Unsupervised Machine Learning and Data Mining

 $\overline{(Course Slides)}$

Javier Béjar

AMLT

Master in Artificial Intelligence

Course 2016/2017 Fall Semester

UPC

 \odot

Esta obra está bajo una licencia Reconocimiento-NoComercial-CompartirIgual de Creative Commons.

Para ver una copia de esta licencia, visite http://creativecommons.org/licenses/by-nc-sa/2.5/es/ o envie una carta a

Creative Commons, 559 Nathan Abbott Way, Stanford, California 94305, USA.

Contents

1	Introduction to Knowledge Discovery	1
2	Data Preprocessing	19
3	Data Clustering	57
4	Cluster Validation	113
5	Clustering of Large Datasets	137
6	Semisupervised Clustering	155
7	Association Rules	169
8	Mining Structured Data	197

Preface

These are the slides for the first half of the course *Advanced Machine Learning Techniques* from the master on Artificial Intelligence of the Barcelona Computer Science School (Facultat d'Informàtica de Barcelona), Technical University of Catalonia (UPC, BarcelonaTech).

This slides are used in class to present the topics of the course and have been prepared using the papers and book references that you can find in the course website http://www.cs.upc.edu/~bejar/amlt/amlt.html).

This document is a complement so you can prepare the classes and use it as a reference, but it is not a substitute for the classes or the material from the webpage of the course.

Javier Béjar Barcelona, 2016

Introduction to Knowledge Discovery











Knowledge Discovery in Databases Definitions of KDD

KDD definitions

"It is the search of valuable information in great volumes of data"

"It is the explorations and analysis, by automatic or semiautomatic tools, of great volumes of data in order to discover patterns and rules"

"It is the nontrivial process of identifying valid, novel, potentially useful, and ultimately understandable patterns in data"

	Knowledge Discovery in Databases Definitions of KDD
Elements	of KDD
Pattern:	Any representation formalism capable to describe the common characteristics of a group if instances
Valid:	A pattern is valid if it is able to predict the behaviour of new information with a degree of <i>certainty</i>
Novelty:	It is novel any knowledge that it is not know respect the domain knowledge and any previous discovered knowledge
Useful:	New knowledge is useful if it allows to perform actions that yield some benefit given a established criteria
Understanda	ble: The knowledge discovered must be analyzed by an expert in the domain, in consequence the interpretability of the result is important
Javier Béjar ©��	O (CS - MIA) Knowledge Discovery in Databases AMLT - 2016/2017 11 / 32 O
	The KDD process
 Knowledg 	e Discovery in Databases
2 The KDD	process
3 Application	ons
4 Tools	
5 Challenge	S



Javier Béjar ©🛈 🕲 (CS - MIA)













- customer churn
- Fraud detection
- Control/analysis of industrial processes
- e-commerce, on-line recommendation
- Financial data (stock market analysis)

• WEB mining

- Text mining, document search/organization
- Social networks analysis
- User behavior



Tools for KDD
 There are a lot of tools available for KDD
 Some tools were developed at universities (C5.0, CART/MARS) and have become a commercial product, others still remain open source (Weka, R, scikit-learn)
• Big fish eats little fish (C5.0 \rightarrow Clementine \rightarrow SPSS-clementine \rightarrow IBM DBMiner)
 Data analysis software companies incorporate KDD techniques inside classical data analysis tools (SPSS, SAS)
 Companies selling databases add KDD tools as an added value (IBM DB2 (intelligent Miner), SQL Server, Oracle)
 Machine Learning as a Service (Amazon, Microsoft, Google, IBM Watson, Big ML,)
Javier Béjar ⊚⊕⊛⊚ (CS - MIA) Knowledge Discovery in Databases AMLT - 2016/2017 29 / 32
Tools
Tools for the course
 Python General programming language, easy to learn numpy, scipy, pandas scikit-learn (http://scikit-learn.org) Data preprocessing Clustering Algorithms Association Rules
 Python General programming language, easy to learn numpy, scipy, pandas scikit-learn (http://scikit-learn.org) Data preprocessing, Clustering Algorithms, Association Rules, R (http://cran.r-project.org/) Statistic analysis oriented language, more steep learning curve Many packages Data preprocessing, Clustering Algorithms, Association Rules,



Open problems Scalability (More data, more attributes) Overfitting (Patterns with low interest) Statistical significance of the results Methods for temporal data/relational data/structured data Methods for data cleaning (Missing data and noise) Pattern comprehensibility Use of domain knowledge Integration with other techniques (OLAP, DataWarehousing, Business Intelligence, Intelligent Decision Support Systems) Privacy

Data Preprocessing




















If you have downloaded the code from the repository you will able to play with the notebooks (run jupyter notebook to open the notebooks)



• Attribute reduction has different solutions









- PCA is a linear transformation, this means that if data is linearly separable, the reduced dataset will be linearly separable (given enough components)
- We can use the *kernel trick* to map the original attribute to a space where non linearly separable data is linearly separable
- Distances among examples are defined as a dot product that can be obtained using a kernel:

$$d(x_i, x_j) = \Phi(x_i)^T \Phi(x_j) = K(x_i, x_j)$$

• Different kernels can be used to perform the transformation to the feature space (polynomial, gaussian, ...)



Kernel PCA

Javier Béjar 🐵 🕤 (CS - MAI)

• The computation of the components is equivalent to PCA but performing the eigen decomposition of the covariance matrix computed for the transformed examples

$$C = rac{1}{M}\sum_{j=1}^M \Phi(x_j) \Phi(x_j)^{ au}$$

• <u>Pro:</u> Helps to discover patterns that are non linearly separable in the original space

Data Preprocessing

<u>Con</u>: Does not give a weight/importance for the new components



AMLT - 2016/2017







Javier Béjar 🐵 🕟 🕥 (CS - MAI)



Nonnegative Matrix Factorization (NMF)

 NMF performs an approximation of a matrix in the product of two matrices

$$V = WH$$

- Similar to PCA, the eigen decomposition transforms the original data matrix in the product of two matrices
- The <u>main difference</u> is that the values of the matrices are constrained to be positive
- The formulation assumes that the data is a sum of unknown positive latent variables
- The positiveness assumption helps to interpret the result
 - Eg.: In text mining, a document is an aggregation of topics
- There are many variants and algorithms









Local linear embedding

- Performs a transformation that preserves local structure
- Assumes that each instance can be reconstructed by a linear combination of its neighbors (weights)
- From this weights a new set of data points that preserve the reconstruction is computed for a lower set of dimensions
- Different variants of the algorithm exist









- Performs a transformation that preserves locality of closer points and puts farther away non neighbor points
- Given a set of pairs of points \mathcal{N} where a pair (i, i') belong to the set if i is among the K neighbors of i' or viceversa
- Minimize the function:

$$S_L(z_1, z_2, \ldots, z_N) = \sum_{(i,i') \in \mathcal{N}} (d_{ii'} - ||z_i - z_{i'}||)^2 - \tau \sum_{(i,i') \notin \mathcal{N}} (||z_i - z_{i'}||)$$

• The parameters τ controls how much the non neighbors are scattered











Dimensionality Reduction Application: Wheelchair control

Wheel chair control characterization

MDS 88 \implies 3 dimensions 2 1 0 -1 -2 -3 -3 -2 $^{-1}$ 0 1 -3 2 -2 $^{-1}$ 3 0 1 4 2 Javier Béjar 📧 (CS - MAI) AMLT - 2016/2017 **>)** Data Preprocessing



Dimensionality Reduction Application: Wheelchair control

Wheel chair control characterization

ISOMAP (k-n=10) 88 \implies 3 dimensions





Dimensionality Reduction Attribute Selection				
Unsupervised feature selection				
 The filter methods order the attributes using measures of the structure of the data Measures of properties of the spatial structure of the data (Entropy, PCA, laplacian matrix) Measures of the relevance of the attributes respect the inherent structure of the data Measures of attribute correlation 				
 The wrapper methods are more diverse Clustering algorithms that compute weights for the attributes Clustering algorithms with an objective function that penalizes the size of the model Consensus clustering 				
Javier Béjar 🐵 🕤 (CS - MAI) Data Preprocessing AMLT - 2016/2017 59 / 71				
Dimensionality Reduction Attribute Selection				
Laplacian Score				

- The Laplacian Score is a filter method that ranks the features respect to their ability of preserving the natural structure of the data.
- This method uses the spectral matrix of the graph computed from the near neighbors of the examples
- Similarity is usually computed using a gaussian kernel (edges not present have a value of 0)

$$S_{ij} = e^{rac{||x_i - x_j||^2}{\sigma}}$$

And all edges not present have a value of 0.

Laplacian Score

• The Laplacian matrix is computed from the similarity matrix *S* and the degree matrix *D* as

Dimensionality Reduction

$$L = S - D$$

Attribute Selection

• The score first computes for each attribute r and their values f_r the transformation \tilde{f}_r as:

$$\tilde{f}_r = f_r - \frac{f_r^T D 1}{1^T D 1} 1$$

• and then the score L_r is computed as:

$$L_r = \frac{\tilde{f}_r^T L \tilde{f}_r}{\tilde{f}_r^T D \tilde{f}_r}$$

• This gives a ranking for the relevance of the attributes

Javier Béjar 📧 🕟 🧿 (CS - MAI)	Data Preprocessing	AMLT - 2016/2017 61 / 71	

Dimensionality Reduction Attribute Selection

Python Notebooks

These two Python Notebooks show some examples dimensionality reduction and feature selection

- Dimensionality reduction and feature selection Notebook (click here to go to the url)
- Linear and non linear dimensionality reduction Notebook (click here to go to the url)

If you have downloaded the code from the repository you will able to play with the notebooks (run jupyter notebook to open the notebooks)









Similarity functions
Python Code
 In the code from the repository inside subdirectory DimReduction you have the Authors python program The code uses the datasets in the directory authors from the datasets zipfile Auth1 has fragments of books that are novels or philosophy works Auth2 has fragments of books written in English and books translated to English
 The code transforms the text to attribute vectors and applies different dimensionality reduction algorithms
 Modifying the code you can process one of the datasets and choose how the text is transformed into vectors
Javier Béjar 🐵 🕤 (CS - MAI) Data Preprocessing AMLT - 2016/2017 71 / 7



Data Clustering











Javier Béjar © 🖲 🕲 (CS - MAI)



AMLT - 2016/2017 12 / 110



Hierarchical Algorithms	Statistical Algorithms
	Statistical Aigorithins

Hierarchical algorithms - Matrices

	2		3	4	ļ	5	
1	6		8	2)	7	
2			1	5	•	3	
3				1	0	9	
4						4	
		2,3	3	4	ŀ	5	
1		7		2)	7	
2,3	3			7.5		6	
4						4	
	I	1,4		ŀ	5		
2	,3	7	'.2!	5	6		
1	,4	5.5					5
		•	1,4,5				
2,3			6	.72	25		








0)=0.50 P(V C)		P(C0)=0.25		P(V C)
Negro	0.0	Color	Negro	0.0
Blanco	1.0	COIOI	Blanco	1.0
Cuadrado	0.0		Cuadrado	0.0
Triángulo	0.0	Forma	Triángulo	1.0
Círculo	1.0		Círculo	0.0

P(V|C)

0.0

1.0

0.0

0.33

0.66

P(C0)=0.75

Color

Form

P(C0)= Ne Color

Bla Cu Forma

Negro

Blanco

Cuadrado

Triángulo

Círculo

COBWEB Category utility (CU)

- Category utility balances:
 - Intra class similarity: $P(A_i = V_{ij}|C_k)$
 - Inter class similarity: $P(C_k | A_i = V_{ij})$
- It measures the difference between a partition of the data and no partition at all
- For qualitative attributes and a set of k categories {C₁, ... C_k} is defined as:

$$\frac{\sum_{k=1}^{K} P(C_k) \sum_{i=1}^{I} \sum_{j=1}^{J} P(A_i = V_{ij} | C_k)^2 - \sum_{i=1}^{I} \sum_{j=1}^{J} P(A_i = V_{ij})^2}{K}$$

(see the derivation on the paper)

Javier Béjar ©€© (CS - MAI)

Unsupervised Learning

AMLT - 2016/2017 21

Hierarchical Algorithms Concept Formation Algorithms

Operators

- **Incorporate:** Put the example inside an existing class
- New class: Create a new class at this level
- **Merge:** Two concepts are merge and the example is incorporated inside the new class
- Divide: A concept is substituted by its children



COBWEB Algorithm

Procedure: Depth-first limited search COBWEB (x: Example, H: Hierarchy)

Update the father with the new example

if we are in a leaf then

Create a new level with this example

else

Compute **CU** of incorporating the example to each class Save the two best **CU**

Compute **CU** of merging the best two classes

Compute **CU** of splitting the best class

Compute **CU** of creating a new class with the example

Recursive call with the best choice

end

Javier Béjar ©��� (CS - MAI)

Unsupervised Learning

Unsupervised Learning
 Hierarchical Algorithms
 Partitional algorithms
 Applications

AMLT - 2016/2017



AMLT - 2016/2017 26 / 110



Javier Béjar @🛈 🏵 🎯 (CS - MAI)

Partitional algorithms Model/Prototype Clustering						
K-means Practical problems						
 The algorithm is sensitive to the initialization (to run the algorithm from several random initializations is a common practice, with the additional computational cost) 						
 Sensitive to clusters with different sizes/densities and <u>outliers</u> 						
 To find the value of k is not an easy problem (experimentation with different values is needed) 						
 A solution is found even if the classes are not hyperspherical (some classes could be split/merged) 						
No guarantee about the quality of the solution						
The spatial complexity makes it not suitable for large datasets						
Javier Báiar @@@@ (CS - MAI) Jacupervised Learning AMLT - 2016/2017 29 / 110						
Partitional algorithms Model/Prototype Clustering						
K-means++ Initialization Strategies						
 K-means++ modifies the initialization strategy of the algorithm 						
The idea is to try to maximize distance among initial centers						

- Algorithm:
 - Choose one center uniformly at random from among the data points
 - 2 For each data point x, compute d(x, c), the distance between x and the nearest center that has already been chosen
 - S Choose one new data point at random as a new center, using a weighted probability distribution where a point x is chosen with probability proportional to $d(x, c)^2$
 - Repeat Steps 2 and 3 until k centers have been chosen
 - Solution Proceed with the standard K-means algorithm



- Keep the clustering with the best objective function as the *C*-clusters solution
- Pro: Reduces the initialization problem/obtains all partitions from 2 to K
- **Con:** Computational cost (runs $K \times N$ the K-means algorithm)

Javier Béjar ©👀 🕲 (CS - MAI)

Unsupervised Learning



• Pro: It is not sensitive to outliers







Javier Béjar @🖲 🏵 🎯 (CS - MAI)



Partitional algorithms Model/Prototype Clustering							
EM Algorithm (K Gaussian)							
The computations depend on the assumptions that we make about the attributes (independent or not, same σ ,)							
 The attributes are independent: μ_i and σ_i have to be computed for each class (O(k) parameters) (model: hyper spheres or ellipsoids parallel to coordinate axis) 							
 The attributes are not independent: μ_i, σ_i and σ_{ij} have to be computed for each class (O(k²) parameters) (model: hyper ellipsoids non parallel to coordinate axis) 							
Javier Béjar @@@@ (CS - MAI) Unsupervised Learning AMLT - 2016/2017 43 / 110							
Partitional algorithms Model/Prototype Clustering							
Partitional algorithms Model/Prototype Clustering EM Algorithm (K Gaussian)							
Partitional algorithms Model/Prototype Clustering EM Algorithm (K Gaussian)							
Partitional algorithms Model/Prototype Clustering EM Algorithm (K Gaussian) Image: Clustering in the second se							
Partitional algorithms Model/Prototype Clustering EM Algorithm (K Gaussian) • • For the case of A independent attributes:							
• For the case of A independent attributes: $P(x \overrightarrow{\mu_i}, \Sigma_i, w_i) = \prod_{j=1}^{A} P(x \mu_{ij}, \sigma_{ij}, w_i)$							
• For the case of A independent attributes: $P(x \overrightarrow{\mu_i}, \Sigma_i, w_i) = \prod_{j=1}^{A} P(x \mu_{ij}, \sigma_{ij}, w_i)$ • The model to fit is							
• For the case of A independent attributes: $P(x \overrightarrow{\mu_i}, \Sigma_i, w_i) = \prod_{j=1}^{A} P(x \mu_{ij}, \sigma_{ij}, w_i)$ • The model to fit is $P(x \overrightarrow{\mu_i}, \overrightarrow{\sigma}) = \sum_{i=1}^{K} P(w_i) \prod_{j=1}^{A} p(x \mu_{ij}, \sigma_{ij}, w_i)$							
• For the case of A independent attributes: $P(x \overrightarrow{\mu_i}, \Sigma_i, w_i) = \prod_{j=1}^{A} P(x \mu_{ij}, \sigma_{ij}, w_i)$ • The model to fit is $P(x \overrightarrow{\mu_i}, \overrightarrow{\sigma}) = \sum_{i=1}^{K} P(w_i) \prod_{j=1}^{A} p(x \mu_{ij}, \sigma_{ij}, w_i)$							







Fuzzy Clustering • Fuzzy clustering relax the hard partition constraint of K-means • Each instance has a membership to each partition • A new optimization function is introduced: $L = \sum_{i=1}^{N} \sum_{k=1}^{K} \delta(C_k, x_i)^b ||x_i - \mu_k||^2$ where $\sum_{k=1}^{K} \delta(C_k, x_i) = 1$ and b is a blending factor • This is an advantage from other algorithms when the groups are overlapped

Fuzzy Clustering

• C-means is the most known fuzzy clustering algorithm, it is the fuzzy version of K-means

Model/Prototype Clustering

- The algorithm performs an iterative optimization of the objective function
- The updating of the cluster centers is computed as:

Partitional algorithms

$$\mu_j = \frac{\sum_{i=1}^N \delta(C_j, x_i)^b x_i}{\sum_{i=1}^N \delta(C_j, x_i)^b}$$

• And the updating of the memberships:

$$\delta(C_j, x_i) = \frac{(1/d_{ij})^{1/(1-b)}}{\sum_{k=1}^{K} (1/d_{ik})^{1/(1-b)}}, d_{ij} = \|x_i - \mu_j\|^2$$

Javier Béjar ©€© (CS - MAI)

Unsupervised Learning

AMLT - 2016/2017

51 / 110

Partitional algorithms Model/Prototype Clustering

Fuzzy Clustering The C-means algorithm looks for spherical clusters, other alternatives: <u>Gustafson-Kessel algorithm</u>: A covariance matrix is introduced for each cluster in the objective function that allows elipsoid shapes and different cluster sizes <u>Gath-Geva algorithm</u>: Adds to the objective function the size and an estimation of the density of the cluster Also different objective functions can be used to detect specific shapes in the data (lines, rectangles, ...) this characteristic is widely used in image recognition









Partitional algorithms Density/Grid Clustering

DBSCAN/OPTICS

Ester, Kriegel, Sander, Xu A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise (DBSCAN) (1996)

Ankerst, Breunig, Kriegel, Sander **OPTICS: Ordering Points To Identify the Clustering Structure** (2000)

- Used in spatial databases, but can be applied to data with more dimensionality
- Based on finding areas of high density, it finds arbitrary shapes
- We define ε-neighbourhood, as the instances that are at a distance less than ε to a given instance

$$N_{\varepsilon}(x) = \{y \in X | d(x, y) \le \varepsilon\}$$

 We define core point as the instances that have a certain number of elements in N_ε(x)

$${\it Core_point} \equiv |N_arepsilon(x)| \geq {\it MinPts}$$



Partitional algorithms Density/Grid Clustering

DBSCAN - Density

 We say that two instances p and q are Direct Density Reachable with respect to ε and MinPts if:

•
$$p \in N_{\varepsilon}(q)$$

- 2 $|N_{\varepsilon}(q)| \geq MinPts$
- We say that two instances p and q are Density Reachable if there are a sequence of instances p = p₁, p₂,..., p_n = q where p_{i+1} is direct density reachable from p_i
- We say that *p* and *q* are **Density connected** if there is an instance *o* such that both *p* and *q* are **Density Reachable** from *o*





.....

DENCLUE

Hinneburg, Keim **An Efficient Approach to Clustering in Large Multimedia Databases with noise** (1998)

- Clustering algorithm based on kernel density estimation
- Defines the influence of an example in a dataset as the sum of a kernel for all the data (for example a gaussian kernel)

$$f_B^{\mathcal{Y}}(x) = f_B(x, y)$$

• Defines the density function of a dataset as the sum of the influences of all examples of the dataset

$$f_B^D(x) = \sum_{i=1}^N f_B^{x_i}(x)$$



- density-attractor x^* as the subset of examples being density-attracted by x^* and with $f_B^D(x^*) > \xi$
- We define an Arbitrary-shape cluster (wrt to σ,ξ) for a set of density-attractors X as the subset of examples being density-attracted to any x* ∈ X with f^D_B(x*) > ξ and with any density-attractor from X connected by a path P with ∀p ∈ P : f^D_B(p) > ξ
- Different algorithms can be reproduced with the right choice of σ and ξ parameters (DBSCAN, K-means, hierarchical clustering)



- For efficiency reasons the density functions are only computed for the neighbours of each point
- The algorithm is divided in two phases:
 - Preclustering: The dataspace is divided in d-dimensional hypercubes, only using the cubes with datapoints. These hypercubes are mapped to a tree structure for efficient search. From this hypercubes given a threshold, only the highly populated ones and their neigbours are considered.
 - Clustering: For each datapoint in the hypercubes a local density function and a local gradient are computed. Using a Hill-Climbing algorithm the density-attractor for each point is computed. For efficiency reasons each point near the path computed during the search of a density-attractor is assigned to that density-attractor.

CLIQUE

Agrawal, Gehrke, Gunopulos, Raghavan Automatic Subspace Clustering of High Dimensional Data for Data Mining Applications (1998)

- It generates CNF descriptions from the groups that discovers
- The goal is to find a space with less dimensions where the groups are easier to identify
- The algorithm is scalable and can be used with data of high dimensionality
- The goal is to obtain interpretability
- It is based on density estimation techniques
 - Each dimension of the dataset is divided using a grid
 - The bins with higher density than certain threshold are identified

Javier Béjar ©€®⊚ (CS - MAI)	Unsupervis	ed Learning	AMLT - 2016/2017	69 / 110
	Partitional algorithms	Density/Grid Clustering		

Density/Grid Clustering

CLIQUE - Algorithm

- 1. Identify the subspaces with clusters
 - It applies a bottom up strategy identifying first the bins of higher density in one dimension and combining them
 - A combination of k + 1 dimensions can only have high density bins if there are high density bins in k dimensions
 - This rule gives the set of candidate to high density bins when we increase the dimensionality by one
 - Some other heuristics are used to reduce the computational cost of the search



CLIQUE - Algorithm

- 2. Identify the clusters
 - We receive the set of dense bins from the previous step
 - We look for contiguous bins and consider them in the same group
 - Two bins are connected if they share a side in one dimension or there is an intermediate bin that connect them
 - This problem is equivalent to find the connected components of a graph









- We define the degree matrix *D* as the matrix with values d_1, d_2, \ldots, d_n as diagonal
- We can define different Laplace matrices:
 - Unnormalized: L = D W
 - Normalized: $L_{sym} = D^{-1/2}LD^{-1/2}$ or also $L_{rw} = D^{-1}L$



Spectral Clustering (Second approach)

- From the graph defined from the similarity matrix and its Laplacian matrix (see previous slides), it is possible to define the problem as a graph partitioning problem, so any graph partitioning algorithm can be used
- Given two disjoint sets of vertex A and B, we define:

$$cut(A,B) = \sum_{i \in A, j \in B} w_{ij}$$

 We can partition the graph solving the mincut problem choosing a partition that minimizes :

$$cut(A_1,\ldots,A_k) = \sum_{i=1}^k cut(A_i,\overline{A_i})$$




Partitional algorithms Other Approaches Affinity Propagation Clustering - Updating All availabilities are initialized to 0 • The responsibilities are updated as: $r(i,k) = r(i,k) - \max_{k' \neq k} \{a(i,k') + s(i,k')\}$ The availabilities are updated as: $a(i,k) = min\{0, r(k,k) + \sum_{i' \notin \{i,k\}} max(0, r(i',k))\}$ • The self availability a(k, k) is updated as: $a(k,k) = \sum_{i' \neq k} max(0,r(i',k))\}$ • The exemplar for a point is identified by the point that maximizes a(i, k) + r(i, k), if this point is the same point, then it is an exemplar. Javier Béjar ©€© (CS - MAI) AMLT - 2016/2017 Unsupervised Learning Partitional algorithms Other Approaches Affinity Propagation Clustering - Algorithm Update the responsibilities given the availabilities Opdate the availabilities given the responsibilities Ompute the exemplars

• Terminate if the exemplars do not change in a number of iterations



Unsupervised Neural Networks

- Self-organizing maps are an unsupervised neural network method
- Can be seen as an on-line constrained version of K-means
- The data transformed to fit in a 1-d or 2-d regular mesh (rectangular or hexagonal)
- The nodes of this mesh are the prototypes
- This algorithm can be used as a dimensionality reduction method (from *N* to 2 dimensions)





Partitional algorithms Other Approaches

Self-Organizing Maps

- During the iterations the mesh is transformed to be closer to the data, but maintaining the bidimensional relationship between prototypes
- The performance of the algorithm depends on the learning rate α , usually is decreased from 1 to 0 during the iterations
- The neighborhood of a prototype is defined by the adjacency of the cells and the distance of the prototypes
- The number of neighbors used in the update is decreased during the iterations from a predefined number to 1 (only the prototype nearest to the observation)
- Different variations of the algorithm give more weight depending on the distance of the prototypes





Applications				
The data				
 The dataset is composed of several hundreds of thousands of events during several months (represents less than 5% of the actual events) 				
 Each event has a geographical position (latitude/longitude) and a time stamp (inside an area of 30 × 30 Km²) 				
 For the purpose of the analysis (geoprofiles) the actual data is difficult to analyze: 				
 The resolution of the coordinates is too fine, low probability of two events in the exact same place 				
 Clustering geographically the events can help to make sense of the data 				
lavier Béjar @@@@ (CS - MAI) Unsupervised Learning AMLT - 2016/2017 97 / 110				
Applications				
Applications Clustering of the events				
Applications Clustering of the events				
Applications Clustering of the events				
Applications Clustering of the events Being the data geographically distributed there are few alternatives (clusters of arbitrary shapes)				
Applications Clustering of the events • Being the data geographically distributed there are few alternatives (clusters of arbitrary shapes) • The more adequate methods seem:				
Applications Clustering of the events • Being the data geographically distributed there are few alternatives (clusters of arbitrary shapes) • The more adequate methods seem: • Density based clustering				
Applications Clustering of the events • Being the data geographically distributed there are few alternatives (clusters of arbitrary shapes) • The more adequate methods seem: • Density based clustering • Grid based clustering				
Applications Clustering of the events • Being the data geographically distributed there are few alternatives (clusters of arbitrary shapes) • The more adequate methods seem: • Density based clustering • Grid based clustering • The size of the dataset (~ 2.5 million events) could arise scalability issues, so a preliminar coarse grained clustering could be helpful (for example k-means, or the leader algorithm)				



Applications
Geographical profiles
 We want to find groups of users that have similar behavior (in place and time) along the period of data collection
 The clusters from the discretization can be the basis for the geographical profiles
 Additionally a discretization of time can help to have a finer grained representation
 Different representation can be used to generate the dataset: Presence/absence of a user in a place at a specific time interval Absolute number of visits of a user to a place in a time interval Normalized frequency of the visits of a user to a place in a time interval
 Different choices of representation and discretization allow for different analysis
Javier Bejar @@@@ (CS - WAT) Onsupervised Learning AMLT - 2010/2017 101 / 110
Applications
Geographical profiles
 It is more difficult to choose what clustering algorithm is adequate We can explore different alternatives and analyze the results Some choices will depend on: If the dataset generated is continuous or discrete The size of the dataset Our assumptions about the model that represents our goals (Shape of the clusters, Separability of the clusters/Distribution of examples) Interpretability/Representability of the clusters Experts on the domain have to validate the results, but some general assumptions can be used to evaluate the alternatives We will explore two alternatives:
 K-means algorithm (simple and efficient, spherical clusters) Affinity propagation clustering (adaptative and non predefined shape)



<section-header>











Applications

Affinity propagation Cluster example



Javier Béjar ©��� (CS - MAI)

Cluster Validation





- Why do we want to evaluate them?
 - To avoid finding patterns in noise
 - To compare clustering algorithms
 - To compare different models/parameters



Ompute the quotient:

$$H = \frac{\sum_{i=1}^{n} d(p_i)}{\sum_{i=1}^{n} d(p_i) + \sum_{i=1}^{n} d(q_i)}$$

• If points are uniformly distributed the value of H will be around 0.5







• Trace criteria (lower overall intracluster distance/higher overall intercluster distance)

$$Tr(S_W) = rac{1}{K} \sum_{i=1}^{K} S_{W_k} \ Tr(S_B) = rac{1}{K} \sum_{i=1}^{K} S_{B_k}$$

• Determinant criteria (higher cross-intercluster distance)

$$Det(S_M) = rac{1}{K^2} \sum_{i=1}^K \sum_{j=1}^K S_{M_{i,j}}$$





Internal criteria - Indices

- In the literature can be found more than 30 different indices
- Several studies and comparisons have been performed
- Recent studies (Arbelatiz et al, 2013) have exhaustively studied these indices, some have a performance significatively better that others
- Some of the indices show a similar performance (not statistically different)
- The study concludes that Silhouette, Davies-Bouldin and Calinski Harabasz perform well in a wide range of situations



Internal criteria - Davies-Bouldin





- It is equivalent to have a labeled dataset (ground thruth)
- If we do not have a model these criteria can be used to compare the results of using different parameters of different algorithms
 - For instance, can be used to assess the sensitivity to initialization
- The main advantage is that these indices are independent of the examples/cluster description
- These means that can be used to assess any clustering algorithm





$$VI(C, C') = H(C) + H(C') - 2I(C, C')$$

• Adjusted Mutual Information:

$$AMI(U, V) = \frac{MI(U, V) - E(MI(U, V))}{\max(H(U), H(V)) - E(MI(U, V))}$$



what number of clusters is more probable



Number of clusters

Number of clusters - The Gap Statistic

- Assess the number of clusters comparing a clustering with the expected distribution of data given the null hypothesis (no clusters)
- Different clusters of the data increasing the number of clusters and compared to datasets (B) generated with a uniform distribution
- The Gap statistic:

$$Gap(k) = (1/B) \sum_{b} log(S_W(k)_b) - log(S_W(k))$$

• From the st. dev. (sd_k) of $\sum_b log(S_W(k)_b)$ is defined s_k as:

$$s_k = sd_k\sqrt{1+1/B}$$

• The probable number of clusters is the smallest number that holds:

$$Gap(k) \geq Gap(k+1) - s_{k+1}$$







Javier Béjar © (CS - MAI) Clustering Evaluation/Model Assessment



Javier Béjar @ 🖲 🕲 (CS - MAI) Clustering Evaluation/Model Assessment

	Cluster Visualization			
1 Cluster Evaluation				
2 Internal criteria				
3 External criteria				
4 Number of clusters				
5 Application				
6 Cluster Visualization				
Javier Béjar @��@ (CS - MAI)	Clustering Evaluation/Model Assessment	AMLT - 2016/2017	37 / 43	
Cluster Visualization				
Cluster visualization				
 Another method for model assessment is to visualize the data and look for clusters 				

- Dimensionality reduction
 - Project the dataset to 2 or 3 dimensions
 - The clusters in the new space could represent clusters in the original space
 - The confidence depends on the reconstruction error of the transformed data and that the transformation maintains the relations in the original space
- Distance matrix visualization
 - The distance matrix represents the examples relationships
 - Can be rearranged so the closer examples appear in adjacent columns
 - Patterns in the rearranged matrix can show cluster tendency
- Both methodologies are computationally expensive



Javier Béjar ��� (CS - MAI) Clustering Evaluation/Model Assessment



Cluster visualization - Distance matrix






Clustering of Large Datasets

Clustering in KDD

Javier Béjar © 🛈 🕲 🎯

CS - MIA

AMLT - 2016/2017

Javier Béjar ©€© (CS - MIA)

 $Clustering \ in \ \mathsf{KDD}$

AMLT - 2016/2017

Outline





- One-pass
 - Process data as a stream
- Summarization/Data compression
 - Compress examples to fit more data in memory
- Sampling/Batch algorithms
 - Process a subset of the data and maintain/compute a global model
- Approximation
 - Avoid expensive computations by approximate estimation
- Paralelization/Distribution
 - Divide the task in several parts and merge models







Scalable Algorithms



Javier Béjar 🎯 🚱 🕲 (CS - MIA)

Scalal	ble Algorithms				
BIRCH					
 Zhang, Ramakrishnan, Livny BIRCH: An Efficient Data Clustering Method for Very Large Databases (1996) Strategy: One-pass + Summarization Hierarchical clustering with limited memory Incremental algorithm Based on probabilistic prototypes and distances We need two pass from the database Based on an specialized data structure named CF-tree (Clustering Feature Tree) 					
Javier Béjar ⊛��&@ (CS - MIA)	Clustering in KDD	AMLT - 2016/2017 15 / 1			
Scalal	ble Algorithms				
BIRCH (CF-tree)					
 It is a balanced n-ary tree that contains groups represented by probabilistic prototypes Leaves can contain as much as <i>L</i> prototypes and its radius can not be more than <i>T</i> Each non terminal node has a fixed branching factor (<i>B</i>), each element is a prototype that summarizes its subtree The choice of the parameters is crucial because we could fill all available space during the process This can be solved by changing the parameters values (basically <i>T</i>) and recompressing the tree. In fact <i>T</i> determines the granularity of the final groups 					



Scalable Algorithms

Javier Béjar © 🛞 🕲 (CS - MIA)

Clustering in KDD

AMLT - 2016/2017 17 / 1





Javier Béjar @🖲 🏵 🕥 (CS - MIA)



Canopy Clustering	
McCallum, Nigam, Ungar Efficient clustering of high-dimensional dat with application to reference matching (2002)	ta sets
 Strategy: Divide&Conquer + Approximation The approach is based on a two stages clustering The first stage can be seen as a preprocess determining the neighborhood of the densities of examples and reducing the n of distances that have to be computed on the second stage This first stage is the called <i>canopy clustering</i>, relies on a che distance and two parameters T₁ > T₂ This parameters are used as two centered spheres that determ to classify the examples. 	number eap nine how
Javier Béjar ⊛⊕⊛⊚ (CS - MIA) Clustering in KDD AMLT - 2016/2	/2017 23 / 1
Scalable Algorithms	
Scalable Algorithms Canopy Clustering - Algorithm	
 Algorithm: One example is picked at random and the cheap distance from example to the rest is computed All the examples that are at less than T2 are deleted and inclut the canopy The points at less than T1 are added to the canopy of this examited the procedure is repeated until the example list is empty Canopies can share examples After this procedure the examples can be clustered using differ algorithms For agglomerative clustering only the distances among the examine the canopies have to be computed 	n this uded in amples erent amples

Scalable Algorithms



Canopy Clustering - Algorithm

• For EM and K-means there are different methods:

Scalable Algorithms

- <u>Method 1</u>: Prototypes are associated with canopies, only distances inside a canopy are computed and only the examples inside a canopy can influence a prototype. The number of prototypes per canopy have to be decided. The EM or K-means only use for each prototypes the examples of the associated canopies
- <u>Method 2</u>: Use a set of prototypes covering the dataset and then distribute them in the canopies. We allow data in other canopies to influence the prototypes but only using the mean of their canopies.
- <u>Method 3</u>: We allow to create or merge prototypes inside a canopy avoiding to have prototypes that cover examples that belong to more than one canopy.



CURE

Guha, Rastogi, Shim CURE: An efficient clustering algorithm for large databases (1998) • Strategy: Sampling + Divide&Conquer • It uses hierarchical agglomerative clustering It is scalable and is capable of treating outliers • Scalability is obtained using sampling techniques and partitioning the dataset • It uses a set of representatives (c) for cluster instead of centroids, this allows to find non spherical groups • The distance is computed as the nearest pair of representatives among groups • The clustering algorithm is agglomerative and merges pairs of groups until k groups are obtained Javier Béjar ©€® (CS - MIA) AMLT - 2016/2017 Clustering in KDD 29 / 1 Scalable Algorithms **CURE - Algorithm** Oraws a random sample from the dataset Partitions the sample in p groups Executes the clustering algorithm on each partition Oeletes outliers Suns the clustering algorithm on the union of all groups until it obtains k groups • Label the data accordingly to the similarity to the k groups



Viswanath, Babu, Rough-DBSCAN: A fast hybrid density based clustering method for large data sets Pattern Recognition Letters, 2009, 30, 1477 - 1488

- Strategy; One-pass + Summarization
- Two stages algorithm:
 - Preprocess using the leader algorithm
 - Determine the instances that belong to the higher densities and their neighbours
 - 2 Apply DBSCAN algorithm
 - Determine the densities for the selected instances
 - Approximate the values of the densities from their distances and the sizes of the neighbor
 - Assign the neighbors accordingly to the found densities



Semisupervised Clustering



Outline



- Sometimes we have available some information about the dataset we are analyzing unsupervisedly
- Could be interesting to incorporate this information to the clustering process in order to:
 - Bias the search of the algorithm toward the solutions more consistent with our knowledge
 - Improve the quality of the result reducing the algorithm natural bias (predictivity/stability)





Semisupervised Clustering/Labeled Examples

Semi supervised clustering using labeled examples

- Assuming that we have some labeled examples, these can be used to obtain an initial model
- We only have to know what examples belong to clusters, the actual clusters are not needed
- We can begin from this model the clustering process
- This means that we decide the starting point of the search using the supervised information
- This initial model changes the search and the final model (bias)
- This differs from semi-supervised learning from a supervised perspective, were the labels of some of the examples are known

Semi supervised clustering using labeled examples

Basu, Banerjee, Mooney, "**Semi supervised clustering by seeding**", ICML 2002

- Algorithm based on K-means (spherical clusters based on prototypes)
- The usual initialization of K-means is by selecting randomly the initial prototype (there are other alternatives)
- Two proposals:

Javier Béjar 🐵 🕟 🧿 (CS - MAI)

- Use the labeled examples to build the initial prototypes (seeding)
- Use the labeled examples to build the initial prototypes and constrain the model so the labeled examples are always in the initial clusters (seed and constraint)
- The initial prototypes give an initial probability distribution for the clustering

Semi-supervised Clustering

Semisupervised Clustering/Labeled Examples

Semi supervised clustering using labeled examples





AMLT - 2016/2017



Constrained-KMeans

Algorithm: Constrained-KMeans

Input: The dataset \mathcal{X} , the number of clusters K, a set \mathcal{S} of labeled instances (k groups)

Output: A partition of \mathcal{X} in K groups

begin

Compute K initial prototypes (μ_i) using the labeled instances **repeat**

Maintain the examples from ${\mathcal S}$ in their initial classes

Assign each example from $\mathcal X$ to their nearest prototype μ_i

Recompute the prototype μ_i with the examples assigned

until Convergence

end





Semi supervised clustering using constraints

- To have labeled examples means that the number of clusters and something about the characteristic of the data is known
- Sometimes it is easier to have information about if two examples have to be in the same or different clusters
- This information can be expressed by means of constraints among examples: *must links* and *cannot links*
- This information can be used to bias the search and only look for models that maintain these constraints



Semi supervised clustering using constraints

Basu, Bilenko, Mooney, "A probabilistic framework for semi-supervised clustering", ICML 2002

- Algorithm based on K-means (spherical clusters based on prototypes)
- A set of must-link cannot-link constraints is defined over a subset of examples
- The quality function of the K-means algorithm is modified to bias the search
- A hidden markov random field is defined using the constraints







HMRF-KMeans

Algorithm: HMRF-KMeans

Input: The dataset \mathcal{X} , the number of clusters K, a set of must and cannot links, a distance function D and a set of weights for violating the constraints

Output: A partition of \mathcal{X} in K groups

begin

Compute K initial prototypes (μ_i) using constraints

repeat

E-step: Reassign the labels of the examples using the prototypes (μ_i) to minimize \mathcal{J}_{obj}

M-step: Given the cluster labels recalculate cluster centroids to minimize $\mathcal{J}_{\textit{obj}}$

until Convergence

end

Javier Béjar 📧 🕟 🕥 (CS - MAI) 🛛 Semi-

Semi-supervised Clustering

AMLT - 2016/2017

21 / 1

Semisupervised Clustering/Constraints

Semi supervised clustering using constraints

- The constraints can also be seen as an indication of the inadequacy of the distance measure between two instances
 - A must-link violation means that the distance function assigns a similarity less than the desired similarity
 - A cannot-link violation means that distance function assigns a similarity larger than the desired similarity
- The previous algorithm can be modified to introduce weights to modify the distances among instances with constraints (an additional maximization step is needed)



Semisupervised Clustering/Distance Learning

Semi supervised clustering with Distance Learning

- Other approaches learn the more adequate distance function to fulfill the constraints
- The constraints are used as a guide to find a distance matrix that represents the relations among examples
- This problem can be defined as an optimization problem that optimizes the distances among examples with respect to the constraints
- This methods are related to kernel methods, the goal is to learn a Kernel matrix that represents a new space where the instances have appropriate distances



Javier Béjar 🐵 🕟 🇿 (CS - MAI)

Association Rules





	Association Rules						
Association rules							
• A Binary database is a database where each row is composed of							
• A Dinary database is a database where each row is composed of binary attributes							
,			P	C	P		
		A	В	C	D	• • •	
	T1	1	0	0	1		
	T2	0	1	1	1		
	Т3	1	0	1	0		
	Τ4	0	0	1	0		
From this database frequent patterns of occurrence can be							
discovered							

Association Rules Introduction					
Association rules					
 This kind of databases appear for example (market basket analysis) An association rule represents coocurence of database <u>TID Items</u> Bread, Chips, Beer, Yogourt, Eggs Flour, Beer, Eggs, Butter, Bread, Ham, Eggs, Butter, Milk Flour, Eggs, Butter, Chocolate 	on transactional data of events in the ${Flour} \rightarrow {Eggs}$ ${Beer} \rightarrow {Chips}$ ${Bacon} \rightarrow {Eggs}$				
5 Beer, Chips, Bacon, Eggs					
Javier Béjar ©�@ (CS - MAI) Association rules	AMLT - 2016/2017 5 / 54				
Association Rules Definitions					
Definitions (I)					
 We define R as the set of attributes of the database We define X as a subset of attributes from R. We say that X is a pattern from DB if there is any row where all the attributes of X are 1. We define the support (frequency) of a pattern X as the function: 					
$fr(X) = \frac{ M(X, r) }{ r }$					
• Where $ M(X, r) $ is the number of times that X appears in the DB and $ r $ is the size of the BD.					






Properties of frequent sets

- We need ways to prune the search space
- Given X and Y with $Y \subseteq X$:
 - If $fr(Y) \ge fr(X)$, if X is frequent, Y also is frequent
 - If any subset Y from X is not frequent then X is not frequent
- This is known as the **anti-monotone** property of support
- A feasible exploration approach is
 - Find the frequent sets starting from size 1 and increasing its size
 - Prune the candidates including infrequent itemsets



candidates of length l + 1 will be those which all subsets are in the frequent sets of length l.

$$C(\mathcal{F}_{l+1}(r)) = \{X \subseteq R/|X| = l+1 \land \\ \forall_Y (Y \subseteq X \land |Y| = l \Rightarrow Y \in \mathcal{F}_l(r))\}$$

• The computation of association rules can be done iteratively starting with the smallest frequent subsets until no more candidates are obtained













Algorithm

Apriori algorithm



Javier Béjar ©€® (CS - MAI)

Apriori algorithm Computational Cost		
Other pruning strategies		
 Some itemsets are redundant because they have identical support as their supersets We could focus only on those itemsets, reducing the number of candidates Maximal frequent itemsets: A frequent itemset for which none of its immediate supersets are frequent Closed frequent itemsets: A frequent itemset for which none of its immediate supersets have the same support, and its support is larger than <i>minsupport</i> 		
Javier Béjar ⊛⊕®⊚ (CS - MAI) Association rules AMLT - 2016/2017 23 / 54		
FP-Growth		
1 Association Rules		
2 Apriori algorithm		
3 FP-Growth		
Measures of interestingness		
5 Twitter - City Connections		









Javier Béjar ©€© (CS - MAI)

4

 \mathbf{f}

AMLT - 2016/2017









FP-Growth FP-tree		
FP-Growth - Example		
Solve the problem for the predecesors of <i>d</i> , in this case <i>b</i> and <i>c</i> (we are looking if <i>bd</i> and <i>cd</i> are frequent)		
a:2 b:4 c:3 b:2 c:2		
If we continue the algorithm, the patterns extracted would be {[d], [bd], [cd], [abd], [bcd]}		
Javier Béjar @ 🖲 🕲 (CS - MAI) Association rules AMLT - 2016/2017 41 / 54		
1 Association Rules		
2 Apriori algorithm		
3 FP-Growth		
 3 FP-Growth 4 Measures of interestingness 		
 ③ FP-Growth ④ Measures of interestingness ⑤ Twitter - City Connections 		
 FP-Growth Measures of interestingness Twitter - City Connections 		



Measures of interestingness

Measures for association rules interestingness

• Other measures are based on probabilistic independence and correlation

Lift/Interest factor	$N imes rac{f_{11}}{f_{1+} imes f_{+1}}$
All Confidence	$\min(rac{f_{11}}{f_{1+}},rac{f_{11}}{f_{+1}})$
Max Confidence	${\sf m}{f {\sf a}}{\sf x}ig(rac{f_{11}}{f_{1+}},rac{f_{11}}{f_{+1}}ig)$
Kulczynski	$\frac{1}{2}(\frac{f_{11}}{f_{1+}} + \frac{f_{11}}{f_{+1}})$
Cosine Measure	$rac{f_{11}}{\sqrt{f_{1+}f_{+1}}}$











Mining Structured Data













Javier Béjar 📧 🕞 🕥 (CS - MAI)





• Return the candidates that appear more than a number of times










Javier Béjar 📧 🕟 🧿 (CS - MAI)





Clustering of Data Streams - Micro clusters

Aggarwal, Han, Wang, Yu, **On Clustering Massive Data Streams: A Summarization Paradigm** Data Streams - Models and Algorithms, Springer, 2007, 31, 9-38

- On-line phase:
 - Maintains/creates microclusters
 - Number of microclusters is larger than the actual number of clusters of the dataset (approximate the densities)
 - When new data arrives, it is incorporated to a microcluster or generates new microclusters
 - The number of microclusters is fixed, so microclusters are merged to maintain the number
 - Periodically the microclusters are stored
- Off-line phase:
 - Given a time window the stored clusters are used to compute the microclusters inside the time frame
 - A k-means algorithm is used to compute the clusters for the time window

		Sequences	instering of Data Strea	ims		
Clustering of	of Data Str	eams - D	ensity Bas	ed		
Cao, Ester, Qia Stream with N Data Mining, 2	n, Zhou, Densi Ioise Proceedin 006	t y-Based C gs of the Siz	lustering over <th inter<="" siam="" td=""><td>an Evolving Data national Conference or</td><td>١</td></th>	<td>an Evolving Data national Conference or</td> <td>١</td>	an Evolving Data national Conference or	١
 On-line p Core The New outlie Off-line p 	hase: -micro-clusters (weight of a poir examples are m core-mc, sets of potential-mc outlier-mc, sets o er-mc dissapear bhase: Modified	a weighted It fades expo erged and n points with v of points with with time I version of	sum of close p onentially with nc are classified veight over a the n weight below a DBSCAN	oints) time (damping model) l as: reshold a threshold		
Javier Béjar 📧 🕟 🧿) (CS - MAI)	Mining of str	uctures	AMLT - 2016/2017 29	/ 63	
		Sequences C	lustering of Data Strea	ims		
Clustering of	of Data Str	eams - D	ensity Bas	ed		
		me n		Time n+t		
	Cluster	S C C C C C C	Clusters			





Application

Spinal Cord Signals Analysis

Shapes dictionary
Goal: Determine an structure for the peaks
Process

Identify the peaks in the signals (peak finding algorithm)
The peaks extracted from each measuring point will be the first dataset
Clean/Transform the data to visualize the structure
Cluster the datasets to obtain patterns of peaks (k-means ⇒ determine the number of clusters)

Results: A tentative dictionary of peaks shapes

34 / 63



Spinal Cord Signals Analysis

Synchronization Patterns • Goal: Find patterns in the synchronizations of the peaks Process • Determine the synchronizations in the signal (synchronization algorithm) • Generate a database of transactions from the synchronizations • Apply a frequent transactions algorithm to the transactions database • Results: A set of frequent sequences of synchronization











40 / 63









Graph mining Introduction











Graph mining Mining large graphs
Graph partitioning (k-way) - Clustering
 Girvan-Newman Algorithm (Social Networks) Is based on the concept of <i>edge betweenness centrality</i>:
$B(e) = rac{NumConstrainedPaths(e, i, j)}{NumShortPahts(i, j)}$
 Random Walk Betweeness: Compute how often a random walk starting on node <i>i</i> passes through node <i>j</i> Detects bridges among dense components (edges that are not in the shortest paths between pairs of nodes) Algorithm:
 Rank edges by B(e) Delete edge with the highest score Iterate until a specific criteria holds (eg. number of components)
Javier Béjar (B) (CS - MAI) Mining of structures AMLT - 2016/2017 51 / 63
Graph mining Mining large graphs
Graph mining Mining large graphs Dense subgraphs Image: State of the state of
Graph mining Mining large graphs Dense subgraphs
Graph mining Mining large graphs Dense subgraphs • • Not always a complete partitioning of the graph is necessary
 Graph mining Mining large graphs Dense subgraphs Not always a complete partitioning of the graph is necessary Dense subgraphs can represent interesting behaviors
 Graph mining Mining large graphs Dense subgraphs Not always a complete partitioning of the graph is necessary Dense subgraphs can represent interesting behaviors Different types of (pseudo)dense subgraphs:
 Graph mining Mining large graphs Dense subgraphs Not always a complete partitioning of the graph is necessary Dense subgraphs can represent interesting behaviors Different types of (pseudo)dense subgraphs: Clique: All nodes are connected
 Graph mining Mining large graphs Dense subgraphs Not always a complete partitioning of the graph is necessary Dense subgraphs can represent interesting behaviors Different types of (pseudo)dense subgraphs: Clique: All nodes are connected Quasi-clique: Almost all nodes are connected (minimum density or minimum degree)
 Graph mining Mining large graphs Dense subgraphs Not always a complete partitioning of the graph is necessary Dense subgraphs can represent interesting behaviors Different types of (pseudo)dense subgraphs: Clique: All nodes are connected Quasi-clique: Almost all nodes are connected (minimum density or minimum degree) K-core: Every node connects to at least k nodes
 Graph mining Mining large graphs Dense subgraphs Not always a complete partitioning of the graph is necessary Dense subgraphs can represent interesting behaviors Different types of (pseudo)dense subgraphs: Clique: All nodes are connected Quasi-clique: Almost all nodes are connected (minimum density or minimum degree) K-core: Every node connects to at least k nodes K-plex: Each node is missing no more than k-1 edges
 Mining large graphs Dense subgraphs Not always a complete partitioning of the graph is necessary Dense subgraphs can represent interesting behaviors Different types of (pseudo)dense subgraphs: Clique: All nodes are connected Quasi-clique: Almost all nodes are connected (minimum density or minimum degree) K-core: Every node connects to at least k nodes K-plex: Each node is missing no more than k-1 edges
 Graph mining Mining large graphs Dense subgraphs Not always a complete partitioning of the graph is necessary Dense subgraphs can represent interesting behaviors Different types of (pseudo)dense subgraphs: Clique: All nodes are connected Quasi-clique: Almost all nodes are connected (minimum density or minimum degree) K-core: Every node connects to at least k nodes K-plex: Each node is missing no more than k-1 edges Different goals: Minimum size, all or the best ranked, overlapping or not





56 / 63





gSpan - Algorithm

Initialization (D, MinSup)

- sort labels of the vertices and edges in D by frequency
- remove infrequent vertices and edges
- S₀=code of all frequent graphs
 with single edge
- sort S₀ in DFS lexicographic order
- **5** $S = S_0$
- for each code s in S_0
 - **1** gSpan(s, D, MinSup, S)
 - O D = D s

Javier Béjar 🐵 🕥 (CS - MAI)

3 if |D| < MinSup then return

gSpan(s,D,MinSup,S)

- if s! = mincode(s) then return
- 2 insert s into S
- $O C = \emptyset$
- scan D
 - find every edge *e* such that *s* can be right-most extended to frequent *s* * *e*
 - insert s * e into C;
- sort C in DFS lexicographic order
- for each s * e in C do

AMLT - 2016/2017

61 / 63

<page-header><page-header><table-container>

Mining of structures

