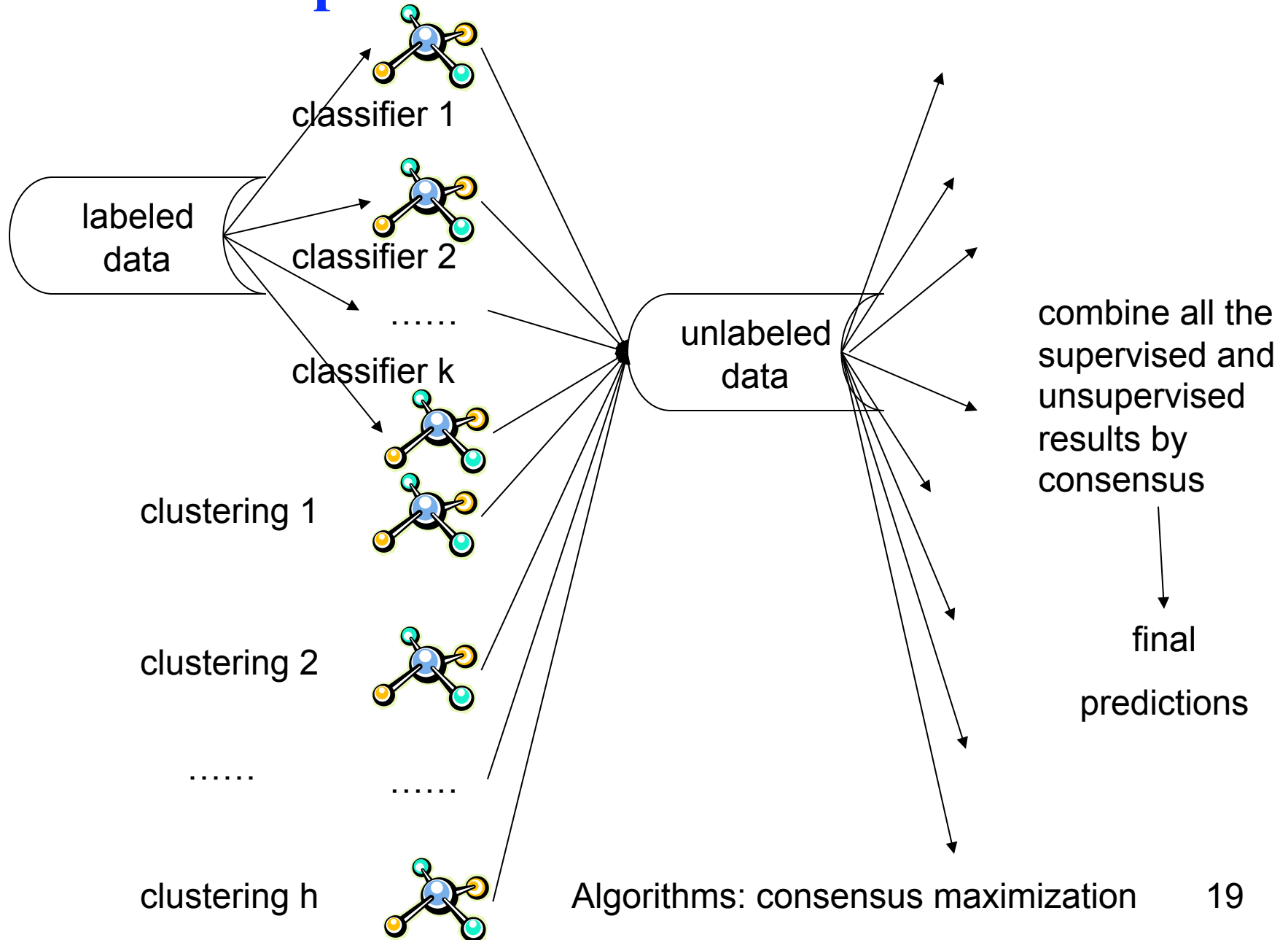
Algorithms: direct approach, object-based, cluster-based, object-cluster-based approaches, generative models

Semi-Supervised Ensemble—Learn to Combine



Algorithms: multi-view learning

Semi-supervised Ensemble—Consensus



Pros and Cons

	Combine by learning	Combine by consensus
Pros	<ul style="list-style-type: none">Get useful feedbacks from labeled dataCan potentially improve accuracy	<ul style="list-style-type: none">Do not need labeled dataCan improve the generalization performance
Cons	<ul style="list-style-type: none">Need to keep the labeled data to train the ensembleMay overfit the labeled dataCannot work when no labels are available	<ul style="list-style-type: none">No feedbacks from the labeled dataRequire the assumption that consensus is better

Outline

- An overview of ensemble methods
 - Motivations
 - Tutorial overview
- Supervised ensemble
- Unsupervised ensemble
- Semi-supervised ensemble
 - Multi-view learning
 - Consensus maximization among supervised and unsupervised models
- Applications
 - Transfer learning, stream classification, anomaly detection

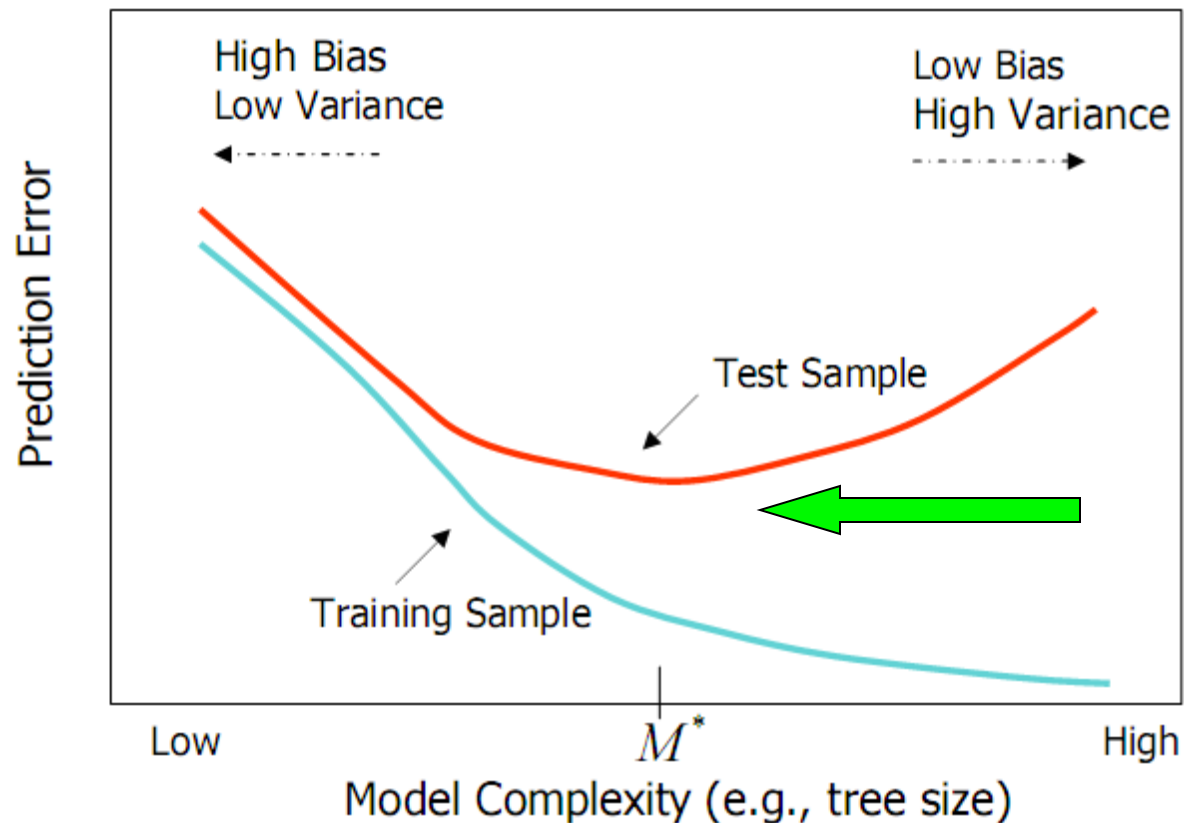
Supervised Ensemble Methods

- **Problem**

- Given a data set $D=\{x_1, x_2, \dots, x_n\}$ and their corresponding labels $L=\{l_1, l_2, \dots, l_n\}$
- An ensemble approach computes:
 - A set of classifiers $\{f_1, f_2, \dots, f_k\}$, each of which maps data to a class label: $f_j(x)=l$
 - A combination of classifiers f^* which minimizes generalization error: $f^*(x)=w_1f_1(x)+w_2f_2(x)+\dots+w_kf_k(x)$

Bias and Variance

- Ensemble methods
 - Combine learners to reduce variance



from Elder, John. From Trees to Forests and Rule Sets - A Unified Overview of Ensemble Methods. 2007.

Generating Base Classifiers

- **Sampling training examples**
 - Train k classifiers on k subsets drawn from the training set
- **Using different learning models**
 - Use all the training examples, but apply different learning algorithms
- **Sampling features**
 - Train k classifiers on k subsets of features drawn from the feature space
- **Learning “randomly”**
 - Introduce randomness into learning procedures

Bagging* (1)

- **Bootstrap**
 - Sampling with replacement
 - Contains around 63.2% original records in each sample
- **Bootstrap Aggregation**
 - Train a classifier on each bootstrap sample
 - Use majority voting to determine the class label of ensemble classifier

*[Breiman96]

Bagging (2)

Original Data:

x	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
y	1	1	1	-1	-1	-1	-1	1	1	1

Bootstrap samples and classifiers:

x	0.1	0.2	0.2	0.3	0.4	0.4	0.5	0.6	0.9	0.9
y	1	1	1	1	-1	-1	-1	-1	1	1

x	0.1	0.2	0.3	0.4	0.5	0.5	0.9	1	1	1
y	1	1	1	-1	-1	-1	1	1	1	1

x	0.1	0.2	0.3	0.4	0.4	0.5	0.7	0.7	0.8	0.9
y	1	1	1	-1	-1	-1	-1	-1	1	1

x	0.1	0.2	0.5	0.5	0.5	0.7	0.7	0.8	0.9	1
y	1	1	-1	-1	-1	-1	-1	1	1	1

Combine predictions by majority voting

Bagging (3)

- **Error Reduction**

- Under mean squared error, bagging reduces variance and leaves bias unchanged
- Consider idealized bagging estimator: $\bar{f}(x) = E(\hat{f}_z(x))$
- The error is

$$\begin{aligned} E[Y - \hat{f}_z(x)]^2 &= E[Y - \bar{f}(x) + \bar{f}(x) - \hat{f}_z(x)]^2 \\ &= E[Y - \bar{f}(x)]^2 + E[\bar{f}(x) - \hat{f}_z(x)]^2 \geq E[Y - \bar{f}(x)]^2 \end{aligned}$$

- Bagging usually decreases MSE

Boosting* (1)

- Principles

- Boost a set of weak learners to a strong learner
- Make records currently misclassified more important

- Example

- Record 4 is hard to classify
- Its weight is increased, therefore it is more likely to be chosen again in subsequent rounds

Original Data	1	2	3	4	5	6	7	8	9	10
Boosting (Round 1)	7	3	2	8	7	9	4	10	6	3
Boosting (Round 2)	5	4	9	4	2	5	1	7	4	2
Boosting (Round 3)	4	4	8	10	4	5	4	6	3	4

*[FrSc97]

Boosting (2)

- AdaBoost

- Initially, set uniform weights on all the records
- At each round
 - Create a bootstrap sample based on the weights
 - Train a classifier on the sample and apply it on the original training set
 - Records that are wrongly classified will have their weights increased
 - Records that are classified correctly will have their weights decreased
 - If the error rate is higher than 50%, start over
- Final prediction is weighted average of all the classifiers with weight representing the training accuracy

Boosting (3)

- Determine the weight

- For classifier i , its error is

$$\varepsilon_i = \frac{\sum_{j=1}^N w_j \delta(C_i(x_j) \neq y_j)}{\sum_{j=1}^N w_j}$$

- The classifier's importance is represented as:

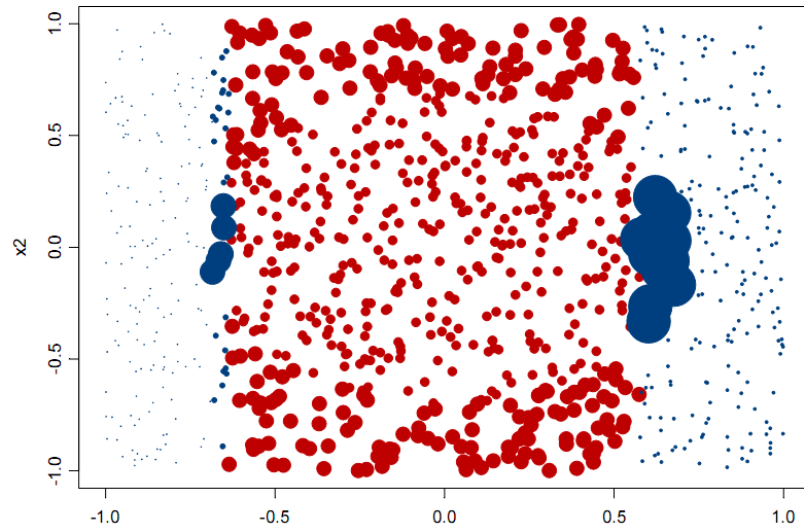
$$\alpha_i = \frac{1}{2} \ln\left(\frac{1 - \varepsilon_i}{\varepsilon_i}\right)$$

- The weight of each record is updated as:

$$w_j^{(i+1)} = \frac{w_j^{(i)} \exp(-\alpha_i y_j C_i(x_j))}{Z^{(i)}}$$

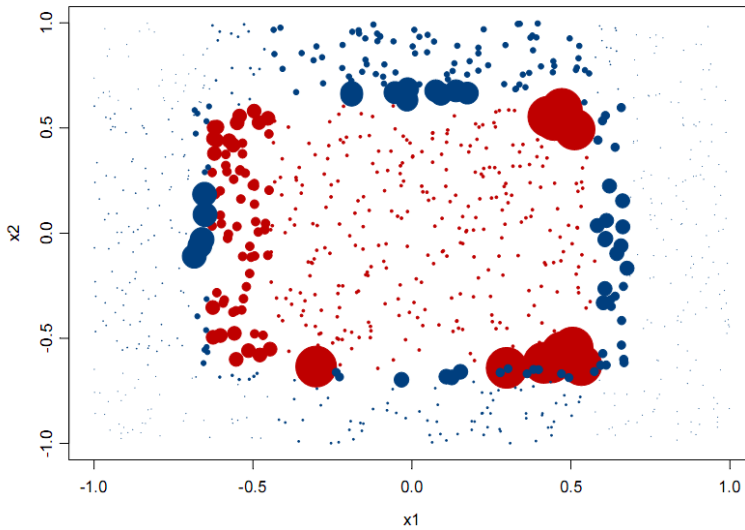
- Final combination:

$$C^*(x) = \arg \max_y \sum_{i=1}^K \alpha_i \delta(C_i(x) = y)$$

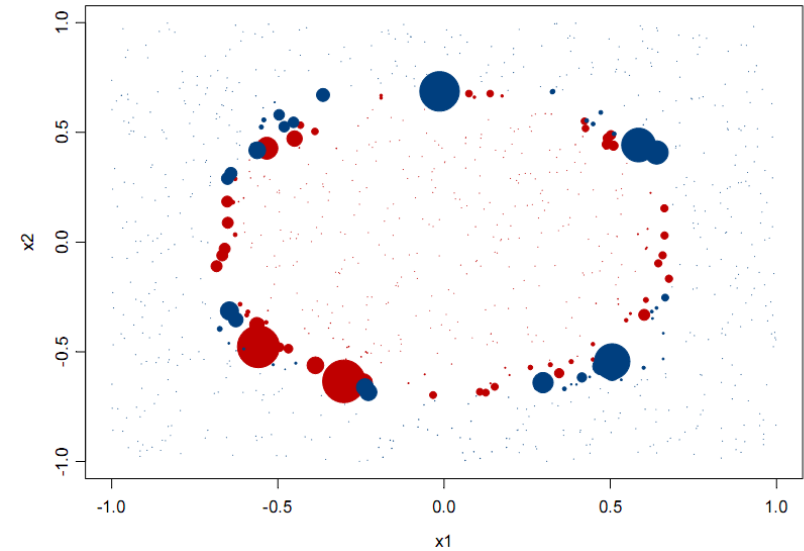


**Classifications (colors) and
Weights (size) after *1 iteration*
Of AdaBoost**

3 iterations



20 iterations



from Elder, John. From Trees to Forests and Rule Sets - A Unified Overview of Ensemble Methods. 2007.

Boosting (4)

- Explanation

- Among the classifiers of the form:

$$f(x) = \sum_{i=1}^K \alpha_i C_i(x)$$

- We seek to minimize the exponential loss function:

$$\sum_{j=1}^N \exp(-y_j f(x_j))$$

- Not robust in noisy settings

Random Forests* (1)

- **Algorithm**
 - Choose T —number of trees to grow
 - Choose $m < M$ (M is the number of total features) — number of features used to calculate the best split at each node (typically 20%)
 - For each tree
 - Choose a training set by choosing N times (N is the number of training examples) with replacement from the training set
 - For each node, randomly choose m features and calculate the best split
 - Fully grown and not pruned
 - Use majority voting among all the trees

*[Breiman01]

Random Forests (2)

- **Discussions**
 - Bagging+random features
 - Improve accuracy
 - Incorporate more diversity and reduce variances
 - Improve efficiency
 - Searching among subsets of features is much faster than searching among the complete set

Random Decision Tree* (1)

- **Single-model learning algorithms**
 - Fix structure of the model, minimize some form of errors, or maximize data likelihood (eg., Logistic regression, Naive Bayes, etc.)
 - Use some “free-form” functions to match the data given some “preference criteria” such as information gain, gini index and MDL. (eg., Decision Tree, Rule-based Classifiers, etc.)
- **Such methods will make mistakes if**
 - Data is insufficient
 - Structure of the model or the preference criteria is inappropriate for the problem
- **Learning as Encoding**
 - Make no assumption about the true model, neither parametric form nor free form

35

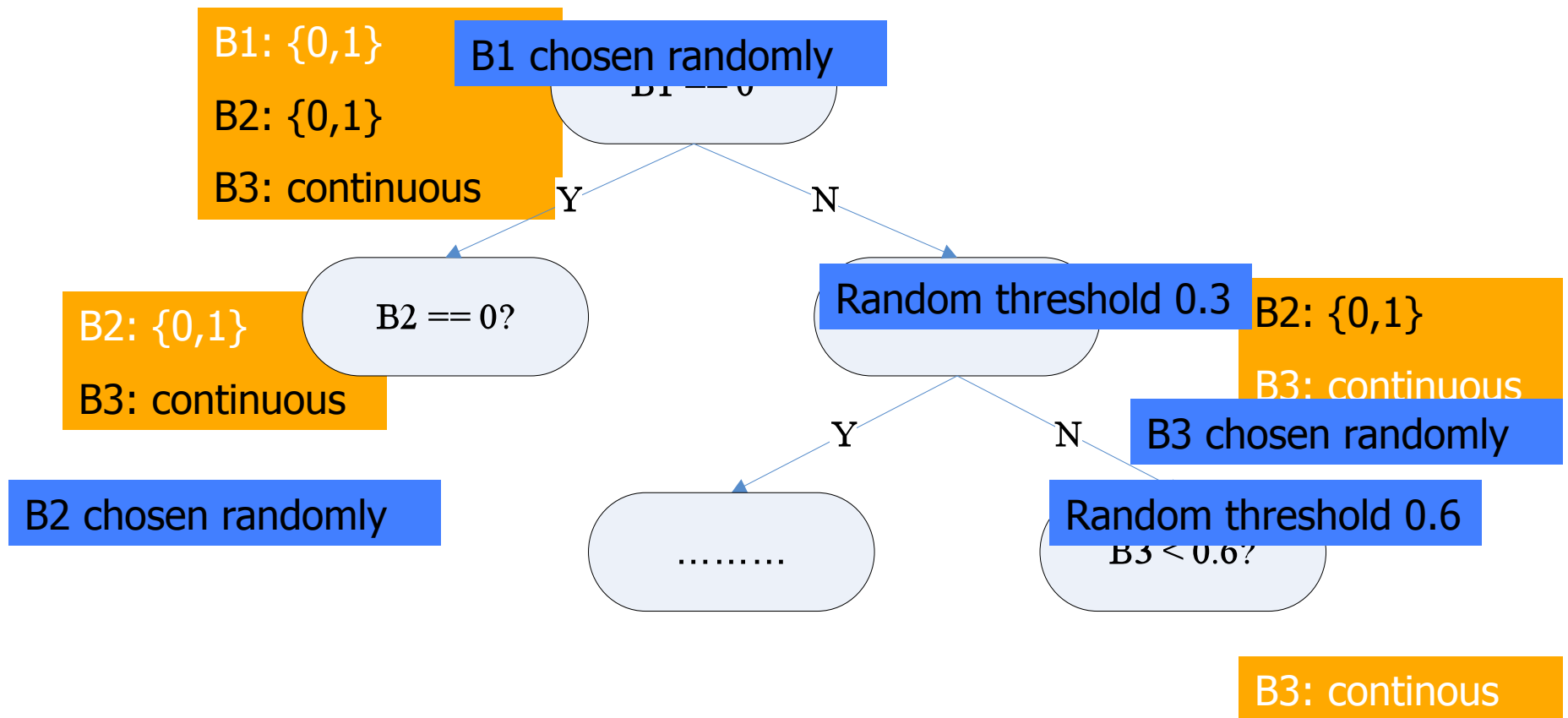
*[FWM+03] Do not prefer one base model over the other, just average them

Random Decision Tree (2)

- **Algorithm**

- At each node, an un-used feature is chosen randomly
 - A discrete feature is un-used if it has never been chosen previously on a given decision path starting from the root to the current node.
 - A continuous feature can be chosen multiple times on the same decision path, but each time a different threshold value is chosen
- We stop when one of the following happens:
 - A node becomes too small (≤ 3 examples).
 - Or the total height of the tree exceeds some limits, such as the total number of features.
- Prediction
 - Simple averaging over multiple trees

Random Decision Tree (3)

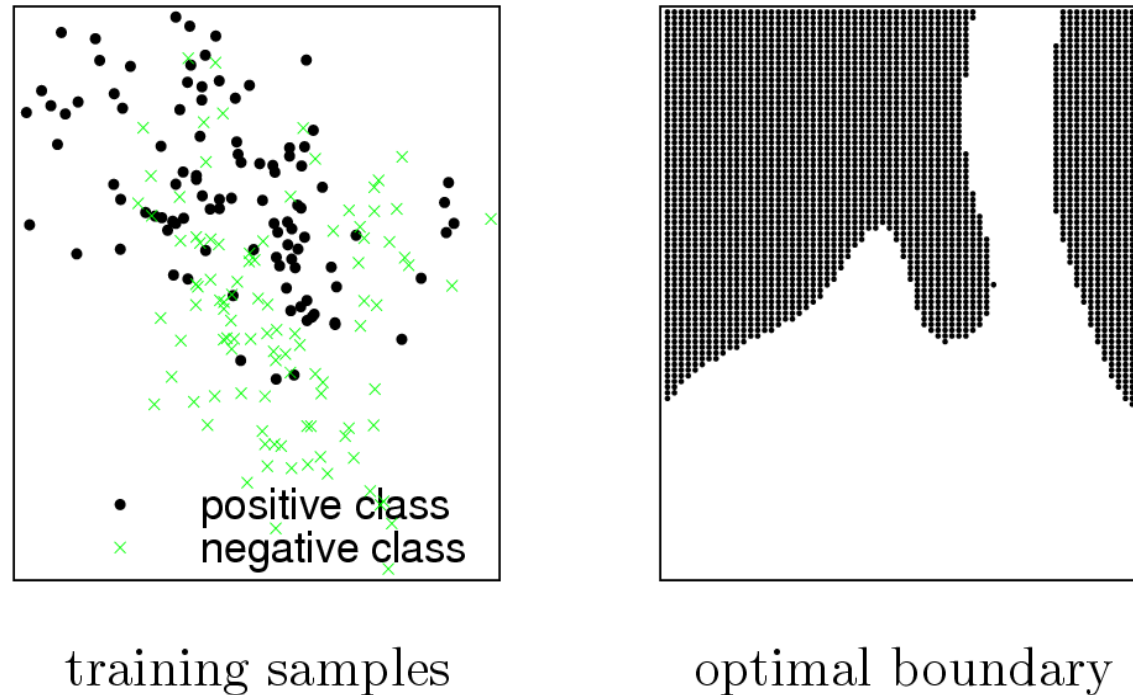


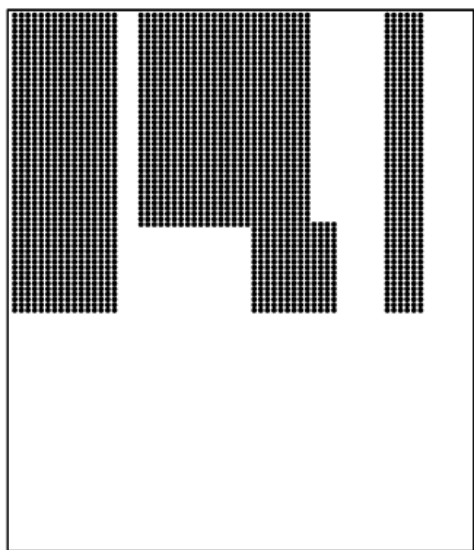
Random Decision Tree (4)

- **Potential Advantages**
 - Training can be very efficient. Particularly true for very large datasets.
 - No cross-validation based estimation of parameters for some parametric methods.
 - Natural multi-class probability.
 - Imposes very little about the structures of the model.

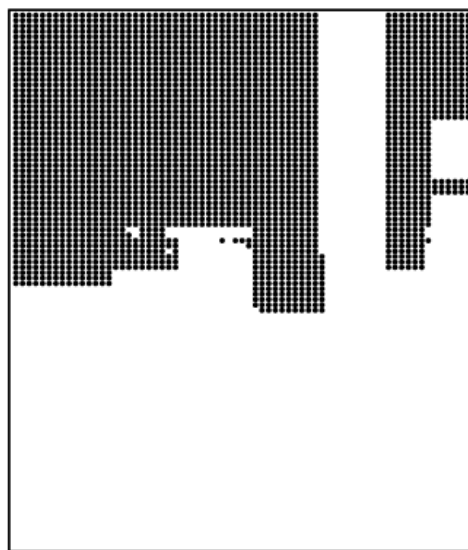
Optimal Decision Boundary

Figure 3.5: Gaussian mixture training samples and optimal boundary.

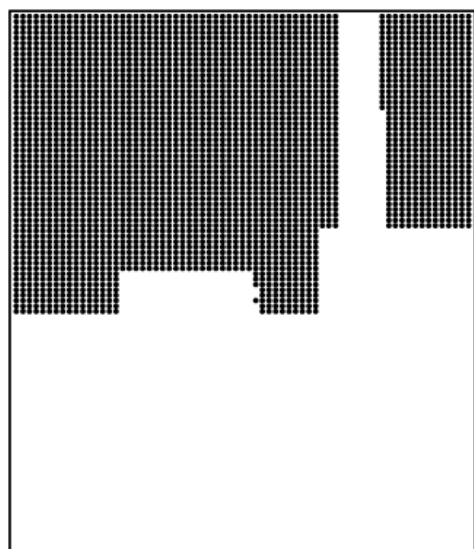
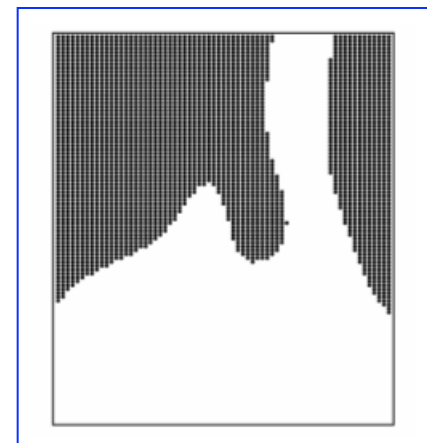




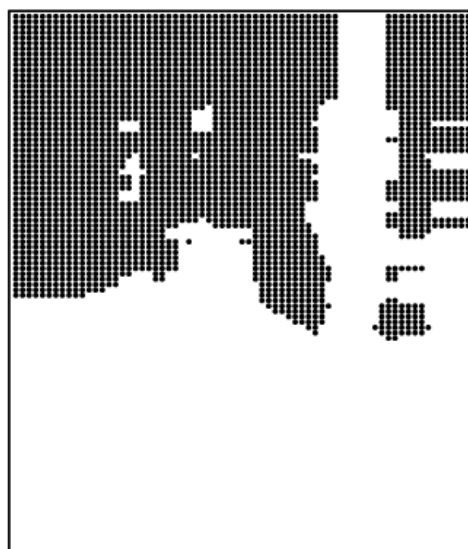
(a) unpruned C4.5



(b) Bagging



(c) Random Forests



(d) Complete-random tree ensemble

RDT looks like the optimal boundary

Outline

- An overview of ensemble methods
 - Motivations
 - Tutorial overview
- Supervised ensemble
- Unsupervised ensemble
- Semi-supervised ensemble
 - Multi-view learning
 - Consensus maximization among supervised and unsupervised models
- Applications
 - Transfer learning, stream classification, anomaly detection

Clustering Ensemble

- **Problem**

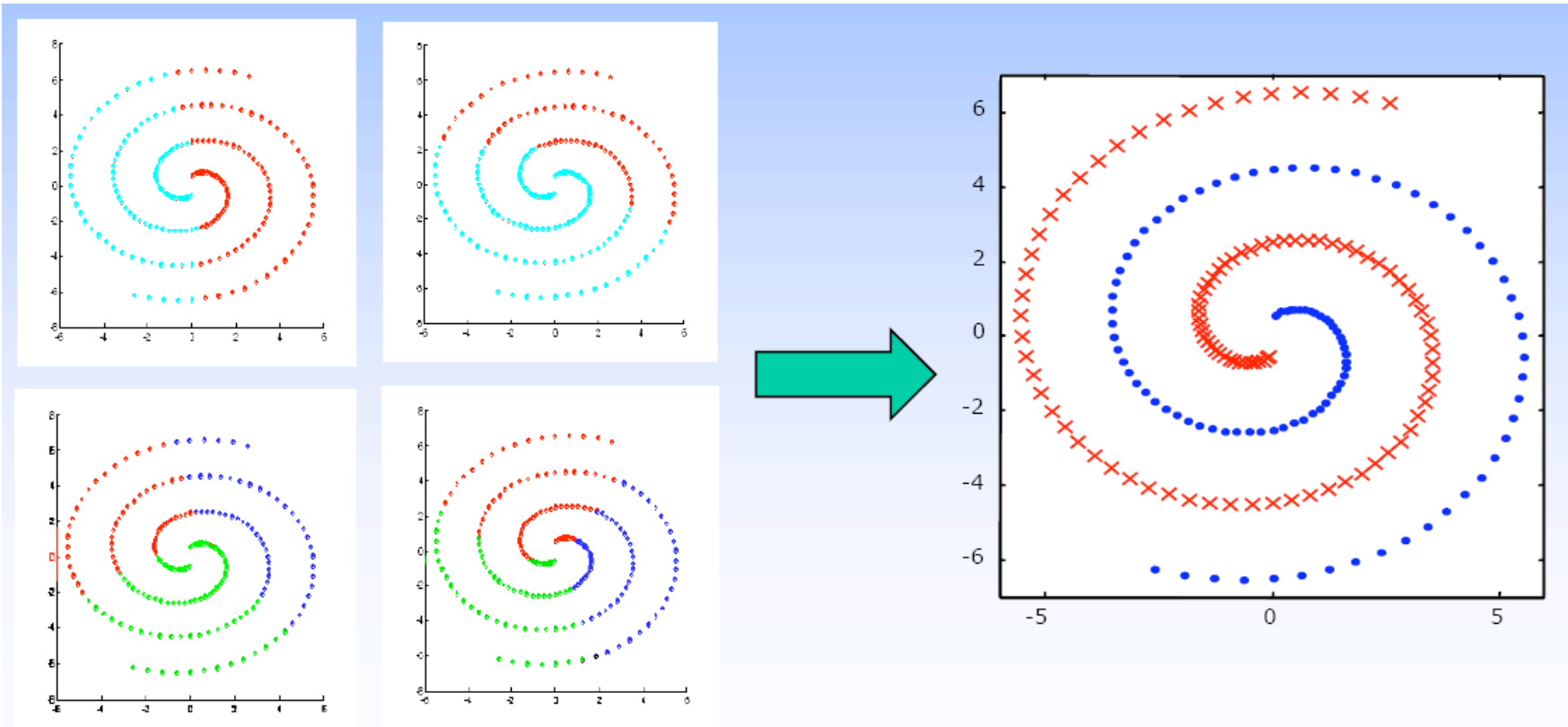
- Given an unlabeled data set $D=\{x_1, x_2, \dots, x_n\}$
- An ensemble approach computes:
 - A set of clustering solutions $\{C_1, C_2, \dots, C_k\}$, each of which maps data to a cluster: $f_j(x)=m$
 - A unified clustering solutions f^* which combines base clustering solutions by their consensus

- **Challenges**

- The correspondence between the clusters in different clustering solutions is unknown
- Unsupervised
- Combinatorial optimization problem-NP-complete

Motivations

- Goal
 - Combine “weak” clusterings to a better one



An Example

base clustering models



objects →

	\mathcal{C}_1	\mathcal{C}_2	\mathcal{C}_3	\mathcal{C}
v_1	1	1	1	1
v_2	1	2	2	2
v_3	2	1	1	1
v_4	2	2	2	2
v_5	3	3	3	3
v_6	3	4	3	3



they may not represent
the same cluster!

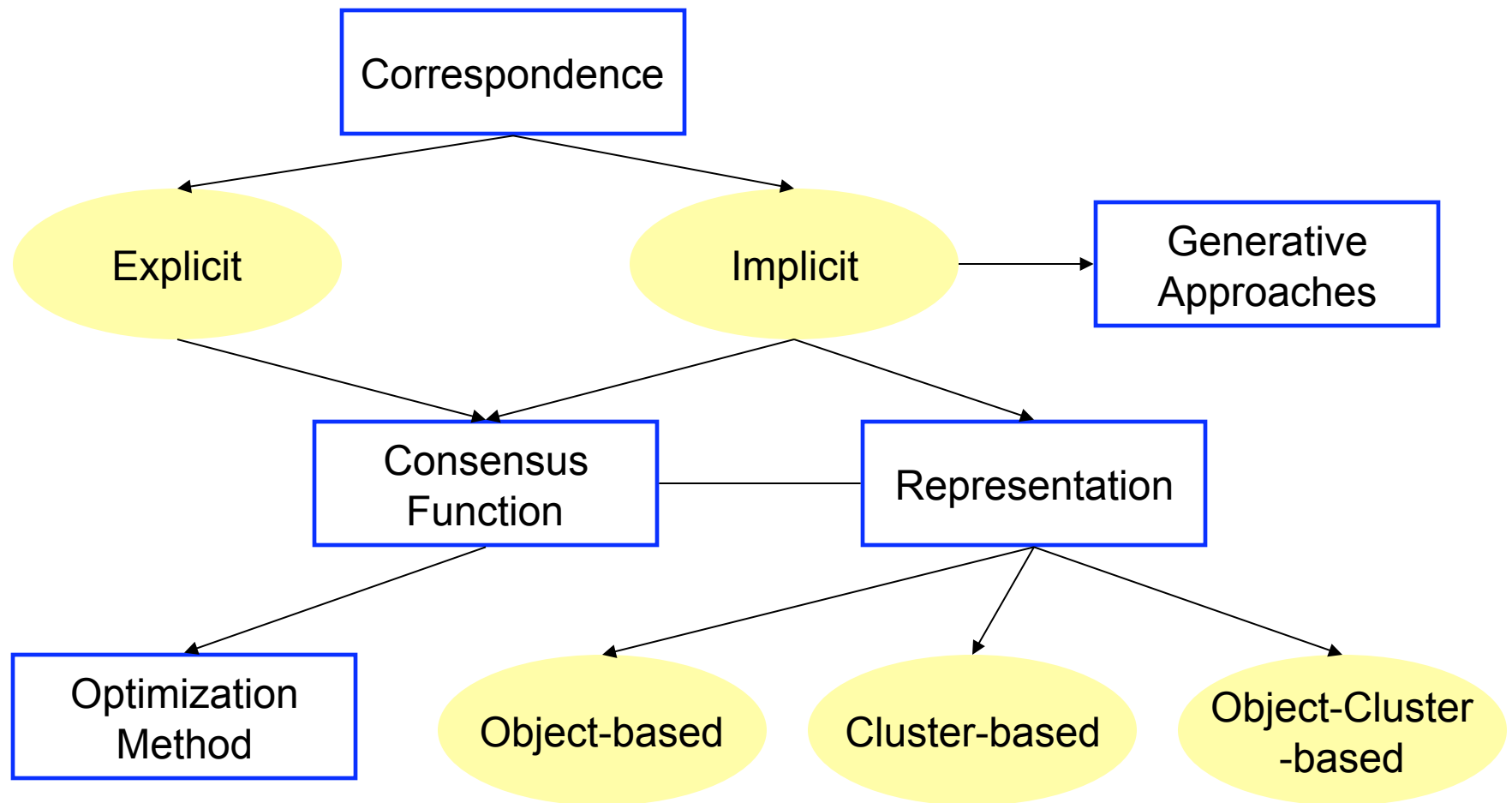
The goal: get the consensus clustering

Methods (1)

- How to get base models?
 - Bootstrap samples
 - Different subsets of features
 - Different clustering algorithms
 - Random number of clusters
 - Random initialization for K-means
 - Incorporating random noises into cluster labels
 - Varying the order of data in on-line methods such as BIRCH

Methods (2)

- How to combine the models?



Hard Correspondence (1)

- Re-labeling+voting

- Find the correspondence between the labels in the partitions and fuse the clusters with the same labels by voting [DuFr03,DWH01]

	Re-labeling			Voting						
	C ₁	C ₂	C ₃		C ₁	C ₂	C ₃	C*		
v ₁	1	3	2		v ₁	1	1	1		
v ₂	1	3	2		v ₂	1	1	1		
v ₃	2	1	2	→	v ₃	2	2	1	→	2
v ₄	2	1	3		v ₄	2	2	2		2
v ₅	3	2	1		v ₅	3	3	3		3
v ₆	3	2	1		v ₆	3	3	3		3

Hard Correspondence (2)

- **Details**

- Hungarian method to match clusters in two different clustering solutions
- Match to a reference clustering or match in a pairwise manner

- **Problems**

- In most cases, clusters do not have one-to-one correspondence

Soft Correspondence* (1)

- Notations

- Membership matrix M_1, M_2, \dots, M_k
- Membership matrix of consensus clustering M
- Correspondence matrix S_1, S_2, \dots, S_k
- $M_i S_i = M$

	C_1	C_2	C_3
v_1	1	3	2
v_2	1	3	2
v_3	2	1	2
v_4	2	1	3
v_5	3	2	1
v_6	3	2	1

$$\begin{array}{c} \mathbf{M}_2 \\ \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} \end{array} \times \begin{array}{c} \mathbf{S}_2 \\ \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \end{array} = \begin{array}{c} \mathbf{M} \\ \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \end{array}$$

*[LZY05]

Soft Correspondence (2)

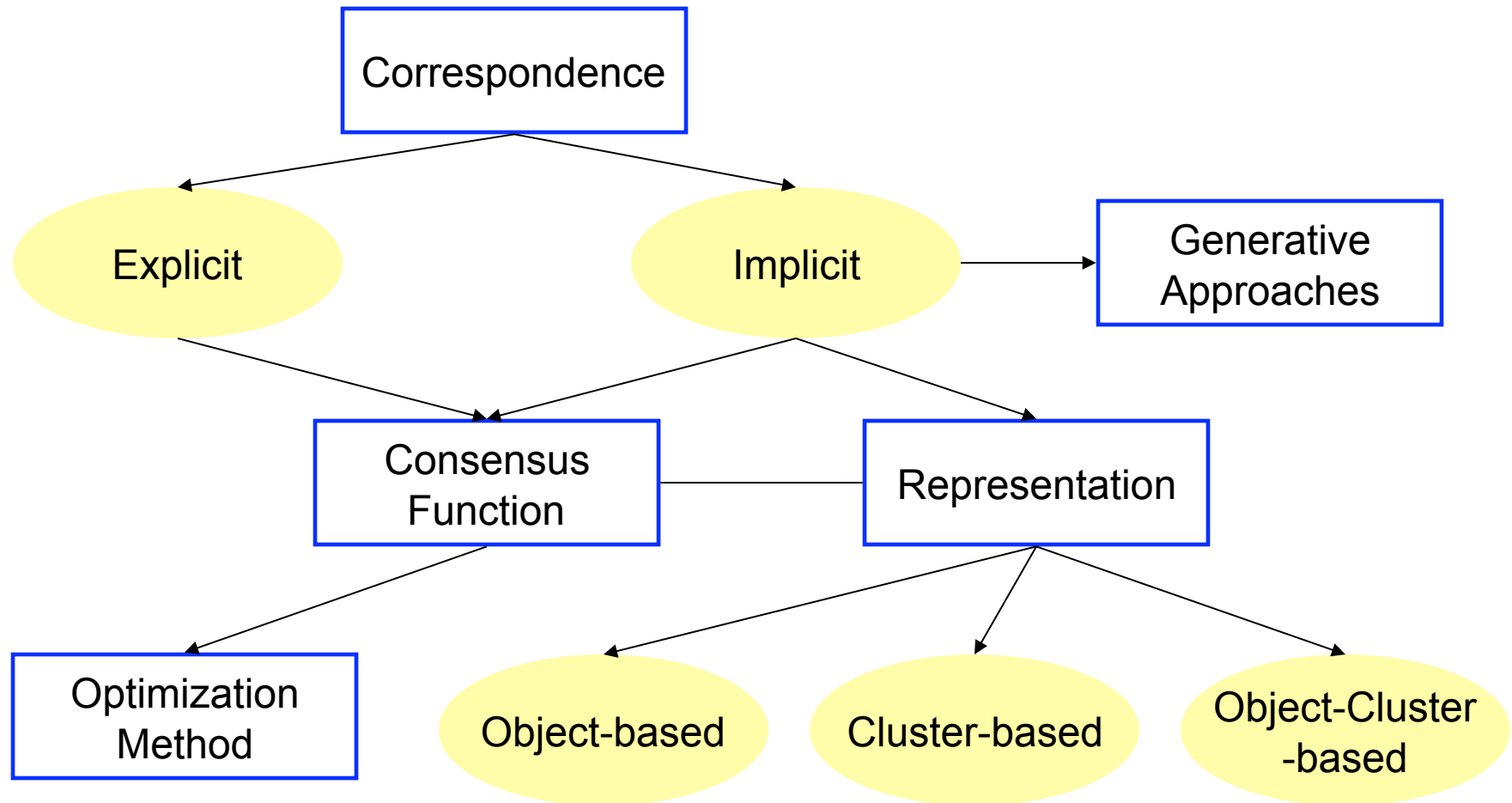
- Consensus function

- Minimize disagreement $\min \sum_{j=1}^k \|M - M_j S_j\|^2$
- Constraint 1: column-sparseness
- Constraint 2: each row sums up to 1
- Variables: M, S_1, S_2, \dots, S_k

- Optimization

- EM-based approach
- Iterate until convergence
 - Update S using gradient descent
 - Update M as $M = \frac{1}{k} \sum_{j=1}^k M_j S_j$

- How to combine the models?

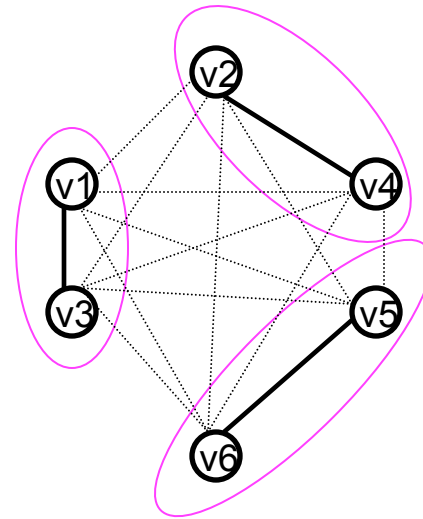


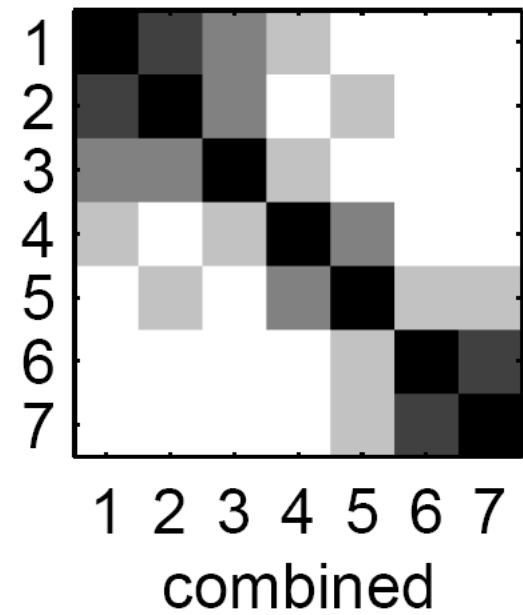
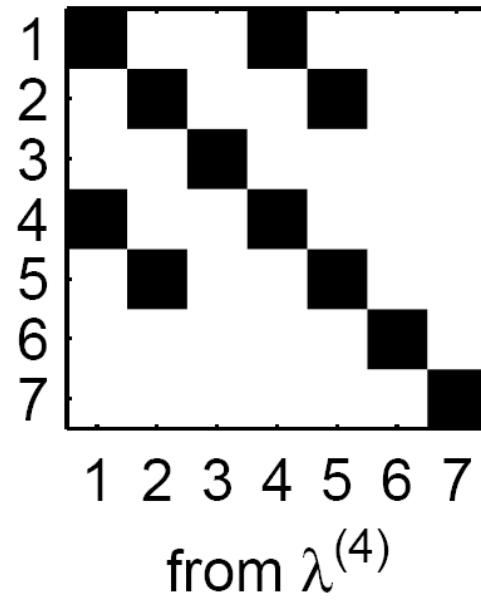
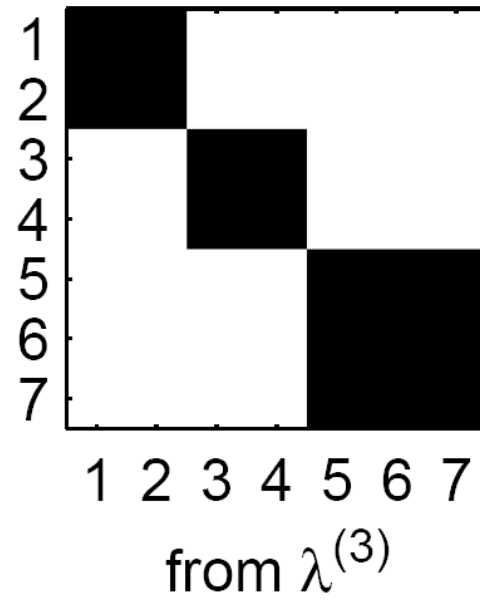
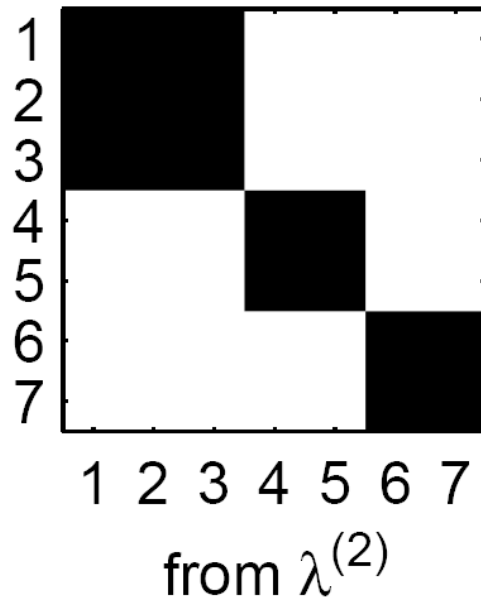
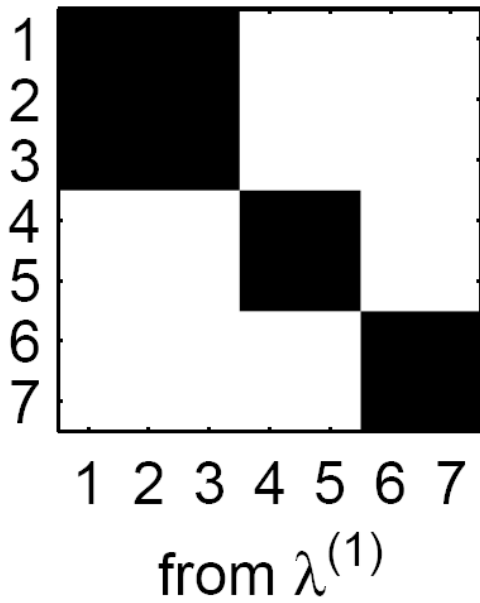
Object-based Methods (1)

- Clustering objects
 - Define a similarity or distance measure:
 - Similarity between two objects can be defined as the percentage of clusterings that assign the two objects into **same** clusters
 - Distance between two objects can be defined as the percentage of clusterings that assign the two objects into **different** clusters
 - Conduct clustering on the new similarity (distance) matrix
 - Result clustering represents the consensus
 - Can view this approach as clustering in the new feature space where clustering results are the categorical features

Object-based Methods (2)

	\mathcal{C}_1	\mathcal{C}_2	\mathcal{C}_3	\mathcal{C}
v_1	1	1	1	1
v_2	1	2	2	2
v_3	2	1	1	1
v_4	2	2	2	2
v_5	3	3	3	3
v_6	3	4	3	3





Co-association
matrix T

Consensus Function

- Minimizing disagreement

- Information-theoretic [StGh03]

$$\max \frac{1}{k} \sum_{j=1}^k NMI(T, T_j) \quad NMI(T, T_j) = \frac{I(T, T_j)}{\sqrt{H(T)H(T_j)}}$$

- Median partition [LDJ07]

$$\bar{T} = \frac{1}{k} \sum_{j=1}^k T_j \quad \min \|\bar{T} - T\|^2$$

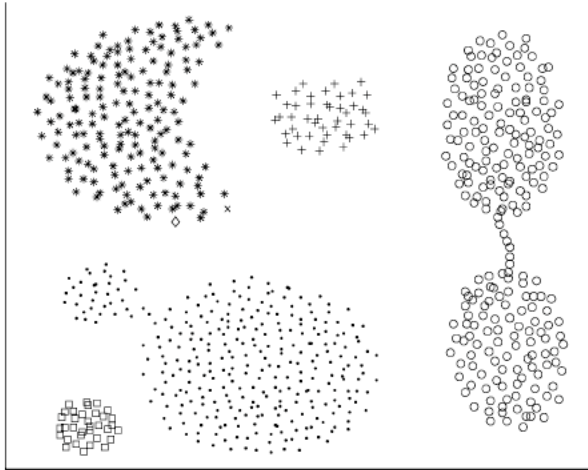
- Correlation clustering [GMT07]

$$\max \sum_{\substack{(u,v) \\ C(u)=C(v)}} \bar{T}_{uv} + \sum_{\substack{(u,v) \\ C(u) \neq C(v)}} (1 - \bar{T}_{uv})$$

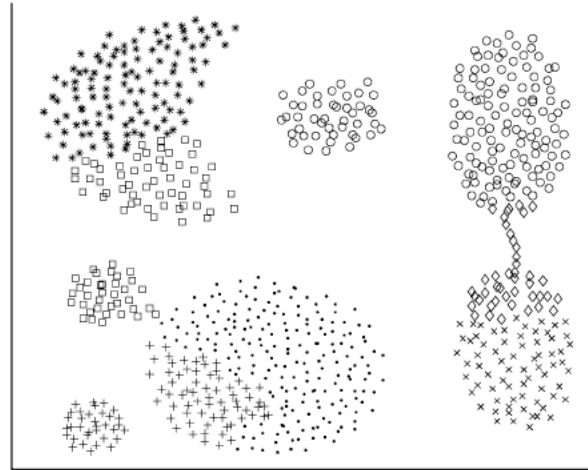
Optimization Method

- **Approximation**
 - Agglomerative clustering (bottom-up) [FrJa02,GMT07]
 - Single link, average link, complete link
 - Divisive clustering (top-down) [GMT07]
 - Furthest
 - LocalSearch [GMT07]
 - Place an object into a different cluster if objective function improved
 - Iterate the above until no improvements can be made
 - BestClustering [GMT07]
 - Select the clustering that maximize (minimize) the objective function
 - Graph partitioning [StGh03]
 - Nonnegative matrix factorization [LDJ07,LiDi08]

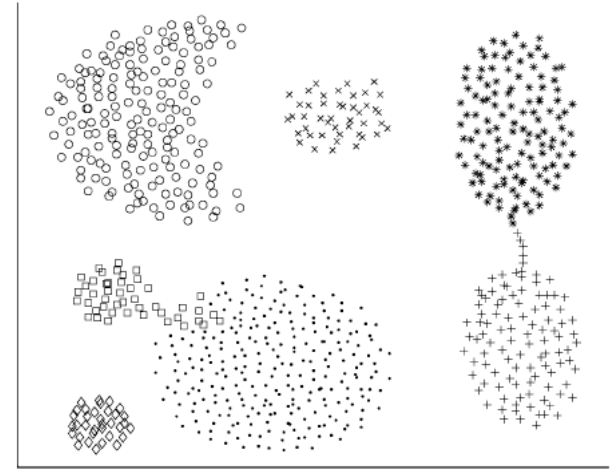
Single linkage



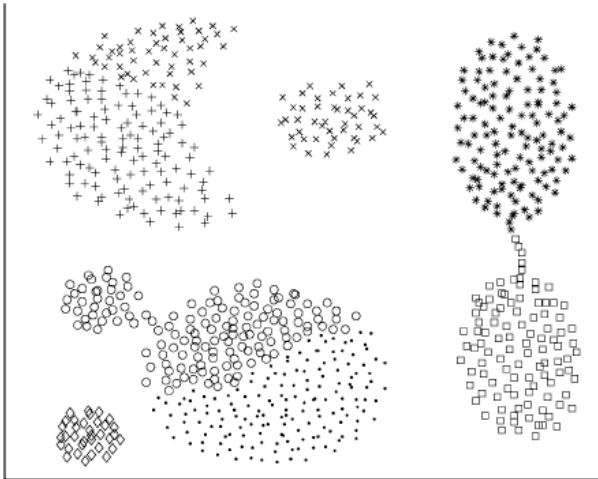
Complete linkage



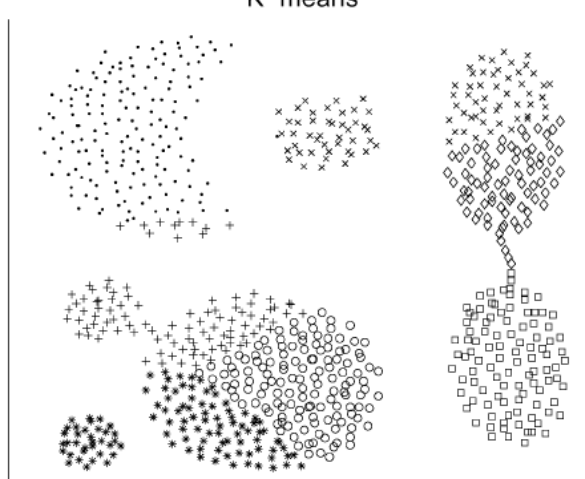
Average linkage



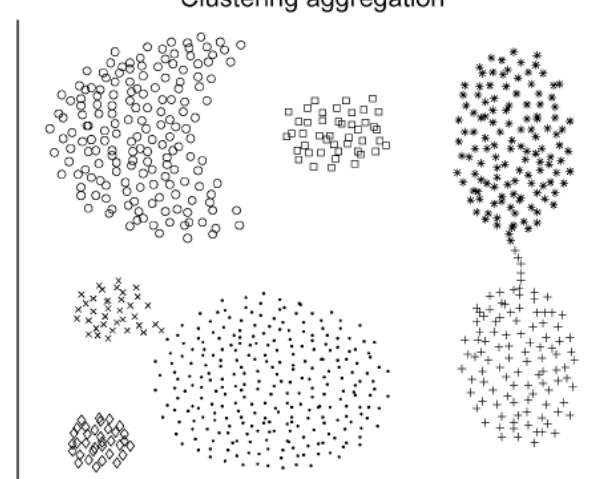
Ward's clustering



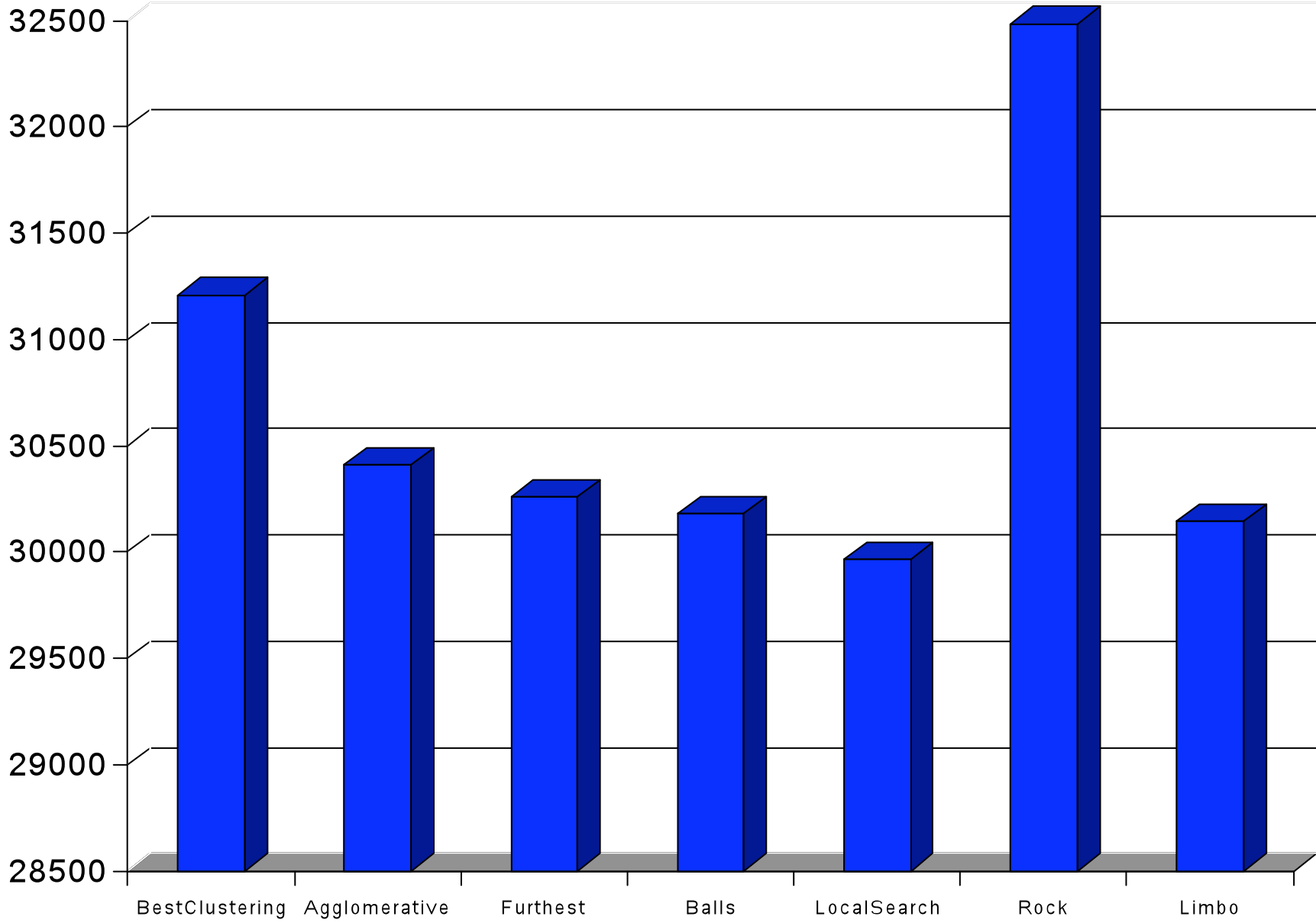
K-means



Clustering aggregation



Overall Distance on Votes data set

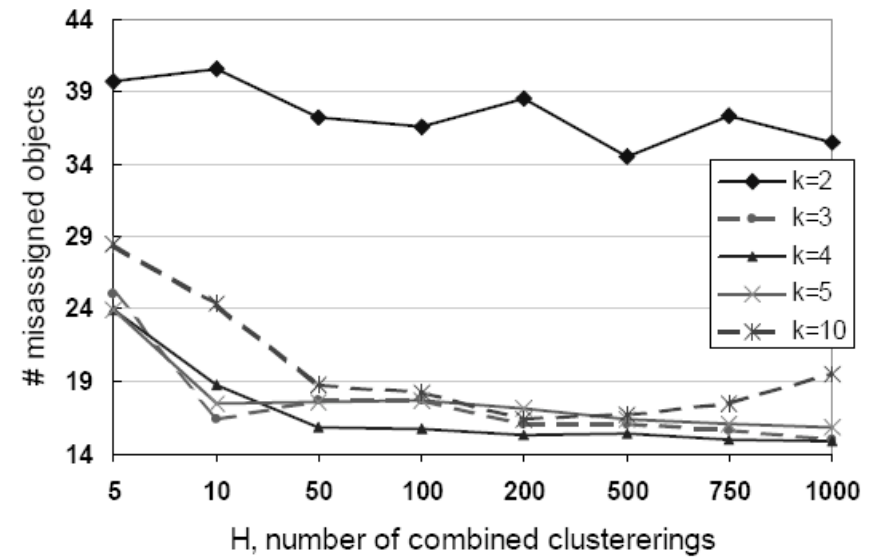
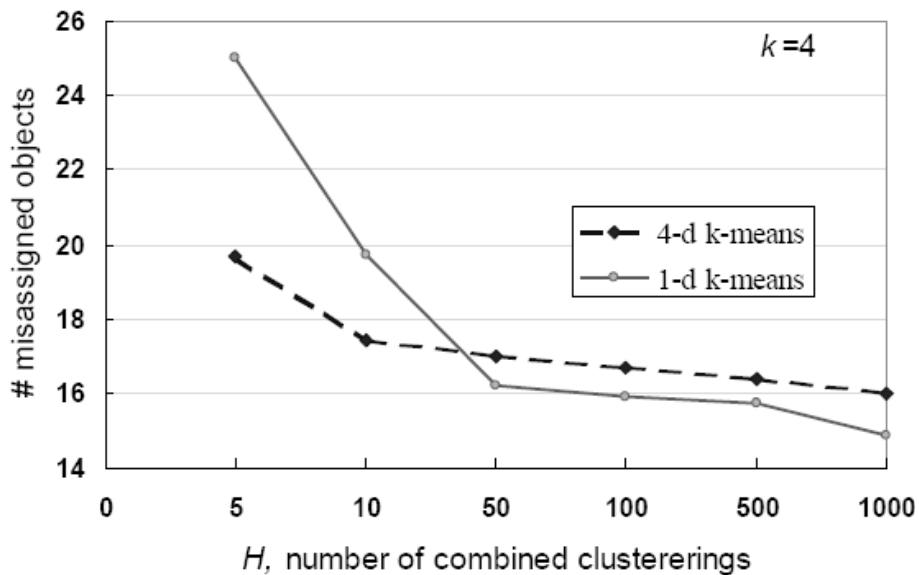
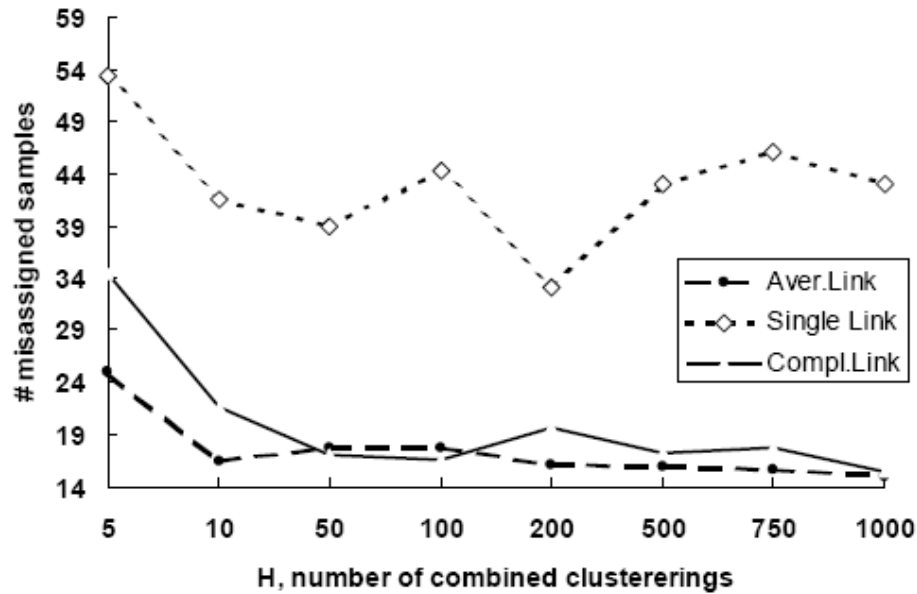


Iris data set

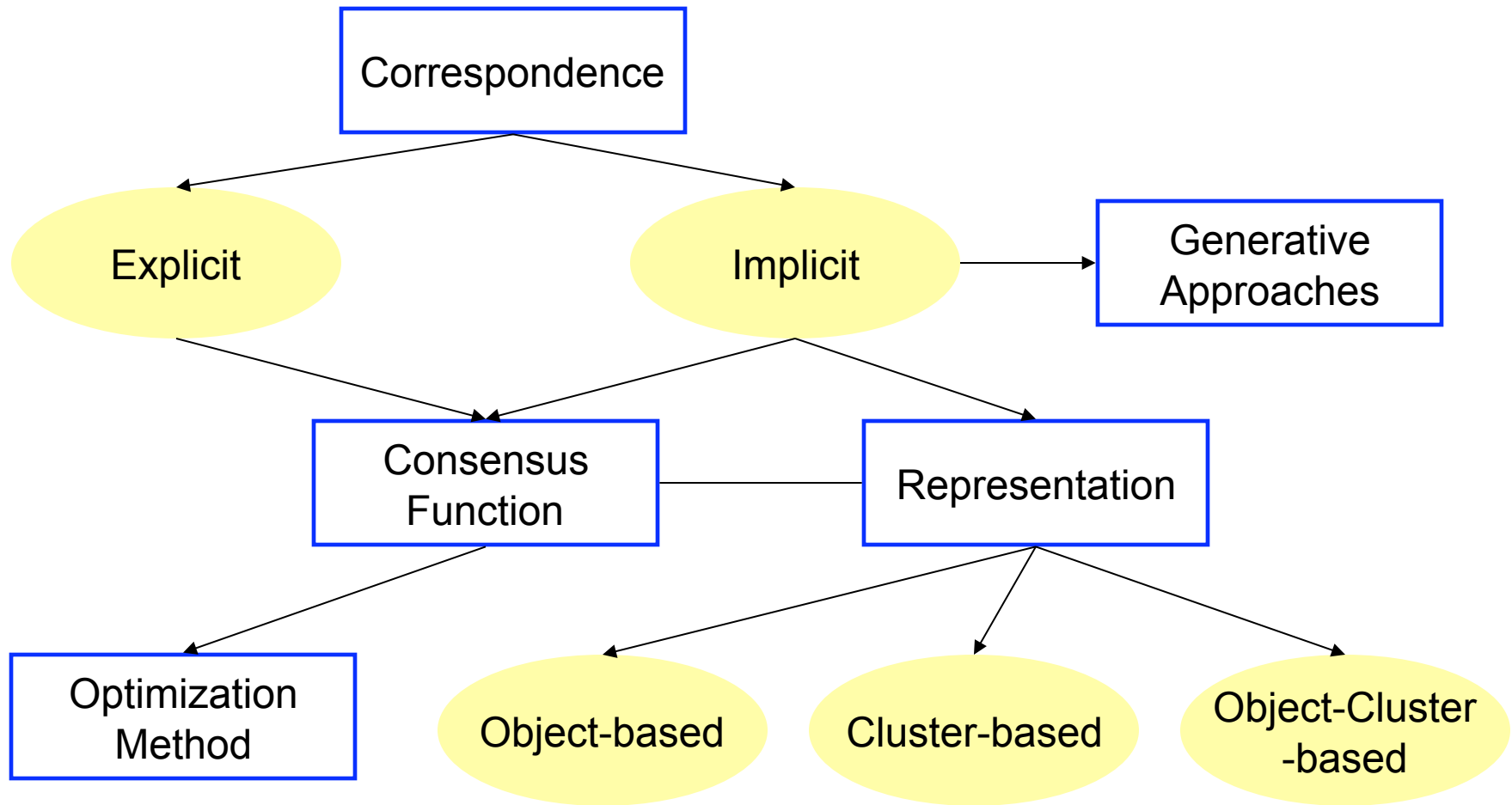
Algorithm: agglomerative clustering

k: number of clusters

H: number of clusterings



- How to combine the models?



Cluster-based Methods

- **Clustering clusters**
 - Regard each cluster from a base model as a record
 - Similarity is defined as the percentage of shared common objects
 - eg. Jaccard measure
 - Conduct clustering on these clusters
 - Assign an object to its most associated consensus cluster

Meta-Clustering Algorithm (MCLA)*

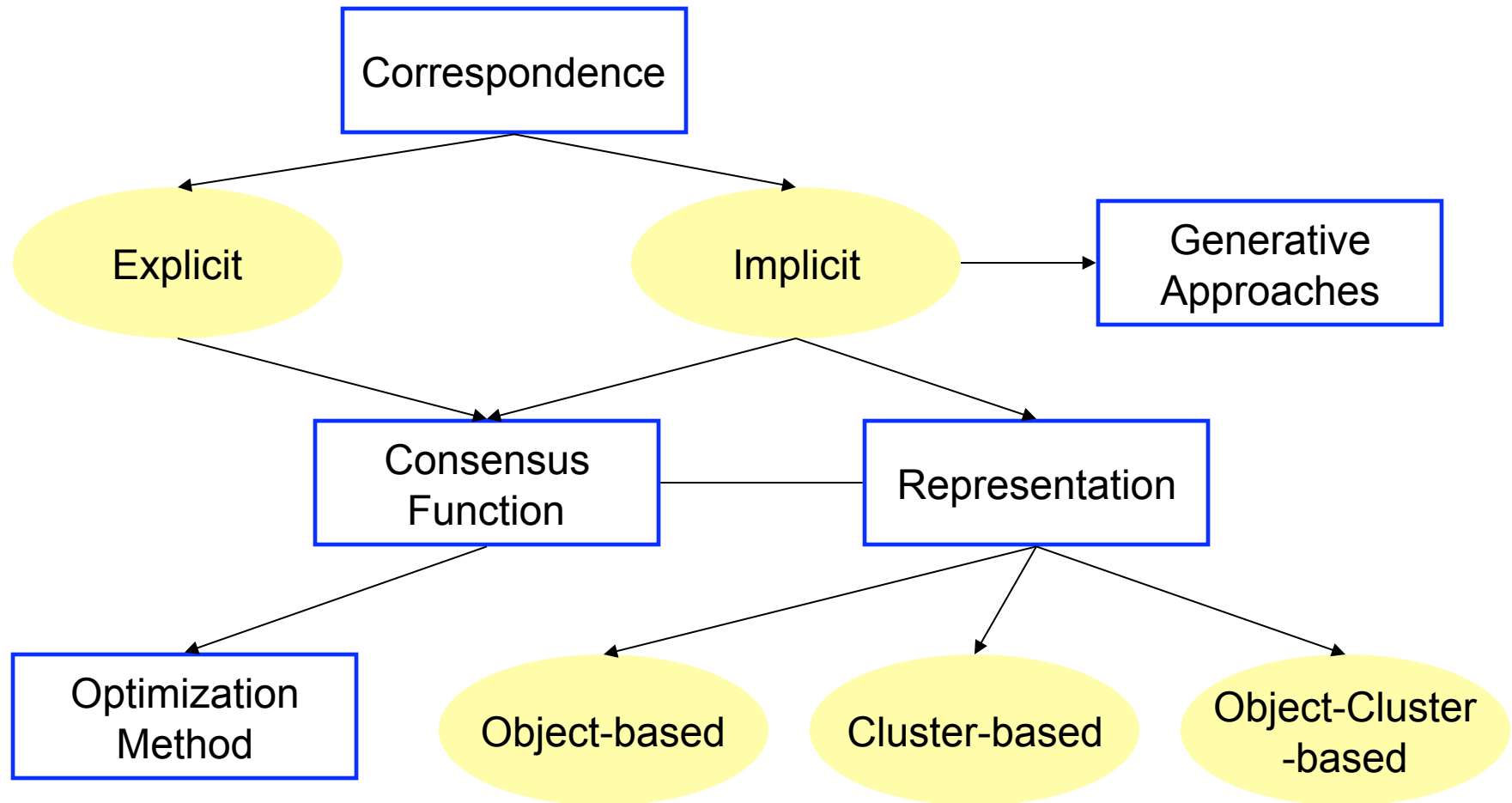
↓

	C_1	C_2	C_3	C	M_1	M_2	M_3
v_1	1	1	1	1	3	0	0
v_2	1	2	2	2	1	2	0
v_3	2	1	1	1	2	1	0
v_4	2	2	2	2	0	3	0
v_5	3	3	3	3	0	0	3
v_6	3	4	3	3	0	0	3

The diagram shows three meta-clusters, M_1 , M_2 , and M_3 , each enclosed in a pink circle. M_1 contains nodes q_1 , q_2 , and q_3 . M_2 contains nodes q_4 , q_5 , and q_6 . M_3 contains nodes q_7 , q_8 , q_9 , and q_{10} .

*[StGh03]

- How to combine the models?



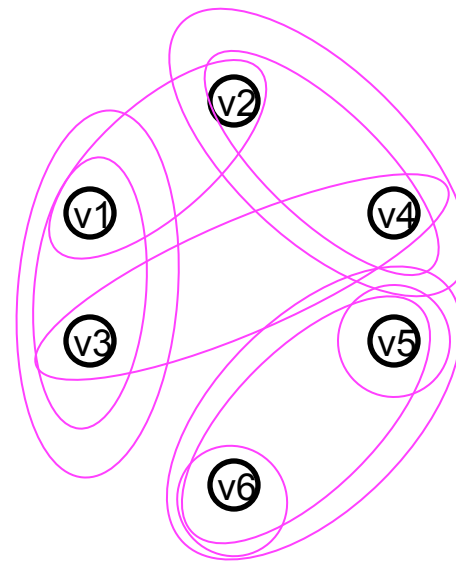
HyperGraph-Partitioning Algorithm (HGPA)*

- **Hypergraph representation and clustering**
 - Each node denotes an object
 - A hyperedge is a generalization of an edge in that it can connect any number of nodes
 - For objects that are put into the same cluster by a clustering algorithm, draw a hyperedge connecting them
 - Partition the hypergraph by minimizing the number of cut hyperedges
 - Each component forms a consensus cluster

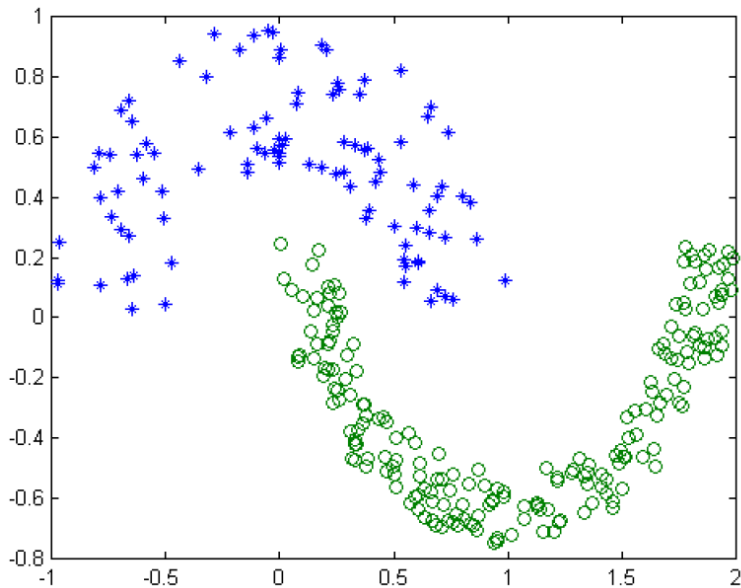
*[StGh03]

HyperGraph-Partitioning Algorithm (HGPA)

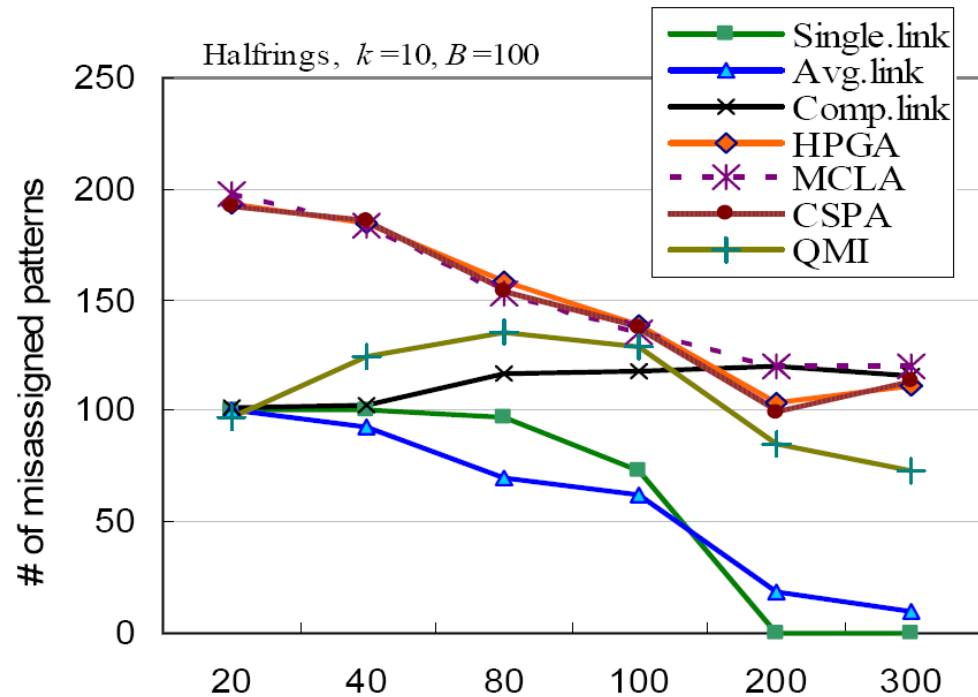
	\mathcal{C}_1	\mathcal{C}_2	\mathcal{C}_3	\mathcal{C}
v_1	1	1	1	1
v_2	1	2	2	2
v_3	2	1	1	1
v_4	2	2	2	2
v_5	3	3	3	3
v_6	3	4	3	3



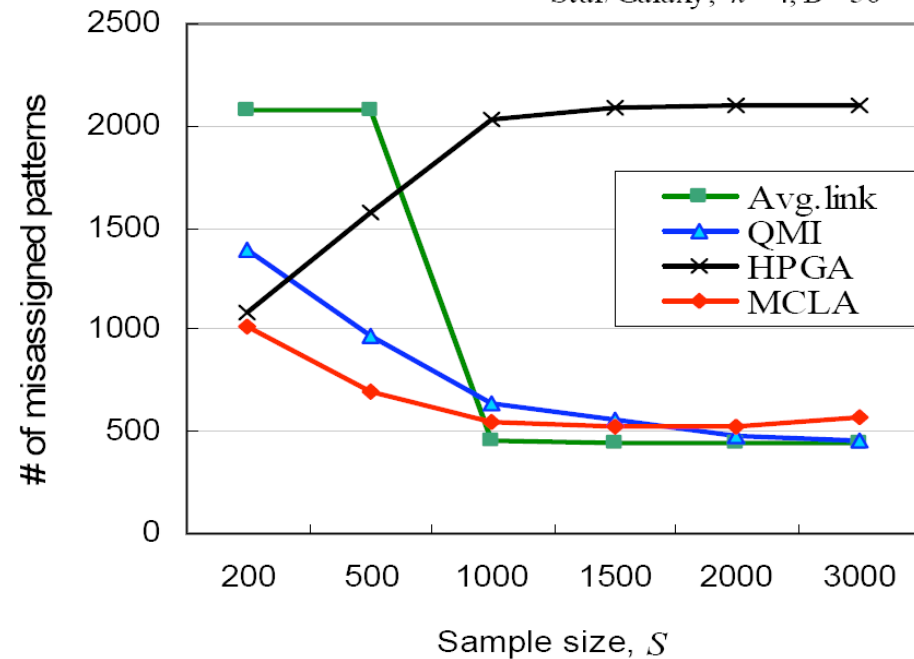
Hypergraph representation– a circle denotes a hyperedge



Halfrings dataset



Star/Galaxy, $k=4, B=50$



Object-based:

Agglomerative: Single link, Avg. link, Comp. link

METIS: CSPA

Quadratic mutual information: QMI

Cluster-based: MCLA

Object-cluster-based: HPGA

[MTP04]

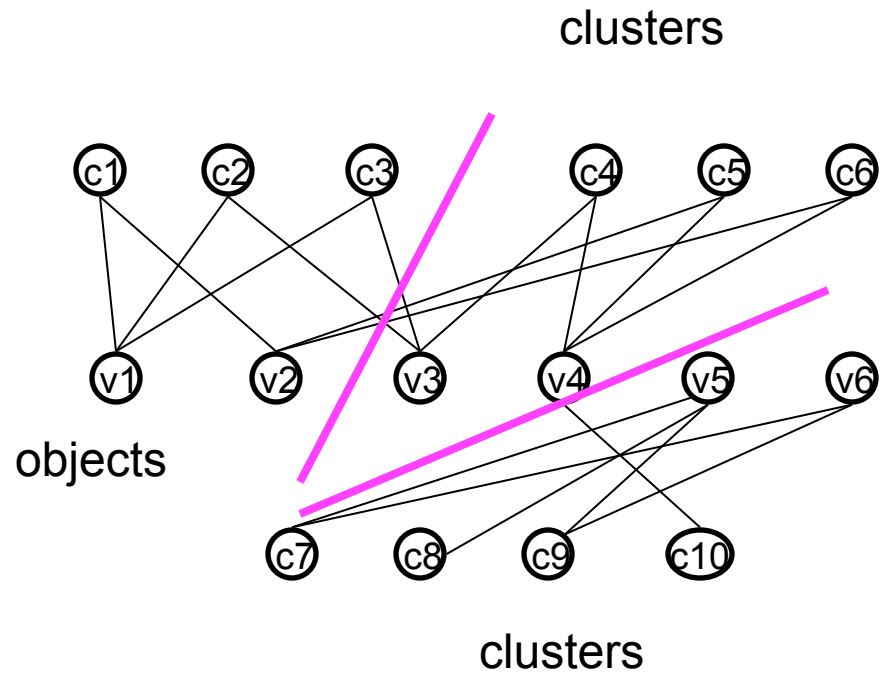
Bipartite Graph Partitioning*

- Hybrid Bipartite Graph Formulation
 - Summarize base model output in a bipartite graph
 - Lossless summarization—base model output can be reconstructed from the bipartite graph
 - Use spectral clustering algorithm to partition the bipartite graph
 - Time complexity $O(nkr)$ —due to the special structure of the bipartite graph
 - Each component represents a consensus cluster

*[FeBr04]

Bipartite Graph Partitioning

	\mathcal{C}_1	\mathcal{C}_2	\mathcal{C}_3	\mathcal{C}
v_1	1	1	1	1
v_2	1	2	2	2
v_3	2	1	1	1
v_4	2	2	2	2
v_5	3	3	3	3
v_6	3	4	3	3



Integer Programming*

- **Three-dimensional representation**

- Object l , cluster i , clustering algorithm j

$A_{lij} = 1$ If object l is assigned to cluster i by algorithm j

$X_{li'} = 1$ If object l is assigned to cluster i' by the consensus output

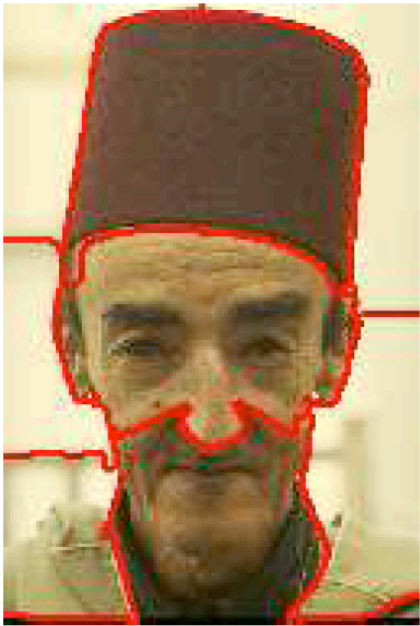
$S_{ij i'} = 1$ For algorithm j , cluster i has the largest overlapping with cluster i' in the consensus output

- **Objective function**

- Median partition

$$\min \sum_{i,j,i'} \left| S_{ij i'} - \frac{\sum_{l=1}^n A_{lij} X_{li'}}{\sum_{l=1}^n X_{li'}} \right|$$

*[SMP+07]



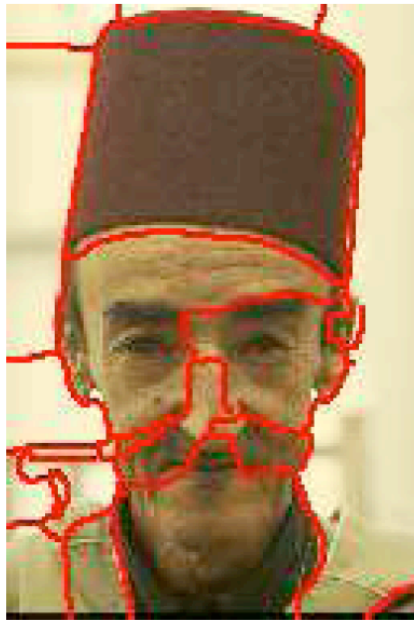
(a)



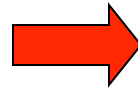
(b)



(c)

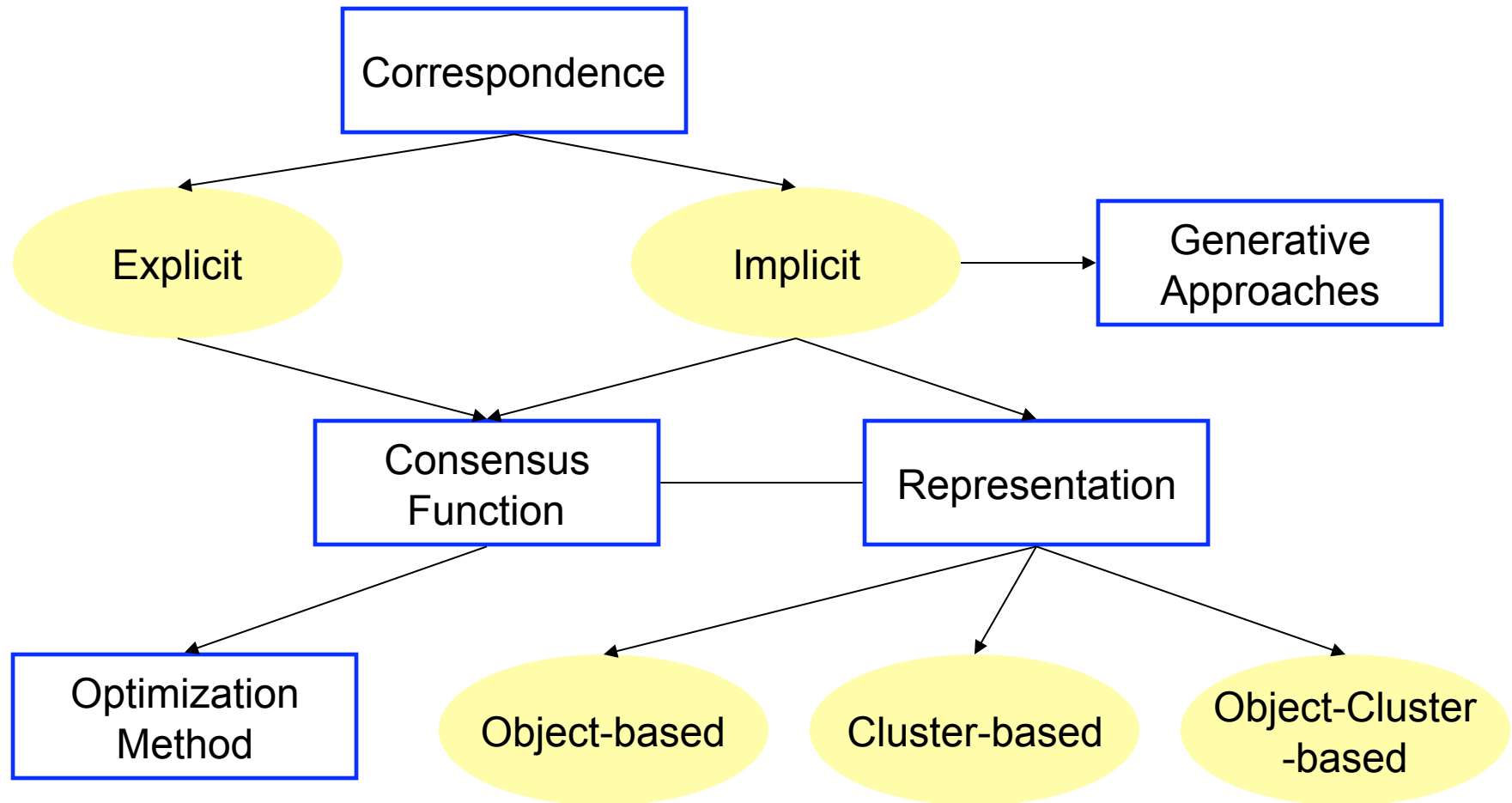


(d)



ensemble

- How to combine the models?



A Mixture Model of Consensus*

- **Probability-based**
 - Assume output comes from a mixture of models
 - Use EM algorithm to learn the model
- **Generative model**
 - The clustering solutions for each object are represented as nominal features-- v_i
 - v_i is described by a mixture of k components, each component follows a multinomial distribution
 - Each component is characterized by distribution parameters θ_j

*[PTJ05]

EM Method

- Maximize log likelihood

$$\sum_{i=1}^n \log \left(\sum_{j=1}^k \alpha_j P(v_i | \theta_j) \right)$$

- Hidden variables
 - z_i denotes which consensus cluster the object belongs to
- EM procedure
 - E-step: compute expectation of z_i
 - M-step: update model parameters to maximize likelihood

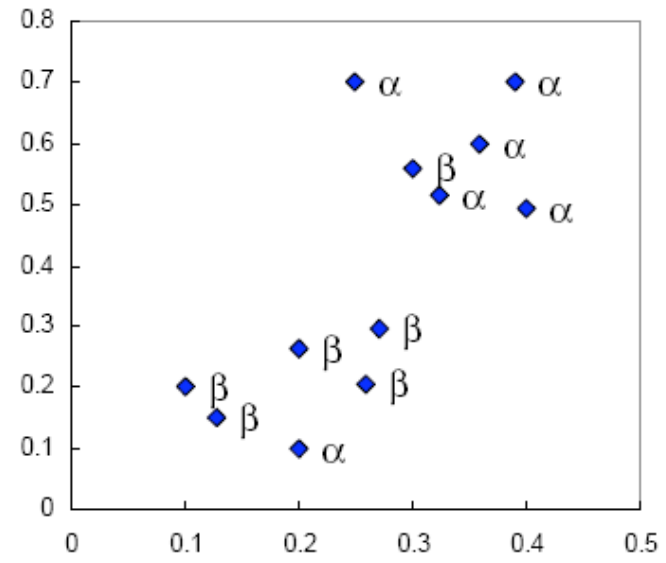
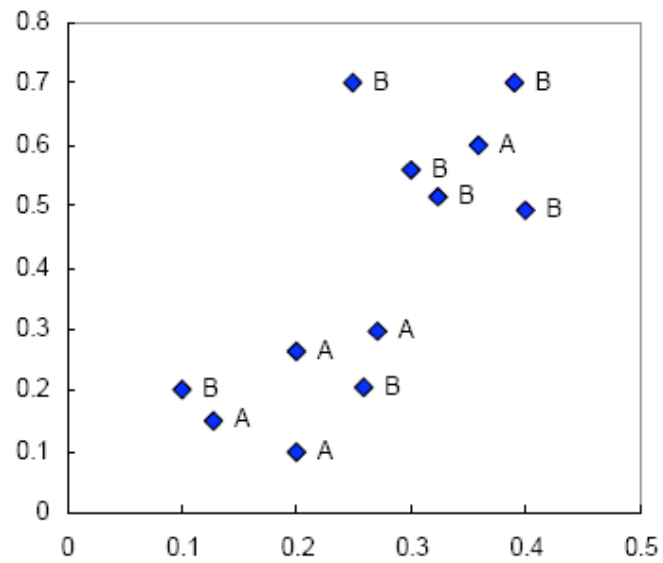
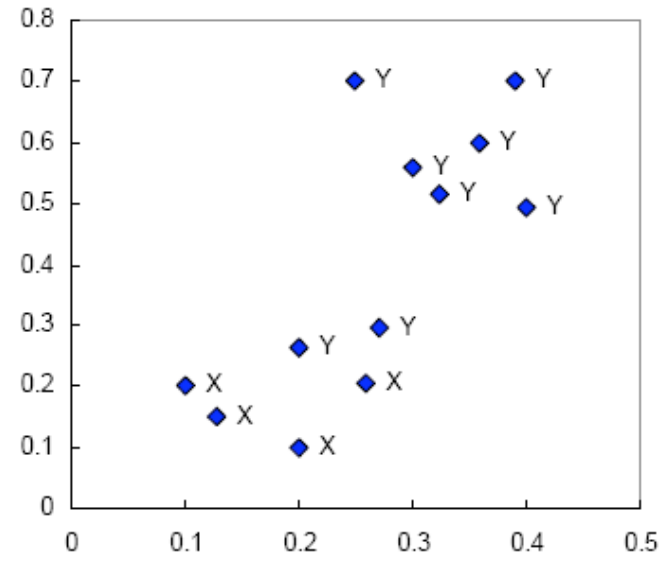
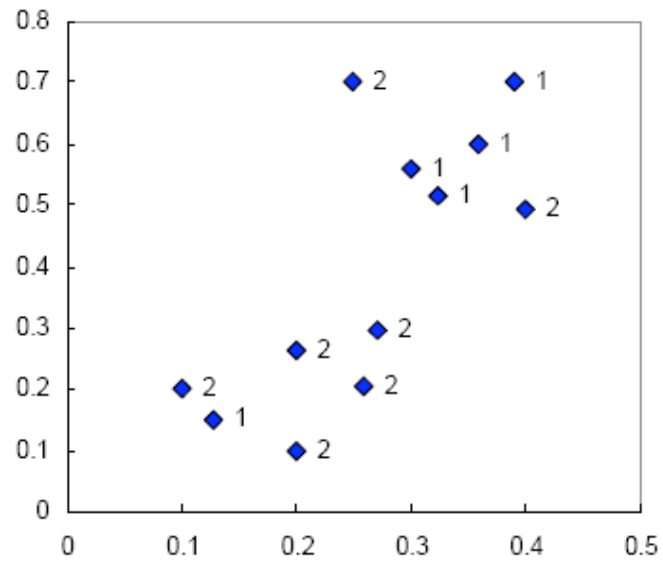
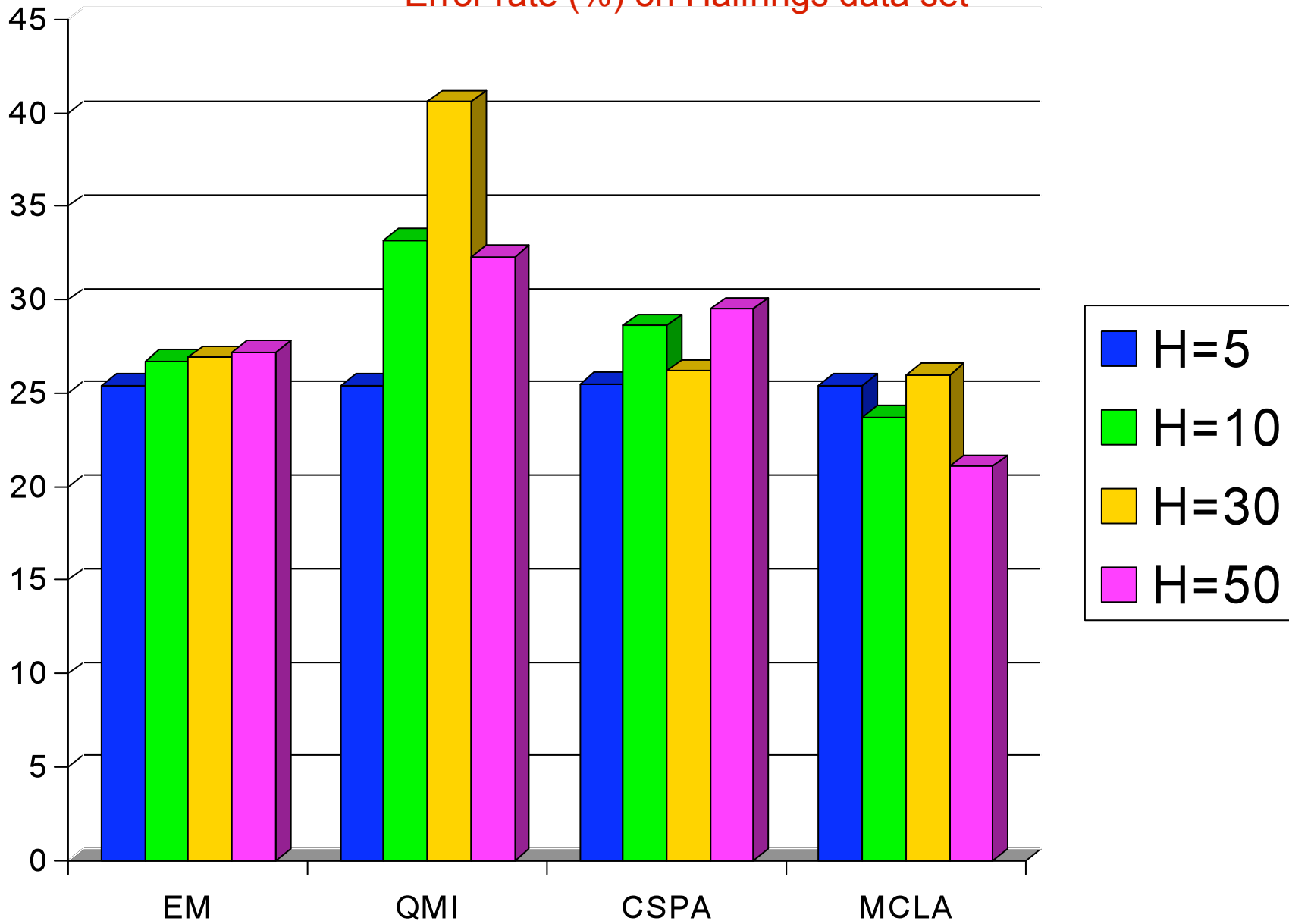


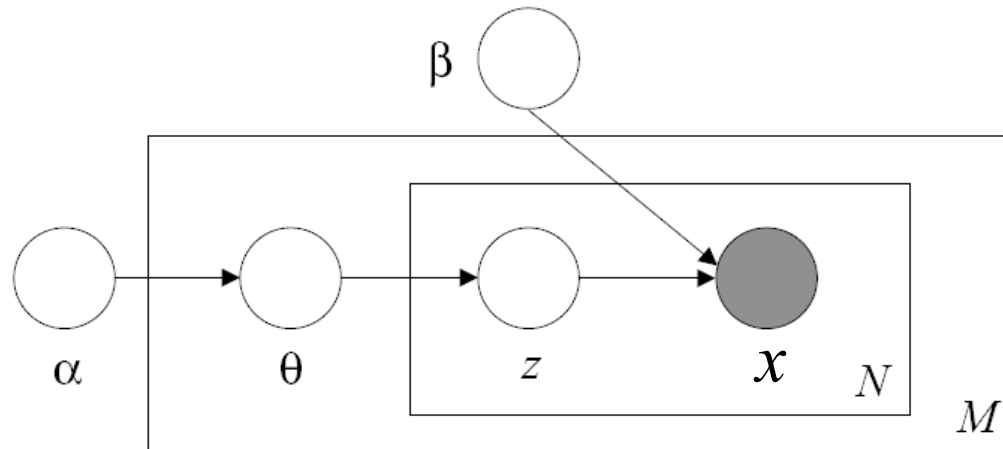
Table 1: Clustering ensemble and consensus solution

	π_1	π_2	π_3	π_4	$E[z_{i1}]$	$E[z_{i2}]$	Consensus
y_1	2	B	X	β	0.999	0.001	1
y_2	2	A	X	α	0.997	0.003	1
y_3	2	A	Y	β	0.943	0.057	1
y_4	2	B	X	β	0.999	0.001	1
y_5	1	A	X	β	0.999	0.001	1
y_6	2	A	Y	β	0.943	0.057	1
y_7	2	B	Y	α	0.124	0.876	2
y_8	1	B	Y	α	0.019	0.981	2
y_9	1	B	Y	β	0.260	0.740	2
y_{10}	1	A	Y	α	0.115	0.885	2
y_{11}	2	B	Y	α	0.124	0.876	2
y_{12}	1	B	Y	α	0.019	0.981	2

Error rate (%) on Halfrings data set



Bayesian Clustering Ensemble*



Consensus cluster—topic

1. Choose $\theta_i \sim \text{Dirichlet}(\alpha)$.

Cluster in base clustering—word

2. For the j^{th} base clustering:

Object—document

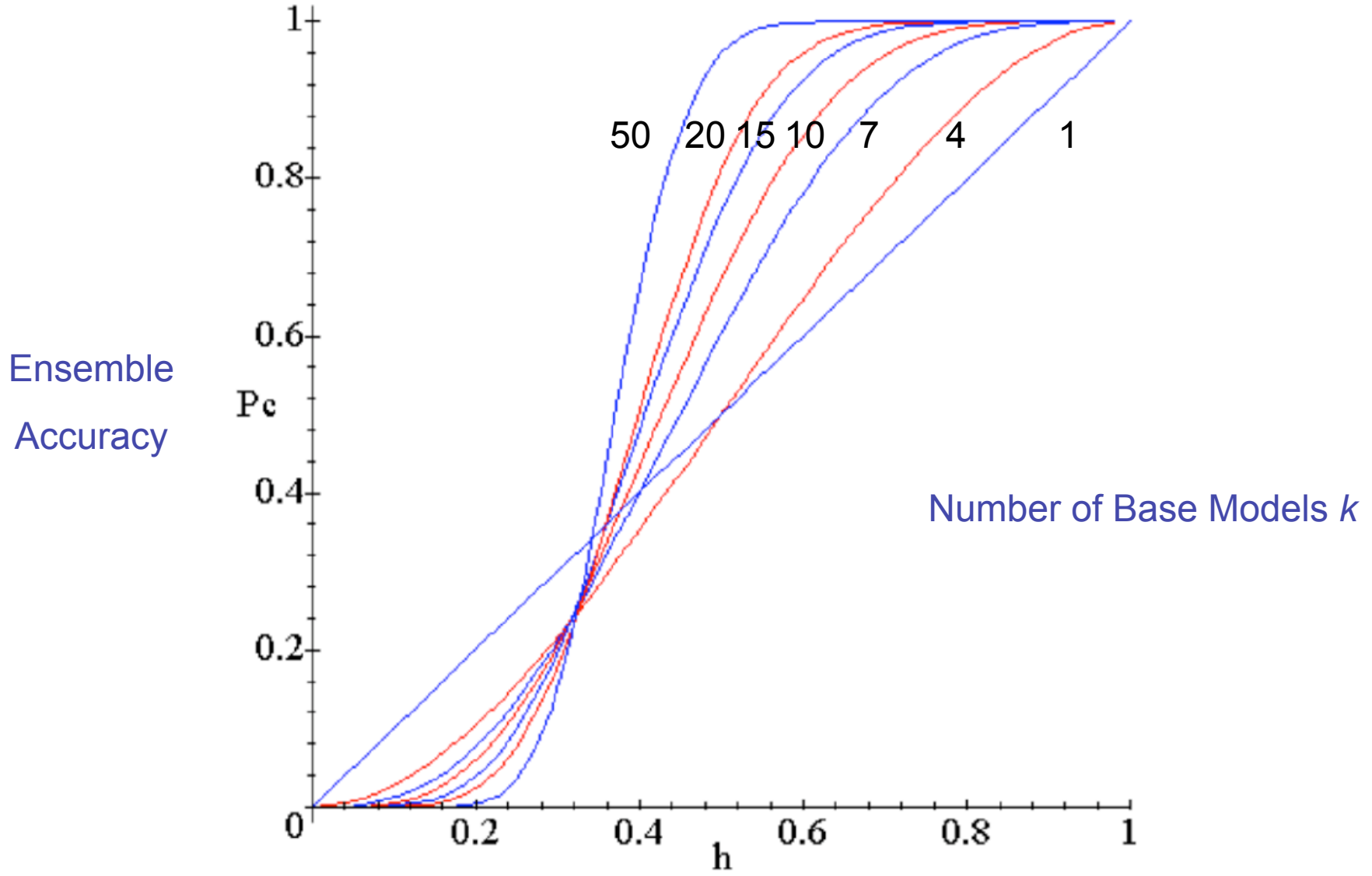
(a) Choose a component $z_{ij} = h \sim \text{discrete}(\theta_i)$;

(b) Choose the base clustering result $x_{ij} \sim \text{discrete}(\beta_{hj})$.

*[WSB09]

Other Research Problems

- **Consensus Clustering Theory**
 - Consensus clustering converges to true clustering as the number of base clustering models increases [TLJ+04]
 - Error incurred by approximation has a lower bound [GMT07,GoFi08]
- **Base model selection**
 - Ensemble selection [FeLi08]
 - Moderate diversity [HKT06,KuWh03]
- **Combining soft clustering**
 - Extend ensemble methods developed for hard clustering [PuGh08]



[TLJ+04]

Base Model Accuracy

Summary of Unsupervised Ensemble

- **Difference from supervised ensemble**
 - The success of clustering ensemble approaches is shown empirically
 - There exist label correspondence problems
- **Characteristics**
 - Experimental results demonstrate that cluster ensembles are better than single models!
 - There is no single, universally successful, cluster ensemble method

Outline

- An overview of ensemble methods
 - Motivations
 - Tutorial overview
- Supervised ensemble
- Unsupervised ensemble
- Semi-supervised ensemble
 - Multi-view learning
 - Consensus maximization among supervised and unsupervised models
- Applications
 - Transfer learning, stream classification, anomaly detection

Multiple Source Classification

flickr

Home The Tour Sign Up Explore

Is there anybody out there?



Actually i'm not a big fan of beach.
It was a sunday afternoon and the summer was going down. I remember i was really excited cause there wasn't anybody over there. Only me and a friend of mine in that desolate beach.
We've smoked a lot and the wind was gentle on our body.
Well, after that day my opinion about beaches is changed.

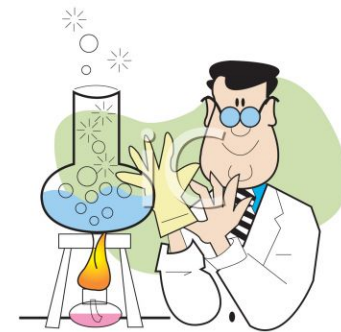


Image Categorization

images, descriptions,
notes, comments,
albums, tags.....

Like? Dislike?

movie genres, cast,
director, plots.....
users viewing history,
movie ratings...

Research Area

publication and co-
authorship network,
published papers,
.....

Multi-view Learning

- **Problem**

- The same set of objects can be described in multiple different views
- Features are naturally separated into K sets:

$$X = (X^1, X^2, \dots, X^K)$$

- Both labeled and unlabeled data are available
- Learning on multiple views:
 - Search for labeling on the unlabeled set and target functions on X : $\{f_1, f_2, \dots, f_k\}$ so that the target functions agree on labeling of unlabeled data

Learning from Two Views

- **Input**

- Features can be split into two sets: $X = X_1 \times X_2$
- The two views are redundant but not completely correlated
- Few labeled examples and relatively large amounts of unlabeled examples are available from the two views

- **Conditions**

- Compatible --- all examples are labeled identically by the target concepts in each view
- Uncorrelated --- given the label of any example, its descriptions in each view are independent

How It Works?

- **Conditions**

- Compatible --- Reduce the search space to where the two classifiers agree on unlabeled data
- Uncorrelated --- If two classifiers always make the same predictions on the unlabeled data, we cannot benefit much from multi-view learning

- **Algorithms**

- Searching for compatible hypotheses
- Canonical correlation analysis
- Co-regularization

- **Theory**

- [DLM01, BBY04, Leskes05]

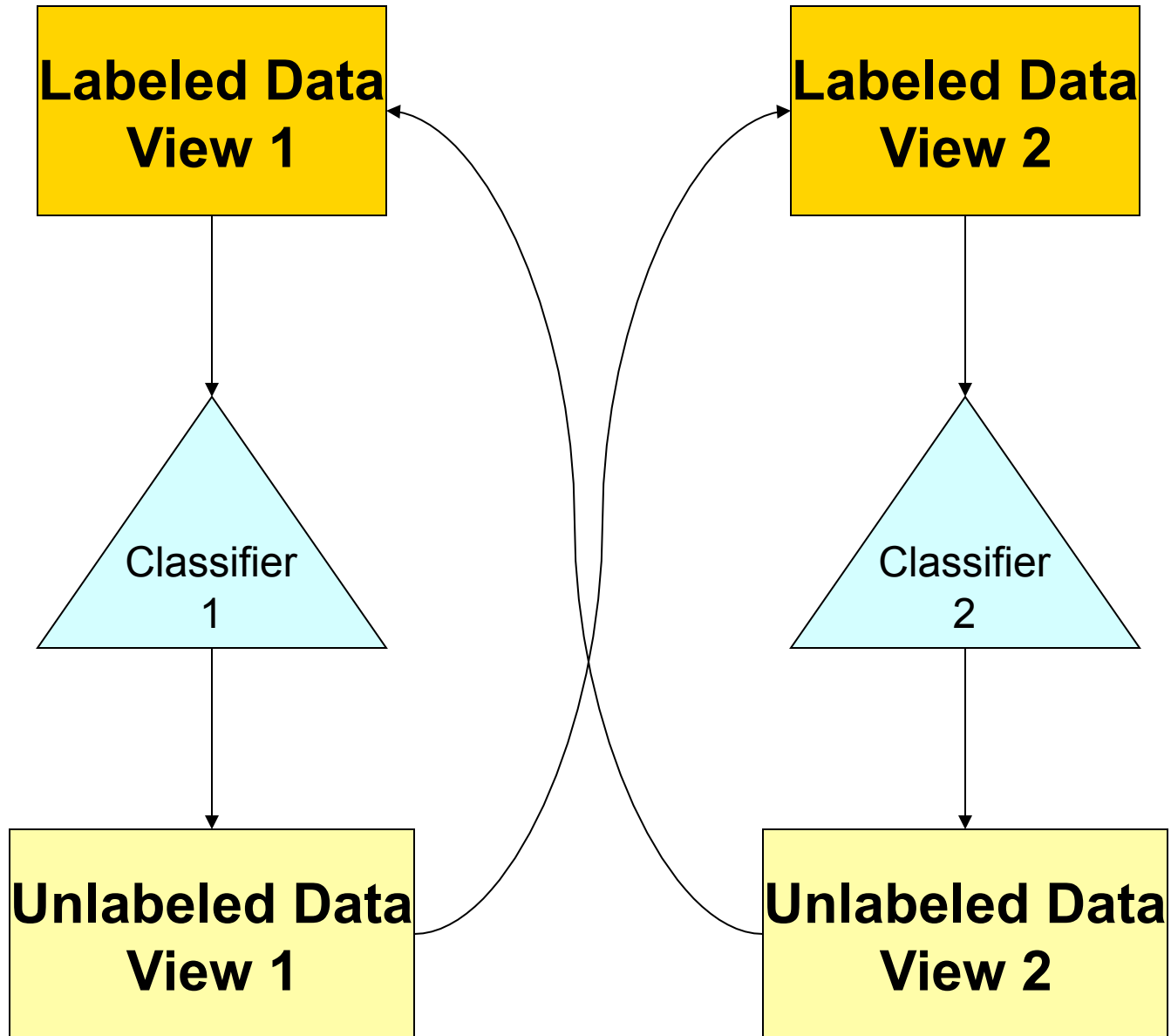
Searching for Compatible Hypotheses

- **Intuitions**

- Two individual classifiers are learnt from the labeled examples of the two views
- The two classifiers' predictions on unlabeled examples are used to enlarge the size of training set
- The algorithm searches for “compatible” target functions

- **Algorithms**

- Co-training [BIMi98]
- Co-EM [NiGh00]
- Variants of Co-training [GoZh00]



Co-Training*

Given:

- a set L of labeled training examples
- a set U of unlabeled examples

Train two classifiers from two views

Create a pool U' of unlabeled examples
Loop for k iterations:
Select the top unlabeled examples with the most confident predictions from the other classifier

Use L to train a classifier h_1 that considers only the x_1 portion of x

Use L to train a classifier h_2 that considers only the x_2 portion of x

Allow h_1 to label p positive and n negative examples from U'

Allow h_2 to label p positive and n negative examples from U'

Add these self-labeled examples to L

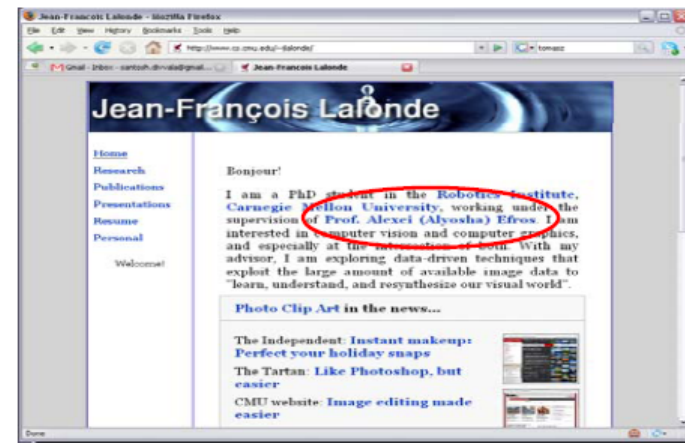
Randomly choose $2p + 2n$ examples from U to replenish U'

Add these self-labeled examples to the training set

Applications: Faculty Webpages Classification



View1: Page Text



View2: Hyperlink Text

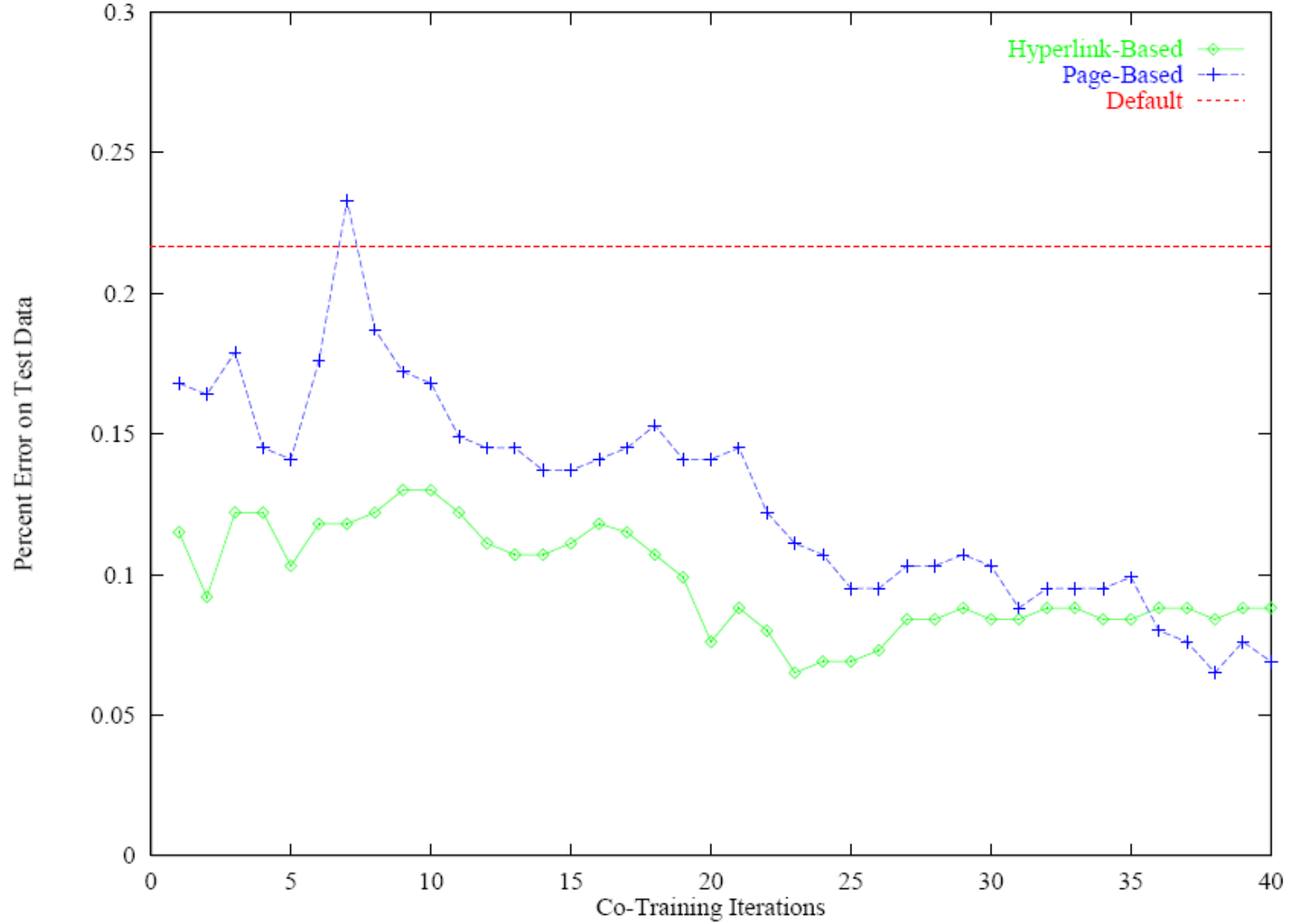


Figure 2: Error versus number of iterations for one run of co-training experiment.

Co-EM*

- **Algorithm**

- Labeled data set L , Unlabeled data set U , Let U_1 be empty, Let $U_2=U$
- Iterate the following
 - Train a classifier h_1 from the feature set X_1 of L and U_1
 - Probabilistically label all the unlabeled data in U_2 using h_1
 - Train a classifier h_2 from the feature set X_2 of L and U_2
 - Let $U_1=U$, probabilistically label all the unlabeled data in U_1 using h_2
- Combine h_1 and h_2

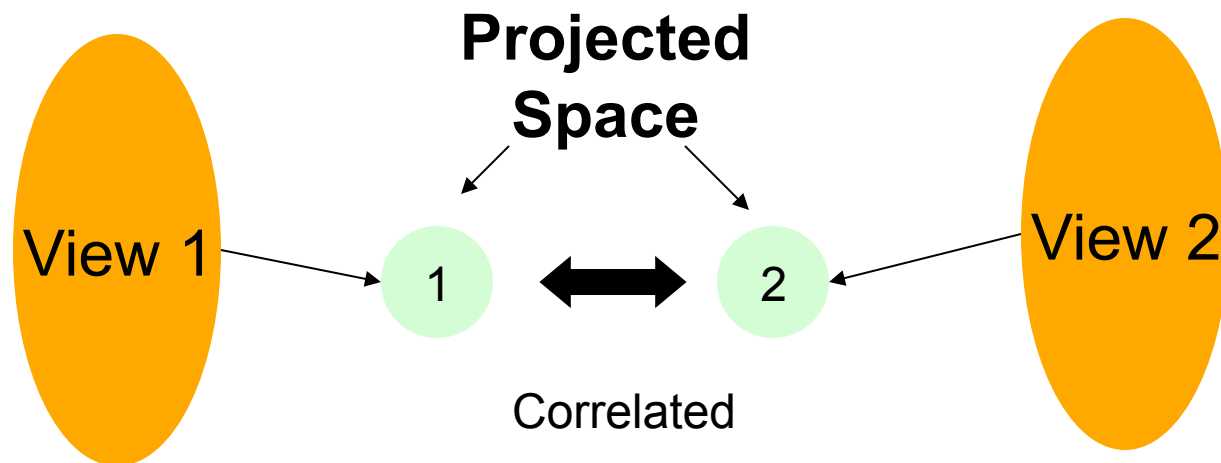
- **Co-EM vs. Co-Training**

- Labeling unlabeled data: soft vs. hard
- Selecting unlabeled data into training set: all vs. the top confident ones

Canonical Correlation Analysis

- **Intuitions**

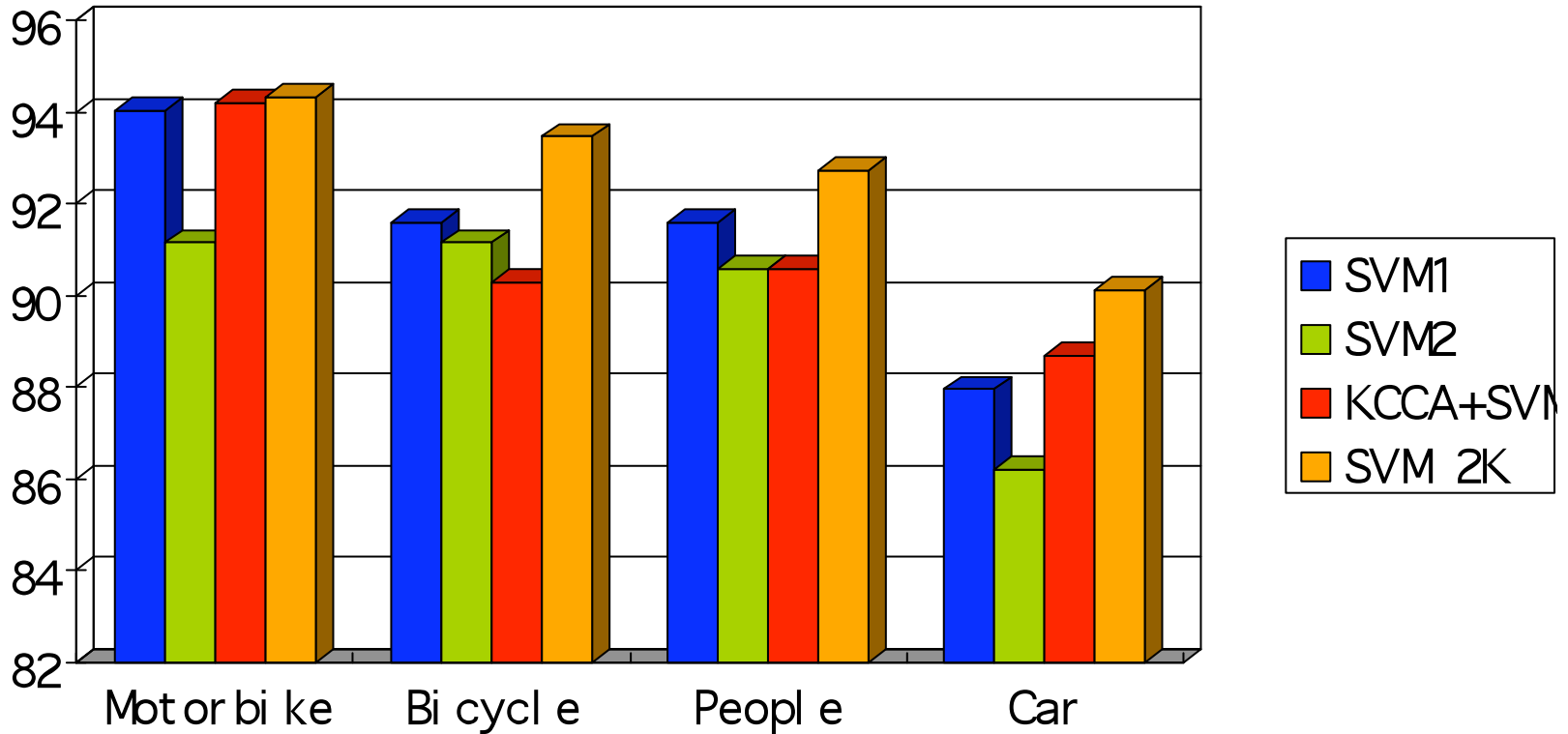
- Reduce the feature space to low-dimensional space containing discriminative information
- With compatible assumption, the discriminative information is contained in the directions that correlate between the two views
- The goal is to maximize the correlation between the data in the two projected spaces



Algorithms

- **Co-training in the reduced spaces [ZZY07]**
 - Project the data into the low-dimensional spaces by maximizing correlations between two views
 - Compute probability of unlabeled data belonging to positive or negative classes using the distance between unlabeled data and labeled data in the new feature spaces
 - Select the top-confident ones to enhance the training set and iterate
- **SVM+Canonical Correlation Analysis [FHM+05]**
 - First reduce dimensions, then train SVM classifiers
 - Combine the two steps together

Experimental Results



Accuracy Comparison on Image Data Set

Co-Regularization Framework

- **Intuitions**

- Train two classifiers from the two views simultaneously
- Add a regularization term to enforce that the two classifiers agree on the predictions of unlabeled data

$$\min R(f_1; L_1) + R(f_2; L_2) + R(f_1, f_2; U_1, U_2)$$

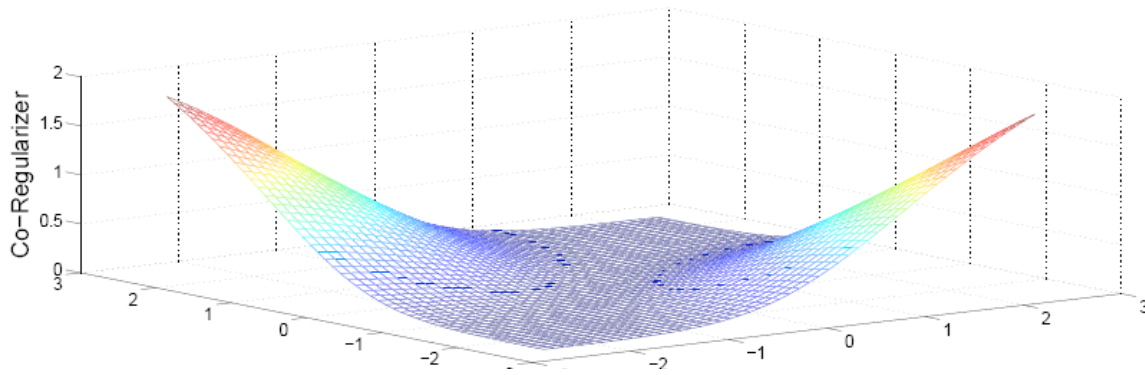
Risk of classifier 1 on view 1 of labeled data

Risk of classifier 2 on view 2 of labeled data

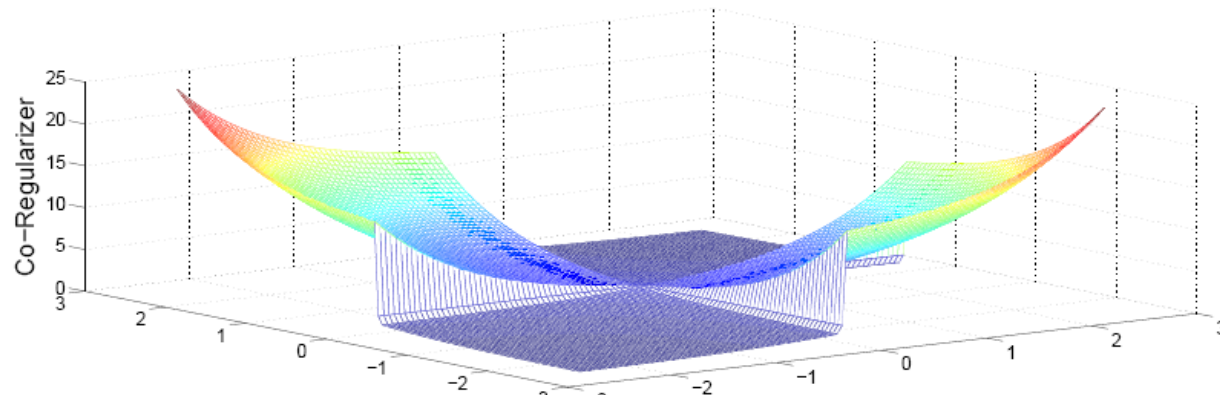
Disagreement between two classifiers on unlabeled data

- **Algorithms**

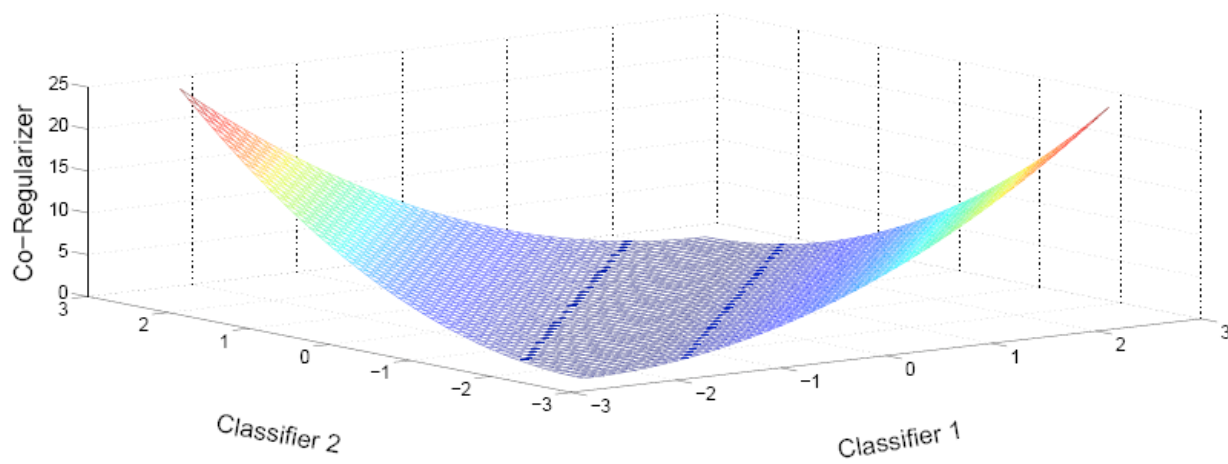
- Co-boosting [CoSi99]
- Co-regularized least squares and SVM [SNB05]
- Bhattacharyya distance regularization [GGB+08]



Bhattacharyya
distance



Exponential loss



Least square

Comparison of Loss Functions

- **Loss functions**

- Exponential: $\sum_{x \in U} \exp(-\tilde{y}_2 f_1(x)) + \exp(-\tilde{y}_1 f_2(x))$

- Least Square: $\sum_{x \in U} (f_1(x) - f_2(x))^2$

- Bhattacharyya distance: $E_U(B(p_1, p_2))$

$$B(p_1, p_2) = -\log \sum_y \sqrt{p_1(y)p_2(y)}$$

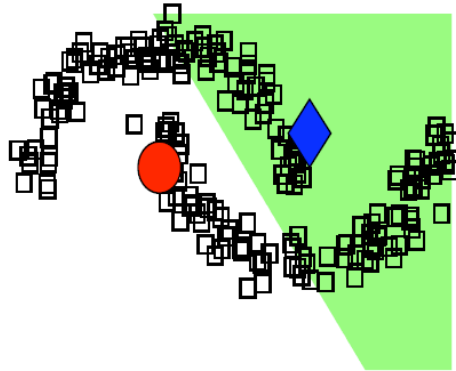
- **When two classifiers don't agree**

- Loss grows exponentially, quadratically, linearly

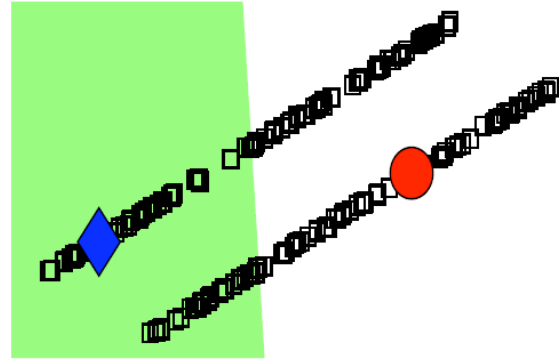
- **When two classifiers agree**

- Little penalty \longrightarrow Penalize the margin

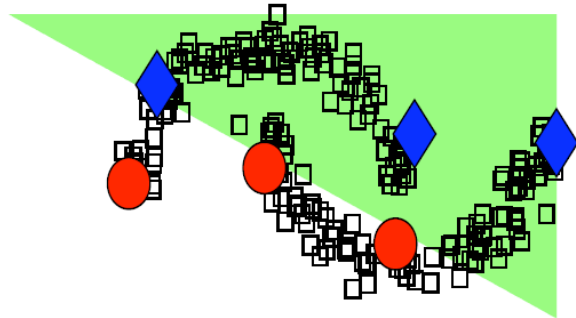
View 1: RLS (2 labeled examples)



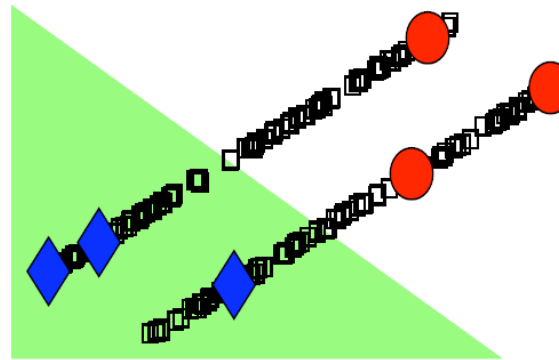
View 2: RLS (2 labeled examples)



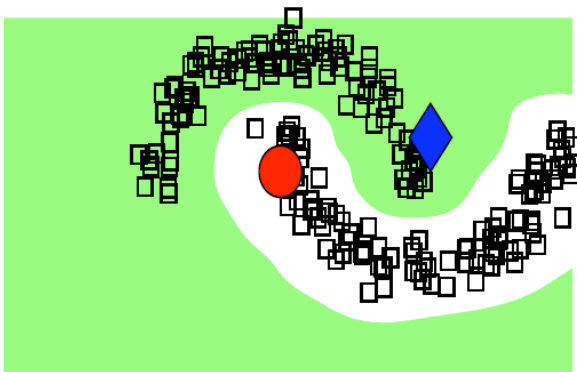
View 1: Co-trained RLS (1 step)



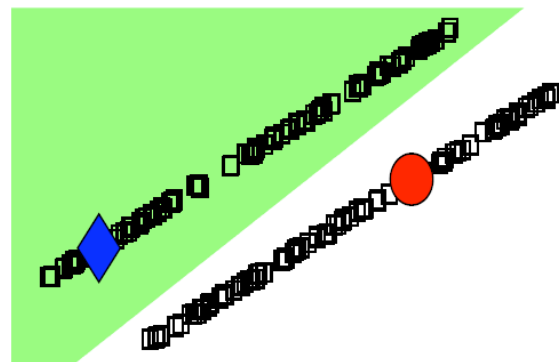
View 2: Co-trained RLS (1 step)



View 1: Co-RLS



View 2: Co-RLS



[SNB05]

Outline

- An overview of ensemble methods
 - Motivations
 - Tutorial overview
- Supervised ensemble
- Unsupervised ensemble
- Semi-supervised ensemble
 - Multi-view learning
 - Consensus maximization among supervised and unsupervised models
- Applications
 - Transfer learning, stream classification, anomaly detection

Consensus Maximization*

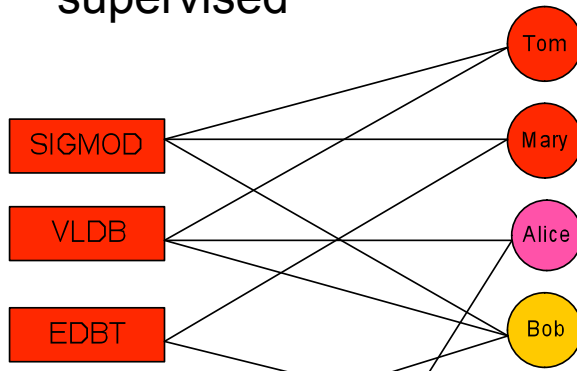
- **Goal**
 - Combine output of multiple supervised and unsupervised models on a set of objects
 - The predicted labels should agree with the base models as much as possible
- **Motivations**
 - Unsupervised models provide useful constraints for classification tasks
 - Model diversity improves prediction accuracy and robustness
 - Model combination at output level is needed due to privacy-preserving or incompatible formats

*[GLF+09]

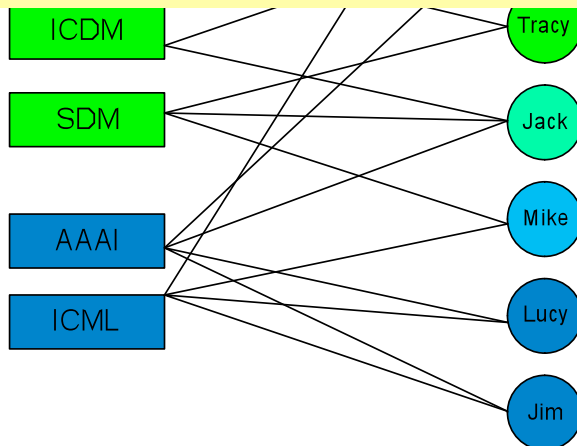
Model Combination helps!

Supervised or
unsupervised

supervised



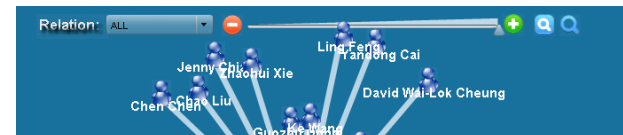
People may publish in relevant
but different areas



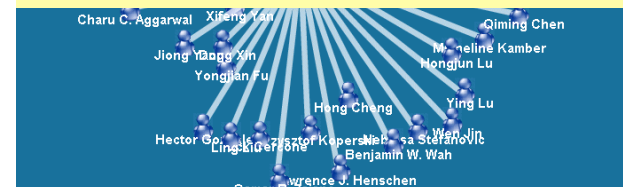
2009	
361	EE Jiawei Han, Xifeng Yan, Philip S. Yu Scalable OLAP and mining of information networks. <i>EDBT 2009</i> : 1159
360	EE Yizhou Sun, Jiawei Han, Pengzeng Zhao, Zhiun Yin, Hong Cheng, Tianyi Wu RankClus: integrating clustering with ranking for heterogeneous information network analysis. <i>EDBT 2009</i> : 565-576
359	EE Bhavani M. Thuraisingham, Latifur Khan, Murat Kantarcioglu, Sonar Chab, Jiawei Han, Sang Son Real-Time Knowledge Discovery and Dissemination for Intelligence Analysis. <i>HICSS 2009</i> : 1-12

Some areas share similar keywords

2008	
355	Deng Cai, Xiaofei He, Jiawei Han Sparse Projections over Graph. <i>AAAI 2008</i> : 610-615
354	EE Chen Chen, Candy Xiao Lin, Xifeng Yan, Jiawei Han On effective presentation of graph patterns: a structural representative approach. <i>CIKM 2008</i> : 299-308
353	EE Deng Cai, Qiaozhu Mei, Jiawei Han, Chengxiang Zhai Modeling hidden topics on document manifold. <i>CIKM 2008</i> : 911-920
352	EE Jiawei Han Data mining for image/video processing: a promising research frontier. <i>CVPR 2008</i> : 1-2

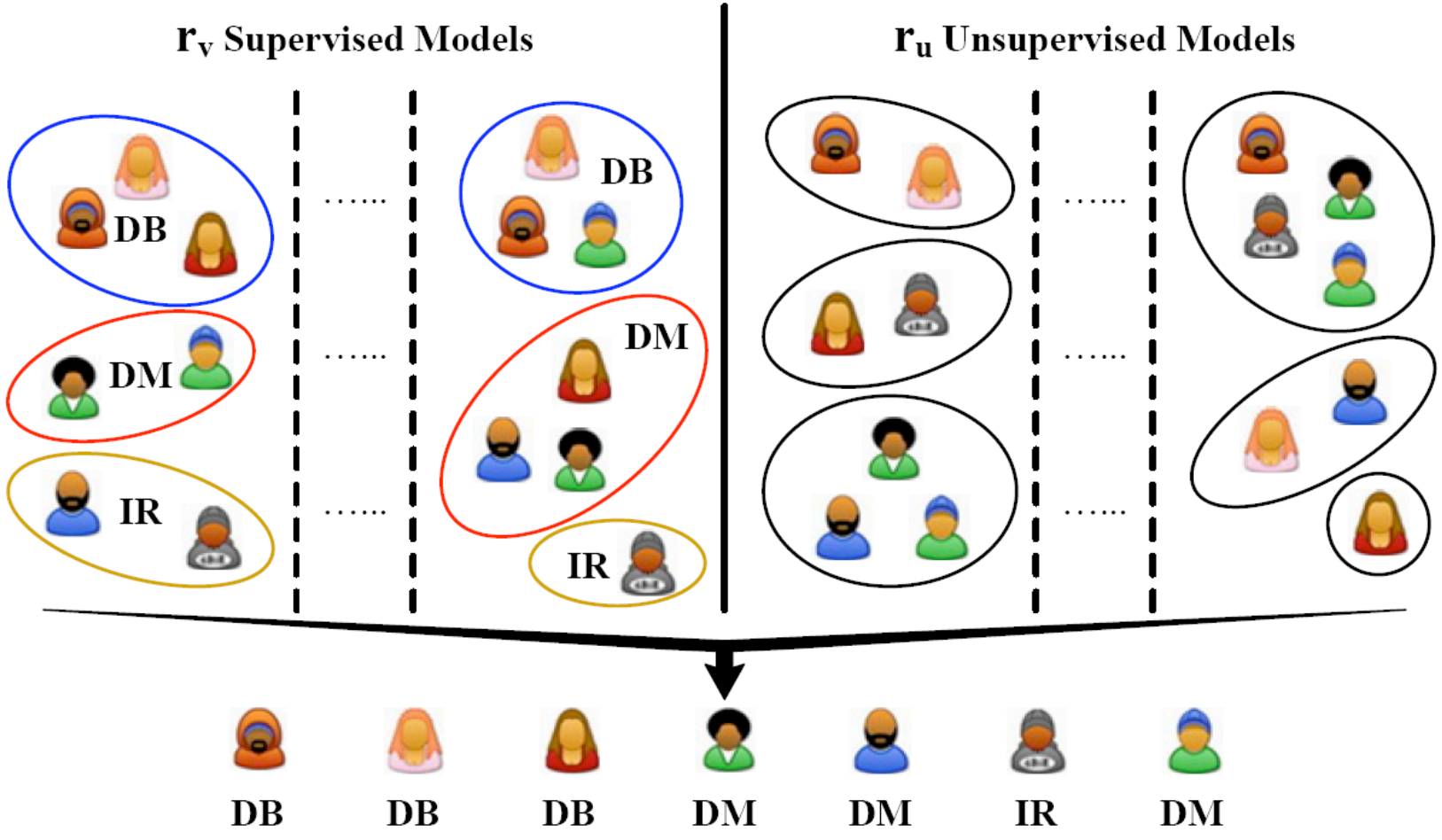


There may be cross-
discipline co-operations

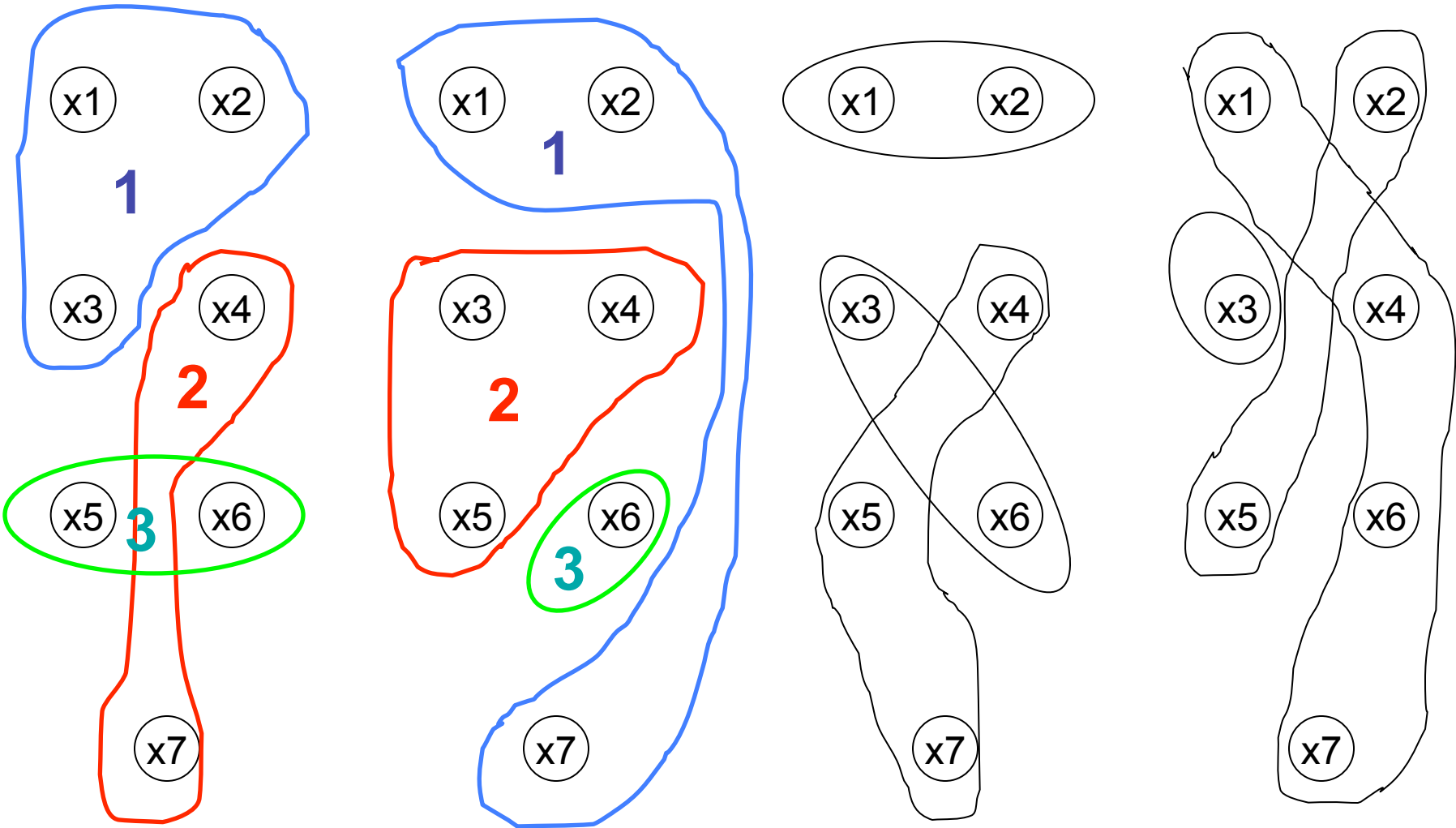


unsupervised

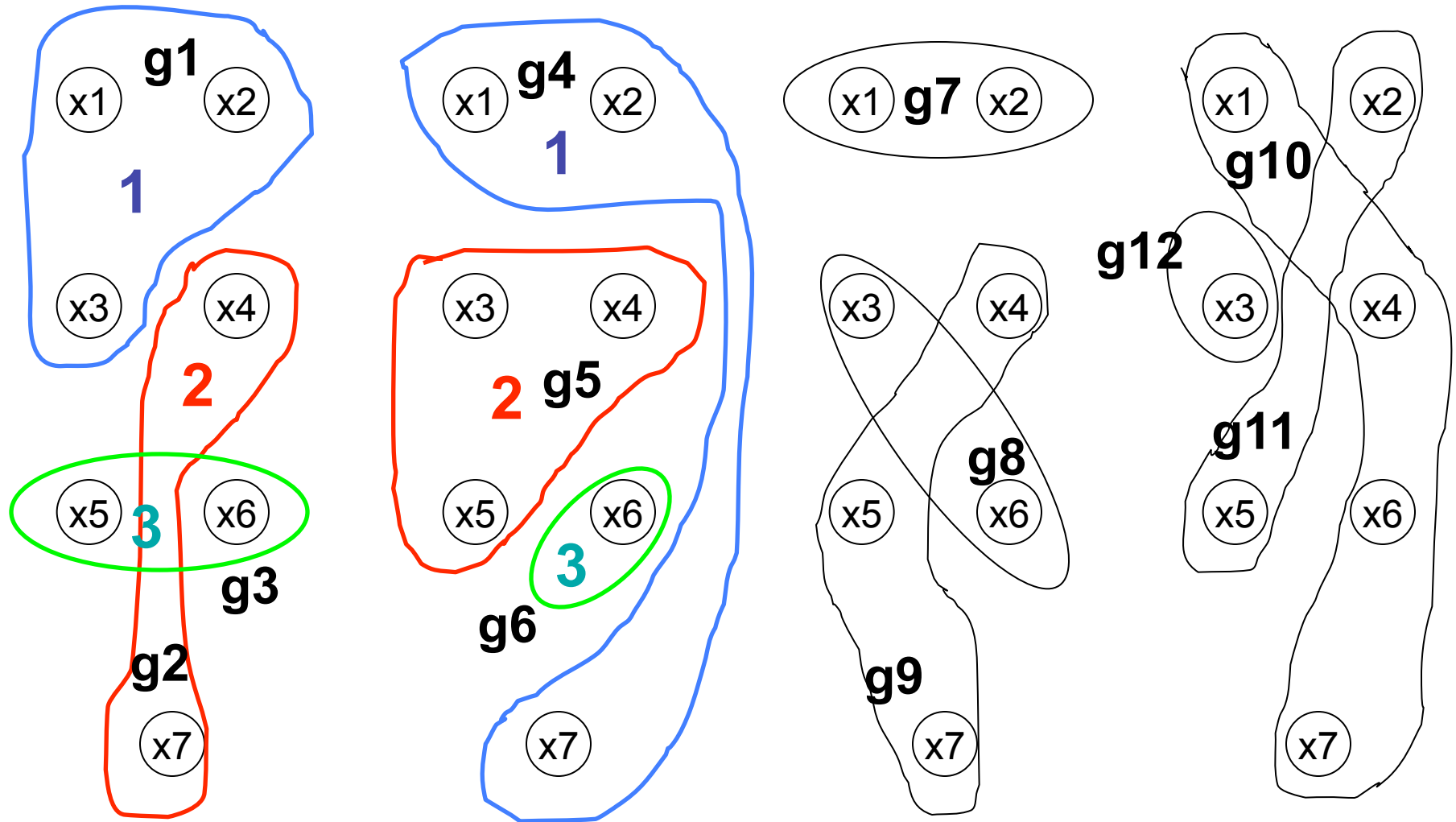
Problem



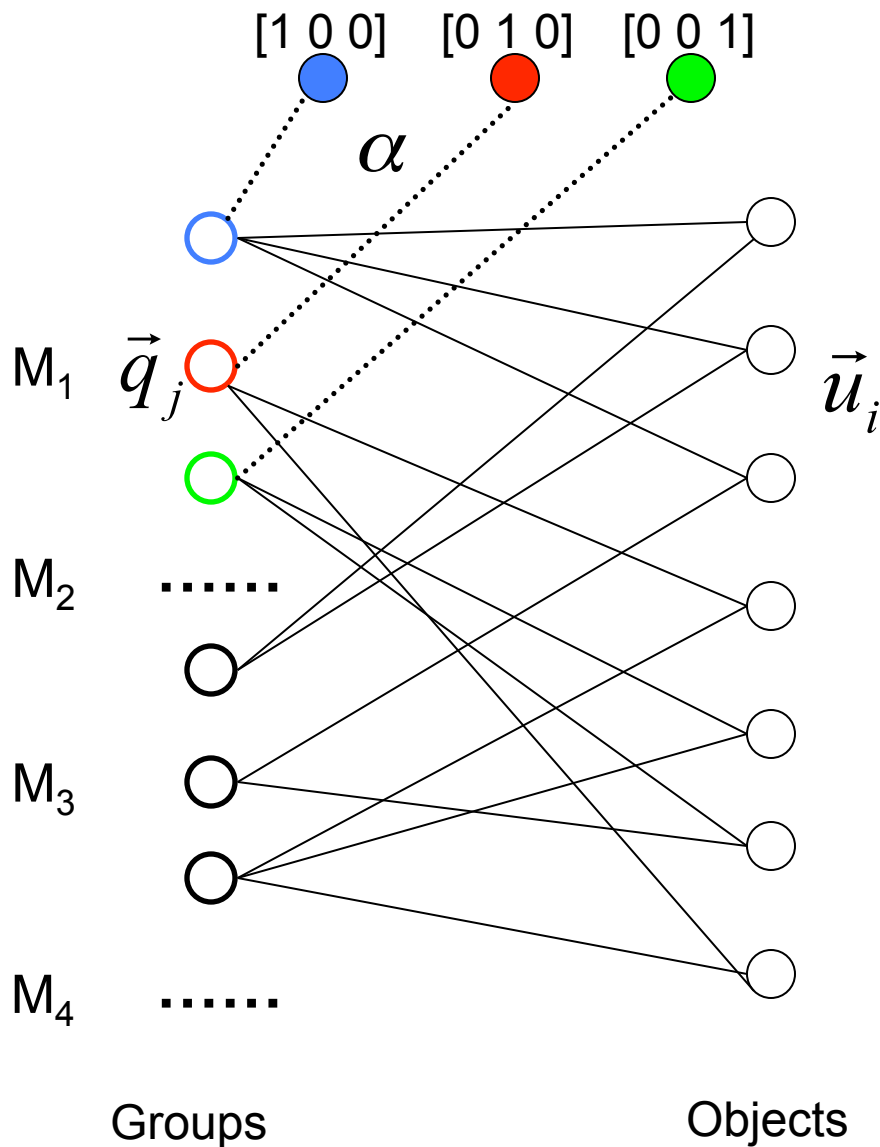
A Toy Example



Groups-Objects



Bipartite Graph



object i $\vec{u}_i = [u_{i1}, \dots, u_{ic}]$

group j $\vec{q}_j = [q_{j1}, \dots, q_{jc}]$

conditional prob vector

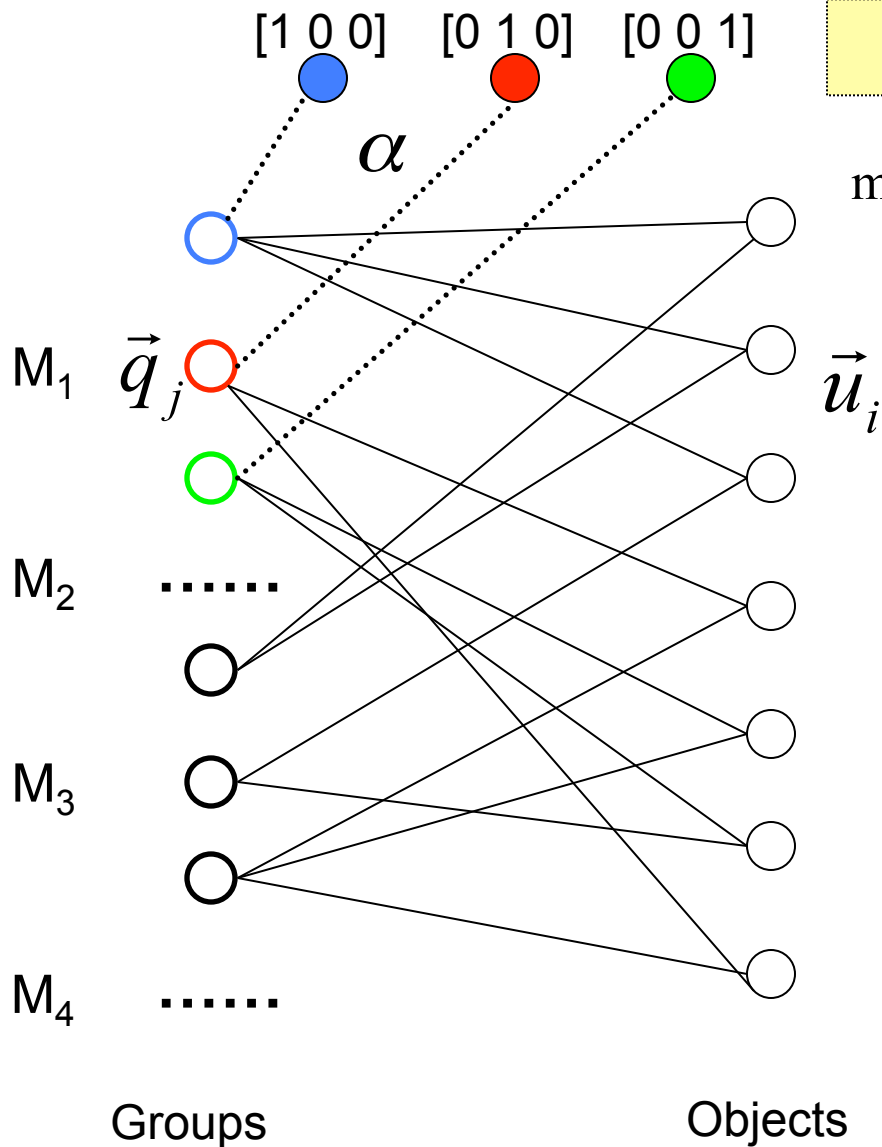
adjacency

$$a_{ij} = \begin{cases} 1 & u_i \leftrightarrow q_j \\ 0 & \text{otherwise} \end{cases}$$

groundtruth probability

$$\vec{y}_j = \begin{cases} [1 \ 0 \dots 0] & g_j \in 1 \\ \dots & \dots \\ [0 \ \dots 0 \ 1] & g_j \in c \end{cases}$$

Objective



minimize disagreement

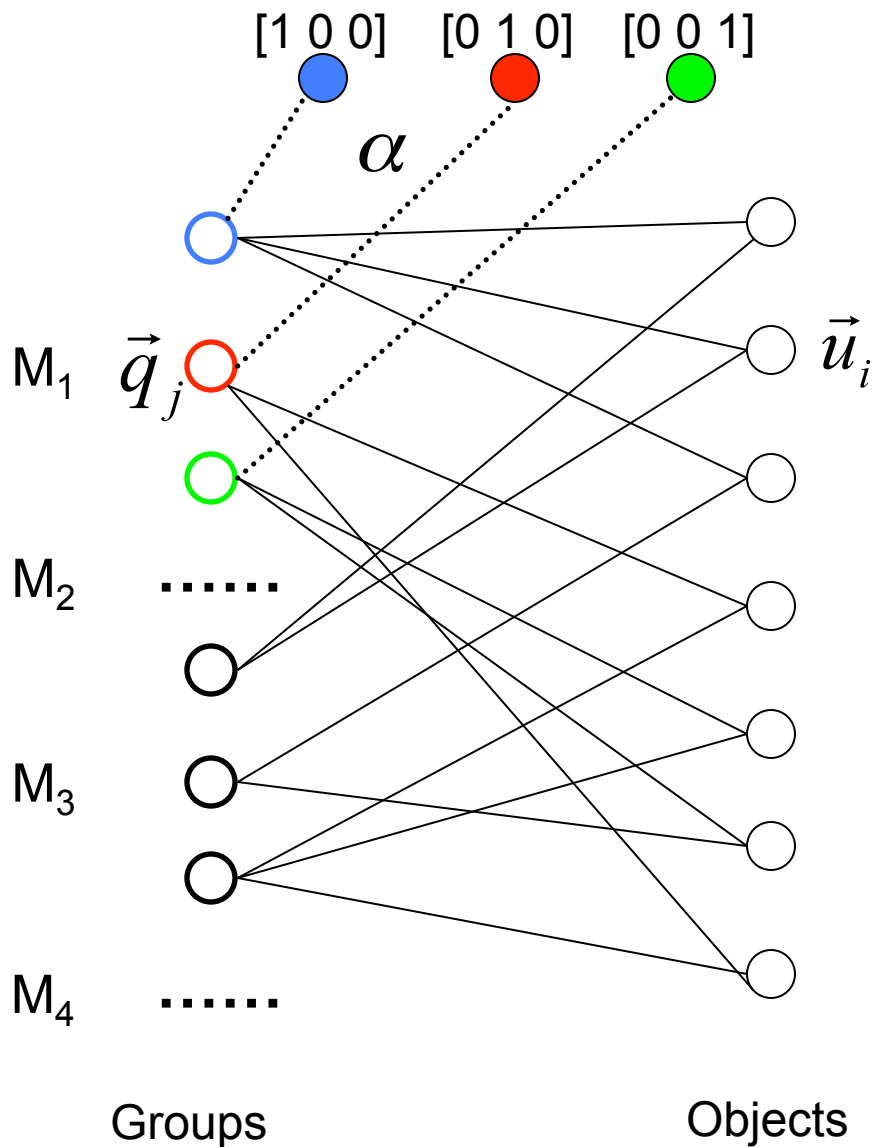
$$\min_{Q,U} \left(\sum_{i=1}^n \sum_{j=1}^v a_{ij} \|\vec{u}_i - \vec{q}_j\|^2 + \alpha \sum_{j=1}^s \|\vec{q}_j - \vec{y}_j\|^2 \right)$$

Similar conditional probability if the object is connected to the group

Do not deviate much from the groundtruth probability

Methodology

Iterate until convergence



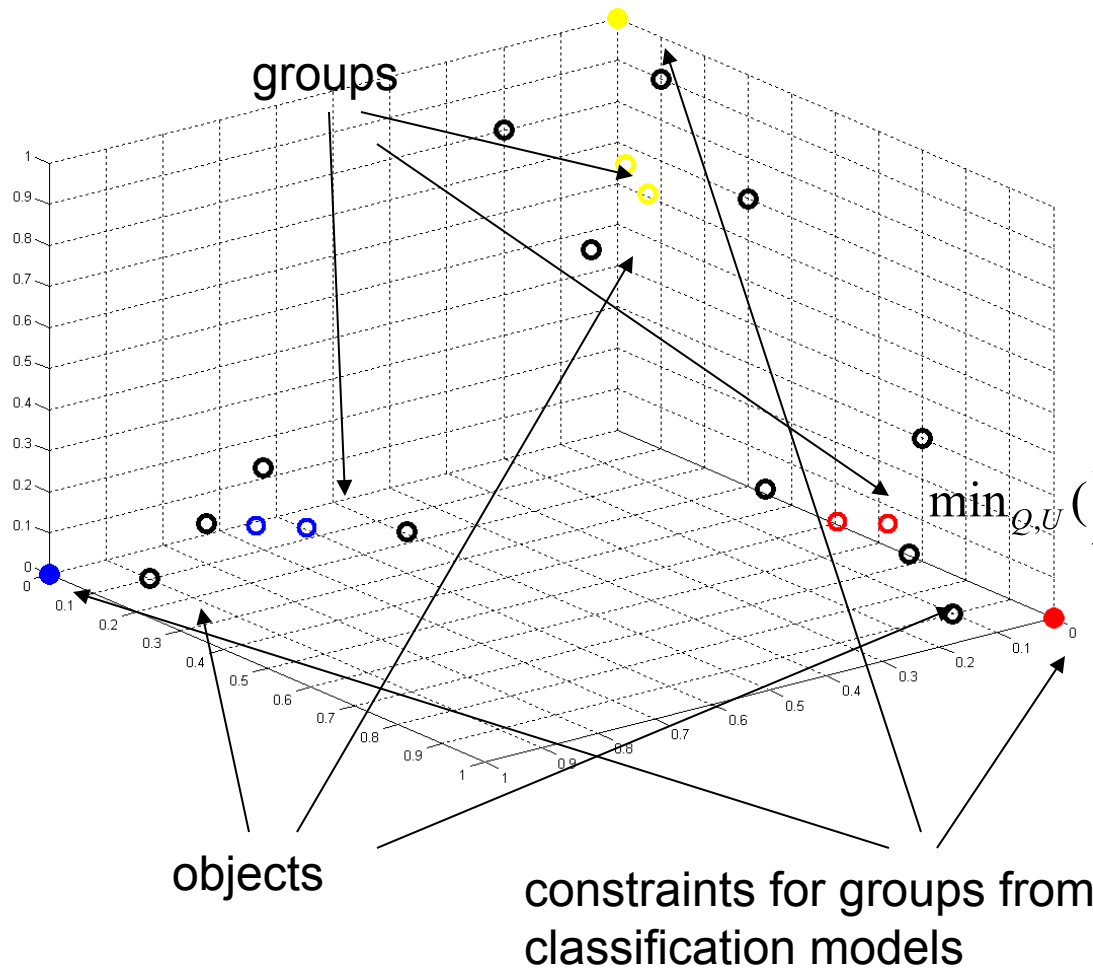
Update probability of a group

$$\vec{q}_j = \frac{\sum_{i=1}^n a_{ij} \vec{u}_i + \alpha \vec{y}_j}{\sum_{i=1}^n a_{ij} + \alpha} \quad \vec{q}_j = \frac{\sum_{i=1}^n a_{ij} \vec{u}_i}{\sum_{i=1}^n a_{ij}}$$

Update probability of an object

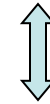
$$\vec{u}_i = \frac{\sum_{j=1}^v a_{ij} \vec{q}_j}{\sum_{j=1}^v a_{ij}}$$

Constrained Embedding



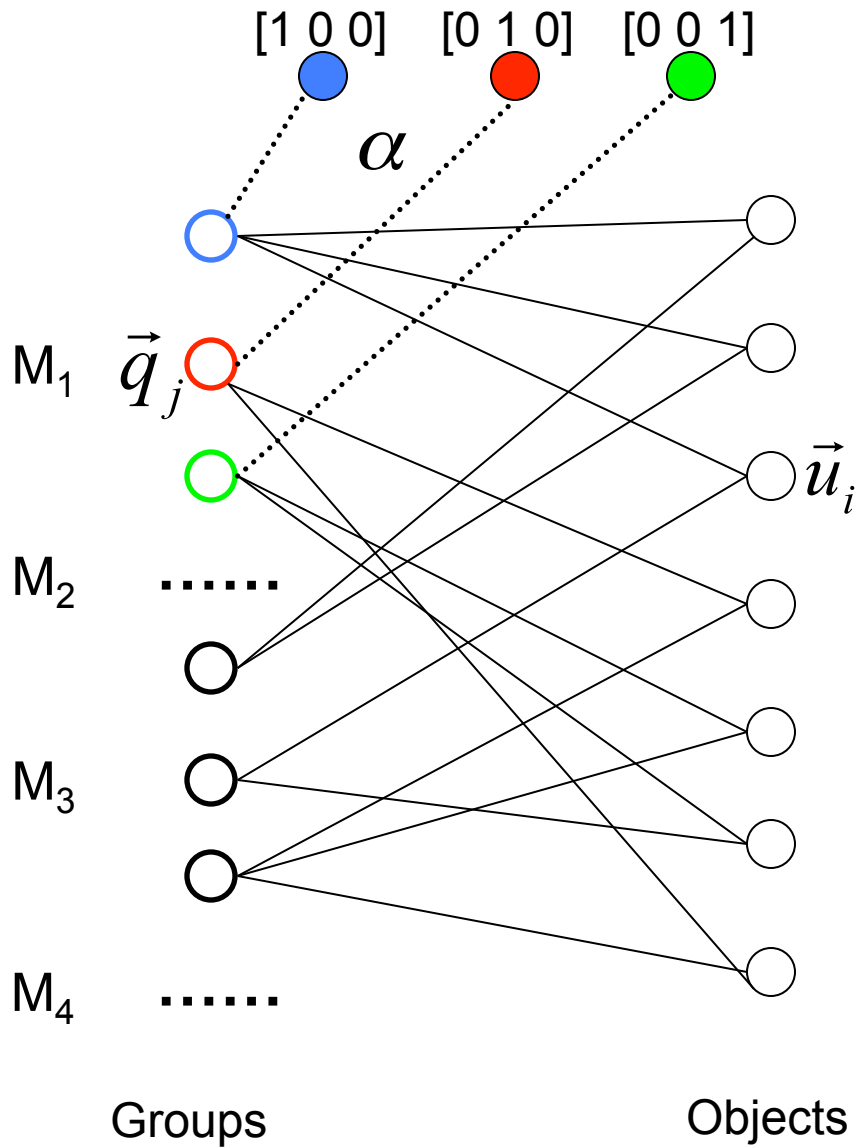
$$\min_{Q,U} \sum_{j=1}^v \sum_{z=1}^c \left| q_{jz} - \frac{\sum_{i=1}^n a_{ij} u_{iz}}{\sum_{i=1}^n a_{ij}} \right|$$

$$q_{jz} = 1 \text{ if } g_j \text{'s label is } z$$



$$\min_{Q,U} \left(\sum_{i=1}^n \sum_{j=1}^v a_{ij} \|\vec{u}_i - \vec{q}_j\|^2 + \alpha \sum_{j=1}^s \|\vec{q}_j - \vec{y}_j\|^2 \right)$$

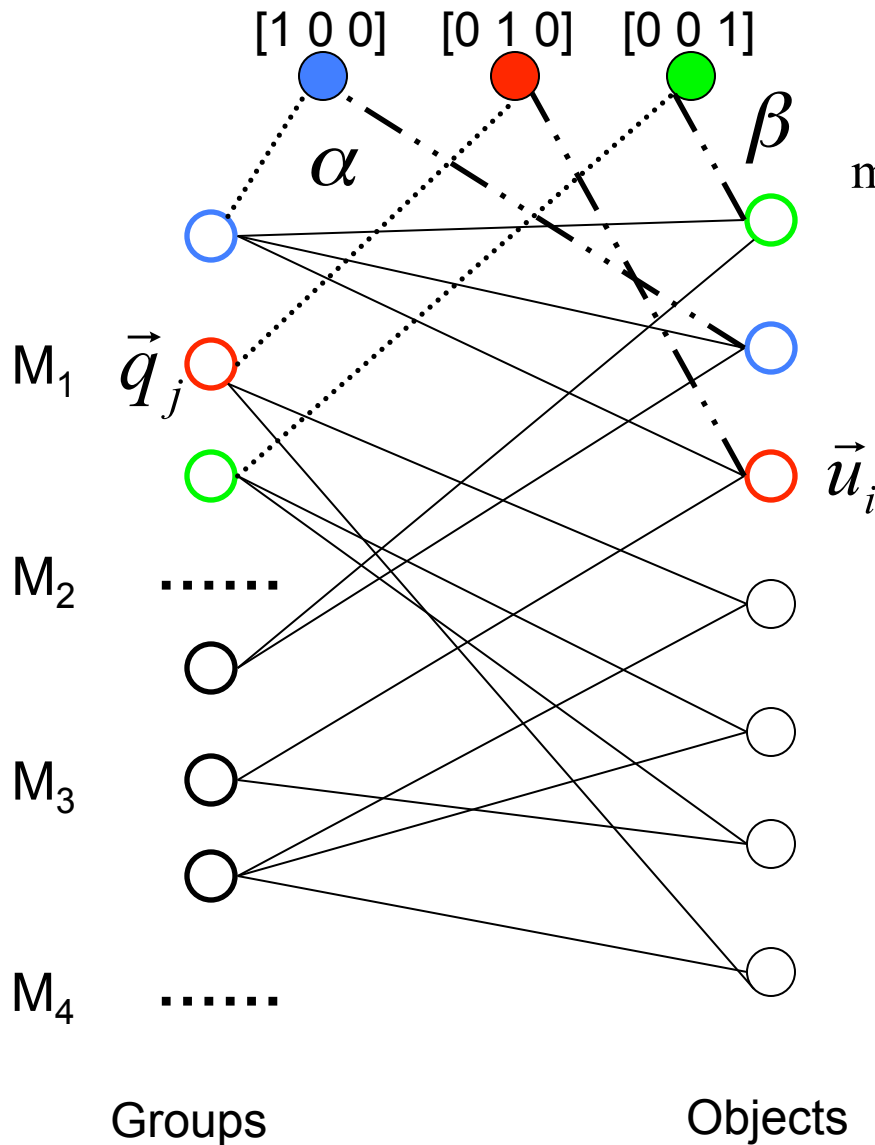
Ranking on Consensus Structure



$$\vec{q}_{.1} = D_\lambda (D_v^{-1} A^T D_n^{-1} A) \vec{q}_{.1} + D_{1-\lambda} \vec{y}_{.1}$$

Diagram illustrating the equation above. The left side of the equation is labeled **adjacency matrix**. The right side is labeled **personalized damping factors**. The right side of the equation is labeled **query**. The diagram shows a bipartite graph between groups and objects, with a query vector \vec{q}_j and a query vector $\vec{y}_{.1}$ shown. The graph is connected to the equation via arrows.

Incorporating Labeled Information



Objective

$$\min_{Q,U} \left(\sum_{i=1}^n \sum_{j=1}^v a_{ij} \|\vec{u}_i - \vec{q}_j\|^2 + \alpha \sum_{j=1}^s \|\vec{q}_j - \vec{y}_j\|^2 \right) + \beta \sum_{i=1}^l \|\vec{u}_i - \vec{f}_i\|^2$$

Update probability of a group

$$\vec{q}_j = \frac{\sum_{i=1}^n a_{ij} \vec{u}_i + \alpha \vec{y}_j}{\sum_{i=1}^n a_{ij} + \alpha} \quad \vec{q}_j = \frac{\sum_{i=1}^n a_{ij} \vec{u}_i}{\sum_{i=1}^n a_{ij}}$$

Update probability of an object

$$\vec{u}_i = \frac{\sum_{j=1}^v a_{ij} \vec{q}_j}{\sum_{j=1}^v a_{ij}} \quad \vec{u}_i = \frac{\sum_{j=1}^v a_{ij} \vec{q}_j + \beta \vec{f}_i}{\sum_{j=1}^v a_{ij} + \beta}$$

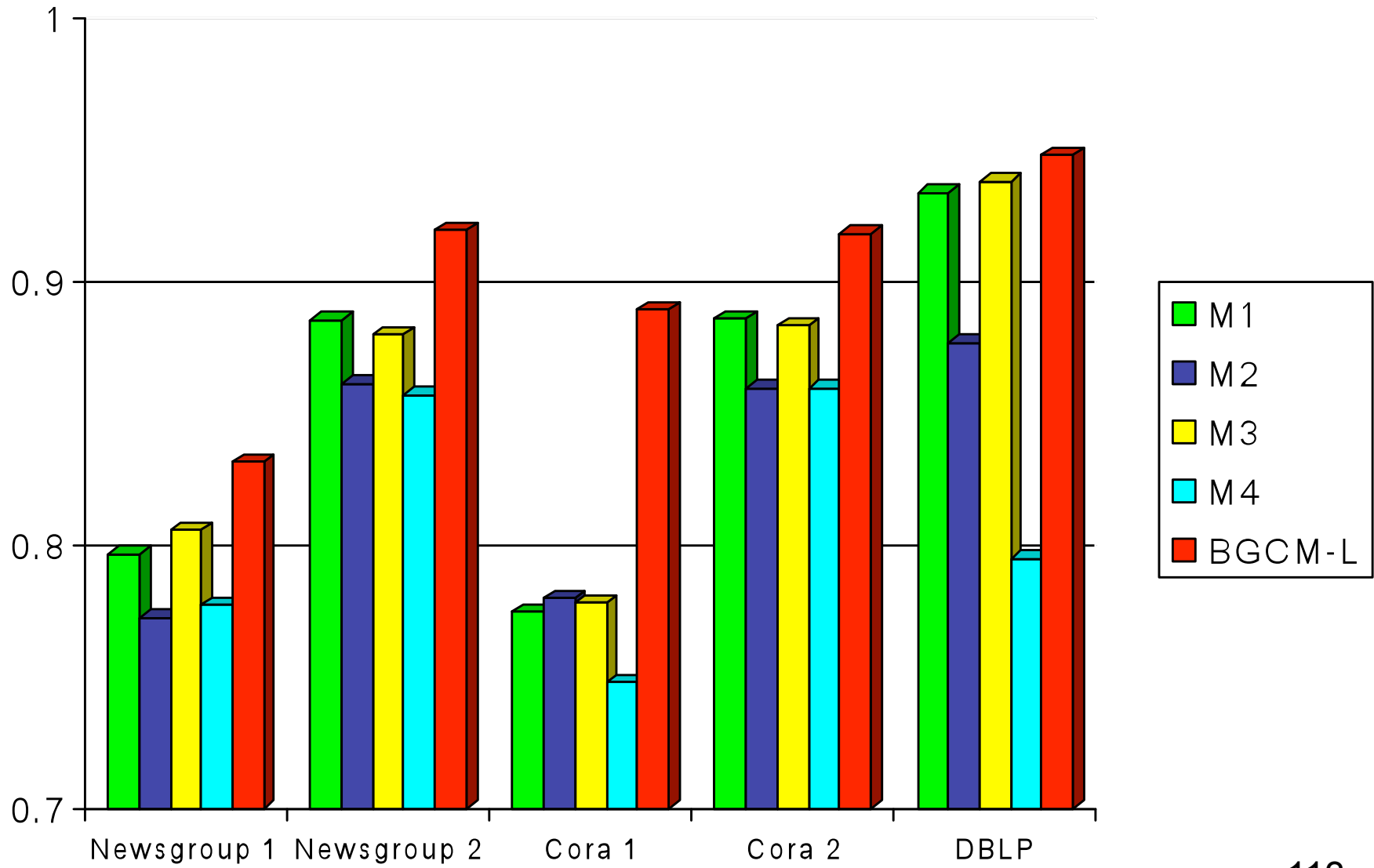
Experiments-Data Sets

- **20 Newsgroup**
 - newsgroup messages categorization
 - only text information available
- **Cora**
 - research paper area categorization
 - paper abstracts and citation information available
- **DBLP**
 - researchers area prediction
 - publication and co-authorship network, and publication content
 - conferences' areas are known

Experiments-Baseline Methods

- **Single models**
 - 20 Newsgroup:
 - logistic regression, SVM, K-means, min-cut
 - Cora
 - abstracts, citations (with or without a labeled set)
 - DBLP
 - publication titles, links (with or without labels from conferences)
- **Proposed method**
 - BGCM
 - BGCM-L: semi-supervised version combining four models

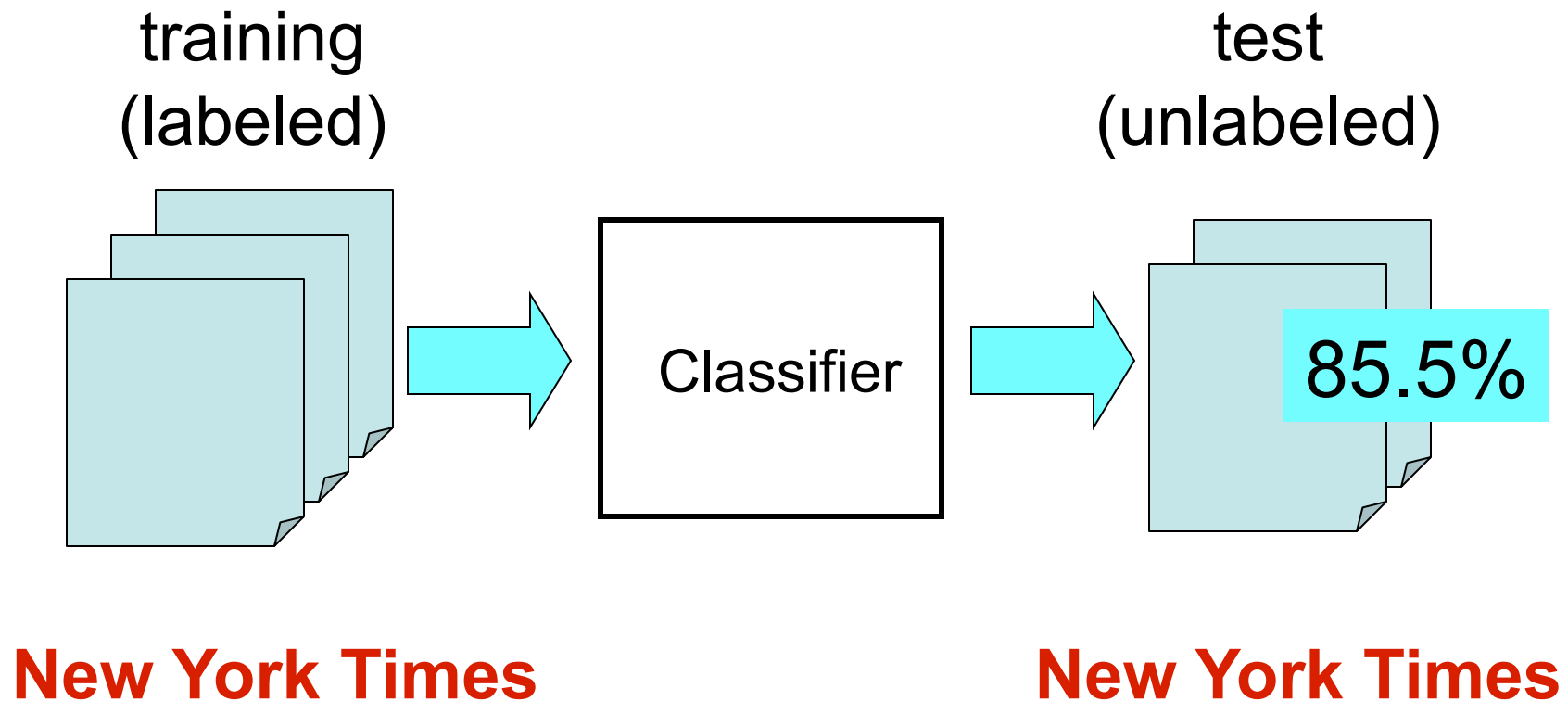
Accuracy



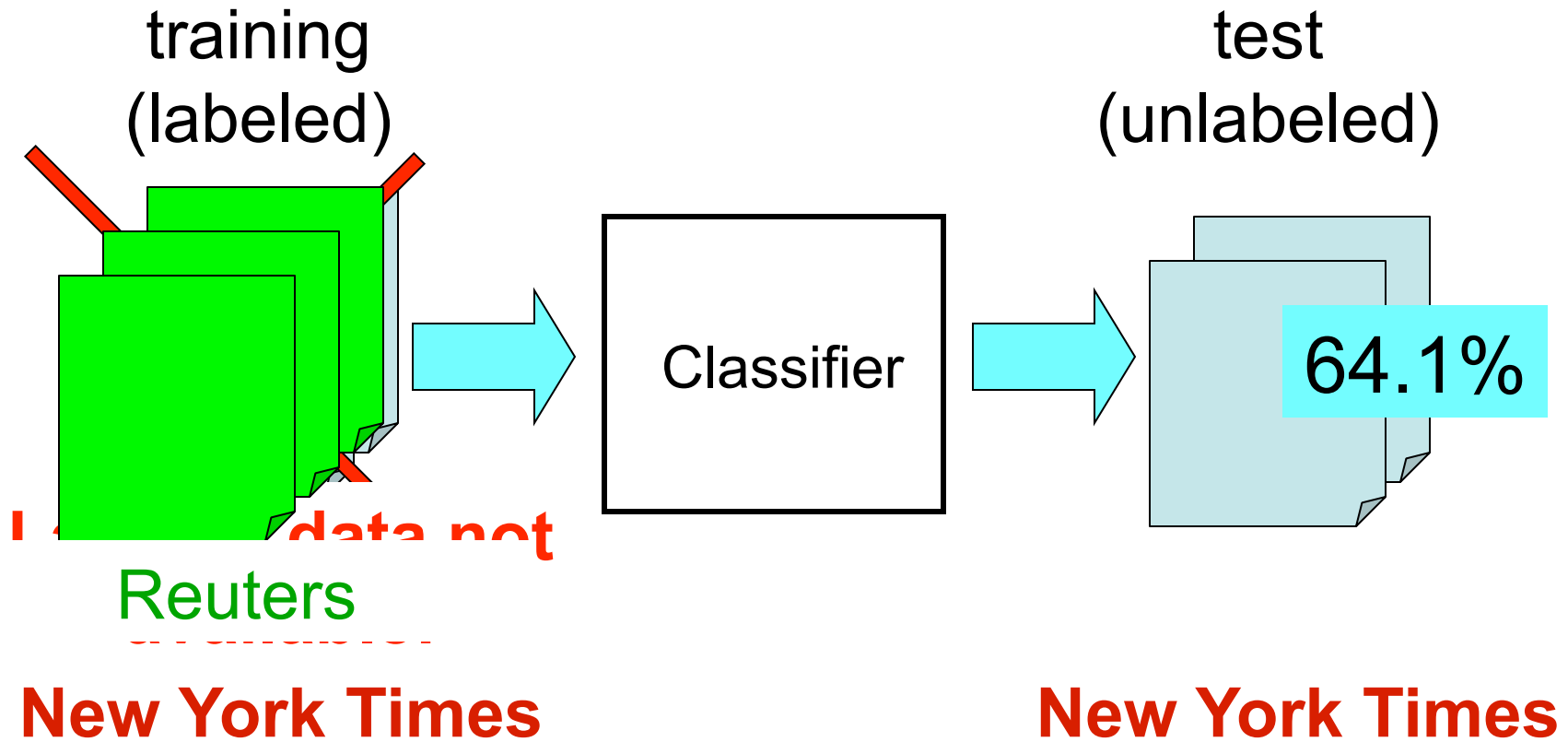
Outline

- An overview of ensemble methods
 - Motivations
 - Tutorial overview
- Supervised ensemble
- Unsupervised ensemble
- Semi-supervised ensemble
 - Multi-view learning
 - Consensus maximization among supervised and unsupervised models
- Applications
 - Transfer learning, stream classification, anomaly detection

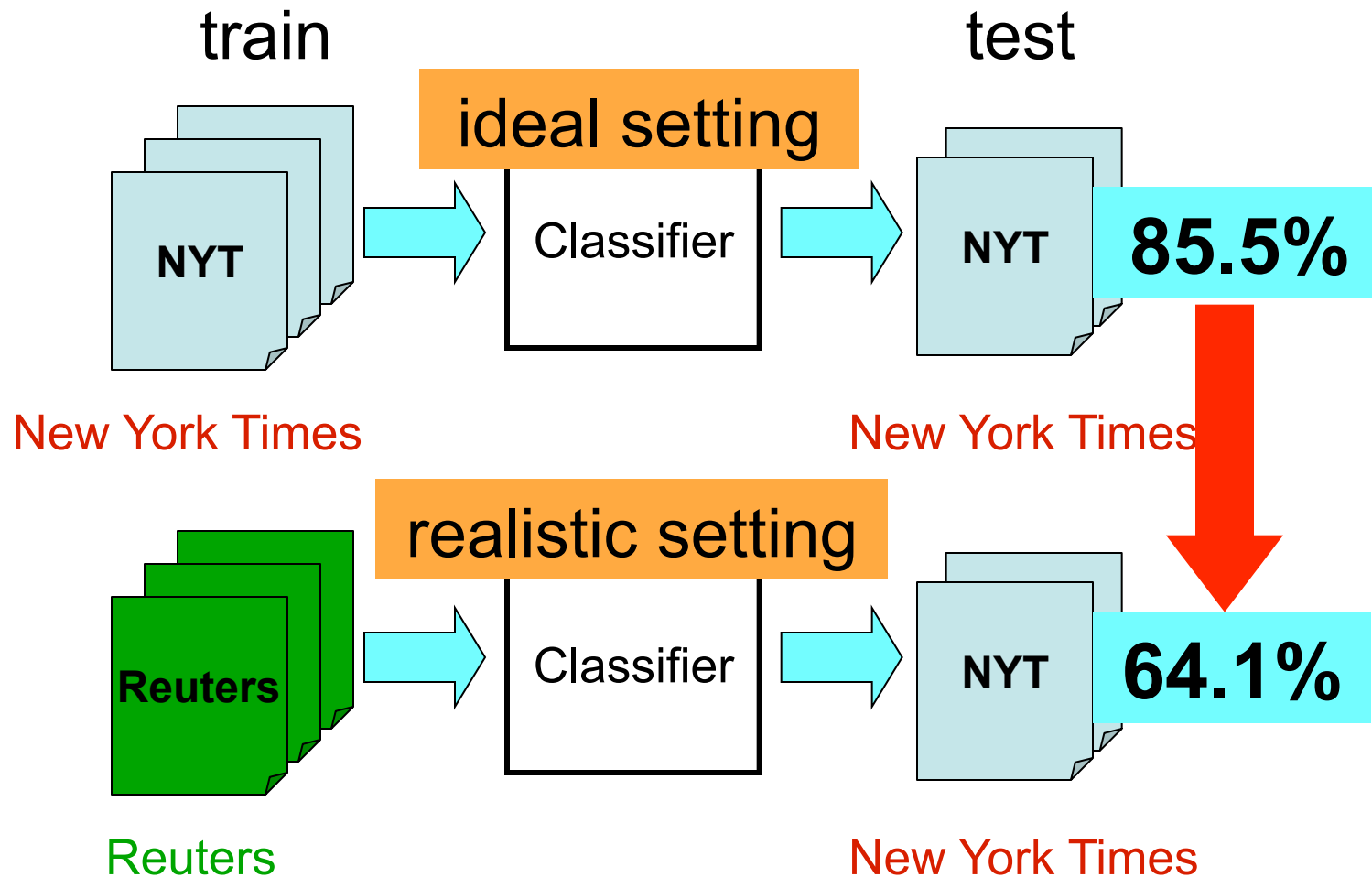
Standard Supervised Learning



In Reality.....



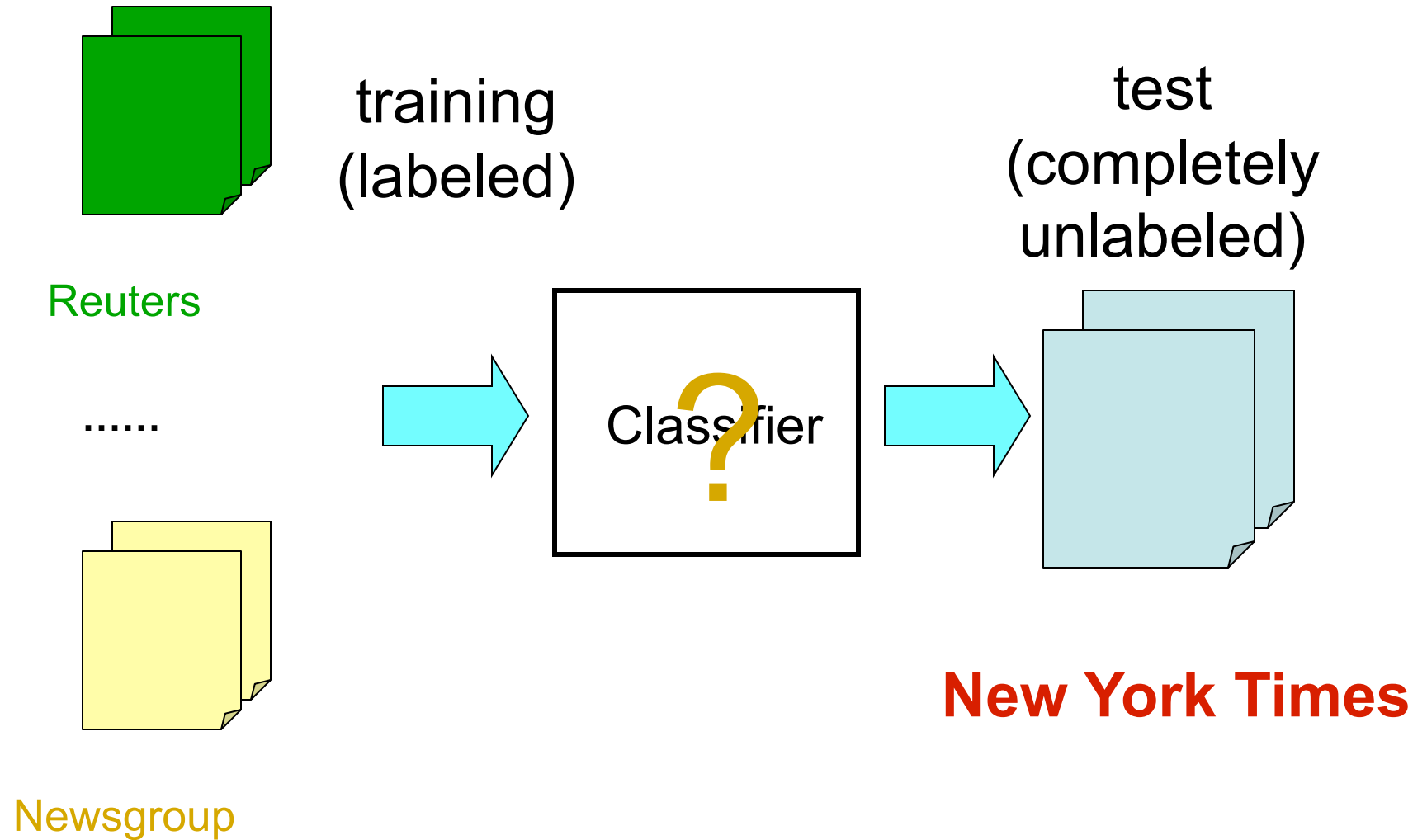
Domain Difference → Performance Drop



Other Examples

- **Spam filtering**
 - Public email collection → personal inboxes
- **Intrusion detection**
 - Existing types of intrusions → unknown types of intrusions
- **Sentiment analysis**
 - Expert review articles → blog review articles
- **The aim**
 - To design learning methods that are aware of the training and test domain difference
- **Transfer learning**
 - Adapt the classifiers learnt from the source domain to the new domain

All Sources of Labeled Information

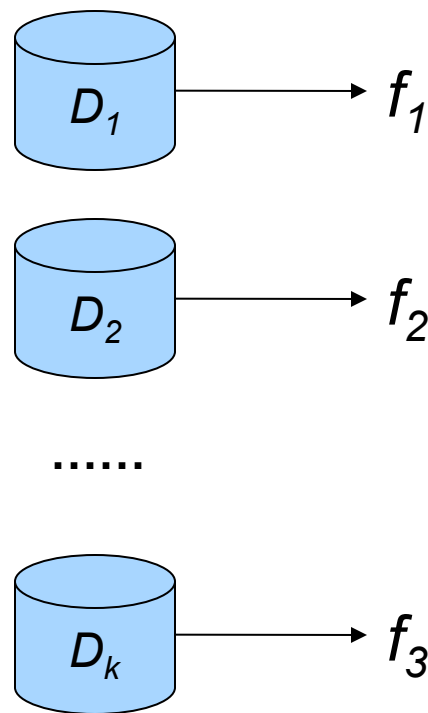


Consensus Regularization Approach* (1)

- Basic idea

- Train k classifiers from k source domains simultaneously
- Incorporate the constraint that the k classifiers reach consensus on the unlabeled data from the target domain

Likelihood



Constraint

	f_1	f_2	f_3
V_1	+	+	+
V_2	-	-	-
V_3	+	+	+
V_4	-	+	-
V_5	+	+	-
V_6	+	-	-



*[LZX+08]

Consensus Regularization Approach (2)

- Optimization framework

- Binary classification
- Base model: logistic regression (on each source domain)

$$\sum_{i=1}^n \log \frac{1}{1 + \exp(-y_i w^T x_i)} - \frac{\lambda}{2} w^T w$$

- Constraint: favoring skewed conditional probability for each object (on target domain)

$$\sum_{i=1}^n (\bar{P}(y = 1 | x) - \bar{P}(y = -1 | x))^2$$

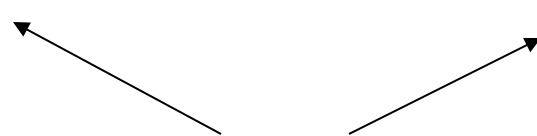
- Maximize: Data Likelihood + constraint violation penalty
- Method: Conjugate gradient

Multiple Model Local Structure Mapping*

- **Locally weighted ensemble framework**
 - transfer useful knowledge from multiple source domains
 - adapt the knowledge to the target domain
- **Graph-based heuristics to compute weights**
 - make the framework practical and effective

*[GFJ+08]

A Synthetic Example



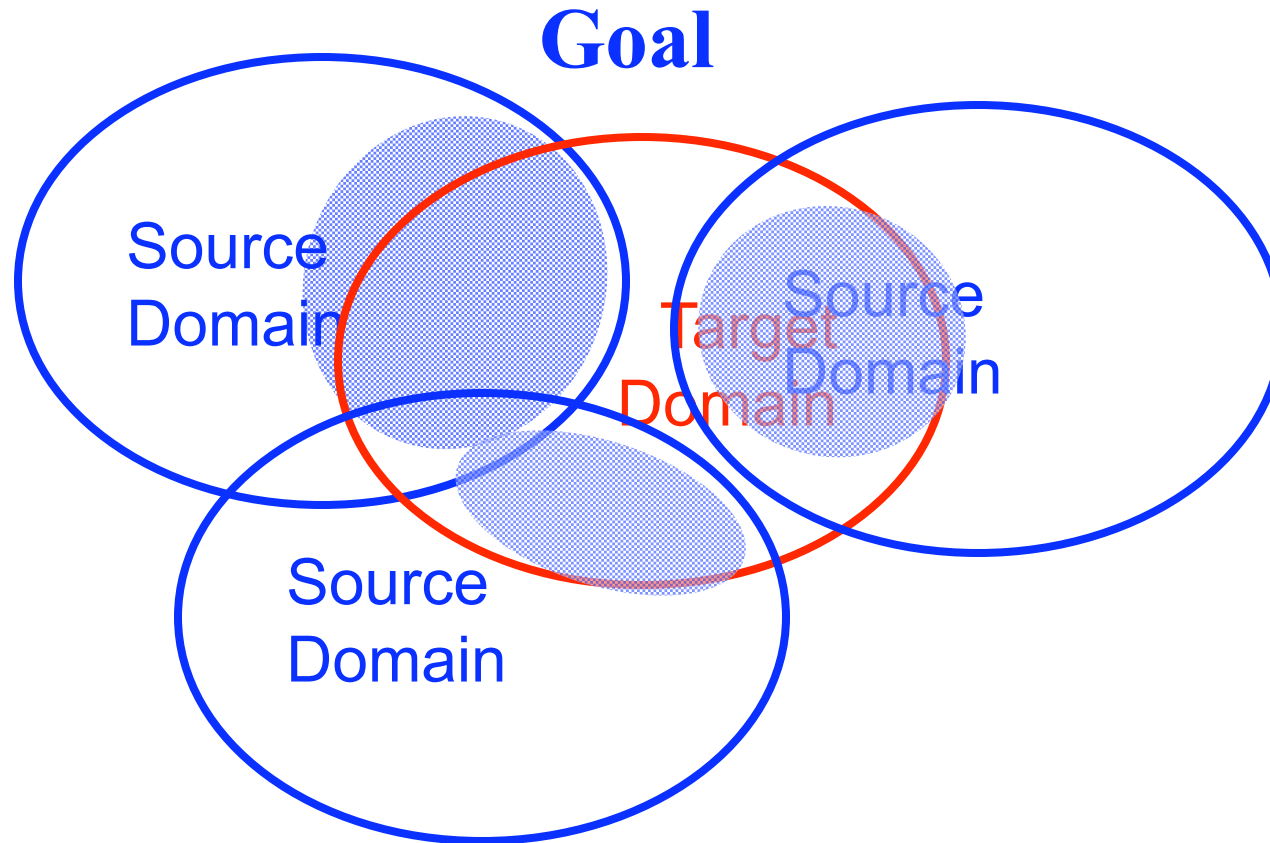
Training



Test



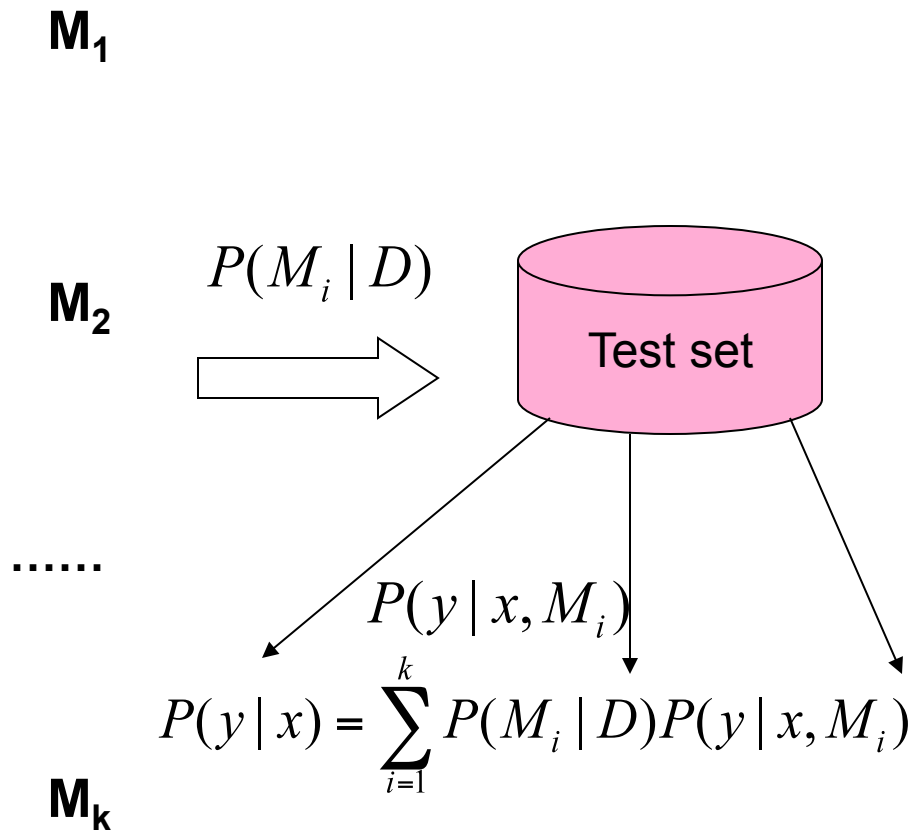
(have conflicting concepts) **Partially overlapping**



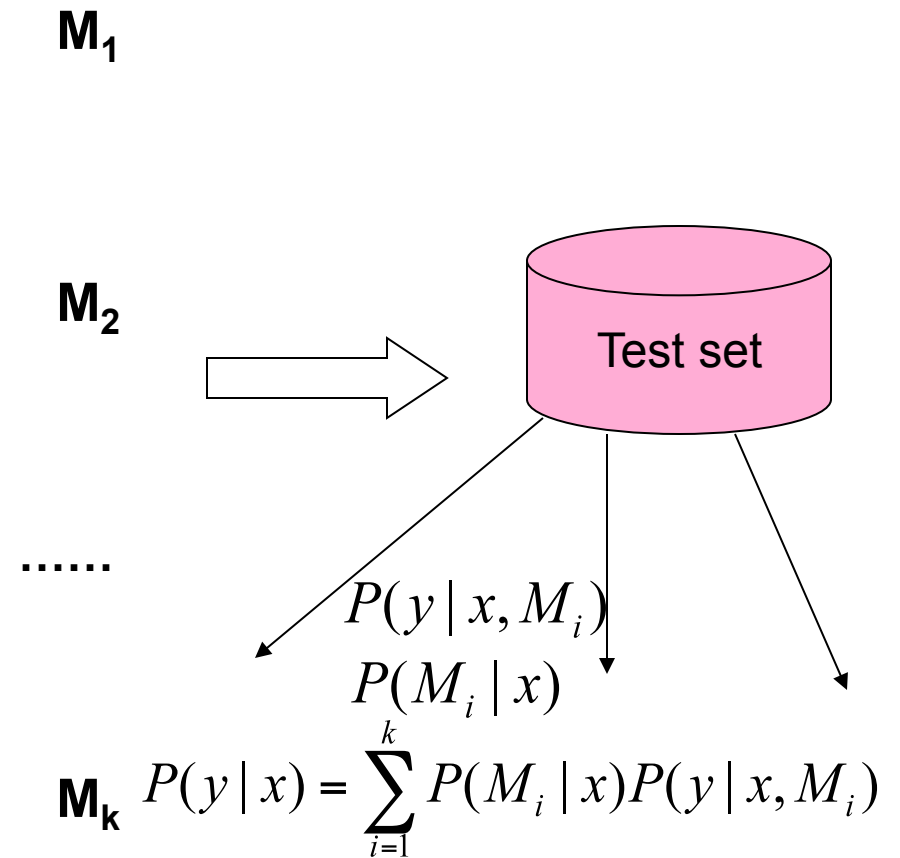
- To unify knowledge that are consistent with the test domain from multiple source domains (models)

Global versus Local Weights (1)

Global weighting



Local weighting



Global versus Local Weights (2)

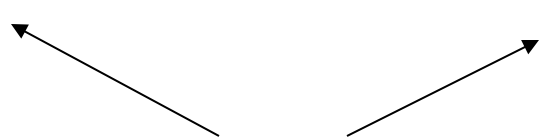
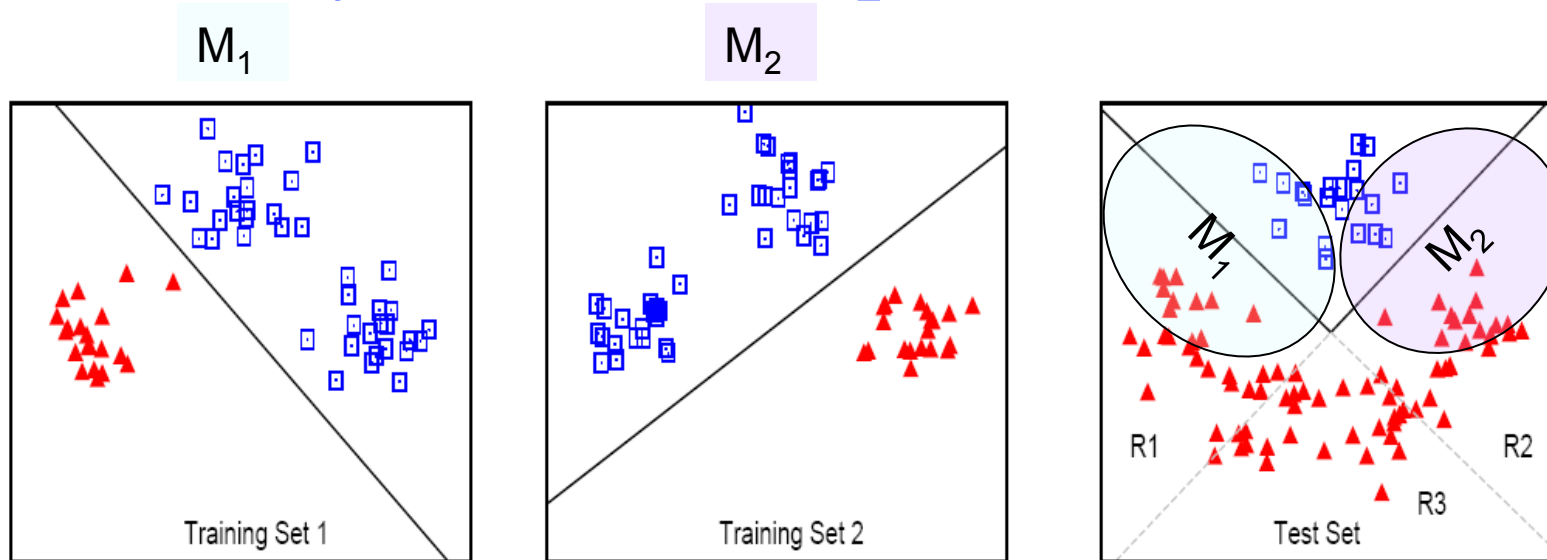
x		y	M ₁	w _g	w _l	M ₂	w _g	w _l
2.40	5.23	1	0.6	0.3	0.2	0.9	0.7	0.8
-2.69	0.55	0	0.4	0.3	0.6	0.6	0.7	0.4
-3.97	-3.62	0	0.2	0.3	0.7	0.4	0.7	0.3
2.08	-3.73	0	0.1	0.3	0.5	0.1	0.7	0.5
5.08	2.15	0	0.6	0.3	0.3	0.3	0.7	0.7
1.43	4.48	1	1	0.3	1	0.2	0.7	0
.....

- **Locally weighting scheme**

- Weight of each model is computed per example
- Weights are determined according to models' performance on the test set, not training set

Training

Synthetic Example Revisited



Training



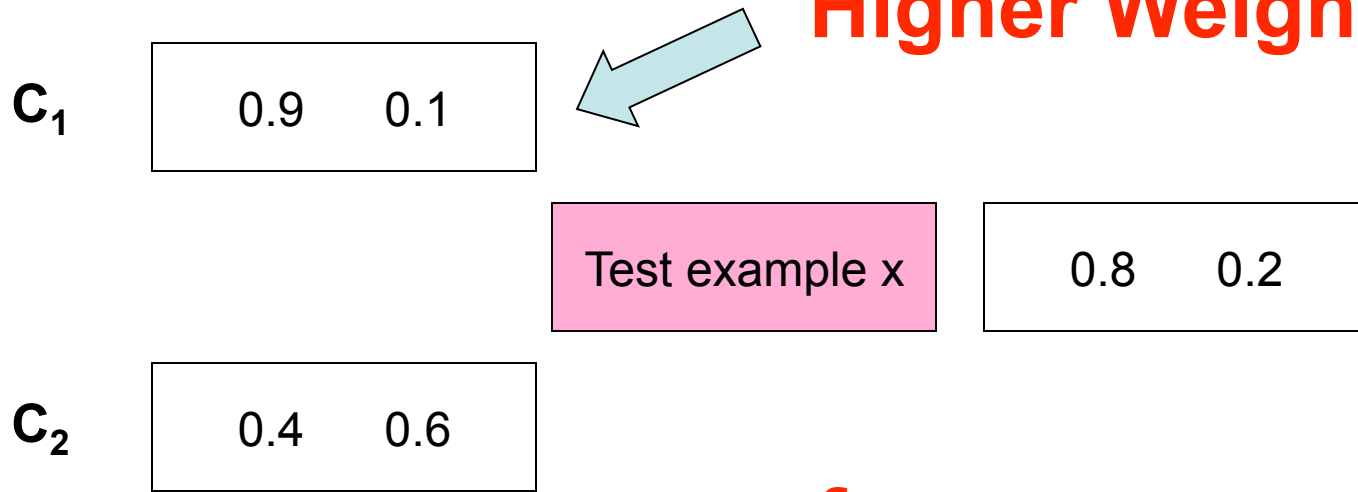
Test



(have conflicting concepts) **Partially overlapping**

Optimal Local Weights

Higher Weight

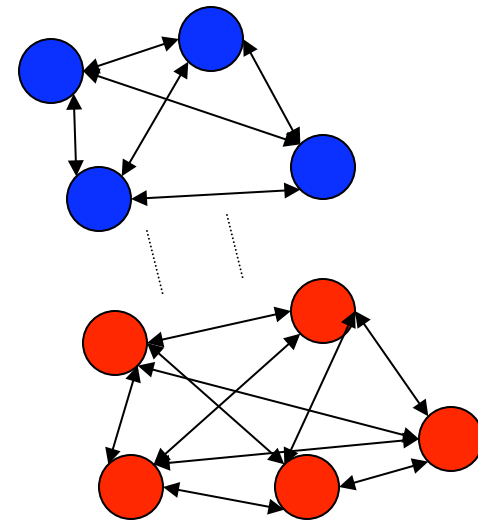
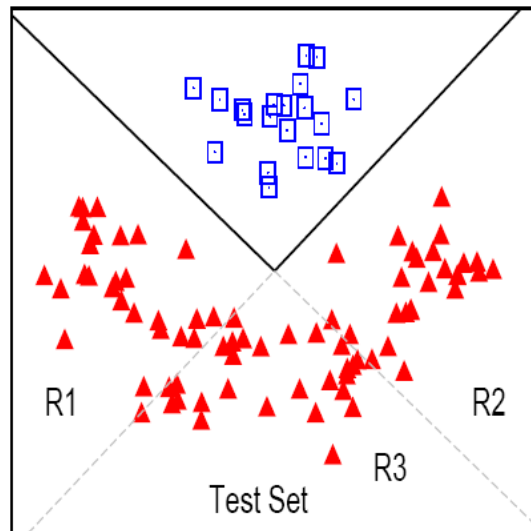


$$\mathbf{H} \begin{pmatrix} 0.9 & 0.4 \\ 0.1 & 0.6 \end{pmatrix} \begin{pmatrix} w^1 \\ w^2 \end{pmatrix} = \begin{pmatrix} 0.8 \\ 0.2 \end{pmatrix} \quad \mathbf{f} \quad \sum_{i=1}^k w^i(x) = 1$$

- **Optimal weights**
 - Solution to a regression problem
 - Impossible to get since \mathbf{f} is unknown!

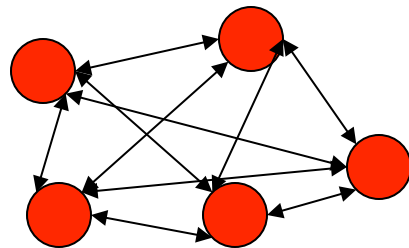
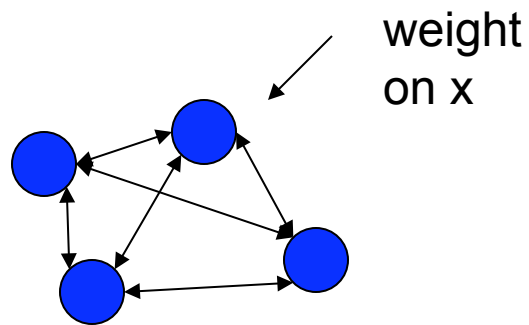
Clustering-Manifold Assumption

Test examples that are closer in feature space are more likely to share the same class label.

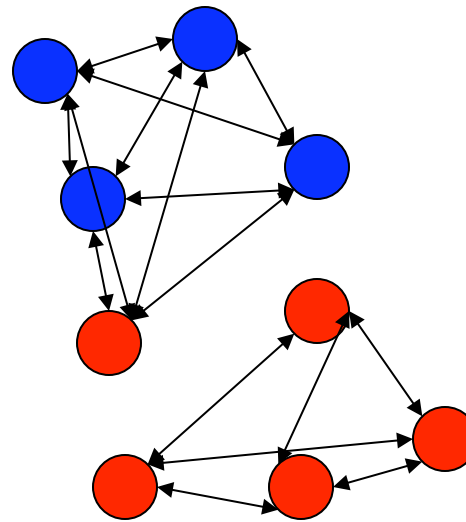


Graph-based Heuristics

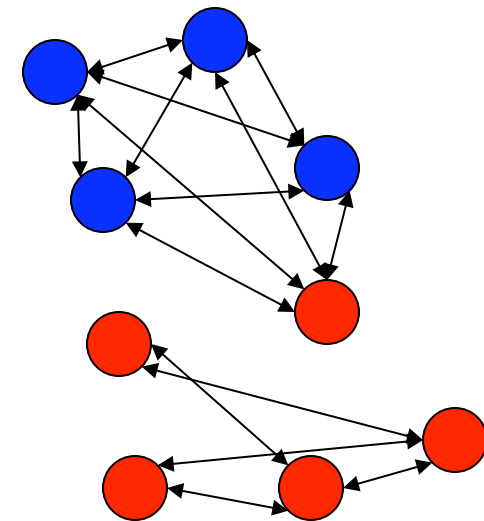
- Graph-based weights approximation
 - Map the structures of models onto test domain



Clustering
Structure

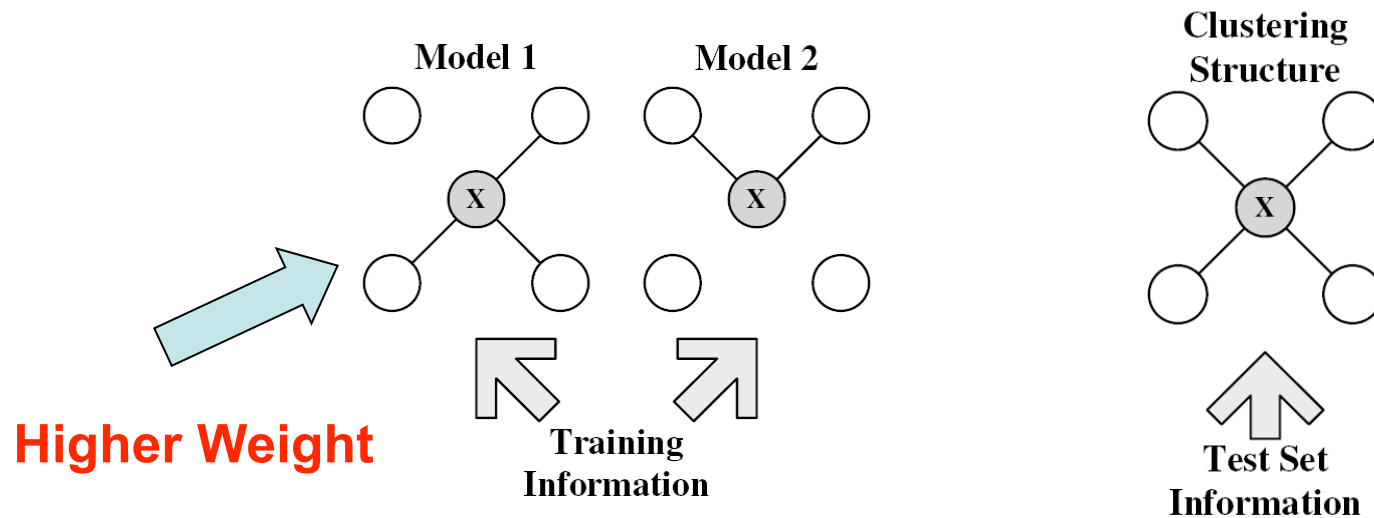


M_1



M_2

Graph-based Heuristics



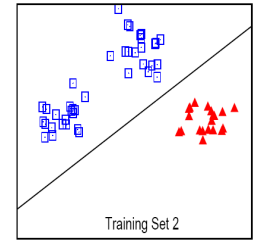
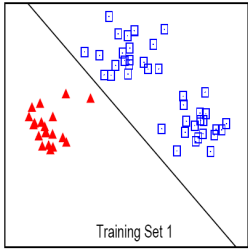
- **Local weights calculation**

- Weight of a model is proportional to the similarity between its neighborhood graph and the clustering structure around x .

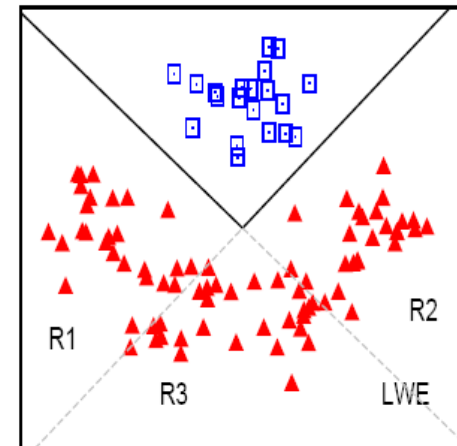
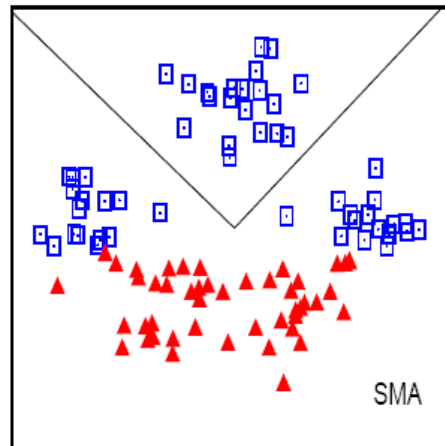
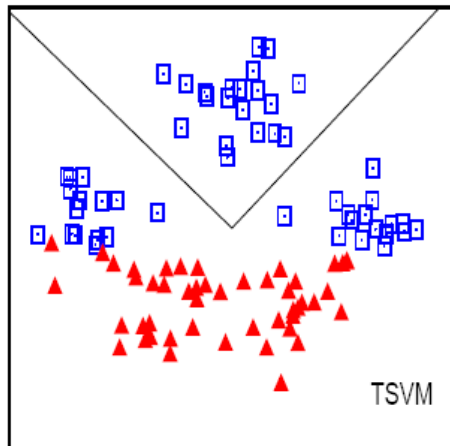
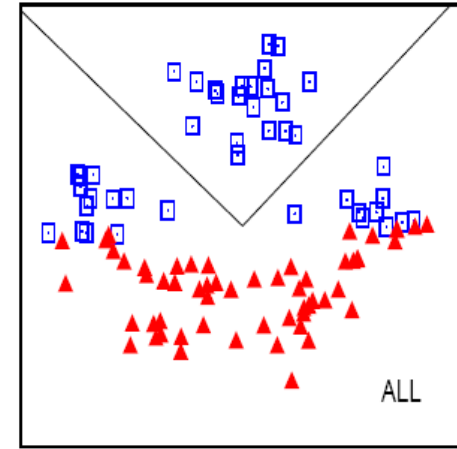
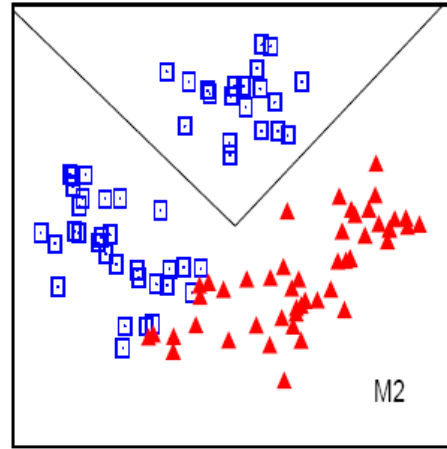
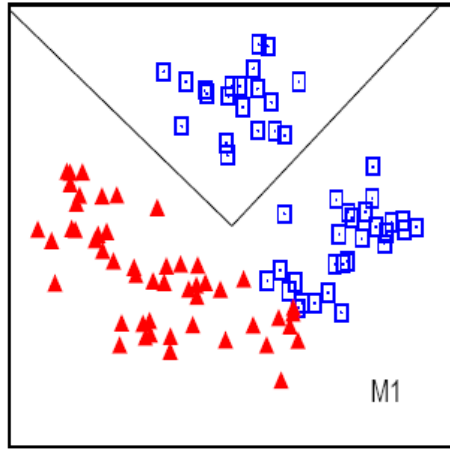
$$w_{M,x} \propto s(G_M, G_T; \mathbf{x}) = \frac{\sum_{v_1 \in V_M} \sum_{v_2 \in V_T} \mathbf{1}\{v_1 = v_2\}}{|V_M| + |V_T|}$$

Experiments Setup

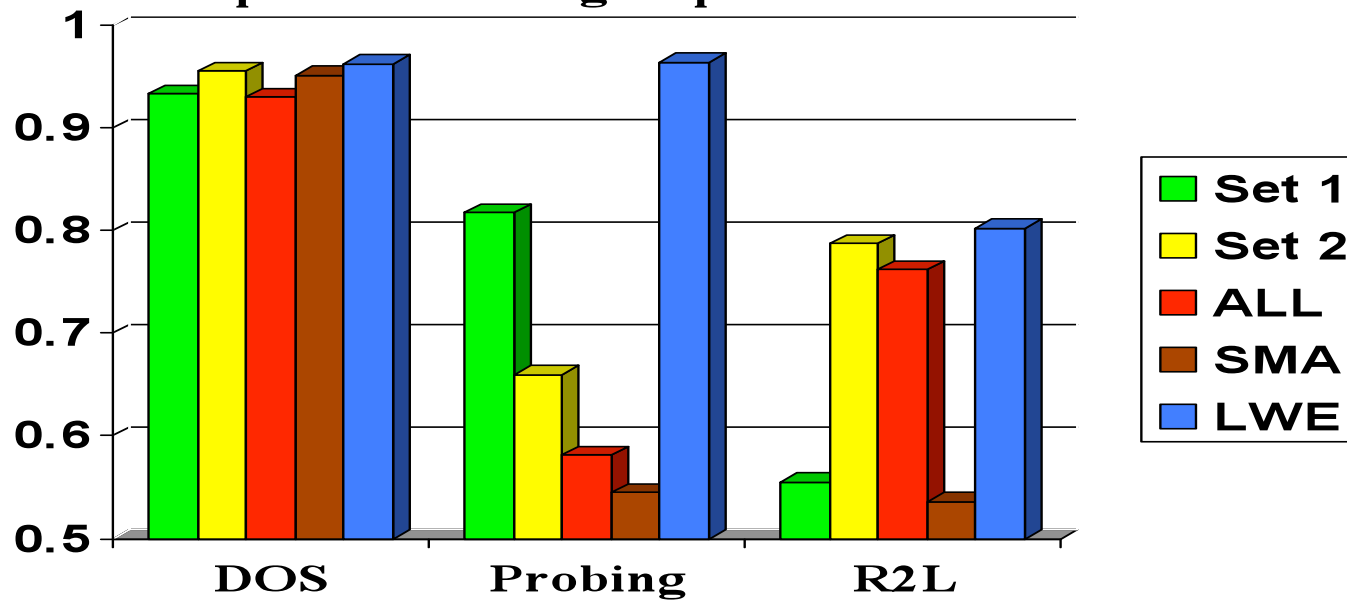
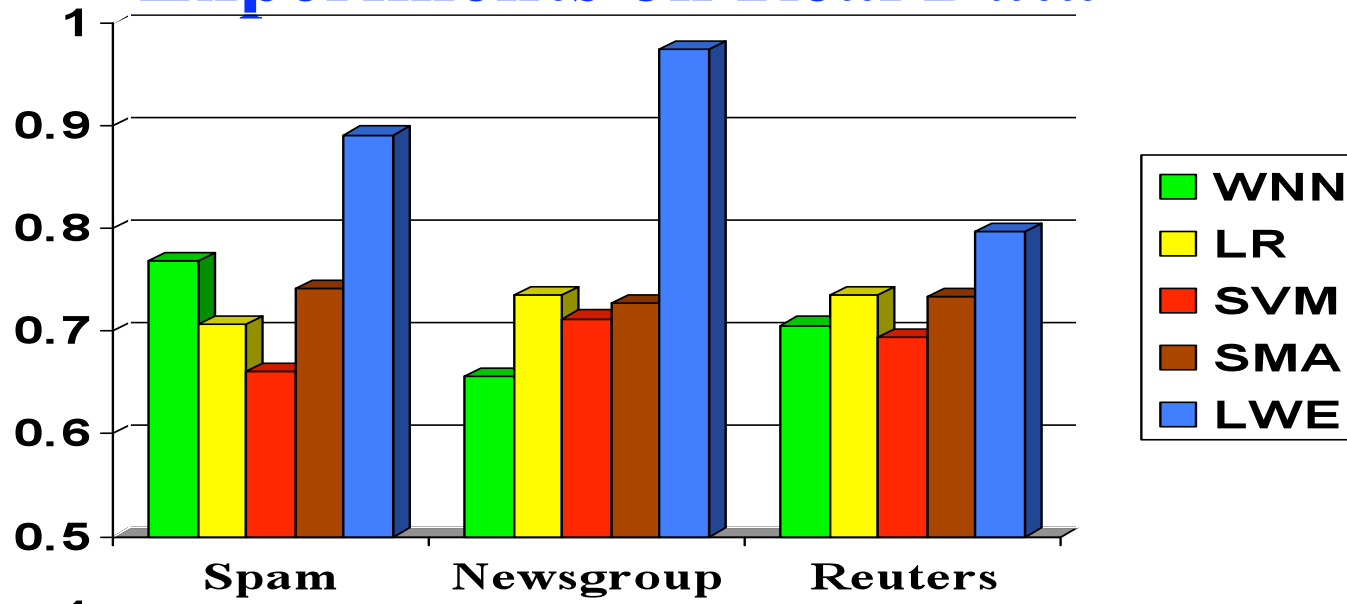
- **Data Sets**
 - Synthetic data sets
 - Spam filtering: public email collection → personal inboxes (u01, u02, u03) (ECML/PKDD 2006)
 - Text classification: same top-level classification problems with different sub-fields in the training and test sets (Newsgroup, Reuters)
 - Intrusion detection data: different types of intrusions in training and test sets.
- **Baseline Methods**
 - One source domain: single models (WNN, LR, SVM)
 - Multiple source domains: SVM on each of the domains
 - Merge all source domains into one: ALL
 - Simple averaging ensemble: SMA
 - Locally weighted ensemble: LWE



Experiments on Synthetic Data



Experiments on Real Data



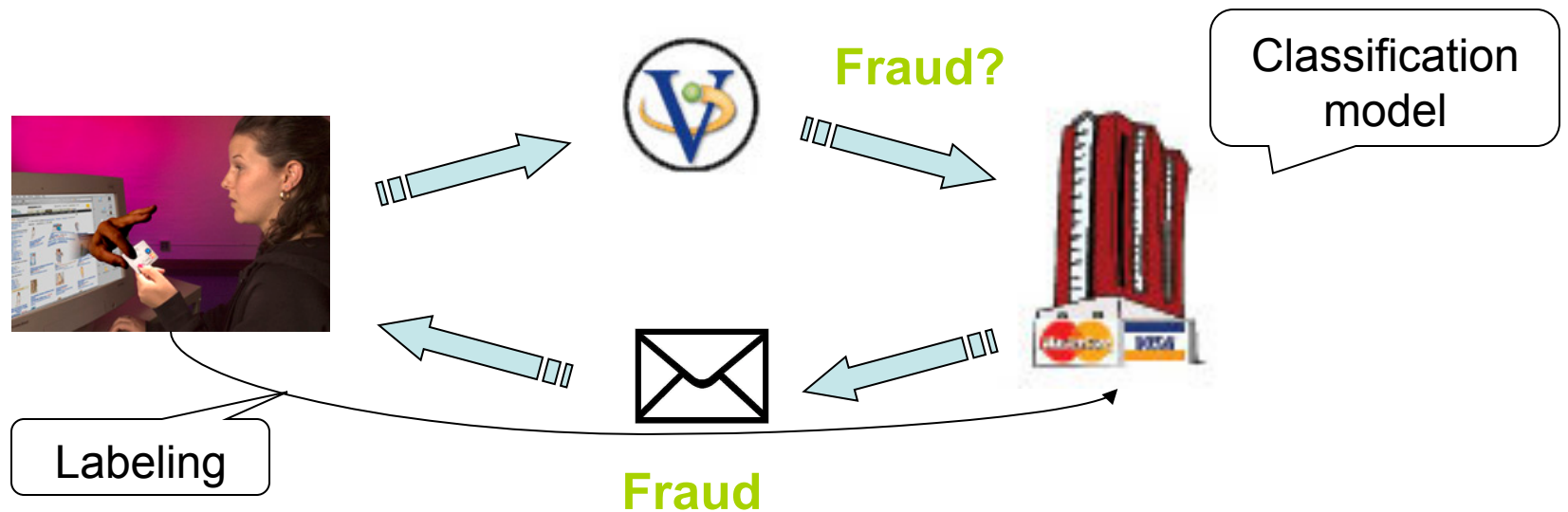
Outline

- An overview of ensemble methods
 - Motivations
 - Tutorial overview
- Supervised ensemble
- Unsupervised ensemble
- Semi-supervised ensemble
 - Multi-view learning
 - Consensus maximization among supervised and unsupervised models
- Applications
 - Transfer learning, stream classification, anomaly detection

Stream Classification*

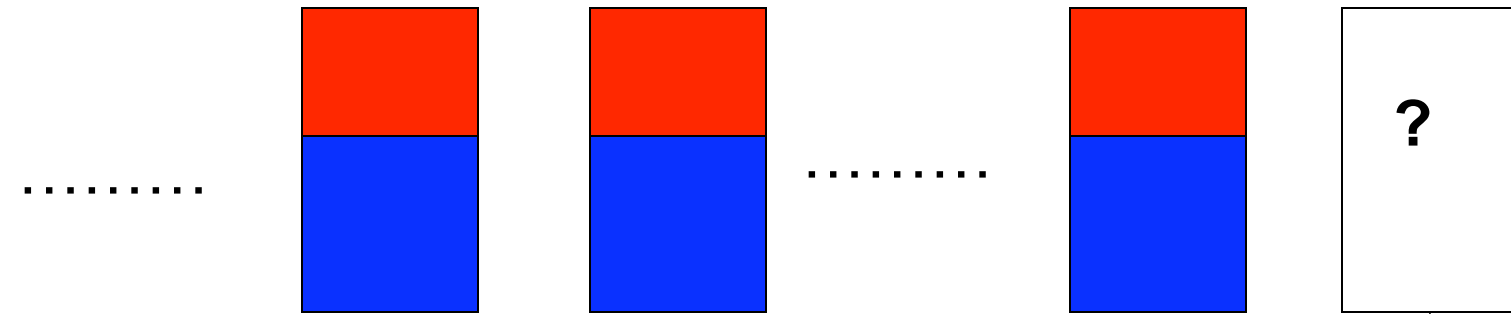
- **Process**

- Construct a classification model based on past records
- Use the model to predict labels for new data
- Help decision making

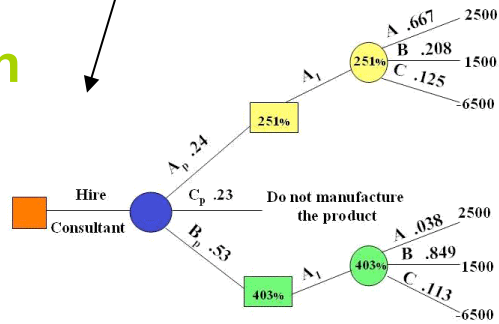


*[GFH07]

Framework



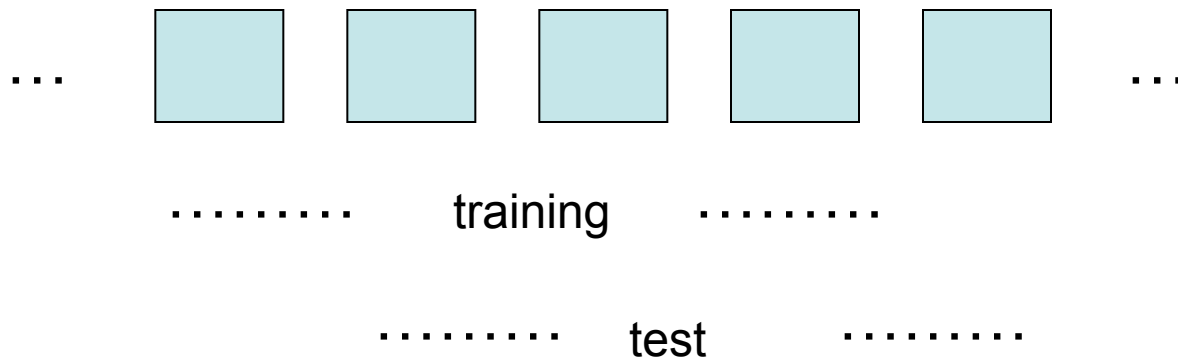
Classification Model



Predict

Existing Stream Mining Methods

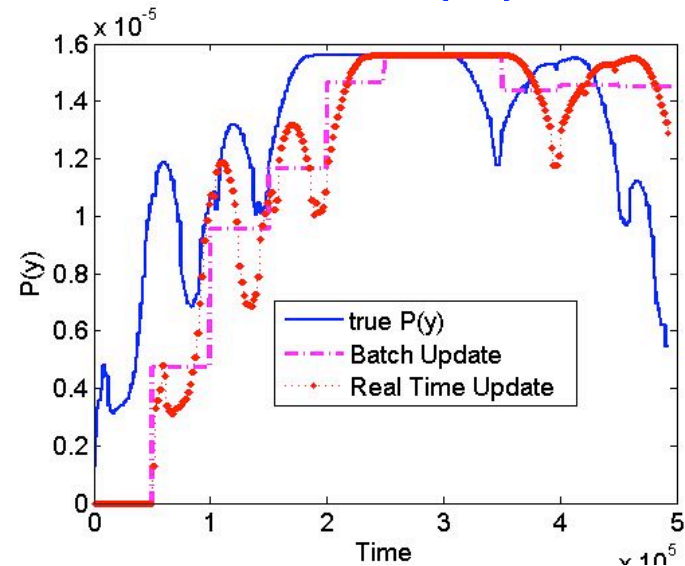
- **Shared distribution assumption**
 - Training and test data are from the same distribution $P(x,y)$ x-feature vector, y-class label
 - Validity of existing work relies on the shared distribution assumption
- **Difference from traditional learning**
 - Both distributions evolve



Evolving Distributions (1)

- An example of stream data

- KDDCUP'99 Intrusion Detection Data
- $P(y)$ evolves

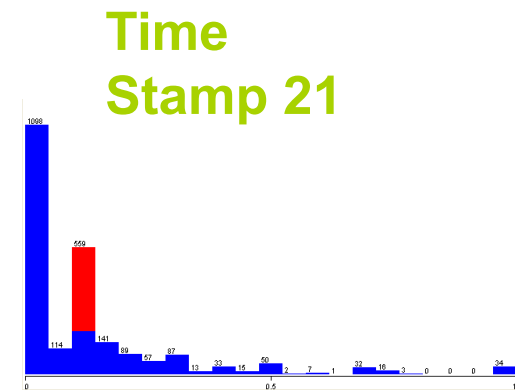
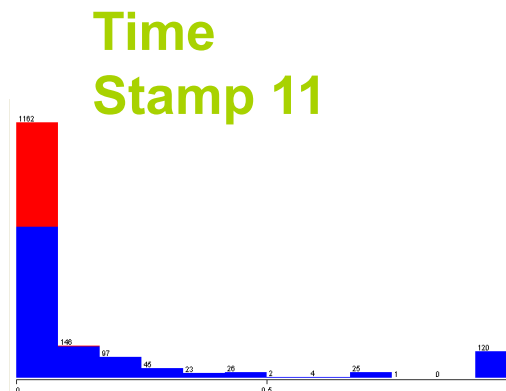
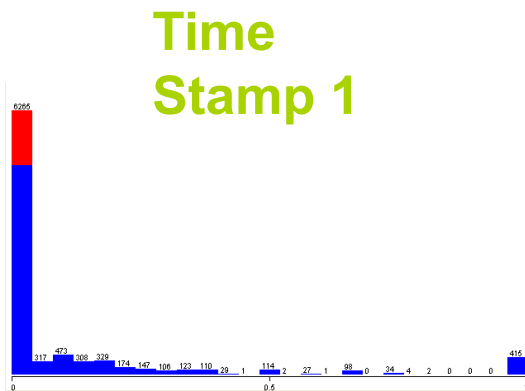


- Shift or delay inevitable

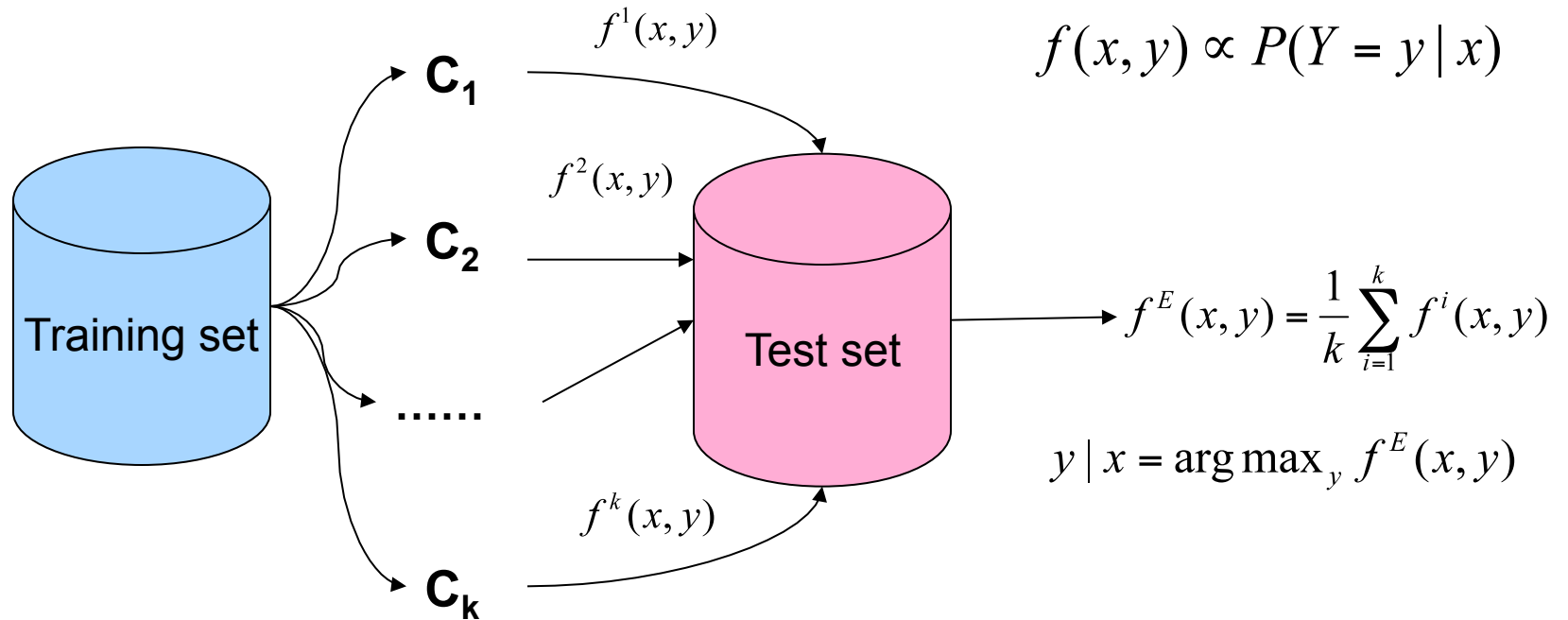
- The future data could be different from current data
- Matching the current distribution to fit the future one is a wrong way
- The shared distribution assumption is inappropriate

Evolving Distributions (2)

- Changes in $P(y)$
 - $P(y) \propto P(x,y)=P(y|x)P(x)$
 - The change in $P(y)$ is attributed to changes in $P(y|x)$ and $P(x)$



Ensemble Method



Simple Voting(SV)

$$f^i(x, y) = \begin{cases} 1 & \text{model } i \text{ predicts } y \\ 0 & \text{otherwise} \end{cases}$$

Averaging Probability(AP)

$$f^i(x, y) = \text{probability of predicting } y \text{ for model } i$$

Why it works?

- **Ensemble**

- Reduce variance caused by single models
- Is more robust than single models when the distribution is evolving

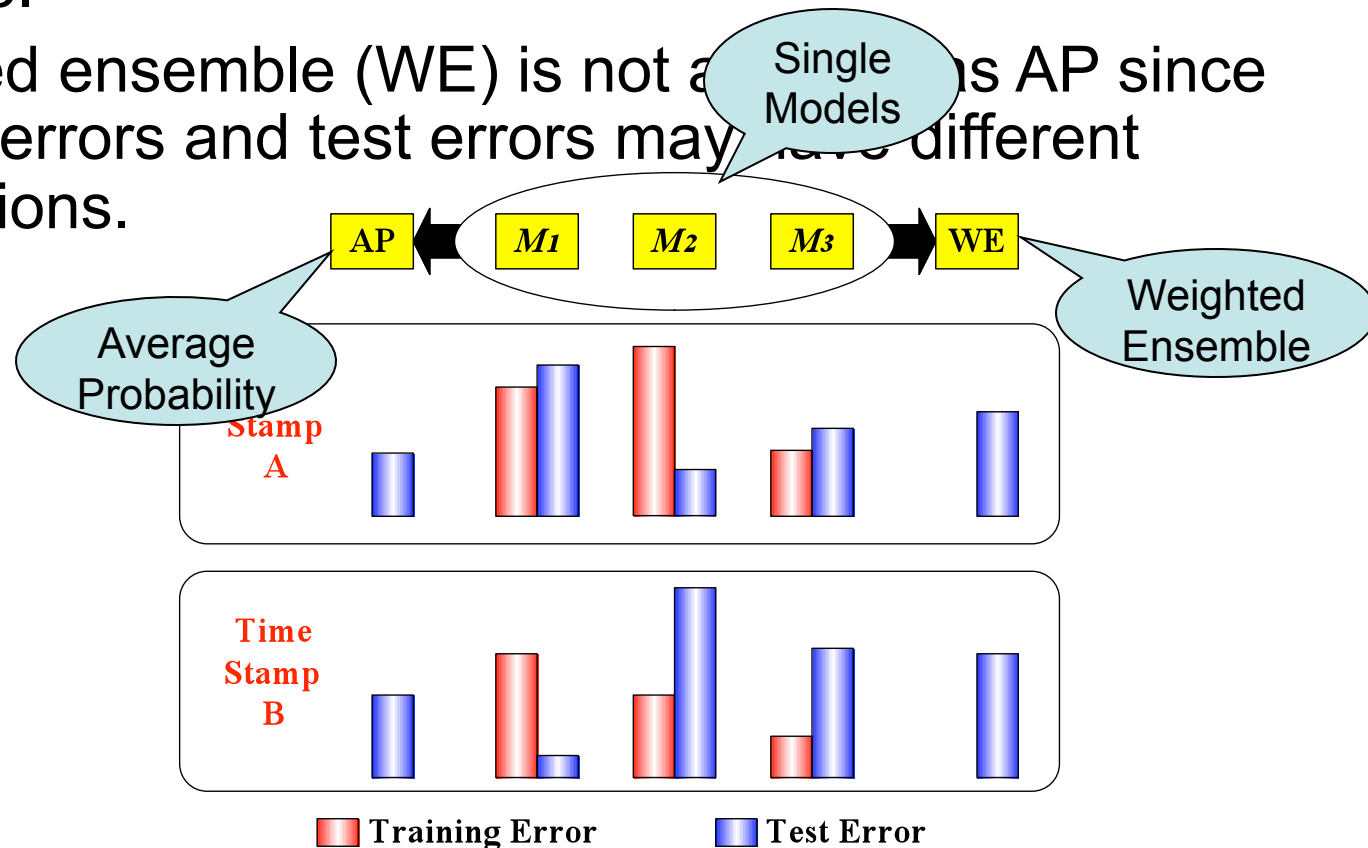
- **Simple averaging**

- Simple averaging: uniform weights $w_i=1/k$
- Weighted ensemble: non-uniform weights
 - w_i is inversely proportional to the training errors
- w_i should reflect $P(M)$, the probability of model M after observing the data
- $P(M)$ is changing and we could never estimate the true $P(M)$ and when and how it changes
- Uniform weights could minimize the expected distance between $P(M)$ and weight vector

$$f^E(x, y) = \sum_{i=1}^k w_i f^i(x, y)$$

An illustration

- Single models (M1, M2, M3) have huge variance.
- Simple averaging ensemble (AP) is more stable and accurate.
- Weighted ensemble (WE) is not as good as AP since training errors and test errors may have different distributions.



Experiments

- **Set up**

- Data streams with chunks T_1, T_2, \dots, T_N
- Use T_i as the training set to classify T_{i+1}

- **Measures**

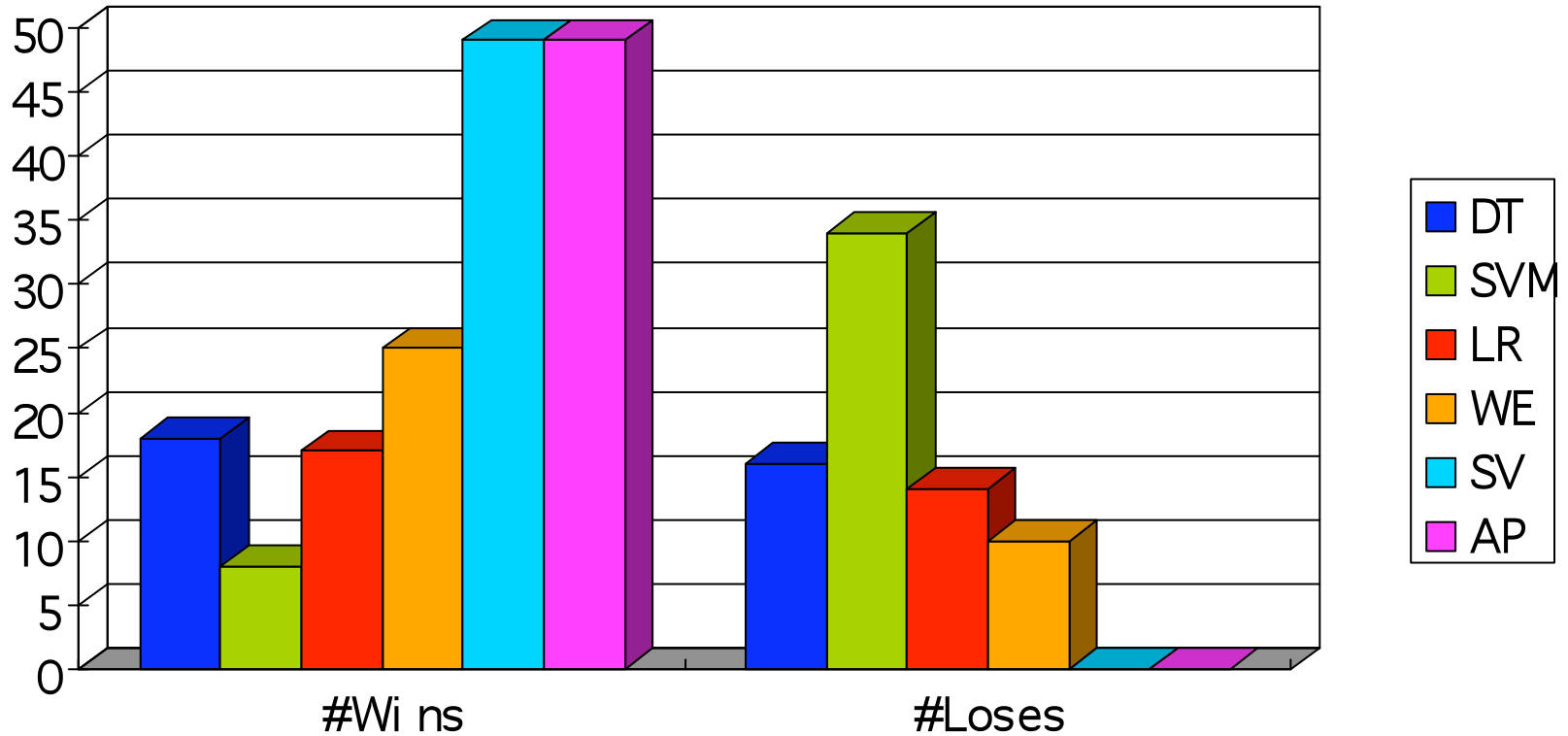
- Mean Squared Error, Accuracy
- Number of Wins, Number of Loses
- Normalized Accuracy, MSE

$$h(A, T) = h(A, T) / \max_A (h(A, T))$$

- **Methods**

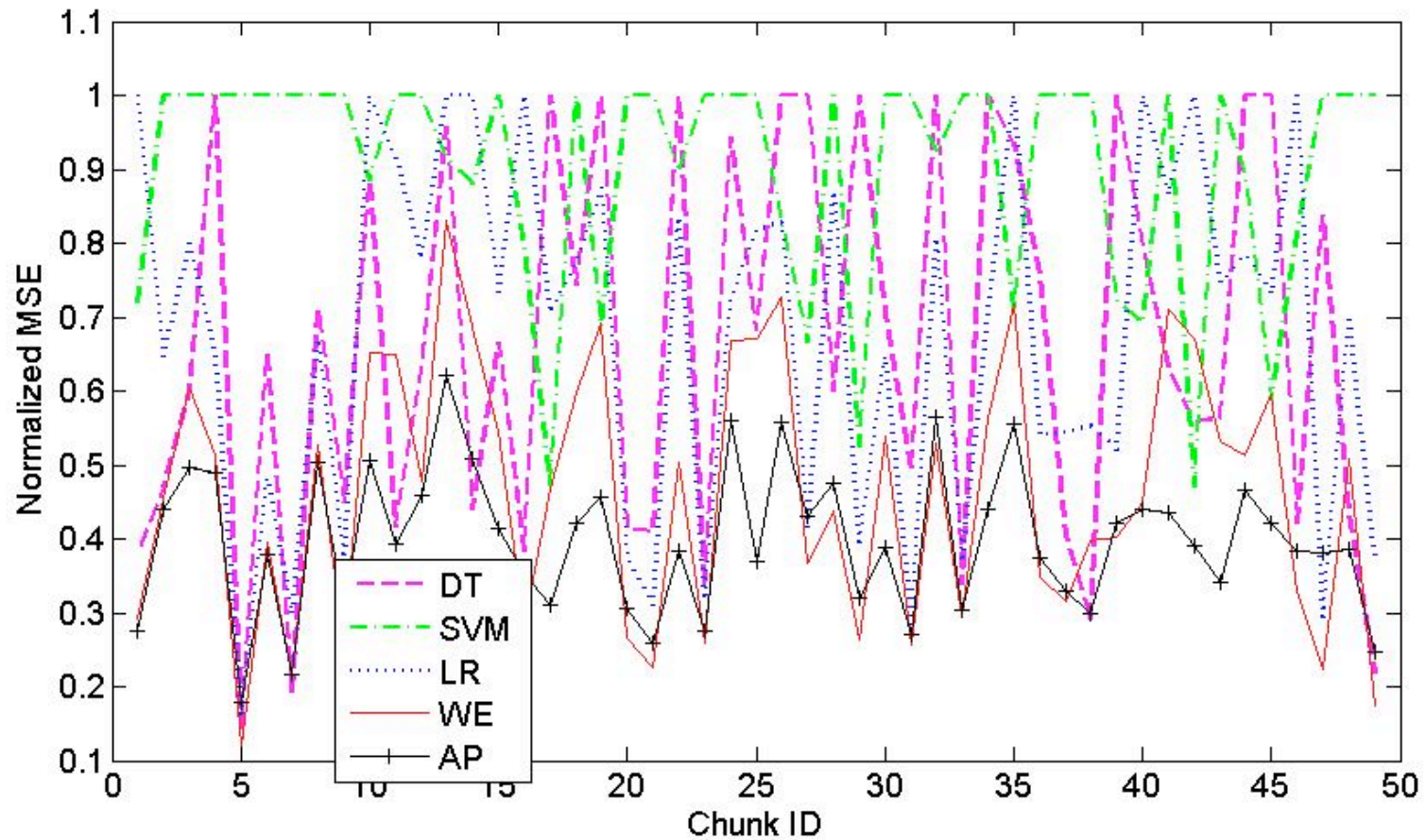
- Single models: Decision tree (DT), SVM, Logistic Regression (LR)
- Weighted ensemble: weights reflect the accuracy on training set (WE)
- **Simple ensemble: voting (SV) or probability averaging (AP)**

Experimental Results (1)



Comparison on Intrusion Data Set

Experimental Results (2)



Mean Squared Error Comparison

Outline

- An overview of ensemble methods
 - Motivations
 - Tutorial overview
- Supervised ensemble
- Unsupervised ensemble
- Semi-supervised ensemble
 - Multi-view learning
 - Consensus maximization among supervised and unsupervised models
- Applications
 - Transfer learning, stream classification, anomaly detection

Combination of Anomaly Detectors

- Simple rules (or atomic rules) are relatively easy to craft.
- Problem:
 - there can be way too many simple rules
 - each rule can have high false alarm or FP rate
- Challenge: can we find their non-trivial combination that significantly improve accuracy?

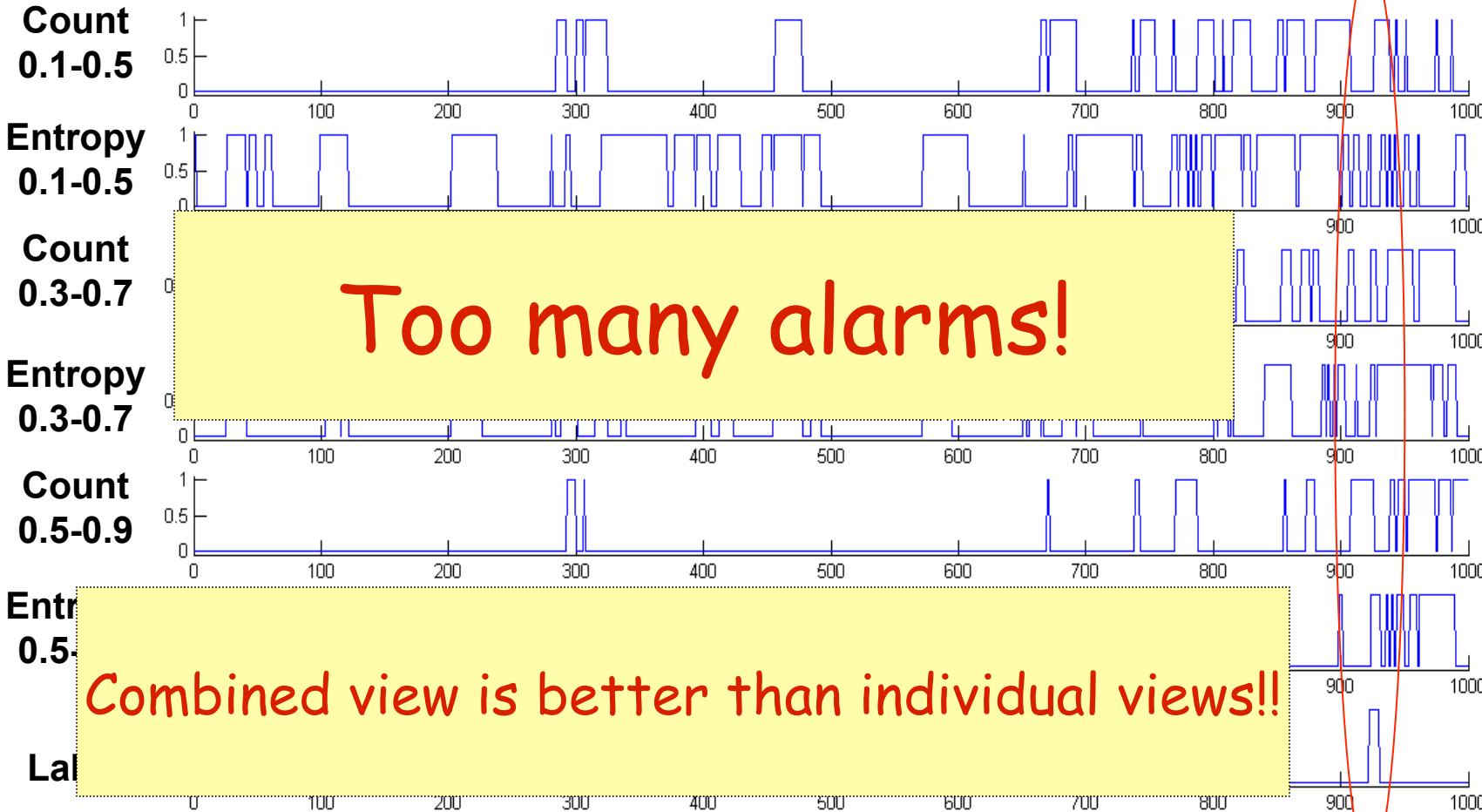
Atomic Anomaly Detectors

Anomaly?

	A_1	A_2	A_{k-1}	A_k
Record 1	Y	N	N	N
Record 2	N	Y	Y	N
Record 3	Y	N	N	N
Record 4	Y	Y	N	Y
Record 5	N	N	Y	Y
Record 6	N	N	N	N
Record 7	N	N	N	N

.....

Why We Need Combine Detectors?



Combining Detectors

- **is non-trivial**
 - We aim at finding a consolidated solution without any knowledge of the true anomalies (**unsupervised**)
 - We don't know which atomic rules are better and which are worse
 - There could be bad base detectors so that majority voting cannot work

How to Combine Atomic Detectors?

- **Basic Assumption:**
 - Base detectors are better than random guessing and systematic flip.
- **Principles**
 - Consensus represents the best we can get from the atomic rules
 - Solution most consistent with atomic detectors
 - Atomic rules should be weighted according to their detection performance
 - We should rank the records according to their probability of being an anomaly
- **Algorithm**
 - Reach consensus among multiple atomic anomaly detectors in an unsupervised way
 - or semi-supervised if we have limited supervision (known botnet site)
 - and incremental in a streaming environment
 - Automatically derive weights of atomic rules and records

Conclusions

- **Ensemble**
 - Combining independent, diversified models improves accuracy
 - No matter in supervised, unsupervised, or semi-supervised scenarios, ensemble methods have demonstrated their strengths
 - Base models are combined by learning from labeled data or by their consensus
- **Beyond accuracy improvements**
 - Information explosion motivates multiple source learning
 - Various learning packages available
 - Combine the complementary predictive powers of multiple models
 - Distributed computing, privacy-preserving applications

Thanks!

- Any questions?

Slides and more references available at <http://ews.uiuc.edu/~jinggao3/sdm10ensemble.htm>

Tutorial on Ensemble of Classifiers

- *Survey of Boosting from an Optimization Perspective*. Manfred K. Warmuth and S.V.N. Vishwanathan. ICML'09, Montreal, Canada, June 2009.
- *Theory and Applications of Boosting*. Robert Schapire. NIPS'07, Vancouver, Canada, December 2007.
- *From Trees to Forests and Rule Sets--A Unified Overview of Ensemble Methods*. Giovanni Seni and John Elder. KDD'07, San Jose, CA, August 2007.

References

- [AUL08] M. Amini, N. Usunier, and F. Laviolette. A transductive bound for the voted classifier with an application to semi-supervised learning. In *Advances in Neural Information Processing Systems 21*, 2008.
- [BBY04] M. Balcan, A. Blum, and K. Yang. Co-training and expansion: Towards bridging theory and practice. In *Advances in Neural Information Processing Systems 17*, 2004.
- [BBM07] A. Banerjee, S. Basu, and S. Merugu. Multi-way clustering on relation graphs. In *Proc. 2007 SIAM Int. Conf. Data Mining (SDM'07)*, 2007.
- [BaKo04] E. Bauer and R. Kohavi. An empirical comparison of voting classification algorithms: Bagging, boosting, and variants. *Machine Learning*, 36:105-139, 2004.
- [BEM05] R. Bekkerman, R. El-Yaniv, and A. McCallum. Multi-way distributional clustering via pairwise interactions. In *Proc. 2005 Int. Conf. Machine Learning (ICML'05)*, pages 41-48, 2005.
- [BDH05] P. N. Bennett, S. T. Dumais, and E. Horvitz. The combination of text classifiers using reliability indicators. *Information Retrieval*, 8(1):67-100, 2005.
- [BiSc04] S. Bickel and T. Scheffer. Multi-view clustering. In *Proc. 2004 Int. Conf. Data Mining (ICDM'04)*, pages 19-26, 2004.
- [BIMi98] A. Blum and T. Mitchell. Combining labeled and unlabeled data with co-training. *Proceedings of the Workshop on Computational Learning Theory*, pages 92-100, 1998.
- [BGS+08] P. Brazdil, C. Giraud-Carrier, C. Soares, and R. Vilalta. *Metalearning: Applications to Data Mining*. Springer, 2008.
- [BBS05] Ulf Brefeld, Christoph Büscher, and Tobias Scheffer. Multi-view discriminative sequential learning. In *Proc. European Conf. Machine Learning (ECML'05)*, pages 60-71, 2005.
- [Breiman96] L. Breiman. Bagging predictors. *Machine Learning*, 26:123-140, 1996.
- [Breiman01] L. Breiman. Random forests. *Machine Learning*, 45(1):5-32, 2001.
- [Caruana97] R. Caruana. Multitask learning. *Machine Learning*, 28(1):41-75, 1997.

References

- [CoSi99] M. Collins and Y. Singer. Unsupervised models for named entity classification. In Proc. 1999 Conf. Empirical Methods in Natural Language Processing (EMNLP'99), 1999.
- [CKW08] K. Crammer, M. Kearns, and J. Wortman. Learning from multiple sources. Journal of Machine Learning Research, 9:1757-1774, 2008.
- [DYX+07] W. Dai, Q. Yang, G.-R. Xue, and Y. Yu. Boosting for transfer learning. In Proc. 2007 Int. Conf. Machine Learning (ICML'07), pages 193-200, 2007.
- [DLM01] S. Dasgupta, M. Littman, and D. McAllester. PAC Generalization Bounds for Co-training. In Advances in Neural Information Processing Systems 14, 2001.
- [DaFa06] I. Davidson and W. Fan. When efficient model averaging out-performs boosting and bagging. In Proc. 2006 European Conf. Principles and Practice of Knowledge Discovery in Databases (PKDD'06), pages 478-486, 2006.
- [DMM03] I. S. Dhillon, S. Mallela, and D. S. Modha. Information-theoretic co-clustering. In Proc. 2003 ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining (KDD'03), pages 89-98, 2003.
- [Dietterich00] T. Dietterich. Ensemble methods in machine learning. In Proc. 2000 Int. Workshop Multiple Classifier Systems, pages 1-15, 2000.
- [DWH01] E. Dimitriadou, A. Weingessel, and K. Homik. Voting-merging: an ensemble method for clustering. In Proc. 2001 Int. Conf. Artificial Neural Networks (ICANN'01), pages 217-224, 2001.
- [DoAI09] C. Domeniconi and M. Al-Razgan. Weighted cluster ensembles: Methods and analysis. ACM Transactions on Knowledge Discovery from Data (TKDD), 2(4):1-40, 2009.
- [Domingos00] P. Domingos. Bayesian averaging of classifiers and the overfitting problem. In Proc. 2000 Int. Conf. Machine Learning (ICML'00), pages 223-230, 2000.
- [DHS01] R. O. Duda, P. E. Hart, and D. G. Stork. Pattern Classification. John Wiley & Sons, second edition, 2001.

References

- [DzZe02] S. Dzeroski and B. Zenko. Is combining classifiers better than selecting the best one. In Proc. 2002 Int. Conf. Machine Learning (ICML'02), pages 123-130, 2002.
- [DuFr03] S. Dudoit and J. Fridlyand. Bagging to improve the accuracy of a clustering procedure. *Bioinformatics*, 19(9): 1090-1099, 2003.
- [FaDa07] W. Fan and I. Davidson. On sample selection bias and its efficient correction via model averaging and unlabeled examples. In Proc. 2007 SIAM Int. Conf. Data Mining (SDM'07), 2007.
- [FGM+05] W. Fan, E. Greengrass, J. McCloskey, P. S. Yu, and K. Drummey. Effective estimation of posterior probabilities: Explaining the accuracy of randomized decision tree approaches. In Proc. 2005 Int. Conf. Data Mining (ICDM'05), pages 154-161, 2005.
- [FHM+05] J. Farquhar, D. Hardoon, H. Meng, J. Shawe-taylor, and S. Szedmak. Two view learning: SVM-2K, theory and practice. In *Advances in Neural Information Processing Systems* 18, 2005.
- [FeBr04] X. Z. Fern and C. E. Brodley. Solving cluster ensemble problems by bipartite graph partitioning. In Proc. 2004 Int. Conf. Machine Learning (ICML'04), pages 281-288, 2004.
- [FeLi08] X. Z. Fern and W. Lin. Cluster ensemble selection. In Proc. 2008 SIAM Int. Conf. Data Mining (SDM'08), 2008.
- [FiSk03] V. Filkov and S. Skiena. Integrating microarray data by consensus clustering. In Proc. 2003 Int. Conf. Tools with Artificial Intelligence, pages 418-426, 2003.
- [FrJa02] A. Fred and A. Jain. Data Clustering using evidence accumulation. In Proc. 2002 Int. Conf. Pattern Recognition (ICPR'02), 2002.
- [FrSc97] Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119-139, 1997.

References

- [FrPo08] J. H. Friedman and B. E. Popescu. Predictive learning via rule ensembles. *Annals of Applied Statistics*, 3(2):916-954, 2008.
- [GGB+08] K. Ganchev, J. Graca, J. Blitzer, and B. Taskar. Multi-view learning over structured and non-identical outputs. In *Proc. 2008 Conf. Uncertainty in Artificial Intelligence (UAI'08)*, pages 204-211, 2008.
- [GFH07] J. Gao, W. Fan, and J. Han. On appropriate assumptions to mine data streams: Analysis and practice. In *Proc. 2007 Int. Conf. Data Mining (ICDM'07)*, pages 143-152, 2007.
- [GFJ+08] J. Gao, W. Fan, J. Jiang, and J. Han. Knowledge transfer via multiple model local structure mapping. In *Proc. 2008 ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining (KDD'08)*, pages 283-291, 2008.
- [GFS+09] J. Gao, W. Fan, Y. Sun, and J. Han. Heterogeneous source consensus learning via decision propagation and negotiation. In *Proc. 2009 ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining (KDD'09)*, pages 339-347, 2009.
- [GLF+09] J. Gao, F. Liang, W. Fan, Y. Sun, and J. Han. Graph-based consensus maximization among multiple supervised and unsupervised models. In *Advances in Neural Information Processing Systems 22*, 2009.
- [GSI+09] R. Ghaemi, M. Sulaiman, H. Ibrahim, and N. Mutspha. A survey: clustering ensembles techniques. *World Academy of Science, Engineering and Technology* 50, 2009.
- [GeTa07] L. Getoor and B. Taskar. *Introduction to statistical relational learning*. MIT Press, 2007.
- [GMT07] A. Gionis, H. Mannila, and P. Tsaparas. Clustering aggregation. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 1(1), 2007.
- [GVB04] C. Giraud-Carrier, R. Vilalta, and P. Brazdil. Introduction to the special issue on meta-learning. *Machine Learning*, 54(3):187-193, 2004.

References

- [GoFi08] A. Goder and V. Filkov. Consensus clustering algorithms: comparison and refinement. In Proc. 2008 Workshop on Algorithm Engineering and Experiments (ALENEX'08), pages 109-117, 2008.
- [GoZh00] S. Goldman and Y. Zhou. Enhancing supervised learning with unlabeled data. In Proc. 2000 Int. Conf. Machine Learning (ICML'00), pages 327-334, 2000.
- [HKT06] S. T. Hadjitodorov, L. I. Kuncheva, and L. P. Todorova. Moderate diversity for better cluster ensembles. *Information Fusion*, 7(3):264-275, 2006.
- [HaKa06] J. Han and M. Kamber. *Data mining: concepts and techniques*. Morgan Kaufmann, second edition, 2006.
- [HTF09] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, second edition, 2009.
- [HMR+99] J. Hoeting, D. Madigan, A. Raftery, and C. Volinsky. Bayesian model averaging: a tutorial. *Statistical Science*, 14:382-417, 1999.
- [JJN+91] R. Jacobs, M. Jordan, S. Nowlan, and G. Hinton. Adaptive mixtures of local experts. *Neural Computation*, 3(1):79-87, 1991.
- [KoMa] J. Kolter and M. Maloof. Using additive expert ensembles to cope with concept drift. In Proc. 2005 Int. Conf. Machine Learning (ICML'05), pages 449-456, 2005.
- [KuWh03] L. I. Kuncheva and C. J. Whitaker. Measures of diversity in classifier ensembles and their relationship with the ensemble accuracy. *Machine Learning*, 51(2):181-207, 2003.
- [Leskes05] B. Leskes. The Value of Agreement, a New Boosting Algorithm. In 2005 Proc. Conf. Learning Theory (COLT'05), pages 95-110, 2005.
- [LiDi08] T. Li and C. Ding. Weighted consensus clustering. In Proc. 2008 SIAM Int. Conf. Data Mining (SDM'08), 2008.

References

- [LDJ07] T. Li, C. Ding, and M. Jordan. Solving consensus and semi-supervised clustering problems using nonnegative matrix factorization. In Proc. 2007 Int. Conf. Data Mining (ICDM'07), pages 577-582, 2007.
- [LiOg05] T. Li and M. Ogihara. Semisupervised learning from different information sources. Knowledge and Information Systems, 7(3):289-309, 2005.
- [LiYa06] C. X. Ling and Q. Yang. Discovering classification from data of multiple sources. Data Mining and Knowledge Discovery, 12(2-3):181-201, 2006.
- [LZY05] B. Long, Z. Zhang, and P. S. Yu. Combining multiple clusterings by soft correspondence. In Proc. 2005 Int. Conf. Data Mining (ICDM'05), pages 282-289, 2005.
- [LZX+08] P. Luo, F. Zhuang, H. Xiong, Y. Xiong, and Q. He. Transfer learning from multiple source domains via consensus regularization. In Proc. 2008 Int. Conf. Information and Knowledge Management (CIKM'08), pages 103-112, 2008.
- [MTP04] B. Minaei-Bidgoli, A. Topchy, and W. Punch: A comparison of resampling methods for clustering ensembles. In Proc. 2004 Int. Conf. Artificial Intelligence (ICAI'04), pages 939-945, 2004.
- [NiGh00] K. Nigam and R. Ghani. Analyzing the effectiveness and applicability of co-training. In Proc. 2000 Int. Conf. Information and Knowledge Management (CIKM'00), pages 86-93, 2000.
- [OkVa08] O. Okun and G. Valentini. Supervised and Unsupervised Ensemble Methods and their Applications. Springer, 2008.
- [Polikar06] R. Polikar. Ensemble based systems in decision making. IEEE Circuits and Systems Magazine, 6(3):21-45, 2006.
- [PrSc08] C. Preisach and L. Schmidt-Thieme. Ensembles of relational classifiers. Knowledge and Information Systems, 14(3):249-272, 2008.

References

- [PTJ05] W. Punch, A. Topchy, and A. K. Jain. Clustering ensembles: Models of consensus and weak partitions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(12): 1866-1881, 2005.
- [PuGh08] K. Punera and J. Ghosh. Consensus based ensembles of soft clusterings. *Applied Artificial Intelligence*, 22(7-8): 780-810, 2008.
- [RoKa07] D. M. Roy and L. P. Kaelbling. Efficient bayesian task-level transfer learning. In *Proc. 2007 Int. Joint Conf. Artificial Intelligence (IJCAI'07)*, pages 2599-2604, 2007.
- [SNB05] V. Sindhwani, P. Niyogi, and M. Belkin. A co-regularization approach to semi-supervised learning with multiple views. In *Proc. 2005 ICML workshop on Learning with Multiple Views*, 2005.
- [SMP+07] V. Singh, L. Mukherjee, J. Peng, and J. Xu. Ensemble clustering using semidefinite programming. In *Advances in Neural Information Processing Systems 20*, 2007.
- [StGh03] A. Strehl and J. Ghosh. Cluster ensembles --a knowledge reuse framework for combining multiple partitions. *Journal of Machine Learning Research*, 3:583-617, 2003.
- [TLJ+04] A. Topchy, M. Law, A. Jain, and A. Fred. Analysis of consensus partition in cluster ensemble. In *Proc. 2004 Int. Conf. Data Mining (ICDM'04)*, pages 225-232, 2004.
- [TuGh96] K. Tumer and J. Ghosh. Analysis of decision boundaries in linearly combined neural classifiers. *Pattern Recognition*, 29, 1996.
- [ViDr02] R. Vilalta and Y. Drissi. A perspective view and survey of meta-learning. *Artificial Intelligence Review*, 18(2):77-95, 2002.
- [WFY+03] H. Wang, W. Fan, P. Yu, and J. Han. Mining concept-drifting data streams using ensemble classifiers. In *Proc. 2003 ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining (KDD'03)*, pages 226-235, 2003.
- [WSB09] H. Wang, H. Shan, and A. Banerjee. Bayesian cluster ensembles. In *Proc. 2009 SIAM Int. Conf. Data Mining (SDM'09)*, 2009.

References

- [Wolpert92] D. H. Wolpert. Stacked generalization. *Neural Networks*, 5:241-259, 1992.
- [WWL09] F. Wang, X. Wang, and T. Li. Generalized Cluster aggregation. In *Proc. 2009 Int. Joint Conf. Artificial Intelligence (IJCAI'09)*, pages 1279-1284, 2009.
- [ZGY05] J. Zhang, Z. Ghahramani, and Y. Yang. Learning multiple related tasks using latent independent component. In *Advances in Neural Information Processing Systems 18*, 2005.
- [ZFY+06] K. Zhang, W. Fan, X. Yuan, I. Davidson, and X. Li. Forecasting skewed biased stochastic ozone days: Analyses and solutions. In *Proc. 2006 Int. Conf. Data Mining (ICDM'06)*, pages 753-764, 2006.
- [ZZY07] Z. Zhou, D. Zhan, and Q. Yang. Semi-Supervised Learning with Very Few Labeled Training Examples. In *Proc. 2007 Conf. Artificial Intelligence (AAAI'07)*, pages 675-680, 2007.