



11th Conference on Intelligent Systems Theory
and Applications.
Mohammedia, 19-20 October 2016

More flexible representations in Data Mining

(An overview to Multiple Instance Learning)

Sebastián Ventura Soto

Knowledge Discovery and Intelligent Systems Research Group

University of Cordoba

Knowledge Discovery and Intelligent Systems

KDISlab
KNOWLEDGE DISCOVERY AND INTELLIGENT SYSTEMS

• CORDOBA



Knowledge Discovery and Intelligent Systems

<http://www.uco.es/grupos/kdis>



KDIS Lab
KNOWLEDGE DISCOVERY AND INTELLIGENT SYSTEMS

- Head of the group:
 - Prof. Sebastian Ventura
- Components:
 - 10 PhD researchers
 - 10 PhD students
- Facts and Figures:
 - 100+ journal papers
 - 200+ conference papers
 - 2 authored and 10 edited books
 - 8 PhD dissertations

Knowledge Discovery and Intelligent Systems

Research Interest



Metaheuristics

- Evolutionary Computation
- Ant Colony Optimization
- Other metaheuristics

Applications

- Search-Based Software Engineering
- Problem solving with Metaheuristics

New Methods in ML / DM

- Association Rule Mining
- Classification
- Regression
- Multiple-instance learning
- Multi-label learning
- Multi-view learning

Scalability in DM

- GPU-based methods
- Big Data Mining (Hadoop and Spark)

Applications

- Educational Data Mining
- Clinical Data Mining

Contents

- Flexible representations in Data Mining
- Multiple-Instance Learning
- Applications of Multiple Instance Classification
- Multiple Instance Classification Algorithms
 - Instance Space Paradigm
 - Bag Space Paradigm
 - Embedded Space Paradigm
- Other Multiple Instance Paradigms

Flexible Representations in Data Mining

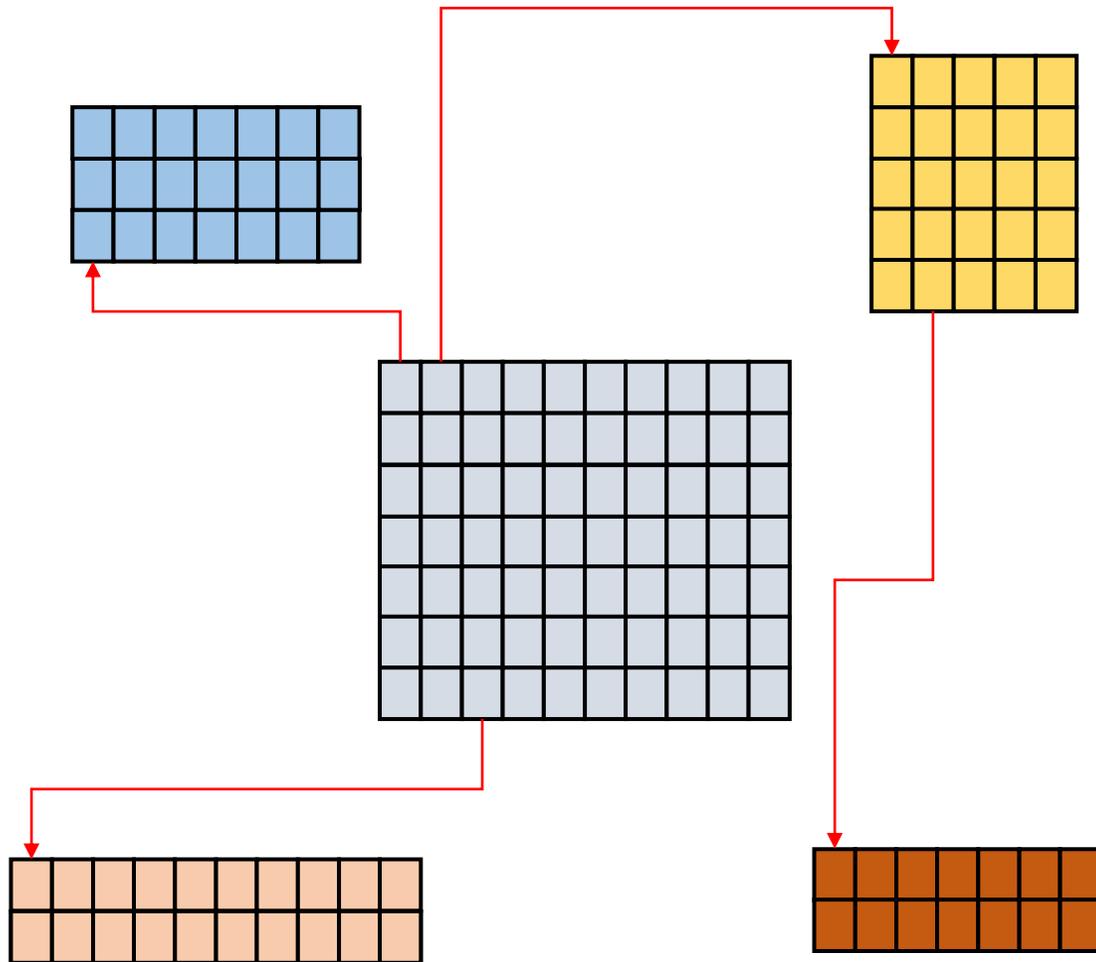
Classical data representation in Machine Learning and Data Mining

	Attribute 1	Attribute 2	Attribute 3							Attribute M
Instance 1										
Instance 2										
Instance 3										
Instance N										

- Table has been the data representation in classical machine learning and data mining
- Both supervised and unsupervised tasks work well with this kind of representation
- There are problems that do not fit to these representations

Alternative representations

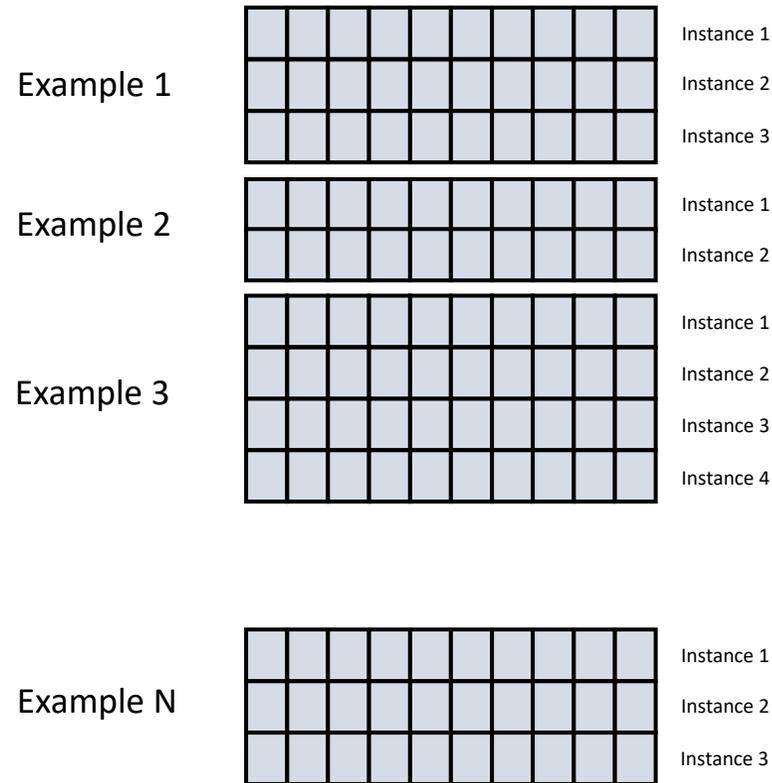
Relational data



- Relational models are a popular data representation, but not very used in ML/DM
- Usually, multi-table data are converted in single-table, conventional data to apply classical ML/DM algorithms
- **Relational learning models** deal with multi-table data representations directly

Alternative representations

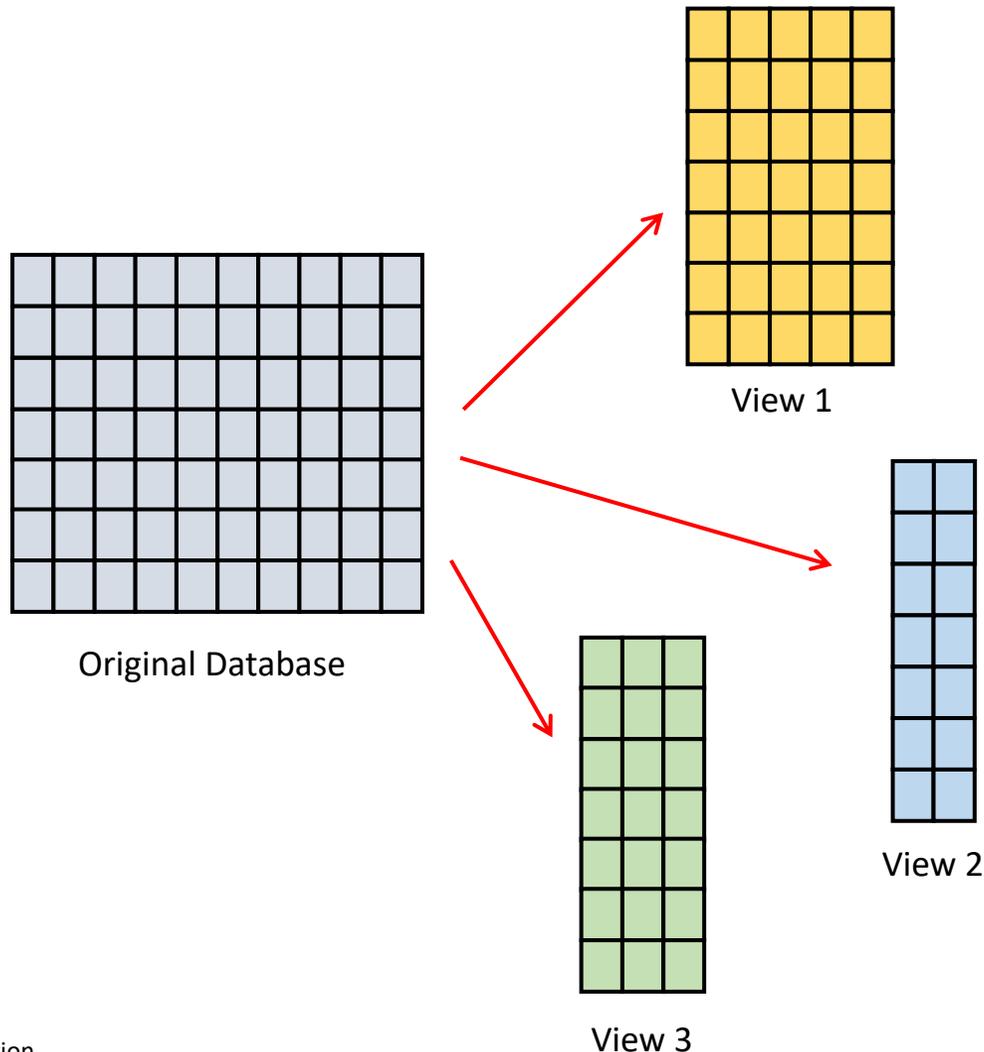
Multi-instance data



- In multiple instance data an object is represented by a variable number of input vectors (*a bag of vectors*)
- Each vector represents a different view or perspective of the object
- **Multiple instance learning** methods deal with this data representation with any kind of preprocessing

Alternative representations

Multiple views



- Sometimes a dataset can be splitted in several *views*, each one related with an attribute subset.
- Generally, attributes in a view keep some relationship.
- Building learning models with each subset is easier that building an overall model, but these models have a partial view of the learned concept
- **Multi-view learning** methods perform a join learning process by combining these partial models in a global one.

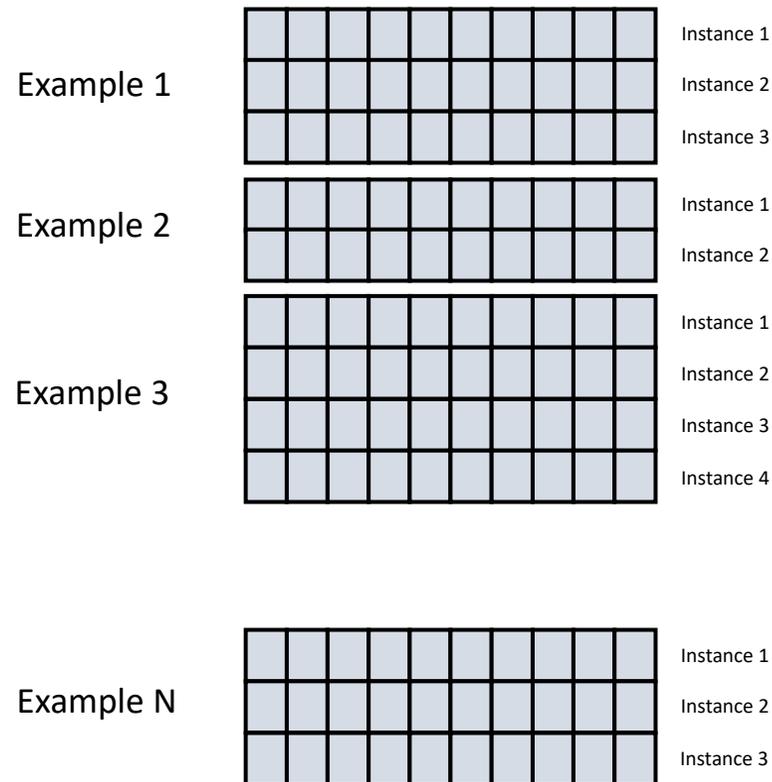
Alternative representations

Flexible representations in data mining

- All these alternative data representations are told **flexible** because can adapt a new problems in a more flexible way that clasical data tables.
- Furthermore, these flexible data representations can combine giving new learning paradigms like:
 - Multi-instance and multi-label classification
 - Multi-view multi-instance learning
 - Multi-view multi-label learning
 - ...
- The rest of this speech will be devoted to the **multiple instance learning** paradigm, due to his recent popularity and the number of applications it has exhibit in the last years.

Multiple Instance Learning

Multiple Instance Learning



- The term Multiple Instance Learning refers in general to solve learning tasks using multiple instances as the input data representation.
- This paradigm appeared at the end of the nineties (paper of Dietterich et al, 1997) and it has become very popular since then.
- There are multiple applications of multiple instance learning in multiple fields:
 - Drug activity prediction
 - Image Classification
 - Text classification
 - ...

Multi-instance Learning Problems/Paradigms

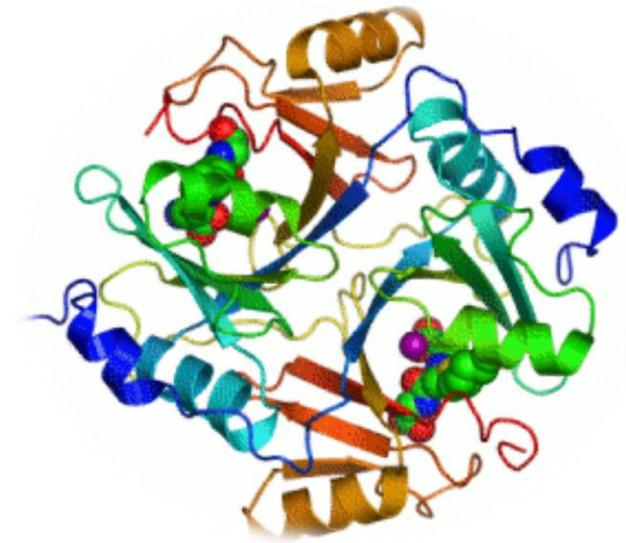
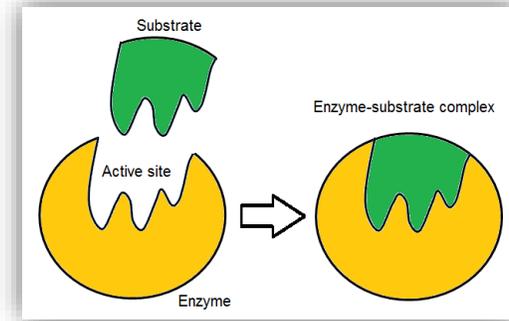
- *Multi-instance Classification*. The objective is to predict unseen bag labels:
 - *Binary Classification*: Binary label
 - *Multiple Classification*: Nominal (non-binary) label
 - *Multi-Label Classification*: Multiple labels (MI-MLL)
- *Multi-instance Regression*. The objective is to predict the continuous label of unseen bags.
- *Multi-instance Clustering*. Grouping similar objects of bags in clusters.
- *Multi-instance Association Rule Mining*. Finding association patterns from bags.

This presentation is focussed on Multi-instance Binary Classification

Applications of Multiple Instance Classification

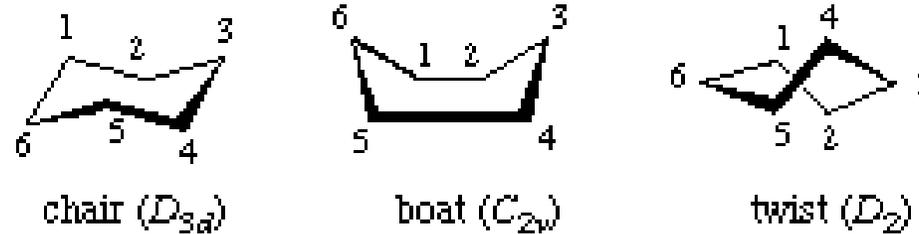
Prediction of Pharmacological Activity

- The first paper on MIL (Dietterich et al., 1997) was motivated by the problem of determining whether a drug molecule exhibits a given activity.
- A molecule presents a given pharmacological activity when it is able to bind with an enzyme or protein. This is only possible if the molecule has certain spatial properties (*key-lock mechanism*).



Prediction of Pharmacological Activity(II)

- A molecule may adopt a wide range of shapes or *conformations*, due to the rotation of its bonds.



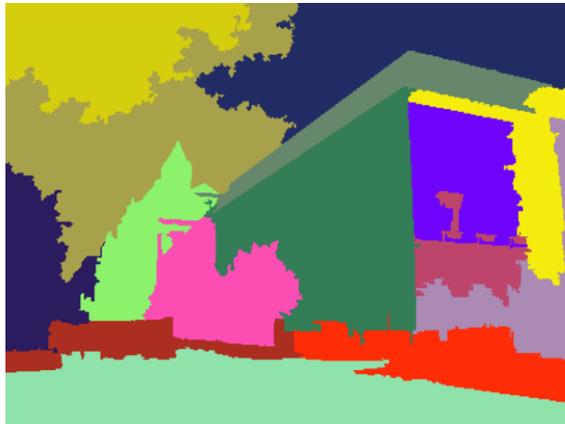
- If a conformation can bind/connect to a pharmacological activity center, the whole molecule exhibits the activity under research. Otherwise, the molecule does not exhibit this activity.
- In Dietterich's paper, the property under study was *musk*. Substances with this property are employed in the manufacture of perfumes and other cosmetic products.

Prediction of Pharmacological Activity(III)

- This problem can be represented by multi-instances in a very natural way:
 - Each molecule is a bag
 - Each conformation is an instance
- Dietterich et al. studied two different datasets:
 - Musk-1: 92 molecules (47 positive y 45 negative), 476 instances and 166 attributes.
 - Musk-2: 102 molecules (39 positive y 63 negative), 6598 instances y 166 attributes.
- There exist other benchmarks related to the pharmacological activity prediction problem. For instance, in mutagenesis dataset, the property under study is the ability to produce mutations.
 - Mutagenesis 1: 188 molecules, 10,468 instances, 7 attributes
 - Mutagenesis 2: 42 molecules, 2,132 instances, 7 attributes

These and other bechmark datasets can be found at <http://www.uco.es/grupos/kdis/mil/dataset.html>

Content-based Image Classification and Retrieval



- The key to the success of image retrieval and classification is the ability to identify the intended target object(s) in images.
- This problem is complicated when the image contains multiple and possibly heterogeneous objects.
- This problem can fit into the MIL setting well:
 - Each image itself is considered as a bag.
 - A region or segment in an image is considered to be an instance.

Text categorization

- Andrews et al (2002) use MIL to categorize documents taken from the TREC9 dataset (a benchmark in text categorization problems).
- They divide each document in 50 word length overlapping sets (authors do not specify what this overlapping is like). In this case, each 50 word set represents an instance and the whole document is a training pattern (bag of instances).

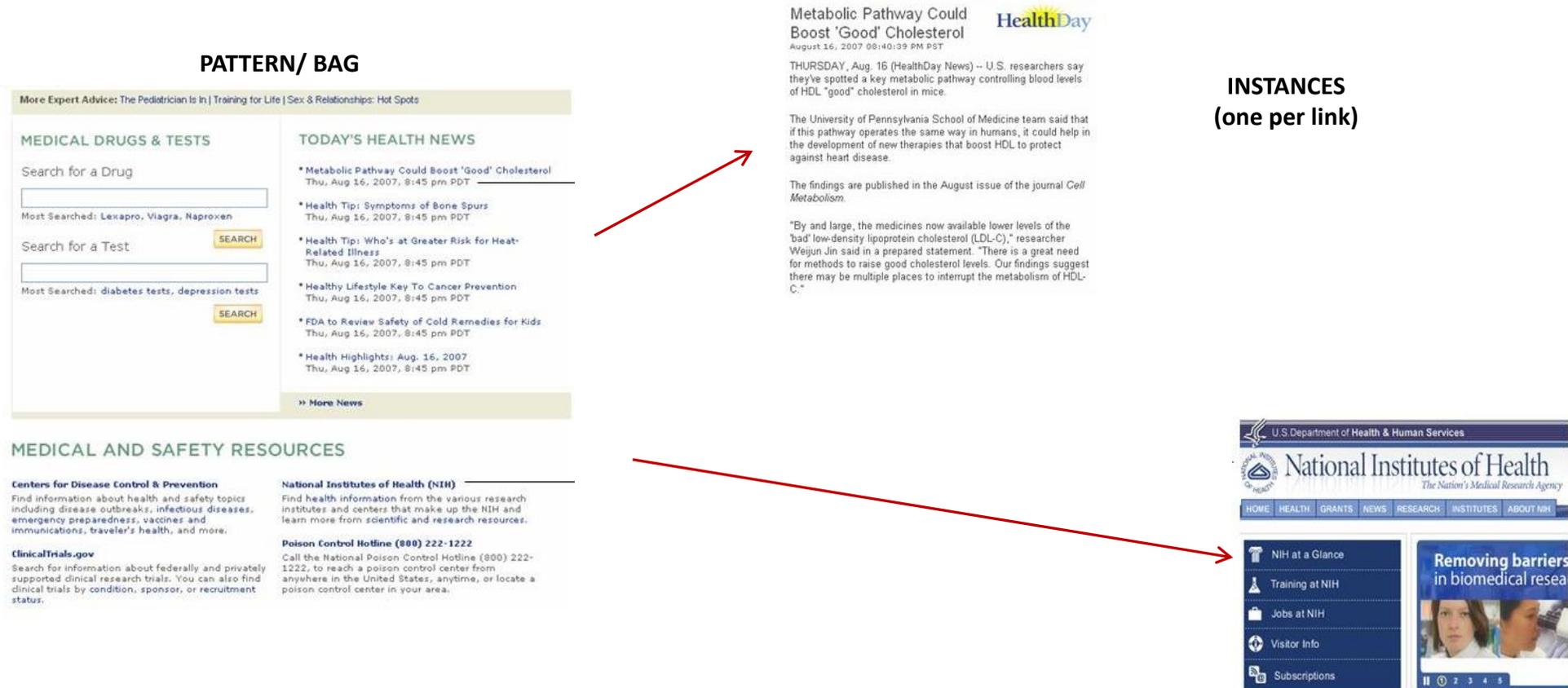
S. Andrews, I. Tsochantaridis & T. Hofmann. Support vector machines for Multiple Instance Learning. In *Advances in Neural Information Processing Systems (NIPS 15)*, pp 1-8, 2002

Web index recommendation



- Web index pages are pages that provide titles or brief summaries and leave the detailed presentation to their linked pages.
- The problem of recommending web index pages consists of determining what pages a given user is interested in.
- In general, if a web index page contain links that the user considers interesting, the user will be attracted to it.
- The problem is that we do not have information about links, but about the page as a whole.

Web index recommendation (II)



A web index page is represented as a bag, and each link in its web index page is represented by an instance.

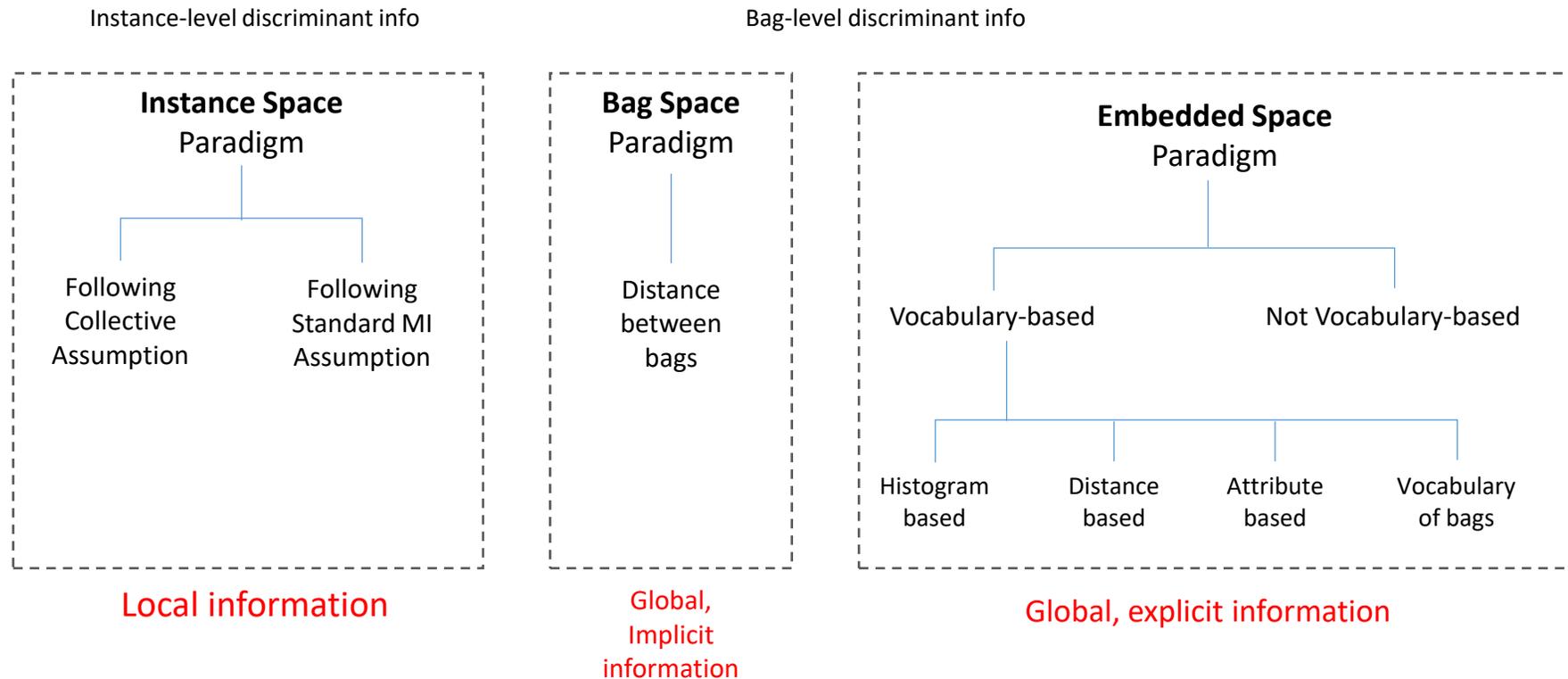
Detecting Fraudulent Users in Personal Banking



- The idea is to detect fraudulent use of credit cards from a transactions dataset
- For each user, we have one or more transactions defined by several attributes:
 - amount
 - time
 - transaction interval
 - service ID
 - merchant type
 - ...
- This problem can also be represented as a multi-instance where a bag represents all the transactions of a given user and the label will tell if the use has been fraudulent or not

Multiple-Instance Classification Algorithms

A Taxonomy For Multi-instance Classification Algorithms



J. Amores. Multiple Instance Classification: Review, taxonomy and comparative study. *Artificial Intelligence*, 201, 81–105 (2013)

Instance-space paradigm

Instance space paradigm

Introduction

- The idea is to infer an *instance based classifier* $f(\vec{x}) \in [0,1]$ from the training data.
- Bag level classification $F(X) \in [0,1]$ is constructed as an aggregation of instance-level responses:

$$F(X) = \frac{f(\vec{x}_1) \circ f(\vec{x}_2) \circ \dots \circ f(\vec{x}_N)}{Z}$$

where \circ represents an *aggregation operator* and Z is a normalization factor

- These methods have to solve the problem of how to infer an instance-level classifier **without** having access to a training set of labelled instances. To do this, some *hypothesis* has to be made about the relationship that exists between the label of the bags and the labels of the instances contained in the bags.
- There are two main hypothesis:
 - Standard or Dietterich's hypothesis
 - Collective hypothesis

Standard Hypothesis

- Every positive bag contains at least one positive instance, while in every negative bag all the instances are negatives.
- The methods following this hypothesis try to identify the type of instance that make the bag positive.
- There are several classical methods that follow this hypothesis:
 - Learning of Axis-Parallel Rectangles (APR)
 - Diverse Density
 - MI-SVM
 - Sparse MIL and Sparse Balanced MIL
 - Adaptations of Single Instance Learning (SIL) algorithms to MIL. According Z.-H. Zhou (2009), SIL algorithms can be adapted to MIL including the standard hypothesis in their development.
 - Decision Trees
 - Rule Based Learning
 - G3P-MI and MO G3P-MI

Learning of Axis Parallel Rectangles (APRs)

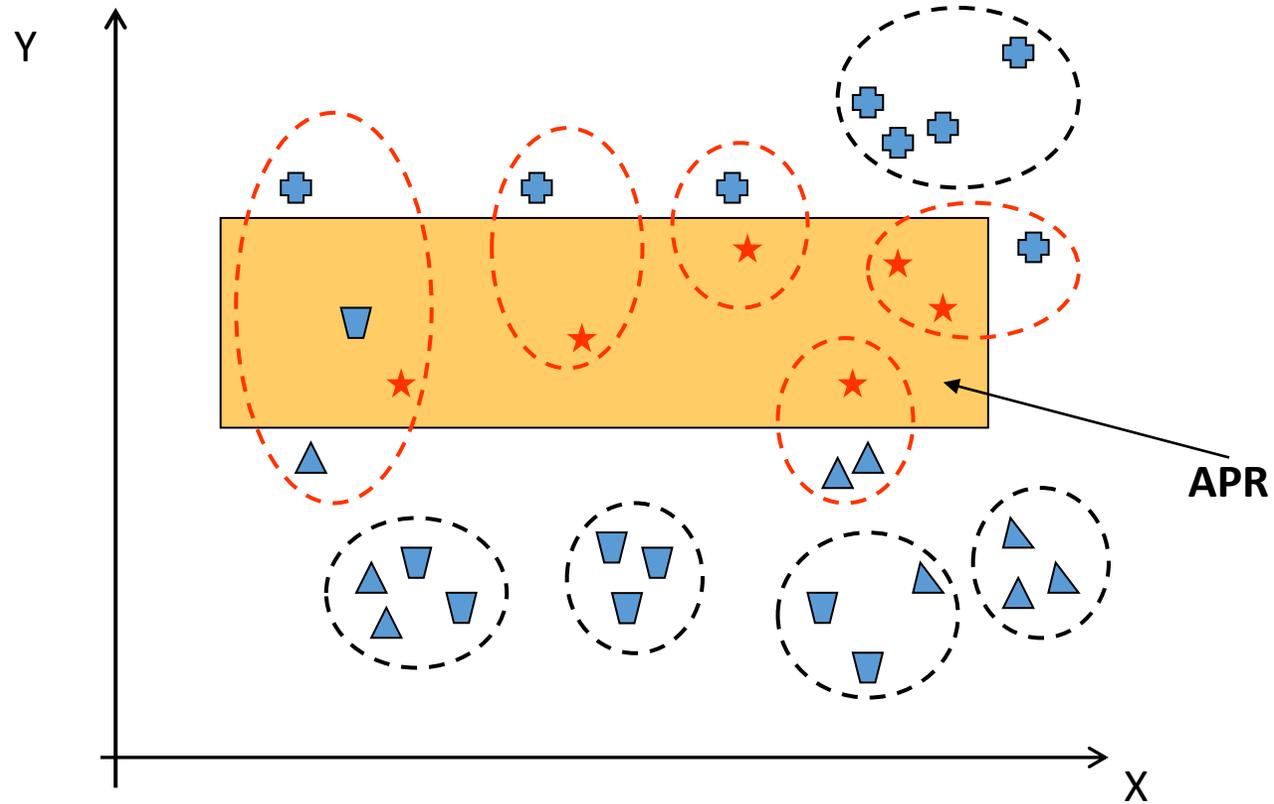
- The first solution to the multiple instance learning problem was proposed by Dietterich et al. 1997

T. G. Dietterich, R.H Lathrop & T. Lozano-Pérez. Solving the multiple instance problem with axis-parallel rectangles.
Artificial Intelligence 89:1-2 (1997), pp 31-71

- They propose representing the concept to be learned by axis parallel rectangles (APR) in the feature space. Intuitively, this APR should contain at least one instance from each positive example and meanwhile exclude all the instances from negative examples.

$$f(\vec{x}; \mathcal{R}) = \begin{cases} 1 & \text{if } \vec{x} \in \mathcal{R} \\ 0 & \text{otherwise} \end{cases} \quad F(X) = \max_{\vec{x} \in X} f(\vec{x})$$

Graphical Description of APR

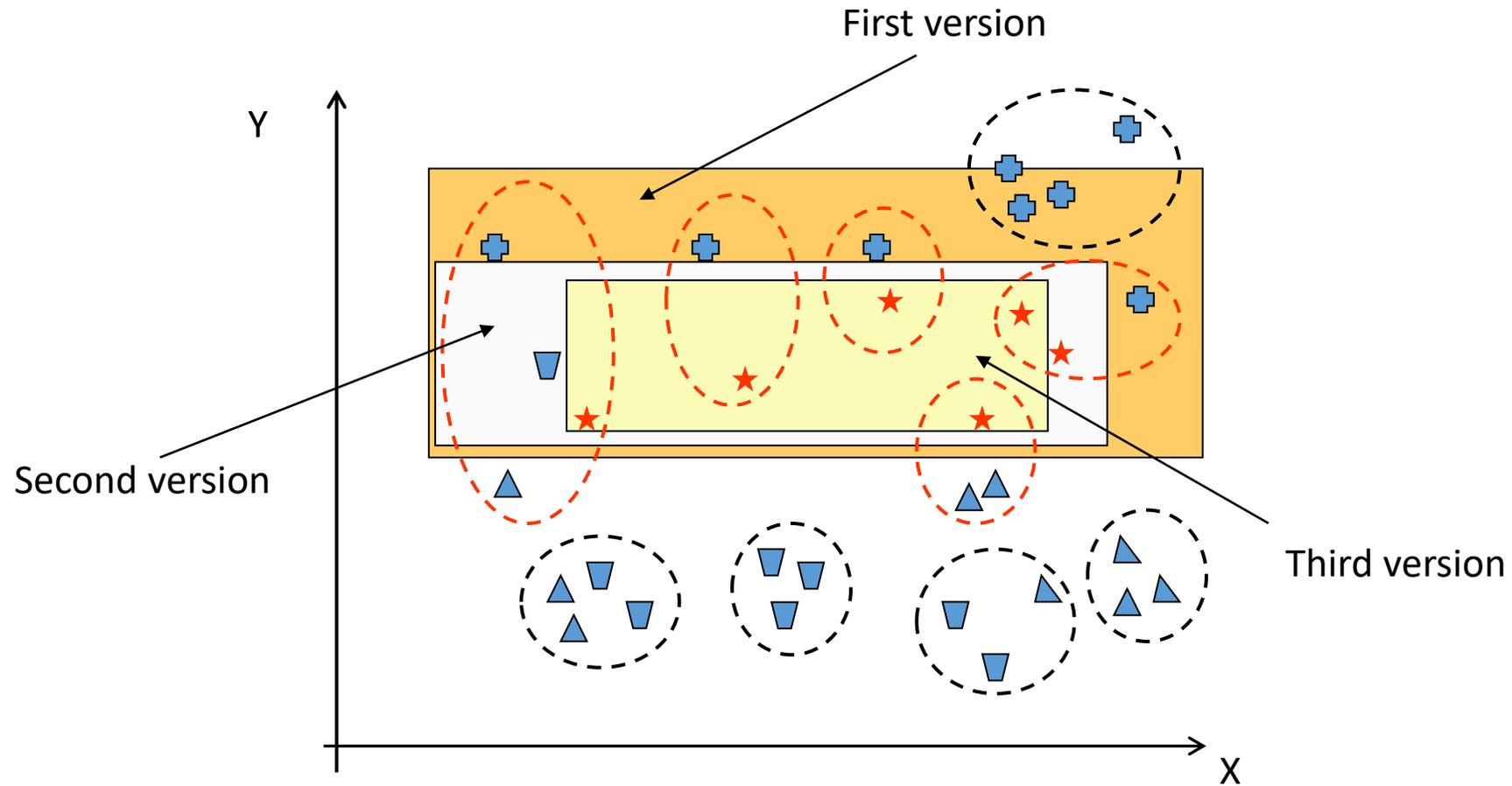


Variants of the APR Algorithm

- Dietterich's article considers three general designs for APR algorithms:
 - A noise-tolerant “standard” algorithm. The naive APR algorithm just forms the smallest APR that bounds the positive examples.
 - An “outside-in” algorithm. This algorithm is a variation on the “standard” algorithm. It constructs the smallest APR that bounds all of the positive examples and then shrinks this APR to exclude false positives.
 - An “inside-out” algorithm. This algorithm starts with a seed point in the feature space and “grows” a rectangle with the goal of finding the smallest rectangle that covers at least one instance of each positive example and no instances of any negative example.
- Authors apply these algorithms and several supervised learning algorithms (C4.5 and BP-ANN) to three datasets: Musk 1, Musk 2 and a synthetic dataset.
- In general, results showed that APR algorithms outperform supervised learning algorithms in this kind of problem.

Variants of APR

(Graphical Description)



Diverse Density

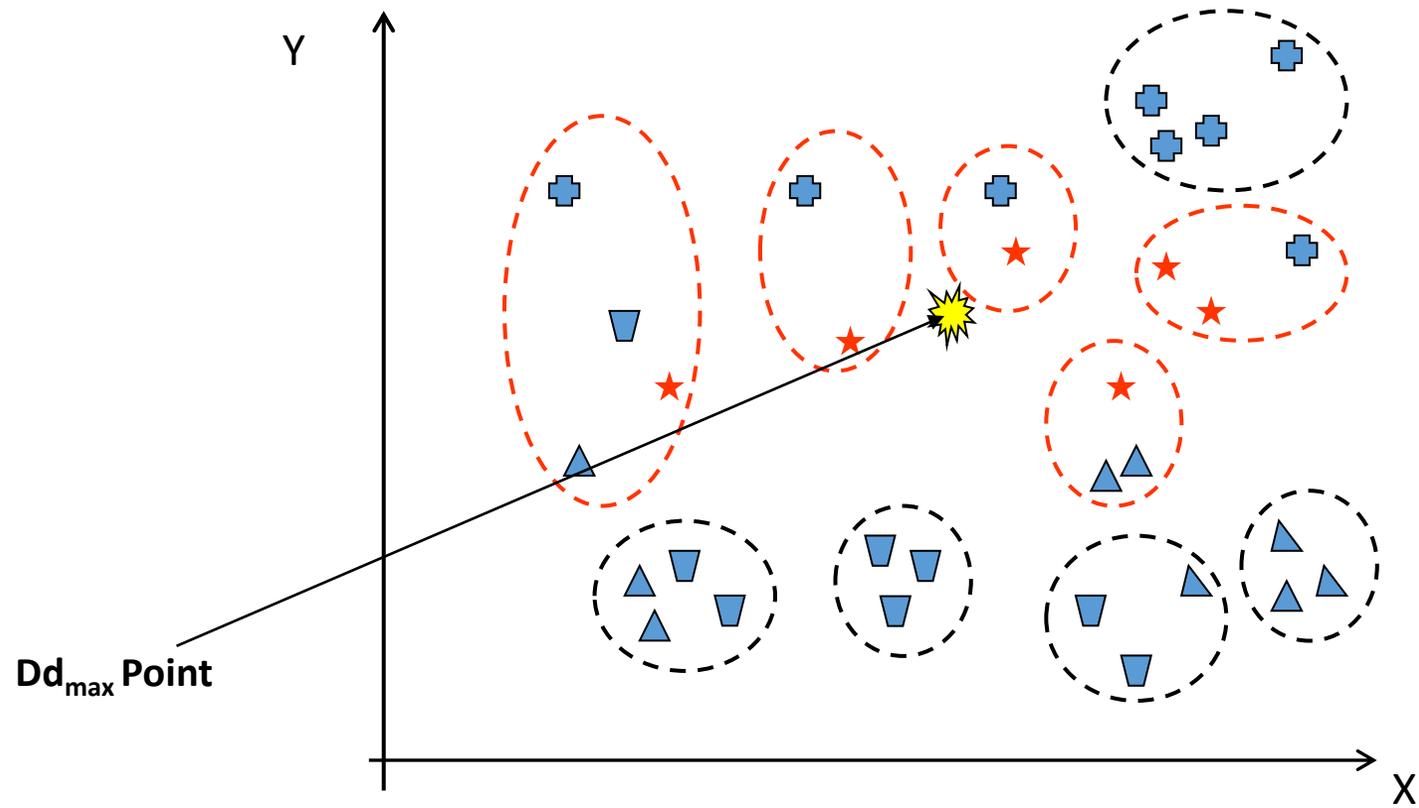
- One of the most popular learning algorithms in multi-instance learning is diverse density (DD), proposed by Maron & Lozano Pérez in 1998.

O. Maron & T. Lozano Pérez. A Framework for Multiple-Instance Learning *In Proc. of the 1997 Conference on Advances in Neural Information Processing Systems* (1998) pp 570-576.

- The main idea of the DD approach is to find a concept point in the feature space that is close to at least one instance from every positive example and meanwhile far away from instances in negative examples.
- The optimal concept point is defined as the one with the maximum diversity density, which is a measure of how many different positive bags have instances near the point, and how far the negative instances are away from that point.

Diverse Density

(Graphical Description)



Decision Trees and Rule Based Systems

- Y. Chevaleyre and J.D. Zucker adapted the algorithms ID3 (decision trees) and RIPPER (rule induction) to the multiple instance learning paradigm.

Y. Chevaleyre & J.-D. Zucker. Solving Multiple-Instance and Multiple-Part Learning Problems with Decision Trees and Rule Sets. Application to the Mutagenesis Problem. In E. Stroulia & S. Matwin (Eds): *AI 2001*, LNAI 2056, pp 204-214, 2001.

Y. Chevaleyre & J.-D. Zucker. A Framework for Learning Rules from Multiple Instance Data. In L. de Raedt & P. Flach (Eds.) *ECML 2001*, LNAI 2167, pp 49-60, 2001.

- These adaptations are based on adapting the concepts of entropy and information gain to the multi-instance context.

Decision Trees and Rule Based Systems (II)

$$Entropy(S) = - \frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

$$Entropy_{multi}(S) = - \frac{u(S)}{u(S)+v(S)} \log_2 \frac{u(S)}{u(S)+v(S)} - \frac{v(S)}{u(S)+v(S)} \log_2 \frac{v(S)}{u(S)+v(S)}$$

u = set of positive bags
v = set of negative bags

$$InfoGain(S, F) = Entropy(S) - \sum_{v \in Values(F)} \frac{p_v + n_v}{p+n} Entropy(S_v)$$

$$InfoGain_{multi}(S, F) = Entropy_{multi}(S) - \sum_{v \in Values(F)} \frac{u(S_v) + v(S_v)}{u(S) + v(S)} Entropy_{multi}(S_v)$$

For adapting the RIPPER, these authors adapt the concept of coverage to the context of multi-instance objects (bags).

Collective Hypothesis

- There are problems where the standard hypothesis does not yield good results.
- The collective hypothesis assumes that **all instances contributes equally to the bag's label**
- Collective hypothesis can also work well in problems like Musk, because all the instances inside a bag might contribute in certain way to concept associated to the bag.



Concept *beach* is associated with the appearance of both *sand* and *water*, not one of them only

Collective Hypothesis (II)

Algorithms based on the collective hypothesis operate as follows:

1. They use a training set where all the instances inherits the label of the bag where it lies.
2. Then, they train a supervised learning classifier $f(\vec{x})$ using this dataset.
3. Finally, they build the bag classifier $F(X)$ aggregating the instance level predictions.

The collective algorithms described in the bibliography only differ in the aggregation method used to build $F(X)$.

SIL and Wrapper MI Algorithms

SIL

- This is the simplest collective algorithm
- Uses the sum as aggregation rule

$$F(X) = \frac{1}{|X|} \sum_{\vec{x} \in X} f(\vec{x})$$

Wrapper MI

- Uses a weighted sum as aggregation method

$$F(X) = \frac{1}{|X|} \sum_{\vec{x} \in X} w(\vec{x}) f(\vec{x}) \quad w(\vec{x}) = \frac{S}{|X|}$$

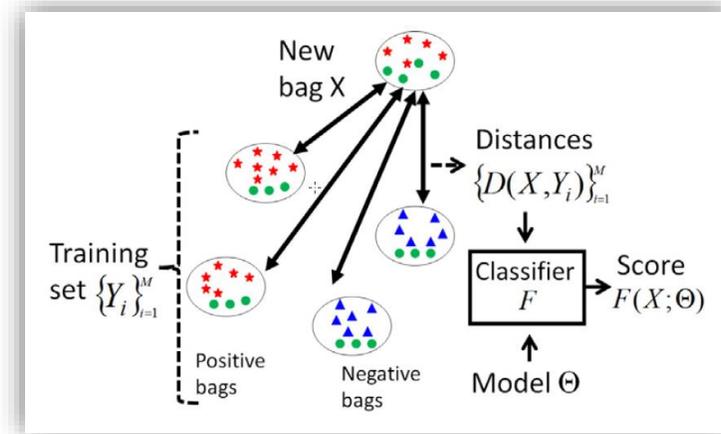
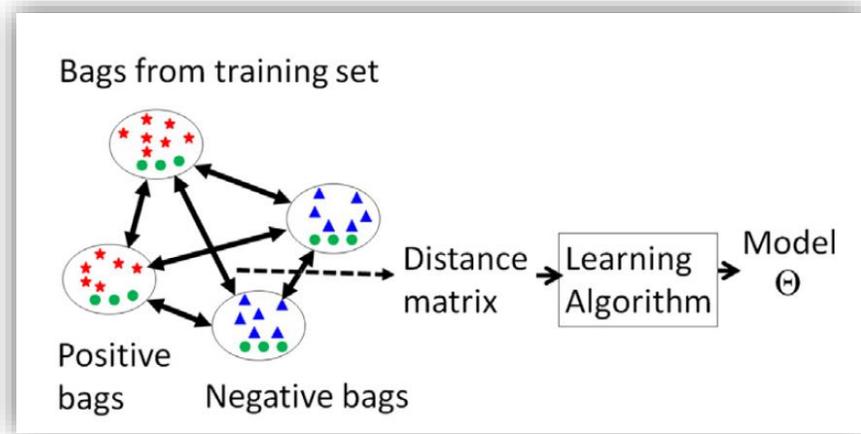
- Authors insist in the importance of these weights, as it makes the different bag of the training set have the same total weight

Bag-Space Paradigm

Mlc Algorithms

Introduction

- This paradigm treat bags as a whole, and the discriminat learning process is performed in the space of bags
- As bags are not vector entities, we have to define a *distance function* $D(X, Y)$ that compare 2 bags X and Y and plug this distance into a standard distance-based classifier like kNN or SVM.



Introduction

Distance commonly used

Minimal Hausdorff distance

$$D(X, Y) = \min_{\vec{x} \in X, \vec{y} \in Y} \|\vec{x} - \vec{y}\|$$

Earth Movers distance

$$D(X, Y) = \frac{\sum_i \sum_j w_{ij} \|\vec{x}_i - \vec{y}_j\|}{\sum_i \sum_j w_{ij}}$$

Chamfer distance

$$D(X, Y) = \frac{1}{|X|} \sum_{\vec{x} \in X} \min_{\vec{y} \in Y} \|\vec{x} - \vec{y}\| + \frac{1}{|Y|} \sum_{\vec{y} \in Y} \min_{\vec{x} \in X} \|\vec{x} - \vec{y}\|$$

Other systems use kernel functions that measure similarity instead of distance

k-NN Algorithms

- Wang y Zucker proposed a k-NN algorithm to MIC problems.

J. Wang & J.-D. Zucker Solving the multiple-instance problem: a lazy learning approach. *In Proc of 17th International Conference on Machine Learning (2000)*, pp 1119-1125

- These authors proposed using the Hausdorff distance as the bag-level distance metric
- The application of k-NN using this metric did not yield good results.

K	K nearest neighbors	# positive	#negative	Total
1	{P}	41	9	50
	{N}	6	36	42
2	{P,P}	41	3	44
	{P,N}	5	15	20
	{N,N}	1	27	28
3	{P,P,P}	40	2	42
	{P,P,N}	5	13	18
	{P,N,N}	2	9	11
	{N,N,N}	0	21	21

- Positive bags also contain negative instances, which attract the negative bags towards themselves.
- There are two ways to solve this problem:
 1. Giving more weight to negative objects.
 2. Defining new ways of combining neighbors to achieve the correct result.

Bayesian k-NN

- Conventional k-NN is based on the votation scheme, that can be represented by

$$\arg \max_{c \in \{\text{positive, negative}\}} \sum_{i=1}^k \delta(c, c_i)$$

- Bayesian k-NN proposes using the probabilities an object has of belonging to the c class, given k nearest neighbors

$$\arg \max_{c \in \{\text{positive, negative}\}} p(c / \{c_1, c_2, \dots, c_k\}) =$$
$$\arg \max_{c \in \{\text{positive, negative}\}} \frac{p(\{c_1, c_2, \dots, c_k\} | c) p(c)}{p(\{c_1, c_2, \dots, c_k\})}$$

$$\arg \max_{c \in \{\text{positive, negative}\}} p(\{c_1, c_2, \dots, c_k\} | c) p(c)$$

These probabilities are calculated from the real distribution of data

Citation k-NN

- This algorithm proposes using, besides the nearest neighbors (called in this case references), the objects that this pattern considers to be a nearest neighbor (called citers).
- Citation k-NN uses R references and C citers and, to decide whether an object is positive or negative, it calculates the following values

$$p = R_p + C_p$$

$$n = R_n + C_n$$

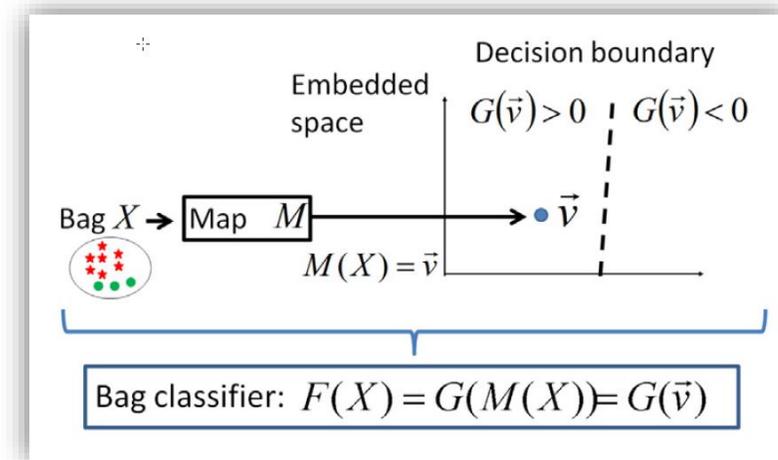
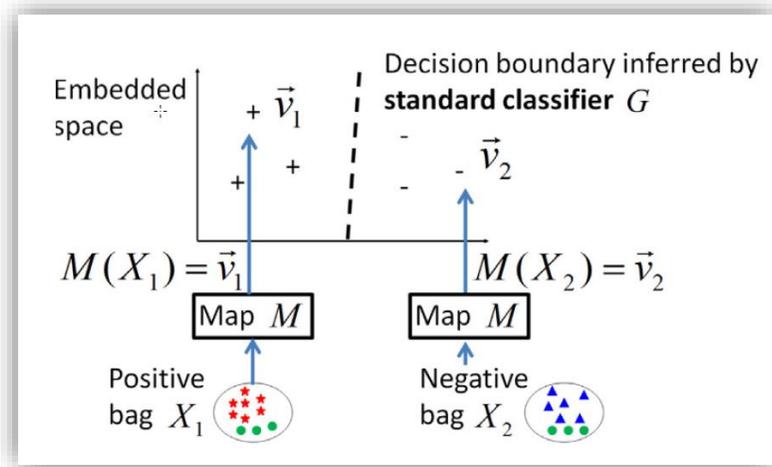
where R_p , C_p , R_n y C_n are, respectively, the number of references and citers that has a positive and negative label.

If $p > n$, the object is labeled as positive. Otherwise, the object is labeled as negative

Embedded-Space Paradigm

Introduction

- Embedded space methods use the information of bags to perform the discriminative process, like bag space methods.
- Instead of using a distance to compare bags, they define a *mapping function* $\mu: X \rightarrow \vec{v}$ from the bag X to a feature vector \vec{v} , which summarizes the characteristics of the whole bag.



Introduction

We can split embedded methods in two categories:

- Methods that simply aggregate statistics of all the instances inside the bag



- Methods that analyze how the instances of the bag match certain prototypes that have been previously discovered in the data (*vocabulary-based methods*).



Methods Without Vocabularies

- These methods aggregate the statistics about the attributes of all the instances without making differentiation among these instances.
- Examples:
 - **Simple MI** (Dong et al, 2011). Maps each bag to the average of the instances inside:

$$\mu(X) = \frac{1}{|X|} \sum_{\vec{x} \in X} \vec{x}$$

- Gartner et al (2002) propose to map each bag X to a *min-max vector*,

$$\mu(X) = (a_1, \dots, a_d, \dots, b_1, \dots, b_d)$$

$$a_j = \min_{\vec{x} \in X} \vec{x} \quad b_j = \max_{\vec{x} \in X} \vec{x} \quad j=1 \dots d \quad (\text{instance dimensionality})$$

Vocabulary-Based Methods

- In these methods bags are related with concepts defined previously (**vocabulary**).
- Embedded space contains information about the relationship between bags and concepts in the vocabulary.
- Many times vocabulary concepts are obtained automatically in a unsupervised way (by clustering).
- There is a mapping from bag space to embedded space where vocabulary concepts play a key role.

Elements in a Vocabulary-based Method

- **Vocabulary**

- Stores K concepts.
- Most of the times the term concept means *class of instances*.

- **Mapping function**

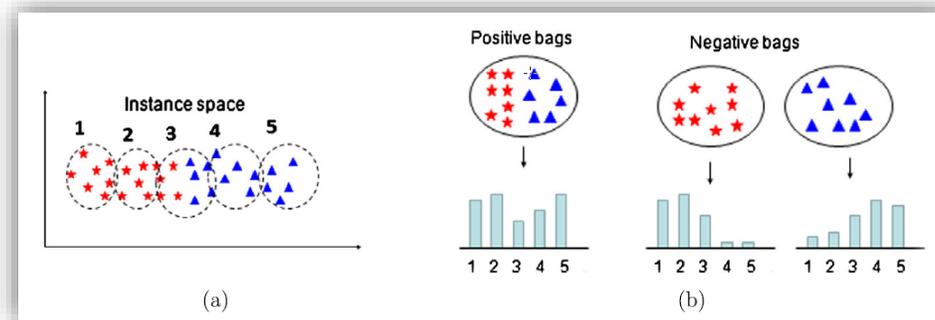
- Given a bag X and a vocabulary V , this mapping function $\mathcal{M}(X, V) = \vec{v}$ obtains a D -dimensional vector that match between the instances $\vec{x}_i \in X$ and the concepts $C_j \in V$

- **Standard supervised classifier.** Classifies the feature vector in the embedded space, using a training set $\mathcal{T}_{\mathcal{M}} = \{(\vec{v}_1, y_1), \dots, (\vec{v}_N, y_N)\}$

Vocabulary-based methods differs in vocabulary and mapping function

Histogram-based Methods

These methods use a function \mathcal{M} that maps each bag X into a histogram $\vec{v} = (v_1, \dots, v_K)$ where the j -th bin v_j counts how many instances of X fall into the j -th vocabulary class C_j .



- Classes (vocabulary) are automatically generated by a clustering algorithm
- Mapping function:

$$\mathcal{M}(X, V) = (v_1, \dots, v_K) \quad v_j = \frac{1}{Z} \sum_{\vec{x}_i \in X} f_j(\vec{x}_i)$$

Examples of Histogram-based Methods

Bag-of-words method (Sivic, 2003)

- Clustering method: K-Means

- Mapping function: $f_j(\vec{x}_i) = \begin{cases} 1 & \text{if } j = \arg \min_{k=1\dots,K} \|\vec{x}_i - \vec{p}_k\| \\ 0 & \text{otherwise} \end{cases}$

YARDS algorithm (Foulds, 2008)

- There are as many clusters as instances

- Mapping function: $f_j(\vec{x}_i) = \exp\left(-\frac{\|\vec{x}_i - \vec{p}_j\|^2}{\sigma^2}\right)$

Distance-Based methods

Instead of counting the number of instances that fall into class C_j , distance-based methods measure the distance $d_j(\vec{x}_i)$ between a given instance $\vec{x}_i \in X$ and the j-th concept. That is:

$$\mathcal{M}(X, V) = (v_1, \dots, v_K) \quad v_j = \min_{\vec{x}_i \in X} d_j(\vec{x}_i) \quad j = 1, \dots, K$$

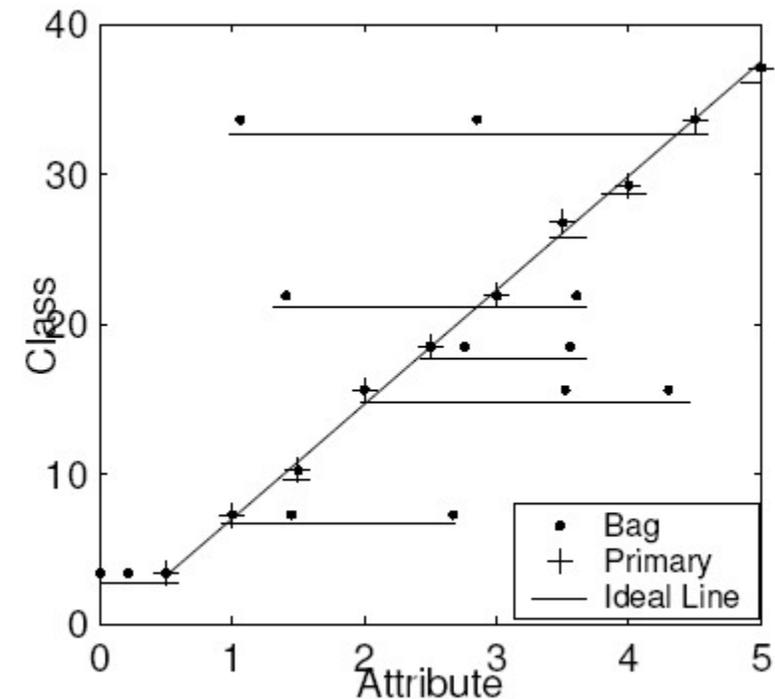
- Several authors developed similar proposals:
 - Auer et al, 2004
 - DD-SVM (Chen & Wang, 2004)
 - MILES (Chen et al, 2006)
- Concepts definition:
 - Hard assignment Clustering, like K-Means
 - Concepts defined explicitly by authors
- Distance functions:
 - Euclidean distance
 - Mahalanobis distance

Distance Based Bag-of-Words

Other Multiple Instance Learning Paradigms

Multiple Instance Regression

- There are few proposals for learning real-valued labels from multi-instance data.
- They assume that one of the instances is responsible of the concept under study (similar to the Dietterich hypothesis).
- The algorithm search the points that best fit to the hyperplane that represents the concept.



S. Ray and D. Page. Multiple instance regression. In Proceedings of the 18th International Conference on Machine Learning, pages 425–432, 2001

Multiple Instance Clustering

- M.-L. Zhang and Z.-H. Zhou published in 2009 the first paper about clustering of multi-instance data

M.-L. Zhang & Z.-H. Zhou. Multi-instance clustering with applications to multi-instance prediction. *Appl. Intell.* 31, 47-68, 2009.

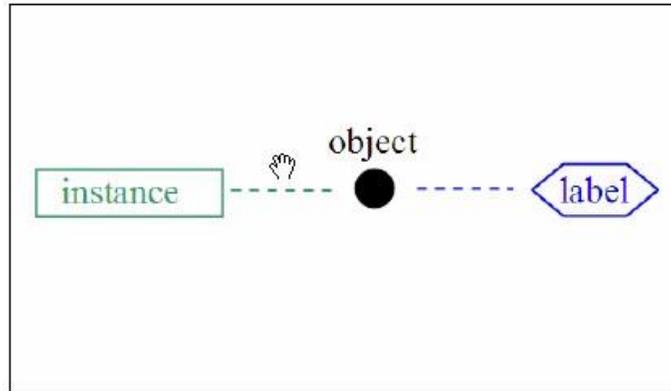
- The algorithm proposed, called BAMIC, is based on the k-medoids algorithm, modifying the metric used to set the distance between examples (bags). They tried three different distances:
 - Minimal Hausdorff distance.
 - Maximal Hausdorff distance
 - **Average Hausdorff distance**

Multi-instance Multi-label Classification

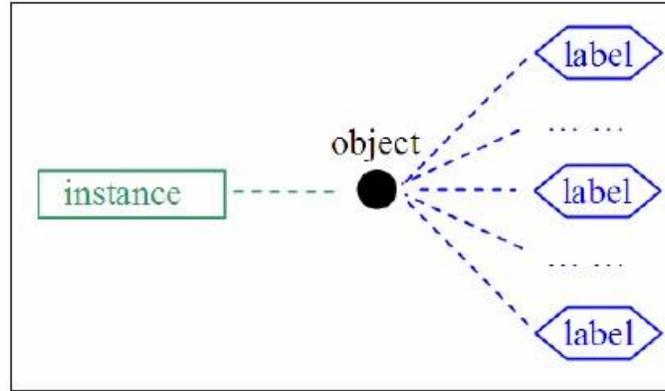
- In this learning paradigm, learner have to learn **a set of labels** for a given object (represented as multi-instances).
- The objective here is dealing with ambiguity both in the input and output spaces.
 - Input space ambiguity. There are several instances per input object
 - Output space ambiguity. There are several labels for the same object
- This paradigm is a generalization of two previous learning paradigms:
 - Multiple instance learning
 - Multi-label learning (or multi-label classification)
- First reference about this topic

Z.-H. Zhou. Mining ambiguous data with multi-instance multi-label representation.
Lecture Notes in Computer Science 4632, 2007

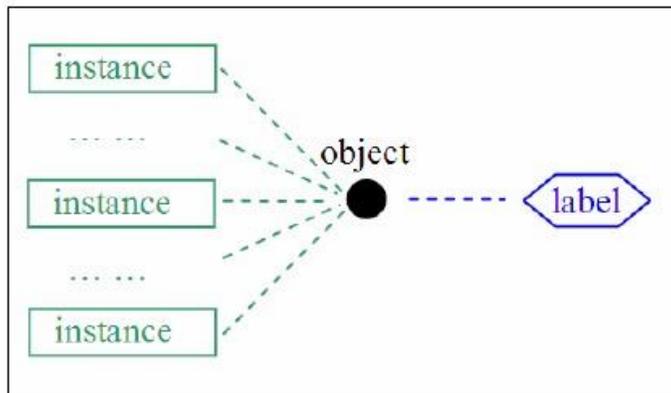
Comparing MIML with other Classification Paradigms



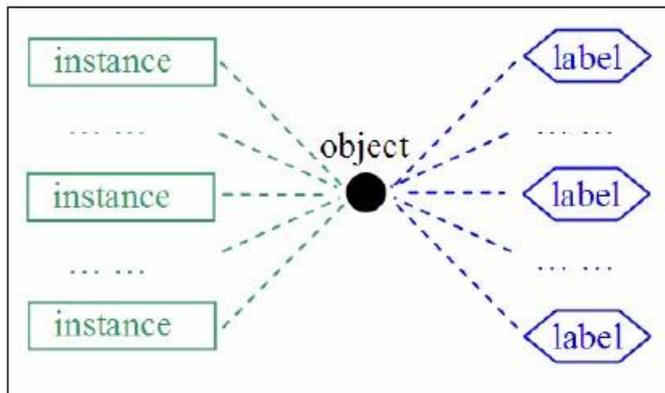
Traditional supervised learning



Multi-label learning

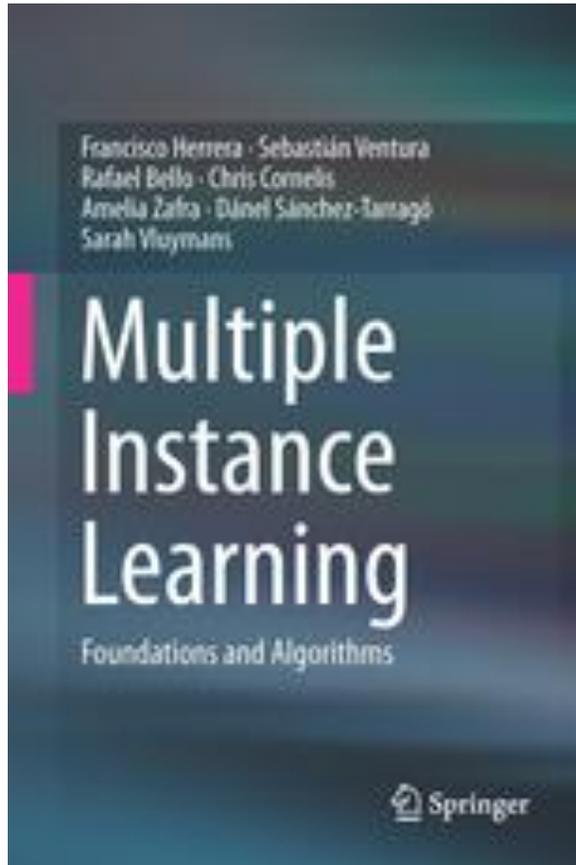


Multi-instance learning



Multi-instance multi-label learning

New Book on Multiple Instance Learning



Publication website

<http://www.springer.com/gp/book/9783319477589>



شكرا

The Mosque of Cordoba (169-633 AH)